# DYNAMIC ENABLEMENT OF LSO (LARGESEND OFFLOAD) IN NETWORK VIRTUALIZED ENVIRONMENT FOR BETTER NETWORK THROUGHPUT

S.Annie Christilla

Associate Professor

Department of Computer Science, St. Francis De Sales College, Bangalore, Karnatka, India

## ABSTRACT

LSO is a very important feature of a GigE/10G NICs providing a good amount of performance benefits. Using this feature Network layers can send bigger packets instead of smaller packets. On virtualized environmentNetwork Interface Card(NIC) will be attached to VIOS(Virutal IO servers) and logical partitions will be sharing this NIC through virtualization technique. In these kind of environment, LSO feature can be turned on/off on the NIC, the bridge layerSEA (Shared Ethernet Adapter) present in VIOS and in the virtual adapter present in logical partition. LSO need to be enabled in all these components to exploit this feature. If LSO is turned off on the bridge(SEA) present in the VIOS, will cause poor network performance. In this paper, method to achieve better network throughput when LSO feature is turned off/on dynamically is being proposed in this draft.

*keywords*: LSO(Largesend Offload), VIOS(Virtual IO Server), SEA(Shared Ethernet Adapter), NIC (Network Interface Card), LPAR ( Logical PARtition)

## 1. INTRODUCTION

Maximum Transmission Unit (MTU) of a network is the maximum protocol data unit that can be transferred on the physical medium. MTU is an inherent property of the physical media. For instance MTU in Ethernet is 1500 bytes. In a Network Protocol Stack, Network Layer or Internet Protocol (IP) layer implements datagram fragmentation so that packets with size larger than the network interface's MTU are fragmented to MTU size before being delivered to the data link layer. Transport protocols such as TCP negotiate MSS (Maximum Segment Size) during connection establishment which is the largest amount of data that TCP is willing to send in a single segment. To avoid IP fragmentation, MSS is always set lower than MTU. Large Send or TCP Segmentation Offload (TSO) is a feature supported by Network Adapters in which the job of fragmenting a larger packets into MTU size is done by the Network Interface Card (NIC) in hardware. network protocol stack can send larger size packets (without having to do the job of fragmentation in software) to NICs and the NIC hardware will do the fragmentation in hardware which would help in improving performance.LSO is a very important feature of a GigE/10G NICs providing a good amount of performance benefits. Using this feature Network layers can send bigger packets instead of smaller packets. However if the network applications have been written in a way to send smaller packets, this hardware feature cannot be exploited well. This is because if applications send smaller packets through the network stack, the stack will end up sending smaller packets to the NIC. With this there would be a lot of packets being sent down from the application to the NIC through the network protocol stack in kernel. This can be avoided by making applications send large size packets so that the lesser number of packets are being sent down by the protocol stack to the NIC

## 2. PROBLEM STATEMENT

On virtualized environment, LSO is turned on NIC, SEA and Virtual Adapter on the logical partition. When TCP connection is established from LPAR to host that resides outside the system, LPAR will be sending bigger packets upto 64k to the VIOS. The NIC on the VIOS will segment this packet based on MTU and send the data out. When the connection is established if LSO is turned off on SEA, LPAR will keep sending large packet with IF_DF flag on. As a result SEA will not be able to fragment this packet and send ICMP back to LPAR. However LPAR will still continue to send bigger packets. When the retransmission timer expires LPAR will send one packet with size of 1500 bytes. This will result poor network throughput as after every retransmission time out one packet will be sent to the remote host.

## 3. EXISTING METHOD

When TCP connection gets establishedLPAR will be sending TCP options0x0E0303 along with SYN packet. If SEA has LSO option turned on while it gets response (SYNACK) back from the destinationit will piggy pack the same TCP option. Upon receiving SYNACK and ifTCP option present, the Transport layer onLPAR will identify LSO is turned on SEA and turns on LSO flag for the connection.Post that transport layer will be sending larger packetsto the VIOS. NIC on VIOS will in turn fragmentthe packet and send it to the destination.TCP connection on the LPAR will keep sending largePacketsupto 64k. The checksum field of the packet will have the MSS value that adapter will use the packet to fragment. Later if LSO is turned off on SEA, and largePacket is received from LPAR(>1500 bytes), packets will be fragmented in the SEA layer and sent down to adapter as small packets. If packet also has IF_DF flag, then those packets will get dropped in SEA layer. In this case SEA willbe sending ICMP "fragmention needed" packetback to the LPAR. On receiving this packet LPAR will turn off LSO feature for the connection. Later if LSO is enabled on SEA, there is no wayTCP connections on the LPAR will be aware of and still will be sending smaller packets. This Results in poor network performance.

## 4. PROPOSED METHOD

In this paper, an algorithm to achieve better network throughput is proposed. On VIOS at SEA layer information (Source IP, Source Port, Destination IP, Destination Port, Largesend capable connection or not) about each connection will be maintained. When LSO is turned off, ICMP message with type 42 will be generated to Source IP, Source Port with information whether LSO is turned off or not. ICMP message will be generated for all the connections. On LPAR once it receives the ICMP message, it turns

**Conference Paper:** International Conference on "Recent Advances in Computing and Communication"
**Organized by:** Department of Computer Science, SSS Shasun Jain College for Women, Chennai, India

ICT ACADEMY
Innovate... Collaborate... Educate...

43

of the LSO feature in Transport layer for the specified connection. ICMP message will also have Destination IP and Destination Port and protocol no in the Data section. This will help the LPAR to find the correct connection and turn off the LSO feature. Later lets assume LSO is turned on on SEA then again the ICMP packets will get generated for the LPARS with the information LSO is turned on. Upon receiving this message LPARs will turn on the LSO for the particular TCP connection. At SEA layer, when it receives SYN Packet the entry about connection will be added. This entry will be removed upon receiving FIN for the same connection.

ICMP message used for communication between SEA and LPAR will be type 42 and code will have value either 0 or 1. 0 represents LSO is turned off on SEA and 1 represents LSO is turned on at the SEA layer. LSO at SEA can only be turn on if the LSO is enabled on the underlying network adapter. Data will have SourceIP, Source Port, Destination IP, Destination port and protocol no. Refer Figure 2 for the ICMP format.
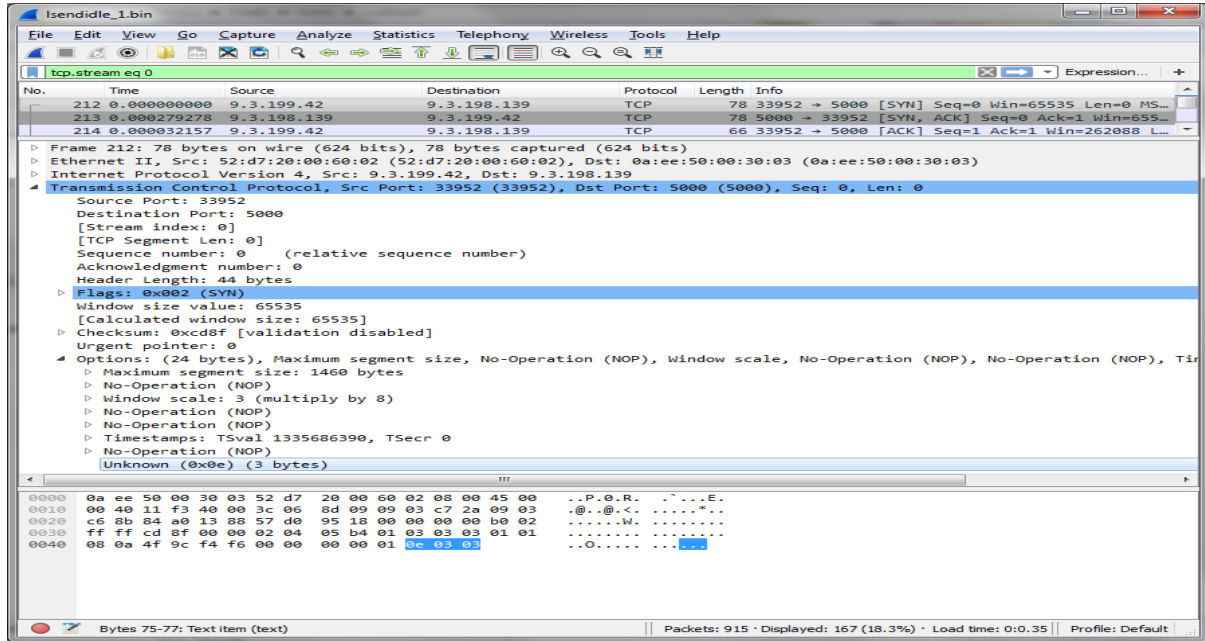


Figure 1. TCP Option to negotiate LSO between LPAR and SEA

**Algorithm:**

1) During connection establishment phase both LPAR and SEA will exchange the TCP option 0x0E0303 if they are LSO capable.

2) If they are LSO capable, TCP layer on LPAR will send large amount of data ( > 1500 bytes) and upto 64K

3) On VIOS entry will be added in the connection cache table with information SourceIP, source Port, Destination IP, Destination port, LSO capable and Protocol number

4) If LSO is turned on the virtual adapter on LPAR, call back will be called and it will search the TCP control block and will turn off LSO for the tcp connections.

5) If LSO is turned off on the SEA, SEA will go through control connection cache table and send ICMP for each entry. ICMP packet type will be 0x42 and code 0 as LSO is turned off. This packet will also have information Source IP, Source Port, Destination IP, Destination port and protocol number.

6) Upon receiving this ICMP packet, LPAR will disable the LARGESEND flag on TCP control block.

7) While sending data, TCP layer on LPAR will not send bigger packet if LARGESEND flag is not set. ( As a result SEA need not fragment the packet).

8) Now lets consider a case when LSO is turned on at SEA, then again SEA layer will go through control connection cache table and send ICMP for each entry. ICMP packet type will be 0x42 and code will be 1 as LSO is enabled. This packet will also have information Source IP, Source Port, Destination IP, Destination port and protocol number.

9) Upon receiving this ICMP packet, LPAR will enable the LARGESEND flag on TCP control block.

10) While sending data, TCP layer on LPAR will send bigger packet if LARGESEND flag is set.

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Type | | | | | | | | Code | | | | | | | | ICMP header checksum | | | | | | | | | | | | | | | |
| Data ::: | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Figure 2. ICMP Protocol format

**5. CONCLUSION**

The objective of the proposed method is to improve the network throughput by transferring the correct sized data to NIC on the VIOS. The proposed method updates the LSO capability of the SEA/underlying NIC to LPAR dynamically as and when it happens, thereby making sure LPAR sends the correct sized data to the NIC on VIOS. This eliminates the fragmentation effort on

VIOS and also when LSO is enabled uses the feature to full extent. Hence the network throughput is improved using this method.

## REFERENCES

1. Retreived from https://tools.ietf.org/html/rfc792.

2. Hai Lin, Lucio Correia, & et. al., IBM PowerVM Virtualization

3. Gary R. Wright &‎ W. Richard Stevens," TCP/IP Illustrated-Implementation", Vol. 2

4. Kumar Reddy, "Network Virtualization".