Sally Floyd David Black K. K. Ramakrishnan December 1999 Expires: June 2000

# **IPsec Interactions with ECN**

Status of this Memo

This document is an Internet-Draft and is in full conformance with all provisions of <u>Section 10 of RFC2026</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 <a href="http://www.ietf.org/ietf/lid-abstracts.txt">http://www.ietf.org/ietf/lid-abstracts.txt</a>

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

## Abstract

IPsec supports secure communication over potentially insecure network components such as intermediate routers. IPsec protocols support two operating modes, transport mode and tunnel mode. Explicit Congestion Notification (ECN) is an experimental addition to the IP architecture that provides notification of onset of congestion to delay- or loss-sensitive applications. ECN provides congestion notifications to enable adaptation to network conditions without the impact of dropped packets [RFC 2481]. The use of two bits in the IP header for ECN experimentation conflicts with header processing at IPsec tunnel endpoints in a manner that makes ECN unusable in the presence of IPsec tunnels. This document considers issues related to this conflict, describes two alternative solutions, and updates the IPsec architecture [RFC 2401] to include these alternatives. Support for

one or the other of these alternatives is REQUIRED to remove the underlying conflict.

# **<u>1</u>**. Introduction.

IPsec supports secure communication over potentially insecure network components such as intermediate routers. IPsec protocols support two operating modes, transport mode and tunnel mode, that span a wide range of security requirements and operating environments. Transport mode security protocol header(s) are inserted between the IP (IPv4 or IPv6) header and higher layer protocol headers (e.g., TCP), and hence transport mode can only be used for end-to-end security on a connection. IPsec tunnel mode is based on adding a new "outer" IP header that encapsulates the original, or "inner" IP header and its associated packet. Tunnel mode security headers are inserted between these two IP headers. In contrast to transport mode, the new "outer" IP header and tunnel mode security headers can be added and removed at intermediate points along a connection, enabling security gateways to secure vulnerable portions of a connection without requiring endpoint participation in the security protocols. An important aspect of tunnel mode security is that the outer header is discarded at tunnel egress, ensuring that security threats based on modifying the IP header do not propagate beyond that tunnel endpoint. Further discussion of IPsec can be found in [RFC 2401].

Explicit Congestion Notification (ECN) is an experimental addition to the IP architecture that provides congestion notifications to delayor loss-sensitive applications to enable them to adapt to network conditions without the impact of dropped packets [RFC 2481]. An ECNcapable router uses the ECN mechanism to signal congestion to connection endpoints by setting a bit in the IP header. These endpoints then react, in terms of congestion control, as if a packet had been dropped (e.g., TCP halves its congestion window). This ability to avoid dropping packets in response to congestion is supported by the use of active queue management mechanisms (e.g., RED) in routers; such mechanisms begin to mark or drop packets as a consequence of congestion before a congested router queue is completely full. ECN is an experimental optimization -- not all routers may be expected to support ECN, and even ECN-capable routers drop packets from ECN-capable connections when necessary. The advantage to routers of not dropping such packets is that ECN can provide a more timely reaction to congestion than reactions based on drop detection via duplicate ACKs or timeout.

ECN as currently specified uses two bits within the IP header in a manner that conflicts with current header processing at IPsec tunnel endpoints. Use of ECN over an IPsec tunnel results in routers

[Page 2]

December 1999

attempting to use the outer IP header to signal congestion to endpoints, but discarding of the outer header at tunnel egress also discards those indications of congestion. <u>RFC 2481</u> recommended that ECN not be used with IPsec tunnels in order to avoid this behavior and its undesirable consequences. This document updates the IPsec architecture to remove that conflict.

In principle, permitting the use of ECN functionality in the outer header of an IPsec tunnel raises security concerns because an adversary could tamper with the information that propagates beyond the tunnel endpoint. Based on an analysis (included in this document) of these concerns and the associated risks, our overall approach is to provide configuration support for the IPsec changes that remove the conflict with ECN. This makes permission to use ECN functionality in the outer header of an IPsec tunnel a configurable part of the corresponding IPsec Security Association (SA), so that it can be disabled in situations where the risks are judged to outweigh the benefits. The result is that an IPsec security administrator is presented with two alternatives for the behavior of ECN-capable connections within an IPsec tunnel:

- A limited-functionality alternative in which the ECN bits are preserved in the inner header, but ECN functionality is disabled in the outer header. The only mechanism available for signaling congestion occurring within the tunnel in this case is dropped packets.

- A full functionality alternative that supports ECN in both the inner and outer headers. This alternative propagates ECN congestion notifications from nodes within the tunnel to endpoints outside the tunnel.

Support for these alternatives involves changes to IP header processing at tunnel ingress and egress. All IPsec implementations MUST implement one of the above two alternatives in order to eliminate the current incompatibility between ECN and IPsec tunnels, but implementers MAY choose to implement either alternative.

The main goal of this document is to provide guidance about the tradeoffs between the limited-functionality and full-functionality alternatives. This includes a full discussion of the potential effects of an adversary's modification to the two bits used by ECN in the IP header.

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 <u>RFC 2119</u> [<u>RFC2119</u>].

[Page 3]

# 2. Architecture.

ECN as specified for experimental purposes uses two bits in the IP header (ECT - ECN Capable Transport, and CE - Congestion Experienced) for signaling between routers and connection endpoints, and uses two flags in the TCP header (ECN-Echo - Echo ECN bit in IP header, CWR - Congestion Window Reduced) for TCP-endpoint to TCP-endpoint signaling. For a TCP connection, a typical sequence of events in an ECN-based reaction to congestion is as follows:

- The ECT bit is set in packets transmitted by the sender to indicate that this TCP connection reacts to ECN congestion notifications for these packets.

- An ECN-capable router detects impending congestion and notices that the ECT bit is set in the packet that the router is about to drop. Instead of dropping the packet, the router sets the CE bit and forwards the packet.

- The packet with the CE bit set arrives at the receiver. The receiver sets the ECN-Echo flag in its next TCP ACK to the sender. - The sender receives the TCP ACK with ECN-Echo set, and reacts to the congestion as if a packet had been dropped.

- The sender sets the CWR flag in the TCP header of the next packet sent to the receiver to acknowledge its receipt of and reaction to the ECN-Echo flag.

Further details on ECN functionality including negotiation of ECNcapability as part of connection setup as well as the responsibilities and requirements of ECN-capable routers and transports can be found in [<u>RFC2481</u>]. These requirements apply only to routers and transports participating in ECN experimentation.

ECN interacts with IPsec tunnels because the bits it uses in the IP header are part of what IPsec refers to as the IPv4 TOS octet or IPv6 Traffic Class octet; this field is copied or mapped from the inner IP header to the outer IP header at IPsec tunnel ingress, and the outer header's copy of this field is discarded at IPsec tunnel egress [RFC2401]. If an ECN-capable router were to set the CE (Congestion Experienced) bit in an IPsec-tunneled packet, this would be discarded at tunnel egress, losing the notification of congestion. As a consequence of this behavior, use of ECN over IPsec tunnels is currently not recommended [RFC 2481].

The IPsec limited-functionality alternative for ECN encapsulation is to always clear (i.e., set to 0) the ECT bit in the outer (encapsulating) header, regardless of the value of the ECT bit in the inner (encapsulated) header. Under this alternative, the ECN bits in the inner header are not altered upon decapsulation. The disadvantage of this approach is that ECN-capable flows do not have

[Page 4]

ECN support for that part of the path that uses IPsec tunneling. That is, if the encapsulated packet arrives at a congested router that is ECN-capable, and the router decides to drop or mark the packet as an indication of congestion to the end nodes, the router has no alternative but to drop the packet.

The IPsec full-functionality alternative for ECN encapsulation copies the ECT bit of the inside header to the outside header on encapsulation, and performs an OR of the CE bits from the outer and inner headers to determine the value of the CE bit on decapsulation. Under the full-functionality alternative, an ECN-capable flow can take advantage of ECN for those parts of the path that use IPsec tunneling. The disadvantage of the full-functionality alternative is that IPsec cannot protect flows from certain modifications to the ECN bits in the IP header within the tunnel. The potential dangers from modifications to the ECN bits in the IP header are described in detail in <u>Section 4</u> below.

This document describes the changes to IPsec that are REQUIRED to enable ECN experimentation over IPsec tunnels without discarding congestion notifications when ECN-capable router or routers are traversed by an IPsec tunnel carrying ECN-capable connections. In summary, two changes to IPsec functionality are involved:

(1) Modify the handling of the IPv4 TOS octet and IPv6 Traffic Class octet at IPsec tunnel endpoints to prevent loss of ECN congestion notifications when an IPsec tunnel traverses an ECN-capable router.

(2) Enable the endpoints of an IPsec tunnel to negotiate enabling ECN functionality in the outer headers of that tunnel based on security policy. ECN is only used in the outer header of packets from ECN-capable connections.

The minimum effort to make ECN compatible with IPsec tunnels is a simplified version of the first change that prevents ECN from being enabled in the outer header of an IPsec tunnel. In contrast, full support for ECN includes the ability to negotiate ECN usage between tunnel endpoints; this enables a security administrator to disable ECN in situations where she believes the risks (e.g., of lost congestion notifications) outweigh the benefits of ECN.

#### **<u>3</u>**. IPsec Changes for ECN usage

This section describes the detailed changes to enable usage of ECN over IPsec tunnels, including the negotiation of ECN support between tunnel endpoints. In order to avoid the loss of congestion notifications at tunnel egress, full ECN functionality for an IPsec

[Page 5]

tunnel supports agreement between both ends of the tunnel that ECN is being used. This is supported by three changes to IPsec:

- A Security Association Database (SAD) field indicating whether tunnel encapsulation and decapsulation processing allows or forbids ECN usage in the outer IP header.

- A new Security Association Attribute that enables negotiation of this SAD field between the two endpoints of an SA that supports tunnel mode.

- Changes to tunnel mode encapsulation and decapsulation processing to allow or forbid ECN usage in the outer IP header based on the value of the SAD field. When ECN usage is allowed in the outer IP header, ECT is set in the outer header for ECNcapable connections and congestion notifications (indicated by the CE bit) from such connections are propagated to the inner header at tunnel egress.

These changes are covered further in the following three subsections.

The first two changes are OPTIONAL, but if negotiation of ECN usage is implemented, then the SAD field SHOULD also be implemented. On the other hand, negotiation of ECN usage is OPTIONAL in all cases, even for implementations that support the SAD field. The encapsulation and decapsulation processing changes are REQUIRED, but MAY be implemented without the other two changes by assuming that ECN usage is always forbidden. The full-functionality alternative for ECN usage over IPsec tunnels consists of the SAD field and the full version of encapsulation and decapsulation processing changes, with or without the OPTIONAL negotiation support. The limitedfunctionality alternative consists of a subset of the encapsulation and decapsulation changes that always forbids ECN usage.

#### **3.1**. ECN Tunnel Security Association Database Field

Full ECN functionality adds a new field to the SAD (see [RFC2401]):

ECN Tunnel: allowed or forbidden.

Indicates whether ECN-capable connections using this SA in tunnel mode are permitted to receive ECN congestion notifications for congestion occurring within the tunnel. The allowed value enables ECN congestion notifications. The forbidden value disables such notifications, causing all congestion to be indicated via dropped packets.

[OPTIONAL. The value of this field SHOULD be assumed to be "forbidden" in implementations that do not support it.]

If this attribute is implemented, then the SA specification in a

[Page 6]

Security Policy Database (SPD) entry MUST support a corresponding attribute, and this SPD attribute MUST be covered by the SPD administrative interface (currently described in <u>Section 4.4.1 of [RFC2401]</u>).

## 3.2. ECN Tunnel Security Association Attribute

A new IPsec Security Association Attribute is defined to enable the support for ECN congestion notifications based on the outer IP header to be negotiated for IPsec tunnels (see [RFC2407]). This attribute is OPTIONAL, although implementations that support it SHOULD also support the SAD field defined in <u>Section 3.1</u>.

Attribute Type

 sic
\$

Class Values

ECN Tunnel

Specifies whether ECN experimental functionality is allowed to be used with Tunnel Encapsulation Mode. This affects tunnel encapsulation and decapsulation processing see Section 3.3.

RESERVED	0
Allowed	1
Forbidden	2

Values 3-61439 are reserved to IANA. Values 61440-65535 are for private use.

If unspecified, the default shall be assumed to be Forbidden.

ECN Tunnel is a new SA attribute, and hence initiators that use it can expect to encounter responders that do not understand it, and therefore reject proposals containing it. For backwards compatibility with such implementations initiators SHOULD always also include a proposal without the ECN Tunnel attribute to enable such a responder to select a transform or proposal that does not contain the ECN Tunnel attribute. <u>RFC 2407</u> currently requires responders to reject all proposals if any proposal contains an unknown attribute; this requirement is expected to be changed to require a responder not to select proposals or transforms containing unknown attributes.

[Page 7]

# 3.3. Changes to IPsec Tunnel Header Processing

Subsequent to the publication of [<u>RFC 2401</u>], the TOS octet of IPv4 and the Traffic Class octet of IPv6 have been superseded by the sixbit DS Field [RFC2474, RFC TBD] and a two-bit "currently unused" (CU) field [RFC TBD]. The two bits in the IP header used for ECN experimentation, ECT and CE, occupy bits 0 and 1 of the CU field.

For full ECN support, the encapsulation and decapsulation processing for the IPv4 TOS field and the IPv6 Traffic Class field are changed from that specified in [RFC2401] to the following:

	< How Outer Hdr Relates to	Inner Hdr>
	Outer Hdr at	Inner Hdr at
IPv4	Encapsulator	Decapsulator
Header fields:		
DS Field	copied from inner hdr (5)	no change
CU Field	constructed (7)	constructed (8)

IPv6

```
Header fields:
DS Field
CU Field
```

ieldcopied from inner hdr (6)no changefieldconstructed (7)constructed (8)

(5)(6) If the packet will immediately enter a domain for which the DSCP value in the outer header is not appropriate, that value MUST be mapped to an appropriate value for the domain [RFC 2474]. Also see [RFC 2475] for further information.

(7) If the value of the ECN Tunnel field in the SAD entry for this SA is "allowed" and the value of ECT (bit 0) is 1 in the inner header, set ECT to 1 in the outer header, else set ECT to 0 in the outer header. Set CE (bit 1) to 0 in the outer header.

(8) If the value of the ECN tunnel field in the SAD entry for this SA is "allowed" and the value of ECT (bit 0) in the inner header is 1, then set the CE bit (bit 1) in the inner header to the logical OR of the CE bit in the inner header with the CE bit in the outer header, else make no change to the CU field.

(5) and (6) are identical to match usage in [<u>RFC2401</u>], although they are different in [<u>RFC2401</u>]. The Differentiated Services Working Group is currently considering interactions between Differentiated Services and tunnels, so implementers are advised to check for additional RFCs that further update the IPsec architecture in this area.

The above description applies to implementations that support the ECN

[Page 8]

Tunnel field in the SAD; such implementations MUST implement this processing of the DS field instead of the processing of the IPv4 TOS octet and IPv6 Traffic Class octet defined in [<u>RFC2401</u>]. This constitutes the full-functionality alternative for ECN usage with IPsec tunnels.

An implementation that does not support the ECN Tunnel field in the SAD MUST implement processing of the DS Field by assuming that the value of the ECN Tunnel field of the SAD is "forbidden" for every SA. In this case, the processing of the CU field reduces to:

- (7) Set the CU field to zero in the outer header.
- (8) Make no change to the CU field.

This constitutes the limited functionality alternative for ECN usage with IPsec tunnels.

In addition, for backwards compatibility, packets with ECT and CE both set to 1 in the outer header SHOULD be dropped if they arrive on an SA that forbids or is assumed to forbid ECN usage in tunnel mode. This applies to both the complete ECN support and partial ECN support implementation approaches. This is discussed further in <u>Section 6</u>.

## **<u>4</u>**. Possible Changes to the ECN Field

This section considers the issues when a router is operating, possibly maliciously, to modify either of the ECN bits in IP header. In this section we represent the ECN bits in the IP header by the tuple (ECT bit, CE bit). The ECT bit, when set to 1, indicates an ECN-Capable Transport. The CE bit, when set to 1, indicates that Congestion was Experienced in the path.

By tampering with the ECN bits, an adversary (or a broken router) could do one or more of the following: erase the ECN congestion indication, falsely report congestion, disable ECN-Capability for an individual packet, or falsely indicate ECN-Capability. We systematically examine the various cases by which the ECN bits could be modified. The important criterion we consider in determining the consequences of such modifications is whether it is likely to lead to worse behavior in any dimension (throughput, delay, fairness or functionality) than if a router were to drop a packet.

## **<u>4.1</u>**. Erasing the Congestion Indication

First, we consider the changes that a router could make that would result in effectively erasing the congestion indication after it had been set by a router upstream. The convention followed is: (ECT, CE) of received packet -> (ECT, CE) of packet transmitted.

[Page 9]

December 1999

 $(1, 1) \rightarrow (1, 0)$ : erase only the CE bit that was set.  $(1, 1) \rightarrow (0, 0)$ : erase both the ECT bit and the CE bit.  $(1, 1) \rightarrow (0, 1)$ : erase the ECT bit

The first change turns off the CE bit after it has been set by some upstream router along the path. The consequence for the upstream router is that there is a potential for congestion to build for a time, because the congestion indication does not reach the source. However, the packet would be received and acknowledged.

The potential effect of erasing the congestion indication is complex, and is discussed in depth in <u>Section 5</u> below. Note that the effect of erasing the congestion indication is different from dropping a packet in the network. When a data packet is dropped, the drop is detected by the TCP sender, and interpreted as an indication of congestion. Similarly, if a sufficient number of consecutive acknowledgement packets are dropped, causing the cumulative acknowledgement field not to be advanced at the sender, the sender is limited by the congestion window from sending additional packets, and ultimately the retransmit timer expires.

In contrast, a systematic erasure of the CE bit by a downstream router can have the effect of causing a queue buildup at an upstream router, including the possible loss of packets due to buffer overflow. There is a potential of unfairness in that another flow that goes through the congested router could react to the CE bit set while the flow that has the CE bit erased could see better performance. The limitations on this potential unfairness are discussed in more detail in <u>Section 5</u> below.

The second change is to turn off both the ECT and the CE bits, thus erasing the congestion indication and disabling ECN-Capability at the same time. The third change turns off only the ECT bit, disabling ECN-Capability. The proposal in this Internet Draft is for the receiver at the end of a tunnel to copy the CE bit, if set, from the outer header to the inner header during decapsulation, if the ECT bit in the inner header is set and the tunnel provides full ECN support. In this case, the third change within an IPsec tunnel would not erase the congestion indication, but would only disable ECN-Capability for that packet within the rest of the tunnel. However, when performed outside of an IPsec tunnel, the third change would also effectively erase the congestion indication, because an ECN field of (0, 1) is undefined.

The `erasure' of the congestion indication is only effective if the packet does not end up being marked or dropped again by a downstream router. With the first change, the packet remains ECN-Capable, and could be either marked or dropped by a downstream router as an

[Page 10]

indication of congestion. With the second and third changes, the packet is no longer ECN-capable, and can therefore be dropped but not marked by a downstream router as an indication of congestion.

## 4.2. Falsely Reporting Congestion

 $(1, 0) \rightarrow (1, 1)$ 

This change is to set the CE bit when the ECT bit was already set, even though there was no congestion. This change does not affect the treatment of that packet along the rest of the path. In particular, a router does not examine the CE bit in deciding whether to drop or mark an arriving packet.

However, this could result in the application unnecessarily invoking end-to-end congestion control, and reducing its arrival rate. By itself, this is no worse (for the application or for the network) than if the tampering router had actually dropped the packet.

# 4.3. Disabling ECN-Capability

$$(1, 0) \rightarrow (0, *)$$

This change is to turn off the ECT bit of a packet that does not have the CE bit set. (Section 4.1 discussed the case of turning off the ECT bit of a packet that does have the CE bit set.) This means that if the packet later encounters congestion (e.g., by arriving to a RED queue with a moderate average queue size), it will be dropped instead of being marked. By itself, this is no worse (for the application) than if the tampering router had actually dropped the packet. The saving grace in this particular case is that there is no congested router upstream expecting a reaction from setting the CE bit.

## **4.4**. Falsely Indicating ECN-Capability

This change is to incorrectly label a packet as ECN-Capable.

 $(0, *) \rightarrow (1, 0);$  $(0, *) \rightarrow (1, 1);$ 

If the packet later encounters moderate congestion at an ECN-Capable router, the router could set the CE bit instead of dropping the packet. If the transport protocol in fact is not ECN-Capable, then the transport will never receive this indication of congestion, and will not reduce its sending rate in response. The potential consequences of falsely indicating ECN-capability are discussed further in <u>Section 5</u> below.

[Page 11]

If the packet never later encounters congestion at an ECN-Capable router, then the first of these two changes would have no effect. The second change, however, would have the effect of giving false reports of congestion to a monitoring device along the path. If the transport protocol is ECN-Capable, then the second of these two changes (when, for example, (0,0) was changed to (1,1)) could also have an effect at the transport level, by combining falsely indicating ECN-Capability with falsely reporting congestion. For an ECN-capable transport, this would cause the transport to unnecessarily react to congestion. In this particular case, the router that is incorrectly changing the ECN field could have dropped the packet. Thus for this case of an ECN-capable transport, the consequence of this change to the ECN field is no worse than dropping the packet.

## 4.5. Changes with No Functional Effect

 $(0, *) \rightarrow (0, *)$ 

The CE bit is ignored in a packet that does not have the ECT bit set. Thus, this change would have no effect, in terms of ECN.

## **<u>4.6</u>**. Information carried in the Transport Header

For TCP, an ECN-capable TCP receiver informs its TCP peer that it is ECN-capable at the TCP level, using information in the TCP header at the time the connection is setup. This document does not consider potential dangers introduced by changes in the transport header because the IPsec tunnel protects the transport header.

## **<u>5</u>**. Implications of Subverting End-to-End Congestion Control

This section focuses on the potential repercussions of subverting end-to-end congestion control by either falsely indicating ECN-Capability, or by erasing the congestion indication in ECN (the CEbit). Subverting end-to-end congestion control by either of these two methods can have consequences both for the application and for the network. We discuss these separately below.

The first method to subvert end-to-end congestion control, falsely indicating ECN-Capability, effectively subverts end-to-end congestion control only if the packet later encounters congestion that results in the setting of the CE bit. In this case, the transport protocol does not receive the indication of congestion from these downstream congested routers.

The second method to subvert end-to-end congestion control, `erasing' the (set) CE bit in a packet, effectively subverts end-to-end

[Page 12]

congestion control only when the CE bit in the packet was set earlier by a congested router. In this case, the transport protocol does not receive the indication of congestion from the upstream congested routers.

Either of these two methods of subverting end-to-end congestion control can potentially introduce more damage to the network (and possibly to the flow itself) than if the adversary had simply dropped packets from that flow. However, as we discuss later in this section and in <u>Section 7</u>, this potential damage is limited.

#### 5.1. Implications for the Network and for Competing Flows

The CE bit of the ECN field is only used by routers as an indication of congestion during periods of \*moderate\* congestion. ECN-capable routers should drop rather than mark packets during heavy congestion even if the router's queue is not yet full. For example, for routers using active queue management based on RED, the router should drop rather than mark packets that arrive while the average queue sizes exceed the RED gueue's maximum threshold.

One consequence for the network of subverting end-to-end congestion control is that flows that do not receive the congestion indications from the network might increase their sending rate until they drive the network into heavier congestion. Then, the congested router could begin to drop rather than mark arriving packets. For flows that are not isolated by some form of per-flow scheduling or other per-flow mechanisms, but that are instead aggregated with other flows in a single queue in an undifferentiated fashion, this packetdropping at the congested router would apply to all flows that share that queue. Thus, the consequences would be to increase the level of congestion in the network.

In some cases, the increase in the level of congestion will lead to a substantial buffer buildup at the congested queue that will be sufficient to drive the congested queue from the packet-marking to the packet-dropping regime. This transition could occur either because of buffer overflow, or because of the active queue management policy described above that drops packets when the average queue is above RED's maximum threshold. At this point, all flows, including the subverted flow, will begin to see packet drops instead of packet marks, and a malicious or broken router will no longer be able to `erase' these indications of congestion in the network. If the end nodes are deploying appropriate end-to-end congestion control, then the subverted flow will reduce its arrival rate in response to congestion. When the level of congestion is sufficiently reduced, the congested queue can return from the packet-dropping regime to the packet-marking regime. The steady-state pattern could be one of the

[Page 13]

congested queue oscillating between these two regimes.

In other cases, the consequences of subverting end-to-end congestion control will not be severe enough to drive the congested link into sufficiently-heavy congestion that packets are dropped instead of being marked. In this case, the implications for competing flows in the network will be a slightly-increased rate of packet marking or dropping, and a corresponding decrease in the bandwidth available to those flows. This can be a stable state if the arrival rate of the subverted flow is sufficiently small, relative to the link bandwidth, that the average queue size at the congested router remains under control. In particular, the subverted flow could have a limited bandwidth demand on the link at this router, while still getting more than its "fair" share of the link. This limited demand could be due to a limited demand from the data source; a limitation from the TCP advertised window; a lower-bandwidth access pipe; or other factors. Thus the subversion of ECN-based congestion control can still lead to unfairness, which we believe is appropriate to note here.

The threat to the network posed by the subversion of ECN-based congestion control in the network is essentially the same as the threat posed by an end-system that intentionally fails to cooperate with end-to-end congestion control. The deployment of mechanisms in routers to address this threat is an open research question, and is discussed further in <u>Section 7</u>.

Let us take the example described in <u>Section 4.1</u>, where the CE bit that was set in a packet is erased:  $\{(1, 1) \rightarrow (1, 0)\}$ . The consequence for the congested upstream router that set the CE bit is that this congestion indication does not reach the end nodes for that flow. The source (even one which is completely cooperative and not malicious) is thus allowed to continue to increase its sending rate (if it is a TCP flow, by increasing its congestion window). The flow potentially achieves better throughput than the other flows that also share the congested router, especially if there are no policing mechanisms or per-flow queueing mechanisms at that router. Consider the behavior of the other flows, especially if they are cooperative: that is, the flows that do not experience subverted end-to-end congestion control. They are likely to reduce their load (e.g., by reducing their window size) on the congested router, thus benefiting our subverted flow. This results in unfairness. As we discussed above, this unfairness could either be transient (because the congested queue is driven into the packet-marking regime), oscillatory (because the congested queue oscillated between the packet marking and the packet dropping regime), or more moderate but a persistent stable state (because the congested queue is never driven to the packet dropping regime).

[Page 14]

The results would be similar if the subverted flow was intentionally avoiding end-to-end congestion control. One difference is that a flow that is intentionally avoiding end-to-end congestion control at the end nodes can avoid end-to-end congestion control even when the congested queue is in packet-dropping mode, by refusing to reduce its sending rate in response to packet drops in the network. Thus the problems for the network of the subversion of ECN-based congestion control are less severe than the problems caused by the intentional avoidance of end-to-end congestion control in the end nodes. It is also the case that it is considerably more difficult to control the behavior of the end nodes than it is to control the behavior of the infrastructure itself. This is not to say that the problems for the network posed by the network's subversion of ECN-based congestion control are small; just that they are dwarfed by the problems for the network posed by the subversion of either ECN-based or packet-based congestion control by the end nodes.

## **<u>5.2</u>**. Implications for the Subverted Flow

When a source indicates that it is ECN-capable, there is an expectation that the routers in the network that are capable of participating in ECN will use the CE bit for indication of congestion. There is the potential benefit of using ECN in reducing the amount of packet loss (in addition to the reduced queueing delays because of active queue management policies). When the packet flows through a tunnel where the nodes that the tunneled packets traverse are untrusted in some way, the expectation is that IPsec will protect the flow from subversion that results in undesirable consequences.

In many cases, a subverted flow will benefit from the subversion of end-to-end congestion control for that flow in the network, by receiving more bandwidth that it would have otherwise, relative to competing non-subverted flows. If the congested queue reaches the packet-dropping stage, then the subversion of end-to-end congestion control might or might not be of overall benefit to the subverted flow, depending on that flow's relative tradeoffs between throughput, loss, and delay.

One form of subverting end-to-end congestion control is to falsely indicate ECN-capability by setting the ECT bit. This has the consequence of downstream congested routers setting the CE bit in vain. However, as we describe in the section below, if the ECT bit is changed in the IPsec tunnel, this can be detected at the egress point of the tunnel.

The second form of subverting end-to-end congestion control is to erase the congestion indication, either by erasing the CE bit directly, or by erasing the ECT bit when the CE bit is already set.

[Page 15]

In this case, it is the upstream congested routers that set the CE bit in vain. There are several possible scenarios for this subversion of end-to-end congestion control within an IPsec tunnel. If the ECT bit is erased within an IPsec tunnel, then this can be detected at the egress point of the tunnel. If the CE bit is set upstream of the IPsec tunnel, then any erasure of the outer header's CE bit within the tunnel will have no effect because the inner header preserves the set value of the CE bit. However, if the CE bit is set within the tunnel, and erased either within or downstream of the tunnel, this is not necessarily detected at the egress point of the tunnel.

With this subversion of end-to-end congestion control, an end-system transport does not respond to the congestion indication. Along with the increased unfairness for the non-subverted flows described in the previous section, the congested router's queue could continue to build, resulting in packet loss at the congested router - which is a means for indicating congestion to the transport in any case. In the interim, the flow might experience higher queueing delays, possibly along with an increased bandwidth relative to other non-subverted flows. But transports do not inherently make assumptions of consistently experiencing carefully managed queueing in the path. We believe that these forms of subverting end-to-end congestion control are no worse for the subverted flow than if the adversary had simply dropped the packets of that flow itself.

#### 5.3. Non-ECN-Based Methods of Subverting End-to-end Congestion Control

We have shown that, in many cases, a malicious or broken router that is able to change the bits in the ECN field can do no more damage than if it had simply dropped the packet in question. However, this is not true in all cases, in particular in the cases where the broken router subverted end-to-end congestion control by either falsely indicating ECN-Capability or by erasing the ECN congestion indication (in the CE-bit). While there are many ways that a router can harm a flow by dropping packets, a router cannot subvert end-to-end congestion control by dropping packets. As an example, a router cannot subvert TCP congestion control by dropping data packets, acknowledgement packets, or control packets.

Even though packet-dropping cannot be used to subvert end-to-end congestion control, there \*are\* non-ECN-based methods for subverting end-to-end congestion control that a broken or malicious router could use. For example, a broken router could duplicate data packets, thus effectively negating the effects of end-to-end congestion control along some portion of the path. (For a router that duplicated packets within an IPsec tunnel, the security administrator can cause the duplicate packets to be discarded by configuring anti-replay

[Page 16]

protection for the tunnel.) This duplication of packets within the network would have similar implications for the network and for the subverted flow as those described in Sections 5.1 and 5.2 above.

## 6. Changes to the ECN Field within an IPsec Tunnel.

The presence of a copy of the ECN field in the inner header of an IPsec tunnel mode packet provides an opportunity for detection of modifications to the ECT bit in the outer header. Comparison of the ECT bits in the inner and outer headers falls into two categories for implementations that conform to this document:

(a) If the SA allows ECN usage within the tunnel, then the values of the ECT bits in the inner and outer headers are expected be identical.

(b) If the SA disallows ECN usage within the tunnel, then the ECT bit in the outer header is expected to be 0.

Receipt of a packet not satisfying the appropriate condition for its SA is an auditable event, but an implementation MAY create audit records with per-SA counts of incorrect packets over some time period rather than creating an audit record for each erroneous packet. Any such audit record SHOULD contain the headers from at least one erroneous packet, but need not contain the headers from every packet represented by the entry.

An important and likely situation involves an IPsec implementation not updated to this document's requirements serving as tunnel ingress for a tunnel egress at an implementation that has been updated. The ECN Tunnel attribute cannot be negotiated in this case because the tunnel ingress implementation does not support it. If packets from an ECN-capable connection use this tunnel, ECT will be set in the outer header. Congestion along the route could then result in ECNcapable routers setting CE in the outer header. All packets arriving at the tunnel egress on this SA will appear to be case (b) errors, but SHOULD be processed according to whether CE was set. Therefore it is RECOMMENDED that packets violating the condition for case (b) above be dropped if CE is set to 1 in the outer header and forwarded if CE is 0 in the outer header.

An IPsec tunnel cannot provide protection against erasure of congestion indications or false reports of congestion based on flipping the value of the CE bit in packets for which ECT is set in the outer header. As described in <u>Section 5</u>, false reports of congestion are equivalent to dropping the packet, an action against which IPsec also provides no protection. On the other hand, erasure of congestion indications could impact the network and other flows in ways that would not be possible in the absence of ECN. It is important to note that erasure of congestion indications can only be

[Page 17]

performed to congestion indications placed by nodes within the tunnel; the copy of the CE bit in the inner header preserves congestion notifications from nodes upstream of the tunnel ingress. If erasure of congestion notifications is judged to be a security risk that exceeds the congestion management benefits of ECN, the security administrator can configure the appropriate tunnel SAs to forbid ECN usage in the outer header.

# 7. Issues Raised by Monitoring and Policing Devices

One possibility is that monitoring and policing devices (or more informally, `penalty boxes') will be installed in the network to monitor whether best-effort flows are appropriately responding to congestion, and to preferentially drop packets from flows determined not to be using adequate end-to-end congestion control procedures. [FF98] proposes three potential classifications for high-bandwidth flows in times in congestion: (1) flows that are not TCP-friendly, in that the arrival rate from that flow exceeds the arrival rate of a conformant TCP connection under the same conditions; (2) flows that are unresponsive, in that they do not decrease their arrival rate appropriately in response to an increase in congestion; and (3) flows using disproportionate bandwidth, defined as flows using a significantly larger share of bandwidth than other flows in times of high congestion. The methods of identifying and classifying flows to be in one of these three categories is outside the scope of this discussion.

[FF98] proposes that flows that are simply determined to be using disproportionate bandwidth could have their bandwidth restricted, in much the same way that a round-robin per-flow scheduling algorithm would restrict the bandwidth received by individual flows, while flows determined to be unresponsive or not TCP-friendly in times of congestion could have their bandwidth even more strongly reduced, as a concrete incentive to end nodes to use end-to-end congestion control.

For an ECN-capable flow, an `ideal' penalty box at a router would be a device that, when it detected that a flow was not responding to ECN indications, would switch to dropping, instead of marking, those packets of a flow that would otherwise have been chosen to carry indications of congestion. In this way, these congestion indications could not be `erased' later in the network, and at the same time there would be no change in the router's treatment of packets of other flows. If a router determines that a flow is still not responding to congestion indications, when the congestion indications consist of packet drops, then the router could take whatever action it deems appropriate for that flow.

[Page 18]

We RECOMMEND that any `penalty box' that detects a flow or an aggregate of flows that is not responding to end-to-end congestion control first change from marking to dropping packets from that flow, before taking any additional action to restrict the bandwidth available to that flow. Thus, initially, the router could drop packets in which the router would otherwise would have set the CE bit. This could include dropping those arriving packets for that flow that are ECN-Capable and that already have the CE bit set. In this way, any congestion indications seen by that router for that flow will be guaranteed to also be seen by the end nodes, even in the presence of malicious or broken routers elsewhere in the path. If we assume that the first action taken at any `penalty box' for an ECNcapable flow will be to drop packets instead of marking them, then there is no way that an adversary that subverts ECN-based end-to-end congestion control can cause a flow to be characterized as being noncooperative and placed into a more severe action within the `penalty box'.

The monitoring and policing devices that are actually deployed could fall short of the `ideal' monitoring device described above, in that the monitoring is applied not to a single flow or to a single IPsec tunnel, but to an aggregate of flows. In this case, the switch from marking to dropping would apply to all of the flows in that aggregate, denying the benefits of ECN to the other flows in the aggregate also. At the highest level of aggregation, another form of the disabling of ECN happens even in the absence of monitoring and policing devices, when ECN-Capable RED queues switch from marking to dropping packets as an indication of congestion when the average queue size has exceeded some threshold.

# 7.1. Complications Introduced by Split Paths

If a router or other network element has access to all of the packets of a flow, then that router could do no more damage to a flow by altering the ECN field that it could by simply dropping all of the packets from that flow. However, in some cases, a malicious or broken router might have access to only a subset of the packets from a flow. The question is as follows: can this router, by altering the ECN field in this subset of the packets, do more damage to that flow than if it has simply dropped that set of the packets?

We will classify the packets in the flow as A packets and B packets, and assume that the adversary only has access to A packets. Assume that the adversary is subverting end-to-end congestion control along the path traveled by A packets only, by either falsely indicating ECN-Capability upstream of the point where congestion occurs, or erasing the congestion indication downstream. Consider also that there exists a monitoring device that sees both the A and B packets,

[Page 19]

and will "punish" both the A and B packets if the total flow is determined not to be properly responding to indications of congestion. Another key characteristic that we believe is likely to be true is that the monitoring device, before `punishing' the A&B flow, will first drop packets instead of setting the CE bit, and will drop arriving packets of that flow that already have the ECT and CE bits set. If the end nodes are in fact using end-to-end congestion control, they will see all of the indications of congestion seen by the monitoring device, and will begin to respond to these indications of congestion. Thus, the monitoring device is successful in providing the indications to the flow at an early stage.

It is true that the adversary that has access only to the A packets might, by subverting ECN-based congestion control, be able to deny the benefits of ECN to the other packets in the A&B aggregate. While this is unfortunate, this is not a reason to disable ECN within an IPsec tunnel.

A variant of falsely reporting congestion occurs when there are two adversaries along a path, where the first adversary falsely reports congestion, and the second adversary `erases' those reports. (Unlike packet drops, ECN congestion reports can be `reversed' later in the network by a malicious or broken router.) While this would be transparent to the end node, it is possible that a monitoring device between the first and second adversaries would see the false indications of congestion. Given our recommendation in this document, before `punishing' a flow for not responding appropriately to congestion, the router will first switch to dropping rather than marking as an indication of congestion, for that flow. When this includes dropping arriving packets from that flow that have the CE bit set, this ensures that these indications of congestion are being seen by the end nodes. Thus, there is no additional harm that we are able to postulate as a result of multiple conflicting adversaries.

# 8. Comments and Rationale

Substantial comments were received on two areas of this document during review by the ipsec working group. This section describes these comments and explains why the proposed changes were not incorporated.

The first comment indicated that per-node configuration is easier to implement than per-SA configuration. After serious thought and despite some initial encouragement of per-node configuration, it no longer seems to be a good idea. The concern is that as IPsec is progressively deployed, many ECN-aware IPsec implementations will find themselves communicating with a mixture of ECN-aware and ECNunaware IPsec tunnel endpoints. In such an environment with per-node

[Page 20]

configuration, the only reasonable thing to do is forbid ECN usage for all IPsec tunnels, which is not the desired outcome.

In the second area, several reviewers noted that SA negotiation is complex, and adding to it is non-trivial. One reviewer suggested using ICMP after tunnel setup as a possible alternative. The addition to SA negotiation in the draft is OPTIONAL and will remain so; implementers are free to ignore it. The authors believe that the assurance it provides can be useful in a number of situations. In practice, if this is not implemented, it can be deleted at a subsequent stage in the standards process. Extending ICMP to negotiate ECN after tunnel setup is more complex than extending SA attribute negotiation. Some tunnels do not permit traffic to be addressed to the egress endpoint, hence the ICMP packet would have to be addressed to somewhere else, scanned for by the egress endpoint, and discarded there or at its actual destination. In addition, ICMP delivery is unreliable, and hence there is a possibility of an ICMP packet being dropped, entailing the invention of yet another ack/retransmit mechanism. It seems better simply to specify an OPTIONAL extension to the existing SA negotiation mechanism.

#### 9. Conclusions.

This document revises the IPsec architecture to remove a conflict between the experimental usage of Explicit Congestion Notification and IPsec tunnels. This revision consists primarily of modifying the IPsec protocol's handling of the bits in the IP header used by ECN during encapsulation and de-capsulation to allow flows that undergo IPsec tunneling to obtain ECN congestion notifications.

## Two alternatives were described:

1) A preferred full-functionality alternative that copies the ECT bit of the inner header to the encapsulating header. At decapsulation, if the ECT bit is set in the inner header, the CE bit from the outer header is ORed with the CE bit of the inner header to update the CE bit of the packet.

2) A limited-functionality alternative that does not permit generation of ECN notifications inside the IPsec tunnel, by setting the ECT bit in the outer header to zero, and not altering the bits used by ECN in inner header upon decapsulation.

This document also specifies a new IPsec SA attribute that enables negotiation of ECN usage within IPsec tunnels and a new field in the Security Association database to indicate whether ECN is permitted in tunnel mode on a SA.

We examined the consequence of modifications of the ECN field within the tunnel, analyzing all the opportunities for an adversary to

[Page 21]

change the ECN field. In many cases, the change to the ECN field is no worse than dropping a packet. However, we noted that some changes have the more serious consequence of subverting end-to-end congestion control. However, we point out that even then the potential damage is limited, and is similar to the threat posed by an end-system intentionally failing to cooperate with end-to-end congestion control.

In order to permit the experimental usage of ECN with IPsec tunnels, all IPsec implementations MUST implement one of the two alternative approaches described above.

#### 10. Acknowledgements

We thank Steve Bellovin and Vern Paxson for discussions of these matters. We thank Derrell Piper and Kero Tivinen for proposing modifications to  $\frac{\text{RFC}}{2407}$  that improve the usability of negotiating the ECN Tunnel SA attribute.

[Page 22]

# **<u>11</u>**. References

[FF98] Floyd, S., and Fall, K., Promoting the Use of End-to-End Congestion Control in the Internet. IEEE/ACM Transactions on Networking, August 1999. URL "http://wwwnrg.ee.lbl.gov/floyd/end2end-paper.html".

[RFC 2119] S. Bradner, Key words for use in RFCs to Indicate Requirement Levels, <u>RFC 2119</u>, March 1997.

[RFC 2401] S. Kent, R. Atkinson, Security Architecture for the Internet Protocol, <u>RFC 2401</u>, November 1998.

[RFC2407] D. Piper, The Internet IP Security Domain of Interpretation for ISAKMP, <u>RFC 2407</u>, November 1998.

[RFC2474] K. Nichols, S. Blake, F. Baker, D. Black, Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers, <u>RFC 2474</u>, December 1998.

[RFC 2475] S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, and W. Weiss, An Architecture for Differentiated Services, <u>RFC 2475</u>, December 1998.

[RFC2481] K. Ramakrishnan, S. Floyd, A Proposal to add Explicit Congestion Notification (ECN) to IP, <u>RFC 2481</u>, January 1999.

[RFC TBD] S. Bradner and V. Paxson, IANA Allocation Guidelines For Values In the Internet Protocol and Related Headers, Internet-Draft (<u>draft-bradner-iana-allocation-03.txt</u>), November 1999.

## **<u>12</u>**. Security Considerations

Security considerations have been addressed in the main body of the document.

AUTHORS' ADDRESSES

Sally Floyd AT&T Center for Internet Research at ICSI (ACIRI) Phone: +1 (510) 642-4274 x189 Email: floyd@aciri.org URL: http://www.aciri.org/floyd/

David L. Black EMC Corporation 42 South St.

[Page 23]

December 1999

Hopkinton, MA 01748 Phone: +1 (508) 435-1000 x75140 Email: black\_david@emc.com

K. K. Ramakrishnan
AT&T Labs. Research
Phone: +1 (973) 360-8766
Email: kkrama@research.att.com
URL: <u>http://www.research.att.com/info/kkrama</u>

This draft was created in December 1999. It expires June 2000.

[Page 24]