INTERNET-DRAFT Intended Status: Informational Expires: September 27, 2015 T. Herbert Facebook L. Yong Huawei O. Zia Microsoft March 26, 2015

Encapsulation Considerations for GUE draft-herbert-gue-encap-considerations-01

Abstract

This document provides a description of how Generic UDP Encapsulation addresses the encapsulation considerations that are described in the "Encapsulation Considerations" Internet Draft.

Status of this Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

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) 2015 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents

Herbert, Yong, Zia Expires September, 2015

(<u>http://trustee.ietf.org/license-info</u>) 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

$\underline{1}$ Introduction	. <u>3</u>
<u>2</u> Entropy	. <u>3</u>
<u>3</u> Next-protocol indication	. <u>4</u>
<u>4</u> MTU and Fragmentation	. <u>4</u>
<u>5</u> OAM	. <u>5</u>
<u>5.1</u> Active OAM	. <u>5</u>
<u>5.2</u> Passive OAM	. <u>5</u>
<u>6</u> Security Considerations	. <u>5</u>
<u>6.1</u> Integrity and authentication of the encapsulation	. <u>5</u>
6.2 Packet level security	. <u>6</u>
<u>7</u> QoS	. <u>6</u>
<u>8</u> Congestion Considerations	. <u>6</u>
9 Header Protection	. <u>7</u>
<u>10</u> Extensibility Considerations	. <u>7</u>
<u>11</u> Layering Considerations	. <u>8</u>
<u>12</u> Service Model	. <u>8</u>
<u>13</u> Hardware Friendly	. <u>8</u>
<u>13.1</u> Switch friendliness	. <u>8</u>
<u>13.2</u> Host friendliness	. <u>9</u>
<u>14</u> Middlebox Considerations	. <u>10</u>
<u>15</u> Network virtualization	. <u>10</u>
<u>16</u> References	. <u>12</u>
<u>16.1</u> Normative References	. <u>12</u>
<u>16.2</u> Informative References	. <u>12</u>
Authors' Addresses	. <u>13</u>

[Page 2]

1 Introduction

This document provides a description of how Generic UDP Encapsulation (GUE) [I.D.herbert-gue] addresses the encapsulation considerations that are described in [I.D.rtg-dt-encap]. That draft is the result of a design team that was chartered by the routing area director to investigate and report on the common issues across different encapsulations.

The organization of this document follows that of the encapsulation considerations draft. There is a section for each of the main areas for consideration with encapsulation. These areas are:

- o Entropy
- o Next-protocol indication
- o MTU and fragmentation
- o OAM
- o Security considerations
- o QoS
- o Congestion considerations
- o Header protection
- o Extensibility considerations
- o Layering considerations
- o Service model
- o Hardware friendly
- o Middleboxes

An additional section covers considerations specific to network virtualization.

2 Entropy

Similar to other UDP encapsulation proposals, the UDP source port of a GUE packet may be set to a value that reflects the inner flow. In GUE parlance, this value based on the inner flow identifier which is often a hash over the 5-tuple of the inner packet's headers. The UDP

Herbert, Yong, Zia Expires September, 2015 [Page 3]

source port provides fourteen bits of entropy assuming that the selected value is restricted to the ephemeral port range. The source port used to indicate a given flow may also change over the lifetime of a flow.

The outer IPv6 flow label may be used to provide additional entropy in the flow identifier when fourteen bits is insufficient.

If the UDP port (and possibly IPv6 flow label) still does not provide enough entropy in the flow classification, then deep parsing of the GUE payload may be performed. See <u>section 14</u>.

<u>3</u> Next-protocol indication

The next protocol indication in GUE is in the proto/ctype field in the GUE header. For GUE data messages (as opposed to control messages, see <u>section 5</u>) the proto field holds an IP protocol number of the next header. An eight bit value is an efficient use of header space, and the lookup of IP protocol can be implemented as a simple 256 entry array (as in Linux).

The IP protocol number allows encapsulation of various layer 2 and layer 3 protocols. In particular, the most common protocols for tunnels are likely to be:

o IPv4: number 4

o IPv6: number 41

o Ethernet: via EtherIP with number 97

The encapsulated protocol may also be GRE (number 47) which allows encapsulation of protocols of any Ethertype in GUE with an additional four bytes of header.

Layer 4 protocols may also be encapsulated within GUE (e.g. ESP, UDP, ICMP, etc.). In this case the UDP and GUE encapsulating headers are considered to be inserted between IP and the transport header. In this way, the outer UDP header for GUE and the encapsulated transport header logically header share the same IP header. For an encapsulated TCP or UDP header checksum calculation, the pseudo checksum is based on the outer IP header ignoring the encapsulation headers.

<u>4</u> MTU and Fragmentation

Similar to other encapsulation protocols, it is recommended in GUE that fragmentation over a tunnel is avoided by configuring tunnel MTUs and using Path MTU Discovery (PMTUD) as necessary. As described

[Page 4]

in [<u>I.D.rtg-dt-encap</u>], detecting mis-configuration causing packets to be dropped due to MTU issues is desirable; having the encapsulator set the don't-fragment (DF) flag in the outer IPv4 header and logging any received ICMP "packet too big" (PTB) are compatible with GUE.

As discussed in [<u>RFC4459</u>] it may not be possible to avoid the need for fragmentation in a all circumstances (for instance a link in a tunnels path may have the minimum MTU for IPv6 which is 1280). To accommodate this, a fragmentation option has been defined for GUE ([<u>I.D.gue-fragmentation</u>]).

5 OAM

Specific OAM support for GUE (and other encapsulation protocols) has not yet been defined. Due to the extensibility model of GUE and definition of control messages, there is a lot of flexibility in how OAM may be supported. GUE includes provisions to support both active OAM (OAM specific messages) and passive OAM (measurements of data messages).

5.1 Active OAM

The GUE header includes a bit that indicates that the payload contains a control message as opposed to a data message. When this bit is set, the proto/ctype field is interpreted to be a control type. Various types of control messages may be defined, including those for OAM.

5.2 Passive OAM

Options in the GUE header may be added to permit passive OAM measurements attached to data messages. This may accomplished by using be single bits of information, or by OAM measurement fields which could contain items such as sequence numbers, timestamps, etc.

<u>6</u> Security Considerations

Security is a very important consideration in GUE. This is particularly motivated by the multi-tentant use case of network virtualization where isolation between tenants is a critical requirement. The GUE security model includes both considerations to protect the headers and the whole packet. For both of these, we assume that a "pluggable" security model is desirable with the assumption that stronger security may be implemented over time in response to changing threats.

6.1 Integrity and authentication of the encapsulation

[Page 5]

Addresses, port numbers, and various elements of a GUE header may need fairly strong assurances of integrity and authentication to protect against corruption or spoofing. This requirement is readily apparent in the use case of network virtualization where ensuring the integrity and authenticity of a virtual network identifier is paramount to guaranteeing isolation between virtual networks even in the presence of users with malicious intent.

To provide for this security, an optional security field is defined in GUE ([I.D.hy-gue-4-secure-transport]). This field has three possible sizes of 64, 128, or 256 bits to allow for different levels of security. The simplest mechanism is a security cookie which is a shared value that is passed in the clear and must matched on receipt. More sophisticated mechanisms may use cryptographic hashes, nonce values, reply detection, etc.

6.2 Packet level security

As GUE is contained in an IP packet, the packet itself may be encapsulated in something like ESP or DTLS to provide security. This is straightforward, however visibility of the encapsulation is lost in the network. This is problematic, for instance, if one wanted to establish firewalls to restrict packets for a certain virtual network.

Security of a GUE payload may be accomplished by applying ESP or DTLS to the payload and encapsulating ESP within GUE. In this model the protocol stack may be something like IP|UDP|GUE|ESP|IP. The GUE next protocol would indicate ESP (number 50), and the UDP and GUE headers would be sent in the clear so that encapsulation is visible to the network. As described above, measures should be taken to ensure the integrity and authentication of addresses and GUE headers. One salient property of this method is that any bits created by an application or virtualization guest are covered by the packet level security mechanism.

7 QoS

There are no specific provisions or options to provide additional QoS facilities in GUE. The provisions of "Diffserv and Tunnels" [<u>RFC2983</u>] are assumed.

<u>8</u> Congestion Considerations

Congestions considerations that are generically specified for tunnels would be applicable to GUE. This would include mechanisms currently being defined such as circuit breaker [I.D.cirtuit-breaker] and common ECN handling for IP tunnels.

[Page 6]

In certain cases, specific congestion control may be necessitated beyond what generic congestion control mechanisms may provide. In particular, this may be required in a data center whose native traffic and resources have been tuned to a very specific congestion control algorithm. When third party network stacks, such as those running in a VM's guest OS, are introduced into such an environment their congestion control model may substantially conflict with that of the native traffic. If traffic and resource isolation is not feasible, the only recourse may be to force third party traffic into compliance with the native congestion control.

This can potentially be accommodated in GUE in two ways:

- Add an additional protocol layer to the encapsulation that provides congestion control. DCCP is a candidate. In this case, the protocol stack for a encapsulated packet may look like IP|UDP|GUE|DCCP|IP.
- 2) Add an option to GUE which provides for congestion control. This would likely include an optional header field that would contain various values needed for congestion control-- sequence numbers, timestamps, ack numbers etc. Note that some of this information may also be in common with passive OAM data.

Extending GUE with finer grained congestion is a topic for further exploration.

9 Header Protection

As with other UDP encapsulation protocols, the UDP checksum may or may not be set transmit. The requirements for setting a zero UDP checksum with IPv6 to be compliant with [RFC6935] and [RFC6936] are enumerated in [I.D.herbert-gue]. In the case that a zero checksum is used (either for IPv4 or IPv6) the GUE specification recommends that the GUE header checksum be used (unless stronger protection such as security are present).

The GUE header checksum is a UDP-lite like option in the GUE header ([I.D.herbert-guecsum]). This checksum covers the entire GUE header and a pseudo header containing the outer IP addresses and UDP port numbers. Optionally, the checksum may cover all or part of the encapsulated GUE payload.

10 Extensibility Considerations

GUE is an extensible protocol that allows a variable length header. The extensibility mechanism in GUE is flag-fields. This is similar to the use of flags and fields in GRE, where if a flag is set a field of

[Page 7]

a specific size is present. See <u>section 13.1</u> for points about the efficiency of the flag-fields solution.

The GUE header includes fifteen bit flags in the primary header and an additional thirty-two in an extension flags field. Up to 124 bytes of optional fields may be present.

All flags are considered mandatory, in the sense that a decapsulator must drop a packet it receives with set flags that are unknown to it. This requirement ensures that new flags with non-trivial semantics can be added without breaking compatibility. A middlebox may ignore flags (see <u>section 14</u>).

<u>11</u> Layering Considerations

GUE does not include any specific provisions for layering of encapsulations other than the fact that it can encapsulate an encapsulated packet represented by an IP protocol. As described in <u>section 3</u>, encapsulation of GRE may be used to encapsulate packets of an arbitrary Ethertype.

Conceptually, the GUE header could be disassociated from UDP and defined as its own IP protocol (similar to GRE being an IP protocol). In this manner GUE could effectively function as an IP extension header and layered encapsulations would essentially be equivalent to multiple extension headers.

12 Service Model

The base service model of GUE is equivalent to that of IP. Packets can be lost, reordered, duplicated, etc. To enact a more elaborate service model over GUE, such as pseudo-wire semantics or reliable tunnels, could be done as a layered encapsulation of a protocol that provides the service.

The exception to the above occurs when encapsulated traffic has the ability to negatively impact unrelated networking traffic. In this case, a service model that provides congestion control or DDOS protection is a candidate to implement within the encapsulation layer (e.g. see section 8).

<u>13</u> Hardware Friendly

<u>13.1</u> Switch friendliness

GUE is mostly intended to be an end to end tunneling protocol. Switches may inspect fields as input to routing operation, however they should not modify GUE headers in flight (checksum and header

[Page 8]

security likely make that infeasible). Encapsulation for the purposes affecting routing per hop or targeting networking services are better left to protocols dedicated to those functions such as BIER or SFC.

In the case that a switch acts as tunnel endpoint (i.e. an NVE in NVO3 parlance), considerations can be made for efficient termination and processing of UDP.

GUE has the following features to be friendly in switches:

- o For a given set of flags, field offsets are fixed
- o The number of possible flag combinations is 2^N for N supported flags. This is a much smaller number than for TLVs which are combinatorial.
- o Minimal overhead. Other than flags, there is no additional overhead associated with fields
- o Flag fields are amenable to hardware parallel parsing mechanisms such as TCAM (based on the above points)
- o The GUE header checksum obviates need for full packet checksum support
- o Hardware support of variable flag-fields has already be demonstrated in GRE

13.2 Host friendliness

The first requirement of encapsulation in the host is that it works with existing NIC offloads. The five common offloads in question are RSS (Receive Side Scaling), TX-csum (transmit checksum offload), RXcsum (receive checksum offload), LSO (Large Segment Offload), and LRO (Large Receive Offload).

RSS is already solved by enabling RSS for UDP which is available on most NICs. This works for any UDP encapsulation that uses source port for flow entropy.

TX-csum and RX-csum offload for encapsulated checksums (no outer UDP checksum) are already generically supported by NICs that provide NETIF_HW_CSUM and CHECKSUM_UNNECESSARY (in Linux parlance). The outer UDP checksum may also be enabled, and this can be leveraged in GUE to provide checksum offload of inner transport checksums for legacy devices (via checksum-unnecessary conversion and remote checksum offload). The ability for NICs to support offload of multiple checksums in a packet may also become pertinent in time.

[Page 9]

For LSO (TSO), the Linux stack already demonstrates a generic solution for segmentation with UDP encapsulation in GSO. This can be implemented in HW as method to provide LSO for any UDP encapsulation.

For LRO, a device needs to do deep parsing of the GUE payload. This can be accomplished by skipping over any options using the Hlen field. Packets should only be considered to match if the GUE encapsulation, including any optional fields and private data, are identical.

The implementation of GUE in a software stack is fairly straight forward and can be efficient. The UDP layer in the software stack should already handle the processing of the UDP/IP headers including that of the checksum. For the GUE headers, processing the flag-fields is likely the most difficult operation. Given the practical constraints on flag-fields, they can be processed without a loop and without the need to check lengths or duplicate fields. The check for unsupported set flags can be implemented with a simple masked comparison on the flags.

<u>14</u> Middlebox Considerations

The following GUE features facilitate middlebox handling:

- o The Hlen field allows middlebox to skip over optional fields to perform deep parsing
- o The meaning of proto/ctype field is invariant regardless of flags
- o Flag-fields permit random access for inspection
- o Middleboxes are not required to understand all possible fields. A principle in GUE is that new fields cannot cause reinterpretation of old fields.

15 Network virtualization

The primary requirement for network virtualization is that a virtual network is indicated as part of the encapsulation (i.e. a virtual network identifier or VNI-ID).

GUE defines a 32 bit VNI-ID in an optional field ([I.D.hy-nvo3-gue-4nvo]). There is no predefined structure to this value. An implemention may apply a hierarchical structure (for instance a tenant might have virtual sub-networks), as well as allocating bits to indicate class or other attributes (such as a bit indicating a trusted or untrusted virtual network).

Herbert, Yong, ZiaExpires September, 2015[Page 10]

A first order requirement of network virtualization is that isolation between virtual networks be ensured. As described in <u>section 6</u>, the GUE security option should be used to provide integrity and authentication.

16 References

<u>16.1</u> Normative References

[I.D.rtg-dt-encap] Nordmark, E., Tian, A., Gross, J.," Hudson, J., Kreeger, L., Garg, P., Thaler, P., Herbert, T., "Encapsulation Considerations", <u>draft-rtg-dt-encap-01</u>.

[I.D.herbert-gue] Herbert, T., and Yong, L., "Generic UNP Encapsulation", <u>draft-herbert-gue-03</u>, work in progress.

<u>16.2</u> Informative References

[RFC4459] Savola, P., "MTU and Fragmentation Issues with In-the-Network Tunneling", <u>RFC 4459</u>, April 2006, <<u>http://www.rfc-editor.org/info/rfc4459</u>>.

[RFC2983] Black, D., "Differentiated Services and Tunnels", <u>RFC 2983</u>, October 2000, <<u>http://www.rfc-</u> editor.org/info/rfc2983>.

[RFC6935] Eubanks, M., Chimento, P., and M. Westerlund, "IPv6 and UDP Checksums for Tunneled Packets", <u>RFC 6935</u>, April 2013, <<u>http://www.rfc-editor.org/info/rfc6935</u>>.

[RFC6936] Fairhurst, G. and M. Westerlund, "Applicability Statement for the Use of IPv6 UDP Datagrams with Zero Checksums", <u>RFC 6936</u>, April 2013, <<u>http://www.rfc-</u> editor.org/info/rfc6936>.

[I.D.gue-fragmentation] Herbert, T. and Templin, F., "Fragmentation option for Generic UDP Encapsulation", <u>draft-herbert-gue-fragmentation-00</u>, work in progress

[I.D.hy-gue-4-secure-transport] Yong, L., Herbert, T., "Generic UDP Encapsulation (GUE) for Secure Transport", <u>draft-hy-gue-4-secure-transport-00</u>, work in progress.

[I.D.tsvwg-circuit-breaker] Fairhurst, G., "Network Transport Circuit Breakers", draft-ietf-tsvwg-circuitbreaker-00

[I.D.herbert-remotecsumoffload] Herbert, T., "Remote checksum offload for encapsulation", <u>draft-herbert-</u> remotecsumoffload-00, work in progress

[I.D.hy-nvo3-gue-4-nvo] Yong, L., Herbert, T., "Generic UDP Encapsulation (GUE) for Network Virtualization

Herbert, Yong, ZiaExpires September, 2015[Page 12]

Overlay", work in progress

Authors' Addresses

Tom Herbert Facebook Menlo Park, CA USA EMail: tom@herbertland.com Lucy Yong Huawei USA 5340 Legacy Dr. Plano, TX 75024 USA Osama Zia Microsoft

EMail: osamaz@microsoft.com

Herbert, Yong, Zia Expires September, 2015 [Page 13]