

# 驗證Cisco IOS XR和BGP上的路徑MTU探索

## 目錄

[簡介](#)

[背景資訊](#)

[TCP PMTUD和TCP MSS](#)

[案例 — TCP PMTUD已停用](#)

[使用預設MTU值](#)

[使用非預設MTU值 — 活動TCP對等體](#)

[使用非預設MTU值 — 被動TCP對等體](#)

[使用TCP選項 — XR活動](#)

[使用TCP選項 — XR被動](#)

[未直接連線的TCP對等體](#)

[未直接連線的TCP對等體 — 使用TCP選項\(MD5\)](#)

[未直接連線的TCP對等點 — 路徑段具有較低的IP MTU](#)

[案例 — 啟用TCP PMTUD](#)

[啟用PMTUD](#)

[PMTUD — 路徑區段的IP MTU更低](#)

[PMTUD - TCP選項\(MD5\)](#)

[PMTUD — 黑孔偵測](#)

## 簡介

本檔案介紹Cisco IOS® XR裝置上的傳輸控制通訊協定(TCP)路徑最大傳輸單元(MTU)探索(PMTUD)。

## 背景資訊

PMTUD機制會嘗試判斷不需要在兩個主機之間的路徑中的任何位置進行分段的最大網際網路通訊協定(IP)封包大小。建立的值是指定的路徑MTU，並等於每一躍點上MTU值的最小值。如果您在傳輸資訊時考慮路徑MTU，便會讓您充分利用網路容量，並避免分段和傳輸效率。使用邊界閘道通訊協定(BGP)作為使用者端通訊協定（其會逐步顯示PMTUD行為），跨各種情形引入PMTUD機制和執行。

## TCP PMTUD和TCP MSS

TCP利用PMTUD結果來影響本機最大區段大小(MSS)，這表示它會動態調整到發現的路徑MTU。因此，在繼續使用PMTUD之前，您可以快速檢視TCP最大區段大小(MSS)，並瞭解其意義和用途。

根據[RFC879](#)中的MSS原始定義：MSS選項的定義如下：在沒有IP報頭選項的IP資料包中傳輸沒有TCP報頭選項的TCP資料段中，此TCP選項的傳送方可以接收的最大資料八位元數。

闡明一些方面並向實施者提供建議，[RFC6691](#) 突出顯示應如何計算MSS值：

當您計算要放入TCP MSS選項的值時，MTU值應只減小固定IP和TCP標頭的大小，而不應減小以考慮任何可能的IP或TCP選項；反之，傳送者必須減少TCP資料長度，以反映其傳送的封包中包含的任何IP或TCP選項。

有關更詳細的MSS定義，請參閱[IOS XR 6.7.x版Cisco ASR 9000系列路由器路由配置指南](#)：

MSS是電腦或通訊裝置可以在單一未分段TCP區段中接收的最大資料量。所有TCP會話都受單個資料包中可傳輸位元組數的限制；此限制為MSS。在將資料包向下傳遞到IP層之前，TCP會將資料包在傳輸隊列中分割成塊。

TCP MSS值取決於介面的MTU，這是在一個情況下通訊協定可以傳輸的最大資料長度。最大TCP資料包長度由源裝置上出站介面的MTU和目的裝置在TCP設定過程中通告的MSS共同決定。MSS越接近MTU，BGP訊息的傳輸就越有效。每個資料流方向可以使用不同的MSS值。

對於給定TCP作業階段上的MSS，TCP應該考慮的值是什麼？如何計算？

根據RFC879，您有**以下預設**值：主機不能傳送大於576個八位位元組的資料包，除非它們知道目的主機準備接受更大的資料包。TCP最大資料段大小是IP最大資料包大小減去40。

預設IP最大資料包大小為576。

預設TCP最大資料段大小為536。

此值會考慮IP MTU值576位元組。但是如果忽略實際IP MTU值，則TCP MSS計算可以總結如下：

- Active Peer — 計算並傳送帶SYN封包的初始MSS。

```
MSS = IPMTU - sizeof(minimum TCPHDR) - sizeof(minimum IPHDR)
```

Where,

```
sizeof(minimum TCPHDR) = 20 bytes.
```

```
sizeof(minimum IPHDR) = 20 bytes.
```

- 被動對等體 — 計算初始MSS，與從主動對等體接收的MSS進行比較，並以這些MSS值中的較低者傳送SYN、ACK。

```
MIN[IPMTU - sizeof(minimum TCPHDR) - sizeof(minimum IPHDR) , Received MSS value]
```

Where,

```
sizeof(minimum TCPHDR) = 20 bytes.
```

```
sizeof(minimum IPHDR) = 20 bytes.
```

```
Received MSS value = MSS value received with Active Peer TCP SYN.
```

沒有關於MSS選項值的交涉。每個節點確定自己的值，並在TCP會話建立時宣佈相同的值。顯然，如果可從PMTUD推匯出MSS計算時考慮的IP MTU值，則對於指定的路徑MTU，MSS值可調整為最有效的值。Cisco IOS XR行為具有一些有關MSS計算和PMTUD角色的具體資訊，概述如下。

Cisco IOS XR上預設停用PMTUD:

- 根據以下內容，本機初始MSS計算會考慮IP MTU: 如果直接連線的對等體 — 請考慮輸出介面IP MTU。如果是非直接連線的對等體 — 請考慮使用1280位元組的IP MTU。MSS值受已設定的TCP選項影響。

在Cisco IOS XR上啟用PMTUD時：

- 根據以下內容，本機初始MSS計算會考慮IP MTU: 無論直接連線/非直接連線對等點如何 — 請考慮輸出介面IP MTU。MSS值受已設定的TCP選項影響。

關於PMTUD機制和實施的其他詳細資訊需要考慮，本檔案將通過下表中總結的實際例子加以介紹。此表還顯示所考慮的每個場景的主動和被動TCP對等體IP MTU以及所選MSS值。

PMTUD	Scenarios	ACTIVE IP MTU	PASSIVE IP MTU	MSS
Disabled	Using default MTU values	1500	1500	1460
	Using non-default MTU value – Active TCP peer	4460	1500	1460
	Using non-default MTU value – Passive TCP peer	1500	4460	1460
	Using TCP Options (MD5) – XR Active	1500	1500	1436
	Using TCP Options (MD5) – XR Passive	1500	1500	1460
	TCP peers not directly connected	1500	1500	1240
	TCP peers not directly connected – Using TCP Options (MD5)	1500	1500	1216
Enabled	Enabling TCP PMTUD	1500	1500	1460
	PMTUD in action – Path segment has lower MTU	1500	1500	1460
	PMTUD in action – TCP Options (MD5)	1500	1500	1436

## 案例 — TCP PMTUD已停用

### 使用預設MTU值

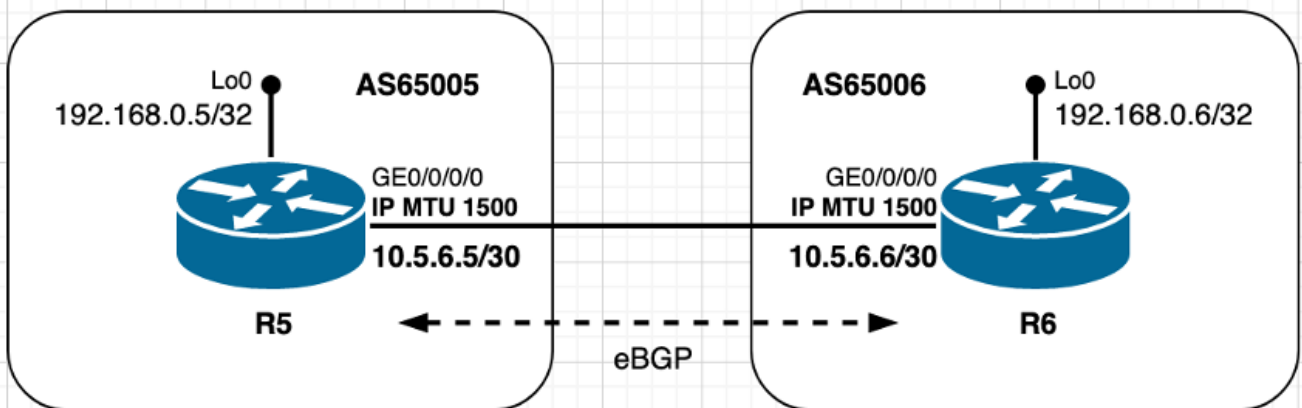


圖2.1.使用預設MTU值

在圖2.1所示的eBGP對等體管理TCP連線的情況下，這意味著，它充當活動角色，並在目標埠179上啟動與R5的TCP會話。對等體是直接連線的，並且兩者都在各自的介面上使用預設IP MTU值。根據本文檔開頭部分共用的資訊，此方案中的MSS計算可以總結如下：

- 兩個節點都使用預設IP MTU 1500位元組
- 預設情況下禁用TCP路徑MTU發現
- TCP對等點直接連線 R6管理BGP連線R6傳送MSS為1460位元組的SYN 1500 ( 介面IP MTU ) — 20(minTCP\_H)- 20(minIP\_H)R5傳送SYN、ACK，MSS為1460位元組 傳送[已接收MSS;本地初始MSS]接收的MSS 1460位元組；本地初始MSS 1460位元組兩個對等點上均使用最低MSS值

## TCP會話詳細資訊，如R6 - ACTIVE上所示：

! - As seen on R6 - ACTIVE

```
RP/0/0/CPU0:R6#show interfaces gigabitEthernet 0/0/0/0
Fri Jan  8 09:35:48.553 UTC
GigabitEthernet0/0/0/0 is up, line protocol is up
Interface state transitions: 1
Hardware is GigabitEthernet, address is fa16.3e85.3dc2 (bia fa16.3e85.3dc2)
Internet address is 10.5.6.6/30
MTU 1514 bytes, BW 1000000 Kbit (Max: 1000000 Kbit)
<snip>
```

```
RP/0/0/CPU0:R6#show tcp brief
Fri Jan  8 09:36:22.491 UTC
PCB      VRF-ID      Recv-Q Send-Q Local Address          Foreign Address         State
<snip>
0x121649fc 0x60000000      0      0 10.5.6.6:24454        10.5.6.5:179           ESTAB
<snip>
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x121649fc
Fri Jan  8 09:37:00.888 UTC
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Fri Jan  8 09:28:28 2021
```

```
PCB 0x121649fc, SO 0x121561b8, TCPCB 0x12156f64, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 78
Local host: 10.5.6.6, Local port: 24454 (Local App PID: 1011918)
Foreign host: 10.5.6.5, Foreign port: 179
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	13	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	10	2	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 3757770712 snduna: 3757770960 sndnxt: 3757770960
sndmax: 3757770960 sndwnd: 32574      sndcwnd: 4380
irs: 1072103647 rcvnxt: 1072103895 rcvwnd: 32593   rcvadv: 1072136488
```

```
SRTT: 155 ms, RTTO: 540 ms, RTV: 385 ms, KRTT: 0 ms
minRTT: 9 ms, maxRTT: 229 ms
```

```
ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE
Connect retries remaining: 30, connect retry interval: 50 secs
```

```
State flags: none
Feature flags: Win Scale, Nagle
Request flags: Win Scale
```

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R6

### TCP會話詳細資訊，如R5上所示 — PASSIVE:

! - As seen on R5 - PASSIVE

RP/0/0/CPU0:R5#show interfaces gigabitEthernet 0/0/0/0  
Fri Jan 8 09:33:04.564 UTC  
GigabitEthernet0/0/0/0 is up, line protocol is up  
Interface state transitions: 1  
Hardware is GigabitEthernet, address is fa16.3ead.518f (bia fa16.3ead.518f)  
Internet address is 10.5.6.5/30  
**MTU 1514 bytes**, BW 1000000 Kbit (Max: 1000000 Kbit)  
<snip>

RP/0/0/CPU0:R5#show tcp brief  
Fri Jan 8 09:33:53.221 UTC

PCB	VRF-ID	Recv-Q	Send-Q	Local Address	Foreign Address	State
<snip>						
0x12155884	0x60000000	0	0	10.5.6.5:179	10.5.6.6:24454	ESTAB
<snip>						

RP/0/0/CPU0:R5#show tcp detail pcb 0x12155884  
Fri Jan 8 09:34:47.317 UTC  
=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Fri Jan 8 09:28:29 2021

PCB 0x12155884, SO 0x1215568c, TCPCB 0x12155a54, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 78  
Local host: 10.5.6.5, Local port: 179 (Local App PID: 1044686)  
Foreign host: 10.5.6.6, Foreign port: 24454

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	9	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	9	7	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 1072103647 snduna: 1072103857 sndnxt: 1072103857
sndmax: 1072103857 sndwnd: 32631 sndcwnd: 4380
irs: 3757770712 rcvnxt: 3757770922 rcvwnd: 32612 rcvadv: 3757803534
```

```
SRTT: 47 ms, RTTO: 300 ms, RTV: 170 ms, KRTT: 0 ms
minRTT: 19 ms, maxRTT: 219 ms
```

```
ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE
Connect retries remaining: 0, connect retry interval: 0 secs
```

```
State flags: none
Feature flags: Win Scale, Nagle
Request flags: Win Scale
```

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

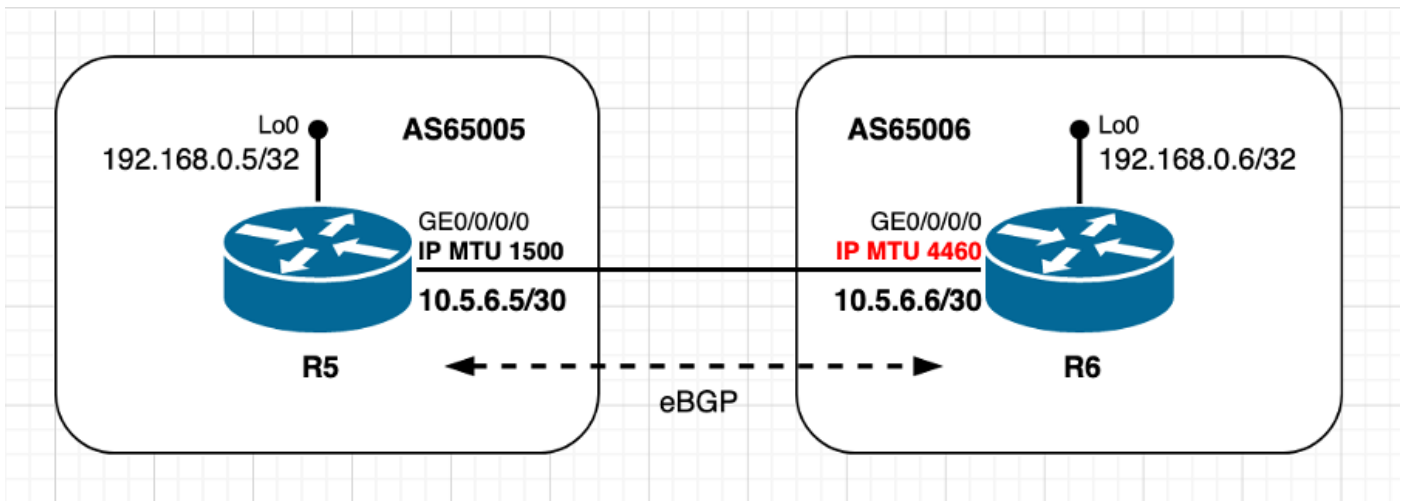
```
Window scales: rcv 0, snd 0, request rcv 0, request snd 0
Timestamp option: recent 0, recent age 0, last ACK sent 0
Sack blocks {start, end}: none
Sack holes {start, end, dups, rxmit}: none
```

```
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0
```

```
PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:
Num Labels: 0 Label Stack:
```

RP/0/0/CPU0:R5#

## 使用非預設MTU值 — 活動TCP對等體



映像2.2 — 活動對等體使用非預設MTU值

此場景與上一場景基本相同，唯一的區別是活動TCP對等體R6現在使用非預設IP MTU值。請注意被動TCP對等體R5是如何對MSS值進行初始計算和決策的。此場景中的TCP MSS計算可以總結如下：

- R6使用非預設IP MTU 4460位元組
- 預設情況下禁用TCP路徑MTU發現
- TCP對等點直接連線 R6管理BGP連線R6傳送MSS為4420位元組的SYN 4460 ( 介面IP MTU )  
— 20(minTCP\_H)- 20(minIP\_H)R5傳送SYN , ACK , MSS為1460位元組 傳送[Received  
MSS;本地初始MSS]已接收4420位元組 ; 本地初始MSS 1460位元組兩個對等點上均使用最低  
MSS值

#### 源自R6的TCP SYN:

! - TCP SYN sourced from R6

```
140      1598.150521      10.5.6.6      10.5.6.5      TCP      62      35502 179 [SYN] Seq=0
Win=16384 Len=0 MSS=4420 WS=1
```

```
Frame 140: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0
Ethernet II, Src: fa:16:3e:85:3d:c2 (fa:16:3e:85:3d:c2), Dst: fa:16:3e:ad:51:8f
(fa:16:3e:ad:51:8f)
```

```
Internet Protocol Version 4, Src: 10.5.6.6, Dst: 10.5.6.5
```

```
Transmission Control Protocol, Src Port: 35502, Dst Port: 179, Seq: 0, Len: 0
```

```
Source Port: 35502
```

```
Destination Port: 179
```

```
[Stream index: 6]
```

```
[TCP Segment Len: 0]
```

```
Sequence number: 0 (relative sequence number)
```

```
Acknowledgment number: 0
```

```
Header Length: 28 bytes
```

```
Flags: 0x002 (SYN)
```

```
Window size value: 16384
```

```
[Calculated window size: 16384]
```

```
Checksum: 0x219d [unverified]
```

```
[Checksum Status: Unverified]
```

```
Urgent pointer: 0
```

```
Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)
```

```
Maximum segment size: 4420 bytes
```

```
Kind: Maximum Segment Size (2)
```

```
Length: 4
```

```
MSS Value: 4420
```

```
Window scale: 0 (multiply by 1)
```

```
End of Option List (EOL)
```

#### TCP SYN、源自R5的ACK:

! - TCP SYN, ACK sourced from R5

```
141      1598.154866      10.5.6.5      10.5.6.6      TCP      62      179 35502 [SYN, ACK] Seq=0
Ack=1 Win=16384 Len=0 MSS=1460 WS=1
```

```
Frame 141: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0
Ethernet II, Src: fa:16:3e:ad:51:8f (fa:16:3e:ad:51:8f), Dst: fa:16:3e:85:3d:c2
(fa:16:3e:85:3d:c2)
```

```
Internet Protocol Version 4, Src: 10.5.6.5, Dst: 10.5.6.6
```

```
Transmission Control Protocol, Src Port: 179, Dst Port: 35502, Seq: 0, Ack: 1, Len: 0
```

```
Source Port: 179
```

```
Destination Port: 35502
```

```
[Stream index: 6]
```

```
[TCP Segment Len: 0]
```

```
Sequence number: 0 (relative sequence number)
```

```
Acknowledgment number: 1 (relative ack number)
```

```
Header Length: 28 bytes
```

```
Flags: 0x012 (SYN, ACK)
```

Window size value: 16384  
[Calculated window size: 16384]  
Checksum: 0xe2b4 [unverified]  
[Checksum Status: Unverified]  
Urgent pointer: 0  
Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)  
Maximum segment size: 1460 bytes  
Kind: Maximum Segment Size (2)  
Length: 4  
**MSS Value: 1460**  
Window scale: 0 (multiply by 1)  
End of Option List (EOL)

## TCP會話詳細資訊，如R6 - ACTIVE上所示：

! - as seen on R6 - Active

```
RP/0/0/CPU0:R6#show interfaces gigabitEthernet 0/0/0/0
Fri Jan  8 09:46:54.138 UTC
GigabitEthernet0/0/0/0 is up, line protocol is up
Interface state transitions: 1
Hardware is GigabitEthernet, address is fa16.3e85.3dc2 (bia fa16.3e85.3dc2)
Internet address is 10.5.6.6/30
MTU 4474 bytes, BW 1000000 Kbit (Max: 1000000 Kbit)
<snip>
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x1215761c
Fri Jan  8 09:56:25.819 UTC
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Fri Jan  8 09:51:46 2021
```

```
PCB 0x1215761c, SO 0x12156f64, TCPCB 0x1216419c, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 886
Local host: 10.5.6.6, Local port: 35502 (Local App PID: 1011918)
Foreign host: 10.5.6.5, Foreign port: 179
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	9	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	6	5	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 764231407  snduna: 764231579  sndnxt: 764231579
sndmax: 764231579  sndwnd: 32650  sndcwnd: 4380
irs: 2712512697  rcvnxt: 2712512869  rcvwnd: 32669  rcvadv: 2712545538
```

```
SRTT: 31 ms, RTTO: 300 ms, RTV: 130 ms, KRTT: 0 ms
minRTT: 9 ms, maxRTT: 239 ms
```

```
ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE
Connect retries remaining: 30, connect retry interval: 50 secs
```

State flags: none



Feature flags: Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 4420, max MSS 4420**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R6#

**TCP會話詳細資訊，如R5上所示 — PASSIVE:**

! - as seen on R5 - Passive

RP/0/0/CPU0:R5#show tcp detail pcb 0x12155a98  
Fri Jan 8 09:55:18.193 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Fri Jan 8 09:51:47 2021

PCB 0x12155a98, SO 0x12153ea0, TCPCB 0x12154e18, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 886  
Local host: 10.5.6.5, Local port: 179 (Local App PID: 1044686)  
Foreign host: 10.5.6.6, Foreign port: 35502

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	6	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	6	1	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 2712512697 snduna: 2712512850 sndnxt: 2712512850  
sndmax: 2712512850 sndwnd: 32688 sndcwnd: 4380  
irs: 764231407 rcvnxt: 764231560 rcvwnd: 32669 rcvadv: 764264229

SRTT: 107 ms, RTTO: 538 ms, RTV: 431 ms, KRTT: 0 ms  
minRTT: 29 ms, maxRTT: 219 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

```

State flags: none
Feature flags: Win Scale, Nagle
Request flags: Win Scale

Datagrams (in bytes): MSS 1460, peer MSS 4420, min MSS 1460, max MSS 1460

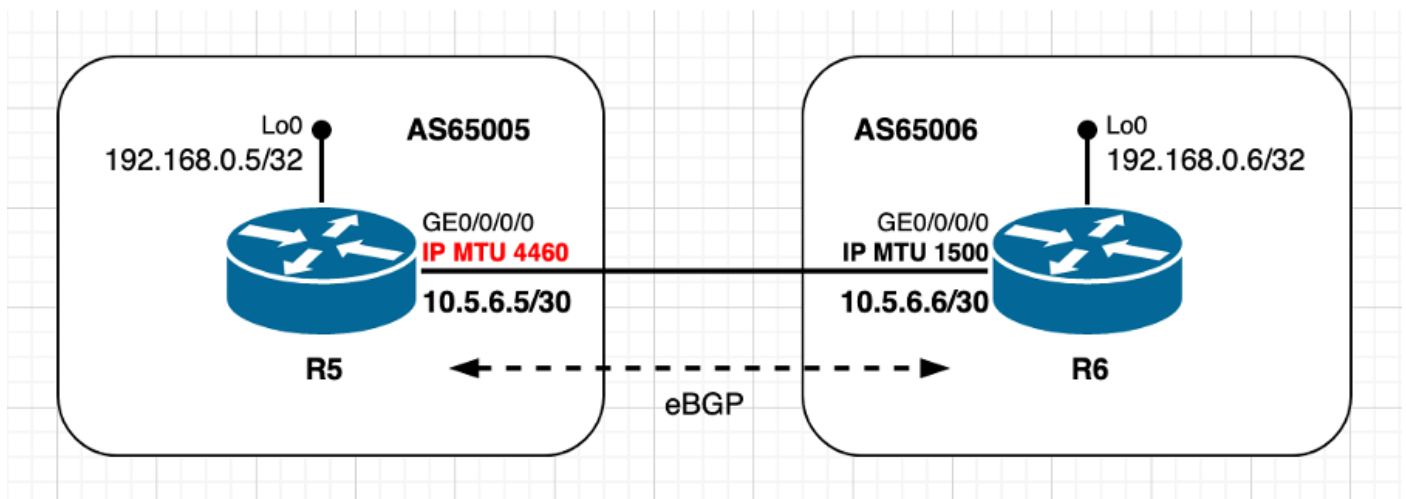
Window scales: rcv 0, snd 0, request rcv 0, request snd 0
Timestamp option: recent 0, recent age 0, last ACK sent 0
Sack blocks {start, end}: none
Sack holes {start, end, dups, rxmit}: none
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R5#

```

## 使用非預設MTU值 — 被動TCP對等體



映像2.3 — 被動對等體使用非預設MTU值。

仍然使用相同的eBGP方案，但現在使用非預設IP MTU配置被動TCP對等體R5，使用預設IP MTU值配置主動TCP對等體R6。與先前的方案一樣，請注意被動對等體R5如何選擇MSS值。此方案中的TCP MSS計算可總結如下：

- R5使用非預設IP MTU 4460位元組
- 預設情況下禁用TCP路徑MTU發現
- TCP對等點直接連線 R6管理BGP連線R6傳送MSS為1460位元組的SYN 1500 ( 介面IP MTU ) —  $20(\text{minTCP\_H}) - 20(\text{minIP\_H})$ R5傳送SYN, ACK, MSS為1460位元組 傳送[Received MSS;本地初始MSS]接收的MSS 1460位元組；本機初始MSS 4420位元組兩個對等點上均使用最低MSS值

源自R6的TCP SYN:

! - TCP SYN sourced from R6

237 2696.666481 10.5.6.6 10.5.6.5 TCP 62 47007 179 [SYN] Seq=0  
Win=16384 Len=0 **MSS=1460** WS=1

Frame 237: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0  
Ethernet II, Src: fa:16:3e:85:3d:c2 (fa:16:3e:85:3d:c2), Dst: fa:16:3e:ad:51:8f  
(fa:16:3e:ad:51:8f)

Internet Protocol Version 4, Src: 10.5.6.6, Dst: 10.5.6.5

Transmission Control Protocol, Src Port: 47007, Dst Port: 179, Seq: 0, Len: 0

Source Port: 47007

Destination Port: 179

[Stream index: 10]

[TCP Segment Len: 0]

Sequence number: 0 (relative sequence number)

Acknowledgment number: 0

Header Length: 28 bytes

Flags: 0x002 (SYN)

Window size value: 16384

[Calculated window size: 16384]

Checksum: 0x2025 [unverified]

[Checksum Status: Unverified]

Urgent pointer: 0

Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)

Maximum segment size: 1460 bytes

Kind: Maximum Segment Size (2)

Length: 4

**MSS Value: 1460**

Window scale: 0 (multiply by 1)

End of Option List (EOL)

## TCP SYN、源自R5的ACK:

! - TCP SYN, ACK sourced from R5

238 2696.702792 10.5.6.5 10.5.6.6 TCP 62 179 47007 [SYN, ACK] Seq=0  
Ack=1 Win=16384 Len=0 **MSS=1460** WS=1

Frame 238: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0  
Ethernet II, Src: fa:16:3e:ad:51:8f (fa:16:3e:ad:51:8f), Dst: fa:16:3e:85:3d:c2  
(fa:16:3e:85:3d:c2)

Internet Protocol Version 4, Src: 10.5.6.5, Dst: 10.5.6.6

Transmission Control Protocol, Src Port: 179, Dst Port: 47007, Seq: 0, Ack: 1, Len: 0

Source Port: 179

Destination Port: 47007

[Stream index: 10]

[TCP Segment Len: 0]

Sequence number: 0 (relative sequence number)

Acknowledgment number: 1 (relative ack number)

Header Length: 28 bytes

Flags: 0x012 (SYN, ACK)

Window size value: 16384

[Calculated window size: 16384]

Checksum: 0x7078 [unverified]

[Checksum Status: Unverified]

Urgent pointer: 0

Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)

Maximum segment size: 1460 bytes

Kind: Maximum Segment Size (2)

Length: 4

**MSS Value: 1460**

Window scale: 0 (multiply by 1)

End of Option List (EOL)

TCP會話詳細資訊，如R6 - ACTIVE上所示：

! - as seen on R6 - Active

RP/0/0/CPU0:R6#show tcp detail pcb 0x1215761c

Fri Jan 8 10:15:20.351 UTC

=====

Connection state is ESTAB, I/O status: 0, socket status: 0

Established at Fri Jan 8 10:10:04 2021

PCB 0x1215761c, SO 0x12162aac, TCPCB 0x12156f64, vrfid 0x60000000,

Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 103

Local host: 10.5.6.6, Local port: 47007 (Local App PID: 1011918)

Foreign host: 10.5.6.5, Foreign port: 179

Current send queue size in bytes: 0 (max 24576)

Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes

Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	10	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	7	5	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 3949093168 snduna: 3949093359 sndnxt: 3949093359  
sndmax: 3949093359 sndwnd: 32631 sndcwnd: 4380  
irs: 54439005 rcvnxt: 54439196 rcvwnd: 32650 rcvadv: 54471846

SRTT: 75 ms, RTTO: 459 ms, RTV: 384 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 30, connect retry interval: 50 secs

State flags: none  
Feature flags: Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none  
Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R6#

TCP會話詳細資訊，如R5上所示 — PASSIVE:

! - as seen on R5 - Passive

```
RP/0/0/CPU0:R5#show interfaces gigabitEthernet 0/0/0/0
Fri Jan  8 10:10:39.110 UTC
GigabitEthernet0/0/0/0 is up, line protocol is up
  Interface state transitions: 1
  Hardware is GigabitEthernet, address is fa16.3ead.518f (bia fa16.3ead.518f)
  Internet address is 10.5.6.5/30
  MTU 4474 bytes, BW 1000000 Kbit (Max: 1000000 Kbit)
<snip>
```

```
RP/0/0/CPU0:R5#show tcp detail pcb 0x121550fc
Fri Jan  8 10:14:20.105 UTC
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Fri Jan  8 10:10:05 2021
```

```
PCB 0x121550fc, SO 0x12154e18, TCPCB 0x12154304, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 103
Local host: 10.5.6.5, Local port: 179 (Local App PID: 1044686)
Foreign host: 10.5.6.6, Foreign port: 47007
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768)  mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	7	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	7	2	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
  iss: 54439005  snduna: 54439177  sndnxt: 54439177
sndmax: 54439177  sndwnd: 32669  sndcwnd: 4380
  irs: 3949093168  rcvnxt: 3949093340  rcvwnd: 32650  rcvadv: 3949125990
```

```
SRTT: 117 ms,  RTTO: 570 ms,  RTV: 453 ms,  KRTT: 0 ms
minRTT: 19 ms,  maxRTT: 229 ms
```

```
ACK hold time: 200 ms,  Keepalive time: 0 sec,  SYN waittime: 30 sec
Giveup time: 0 ms,  Retransmission retries: 0,  Retransmit forever: FALSE
Connect retries remaining: 0,  connect retry interval: 0 secs
```

```
State flags: none
Feature flags: Win Scale, Nagle
Request flags: Win Scale
```

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 4420, max MSS 4420**

```
Window scales: rcv 0, snd 0, request rcv 0, request snd 0
Timestamp option: recent 0, recent age 0, last ACK sent 0
Sack blocks {start, end}: none
Sack holes {start, end, rxmit}: none
```

```
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
```

```
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0
```

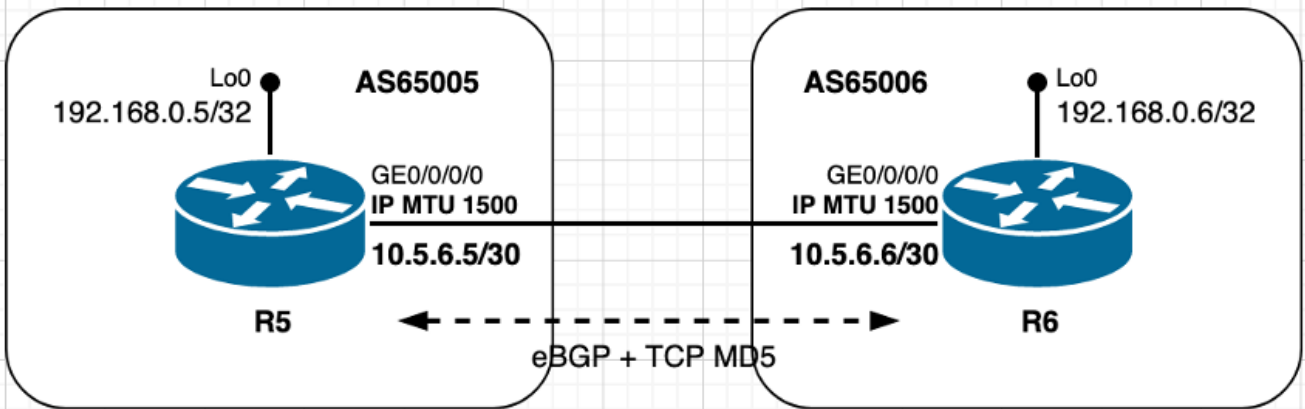
```
PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:
Num Labels: 0 Label Stack:
```

```
RP/0/0/CPU0:R5#
```

## 使用TCP選項 — XR活動

如本檔案前面所述，TCP選項(例如[TCP MD5](#)、[TCP selective-ack](#)或[TCP timestamps](#))的使用會影響MSS計算，因為這些選項會導致在MSS計算中考慮其他位元組。

本節以及下一個目的是說明存在TCP選項時對等點進行的MSS計算。例如，使用TCP MD5身份驗證選項。請參閱映像2.4中的參考案例，如下圖所示。



影象2.4 — 使用TCP選項(MD5)- XR活動。

在此方案中，兩個對等體都使用預設IP MTU值、直接連線，並且對等體R6扮演TCP主動角色。已共用TCP MD5身份驗證帳戶的配置和使用，從而產生額外開銷。在此特定案例中，TCP MSS計算可以總結如下：

- 兩個節點都使用預設IP MTU 1500位元組
- 預設情況下禁用TCP路徑MTU發現
- TCP對等點直接連線
- 兩個節點上都啟用了TCP MD5身份驗證 R6管理BGP連線R6傳送MSS為1436位元組的SYN 1500 ( 介面IP MTU ) — 20(minTCP\_H)- 20(minIP\_H)- 24位元組 ( IOS XR TCP選項額外負荷 ) R5傳送SYN, ACK, MSS為1436位元組 傳送[Received MSS;本地初始MSS]收到1436位元組；本地初始MSS 1460位元組兩個對等點上均使用最低MSS值

從摘要中可看出，Cisco IOS XR的行為方式並不嚴格符合[RFC 879](#)和[RFC 6691](#)，後者指出TCP選項不應計入MSS計算。

tcp標頭長度的額外因素的Cisco IOS XR帳戶進一步記錄在Cisco錯誤ID [CSCvf20166](#)中：

"(...)XR啟動BGP連線時，BGP首先建立套接字，然後設定套接字選項，包括MD5。這會使tcp選項標頭長度= 24。因此，初始MSS變為1500 - 40 - 24 = 1436。這被傳送到對等項和對等項使用 min(1436, 1460)= 1436.(..)

## 源自R6的TCP SYN:

! - TCP SYN sourced from R6

```
430      5775.839420    10.5.6.6      10.5.6.5      TCP      82      24785  179 [SYN] Seq=0
Win=16384 Len=0 MSS=1436 WS=1
```

Frame 430: 82 bytes on wire (656 bits), 82 bytes captured (656 bits) on interface 0  
Ethernet II, Src: fa:16:3e:85:3d:c2 (fa:16:3e:85:3d:c2), Dst: fa:16:3e:ad:51:8f  
(fa:16:3e:ad:51:8f)

Internet Protocol Version 4, Src: 10.5.6.6, Dst: 10.5.6.5

Transmission Control Protocol, Src Port: 24785, Dst Port: 179, Seq: 0, Len: 0

Source Port: 24785

Destination Port: 179

[Stream index: 14]

[TCP Segment Len: 0]

Sequence number: 0 (relative sequence number)

Acknowledgment number: 0

Header Length: 48 bytes

Flags: 0x002 (SYN)

Window size value: 16384

[Calculated window size: 16384]

Checksum: 0xd62b [unverified]

[Checksum Status: Unverified]

Urgent pointer: 0

Options: (28 bytes), Maximum segment size, Window scale, No-Operation (NOP), **TCP MD5**

**signature**, End of Option List (EOL)

Maximum segment size: 1436 bytes

Kind: Maximum Segment Size (2)

Length: 4

**MSS Value: 1436**

Window scale: 0 (multiply by 1)

No-Operation (NOP)

TCP MD5 signature

End of Option List (EOL)

## TCP SYN、源自R5的ACK:

! - TCP SYN, ACK sourced from R5

```
431      5775.845744    10.5.6.5      10.5.6.6      TCP      82      179  24785 [SYN, ACK] Seq=0
Ack=1 Win=16384 Len=0 MSS=1436 WS=1
```

Frame 431: 82 bytes on wire (656 bits), 82 bytes captured (656 bits) on interface 0  
Ethernet II, Src: fa:16:3e:ad:51:8f (fa:16:3e:ad:51:8f), Dst: fa:16:3e:85:3d:c2  
(fa:16:3e:85:3d:c2)

Internet Protocol Version 4, Src: 10.5.6.5, Dst: 10.5.6.6

Transmission Control Protocol, Src Port: 179, Dst Port: 24785, Seq: 0, Ack: 1, Len: 0

Source Port: 179

Destination Port: 24785

[Stream index: 14]

[TCP Segment Len: 0]

Sequence number: 0 (relative sequence number)

Acknowledgment number: 1 (relative ack number)

Header Length: 48 bytes

Flags: 0x012 (SYN, ACK)

Window size value: 16384

[Calculated window size: 16384]

Checksum: 0xe83d [unverified]

[Checksum Status: Unverified]

Urgent pointer: 0

Options: (28 bytes), Maximum segment size, Window scale, No-Operation (NOP), **TCP MD5 signature**, End of Option List (EOL)  
Maximum segment size: 1436 bytes  
Kind: Maximum Segment Size (2)  
Length: 4  
**MSS Value: 1436**  
Window scale: 0 (multiply by 1)  
No-Operation (NOP)  
TCP MD5 signature  
End of Option List (EOL)

## TCP會話詳細資訊，如R6 - ACTIVE上所示：

! - as seen on R6 - Active

RP/0/0/CPU0:R6#show tcp detail pcb 0x1215761c

Fri Jan 8 11:14:13.599 UTC

=====

Connection state is ESTAB, I/O status: 0, socket status: 0

Established at Fri Jan 8 11:01:21 2021

PCB 0x1215761c, SO 0x1216419c, TCPCB 0x121649fc, vrfid 0x60000000,

Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 409

Local host: 10.5.6.6, Local port: 24785 (Local App PID: 1011918)

Foreign host: 10.5.6.5, Foreign port: 179

Current send queue size in bytes: 0 (max 24576)

Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes

Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	17	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	14	13	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 1379482495 snduna: 1379482819 sndnxt: 1379482819

sndmax: 1379482819 sndwnd: 32498 sndcwnd: 4308

irs: 3750694052 rcvnx: 3750694376 rcvwnd: 32517 rcvadv: 3750726893

SRTT: 55 ms, RTTO: 300 ms, RTV: 176 ms, KRTT: 0 ms

minRTT: 9 ms, maxRTT: 259 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec

Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE

Connect retries remaining: 30, connect retry interval: 50 secs

State flags: none

Feature flags: **MD5**, Win Scale, Nagle

Request flags: Win Scale

**Datagrams (in bytes): MSS 1436, peer MSS 1436, min MSS 1436, max MSS 1436**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0

Timestamp option: recent 0, recent age 0, last ACK sent 0

Sack blocks {start, end}: none

Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO



Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:

#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R6#

TCP會話詳細資訊，如R5上所示 — PASSIVE:

! - as seen on R5 - Passive

RP/0/0/CPU0:R5#show tcp detail pcb 0x12155d04

Fri Jan 8 11:12:51.984 UTC

=====

Connection state is ESTAB, I/O status: 0, socket status: 0

Established at Fri Jan 8 11:01:22 2021

PCB 0x12155d04, SO 0x12154e18, TCPCB 0x12154304, vrfid 0x60000000,

Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 409

Local host: 10.5.6.5, Local port: 179 (Local App PID: 1044686)

Foreign host: 10.5.6.6, Foreign port: 24785

Current send queue size in bytes: 0 (max 24576)

Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes

Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	14	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	14	3	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 3750694052 snduna: 3750694357 sndnxt: 3750694357

sndmax: 3750694357 sndwnd: 32536 sndcwnd: 4308

irs: 1379482495 rcvnxt: 1379482800 rcvwnd: 32517 rcvadv: 1379515317

SRTT: 181 ms, RTTO: 443 ms, RTV: 262 ms, KRTT: 0 ms

minRTT: 29 ms, maxRTT: 219 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec

Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE

Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none

Feature flags: MD5, Win Scale, Nagle

Request flags: Win Scale

**Datagrams (in bytes): MSS 1436, peer MSS 1436, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0

Timestamp option: recent 0, recent age 0, last ACK sent 0

Sack blocks {start, end}: none

Sack holes {start, end, dups, rxmit}: none

```
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0
```

```
PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:
Num Labels: 0 Label Stack:
```

```
RP/0/0/CPU0:R5#
```

其他TCP選項也存在類似行為，在配置時這些選項會增加額外開銷並影響Cisco IOS XR中的MSS計算。請考慮配置TCP時間戳和TCP選擇性確認選項時記錄MSS計算的相同方案和這些示例。

TCP會話詳細資訊，如R6上所示 — ACTIVE — 已配置TCP選項時間戳和選擇性確認選項：

```
! - as seen on R6 - Active
! -- tcp timestamp configured
! -- 12 bytes of additional overhead
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x1539c844
```

```
<snip>
Feature flags: Timestamp, Win Scale, Nagle
Request flags: Timestamp, Win Scale
```

```
Datagrams (in bytes): MSS 1448, peer MSS 1448, min MSS 1448, max MSS 1448
<snip>
```

```
! - as seen on R6 - Active
! -- tcp selective-ack configured
! -- 36 bytes of additional overhead
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x1539df38
```

```
<snip>
Feature flags: Sack, Win Scale, Nagle
Request flags: Sack, Win Scale
```

```
Datagrams (in bytes): MSS 1424, peer MSS 1424, min MSS 1424, max MSS 1424
<snip>
```

```
! - as seen on R6 - Active
! -- tcp selective-ack and tcp timestamp configured
! -- 40 bytes of additional overhead
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x1539e130
```

```
<snip>
State flags: none
Feature flags: Sack, Timestamp, Win Scale, Nagle
Request flags: Sack, Timestamp, Win Scale
```

```
Datagrams (in bytes): MSS 1420, peer MSS 1420, min MSS 1420, max MSS 1420
<snip>
```

```
! - as seen on R6 - Active
! -- MD5 and tcp selective-ack configured
! -- 36 bytes of additional overhead
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x1539b3cc
```

```
<snip>
```

```
Feature flags: Sack, MD5, Win Scale, Nagle
Request flags: Sack, Win Scale
```

```
Datagrams (in bytes): MSS 1424, peer MSS 1424, min MSS 1424, max MSS 1424
<snip>
```

```
! - as seen on R6 - Active
! -- MD5 and tcp timestamp configured
! -- 36 bytes of additional overhead
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x15397b4c
<snip>
```

```
Feature flags: MD5, Timestamp, Win Scale, Nagle
Request flags: Timestamp, Win Scale
```

```
Datagrams (in bytes): MSS 1424, peer MSS 1424, min MSS 1424, max MSS 1424
<snip>
```

```
! - as seen on R6 - Active
! -- MD5, tcp timestamp, and tcp selective-ack configured
! -- 40 bytes of additional overhead
```

```
RP/0/0/CPU0:R6#show tcp detail pcb 0x1539a4cc
<snip>
```

```
State flags: none
Feature flags: MD5, Timestamp, Win Scale, Nagle
Request flags: Timestamp, Win Scale
```

```
Datagrams (in bytes): MSS 1420, peer MSS 1420, min MSS 1420, max MSS 1420
<snip>
```

## 使用TCP選項 — XR被動

從先前的場景中，您可能已經注意到Cisco IOS XR節點在初始MSS計算方面處於被動角色時的不同行為。節點不考慮**tcp選項標頭長度**。此案例旨在突出顯示此截然不同的行為，思科錯誤ID也對此進行了描述：

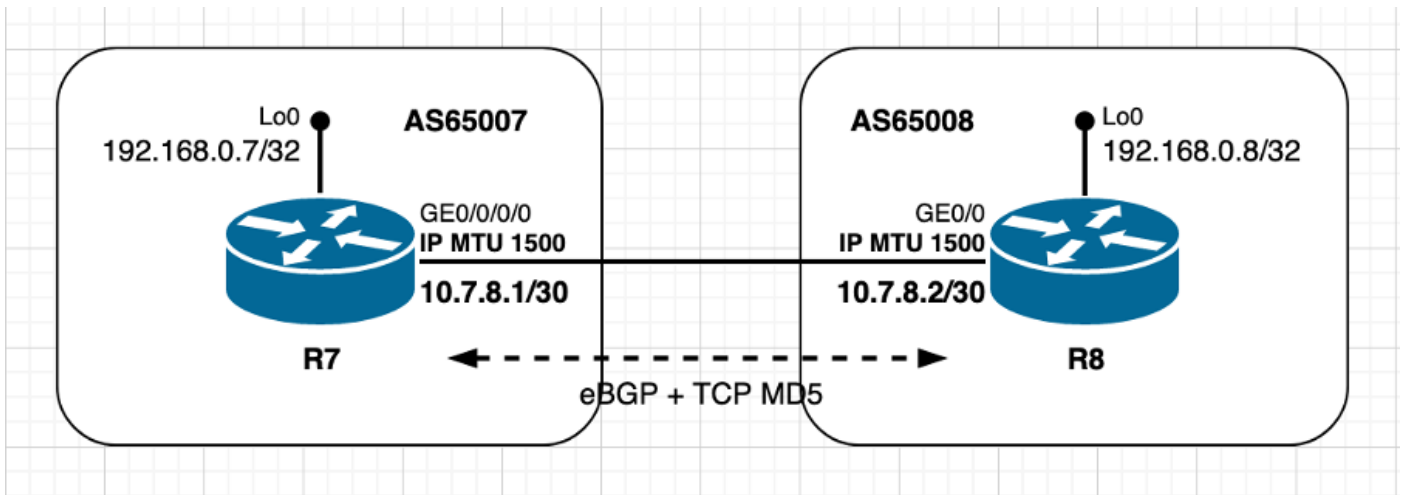
"(...) — 當對等體發起連線時，它會將初始MSS作為1460傳送。XR TCP建立socket、pcb等，然後按指定順序採取以下兩個操作：

— 首先，它在減去**tcp選項標頭長度**後計算初始MSS。這是「0」，因為MD5選項尚未從偵聽套接字繼承到此套接字。

— 然後，它繼承「MD5」和其他選項，這使得「選項標頭位元組長度」為24。

因此，在此案例中，XR TCP將初始MSS傳送為1460，因此供兩者使用。 (...)」

在此案例中，雖然活動的TCP對等體R8是Cisco IOS節點，但此事實不會帶來任何差異，也不會帶來任何場景要強調的細節。但是，有趣的是，請注意，與上一節場景所顯示的Cisco IOS XR不同，此處活動TCP對等體R8在初始MSS計算時不考慮TCP選項。



影象2.5 — 使用TCP選項(MD5)- XR被動。

兩個對等點使用預設的IP MTU值，且均已直接連線。Cisco IOS對等體R8扮演主動角色。此案例中的TCP MSS計算可總結如下：

- 兩個節點都使用預設IP MTU 1500位元組
- Cisco IOS XR R7預設禁用TCP路徑MTU發現
- Cisco IOS R8上預設啟用TCP路徑MTU發現
- TCP對等點直接連線
- 兩個節點上都啟用了TCP MD5身份驗證 IOS R8管理BGP連線IOS R8傳送MSS為1460位元組的SYN 1500 ( 介面IP MTU ) — 20(minTCP\_H)- 20(minIP\_H)IOS XR R7傳送SYN、ACK，MSS為1460位元組 傳送[Received MSS;本地初始MSS]接收的MSS 1460位元組；本地初始MSS 1460位元組兩個對等點上均使用最低MSS值

源自R8的TCP SYN - Cisco IOS:

! - TCP SYN sourced from R8

```
96      5.907127      10.7.8.2      10.7.8.1      TCP      78      52975  179 [SYN] Seq=0
Win=16384 Len=0  MSS=1460
```

```
Frame 96: 78 bytes on wire (624 bits), 78 bytes captured (624 bits) on interface 0
Ethernet II, Src: fa:16:3e:58:21:ba (fa:16:3e:58:21:ba), Dst: fa:16:3e:68:d9:e5
(fa:16:3e:68:d9:e5)
```

```
Internet Protocol Version 4, Src: 10.7.8.2, Dst: 10.7.8.1
```

```
Transmission Control Protocol, Src Port: 52975, Dst Port: 179, Seq: 0, Len: 0
```

```
Source Port: 52975
```

```
Destination Port: 179
```

```
[Stream index: 3]
```

```
[TCP Segment Len: 0]
```

```
Sequence number: 0 (relative sequence number)
```

```
Acknowledgment number: 0
```

```
Header Length: 44 bytes
```

```
Flags: 0x002 (SYN)
```

```
Window size value: 16384
```

```
[Calculated window size: 16384]
```

```
Checksum: 0xb612 [unverified]
```

```
[Checksum Status: Unverified]
```

```
Urgent pointer: 0
```

```
Options: (24 bytes), Maximum segment size, TCP MD5 signature, End of Option List (EOL)
```

```
Maximum segment size: 1460 bytes
```

```
Kind: Maximum Segment Size (2)
```

```
Length: 4
```

**MSS Value: 1460**

TCP MD5 signature

End of Option List (EOL)

## TCP SYN、源自R7的ACK - Cisco IOS XR:

! - TCP SYN,ACK sourced from R7

```
97      0.003446      10.7.8.1      10.7.8.2      TCP      78      179 52975 [SYN, ACK] Seq=0
Ack=1 Win=16384 Len=0 MSS=1460
```

Frame 97: 78 bytes on wire (624 bits), 78 bytes captured (624 bits) on interface 0  
Ethernet II, Src: fa:16:3e:68:d9:e5 (fa:16:3e:68:d9:e5), Dst: fa:16:3e:58:21:ba  
(fa:16:3e:58:21:ba)

Internet Protocol Version 4, Src: 10.7.8.1, Dst: 10.7.8.2

Transmission Control Protocol, Src Port: 179, Dst Port: 52975, Seq: 0, Ack: 1, Len: 0

Source Port: 179

Destination Port: 52975

[Stream index: 3]

[TCP Segment Len: 0]

Sequence number: 0 (relative sequence number)

Acknowledgment number: 1 (relative ack number)

Header Length: 44 bytes

Flags: 0x012 (SYN, ACK)

Window size value: 16384

[Calculated window size: 16384]

Checksum: 0xfb47 [unverified]

[Checksum Status: Unverified]

Urgent pointer: 0

Options: (24 bytes), Maximum segment size, **TCP MD5 signature**, End of Option List (EOL)

Maximum segment size: 1460 bytes

Kind: Maximum Segment Size (2)

Length: 4

**MSS Value: 1460**

TCP MD5 signature

End of Option List (EOL)

## TCP會話詳細資訊，如R8上所示 — Cisco IOS - ACTIVE:

! - as seen from R8 - Cisco IOS

R8#show ip bgp neighbors

BGP neighbor is 10.7.8.1, remote AS 65007, external link

BGP version 4, remote router ID 192.168.0.7

BGP state = Established, up for 00:06:12

Last read 00:00:16, last write 00:00:16, hold time is 180, keepalive interval is 60 seconds

Neighbor sessions:

1 active, is not multiseession capable (disabled)

Neighbor capabilities:

Route refresh: advertised and received(new)

Four-octets ASN Capability: advertised and received

Address family IPv4 Unicast: advertised and received

Enhanced Refresh Capability: advertised

Multiseession Capability:

Stateful switchover support enabled: NO for session 1

Message statistics:

InQ depth is 0

OutQ depth is 0

	Sent	Rcvd
Opens:	1	1
Notifications:	0	0

```

Updates:          1          1
Keepalives:      7          7
Route Refresh:   0          0
Total:           9          9

```

Do log neighbor state changes (via global configuration)  
Default minimum time between advertisement runs is 30 seconds

For address family: IPv4 Unicast  
Session: 10.7.8.1  
BGP table version 1, neighbor version 1/0  
Output queue size : 0  
Index 6, Advertise bit 0  
6 update-group member  
Slow-peer detection is disabled  
Slow-peer split-update-group dynamic is disabled

	Sent	Rcvd
Prefix activity:	----	----
Prefixes Current:	0	0
Prefixes Total:	0	0
Implicit Withdraw:	0	0
Explicit Withdraw:	0	0
Used as bestpath:	n/a	0
Used as multipath:	n/a	0
Used as secondary:	n/a	0

	Outbound	Inbound
Local Policy Denied Prefixes:	-----	-----
Total:	0	0

Number of NLRI in the update sent: max 0, min 0

Last detected as dynamic slow peer: never  
Dynamic slow peer recovered: never  
Refresh Epoch: 1  
Last Sent Refresh Start-of-rib: never  
Last Sent Refresh End-of-rib: never  
Last Received Refresh Start-of-rib: never  
Last Received Refresh End-of-rib: never

	Sent	Rcvd
Refresh activity:	----	----
Refresh Start-of-RIB	0	0
Refresh End-of-RIB	0	0

Address tracking is enabled, the RIB does have a route to 10.7.8.1  
Connections established 6; dropped 5  
Last reset 00:06:18, due to BGP Notification received of session 1, Administrative Reset  
External BGP neighbor configured for connected checks (single-hop no-disable-connected-check)  
Interface associated: GigabitEthernet0/1 (peering address in same link)

**Transport(tcp) path-mtu-discovery is enabled**

Graceful-Restart is disabled  
SSO is disabled

Connection state is ESTAB, I/O status: 1, unread input bytes: 0  
Connection is ECN Disabled, Minimum incoming TTL 0, Outgoing TTL 1  
Local host: 10.7.8.2, Local port: 52975  
Foreign host: 10.7.8.1, Foreign port: 179  
Connection tableid (VRF): 0  
Maximum output segment queue size: 50

Enqueued packets for retransmit: 0, input: 0 mis-ordered: 0 (0 bytes)

Event Timers (current time is 0x15DD97):

Timer	Starts	Wakeups	Next
Retrans	10	0	0x0
TimeWait	0	0	0x0
AckHold	9	5	0x0
SendWnd	0	0	0x0

```
KeepAlive      0          0          0x0
GiveUp         0          0          0x0
PmtuAger      1          0        0x195465
DeadWait      0          0          0x0
Linger        0          0          0x0
ProcessQ      0          0          0x0
```

```
iss: 1154289541  snduna: 1154289755  sndnxt: 1154289755
irs: 2149897425  rcvnxt: 2149897635
```

```
sndwnd: 32612  scale:      0  maxrcvwnd: 16384
rcvwnd: 16175  scale:      0  delrcvwnd: 209
```

```
SRTT: 737 ms, RTTO: 2506 ms, RTV: 1769 ms, KRTT: 0 ms
minRTT: 7 ms, maxRTT: 1000 ms, ACK hold: 200 ms
uptime: 372981 ms, Sent idletime: 16648 ms, Receive idletime: 16431 ms
Status Flags: active open
Option Flags: nagle, path mtu capable, md5
IP Precedence value : 6
```

**Datagrams (max data segment is 1460 bytes):**

```
Rcvd: 18 (out of order: 0), with data: 8, total data bytes: 209
Sent: 16 (retransmit: 0, fastretransmit: 0, partialack: 0, Second Congestion: 0), with data: 9,
total data bytes: 213
```

```
Packets received in fast path: 0, fast processed: 0, slow path: 0
fast lock acquisition failures: 0, slow path: 0
TCP Semaphore      0x0FBFA8A4  FREE
```

R8#

**TCP會話詳細資訊，如R7上所示 — Cisco IOS XR - PASSIVE:**

! - as seen from R7 - Cisco IOS XR

```
RP/0/0/CPU0:R7#show tcp detail pcb 0x12152e48
Wed Jan 13 13:03:43.363 UTC
```

```
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Wed Jan 13 12:58:16 2021
```

```
PCB 0x12152e48, SO 0x1213c130, TCPCB 0x12156060, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 1, Hash index: 947
Local host: 10.7.8.1, Local port: 179 (Local App PID: 983244)
Foreign host: 10.7.8.2, Foreign port: 52975
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768)  mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	8	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	8	7	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 2149897425  snduna: 2149897616  sndnxt: 2149897616
sndmax: 2149897616  sndwnd: 16194  sndcwnd: 4380
irs: 1154289541  rcvnxt: 1154289736  rcvwnd: 32631  rcvadv: 1154322367
```

SRTT: 125 ms, RTTO: 552 ms, RTV: 427 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 229 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: MD5, Nagle  
Request flags: none

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

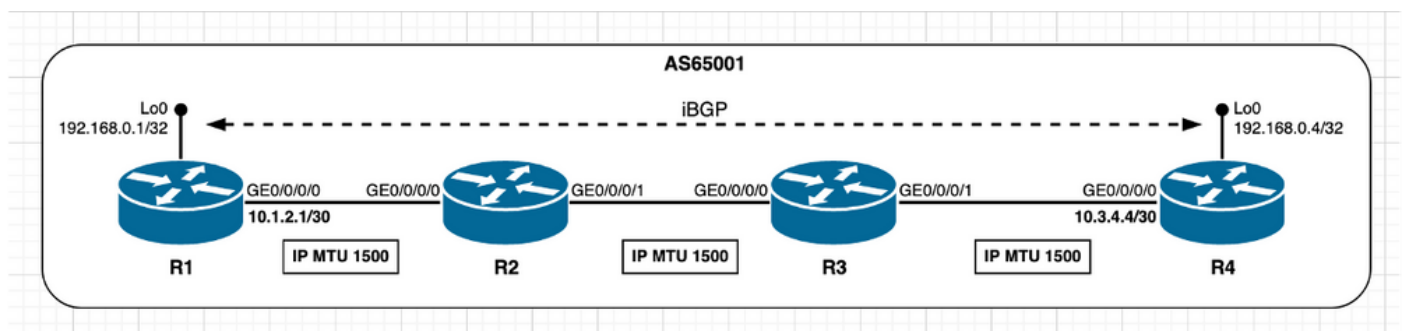
Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R7#

## 未直接連線的TCP對等體

當對等點沒有直接連線時，完成TCP MSS初始計算的方式將更改，如本文檔的介紹部分中所述。使用設定為預設IP MTU值的所有對等點的iBGP作業階段的案例，來逐步執行MSS計算。



映像2.6 — 未直接連線的TCP對等體 — iBGP。

需要注意的重要一點是，當禁用TCP路徑MTU發現且對等體未直接連線時，根據設計，Cisco IOS XR使用固定IP MTU值1280位元組。

在上圖中，R4扮演活動角色並管理TCP連線，R4在目的地埠179上開啟與R1的TCP會話。兩個節點在其介面上均使用預設IP MTU值。此方案中的MSS計算可概述如下：

- 所有節點均使用預設IP MTU 1500位元組
- 預設情況下禁用TCP路徑MTU發現
- TCP對等點未直接連線 R4管理BGP連線R4傳送MSS為1240位元組的SYN 當對等點未直接連線



且TCP路徑MTU探索已停用時，不會考慮介面MTU根據Cisco IOS XR設計，1280位元組被認為是TCP\_DEFAULT\_MTU1280(TCP\_DEFAULT\_MTU)- 20(minTCP\_H)- 20(minIP\_H)R1傳送SYN，ACK，MSS為1240位元組 傳送[已接收MSS;本地初始MSS]接收的MSS 1240位元組；本地初始MSS 1240位元組兩個對等點上均使用最低MSS值

來源為R4的TCP SYN:

! - TCP SYN sourced from R4

```
194      434.274181      192.168.0.4 192.168.0.1 TCP      62      37740 179 [SYN] Seq=0 Win=16384
Len=0 MSS=1240 WS=1
```

Frame 194: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0  
Ethernet II, Src: fa:16:3e:d7:7e:f6 (fa:16:3e:d7:7e:f6), Dst: fa:16:3e:8f:8f:54  
(fa:16:3e:8f:8f:54)

Internet Protocol Version 4, Src: 192.168.0.4, Dst: 192.168.0.1

Transmission Control Protocol, Src Port: 37740, Dst Port: 179, Seq: 0, Len: 0

Source Port: 37740

Destination Port: 179

[Stream index: 7]

[TCP Segment Len: 0]

Sequence number: 0 (relative sequence number)

Acknowledgment number: 0

Header Length: 28 bytes

Flags: 0x002 (SYN)

Window size value: 16384

[Calculated window size: 16384]

Checksum: 0x8643 [unverified]

[Checksum Status: Unverified]

Urgent pointer: 0

Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)

Maximum segment size: 1240 bytes

Kind: Maximum Segment Size (2)

Length: 4

**MSS Value: 1240**

Window scale: 0 (multiply by 1)

End of Option List (EOL)

源自R1的TCP SYN、ACK:

! - TCP SYN,ACK sourced from R1

```
195      434.277985      192.168.0.1 192.168.0.4 TCP      62      179 37740 [SYN, ACK] Seq=0 Ack=1
Win=16384 Len=0 MSS=1240 WS=1
```

Frame 195: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0  
Ethernet II, Src: fa:16:3e:8f:8f:54 (fa:16:3e:8f:8f:54), Dst: fa:16:3e:d7:7e:f6  
(fa:16:3e:d7:7e:f6)

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

Transmission Control Protocol, Src Port: 179, Dst Port: 37740, Seq: 0, Ack: 1, Len: 0

Source Port: 179

Destination Port: 37740

[Stream index: 7]

[TCP Segment Len: 0]

Sequence number: 0 (relative sequence number)

Acknowledgment number: 1 (relative ack number)

Header Length: 28 bytes

Flags: 0x012 (SYN, ACK)

Window size value: 16384

[Calculated window size: 16384]

Checksum: 0xd8f7 [unverified]

```
[Checksum Status: Unverified]
Urgent pointer: 0
Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)
  Maximum segment size: 1240 bytes
    Kind: Maximum Segment Size (2)
    Length: 4
    MSS Value: 1240
  Window scale: 0 (multiply by 1)
  End of Option List (EOL)
```

## R4上看到的TCP會話詳細資訊 — ACTIVE:

! - as seen on R4 - Active

```
RP/0/0/CPU0:R4#show tcp detail pcb 0x12154d3c
Fri Jan  8 12:32:41.096 UTC
```

```
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Fri Jan  8 12:17:46 2021
```

```
PCB 0x12154d3c, SO 0x12154460, TCPCB 0x1215486c, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 1577
Local host: 192.168.0.4, Local port: 37740 (Local App PID: 1052958)
Foreign host: 192.168.0.1, Foreign port: 179
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768)  mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	19	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	16	15	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 2075436506  snduna: 2075436868  sndnxt: 2075436868
sndmax: 2075436868  sndwnd: 32460  sndcwnd: 3720
irs: 4238127261  rcvnx: 4238127623  rcvwnd: 32479  rcvadv: 4238160102
```

```
SRTT: 65 ms,  RTTO: 300 ms,  RTV: 40 ms,  KRTT: 0 ms
minRTT: 9 ms,  maxRTT: 229 ms
```

```
ACK hold time: 200 ms,  Keepalive time: 0 sec,  SYN waittime: 30 sec
Giveup time: 0 ms,  Retransmission retries: 0,  Retransmit forever: FALSE
Connect retries remaining: 30,  connect retry interval: 30 secs
```

```
State flags: none
Feature flags: Win Scale, Nagle
Request flags: Win Scale
```

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**

```
Window scales: rcv 0, snd 0, request rcv 0, request snd 0
Timestamp option: recent 0, recent age 0, last ACK sent 0
Sack blocks {start, end}: none
Sack holes {start, end, dups, rxmit}: none
```

```
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
```

Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:

#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R4#

## R1上看到的TCP會話詳細資訊 — PASSIVE:

! - as seen on R1 - Passive

RP/0/0/CPU0:R1#show tcp detail pcb 0x12155390

Fri Jan 8 12:23:52.041 UTC

=====

Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Fri Jan 8 12:17:43 2021

PCB 0x12155390, SO 0x121573e4, TCPCB 0x12156948, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 1577  
Local host: 192.168.0.1, Local port: 179 (Local App PID: 983326)  
Foreign host: 192.168.0.4, Foreign port: 37740

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	9	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	9	1	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 4238127261 snduna: 4238127471 sndnxt: 4238127471  
sndmax: 4238127471 sndwnd: 32631 sndcwnd: 3720  
irs: 2075436506 rcvnxt: 2075436716 rcvwnd: 32612 rcvadv: 2075469328

SRTT: 144 ms, RTTO: 578 ms, RTV: 434 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

```

Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

```

```

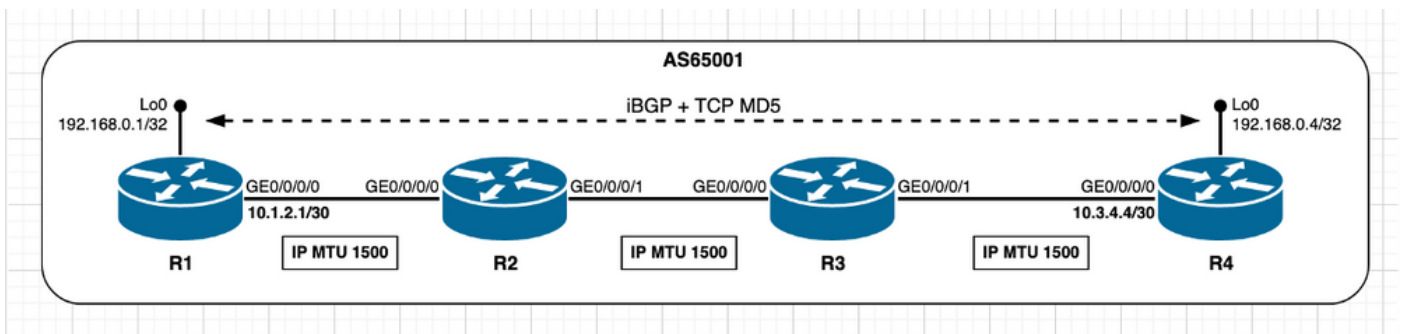
PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:
Num Labels: 0 Label Stack:

```

```
RP/0/0/CPU0:R1#
```

## 未直接連線的TCP對等體 — 使用TCP選項(MD5)

對於非直接連線對等體場景和使用TCP MD5身份驗證時，與前面所述的測試案例或場景沒有根本的差異。如之前使用TCP MD5驗證時所示，Cisco IOS XR會考慮其他額外負荷，且初始MSS值反映相同。請參閱前面的使用TCP選項 — XR主動和使用TCP選項 — XR被動部分以瞭解其他有關TCP選項對TCP MSS計算影響的詳細資訊。



映像2.7 — 未直接連線的TCP對等體 — iBGP + TCP MD5。

此案例中的TCP MSS計算可總結如下：

- 所有節點均使用預設IP MTU 1500位元組
- 預設情況下禁用TCP路徑MTU發現
- TCP對等點未直接連線 R4管理BGP連線目標R1未直接連線R4傳送MSS為1216位元組的SYN 當對等點未直接連線且TCP路徑MTU探索已停用時，不會考慮介面MTU根據設計，1280位元組被視為TCP\_DEFAULT\_MTU1280(TCP\_DEFAULT\_MTU)- 20(minTCP\_H)- 20(minIP\_H)- 24位元組 (IOS XR TCP選項額外負荷) R1傳送SYN, ACK, MSS為1216位元組 傳送[已接收MSS;本地初始MSS]接收的MSS 1216位元組；本地初始MSS 1240位元組兩個對等點上均使用最低MSS值

來源為R4的TCP SYN:

```
! - TCP SYN sourced from R4
```

```
3425 3.691042 192.168.0.4 192.168.0.1 TCP 82 42135 179 [SYN] Seq=0 Win=16384
Len=0 MSS=1216 WS=1
```

```

Frame 3425: 82 bytes on wire (656 bits), 82 bytes captured (656 bits) on interface 0
Ethernet II, Src: fa:16:3e:d7:7e:f6 (fa:16:3e:d7:7e:f6), Dst: fa:16:3e:8f:8f:54
(fa:16:3e:8f:8f:54)
Internet Protocol Version 4, Src: 192.168.0.4, Dst: 192.168.0.1
Transmission Control Protocol, Src Port: 42135, Dst Port: 179, Seq: 0, Len: 0
Source Port: 42135

```

```
Destination Port: 179
[Stream index: 10]
[TCP Segment Len: 0]
Sequence number: 0 (relative sequence number)
Acknowledgment number: 0
Header Length: 48 bytes
Flags: 0x002 (SYN)
Window size value: 16384
[Calculated window size: 16384]
Checksum: 0xc503 [unverified]
[Checksum Status: Unverified]
Urgent pointer: 0
Options: (28 bytes), Maximum segment size, Window scale, No-Operation (NOP), TCP MD5
signature, End of Option List (EOL)
    Maximum segment size: 1216 bytes
        Kind: Maximum Segment Size (2)
        Length: 4
        MSS Value: 1216
    Window scale: 0 (multiply by 1)
    No-Operation (NOP)
    TCP MD5 signature
    End of Option List (EOL)
```

### 源自R1的TCP SYN、ACK:

! - TCP SYN,ACK sourced from R1

```
3426 0.004186 192.168.0.1 192.168.0.4 TCP 82 179 42135 [SYN, ACK] Seq=0 Ack=1
Win=16384 Len=0 MSS=1216 WS=1
```

```
Frame 3426: 82 bytes on wire (656 bits), 82 bytes captured (656 bits) on interface 0
Ethernet II, Src: fa:16:3e:8f:8f:54 (fa:16:3e:8f:8f:54), Dst: fa:16:3e:d7:7e:f6
(fa:16:3e:d7:7e:f6)
Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4
Transmission Control Protocol, Src Port: 179, Dst Port: 42135, Seq: 0, Ack: 1, Len: 0
    Source Port: 179
    Destination Port: 42135
    [Stream index: 10]
    [TCP Segment Len: 0]
    Sequence number: 0 (relative sequence number)
    Acknowledgment number: 1 (relative ack number)
    Header Length: 48 bytes
    Flags: 0x012 (SYN, ACK)
    Window size value: 16384
    [Calculated window size: 16384]
    Checksum: 0xbb05 [unverified]
    [Checksum Status: Unverified]
    Urgent pointer: 0
    Options: (28 bytes), Maximum segment size, Window scale, No-Operation (NOP), TCP MD5
signature, End of Option List (EOL)
        Maximum segment size: 1216 bytes
            Kind: Maximum Segment Size (2)
            Length: 4
            MSS Value: 1216
        Window scale: 0 (multiply by 1)
        No-Operation (NOP)
        TCP MD5 signature
        End of Option List (EOL)
```

### R4上看到的TCP會話詳細資訊 — ACTIVE:

! - as seen from R4 - Active

RP/0/0/CPU0:R4#show tcp detail pcb 0x12154490

Tue Jan 12 14:37:32.097 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Tue Jan 12 14:27:42 2021

PCB 0x12154490, SO 0x12155014, TCPCB 0x12155a84, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 1876  
Local host: 192.168.0.4, Local port: 42135 (Local App PID: 1052958)  
Foreign host: 192.168.0.1, Foreign port: 179

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	14	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	11	9	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 3124761989 snduna: 3124763317 sndnxt: 3124763317  
sndmax: 3124763317 sndwnd: 32711 sndcwnd: 3648  
irs: 1090344992 rcvnxt: 1090346320 rcvwnd: 32730 rcvadv: 1090379050

SRTT: 28 ms, RTTO: 300 ms, RTV: 57 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 229 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 30, connect retry interval: 30 secs

State flags: none  
Feature flags: MD5, Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1216, peer MSS 1216, min MSS 1216, max MSS 1216**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R4#

R1上看到的TCP會話詳細資訊 — PASSIVE:

! - as seen from R1 - Passive

RP/0/0/CPU0:R1#show tcp detail pcb 0x12168df4

Tue Jan 12 14:36:38.860 UTC

=====

Connection state is ESTAB, I/O status: 0, socket status: 0

Established at Tue Jan 12 14:27:32 2021

PCB 0x12168df4, SO 0x12156bf8, TCPCB 0x12157a44, vrfid 0x60000000,

Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 1876

Local host: 192.168.0.1, Local port: 179 (Local App PID: 983326)

Foreign host: 192.168.0.4, Foreign port: 42135

Current send queue size in bytes: 0 (max 24576)

Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes

Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	12	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	12	1	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 1090344992 snduna: 1090346320 sndnxt: 1090346320  
sndmax: 1090346320 sndwnd: 32730 sndcwnd: 3648  
irs: 3124761989 rcvnxt: 3124763317 rcvwnd: 32711 rcvadv: 3124796028

SRTT: 150 ms, RTTO: 558 ms, RTV: 408 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: MD5, Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1216, peer MSS 1216, min MSS 1240, max MSS 1240**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

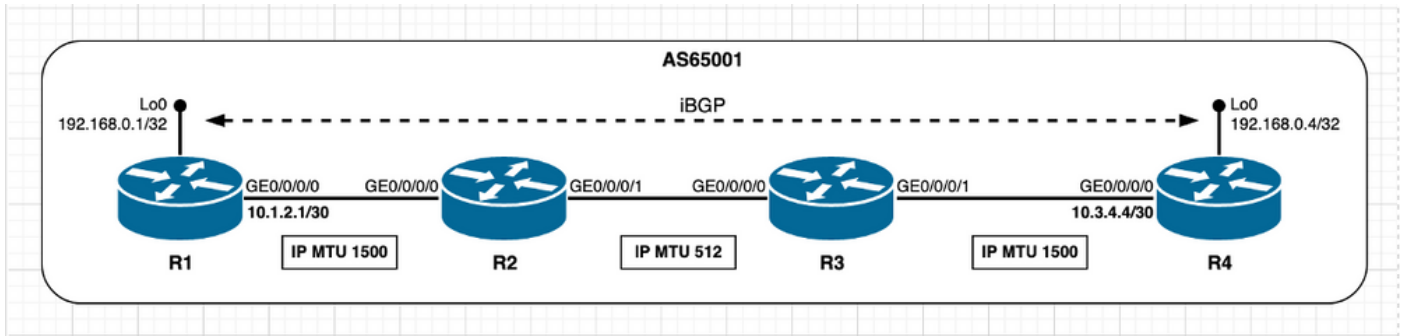
Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R1#

## 未直接連線的TCP對等點 — 路徑段具有較低的IP MTU

在下一個場景中，目標是觀察並總結在預設條件下如果存在具有較低IP MTU的中間路徑區段會發生的情況，這表示已停用TCP PMTUD。請參閱此映像。



映像2.8 - R2/R3路徑段的IP MTU較低。

最初的情況是認為BGP資訊是最小的，也就是說，BGP對等點之間需要交換的任何內容，都可以使用適合在512位元組的最小路徑MTU之下的IP封包完成。在此假設下，MSS計算會如未直接連線的TCP對等點一節所述。R1和R4都選擇1240位元組的MSS值。

### R4上看到的TCP會話詳細資訊 — ACTIVE:

! - as seen from R4 - Active

```
RP/0/0/CPU0:R4#show tcp detail pcb 0x15390fe8
```

```
=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Wed May 12 12:09:48 2021
```

```
PCB 0x15390fe8, SO 0x15391a7c, TCPCB 0x15391368, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 835  
Local host: 192.168.0.4, Local port: 39046 (Local App PID: 1196319)  
Foreign host: 192.168.0.1, Foreign port: 179  
(Local App PID/instance/SPL_APP_ID: 1196319/1/0)
```

```
Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	1267	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	1280	1235	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 1991226354 snduna: 1991250450 sndnxt: 1991250450  
sndmax: 1991250450 sndwnd: 32578 sndcwnd: 2480  
irs: 4276699304 rcvnxt: 4276746737 rcvwnd: 31568 rcvadp: 4276778305
```

```
SRTT: 213 ms, RTTO: 300 ms, RTV: 54 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 269 ms
```

```
ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec
```



Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 10, connect retry interval: 30 secs

State flags: none  
Feature flags: Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**  
<snip>

## R1上看到的TCP會話詳細資訊 — PASSIVE:

! - as seen from R1 - Passive

RP/0/0/CPU0:R1#show tcp detail pcb 0x15393770

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Wed May 12 12:09:46 2021

PCB 0x15393770, SO 0x15392224, TCPCB 0x153928cc, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 835  
Local host: 192.168.0.1, Local port: 179 (Local App PID: 1192224)  
Foreign host: 192.168.0.4, Foreign port: 39046  
(Local App PID/instance/SPL\_APP\_ID: 1192224/1/0)

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	1280	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	1264	1213	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 4276699304 snduna: 4276746718 sndnxt: 4276746718  
sndmax: 4276746718 sndwnd: 31587 sndcwnd: 3720  
irs: 1991226354 rcvnx: 1991250431 rcvwnd: 32597 rcvadv: 1991283028

SRTT: 202 ms, RTTO: 355 ms, RTV: 153 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 309 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**  
<snip>

現在建立BGP作業階段時，請考慮觸發大小高於最小路徑MTU 512位元組的BGP更新訊息。從輸出中可看出，Cisco IOS XR不會使用BGP更新消息設定df位元，這表示傳輸BGP資訊時，會犧牲中間節點上的封包分段。

源自R1的BGP更新 — 被動：

! - as seen from R1 - Passive - BGP UPDATE  
! - Note Total Length of 1097 bytes higher than the IP MTU value of 512 bytes at R2-R3 path segment

23 3.450878 192.168.0.1 192.168.0.4 BGP 1111 UPDATE Message

Frame 23: 1111 bytes on wire (8888 bits), 1111 bytes captured (8888 bits) on interface 0  
Ethernet II, Src: fa:16:3e:42:18:05 (fa:16:3e:42:18:05), Dst: fa:16:3e:5c:f1:80  
(fa:16:3e:5c:f1:80)

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

0100 .... = Version: 4

.... 0101 = Header Length: 20 bytes (5)

Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)

**Total Length: 1097**

Identification: 0x5841 (22593)

Flags: 0x00

0... .... = Reserved bit: Not set

.0.. .... = Don't fragment: Not set

..0. .... = More fragments: Not set

Fragment offset: 0

Time to live: 255

Protocol: TCP (6)

Header checksum: 0x54a4 [validation disabled]

[Header checksum status: Unverified]

Source: 192.168.0.1

Destination: 192.168.0.4

[Source GeoIP: Unknown]

[Destination GeoIP: Unknown]

Transmission Control Protocol, Src Port: 179, Dst Port: 39046, Seq: 20, Ack: 20, Len: 1057

Border Gateway Protocol - UPDATE Message

Marker: ff

Length: 1057

Type: UPDATE Message (2)

Withdrawn Routes Length: 0

Total Path Attribute Length: 1034

Path attributes

Path Attribute - MP\_REACH\_NLRI

Path Attribute - ORIGIN: INCOMPLETE

Path Attribute - AS\_PATH: empty

Path Attribute - MULTI\_EXIT\_DISC: 0

Path Attribute - LOCAL\_PREF: 100

源自R1的BGP更新消息的分段發生在節點R2,R2介面GE0/0/0/1上的流量捕獲可以觀察到這一點。

節點R2上的IP分段：

! - as seen from R2 - GE0/0/0/1

! - Node R2 fragments original packet in three distinct packets

4 1.334852 192.168.0.1 192.168.0.4 BGP 522 UPDATE Message

5 0.000289 192.168.0.1 192.168.0.4 IPv4 522 Fragmented IP protocol (proto=TCP 6, off=488, ID=7b41)

6 0.000122 192.168.0.1 192.168.0.4 IPv4 135 Fragmented IP protocol (proto=TCP 6, off=976, ID=7b41)

! - Captured frame details

Frame 4: 522 bytes on wire (4176 bits), 522 bytes captured (4176 bits) on interface 0  
Ethernet II, Src: fa:16:3e:61:25:f0 (fa:16:3e:61:25:f0), Dst: fa:16:3e:23:ab:27  
(fa:16:3e:23:ab:27)

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)  
**Total Length: 508**  
**Identification: 0x7b41 (31553)**  
Flags: 0x01 (More Fragments)  
  0... .... = Reserved bit: Not set  
  .0.. .... = Don't fragment: Not set  
  ..1. .... = **More fragments: Set**

**Fragment offset: 0**  
Time to live: 254  
Protocol: TCP (6)  
Header checksum: 0x14f1 [validation disabled]  
[Header checksum status: Unverified]  
Source: 192.168.0.1  
Destination: 192.168.0.4  
[Source GeoIP: Unknown]  
[Destination GeoIP: Unknown]

Transmission Control Protocol, Src Port: 179, Dst Port: 39046, Seq: 4276759681, Ack: 1991250830  
Border Gateway Protocol - UPDATE Message  
<snip>

Frame 5: 522 bytes on wire (4176 bits), 522 bytes captured (4176 bits) on interface 0  
Ethernet II, Src: fa:16:3e:61:25:f0 (fa:16:3e:61:25:f0), Dst: fa:16:3e:23:ab:27  
(fa:16:3e:23:ab:27)

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)  
**Total Length: 508**  
**Identification: 0x7b41 (31553)**  
Flags: 0x01 (More Fragments)  
  0... .... = Reserved bit: Not set  
  .0.. .... = Don't fragment: Not set  
  ..1. .... = **More fragments: Set**

**Fragment offset: 488**  
Time to live: 254  
Protocol: TCP (6)  
Header checksum: 0x14b4 [validation disabled]  
[Header checksum status: Unverified]  
Source: 192.168.0.1  
Destination: 192.168.0.4  
[Source GeoIP: Unknown]  
[Destination GeoIP: Unknown]

Data (488 bytes)

<snip>

Frame 6: 135 bytes on wire (1080 bits), 135 bytes captured (1080 bits) on interface 0  
Ethernet II, Src: fa:16:3e:61:25:f0 (fa:16:3e:61:25:f0), Dst: fa:16:3e:23:ab:27  
(fa:16:3e:23:ab:27)

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)  
**Total Length: 121**  
**Identification: 0x7b41 (31553)**  
Flags: 0x00  
  0... .... = Reserved bit: Not set  
  .0.. .... = Don't fragment: Not set  
  ..0. .... = **More fragments: Not set**

**Fragment offset: 976**  
Time to live: 254  
Protocol: TCP (6)  
Header checksum: 0x35fa [validation disabled]

```

[Header checksum status: Unverified]
Source: 192.168.0.1
Destination: 192.168.0.4
[Source GeoIP: Unknown]
[Destination GeoIP: Unknown]
Data (101 bytes)
<snip>

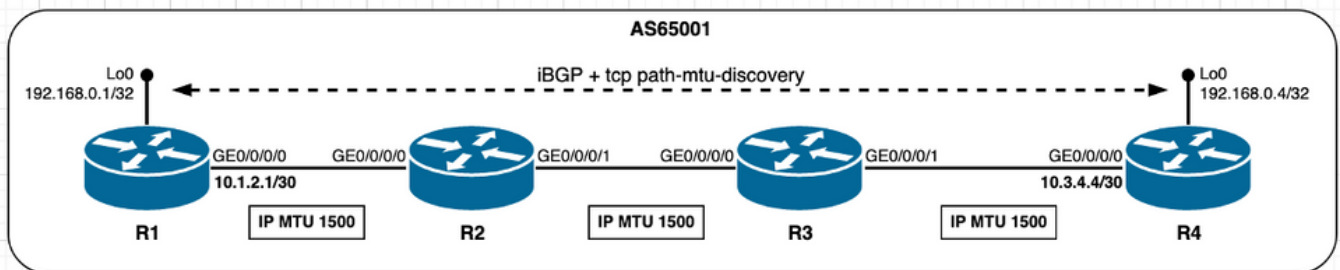
```

## 案例 — 啟用TCP PMTUD

### 啟用PMTUD

啟用PMTUD後，無論對等點是直接連線還是非直接連線，MSS初始計算一律會考慮輸出介面IP MTU。

此案例可深入瞭解啟用PMTUD時的預期行為。此處，Cisco IOS XR節點R4扮演主動角色，管理TCP連線，並在目的地埠179上開啟與Cisco IOS XR節點R1的TCP會話。兩個節點在其介面上使用預設IP MTU值。



映像3.1 — 已啟用TCP PMTUD。

此方案中的MSS計算可概述如下：

- 所有節點均使用預設IP MTU 1500位元組
- 已啟用TCP路徑MTU探索
- TCP對等點未直接連線 R4管理BGP連線R4傳送MSS為1460位元組的SYN 1500 ( 介面IP MTU ) — 20(minTCP\_H)- 20(minIP\_H)R1傳送SYN，ACK，MSS為1460位元組 傳送[已接收MSS;本地初始MSS]接收的MSS 1460位元組；本地初始MSS 1460位元組兩個對等點上均使用最低MSS值

為了突出啟用PMTUD所產生的行為變更，接下來的輸出說明事件的順序：

1. 在預設停用PMTUD的情況下已建立TCP作業階段的初始狀態；
2. 在TCP對等路由器R4和R1上設定和啟用PMTUD;
3. TCP作業階段會重新啟動，MSS計算會發生，並受到TCP PMTUD的影響。

如R4上所示 — ACTIVE - TCP PMTUD已禁用 ( 預設 )：

```

! - as seen on R4 - Active
! - TCP path mtu discovery disabled (default)
! - TCP session initial state

```

```

RP/0/0/CPU0:R4#show tcp detail pcb 0x121536c8
Fri Jan 8 16:06:30.237 UTC

```

```

=====

```

Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Fri Jan 8 16:05:15 2021

PCB 0x121536c8, SO 0x12155370, TCPCB 0x12154f64, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 376  
Local host: 192.168.0.4, Local port: 20155 (Local App PID: 1052958)  
Foreign host: 192.168.0.1, Foreign port: 179

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	6	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	3	2	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 357400981 snduna: 357401257 sndnxt: 357401257  
sndmax: 357401257 sndwnd: 32546 sndcwnd: 3720  
irs: 524019443 rcvnxt: 524019719 rcvwnd: 32565 rcvadv: 524052284

SRTT: 72 ms, RTTO: 416 ms, RTV: 344 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 229 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 30, connect retry interval: 30 secs

State flags: none  
Feature flags: Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R4#  
**如R1上所示 — 被動 — TCP PMTUD已禁用 ( 預設 ) :**

! - as seen on R1 - Passive  
! - TCP path mtu discovery disabled (default)  
! - TCP session initial state

RP/0/0/CPU0:R1#show tcp detail pcb 0x12157020

Fri Jan 8 16:05:52.868 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Fri Jan 8 16:05:12 2021

PCB 0x12157020, SO 0x121565ac, TCPCB 0x121560ec, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 376  
Local host: 192.168.0.1, Local port: 179 (Local App PID: 983326)  
Foreign host: 192.168.0.4, Foreign port: 20155

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	3	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	3	1	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 524019443 snduna: 524019700 sndnxt: 524019700  
sndmax: 524019700 sndwnd: 32584 sndcwnd: 3720  
irs: 357400981 rcvnxt: 357401238 rcvwnd: 32565 rcvadv: 357433803

SRTT: 46 ms, RTTO: 300 ms, RTV: 249 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: Win Scale, Nagle  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R1#  
如R4上所示 — ACTIVE - TCP PMTUD已啟用 :

! - 'debug tcp pmtud' output on R4  
! - tcp path mtu discovery enabled and uses default Path MTU aging timer (10 min / 600000 msec)

RP/0/0/CPU0:Jan 8 16:09:28.285 : tcp[399]: [t21] Try to enable path MTU discovery(neww age timer: 10 min)

RP/0/0/CPU0:Jan 8 16:09:28.285 : tcp[399]: [t21] Path mtu is ON (age-timer: 10)

! - as seen on R4 - Active  
! - TCP PMTUD is enabled

RP/0/0/CPU0:R4#show tcp detail pcb 0x121536c8

Fri Jan 8 16:11:00.138 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Fri Jan 8 16:05:15 2021

PCB 0x121536c8, SO 0x12155370, TCPCB 0x12154f64, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 376  
Local host: 192.168.0.4, Local port: 20155 (Local App PID: 1052958)  
Foreign host: 192.168.0.1, Foreign port: 179

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	10	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	7	4	0
KeepAlive	1	0	0
<b>PmtuAger</b>	<b>1</b>	<b>0</b>	<b>508096</b>
GiveUp	0	0	0
Throttle	0	0	0

iss: 357400981 snduna: 357401333 sndnxt: 357401333  
sndmax: 357401333 sndwnd: 32470 sndcwnd: 3720  
irs: 524019443 rcvnxt: 524019795 rcvwnd: 32489 rcvadv: 524052284

SRTT: 116 ms, RTTO: 578 ms, RTV: 462 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 229 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 30, connect retry interval: 30 secs

State flags: PMTU ager  
Feature flags: Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768

Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:

#PDU's in buffer: 0

FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:

Num Labels: 0 Label Stack:

RP/0/0/CPU0:R4#

如R1上所示 — 被動 — 已啟用TCP PMTUD:

! - 'debug tcp pmtud' output on R1

! - tcp path mtu discovery is enabled and uses default Path MTU aging timer (10 min / 60000 msec)

RP/0/0/CPU0:Jan 8 16:09:25.214 : tcp[399]: [t21] Try to enable path MTU discovery(neww age timer: 10 min)

RP/0/0/CPU0:Jan 8 16:09:25.214 : tcp[399]: [t21] Path mtu is ON (age-timer: 10)

! - as seen on R1 - Passive

! - TCP PMTUD is enabled

RP/0/0/CPU0:R1#show tcp detail pcb 0x12157020

Fri Jan 8 16:10:03.101 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Fri Jan 8 16:05:12 2021

PCB 0x12157020, SO 0x121565ac, TCPCB 0x121560ec, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 376  
Local host: 192.168.0.1, Local port: 179 (Local App PID: 983326)  
Foreign host: 192.168.0.4, Foreign port: 20155

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	7	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	7	4	0
KeepAlive	1	0	0
<b>PmtuAger</b>	<b>1</b>	<b>0</b>	<b>562042</b>
GiveUp	0	0	0
Throttle	0	0	0

iss: 524019443 snduna: 524019776 sndnxt: 524019776  
sndmax: 524019776 sndwnd: 32508 sndcwnd: 3720  
irs: 357400981 rcvnx: 357401314 rcvwnd: 32489 rcvadv: 357433803

SRTT: 95 ms, RTTO: 528 ms, RTV: 433 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: PMTU ager  
Feature flags: Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1240, peer MSS 1240, min MSS 1240, max MSS 1240**



Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R1#

請注意PMTU老化器計時器行為：

! - Note PmtuAger timer initial value is 10min  
! - but after initial interval expires then it expires every 2min  
! - As seen from 'debug tcp pmtud' output  
! - TCP PMTUD is enabled

RP/0/0/CPU0:Jan 8 16:09:25.214 : tcp[399]: [t21] Try to enable path MTU discovery(neww age timer: 10 min)  
RP/0/0/CPU0:Jan 8 16:09:25.214 : tcp[399]: [t21] Path mtu is ON (age-timer: 10)  
RP/0/0/CPU0:Jan 8 16:19:25.233 : tcp[399]: [t21] PCB 0x12157020: Trying next higher MTU: 1240  
RP/0/0/CPU0:Jan 8 16:21:25.245 : tcp[399]: [t21] PCB 0x12157020: Trying next higher MTU: 1240  
RP/0/0/CPU0:Jan 8 16:23:25.256 : tcp[399]: [t21] PCB 0x12157020: Trying next higher MTU: 1240

如R4上所示 — ACTIVE - BGP Session restart - TCP SYN:

! - Once BGP session is cleared  
! - TCP SYN sourced from R4 - Active  
! - MSS calculation takes place and is influenced by TCP PMTUD

2734 4.810311 192.168.0.4 192.168.0.1 TCP 62 32077 179 [SYN] Seq=0 Win=16384  
Len=0 **MSS=1460** WS=1

Frame 2734: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0  
Ethernet II, Src: fa:16:3e:d7:7e:f6 (fa:16:3e:d7:7e:f6), Dst: fa:16:3e:8f:8f:54  
(fa:16:3e:8f:8f:54)  
Internet Protocol Version 4, Src: 192.168.0.4, Dst: 192.168.0.1  
Transmission Control Protocol, Src Port: 32077, Dst Port: 179, Seq: 0, Len: 0  
Source Port: 32077  
Destination Port: 179  
[Stream index: 25]  
[TCP Segment Len: 0]  
Sequence number: 0 (relative sequence number)  
Acknowledgment number: 0  
Header Length: 28 bytes  
Flags: 0x002 (SYN)  
Window size value: 16384  
[Calculated window size: 16384]  
Checksum: 0x6398 [unverified]  
[Checksum Status: Unverified]  
Urgent pointer: 0  
Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)

```
Maximum segment size: 1460 bytes
  Kind: Maximum Segment Size (2)
  Length: 4
  MSS Value: 1460
Window scale: 0 (multiply by 1)
End of Option List (EOL)
```

如R1上所示 — 被動 — BGP會話重新啟動 — TCP SYN , ACK。

```
! - Once BGP session is cleared
! - TCP SYN,ACK sourced from R1 - Passive
! - MSS calculation takes place and is influenced by TCP PMTUD
```

```
2735  0.003879      192.168.0.1 192.168.0.4 TCP    62      179  32077 [SYN, ACK] Seq=0 Ack=1
Win=16384 Len=0 MSS=1460 WS=1
```

```
Frame 2735: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0
Ethernet II, Src: fa:16:3e:8f:8f:54 (fa:16:3e:8f:8f:54), Dst: fa:16:3e:d7:7e:f6
(fa:16:3e:d7:7e:f6)
Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4
Transmission Control Protocol, Src Port: 179, Dst Port: 32077, Seq: 0, Ack: 1, Len: 0
  Source Port: 179
  Destination Port: 32077
  [Stream index: 25]
  [TCP Segment Len: 0]
  Sequence number: 0      (relative sequence number)
  Acknowledgment number: 1    (relative ack number)
  Header Length: 28 bytes
  Flags: 0x012 (SYN, ACK)
  Window size value: 16384
  [Calculated window size: 16384]
  Checksum: 0xbf77 [unverified]
  [Checksum Status: Unverified]
  Urgent pointer: 0
  Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)
    Maximum segment size: 1460 bytes
      Kind: Maximum Segment Size (2)
      Length: 4
      MSS Value: 1460
    Window scale: 0 (multiply by 1)
    End of Option List (EOL)
```

啟用TCP PMTUD並清除BGP作業階段後，在R4上看到的TCP作業階段詳細資訊 — ACTIVE:

```
! - BGP session re-established
! - as seen on R4 - Active
```

```
RP/0/0/CPU0:R4#show tcp detail pcb 0x121567f4
Fri Jan  8 16:45:13.928 UTC
```

```
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Fri Jan  8 16:41:49 2021
```

```
PCB 0x121567f4, SO 0x12154460, TCPCB 0x12156190, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 10
Local host: 192.168.0.4, Local port: 32077 (Local App PID: 1052958)
Foreign host: 192.168.0.1, Foreign port: 179
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768)  mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	8	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	5	3	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```

iss: 1254100669  snduna: 1254100983  sndnxt: 1254100983
sndmax: 1254100983  sndwnd: 32508      sndcwnd: 4380
irs: 839938559   rcvnxt: 839938873   rcvwnd: 32527   rcvadv: 839971400

```

```

SRTT: 79 ms,  RTTO: 485 ms,  RTV: 406 ms,  KRTT: 0 ms
minRTT: 9 ms,  maxRTT: 229 ms

```

```

ACK hold time: 200 ms,  Keepalive time: 0 sec,  SYN waittime: 30 sec
Giveup time: 0 ms,  Retransmission retries: 0,  Retransmit forever: FALSE
Connect retries remaining: 30,  connect retry interval: 30 secs

```

```

State flags: none
Feature flags: Win Scale, Nagle, Path MTU
Request flags: Win Scale

```

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

```

Window scales: rcv 0, snd 0, request rcv 0, request snd 0
Timestamp option: recent 0, recent age 0, last ACK sent 0
Sack blocks {start, end}: none
Sack holes {start, end, dups, rxmit}: none

```

```

Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer   : Low/High watermark 2048/24576, Notify threshold 0

```

```

PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x40  PD ctx: size: 0  data:
Num Labels: 0  Label Stack:

```

RP/0/0/CPU0:R4#

啟用TCP PMTUD並清除BGP作業階段後，在R1上看到的TCP作業階段詳細資訊 — 被動。

```

! - BGP session re-established
! - as seen on R1 - Passive

```

RP/0/0/CPU0:R1#show tcp detail pcb 0x121558cc

Fri Jan 8 16:44:59.448 UTC

```

=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Fri Jan 8 16:41:46 2021

```

```

PCB 0x121558cc, SO 0x121556d4, TCPCB 0x121575bc, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 10
Local host: 192.168.0.1, Local port: 179 (Local App PID: 983326)
Foreign host: 192.168.0.4, Foreign port: 32077

```

Current send queue size in bytes: 0 (max 24576)

Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	6	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	6	3	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 839938559 snduna: 839938873 sndnxt: 839938873  
sndmax: 839938873 sndwnd: 32527 sndcwnd: 4380  
irs: 1254100669 rcvnxt: 1254100983 rcvwnd: 32508 rcvadp: 1254133491

SRTT: 76 ms, RTTO: 454 ms, RTV: 378 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 219 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

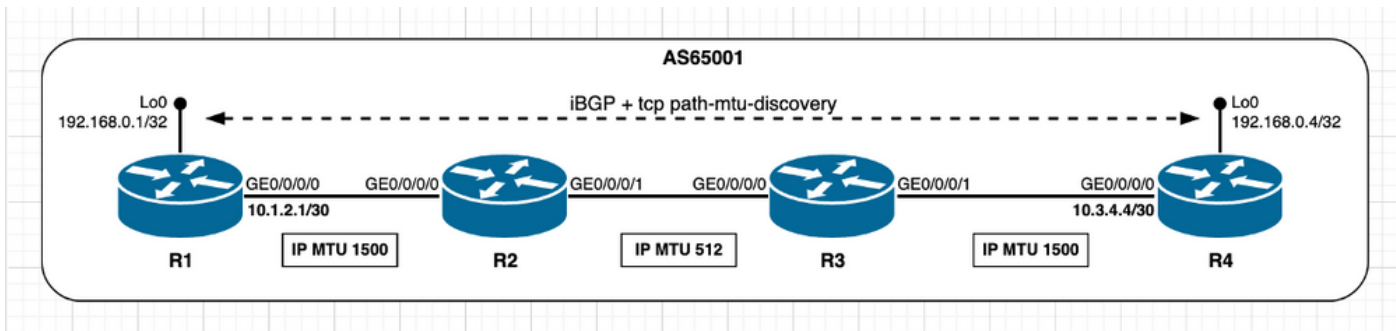
Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R1#

## PMTUD — 路徑區段的IP MTU更低

上一個案例有助於瞭解在啟用PMTUD的情況下建立初始TCP作業階段時發生的情況。此案例建立在上面，有助於瞭解TCP PMTUD的運作方式及其對已建立TCP作業階段的影響。



映像3.2 — 已啟用PMTUD，但路徑區段的IP MTU較低。

將上一個映像視為參考，假設BGP會話已建立，並且R1會傳送大小大於512位元組的IP資料包所攜帶的BGP更新消息。啟用PMTUD後，現在已設定DF位元（不分段）。因此，節點R2丟棄IP資料包並傳送ICMP（網際網路控制消息協定）消息（無法到達目的地 — 型別3；需要分段 — 代碼4）傳回R1。在節點R1收到ICMP訊息後，會觸發PMTUD，並嘗試建立路徑最小的IP MTU。它使用一組定義良好的平台級別中的下一個較低值，即視為新的TCP會話MSS值。接著，TCP使用新的MSS值重新傳輸原始BGP更新，此程式會根據需要重複多次，直到ICMP訊息（無法到達目的地 — 型別3；不再接收需要分段 — 代碼4）。這表示直到使用的MSS值使得每個傳送的封包都位於最低路徑區段IP MTU之下。隨著時間的流逝，由PmtuAger計時器控制的PMTUD會反向經過平台級別，並將MSS提升回其最大值。在任何指定時間，如果ICMP消息（無法到達目的地 — 型別3；需要分段 — 再次收到代碼4），然後PMTUD會如前所述運作。

接下來的輸出會瀏覽剛才介紹的PMTUD行為，並從已建立TCP作業階段的案例開始。這裡，Cisco IOS XR節點R4扮演主動角色，因此管理TCP連線，並在目的地埠179上開啟與R1的TCP會話。兩個節點在其介面上使用預設IP MTU值。此方案中的初始MSS計算可概述如下：

- R2和R3節點之間的中間網段使用非預設IP MTU 512位元組。
- R1和R4在其介面上使用預設MTU值。
- 已啟用TCP路徑MTU發現。
- TCP對等點沒有直接連線。R4管理BGP連線。R4傳送MSS為1460位元組的SYN。1500（介面IP MTU）— 20(minTCP\_H)- 20(minIP\_H)。R1傳送SYN和MSS為1460位元組的ACK。傳送[Received MSS ;本地初始MSS]。接收的MSS 1460位元組；本地初始MSS 1460位元組。兩個對等點上使用最小MSS值。

來源為R4的TCP SYN:

```
! - Initial TCP session establishment
! - TCP SYN sourced from R4

392      6.752774      192.168.0.4 192.168.0.1 TCP      62      32449 179 [SYN] Seq=0 Win=16384
Len=0 MSS=1460 WS=1

Frame 392: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0
Ethernet II, Src: fa:16:3e:5c:f1:80 (fa:16:3e:5c:f1:80), Dst: fa:16:3e:42:18:05
(fa:16:3e:42:18:05)
Internet Protocol Version 4, Src: 192.168.0.4, Dst: 192.168.0.1
Transmission Control Protocol, Src Port: 32449, Dst Port: 179, Seq: 0, Len: 0
  Source Port: 32449
  Destination Port: 179
  [Stream index: 10]
  [TCP Segment Len: 0]
  Sequence number: 0 (relative sequence number)
  Acknowledgment number: 0
  Header Length: 28 bytes
  Flags: 0x002 (SYN)
```

```

Window size value: 16384
[Calculated window size: 16384]
Checksum: 0x6858 [unverified]
[Checksum Status: Unverified]
Urgent pointer: 0
Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)
  Maximum segment size: 1460 bytes
    Kind: Maximum Segment Size (2)
    Length: 4
    MSS Value: 1460
  Window scale: 0 (multiply by 1)
  End of Option List (EOL)

```

## 源自R1的TCP SYN、ACK:

```

! - Initial TCP session establishment
! - TCP SYN,ACK sourced from R1

```

```

393      0.003628      192.168.0.1 192.168.0.4 TCP      62      179  32449 [SYN, ACK] Seq=0 Ack=1
Win=16384 Len=0 MSS=1460 WS=1

```

```

Frame 393: 62 bytes on wire (496 bits), 62 bytes captured (496 bits) on interface 0
Ethernet II, Src: fa:16:3e:42:18:05 (fa:16:3e:42:18:05), Dst: fa:16:3e:5c:f1:80
(fa:16:3e:5c:f1:80)
Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4
Transmission Control Protocol, Src Port: 179, Dst Port: 32449, Seq: 0, Ack: 1, Len: 0
  Source Port: 179
  Destination Port: 32449
  [Stream index: 10]
  [TCP Segment Len: 0]
  Sequence number: 0      (relative sequence number)
  Acknowledgment number: 1      (relative ack number)
  Header Length: 28 bytes
  Flags: 0x012 (SYN, ACK)
  Window size value: 16384
  [Calculated window size: 16384]
  Checksum: 0x509e [unverified]
  [Checksum Status: Unverified]
  Urgent pointer: 0
  Options: (8 bytes), Maximum segment size, Window scale, End of Option List (EOL)
    Maximum segment size: 1460 bytes
      Kind: Maximum Segment Size (2)
      Length: 4
      MSS Value: 1460
    Window scale: 0 (multiply by 1)
    End of Option List (EOL)

```

建立BGP作業階段後，節點R1會傳送BGP更新訊息並接收ICMP訊息(無法到達目的地 — 型別3 ;需要分段 — 代碼4)返回源自R2節點。

之所以會出現這種情況，是因為傳送BGP更新消息的IP封包已設定DF位元，且R2/R3區段上使用的512位元組的IP MTU低於1116位元組的IP封包大小。如前所述，接收ICMP訊息會觸發PMTUD。

在R1 ICMP處，收到型別3/代碼4消息：

```

! - as seen from R1 - Passive
! - After session is established R1 sends BGP Update message with IP length of 1116 Bytes
! - note IP Header Flags shows DF bit set

```

```

528      5.893055      192.168.0.1 192.168.0.4 BGP      1130    UPDATE Message, KEEPALIVE Message

```

Frame 528: 1130 bytes on wire (9040 bits), 1130 bytes captured (9040 bits) on interface 0  
Ethernet II, Src: fa:16:3e:42:18:05 (fa:16:3e:42:18:05), Dst: fa:16:3e:5c:f1:80  
(fa:16:3e:5c:f1:80)

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)

**Total Length: 1116**

Identification: 0x8c37 (35895)

**Flags: 0x02 (Don't Fragment)**

Fragment offset: 0

Time to live: 255

Protocol: TCP (6)

Header checksum: 0xe09a [validation disabled]

[Header checksum status: Unverified]

Source: 192.168.0.1

Destination: 192.168.0.4

[Source GeoIP: Unknown]

[Destination GeoIP: Unknown]

Transmission Control Protocol, Src Port: 179, Dst Port: 32449, Seq: 318, Ack: 251, Len: 1076

Border Gateway Protocol - UPDATE Message

Border Gateway Protocol - KEEPALIVE Message

<snip>

! - as seen from R1 - Passive

! - IP MTU on R2/R3 is lower than IP packet length and DF bit is set

! - R1 receives ICMP error message from R2

! - note R2 ICMP error message carries Next-Hop MTU

! - "The size in octets of the largest datagram that could be forwarded, along the path of

! the original datagram, without being fragmented at this router. The size includes the

! IP header and IP data, and does not include any lower-level headers."

529 0.002423 10.2.3.1 192.168.0.1 ICMP 110 **Destination unreachable  
(Fragmentation needed)**

Frame 529: 110 bytes on wire (880 bits), 110 bytes captured (880 bits) on interface 0  
Ethernet II, Src: fa:16:3e:5c:f1:80 (fa:16:3e:5c:f1:80), Dst: fa:16:3e:42:18:05  
(fa:16:3e:42:18:05)

Internet Protocol Version 4, Src: 10.2.3.1, Dst: 192.168.0.1

0100 .... = Version: 4

.... 0101 = Header Length: 20 bytes (5)

Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)

Total Length: 96

Identification: 0x0001 (1)

Flags: 0x00

Fragment offset: 0

Time to live: 255

**Protocol: ICMP (1)**

Header checksum: 0xac97 [validation disabled]

[Header checksum status: Unverified]

Source: 10.2.3.1

Destination: 192.168.0.1

[Source GeoIP: Unknown]

[Destination GeoIP: Unknown]

Internet Control Message Protocol

**Type: 3 (Destination unreachable)**

**Code: 4 (Fragmentation needed)**

Checksum: 0x2d52 [correct]

[Checksum Status: Good]

Length: 17

[Length of original datagram: 68]

Unused: 0011

**MTU of next hop: 512**

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

```
0100 .... = Version: 4
.... 0101 = Header Length: 20 bytes (5)
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)
Total Length: 1116
Identification: 0x8c37 (35895)
Flags: 0x02 (Don't Fragment)
Fragment offset: 0
Time to live: 254
Protocol: TCP (6)
Header checksum: 0xe19a [validation disabled]
[Header checksum status: Unverified]
Source: 192.168.0.1
Destination: 192.168.0.4
[Source GeoIP: Unknown]
[Destination GeoIP: Unknown]
```

```
Transmission Control Protocol, Src Port: 179, Dst Port: 32449, Seq: 2847698730, Ack:
2130367817
```

```
Border Gateway Protocol - UPDATE Message
```

```
[Packet size limited during capture: IPv4 truncated]
```

在節點R1 ( 由ICMP訊息觸發 ) , TCP PMTUD會嘗試使用定義良好的平台(IP MTU)層中的下一個較低值來建立端對端的最低IP MTU。這些平台層級記錄在[RFC1191 — 路徑MTU探索中](#)。

```
MTU plateaus from RFC 1191
```

```
- values include both TCP and IP headers
```

```
65535
32000
17914
8166
4352
2002
1492
1006
508
296
68
```

但由於ICMP(目的地無法連線 — 型別3;需要分段 — 代碼4)節點R1收到的消息傳遞下一躍點的MTU, 然後如圖所示, 節點R1使用此值 ( 在我們的示例中為512位元組 ) , 並調整TCP會話MSS值。請注意, 原始TCP區段的長度為1076位元組, 因此重新傳輸原始TCP區段需要三個封包。

如R1上所示 — 被動 — PMTUD操作 :

```
! - As seen from R1 - Passive
! - Hint is provided by ICMP unreachable message MTU of next-hop field: 512 bytes
! - R1 then considers this value and retransmits BGP Update split in three distinct packets
! - Sum of TCP length = 472 + 472 + 132 = 1076 bytes
```

```
530    0.007497      192.168.0.1 192.168.0.4 TCP    526    [TCP Out-Of-Order] 179  32449 [ACK]
Seq=318 Ack=251 Win=32593 Len=472
532    0.015374      192.168.0.1 192.168.0.4 TCP    526    [TCP Retransmission] 179  32449
[ACK] Seq=790 Ack=251 Win=32593 Len=472
533    0.004129      192.168.0.1 192.168.0.4 TCP    186    [TCP Retransmission] 179  32449
[PSH, ACK] Seq=1262 Ack=251 Win=32593 Len=132
```

如前所述, 在一段時間內, 傳送所有封包後, PMTUD會沿PmtuAger計時器所界定的相反方向走過平台層級, 並嘗試根據適當的情境將MSS提升為其最大值。

如R1上所示 — PMTUD跨定義的平台 :



! - As seen from R1 - Passive - 'debug tcp pmtud' and 'debug icmp' active  
! - TCP PMTUD is triggered once ICMP unreachable received

```
RP/0/0/CPU0:May 12 09:09:22.763 UTC: ipv4_io[266]: IPv4 ICMP: Received ICMP too big from
192.168.0.1 about 192.168.0.4, MTU=512
RP/0/0/CPU0:May 12 09:09:22.763 UTC: ipv4_io[266]: ipv4_icmp_unreachable_rcvd ICMP unreach
recvd: sending pak(0xb0c07d8f) to transport: 6, tid: 5
RP/0/0/CPU0:May 12 09:09:22.763 UTC: ipv4_io[266]: ip_icmp_lib_ipv4_receive: sending
pak(0xb0c07d8f) to transport: 1, tid: 5
RP/0/0/CPU0:May 12 09:09:22.763 UTC: tcp[399]: [t4] PCB 0x15393770: Process ICMP Dest-unreach
(next hop mtu: 512)
```

! - attempt new MSS 472 = MTU of next-hop(512) - TCP\_H(20) - IP\_H(20)

```
RP/0/0/CPU0:May 12 09:09:22.763 UTC: tcp[399]: [t4] PCB 0x15393770: Process ICMP Dest-unreach
(next hop mtu: 512)
RP/0/0/CPU0:May 12 09:09:22.763 UTC: tcp[399]: [t4] PCB 0x15393770: Try to use new MSS: 472
RP/0/0/CPU0:May 12 09:09:22.763 UTC: tcp[399]: [t4] PCB 0x15393770, New path MTU decided to use:
472 configured tp_user_mss 0
```

! - over time PMTUD attempts to raise MSS as per egress interface configured MTU

```
RP/0/0/CPU0:May 12 09:19:22.782 UTC: tcp[399]: [t23] PCB 0x15393770: Trying next higher MTU: 966
RP/0/0/CPU0:May 12 09:21:22.793 UTC: tcp[399]: [t23] PCB 0x15393770: Trying next higher MTU:
1452
RP/0/0/CPU0:May 12 09:23:22.805 UTC: tcp[399]: [t23] PCB 0x15393770: Trying next higher MTU:
1460
```

在這些輸出上可觀察到最終狀態。特別要注意節點R1顯示的最小MSS值和最大MSS值，它突出顯示並顯示PMTUD已觸發。

## R4上看到的TCP會話詳細資訊 — ACTIVE:

! - Final stage as seen from R4 - Active

```
RP/0/0/CPU0:R4#show tcp detail pcb 0x153913b8
Wed May 12 10:09:43.246 UTC
```

```
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Wed May 12 09:02:07 2021
```

```
PCB 0x153913b8, SO 0x153917f0, TCPCB 0x1538fb58, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 382
Local host: 192.168.0.4, Local port: 32449 (Local App PID: 1196319)
Foreign host: 192.168.0.1, Foreign port: 179
(Local App PID/instance/SPL_APP_ID: 1196319/1/0)
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	72	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	71	69	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 2130367566 snduna: 2130368957 sndnxt: 2130368957
```

sndmax: 2130368957 sndwnd: 31453 sndcwnd: 2920  
irs: 2847698412 rcvnxt: 2847700946 rcvwnd: 31799 rcvadv: 2847732745

SRTT: 220 ms, RTTO: 300 ms, RTV: 12 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 10, connect retry interval: 30 secs

State flags: none  
Feature flags: Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0  
Socket misc info : Rcv data size (sb\_cc) 0, so\_qlen 0,  
so\_q0len 0, so\_qlimit 0, so\_error 0  
so\_auto\_rearm 1

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:  
Num of peers with authentication info: 0

RP/0/0/CPU0:R4#

## R1上看到的TCP會話詳細資訊 — PASSIVE:

! - Final stage as seen from R1 - Passive

RP/0/0/CPU0:R1#show tcp detail pcb 0x15393770  
Wed May 12 10:12:41.432 UTC

=====  
Connection state is ESTAB, I/O status: 240, socket status: 0  
Established at Wed May 12 09:02:05 2021

PCB 0x15393770, SO 0x15394ea0, TCPCB 0x15391c0c, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 382  
Local host: 192.168.0.1, Local port: 179 (Local App PID: 1192224)  
Foreign host: 192.168.0.4, Foreign port: 32449  
(Local App PID/instance/SPL\_APP\_ID: 1192224/1/0)

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	75	0	0
SendWnd	0	0	0
TimeWait	0	0	0

```
AckHold          73          71          0
KeepAlive        1           0           0
PmtuAger       28         27         41595
GiveUp           0           0           0
Throttle         0           0           0
```

```
iss: 2847698412  snduna: 2847701003  sndnxt: 2847701003
sndmax: 2847701003  sndwnd: 31742          sndcwnd: 4380
irs: 2130367566  rcvnxt: 2130369014  rcvwnd: 31396  rcvadp: 2130400410
```

```
SRTT: 224 ms,  RTTO: 300 ms,  RTV: 23 ms,  KRTT: 0 ms
minRTT: 9 ms,  maxRTT: 259 ms
```

```
ACK hold time: 200 ms,  Keepalive time: 0 sec,  SYN waittime: 30 sec
Giveup time: 0 ms,  Retransmission retries: 0,  Retransmit forever: FALSE
Connect retries remaining: 0,  connect retry interval: 0 secs
```

```
State flags: PMTU ager
Feature flags: Win Scale, Nagle, Path MTU
Request flags: Win Scale
```

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 472, max MSS 1460**

```
Window scales: rcv 0, snd 0, request rcv 0, request snd 0
Timestamp option: recent 0, recent age 0, last ACK sent 0
Sack blocks {start, end}: none
Sack holes {start, end, dups, rxmit}: none
```

```
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer   : Low/High watermark 2048/24576, Notify threshold 0
Socket misc info     : Rcv data size (sb_cc) 0, so_qlen 0,
                      so_q0len 0, so_qlimit 0, so_error 0
                      so_auto_rearm 1
```

```
PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x20  PD ctx: size: 0  data:
Num Labels: 0  Label Stack:
Num of peers with authentication info: 0
```

```
RP/0/0/CPU0:R1#
```

最後，如果在任何指定時間存在ICMP(目的地無法連線 — 型別3 ;需要分段 — 代碼4)訊息再次收到，然後PMTUD再次按之前所述運作。

如R1所示 — PMTUD已再次觸發：

```
! - As seen from R1 - Passive
! - TCP PMTUD is again triggered upon new ICMP unreachable received
! - Behavior can be triggered via clearing redistributed, network and aggregate routes
originated
```

```
RP/0/0/CPU0:R1#clear bgp ipv4 all self-originated
Wed May 12 10:19:06.836 UTC
RP/0/0/CPU0:R1#
```

```
! - New BGP update message is sourced from R1 after clear bgp command
```

1707 1.712657 192.168.0.1 192.168.0.4 BGP 1121 UPDATE Message

Frame 1707: 1121 bytes on wire (8968 bits), 1121 bytes captured (8968 bits) on interface 0  
Ethernet II, Src: fa:16:3e:42:18:05 (fa:16:3e:42:18:05), Dst: fa:16:3e:5c:f1:80  
(fa:16:3e:5c:f1:80)

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)  
Total Length: 1107  
Identification: 0x1a38 (6712)  
Flags: 0x02 (Don't Fragment)  
Fragment offset: 0  
Time to live: 255  
Protocol: TCP (6)  
Header checksum: 0x52a3 [validation disabled]  
[Header checksum status: Unverified]  
Source: 192.168.0.1  
Destination: 192.168.0.4  
[Source GeoIP: Unknown]  
[Destination GeoIP: Unknown]

Transmission Control Protocol, Src Port: 179, Dst Port: 32449, Seq: 2705, Ack: 1562, Len: 1067  
Border Gateway Protocol - UPDATE Message

! - ICMP Destination Unreachable / Fragmentation needed is received and triggers PMTUD

1708 0.001614 10.2.3.1 192.168.0.1 ICMP 110 **Destination unreachable  
(Fragmentation needed)**

Frame 1708: 110 bytes on wire (880 bits), 110 bytes captured (880 bits) on interface 0  
Ethernet II, Src: fa:16:3e:5c:f1:80 (fa:16:3e:5c:f1:80), Dst: fa:16:3e:42:18:05  
(fa:16:3e:42:18:05)

Internet Protocol Version 4, Src: 10.2.3.1, Dst: 192.168.0.1

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)  
Total Length: 96  
Identification: 0x0002 (2)  
Flags: 0x00  
Fragment offset: 0  
Time to live: 255  
**Protocol: ICMP (1)**  
Header checksum: 0xac96 [validation disabled]  
[Header checksum status: Unverified]  
Source: 10.2.3.1  
Destination: 192.168.0.1  
[Source GeoIP: Unknown]  
[Destination GeoIP: Unknown]

Internet Control Message Protocol

**Type: 3 (Destination unreachable)**

**Code: 4 (Fragmentation needed)**

Checksum: 0x3b73 [correct]  
[Checksum Status: Good]  
Length: 17  
[Length of original datagram: 68]  
Unused: 0011

**MTU of next hop: 512**

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)  
Total Length: 1107  
Identification: 0x1a38 (6712)

Flags: 0x02 (Don't Fragment)  
Fragment offset: 0  
Time to live: 254  
Protocol: TCP (6)  
Header checksum: 0x53a3 [validation disabled]  
[Header checksum status: Unverified]  
Source: 192.168.0.1  
Destination: 192.168.0.4  
[Source GeoIP: Unknown]  
[Destination GeoIP: Unknown]

Transmission Control Protocol, Src Port: 179, Dst Port: 32449, Seq: 2847701117, Ack: 2130369128

Border Gateway Protocol - UPDATE Message

! - Note new/updated MSS value and PmtuAger  
! - MSS 472 ; Aligned with "MTU of next hop" value contained in ICMP message

RP/0/0/CPU0:R1#show tcp detail pcb 0x15393770

Wed May 12 10:19:31.494 UTC

=====

Connection state is ESTAB, I/O status: 240, socket status: 0

Established at Wed May 12 09:02:05 2021

PCB 0x15393770, SO 0x15394ea0, TCPCB 0x15391c0c, vrfid 0x60000000,

Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 382

Local host: 192.168.0.1, Local port: 179 (Local App PID: 1192224)

Foreign host: 192.168.0.4, Foreign port: 32449

(Local App PID/instance/SPL\_APP\_ID: 1192224/1/0)

Current send queue size in bytes: 0 (max 24576)

Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes

Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	83	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	80	77	0
KeepAlive	1	0	0
<b>PmtuAger</b>	<b>32</b>	<b>30</b>	<b>575401</b>
GiveUp	0	0	0
Throttle	0	0	0

iss: 2847698412 snduna: 2847702184 sndnxt: 2847702184

sndmax: 2847702184 sndwnd: 32173 sndcwnd: 944

irs: 2130367566 rcvnxt: 2130369147 rcvwnd: 32730 rcvadv: 2130401877

SRTT: 221 ms, RTTO: 300 ms, RTV: 16 ms, KRTT: 0 ms

minRTT: 9 ms, maxRTT: 259 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec

Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE

Connect retries remaining: 0, connect retry interval: 0 secs

State flags: PMTU ager

Feature flags: Win Scale, Nagle, **Path MTU**

Request flags: Win Scale

**Datagrams (in bytes): MSS 472, peer MSS 1460, min MSS 472, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0

Timestamp option: recent 0, recent age 0, last ACK sent 0

Sack blocks {start, end}: none

Sack holes {start, end, dups, rxmit}: none

```
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0
Socket misc info : Rcv data size (sb_cc) 0, so_qlen 0,
                  so_q0len 0, so_qlimit 0, so_error 0
                  so_auto_rearm 1
```

```
PDU information:
 #PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x20 PD ctx: size: 0 data:
 Num Labels: 0 Label Stack:
Num of peers with authentication info: 0
```

```
RP/0/0/CPU0:R1#
```

在受Cisco錯誤ID [CSCvf10395](#)影響的Cisco IOS XR版本上，會略過ICMP錯誤訊息中包含的下一個躍點，且節點會嘗試使用之前提到且由[RFC1191 - Path MTU discovery](#)所記錄的一組定義良好的平台(IP MTU)層級中的下一個較低值來建立端到端最低IP MTU。這些嘗試一直進行到成功傳輸，這表示一直到ICMP(目的地無法連線 — 型別3 ;需要分段 — 代碼4)不再接收消息。

從受Cisco錯誤ID [CSCvf10395](#)影響的Cisco IOS XR版本節點中可看出：

```
! - As seen from IOX XR node with a release impacted by Cisco bug ID CSCvf10395
! - Node ignores "MTU of next hop" and tries next lower plateau
! - This is observed till ICMP error messages are no longer received
! - Practical consequence is extra retransmissions occurrence
```

```
RP/0/0/CPU0:Feb 23 17:05:32.929 : tcp[399]: [t4] PCB 0x12152adc: Process ICMP Dest-unreach (next hop mtu: 33554432)
```

```
RP/0/0/CPU0:Feb 23 17:05:32.929 : tcp[399]: [t4] PCB 0x12152adc: Invalid next hop mtu (33554432), ignore it
```

```
RP/0/0/CPU0:Feb 23 17:05:34.649 : tcp[399]: [t27] PCB 0x12152adc: Trying next lower MTU: 1452
<<<<<<< HERE: Plateau 1492
```

```
RP/0/0/CPU0:Feb 23 17:05:35.519 : tcp[399]: [t4] PCB 0x12152adc: Process ICMP Dest-unreach (next hop mtu: 33554432)
```

```
RP/0/0/CPU0:Feb 23 17:05:35.519 : tcp[399]: [t4] PCB 0x12152adc: Invalid next hop mtu (33554432), ignore it
```

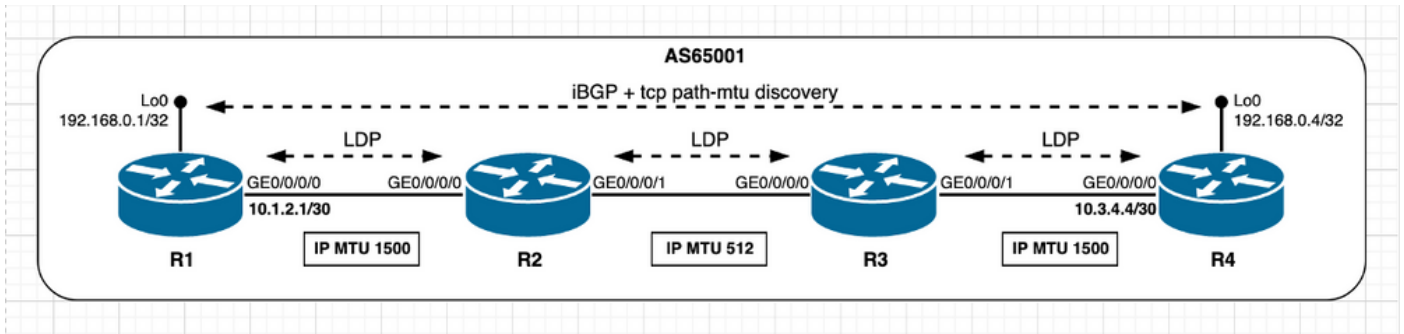
```
RP/0/0/CPU0:Feb 23 17:05:37.239 : tcp[399]: [t27] PCB 0x12152adc: Trying next lower MTU: 966
<<<<<<< HERE: Plateau 1006
```

```
RP/0/0/CPU0:Feb 23 17:05:38.109 : tcp[399]: [t4] PCB 0x12152adc: Process ICMP Dest-unreach (next hop mtu: 33554432)
```

```
RP/0/0/CPU0:Feb 23 17:05:38.109 : tcp[399]: [t4] PCB 0x12152adc: Invalid next hop mtu (33554432), ignore it
```

```
RP/0/0/CPU0:Feb 23 17:05:39.829 : tcp[399]: [t27] PCB 0x12152adc: Trying next lower MTU: 468
<<<<<<< HERE: Plateau 508
```

作為下一步，請考慮所有介面上使用標籤分發協定(LDP)的相同方案。此處的目標是瞭解在支援MPLS的環境中，從以前的場景中可觀察到哪些差異。



映像3.3 — 啟用PMTUD，且路徑段具有較低的IP MTU - MPLS案例。

首先，考慮在PMTUD觸發之前建立的BGP作業階段的初始階段，如此處所示。

在R4上看到的TCP(BGP)初始狀態 — 活動 — 啟用MPLS的方案：

```
! - as seen on R4 - Active
! - TCP path MTU discovery enabled
! - MPLS LDP enabled
! - TCP session initial state

RP/0/0/CPU0:R4#show tcp detail pcb 0x153bdaf0
Mon May 17 08:32:16.673 UTC
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Mon May 17 08:31:57 2021

PCB 0x153bdaf0, SO 0x153acc80, TCPCB 0x153acea8, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 757
Local host: 192.168.0.4, Local port: 57400 (Local App PID: 1196319)
Foreign host: 192.168.0.1, Foreign port: 179
(Local App PID/instance/SPL_APP_ID: 1196319/1/0)

Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768)  mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)

Timer           Starts      Wakeups      Next(msec)
Retrans         5           1            0
SendWnd         0           0            0
TimeWait        0           0            0
AckHold         2           1            0
KeepAlive       1           0            0
PmtuAger        0           0            0
GiveUp          0           0            0
Throttle        0           0            0

  iss: 1386459919  snduna: 1386460037  sndnxt: 1386460037
sndmax: 1386460037  sndwnd: 32726      sndcwnd: 4380
  irs: 3874414679  rcvnxt: 3874414864  rcvwnd: 32678   rcvadp: 3874447542

SRTT: 48 ms,  RTTO: 300 ms,  RTV: 228 ms,  KRTT: 0 ms
minRTT: 9 ms,  maxRTT: 229 ms

ACK hold time: 200 ms,  Keepalive time: 0 sec,  SYN waittime: 30 sec
Giveup time: 0 ms,  Retransmission retries: 0,  Retransmit forever: FALSE
Connect retries remaining: 10,  connect retry interval: 30 secs

State flags: none
Feature flags: Win Scale, Nagle, Path MTU
```

Request flags: Win Scale

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0

Timestamp option: recent 0, recent age 0, last ACK sent 0

Sack blocks {start, end}: none

Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO

Socket states: SS\_ISCONNECTED, SS\_PRIV

Socket receive buffer states: SB\_DEL\_WAKEUP

Socket send buffer states: SB\_DEL\_WAKEUP

Socket receive buffer: Low/High watermark 1/32768

Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

Socket misc info : Rcv data size (sb\_cc) 0, so\_qlen 0,  
so\_q0len 0, so\_qlimit 0, so\_error 0  
so\_auto\_rearm 1

PDU information:

#PDU's in buffer: 0

FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:

Num Labels: 1 Label Stack: 0x5dc2

Num of peers with authentication info: 0

RP/0/0/CPU0:R4#

**R1上看到的TCP(BGP)初始狀態 — 被動 — 啟用MPLS的情況 :**

! - as seen on R1 - Passive

! - TCP path MTU discovery enabled

! - MPLS LDP enabled

! - TCP session initial state

RP/0/0/CPU0:R1#show tcp detail pcb 0x153acc8c

Mon May 17 08:32:56.618 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0

Established at Mon May 17 08:31:55 2021

PCB 0x153acc8c, SO 0x153adad4, TCPCB 0x153adcfc, vrfid 0x60000000,

Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 757

Local host: 192.168.0.1, Local port: 179 (Local App PID: 1192224)

Foreign host: 192.168.0.4, Foreign port: 57400

(Local App PID/instance/SPL\_APP\_ID: 1192224/1/0)

Current send queue size in bytes: 0 (max 24576)

Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes

Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	3	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	3	1	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 3874414679 snduna: 3874414864 sndnxt: 3874414864

sndmax: 3874414864 sndwnd: 32678 sndcwnd: 4380

irs: 1386459919 rcvnxt: 1386460037 rcvwnd: 32726 rcvadv: 1386492763



SRTT: 45 ms, RTTO: 300 ms, RTV: 239 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 229 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1460, peer MSS 1460, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0  
Socket misc info : Rcv data size (sb\_cc) 0, so\_qlen 0,  
so\_q0len 0, so\_qlimit 0, so\_error 0  
so\_auto\_rearm 1

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x20 PD ctx: size: 0 data:  
Num Labels: 1 Label Stack: 0x5dc3  
Num of peers with authentication info: 0

RP/0/0/CPU0:R1#

在此啟用MPLS的場景中，觀察到已建立TCP(LDP)作業階段的詳細資訊。請注意，之前描述的所有有關TCP(BGP)會話的MSS計算的內容，同樣適用於TCP(LDP)會話。例如，節點R3和R2 TCP(LDP)會話MSS計算可以總結如下：

- R2和R3都使用512位元組的非預設IP MTU。
- 已啟用路徑MTU發現。
- TCP對等點沒有直接連線 ( TCP會話在環回介面之間建立 )。 R3管理LDP連線。R3傳送MSS為472位元組的SYN。 512 ( 介面IP MTU ) — 20(minTCP\_H)- 20(minIP\_H)。R2傳送SYN，ACK的MSS為472位元組。 傳送[已接收MSS;本地初始MSS]。已接收472位元組；本地初始MSS 472位元組。兩個對等點上使用最小MSS值。

TCP(LDP)會話詳細資訊，如在R3上所示 — 活動 — 啟用MPLS的方案：

```
! - as seen on R3 - Active
! - TCP path MTU discovery enabled
! - MPLS LDP enabled
! - TCP session initial state
```

RP/0/0/CPU0:R3#show tcp detail pcb 0x15393fbc  
Mon May 17 08:33:30.627 UTC

```
=====
Connection state is ESTAB, I/O status: 0, socket status: 0
Established at Mon May 17 08:30:04 2021
```

PCB 0x15393fbc, SO 0x15393d94, TCPCB 0x153941b4, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 970  
Local host: 192.168.0.3, Local port: 57146 (Local App PID: 1151216)  
Foreign host: 192.168.0.2, Foreign port: 646  
(Local App PID/instance/SPL\_APP\_ID: 1151216/0/0)

Current send queue size in bytes: 0 (max 16384)  
Current receive queue size in bytes: 0 (max 16384) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 60)

Timer	Starts	Wakeups	Next(msec)
Retrans	8	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	6	4	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 2917752466 snduna: 2917752838 sndnxt: 2917752838  
sndmax: 2917752838 sndwnd: 16013 sndcwnd: 944  
irs: 228184383 rcvnxt: 228184763 rcvwnd: 16005 rcvadv: 228200768

SRTT: 103 ms, RTTO: 580 ms, RTV: 477 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 279 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 1, connect retry interval: 3 secs

State flags: none  
Feature flags: Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 472, peer MSS 472, min MSS 472, max MSS 472**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_SEL, SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/16384  
Socket send buffer : Low/High watermark 2048/16384, Notify threshold 0  
Socket misc info : Rcv data size (sb\_cc) 0, so\_qlen 0,  
so\_q0len 0, so\_qlimit 0, so\_error 0  
so\_auto\_rearm 1

PDU information:  
#PDU's in buffer: 0  
FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 1 Label Stack: 0x5dc2  
Num of peers with authentication info: 0

RP/0/0/CPU0:R3#

R2上看到的TCP(LDP)會話詳細資訊 — 被動 — 啟用MPLS的方案 :

! - as seen on R2 - Passive

! - TCP path MTU discovery enabled  
! - MPLS LDP enabled  
! - TCP session initial state

RP/0/0/CPU0:R2#show tcp detail pcb 0x153a1f44

Mon May 17 08:34:28.843 UTC

=====

Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Mon May 17 08:30:31 2021

PCB 0x153a1f44, SO 0x153a1d1c, TCPCB 0x153a213c, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 970  
Local host: 192.168.0.2, Local port: 646 (Local App PID: 1151216)  
Foreign host: 192.168.0.3, Foreign port: 57146  
(Local App PID/instance/SPL\_APP\_ID: 1151216/0/0)

Current send queue size in bytes: 0 (max 16384)  
Current receive queue size in bytes: 0 (max 16384) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 60)

Timer	Starts	Wakeups	Next(msec)
Retrans	7	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	7	5	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 228184383 snduna: 228184763 sndnxt: 228184763  
sndmax: 228184763 sndwnd: 16005 sndcwnd: 944  
irs: 2917752466 rcvnxt: 2917752856 rcvwnd: 15995 rcvadv: 2917768851

SRTT: 95 ms, RTTO: 561 ms, RTV: 466 ms, KRTT: 0 ms  
minRTT: 0 ms, maxRTT: 219 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 472, peer MSS 472, min MSS 472, max MSS 472**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_SEL, SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/16384  
Socket send buffer : Low/High watermark 2048/16384, Notify threshold 0  
Socket misc info : Rcv data size (sb\_cc) 0, so\_qlen 0,  
so\_q0len 0, so\_qlimit 0, so\_error 0  
so\_auto\_rearm 1

PDU information:  
#PDU's in buffer: 0

```
FIB Lookup Cache:  IFH: 0x60  PD ctx: size: 0  data:
  Num Labels: 1  Label Stack: 0x5dc1
Num of peers with authentication info: 0
```

```
RP/0/0/CPU0:R2#
```

建立BGP作業階段後，R1會傳送BGP更新訊息並接收ICMP訊息(無法到達目的地 — 型別3;需要分段 — 代碼4)，返回源自R2的節點R2，該節點在節點R1處觸發TCP PMTUD。出現這種情況是因為傳送BGP更新消息的IP資料包設定了DF位元，並且R2/R3網段上使用的512位元組的IP MTU低於IP資料包大小1116位元組。與之前一樣，收到此ICMP訊息時會觸發PMTUD。與先前的非MPLS方案相比，啟用MPLS方案的差異在於節點R2 ICMP消息中包含的下一跳值(無法到達目的地 — 型別3;需要分段 — 代碼4)。在此啟用MPLS的情況下，下一躍點的MTU值佔用4位元組的額外MPLS額外負荷，這意味著它佔用R2處的輸出MPLS標籤堆疊，如以下輸出所示。

在R1上看到的TCP路徑MTU發現 — 被動 — 啟用MPLS的情況：

```
! - as seen from R1 - Passive
! - R1 sends BGP Update message with IP length of 1116 Bytes
! - Note MPLS Header as packet is to be label-switched (single label ; IGP label)
! - note IP Header Flags shows DF bit set

455      0.044859      192.168.0.1 192.168.0.4 BGP      1134      UPDATE Message, KEEPALIVE Message

Frame 455: 1134 bytes on wire (9072 bits), 1134 bytes captured (9072 bits) on interface 0
Ethernet II, Src: fa:16:3e:42:18:05 (fa:16:3e:42:18:05), Dst: fa:16:3e:5c:f1:80
(fa:16:3e:5c:f1:80)
MultiProtocol Label Switching Header, Label: 24002, Exp: 6, S: 1, TTL: 255
Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4
  0100 .... = Version: 4
  .... 0101 = Header Length: 20 bytes (5)
  Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)
  Total Length: 1116
  Identification: 0xc6dd (50909)
  Flags: 0x02 (Don't Fragment)
    0... .... = Reserved bit: Not set
    .1.. .... = Don't fragment: Set
    ..0. .... = More fragments: Not set
  Fragment offset: 0
  Time to live: 255
  Protocol: TCP (6)
  Header checksum: 0xa5f4 [validation disabled]
  [Header checksum status: Unverified]
  Source: 192.168.0.1
  Destination: 192.168.0.4
  [Source GeoIP: Unknown]
  [Destination GeoIP: Unknown]
Transmission Control Protocol, Src Port: 179, Dst Port: 57400, Seq: 242, Ack: 175, Len: 1076
Border Gateway Protocol - UPDATE Message
Border Gateway Protocol - KEEPALIVE Message
<snip>

! - as seen from R1 - Passive
! - IP MTU on R2/R3 of 512 bytes is lower than IP packet length and DF bit is set
! - R1 receives ICMP error message from R2
! - note R2 ICMP error message carries Next-Hop MTU
! - "The size in octets of the largest datagram that could be forwarded, along the path of
!   the original datagram, without being fragmented at this router. The size includes the
!   IP header and IP data, and does not include any lower-level headers."
! - In present MPLS-enabled scenario Next-Hop MTU value is 508 bytes
! - In previous non-MPLS scenario Next-Hop MTU value was 512 bytes
```

456 0.014117 10.2.3.1 192.168.0.1 ICMP 182 **Destination unreachable**  
**(Fragmentation needed)**

Frame 456: 182 bytes on wire (1456 bits), 182 bytes captured (1456 bits) on interface 0  
Ethernet II, Src: fa:16:3e:5c:f1:80 (fa:16:3e:5c:f1:80), Dst: fa:16:3e:42:18:05  
(fa:16:3e:42:18:05)

Internet Protocol Version 4, Src: 10.2.3.1, Dst: 192.168.0.1

0100 .... = Version: 4  
.... 0101 = Header Length: 20 bytes (5)  
Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)  
Total Length: 168  
Identification: 0x001f (31)  
Flags: 0x00  
0... .... = Reserved bit: Not set  
.0.. .... = Don't fragment: Not set  
..0. .... = More fragments: Not se

Fragment offset: 0

Time to live: 251

**Protocol: ICMP (1)**

Header checksum: 0xb031 [validation disabled]

[Header checksum status: Unverified]

Source: 10.2.3.1

Destination: 192.168.0.1

[Source GeoIP: Unknown]

[Destination GeoIP: Unknown]

Internet Control Message Protocol

**Type: 3 (Destination unreachable)**

**Code: 4 (Fragmentation needed)**

Checksum: 0x5199 [correct]

[Checksum Status: Good]

Length: 17

[Length of original datagram: 68]

Unused: 0011

**MTU of next hop: 508**

Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4

Transmission Control Protocol, Src Port: 179, Dst Port: 57400, Seq: 3874414921, Ack:

1386460094

Border Gateway Protocol - UPDATE Message

! - As seen from R1 - Passive

! - Hint is provided by ICMP unreachable message MTU of next-hop field: 508 bytes

! - R1 then considers this value and retransmits BGP Update split in three distinct packets

! - Sum of TCP length = 468 + 468 + 140 = 1076 bytes

457 0.006689 192.168.0.1 192.168.0.4 TCP 526 [TCP Retransmission] 179 57400

[ACK] Seq=242 Ack=175 Win=32669 **Len=468**

460 0.004001 192.168.0.1 192.168.0.4 TCP 526 [TCP Retransmission] 179 57400

[ACK] Seq=710 Ack=175 Win=32669 **Len=468**

461 0.001788 192.168.0.1 192.168.0.4 TCP 198 [TCP Retransmission] 179 57400

[PSH, ACK] Seq=1178 Ack=175 Win=32669 **Len=140**

463 0.056695 192.168.0.4 192.168.0.1 TCP 54 57400 179 [ACK] Seq=175 Ack=1318

Win=31545 Len=0

! - As seen from R1 - Passive - 'debug tcp pmtud' and 'debug icmp' active

! - TCP PMTUD is triggered once ICMP unreachable received

RP/0/0/CPU0:May 17 08:29:56.131 UTC: tcp[399]: [t1] Try to enable path MTU discovery(neww age  
timer: 10 min)

RP/0/0/CPU0:May 17 08:29:56.131 UTC: tcp[399]: [t1] Path mtu is ON (age-timer: 10)

RP/0/0/CPU0:May 17 08:35:51.726 UTC: ipv4\_io[266]: ip\_icmp\_lib\_ipv4\_receive: Receiving  
pak(0xb0c07d8f) tid: 5

RP/0/0/CPU0:May 17 08:35:51.726 UTC: ipv4\_io[266]: Entering ipv4\_mtu\_update\_cb

RP/0/0/CPU0:May 17 08:35:51.726 UTC: ipv4\_io[266]: IPv4 ICMP: Received ICMP too big from  
192.168.0.1 about 192.168.0.4, MTU=508

```
RP/0/0/CPU0:May 17 08:35:51.726 UTC: ipv4_io[266]: ipv4_icmp_unreachable_rcvd ICMP unreach
recvd: sending pak(0xb0c07d8f) to transport: 6, tid: 5
RP/0/0/CPU0:May 17 08:35:51.726 UTC: ipv4_io[266]: ip_icmp_lib_ipv4_receive: sending
pak(0xb0c07d8f) to transport: 1, tid: 5
RP/0/0/CPU0:May 17 08:35:51.726 UTC: tcp[399]: [t4] PCB 0x153acc8c: Process ICMP Dest-unreach
(next hop mtu: 508)
```

```
! - attempt new MSS 468 = MTU of next-hop(508) - TCP_H(20) - IP_H(20)
```

```
RP/0/0/CPU0:May 17 08:35:51.726 UTC: tcp[399]: [t4] PCB 0x153acc8c: Try to use new MSS: 468
RP/0/0/CPU0:May 17 08:35:51.726 UTC: tcp[399]: [t4] PCB 0x153acc8c, New path MTU decided to use:
468 configured tp_user_mss 0
```

```
! - over time PMTUD attempts to raise MSS as per egress interface configured MTU
```

```
RP/0/0/CPU0:May 17 08:45:51.745 UTC: tcp[399]: [t29] PCB 0x153acc8c: Trying next higher MTU: 966
RP/0/0/CPU0:May 17 08:47:51.757 UTC: tcp[399]: [t29] PCB 0x153acc8c: Trying next higher MTU:
1452
RP/0/0/CPU0:May 17 08:49:51.769 UTC: tcp[399]: [t29] PCB 0x153acc8c: Trying next higher MTU:
1460
```

如R1所示 — 被動 — TCP PMTUD觸發 — 啟用MPLS的情況：

```
! - as seen on R1 - Passive
! - R1 session details after TCP PMTUD trigger
```

```
RP/0/0/CPU0:R1#show tcp detail pcb 0x153acc8c
Mon May 17 08:43:07.077 UTC
```

```
=====
Connection state is ESTAB, I/O status: 240, socket status: 0
Established at Mon May 17 08:31:55 2021
```

```
PCB 0x153acc8c, SO 0x153adad4, TCPCB 0x153adcfc, vrfid 0x60000000,
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 757
Local host: 192.168.0.1, Local port: 179 (Local App PID: 1192224)
Foreign host: 192.168.0.4, Foreign port: 57400
(Local App PID/instance/SPL_APP_ID: 1192224/1/0)
```

```
Current send queue size in bytes: 0 (max 24576)
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes
Current receive queue size in packets: 0 (max 0)
```

Timer	Starts	Wakeups	Next(msec)
Retrans	15	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	14	9	0
KeepAlive	1	0	0
<b>PmtuAger</b>	<b>1</b>	<b>0</b>	<b>164599</b>
GiveUp	0	0	0
Throttle	0	0	0

```
iss: 3874414679 snduna: 3874416130 sndnxt: 3874416130
sndmax: 3874416130 sndwnd: 31412 sndcwnd: 936
irs: 1386459919 rcvnxt: 1386460246 rcvwnd: 32517 rcvadv: 1386492763
```

```
SRTT: 180 ms, RTTO: 509 ms, RTV: 329 ms, KRTT: 0 ms
minRTT: 19 ms, maxRTT: 239 ms
```

```
ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE
Connect retries remaining: 0, connect retry interval: 0 secs
```

```
State flags: PMTU ager
Feature flags: Win Scale, Nagle, Path MTU
Request flags: Win Scale
```

```
Datagrams (in bytes): MSS 468, peer MSS 1460, min MSS 468, max MSS 1460
```

```
Window scales: rcv 0, snd 0, request rcv 0, request snd 0
Timestamp option: recent 0, recent age 0, last ACK sent 0
Sack blocks {start, end}: none
Sack holes {start, end, dups, rxmit}: none
```

```
Socket options: SO_REUSEADDR, SO_REUSEPORT, SO_NBIO
Socket states: SS_ISCONNECTED, SS_PRIV
Socket receive buffer states: SB_DEL_WAKEUP
Socket send buffer states: SB_DEL_WAKEUP
Socket receive buffer: Low/High watermark 1/32768
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0
Socket misc info : Rcv data size (sb_cc) 0, so_qlen 0,
                  so_q0len 0, so_qlimit 0, so_error 0
                  so_auto_rearm 1
```

```
PDU information:
#PDU's in buffer: 0
FIB Lookup Cache: IFH: 0x20 PD ctx: size: 0 data:
Num Labels: 1 Label Stack: 0x5dc3
Num of peers with authentication info: 0
```

```
RP/0/0/CPU0:R1#
```

請注意，在啟用MPLS的情況下，節點R2 ICMP消息中包含的下一跳MTU的值將計入輸出MPLS標籤堆疊。要進一步強化這一方面，請考慮下一個示例。如果在R2過濾的IP資料包與L3VPN服務相關聯，則表示乙太網幀現在帶有兩個標籤（IGP標籤和VPN標籤）。接下來躍點的MTU會反映所需的標籤堆疊大小。請參閱這些輸出。

如R1上所示 — PASSIVE - L3 VPN服務資料包：

```
! - as seen from R1 - Passive
! - L3 VPN service packet is sourced by node R1 and destined to node R4
! - Note presence of MPLS label stack - both IGP and VPN label are present
! - Note IP Total Length of 610 bytes higher than the IP MTU on R2/R3 segment
! - note IP Header Flags shows DF bit set
```

```
2024 0.302370 10.1.14.1 10.1.14.14 TELNET 632 Telnet Data ...
```

```
Frame 2024: 632 bytes on wire (5056 bits), 632 bytes captured (5056 bits) on interface 0
Ethernet II, Src: fa:16:3e:42:18:05 (fa:16:3e:42:18:05), Dst: fa:16:3e:5c:f1:80
(fa:16:3e:5c:f1:80)
```

```
MultiProtocol Label Switching Header, Label: 24002, Exp: 0, S: 0, TTL: 255
```

```
0000 0101 1101 1100 0010 .... = MPLS Label: 24002
.... 000. .... = MPLS Experimental Bits: 0
.... 0 .... = MPLS Bottom Of Label Stack: 0
.... 1111 1111 = MPLS TTL: 255
```

```
MultiProtocol Label Switching Header, Label: 24005, Exp: 0, S: 1, TTL: 255
```

```
0000 0101 1101 1100 0101 .... = MPLS Label: 24005
.... 000. .... = MPLS Experimental Bits: 0
.... 1 .... = MPLS Bottom Of Label Stack: 1
.... 1111 1111 = MPLS TTL: 255
```

```
Internet Protocol Version 4, Src: 10.1.14.1, Dst: 10.1.14.14
```

```
0100 .... = Version: 4
.... 0101 = Header Length: 20 bytes (5)
Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)
Total Length: 610
```

```
Identification: 0x7c9f (31903)
Flags: 0x02 (Don't Fragment)
  0... .... = Reserved bit: Not set
  .1.. .... = Don't fragment: Set
  ..0. .... = More fragments: Not set
Fragment offset: 0
Time to live: 255
Protocol: TCP (6)
Header checksum: 0xcce5 [validation disabled]
[Header checksum status: Unverified]
Source: 10.1.14.1
Destination: 10.1.14.14
[Source GeoIP: Unknown]
[Destination GeoIP: Unknown]
```

Transmission Control Protocol, Src Port: 22008, Dst Port: 23, Seq: 34755, Ack: 93250, Len: 570  
如在R1上所示 — PASSIVE - L3 VPN服務 — ICMP型別3/代碼4:

```
! - as seen from R1 - Passive
! - IP MTU on R2/R3 of 512 bytes is lower than IP packet length and DF bit is set
! - R1 receives ICMP error message from R2
! - note R2 ICMP error message carries Next-Hop MTU
! - "The size in octets of the largest datagram that could be forwarded, along the path of
!   the original datagram, without being fragmented at this router. The size includes the
!   IP header and IP data, and does not include any lower-level headers."
! - In present L3VPN MPLS-enabled scenario (dual-label) Next-Hop MTU value is 504 bytes
! - In previous MPLS scenario (single-label) Next-Hop MTU value was 508 bytes
```

```
2030  0.020299      10.2.3.1      10.1.14.1      ICMP  190      Destination unreachable
(Fragmentation needed)
```

```
Frame 2030: 190 bytes on wire (1520 bits), 190 bytes captured (1520 bits) on interface 0
Ethernet II, Src: fa:16:3e:5c:f1:80 (fa:16:3e:5c:f1:80), Dst: fa:16:3e:42:18:05
(fa:16:3e:42:18:05)
```

```
MultiProtocol Label Switching Header, Label: 24005, Exp: 0, S: 1, TTL: 251
  0000 0101 1101 1100 0101 .... .... .... = MPLS Label: 24005
  .... .... .... .... .... 000. .... .... = MPLS Experimental Bits: 0
  .... .... .... .... .... ...1 .... .... = MPLS Bottom Of Label Stack: 1
  .... .... .... .... .... .... 1111 1011 = MPLS TTL: 251
```

```
Internet Protocol Version 4, Src: 10.2.3.1, Dst: 10.1.14.1
0100 .... = Version: 4
.... 0101 = Header Length: 20 bytes (5)
Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)
Total Length: 172
Identification: 0x002b (43)
Flags: 0x00
  0... .... = Reserved bit: Not set
  .0.. .... = Don't fragment: Not set
  ..0. .... = More fragments: Not set
```

```
Fragment offset: 0
Time to live: 253
Protocol: ICMP (1)
Header checksum: 0x9821 [validation disabled]
[Header checksum status: Unverified]
Source: 10.2.3.1
Destination: 10.1.14.1
[Source GeoIP: Unknown]
[Destination GeoIP: Unknown]
```

```
Internet Control Message Protocol
  Type: 3 (Destination unreachable)
  Code: 4 (Fragmentation needed)
Checksum: 0xbbac [correct]
[Checksum Status: Good]
```

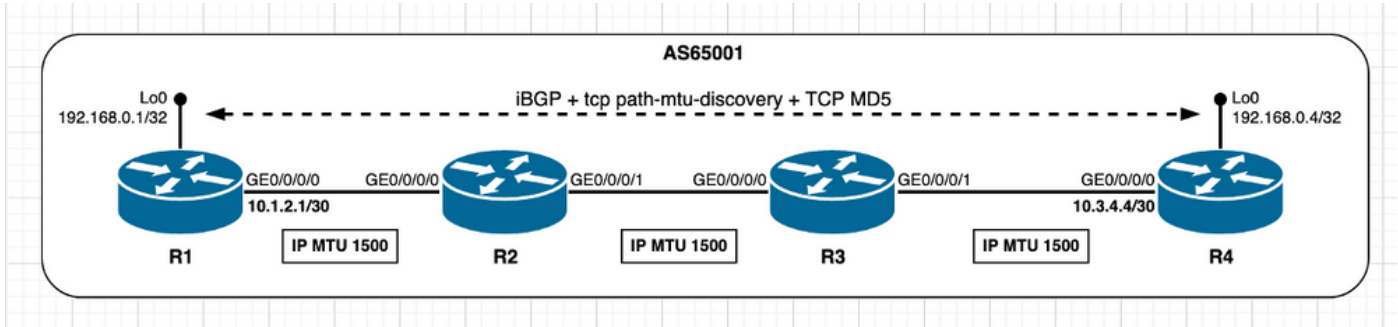


```

Length: 17
[Length of original datagram: 68]
Unused: 0011
MTU of next hop: 504
Internet Protocol Version 4, Src: 10.1.14.1, Dst: 10.1.14.14
  0100 .... = Version: 4
  .... 0101 = Header Length: 20 bytes (5)
  Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)
  Total Length: 610
  Identification: 0x7c9f (31903)
  Flags: 0x02 (Don't Fragment)
    0... .... = Reserved bit: Not set
    .1.. .... = Don't fragment: Set
    ..0. .... = More fragments: Not set
  Fragment offset: 0
  Time to live: 255
  Protocol: TCP (6)
  Header checksum: 0xcce5 [validation disabled]
  [Header checksum status: Unverified]
  Source: 10.1.14.1
  Destination: 10.1.14.14
  [Source GeoIP: Unknown]
  [Destination GeoIP: Unknown]
Transmission Control Protocol, Src Port: 22008, Dst Port: 23, Seq: 586828435, Ack: 754580617

```

### PMTUD - TCP選項(MD5)



映像3.4 — 啟用PMTUD和TCP MD5身份驗證。

在啟用TCP MD5驗證的情況下，不會從先前案例中說明的內容區分PMTUD行為。與先前在使用中的TCP MD5驗證共用一樣，Cisco IOS XR會考慮其他額外負荷，而作用中TCP對等體的初始MSS值反映相同。請參閱前面的使用TCP選項 — XR主動和使用TCP選項 — XR被動小節，以瞭解有關使用TCP選項的影響的其他詳細資訊。此案例中的TCP MSS計算可總結如下：

- 所有節點均使用預設IP MTU 1500位元組。
- 已啟用TCP路徑MTU發現。
- TCP對等點沒有直接連線。
- 在R1和R4上啟用了TCP MD5身份驗證。R4管理BGP連線。R4傳送MSS為1436位元組的SYN。1500 (介面IP MTU) — 20(minTCP\_H)- 20(minIP\_H)- 24位元組 (IOS XR TCP選項額外負荷)。R1傳送SYN, ACK, MSS為1436位元組。傳送[Received MSS ;本地初始MSS]。收到1436位元組；本地初始MSS 1460位元組。兩個對等點上使用最低的MSS值。

來源為R4的TCP SYN:

```
! - TCP SYN sourced from R4
```

```
2408 5.695076 192.168.0.4 192.168.0.1 TCP 82 59050 179 [SYN] Seq=0 Win=16384
Len=0 MSS=1436 WS=1
```

```
Frame 2408: 82 bytes on wire (656 bits), 82 bytes captured (656 bits) on interface 0
Ethernet II, Src: fa:16:3e:d7:7e:f6 (fa:16:3e:d7:7e:f6), Dst: fa:16:3e:8f:8f:54
(fa:16:3e:8f:8f:54)
Internet Protocol Version 4, Src: 192.168.0.4, Dst: 192.168.0.1
Transmission Control Protocol, Src Port: 59050, Dst Port: 179, Seq: 0, Len: 0
  Source Port: 59050
  Destination Port: 179
  [Stream index: 8]
  [TCP Segment Len: 0]
  Sequence number: 0 (relative sequence number)
  Acknowledgment number: 0
  Header Length: 48 bytes
  Flags: 0x002 (SYN)
  Window size value: 16384
  [Calculated window size: 16384]
  Checksum: 0x20d7 [unverified]
  [Checksum Status: Unverified]
  Urgent pointer: 0
  Options: (28 bytes), Maximum segment size, Window scale, No-Operation (NOP), TCP MD5
signature, End of Option List (EOL)
    Maximum segment size: 1436 bytes
      Kind: Maximum Segment Size (2)
      Length: 4
      MSS Value: 1436
    Window scale: 0 (multiply by 1)
    No-Operation (NOP)
    TCP MD5 signature
    End of Option List (EOL)
```

### 源自R1的TCP SYN、ACK:

! - TCP SYN,ACK sourced from R1

```
2409 0.004352 192.168.0.1 192.168.0.4 TCP 82 179 59050 [SYN, ACK] Seq=0 Ack=1
Win=16384 Len=0 MSS=1436 WS=1
```

```
Frame 2409: 82 bytes on wire (656 bits), 82 bytes captured (656 bits) on interface 0
Ethernet II, Src: fa:16:3e:8f:8f:54 (fa:16:3e:8f:8f:54), Dst: fa:16:3e:d7:7e:f6
(fa:16:3e:d7:7e:f6)
Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4
Transmission Control Protocol, Src Port: 179, Dst Port: 59050, Seq: 0, Ack: 1, Len: 0
  Source Port: 179
  Destination Port: 59050
  [Stream index: 8]
  [TCP Segment Len: 0]
  Sequence number: 0 (relative sequence number)
  Acknowledgment number: 1 (relative ack number)
  Header Length: 48 bytes
  Flags: 0x012 (SYN, ACK)
  Window size value: 16384
  [Calculated window size: 16384]
  Checksum: 0xcbf8 [unverified]
  [Checksum Status: Unverified]
  Urgent pointer: 0
  Options: (28 bytes), Maximum segment size, Window scale, No-Operation (NOP), TCP MD5
signature, End of Option List (EOL)
    Maximum segment size: 1436 bytes
      Kind: Maximum Segment Size (2)
      Length: 4
      MSS Value: 1436
    Window scale: 0 (multiply by 1)
    No-Operation (NOP)
```

TCP MD5 signature  
End of Option List (EOL)

## R4上看到的TCP會話詳細資訊 — ACTIVE:

! - as seen from R4 - Active

RP/0/0/CPU0:R4#show tcp detail pcb 0x121542c0  
Tue Jan 12 13:27:23.526 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Tue Jan 12 13:25:41 2021

PCB 0x121542c0, SO 0x1213c0e4, TCPCB 0x12156010, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 359  
Local host: 192.168.0.4, Local port: 59050 (Local App PID: 1052958)  
Foreign host: 192.168.0.1, Foreign port: 179

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	6	1	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	3	2	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 3299472269 snduna: 3299473445 sndnxt: 3299473445  
sndmax: 3299473445 sndwnd: 31646 sndcwnd: 4308  
irs: 3225544359 rcvnx: 3225545535 rcvwnd: 31665 rcvadv: 3225577200

SRTT: 89 ms, RTTO: 530 ms, RTV: 441 ms, KRTT: 0 ms  
minRTT: 19 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 30, connect retry interval: 30 secs

State flags: none  
Feature flags: **MD5**, Win Scale, Nagle, **Path MTU**  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1436, peer MSS 1436, min MSS 1436, max MSS 1436**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

PDU information:  
#PDU's in buffer: 0

FIB Lookup Cache: IFH: 0x40 PD ctx: size: 0 data:  
Num Labels: 0 Label Stack:

RP/0/0/CPU0:R4#

## R1上看到的TCP會話詳細資訊 — PASSIVE:

! - as seen from R1 - Passive

RP/0/0/CPU0:R1#show tcp detail pcb 0x121560ec  
Tue Jan 12 13:25:59.310 UTC

=====  
Connection state is ESTAB, I/O status: 0, socket status: 0  
Established at Tue Jan 12 13:25:31 2021

PCB 0x121560ec, SO 0x121556d4, TCPCB 0x121575bc, vrfid 0x60000000,  
Pak Prio: Medium, TOS: 192, TTL: 255, Hash index: 359  
Local host: 192.168.0.1, Local port: 179 (Local App PID: 983326)  
Foreign host: 192.168.0.4, Foreign port: 59050

Current send queue size in bytes: 0 (max 24576)  
Current receive queue size in bytes: 0 (max 32768) mis-ordered: 0 bytes  
Current receive queue size in packets: 0 (max 0)

Timer	Starts	Wakeups	Next(msec)
Retrans	3	0	0
SendWnd	0	0	0
TimeWait	0	0	0
AckHold	3	2	0
KeepAlive	1	0	0
PmtuAger	0	0	0
GiveUp	0	0	0
Throttle	0	0	0

iss: 3225544359 snduna: 3225545516 sndnxt: 3225545516  
sndmax: 3225545516 sndwnd: 31684 sndcwnd: 4308  
irs: 3299472269 rcvnxt: 3299473426 rcvwnd: 31665 rcvadv: 3299505091

SRTT: 37 ms, RTTO: 300 ms, RTV: 244 ms, KRTT: 0 ms  
minRTT: 9 ms, maxRTT: 239 ms

ACK hold time: 200 ms, Keepalive time: 0 sec, SYN waittime: 30 sec  
Giveup time: 0 ms, Retransmission retries: 0, Retransmit forever: FALSE  
Connect retries remaining: 0, connect retry interval: 0 secs

State flags: none  
Feature flags: MD5, Win Scale, Nagle, Path MTU  
Request flags: Win Scale

**Datagrams (in bytes): MSS 1436, peer MSS 1436, min MSS 1460, max MSS 1460**

Window scales: rcv 0, snd 0, request rcv 0, request snd 0  
Timestamp option: recent 0, recent age 0, last ACK sent 0  
Sack blocks {start, end}: none  
Sack holes {start, end, dups, rxmit}: none

Socket options: SO\_REUSEADDR, SO\_REUSEPORT, SO\_NBIO  
Socket states: SS\_ISCONNECTED, SS\_PRIV  
Socket receive buffer states: SB\_DEL\_WAKEUP  
Socket send buffer states: SB\_DEL\_WAKEUP  
Socket receive buffer: Low/High watermark 1/32768  
Socket send buffer : Low/High watermark 2048/24576, Notify threshold 0

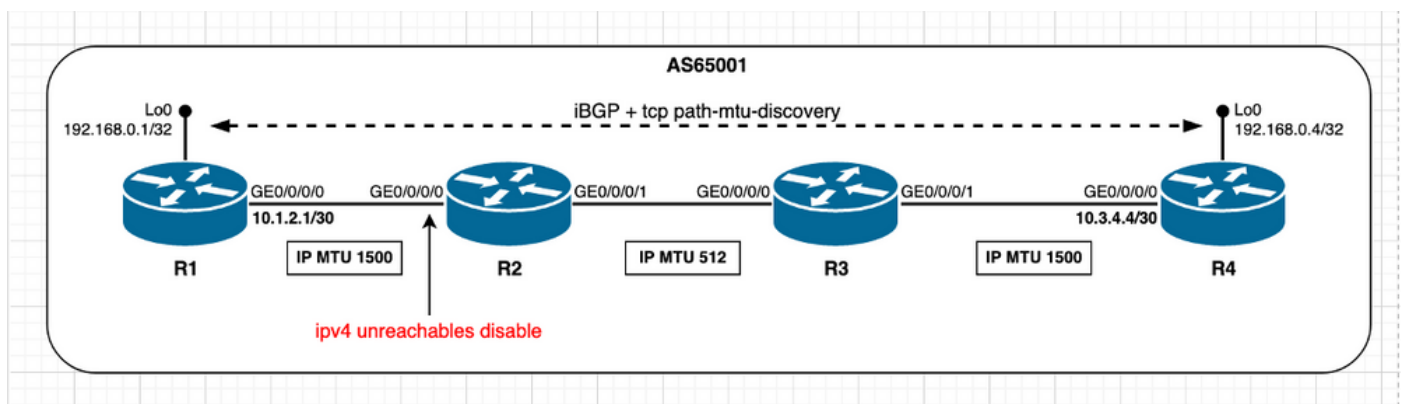
```
PDU information:
#PDU's in buffer: 0
FIB Lookup Cache:  IFH: 0x40  PD ctx: size: 0  data:
Num Labels: 0  Label Stack:
```

```
RP/0/0/CPU0:R1#
```

## PMTUD — 黑孔偵測

如前文PMTUD — 路徑區段的IP MTU更低一節所述，當啟用時，TCP PMTUD會透過接收ICMP(目的地無法連線 — 型別3;需要分段 — 代碼4)消息。在某些情況下，由於某些原因未收到這些消息，因此不會觸發PMTUD。在這種情況下，不會得知TCP對等路由器之間路徑的最小IP MTU。如果IP封包已設定DF位元，且它們的大小高於最低IP MTU路徑區段，則此類情況會引入潛在的黑洞。這些資料包將被靜默丟棄。

本節旨在重點介紹Cisco IOS XR如何檢測此類潛在黑洞場景並對其執行操作。為此，在R2介面GE0/0/0/0上禁用IPv4不可達功能，如下一映像和CLI輸出所示。



映像3.5 — 在R1/R4和R2 IPv4上啟用PMTUD且無法連線。

R2上禁用的IPv4無法連線：

```
!- R2 - IP unreachable is disabled
```

```
RP/0/0/CPU0:R2#show run interface gigabitEthernet 0/0/0/0
Thu May 13 12:09:45.483 UTC
interface GigabitEthernet0/0/0/0
 ipv4 address 10.1.2.2 255.255.255.252
 ipv4 unreachable disable
!
```

```
RP/0/0/CPU0:R2#show ipv4 interface gigabitEthernet 0/0/0/0
Thu May 13 12:10:04.112 UTC
GigabitEthernet0/0/0/0 is Up, ipv4 protocol is Up
 Vrf is default (vrfid 0x60000000)
 Internet address is 10.1.2.2/30
 MTU is 1514 (1500 is available to IP)
 Helper address is not set
 Multicast reserved groups joined: 224.0.0.2 224.0.0.1 224.0.0.5
 224.0.0.6
 Directed broadcast forwarding is disabled
 Outgoing access list is not set
 Inbound common access list is not set, access list is not set
 Proxy ARP is disabled
 ICMP redirects are never sent
 ICMP unreachable are never sent
```

ICMP mask replies are never sent  
Table Id is 0xe0000000

Cisco IOS XR處理此黑洞情況的方式是重新傳輸相同封包兩次，如果仍不成功，即未收到預期的TCP ACK，則重試但使用下一個較低的、定義良好的平台值，如[RFC1191 — 路徑MTU探索](#)中所述（請參閱[PMTUD — 路徑片段具有較低的IP MTU](#)一節，瞭解平台清單）。總而言之，Cisco IOS XR假設在到達目的地的路徑內因資料包大小而丟棄資料包，並嘗試通過資料包重新傳輸來解決這種情況。通過節點R1介面處捕獲的資料包以及debug tcp pmtud的輸出中的下一個示例，可以觀察到此行為。

R1上的IOS-XR黑洞檢測：

```
! - at R1
! - Original BGP Update message is sent
! - Note IP Total Length of 1116 bytes and TCP Segment Length of 1076 bytes
! - R2 filters such packet and send and ICMP error message towards R1 which triggers PMTUD
! - But because IPv4 unreachable are disabled at R2 GE0/0/0/0 ICMP message is not sent
! - Hence BGP message is silently filtered at R2

562      7.638774          192.168.0.1 192.168.0.4 BGP      1130    UPDATE Message, KEEPALIVE Message

Frame 562: 1130 bytes on wire (9040 bits), 1130 bytes captured (9040 bits) on interface 0
Ethernet II, Src: fa:16:3e:42:18:05 (fa:16:3e:42:18:05), Dst: fa:16:3e:5c:f1:80
(fa:16:3e:5c:f1:80)
Internet Protocol Version 4, Src: 192.168.0.1, Dst: 192.168.0.4
 0100 .... = Version: 4
  .... 0101 = Header Length: 20 bytes (5)
Differentiated Services Field: 0xc0 (DSCP: CS6, ECN: Not-ECT)
Total Length: 1116
Identification: 0x4a37 (18999)
Flags: 0x02 (Don't Fragment)
 0... .... = Reserved bit: Not set
 .1.. .... = Don't fragment: Set
 ..0. .... = More fragments: Not set
Fragment offset: 0
Time to live: 255
Protocol: TCP (6)
Header checksum: 0x229b [validation disabled]
[Header checksum status: Unverified]
Source: 192.168.0.1
Destination: 192.168.0.4
[Source GeoIP: Unknown]
[Destination GeoIP: Unknown]
Transmission Control Protocol, Src Port: 179, Dst Port: 57082, Seq: 318, Ack: 251, Len: 1076
Border Gateway Protocol - UPDATE Message
Border Gateway Protocol - KEEPALIVE Message
<snip>

! - at R1
! - No TCP ACK is received
! - Packet retransmission is attempted (2 attempts)
! - Note initial MSS value is of 1460 bytes

563      0.560058          192.168.0.1 192.168.0.4 TCP      1130    [TCP Retransmission] 179  57082
[PSH, ACK] Seq=318 Ack=251 Win=32593 Len=1076
564      1.101367          192.168.0.1 192.168.0.4 TCP      1130    [TCP Retransmission] 179  57082
[PSH, ACK] Seq=318 Ack=251 Win=32593 Len=1076

! - at R1
! - Still no TCP ACK received; previous retransmissions failed
! - Next lower plateau value is attempted - 1492 bytes
! - Packet retransmission is attempted (2 attempts)
```

RP/0/0/CPU0:May 13 10:20:44.251 UTC: tcp[399]: [t1] PCB 0x15392224: Trying next lower MTU: 1452

567 1.850294 192.168.0.1 192.168.0.4 TCP 1130 [TCP Retransmission] 179 57082  
[PSH, ACK] Seq=318 Ack=251 Win=32593 Len=1076  
568 1.111361 192.168.0.1 192.168.0.4 TCP 1130 [TCP Retransmission] 179 57082  
[PSH, ACK] Seq=318 Ack=251 Win=32593 Len=1076

! - at R1  
! - Still no TCP ACK received; previous retransmissions failed  
! - Next lower plateau value is attempted - 1006 bytes  
! - Packet retransmission is attempted (2 attempts)

RP/0/0/CPU0:May 13 10:20:47.560 UTC: tcp[399]: [t1] PCB 0x15392224: Trying next lower MTU: 966

569 2.198327 192.168.0.1 192.168.0.4 TCP 1020 [TCP Retransmission] 179 57082  
[ACK] Seq=318 Ack=251 Win=32593 Len=966  
570 1.109602 192.168.0.1 192.168.0.4 TCP 1020 [TCP Retransmission] 179 57082  
[ACK] Seq=318 Ack=251 Win=32593 Len=966

! - at R1  
! - Still no TCP ACK received; previous retransmissions failed  
! - Next lower plateau value is attempted - 508 bytes  
! - Original information (TCP Length of 1076 bytes) is split in three distinct packets  
! - TCP Segment Lengths 468 + 468 + 140 = 1076  
! - TCP ACK is received from peer R4

RP/0/0/CPU0:May 13 10:20:50.870 UTC: tcp[399]: [t1] PCB 0x15392224: Trying next lower MTU: 468

571 2.205552 192.168.0.1 192.168.0.4 TCP 522 [TCP Retransmission] 179 57082  
[ACK] Seq=318 Ack=251 Win=32593 **Len=468**  
573 0.004254 192.168.0.1 192.168.0.4 TCP 522 [TCP Retransmission] 179 57082  
[ACK] Seq=786 Ack=251 Win=32593 **Len=468**  
574 0.002724 192.168.0.1 192.168.0.4 TCP 194 [TCP Retransmission] 179 57082  
[PSH, ACK] Seq=1254 Ack=251 Win=32593 **Len=140**

! - Peer R4 TCP ACK is received

575 0.223172 192.168.0.4 192.168.0.1 TCP 54 57082 179 [ACK] Seq=251 Ack=1394  
Win=31469 Len=0