                     BIER Use Case in Data Centers
                    draft-wang-bier-vxlan-use-case-02

Abstract

   Bit Index Explicit Replication (BIER) is an architecture that
   provides optimal multicast forwarding through a "BIER domain" without
   requiring intermediate routers to maintain any multicast related per-
   flow state.  BIER also does not require any explicit tree-building
   protocol for its operation.  A multicast data packet enters a BIER
   domain at a "Bit-Forwarding Ingress Router" (BFIR), and leaves the
   BIER domain at one or more "Bit-Forwarding Egress Routers" (BFERs).
   The BFIR router adds a BIER header to the packet.  The BIER header
   contains a bit-string in which each bit represents exactly one BFER
   to forward the packet to.  The set of BFERs to which the multicast
   packet needs to be forwarded is expressed by setting the bits that
   correspond to those routers in the BIER header.

   This document tries to describe the drawbacks of how BUM services are
   deployed in current data centers, and proposes how to take full
   advantage of BIER to implement BUM services in data centers.

Status of This Memo

Copyright Notice

Table of Contents

## 1.  Introduction

   This document is motivated by [I-D.ietf-bier-use-cases].

   In current data center virtualization, virtual eXtensible Local Area
   Network (VXLAN) [RFC7348] is a kind of network virtualization overlay
   technology which is overlaid between NVEs and is intended for multi-
   tenancy data center networks, whose reference architecture is
   illustrated as per Figure 1.

```
    +--------+                                          +--------+
    | Tenant +--+                               +----| Tenant |
    | System |  |                               (')   | System |
    +--------+  |        ...............         (   )  +--------+
            |  +-+--+  .                 .  +--+-+  (_)
            |  | NVE|--.               .--| NVE|    |
            +--|    |  .               .  |    |---+
              +-+--+  .               .  +--+-+
               /       .             .
              /        . L3 Overlay  .  +--+-++--------+
    +--------+  /       .    Network   .  | NVE|| Tenant |
    | Tenant +--+        .             .--|    || System |
    | System |           .             .  +--+-++--------+
    +--------+           ...............
```
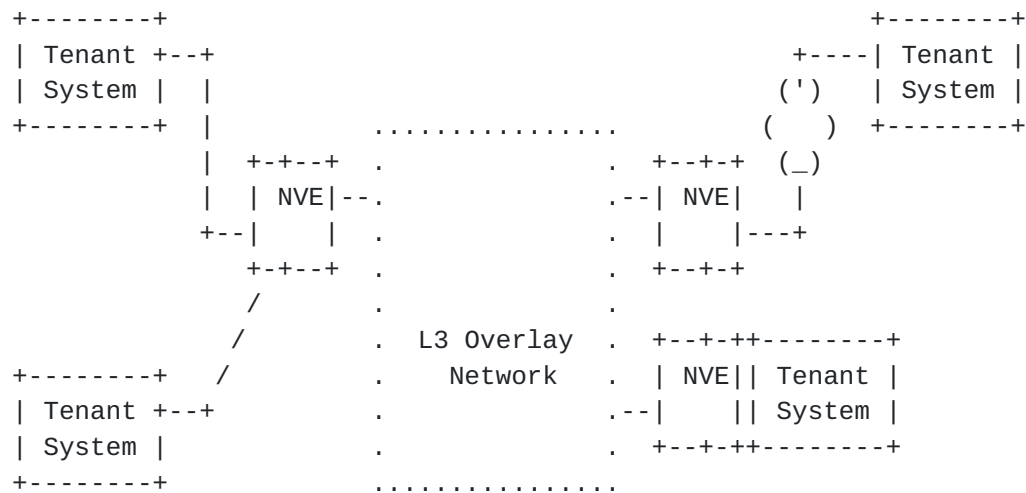

                       Figure 1: NVO3 Architecture

   And there are two kinds of most common methods about how to forward
   BUM packets in this virtualization overlay network.  One is using PIM
   as underlay multicast routing protocol to build explicit multicast
   distribution tree, such as PIM-SM[RFC4601] or PIM-BIDIR
   [RFC5015]multicast routing protocol.  Then, when BUM packets arrive
   at NVE, it requires NVE to have a mapping between the VXLAN Network
   Identifier and the IP multicast group.  According to the mapping, NVE
   can encapsulate BUM packets in a multicast packet which group address
   is the mapping IP multicast group address and steer them through
   explicit multicast distribution tree to the destination NVEs.  This
   method has two serious drawbacks.  It need the underlay network
   supports complicated multicast routing protocol and maintains
   multicast related per-flow state in every transit nodes.  What is
   more, how to configure the ratio of the mapping between VNI and IP
   multicast group is also an issue.  If the ratio is 1:1, there should
   be 16M multicast groups in the underlay network at maximum to map to
   the 16M VNIs, which is really a significant challenge for the data
   center devices.  If the ratio is n:1, it would result in inefficiency
   bandwidth utilization which is not optimal in data center networks.

   The other method is using ingress replication to require each NVE to
   create a mapping between the VXLAN Network Identifier and the remote
   addresses of NVEs which belong to the same virtual network.  When NVE
   receives BUM traffic from the attached tenant, NVE can encapsulate
   these BUM packets in unicast packets and replicate them and tunnel
   them to different remote NVEs respectively.  Although this method can
   eliminate the burden of running multicast protocol in the underlay
   network, it has a significant disadvantage: large waste of bandwidth,
   especially in big-sized data center where there are many receivers.

Bit Index Explicit Replication (BIER) [I-D.ietf-bier-architecture] is
an architecture that provides optimal multicast forwarding through a
"BIER domain" without requiring intermediate routers to maintain any
multicast related per-flow state.  BIER also does not require any
explicit tree-building protocol for its operation.  A multicast data
packet enters a BIER domain at a "Bit-Forwarding Ingress Router"
(BFIR), and leaves the BIER domain at one or more "Bit-Forwarding
Egress Routers" (BFERs).  The BFIR router adds a BIER header to the
packet.  The BIER header contains a bit-string in which each bit
represents exactly one BFER to forward the packet to.  The set of
BFERs to which the multicast packet needs to be forwarded is
expressed by setting the bits that correspond to those routers in the
BIER header.  Specifically, for BIER-TE, the BIER header may also
contain a bit-string in which each bit indicates the link the flow
passes through.

The following sections try to propose how to take full advantage of
overlay multicast protocol to carry virtual network information, and
create a mapping between the virtual network information and the bit-
string to implement BUM services in data centers.

## 2.  Convention and Terminology

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 [RFC2119].

The terms about BIER are defined in [I-D.ietf-bier-architecture].

The terms about NVO3 are defined in [RFC7365].

Here tries to list the most common terminology mentioned in this
draft.

BIER: Bit Index Explicit Replication(Bit Index Explicit Replication
(The overall architecture of forwarding multicast using a Bit
Position).

NVE: Network Virtualization Edge, which is the entity that implements
the overlay functionality.  An NVE resides at the boundary between a
Tenant System and the overlay network.

VXLAN: Virtual eXtensible Local Area Network

VNI: VXLAN Network Identifier

Virtal Network Context Identifier: Field in an overlay encapsulation
header that identifies the specific VN the packet belongs to.

3.  **BIER in data centers**

   This section tries to describe how to use BIER as an optimal scheme
   to forward the broadcast, unknown and multicast (BUM) packets when
   they arrive at the ingress NVE in data centers.

   The principle of using BIER to forward BUM traffic is that: firstly,
   it requires each ingress NVE to have a mapping between the Virtual
   Network Context Identifier and the bit-string in which each bit
   represents exactly one egress NVE to forward the packet to.  And
   then, when receiving the BUM traffic, the BFIR/Ingree NVE maps the
   receiving BUM traffic to the mapping bit-string, encapsulates the
   BIER header, and forwards the encapsulated BUM traffic into the BIER
   domain to the other BFERs/Egress NVEs indicated by the bit-string.

   Furthermore, as for how each ingress NVE knows the other egress NVEs
   that belong to the same virtual network and creates the mapping is
   the main issue discussed below.  Basically, BIER Multicast Listener
   Discovery is an overlay solution to support ingress routers to keep
   per-egress-router state to construct the BIER bit-string associated
   with IP multicast packets entering the BIER domain.  The following
   section tries to extend BIER MLD to carry virtual network
   information(such as Virtual Network Context identifier), and
   advertise them between NVEs.  When each NVE receive these
   information, they create the mapping between the virtual network
   information and the bit-string representing the other NVEs belonged
   to the same virtual network.

4.  **BIER MLD extension for Virtual Network information**

   Figure 2 draws the extension of the MLD report message format.  In
   order to support Virtual Network information advertisement, one bit
   of the reserved field need be defined to indicate that there is a
   extended MLD Virtual Network information TLV immediately following in
   this message.

```
    0                   1                   2                   3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |     Type      |     Code      |          Checksum             |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |     Maximum Response Delay     |          Reserved           |1|
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |                                                               |
   +                                                               +
   |                                                               |
   +                     Multicast Address                         +
   |                                                               |
   +                                                               +
   |                                                               |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   | VN Info type  |    length     |            Value              |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |                          ......                               |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```
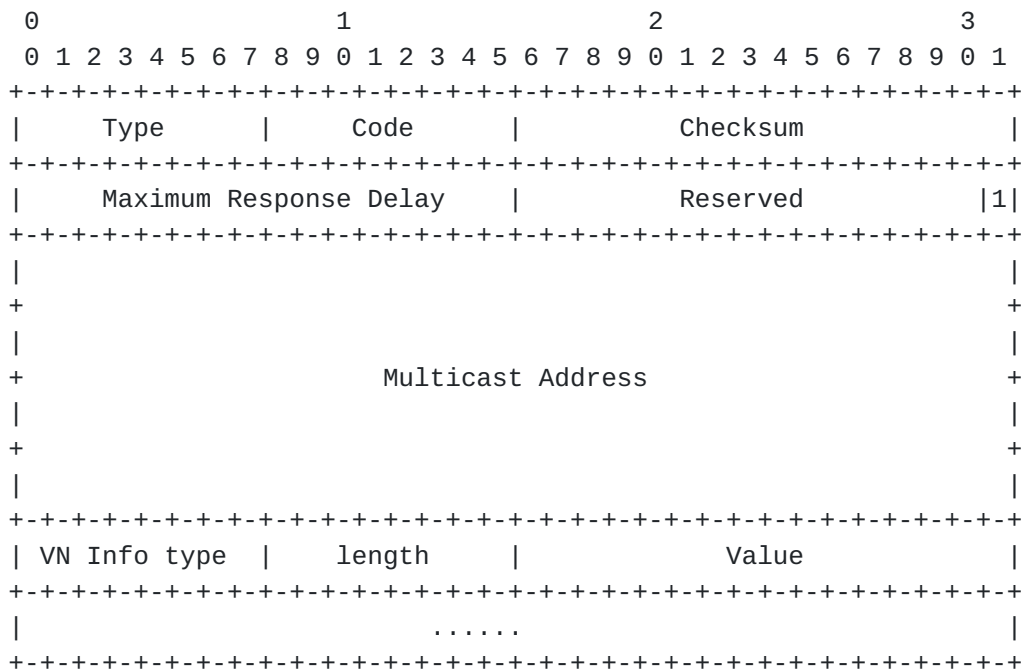
                    Figure 2: MLD message format

Specifically, Figure 3 illustrates the detailed Virtual Network
information TLV format in MLD report message to carry Virtual Network
information.

```
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |     type      |    length     |  Encap Type   |   Reserved    |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |      Virtual Network Context Identifier        |   Reserved    |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |      Virtual Network Context Identifier        |   Reserved    |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |                          ......                               |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```
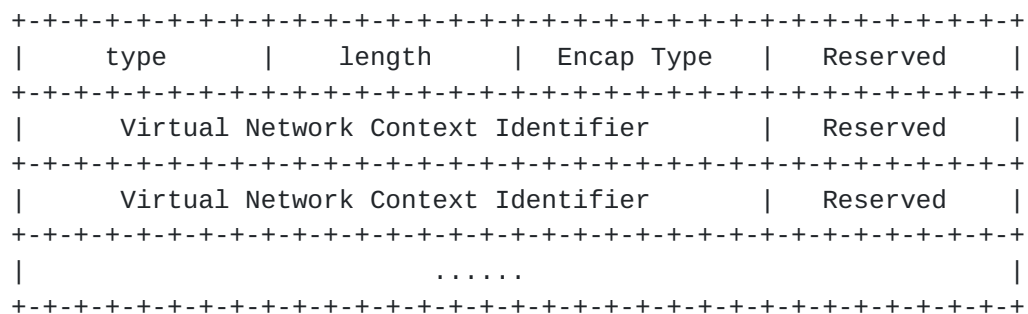
              Figure 3: MLD Virtual Network information TLV

   Type:

   indicates Virtual Network information TLV.  Value 1 indicates Virtual
   Network information advertisement.  Value 2 indicates Virtual Network
   information withdraw.

   Length: 1 cotet.

Encap Type:

indicates the encapsulation type the virtual netowrk using, such as
VxLAN,NVGRE,Geneve and so on.

Virtual Network Context Idenfifier:

indicates the identifier of the virtual network.  Different
encapsulation type has different meaning for this field.  In VxLAN
encasulation type, this field indicates VxLAN Network Identifier.  In
NVGRE encapsulation type, this field indicates Virtual Subnet
ID(VSID).In Geneve encapsulation type, this field indicates Virtual
Network Identifier.

Each NVE acquires the Virtual Network information, and advertises
this Virtual Network information to other NVEs through the MLD
messages.  If the NVE attaches to several virtual networks, it will
carry several Virtual Network Context Identifiers in the Virtual
Network advertisement message.  if the NVE supports several
encapsulation types, it will carry the Virtual Network Context
Identifiers belonging to one encapsulation type in one Virtual
Network information TLV.  if one attached virtual network is
migrated, the NVE will withdraw the Virtual Network information.

When ingress NVE receives the Virtual Network information
advertisement message, it builds a mapping between the receiving
Virtual Network Context Identifier in this message and the bit-string
in which each bit represents one egress NVE who sends the same
Virtual Network information.  Subsequently, once this ingress NVE
receives some other MLD advertisements which include the same Virtual
Network information from some other NVEs , it updates the bit-string
in the mapping and adds the corresponding sending NVE to the updated
bit-string.  Once the ingress NVE removes one virtual network, it
will delete the mapping corresponding to this virtual network as well
as send withdraw message to other NVEs.

After finishing the above interaction of MLD messages, each ingress
NVE knows where the other egress NVEs are in the same virtual
network.  When receiving BUM traffic from the attached virtual
network, each ingress NVE knows exactly how to encapsulate this
traffic and where to forward them to.

This can be used in both IPv4 network and IPv6 network.  In IPv4,
IGMP protocol does the similar extension for carrying Virtual Network
information TLV in Version 2 membership report message.

## 5.  Security Considerations

   It will be considered in a future revision.

## 6.  IANA Considerations

   There need one new Type for Virtual Network information TLV in MLD
   protocol and one in IGMP protocol.

## 7.  References

### 7.1.  Normative References

   [RFC2119]  Bradner, S., "Key words for use in RFCs to Indicate
              Requirement Levels", BCP 14, RFC 2119,
              DOI 10.17487/RFC2119, March 1997,
              <http://www.rfc-editor.org/info/rfc2119>.

   [RFC4601]  Fenner, B., Handley, M., Holbrook, H., and I. Kouvelas,
              "Protocol Independent Multicast - Sparse Mode (PIM-SM):
              Protocol Specification (Revised)", RFC 4601,
              DOI 10.17487/RFC4601, August 2006,
              <http://www.rfc-editor.org/info/rfc4601>.

   [RFC5015]  Handley, M., Kouvelas, I., Speakman, T., and L. Vicisano,
              "Bidirectional Protocol Independent Multicast (BIDIR-
              PIM)", RFC 5015, DOI 10.17487/RFC5015, October 2007,
              <http://www.rfc-editor.org/info/rfc5015>.

   [RFC7348]  Mahalingam, M., Dutt, D., Duda, K., Agarwal, P., Kreeger,
              L., Sridhar, T., Bursell, M., and C. Wright, "Virtual
              eXtensible Local Area Network (VXLAN): A Framework for
              Overlaying Virtualized Layer 2 Networks over Layer 3
              Networks", RFC 7348, DOI 10.17487/RFC7348, August 2014,
              <http://www.rfc-editor.org/info/rfc7348>.

   [RFC7365]  Lasserre, M., Balus, F., Morin, T., Bitar, N., and Y.
              Rekhter, "Framework for Data Center (DC) Network
              Virtualization", RFC 7365, DOI 10.17487/RFC7365, October
              2014, <http://www.rfc-editor.org/info/rfc7365>.

### 7.2.  Informative References

   [I-D.ietf-bier-architecture]
              Wijnands, I., Rosen, E., Dolganow, A., Przygienda, T., and
              S. Aldrin, "Multicast using Bit Index Explicit
              Replication", draft-ietf-bier-architecture-04 (work in
              progress), July 2016.

   [I-D.ietf-bier-idr-extensions]
              Xu, X., Chen, M., Patel, K., Wijnands, I., and T.
              Przygienda, "BGP Extensions for BIER", draft-ietf-bier-
              idr-extensions-01 (work in progress), June 2016.

   [I-D.ietf-bier-isis-extensions]
              Ginsberg, L., Przygienda, T., Aldrin, S., and Z. Zhang,
              "BIER support via ISIS", draft-ietf-bier-isis-
              extensions-02 (work in progress), March 2016.

   [I-D.ietf-bier-ospf-bier-extensions]
              Psenak, P., Kumar, N., Wijnands, I., Dolganow, A.,
              Przygienda, T., Zhang, Z., and S. Aldrin, "OSPF Extensions
              for BIER", draft-ietf-bier-ospf-bier-extensions-02 (work
              in progress), March 2016.

   [I-D.ietf-bier-use-cases]
              Kumar, N., Asati, R., Chen, M., Xu, X., Dolganow, A.,
              Przygienda, T., arkadiy.gulko@thomsonreuters.com, a.,
              Robinson, D., Arya, V., and C. Bestler, "BIER Use Cases
              draft-ietf-bier-use-cases-03 .txt", draft-ietf-bier-use-
              cases-03 (work in progress), July 2016.

Authors' Addresses

   Cui(Linda) Wang
   ZTE Corporation
   No.50 Software Avenue, Yuhuatai District
   Nanjing
   China


   Email: wang.cui1@zte.com.cn


   Zheng(Sandy) Zhang
   ZTE Corporation
   No.50 Software Avenue, Yuhuatai District
   Nanjing
   China


   Email: zhang.zheng@zte.com.cn