

Network Working Group
Internet-Draft
Intended Status: Experimental
Expires: September 14, 2014

H. Zhou
C. Li
eBay Inc.
March 13, 2014

Segmentation Offloading Extension for VxLAN
draft-zhou-li-vxlan-soe-00

Abstract

Segmentation offloading is nowadays common in network stack implementation and well supported by para-virtualized network device drivers for virtual machine (VM)s. This draft describes an extension to Virtual eXtensible Local Area Network (VXLAN) so that segmentation can be decoupled from physical/underlay networks and offloaded further to the remote end-point thus improving data-plane performance for VMs running on top of overlay networks.

Status of this Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/lid-abstracts.html>

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>

Copyright and License Notice

Copyright (c) 2014 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal

Internet-Draft

VXLAN-soe

March 2014

Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1	Introduction	3
1.1	Requirements Notation	4
1.2	Definition of Terms	4
2	Approach	4
2.1	VXLAN Header Extension	4
2.2	TX VTEP	5
2.3	RX VTEP - Hypervisors	6
2.4	RX VTEP - Gateways	6
3	Interoperability	6
4	Security Considerations	6
5	IANA Considerations	6
6	References	7
6.1	Normative References	7
6.2	Informative References	7
	Authors' Addresses	7

Internet-Draft

VXLAN-soe

March 2014

1 Introduction

Network virtualization over L3 transport is evolved along with server virtualization in data-centers, and data plane performance is one of the keys to the success of this combination. One of the most critical improvements in OS kernel TCP/IP stack in recent years is segmentation offloading, and now hypervisor providers support same mechanism in para-virtualized Ethernet drivers so that virtual servers can benefit from the same mechanism in virtualized world by offloading segmentation tasks to the lowest layer on hypervisors or NICs (if TSO is supported by the NICs equipped in the hypervisor).

Essentially, overlay networks has its own advantage comparing with physical underlay networks in that it does not have a hard MTU limitation. Therefore, segmentation offloading can be pushed to the remote end-point of the transport tunnel, where segmentation can be completely omitted if this remote end-point is on a hypervisor. However, this advantage is not utilized when the transport of the overlay is based on the Virtual eXtensible Local Area Network [I-D.mahalingam-dutt-dcops-vxlan], which provides a transport mechanism for logically isolated L2 overlay networks between hypervisors. Lacking segmentation information in the VXLAN header, hypervisor implementations have to make pessimistic decisions to always segment the packet in the size specified by VMs before delivering to hypervisors' IP stack, because it does not know whether the remote end-point is bridged to a physical network with hard MTU limitations. It is worth noting that the segmentation here is not the IP fragmentation in terms of the physical network MTU, which may still follow if the segment size resulting from the process above plus the tunnel outer header is bigger than the physical network MTU.

To fulfill the potential of segmentation offloading on overlay, this draft introduces segmentation metadata in VXLAN header. With the capability of carrying segmentation metadata in packets, hypervisors can offload the segmentation decision further to the remote tunnel end-point, thus decoupling the segmentation for overlay from physical

limitations of underlay, providing higher flexibility to hypervisor implementations to achieve significant performance gains in a major part of VXLAN deployment scenarios.

Although the performance gains can be achieved is affected by the physical network MTU, there is inherently no mandatory requirement to physical layer:

1) When physical network MTU is far bigger than overlay MTU, the offloading reduces the number of packets being transmitted by TX hypervisors and received in RX hypervisors and RX VMs.

2) When physical network MTU is close to overlay MTU, the number of packets being transmitted in physical network (resulted in IP fragmentation) may not be reduced significantly, but on RX side after IP reassembling, the number of packets being delivered from the hypervisor to the receiving VM is largely reduced, thus saving the cost of hypervisor <-> VM interaction and protocol stack of the receiving VM. Furthermore, a minor cost saving is that the bytes being transmitted over physical network is slightly reduced because only one copy of headers (inner L2-L4 header, VXLAN header and outer UDP header) is transmitted for a large overlay packet.

In addition, offloading features support from NIC hardware is NOT a requirement, either, to the performance gains discussed above.

1.1 Requirements Notation

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#) [[RFC2119](#)].

1.2 Definition of Terms

GS0: Generic Segmentation Offload.

TS0: TCP Segmentation Offload.

NIC: Network Interface Card.

VM: Virtual Machine.

TX: Sending side.

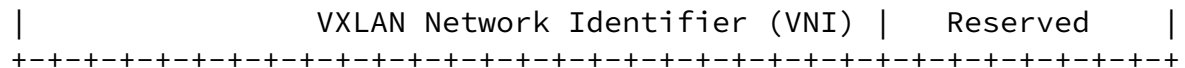
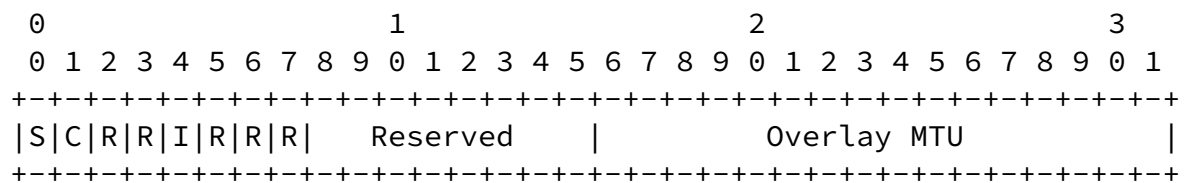
RX: Receiving side.

VTEP: Virtual Tunnel End Point

2. Approach

2.1 VXLAN Header Extension

The new VXLAN Segmentation Offloading Extension (VXLAN-soe) header is defined as:



The changes to VXLAN are:

S Bit: Flag bit 0 is defined as the S (Segmentation Offloading Extension) bit.

S = 1 indicates that VXLAN-soe is applied to the encapsulated overlay packet, and the C Bit and Overlay MTU field (see below) are valid.

S = 0 indicates that VXLAN-soe is NOT applied, and the C Bit and Segment Size field MUST be set to 0 in accordance with VXLAN.

C Bit: Flag bit 1 is defined as the C (Checksum) bit. This bit is valid only if the S bit is set to 1.

C = 1 indicates that the checksum to the encapsulated packet is required, and SHALL be re-calculated when the segmentation is being performed.

C = 0 indicates that the checksum to the encapsulated packet is NOT needed.

Overlay MTU: bit 16 - 31 is defined as the MTU desired by TX VM for the segmentation being offloaded.

Its value indicates the max size of an overlay segment including its L3 header, but NOT including Ethernet header. This field is valid only if the S bit is set.

[2.2 TX VTEP](#)

VTEP at TX side MUST set the S bit to 1 if the packet to be encapsulated is NOT segmented and it decides to offload the segmentation to the remote end-point. In such case the C bit and Overlay MTU field MUST be set accordingly. This is the typical use case when the TX VTEP is a hypervisor transmitting TCP stream of VMs with large sliding windows.

VTEP at TX side MUST clear the S bit if the packet to be encapsulated is segmented already or does NOT need to be segmented in terms of the overlay MTU. In such case, the encapsulation is in the same format as specified in VXLAN. This is the typical use case when the TX VTEP is a hypervisor transmitting small size overlay packets, or a gateway forwarding overlay packets without

offloading requirements.

[2.3 RX VTEP - Hypervisors](#)

When a VTEP at RX side is on a hypervisor, checking of the S bit is OPTIONAL.

[2.4 RX VTEP - Gateways](#)

When a VTEP at RX side is on a gateway node that connects overlay networks and physical networks, the S bit MUST be checked and the VTEP MUST ensure the segmentation specified by the header fields is performed by the VTEP itself or offloaded further - it MAY offload the segmentation again to the subsequent transmission mechanisms: such as GSO and TSO, or, if the link to the next hop

is also an overlay based on VXLAN-soe (or other tunneling protocols that supports segmentation offloading), pass the segmentation metadata to the next hop.

3 Interoperability

In addition to offload segmentation requests from VMs, VXLAN-soe enabled VTEP is able to offload segmentation requests from STT [I-D.davie-stt] overlay, because the metadata required in VXLAN-soe header is a subset of STT metadata. The additional segmentation offloading information carried in STT metadata such as L4 offset can be obtained by examine inner headers of the packets.

VXLAN-soe defines Overlay MTU at the same position of Protocol Type field in VXLAN-gpe [I-D.quinn-vxlan-gpe], another extension of VXLAN. This is not a problem because VXLAN-soe is introduced for segmentation offloading use cases where Ethernet header is always encapsulated, and it uses different flag bits to be distinguished from VXLAN-gpe.

4 Security Considerations

There is no special security issues introduced by this extension to VXLAN.

5 IANA Considerations

This document creates no new requirements on IANA namespaces [RFC5226].

6 References

6.1 Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

[6.2](#) Informative References

[I-D.mahalingam-dutt-dcops-vxlan]

Mahalingam, M., Dutt, D., Duda, K., Agarwal, P., Kreeger, L., Sridhar, T., Bursell, M., and C. Wright, "VXLAN: A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks", [draft-mahalingam-dutt-dcops-vxlan-08](#) (work in progress), February 2014.

[I-D.davie-stt]

Davie, B. and J. Gross, "A Stateless Transport Tunneling Protocol for Network Virtualization (STT)", [draft-davie-stt-05](#)(work in progress), March 2014.

[I-D.quinn-vxlan-gpe]

Agarwal, P., Fernando, R., Kreeger, L., Lewis, D., Maino, F., Quinn, P., Yong, L., Xu, X., Smith, M., Yadav, N., and U. Elzur, "Generic Protocol Extension for VXLAN", [draft-quinn-vxlan-gpe-02](#) (work in progress), December 2013.

Authors' Addresses

Han Zhou
eBay, Inc.

EMail: hzhou8@ebay.com

Chengyuan Li
eBay, Inc.

Email: chengyli@ebay.com