

Internet Congestion Control Research Group
Internet-Draft
Intended status: Experimental
Expires: May 4, 2017

M. Welzl
S. Islam
K. Hiorth
University of Oslo
J. You
Huawei
October 31, 2016

**TCP-CCC: single-path TCP congestion control coupling
draft-welzl-tcp-ccc-00**

Abstract

This document specifies a method, TCP-CCC, to combine the congestion controls of multiple TCP connections between the same pair of hosts. This can have several performance benefits, and it makes it possible to precisely assign a share of the congestion window to the connections based on priorities. This document also addresses the problem that TCP connections between the same pair of hosts may not share the same path. We discuss methods to detect if, or enforce that connections traverse a common bottleneck.

Requirements Language

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](#)].

Status of This Memo

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

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

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."

This Internet-Draft will expire on May 4, 2017.

Copyright Notice

Copyright (c) 2016 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

- [1.](#) Introduction [2](#)
- [2.](#) Coupled Congestion Control [3](#)
- [3.](#) Ensuring a Common Bottleneck [6](#)
 - [3.1.](#) Encapsulation [7](#)
 - [3.1.1.](#) TCP in UDP [7](#)
 - [3.1.2.](#) Other Methods [14](#)
- [4.](#) Related Work [14](#)
- [5.](#) Implementation Status [15](#)
- [6.](#) IANA Considerations [15](#)
- [7.](#) Security Considerations [15](#)
- [8.](#) Acknowledgements [16](#)
- [9.](#) References [16](#)
 - [9.1.](#) Normative References [16](#)
 - [9.2.](#) Informative References [16](#)
- Authors' Addresses [19](#)

1. Introduction

When multiple TCP connections between the same host pair compete on the same bottleneck, they often incur more delay and losses than a single TCP connection. Moreover, it is often not possible to precisely divide the available capacity among the connections. To address this problem, this document presents TCP-CCC, a method to combine the congestion controls of multiple TCP connections between the same pair of hosts. This can have several performance benefits:

- o Reduced average loss and queuing delay (because the competition between the encapsulated TCP connections is avoided)
- o Assign a precise capacity share based on a priority.

- o Even in the absence of prioritization, better fairness between the TCP connections.
- o No need for new connections to slow start up to a reasonable cwnd value that ongoing connections already have: a connection can immediately be assigned its share of the aggregate's total cwnd. This can significantly reduce the completion time of short connections.

All of these benefits only play out when there are more than one TCP connections. Some of the benefits in the list above are more significant when some transfers are short. This makes the usage of TCP-CCC especially attractive in situations where some transfers are short.

We discuss methods to determine if connections traverse the same bottleneck as well as methods to ensure this. To this end, we propose a light-weight, dynamically configured TCP-in-UDP (TiU) encapsulation scheme. TiU is optional, as our coupled congestion control strategy is applicable wherever overlapping TCP flows must follow the same path (such as when routed over a VPN tunnel).

2. Coupled Congestion Control

For each TCP connection c , the algorithm described below receives cwnd and ssthresh as input and stores the following information:

- o the Connection ID.
- o a priority $P(c)$ -- e.g., an integer value in the range from 1 (unimportant) to 10 (very important).
- o The previously used cwnd used by the connection c , $ccc_cwnd(c)$.
- o The previously used ssthresh used by the connection c , $ccc_ssthresh(c)$.

Three global variables sum_cwnd , $sum_ssthresh$ and sum_p are used to represent the sum of all the ccc_cwnd values, $ccc_ssthresh$ values and priorities of all TCP connections, respectively. sum_cwnd and $sum_ssthresh$ are used to update the cwnd and ssthresh values for all connections.

This algorithm emulates the behavior of a single TCP connection by choosing one connection as the connection that dictates the increase / decrease behavior for the aggregate. We call it the "Coordinating Connection" (CoCo). The algorithm was designed to be as simple as possible. Below, abbreviations are used to refer to the phases of

TCP congestion control as defined in [[RFC5681](#)]: SS refers to Slow Start, CA refers to Congestion Avoidance and FR refers to Fast Recovery.

For simplicity, this algorithm refrains from changing cwnd when a connection is in FR. SS should not happen as long as ACKs arrive. Hence, the algorithm ensures that the aggregate's behavior is only dictated by SS when all connections are in the SS phase. We use a bit array, *ssbits*, with a bit for each connection in the group. We set the bit if the connection state is SS due to an RT0.

- (1) When a connection *c* starts, it adds its priority $P(c)$ to *sum_p*. If it is the very first connection, it sets *sum_cwnd* to its own cwnd. After that, the connection's globally known cwnd and ssthresh values (*ccc_cwnd(c)* and *ccc_ssthresh(c)*) are updated, and the connection updates its own cwnd and ssthresh values to be equal to *ccc_cwnd(c)* and *ccc_ssthresh(c)*.

```
ccc_P(c) = P
sum_P = sum_P + P
sum_cwnd = sum_cwnd + cwnd
ccc_cwnd(c) = sum_cwnd / sum_P
ccc_ssthresh(c) = ssthresh
if sum_ssthresh > 0 then
    ccc_ssthresh(c) = sum_ssthresh / sum_P
end if
// Update c's own cwnd and ssthresh for immediate use:
Send ccc_cwnd(c) and ccc_ssthresh(c) to c
```

- (2) When a connection *c* stops, its entry is removed. *sum_p* is recalculated.

```
if c = CoCo then
    CoCo = the next connection
end if
sum_p = sum_p - ccc_P(c)
Remove ccc_P(c), ccc_cwnd(c), ccc_ssthresh(c)
```

- (3) Every time the congestion controller of a connection *c* calculates a new cwnd, the connection calls UPDATE, which carries out the tasks listed below to derive the new cwnd and ssthresh values. Whenever the CoCo calls UPDATE, *sum_cwnd* and *sum_ssthresh* are additionally updated to reflect the current sum of all stored *ccc_cwnd* and *ccc_ssthresh* values. Initially,

there is only one connection and this connection automatically becomes the CoCo. It updates `sum_cwnd` to its own `cwnd` and sets `sum_ssthresh` to 0.

(4) WHEN a non-CoCo connection `c` CALLS UPDATE.....

```
if(all of the connections including CoCo are in CA but c is in FR)
  c becomes the new CoCo.
```

```
else
```

```
  if(c is in CA or SS)
```

```
    c's cwnd is assigned its previously stored ccc_cwnd value.
```

(5) WHEN `c`(CoCo) CALLS UPDATE.....

```
if CoCo == c then
```

```
  if state == CA and ssbits(c) == 0 then
```

```
    if cwnd >= ccc_cwnd(c) then // increased cwnd
```

```
      sum_cwnd = sum_cwnd + cwnd - ccc_cwnd(c)
```

```
    else
```

```
      sum_cwnd = sum_cwnd * cwnd / ccc_cwnd(c)
```

```
    end if
```

```
    ccc_cwnd(c) = ccc_P(c) * sum_cwnd / sum_p
```

```
    ccc_ssthresh(c) ssthresh
```

```
    if sum_ssthresh > 0 then
```

```
      ccc_ssthresh(c) ccc_P(c) * sum_ssthresh/sum_p
```

```
    end if
```

```
  else if state == FR then
```

```
    sum_ssthresh = sum_cwnd/2
```

```
  else if state == SS then
```

```
    if c experienced a timeout then
```

```
      ssbits(c) = 1
```

```
    end if
```

```
    if ssbits(x) == 1 for all x then
```

```
      ssbits(x) = 0 // for all x
```

```
      sum_cwnd = sum_cwnd * cwnd / ccc_cwnd(c)
```

```
      ccc_cwnd(c) = ccc_P(c) * sum_cwnd / sum_p
```

```
      sum_ssthresh = sum_cwnd/2
```

```
    else
```

```
      CoCo = first connection where ccc_state == SS
```

```
    end if
```

```
  end if
```

```
end if
```

(6) After that, if the `ccc_state(c)` is not equal to FR


```
if state != FR then
    Send ccc_cwnd(c) and ccc_ssthresh(c) to c
end if
```

When a flow gets a large share of the aggregate immediately after joining, it can potentially create a burst in the network. We propose a mechanism [[anrw2016](#)] to clock the packet transmission out by using the ack-clock of TCP. Our algorithm achieves a form of "pacing", but it does not rely on any timers.

When a connection *c* joins, it turns on the ack-clock feature and calculates the share of the aggregate, `clocked_cwnd c`. Below, we illustrate the ack-clock mechanism that is used to distribute the share of the `cwnd` based on the acknowledgements received from other flows.

```
if clocked_cwnd(c) <= 0 then
    return // alg. ends; other connections can increase cwnd again
end if
if number_of_acks c % N = 0 then
    send a new segment for connection c
    clocked_cwnd(c) = clocked_cwnd(c) - 1
end if
number_of_acks(c) = number_of_acks(c) + 1
```

3. Ensuring a Common Bottleneck

Our algorithm, as well as EFCM [[EFCM](#)], E-TCP [[EFCM](#)] and the CM [[RFC3124](#)] assume that multiple TCP connections between the same host pair traverse the same bottleneck. This is not always true: load-balancing mechanisms such as Link Aggregation Group (LAG) and Equal-Cost Multi-Path (ECMP) may force them to take different paths [[RFC7424](#)]. If this leads to the connections seeing different bottlenecks, combining the congestion controllers would incur wrong behavior. There are, however, several application scenarios where the single-bottleneck assumption is correct.

Sometimes, the network configuration is known, and it is known that mechanisms such as ECMP and LAG do not operate on the bottleneck or are simply not in use. Alternatively, measurements can infer whether flows traverse the same bottleneck [[I-D.ietf-rmcat-sbd](#)]. When IPv6 is available, the TCP connections could be assigned the same IPv6 flow label. According to [[RFC6437](#)], "The usage of the 3-tuple of the Flow Label, Source Address, and Destination Address fields enables efficient IPv6 flow classification, where only IPv6 main header

fields in fixed positions are used" - this would be favorable for TCP congestion control coupling. However, this [[RFC6437](#)] does not make a clear recommendation about either using the 3-tuple or 5-tuple (which includes the port numbers) - both methods are valid. Thus, whether it works to use the flow label as the sole means to put connections on the same path depends on router configuration. When it works, it is an attractive option because it does not require changing the receiver.

Finally, encapsulating packets with a header that ensures a common path is another possibility to make connections traverse the same bottleneck. We will discuss encapsulation in the next section.

[3.1.](#) Encapsulation

[3.1.1.](#) TCP in UDP

[3.1.1.1.](#) Introduction

We want to be able to ensure that TCP congestion control coupling can always work, provided that the required code is available at the receiver - and be able to efficiently fall back to the standard behaviour in case it is not. To achieve this, we present a method, TCP-in-UDP (TiU), to encapsulate multiple TCP connections using the same UDP port pair.

TCP-in-UDP (TiU) is based on [[Che13](#)]. It differs from it in that:

- o Other than [[Che13](#)], TiU encapsulates multiple TCP connections using the same UDP port number pair. TCP port numbers are preserved; a single well-known UDP port is used for TiU. If TiU is implemented in the kernel, this allows using normal TCP sockets, where enabling the usage of TiU could be done via a socket option, for example.
- o The header format is slightly different to allow representing a TCP connection with a few bits that are encoded across the original TCP header's "Reserved" field and the URG (Urgent) flag to encode a Connection ID. With this encoding, similar to the encapsulation in [[Che13](#)], the total TiU header size does not exceed the original TCP header size.
- o A (TiU-encapsulated) TCP SYN uses a newly defined TCP option to establish the mapping between a Connection ID and the original TCP port number pair.

TiU inherits all the benefits of [[Che13](#)] and a preceding similar proposal, [[Den08](#)]. It enables TCP-CCC coupled congestion control,

and it adds the potential disadvantage of not being able to benefit from ECMP. In short, the benefits and features of TiU that are already explained in detail in [Che13] and [Den08] are:

- o To establish direct communication between two devices that are both behind NAT gateways, Interactive Connectivity Establishment (ICE) [RFC5245] is used to create the necessary mappings in both NAT gateways, and ICE can have higher success rates using UDP [RFC5128].
- o TCP options, as required for Multipath TCP [RFC6824], for example, are expected to work more reliably because middleboxes will be less able to interfere with them.
- o Because the packet format allows the first octet to be in the range 0x0-0x3 (as is the case for a STUN [RFC5389] packet, where the most significant two bits are always zero), the UDP port number pair used by TiU can be used to exchange STUN packets with a STUN server that is unaware of TiU.
- o Following the method described in [Che13] and [Den08], other transport protocols than TCP (e.g., SCTP) could be UDP-encapsulated in a similar fashion. With TiU, the same outer UDP port number pair could be used for different encapsulated protocols at the same time.

[Che13] also lists a disadvantage of UDP-encapsulating TCP packets: because NAT gateways typically use shorter timeouts for UDP port mappings than they do for TCP port mappings, long-lived UDP-encapsulated TCP connections will need to send more frequent keepalive packets than native TCP connections. TiU inherits this problem too, although using a single five-tuple for multiple TCP connections alleviates it by reducing the chance of experiencing long periods of silence.

3.1.1.2. Specification

TiU uses a header that is very similar to the header format in [Den08] and [Che13], where it is explained in greater detail. It consists of a UDP header that is followed by a slightly altered TCP header. The UDP source and destination ports are semantically different from [Den08] and [Che13]: TiU uses a single well-known UDP port, and multiple TCP connections use the same UDP port number pair. The encapsulated TCP header is changed to fit into a UDP packet without increasing the MSS; this is achieved by removing the TCP source and destination ports, the Urgent Pointer and the (now unnecessary) TCP checksum. Moreover, the order of fields is changed to move the Data Offset field to the beginning of the UDP payload.

This allows using it to identify other encapsulated content such as a STUN packet: for TCP, the Data Offset must be at least 5, i.e. the most-significant four bits of the first octet of the UDP payload are in the range 0x5-0xF, whereas this is not the case for other protocols (e.g., STUN requires these bits to be 0). The altered TCP header for TiU is shown below:

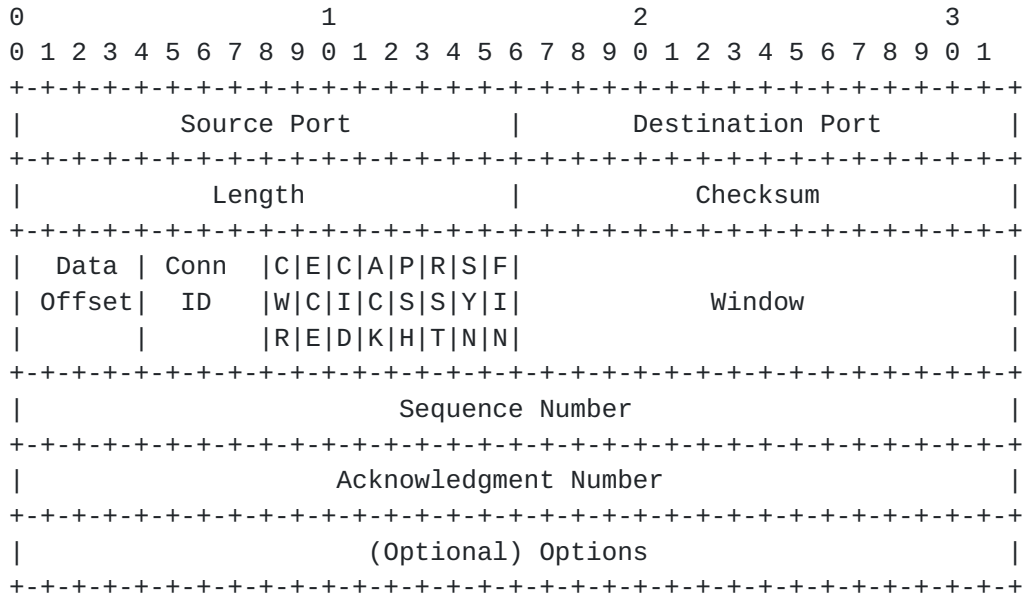


Figure 1: Encapsulated TCP-in-UDP Header Format (the first 8 bytes are the UDP header)

Different from [Den08] and [Che13], the least-significant four bits of the first octet and a bit that replaces the URG bit in the next octet together form a five-bit "Connection ID" (Conn ID). TiU maintains the port numbers of the TCP connections that it encapsulates; the Connection ID is a way to encode the port number information with a few unused header bits. It uniquely identifies a port number pair of a TCP connection that is encapsulated with TiU. Using these five bits, TiU can combine up to 32 TCP connections with one UDP port number pair.

The TiU-TCP SYN and SYN/ACK packets look slightly little different, because they need to establish the mapping between the Connection ID and the port numbers that are used by TiU-encapsulated TCP connections:

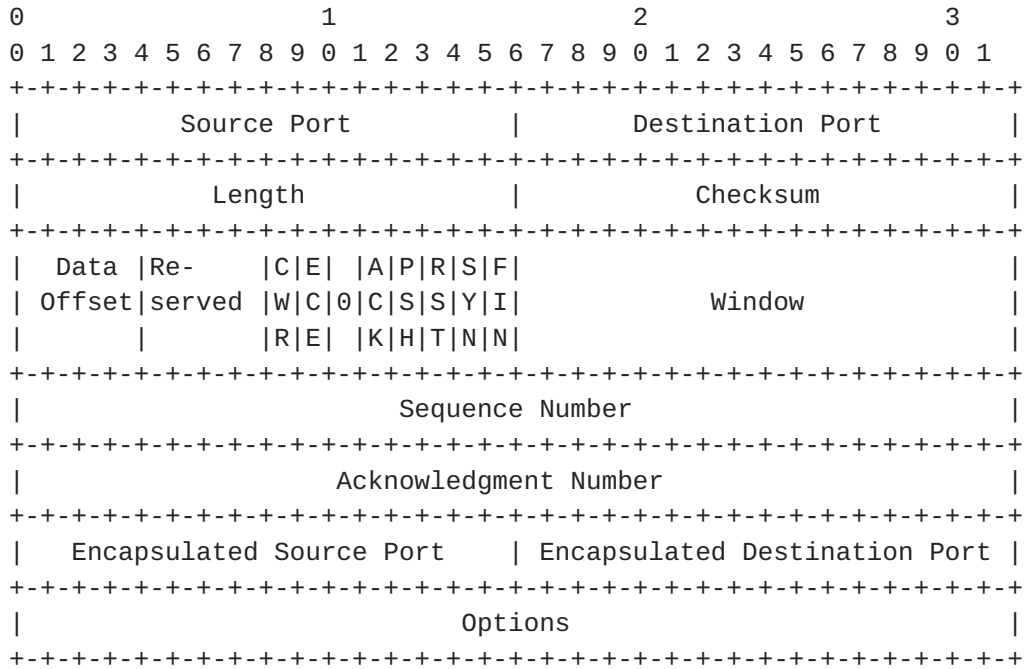


Figure 2: Encapsulated TCP-in-UDP SYN and SYN/ACK Packet Header Format

The Encapsulated Source Port and Encapsulated Destination Port are the port numbers of the TCP connection. To create this header, an implementation can simply swap the position of the original TCP header's port number fields with the position of the Data Offset / Reserved / Flags / Window fields.

Every TiU SYN or TiU SYN-ACK packet also carries at least the TiU-Setup TCP option. This option contains a Connection ID number. On a SYN packet, it is the Connection ID that the sender intends to use in future packets to represent the Encapsulated Source Port and Encapsulated Destination Port. On a SYN/ACK packet, it confirms that such usage is accepted by the recipient of the SYN. A special value of 255 is used to signify an error, upon which TiU will no longer be used (i.e., the next packet is expected to be a non-encapsulated TCP packet). The TiU-Setup TCP option is defined as follows:

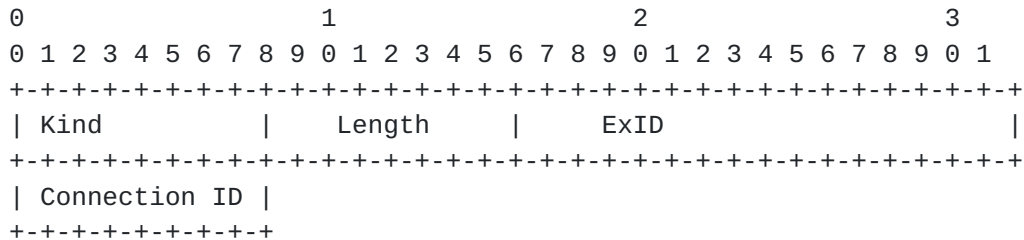


Figure 3: TiU Setup TCP Option

The option follows the format for Experimental TCP Options defined in [RFC6994]. It has Kind=253, Length=5, an ExID that is with value TBD (see [Section 6](#)) and the Connection ID. The Connection ID is an 8-bit field for easier parsing, but only values 0-31 are valid Connection IDs (because the Connection ID in non - SYN or SYN/ACK TiU packets is only 5 bit long).

3.1.1.3. Protocol Operation and Implementation Notes

There can be several ways to implement TCP-in-UDP. The following gives an overview of how a TiU implementation can operate. This description matches the implementation described in [Section 5](#).

A goal of TiU is to achieve congestion control coupling with a simple implementation that minimizes changes to existing code. It is thus recommendable to implement TiU in the kernel, as a change to the existing kernel TCP code. The changes fall in two basic categories:

- o Encapsulation and decapsulation: this is code that should, in the simplest case, operate just before a TCP segment is transmitted. Based on e.g. a socket option that enables/disables TiU, the TCP segment is changed into the TiU header format (Figure 1). In case it is a TCP SYN or TCP SYN/ACK packet, the header format is defined as in Figure 2, and the TiU-Setup TCP option is appended. This packet is then transmitted. For decapsulation, the reverse mechanism applies, upon reception of a UDP packet that uses destination port XXX (TBD, see [Section 6](#)). Both hosts keep a list of encapsulated TCP port numbers and their corresponding Connection IDs. In case a SYN packet requests using a Connection ID that is already reserved, an error (Connection ID value 255 in the TiU Setup TCP option) must be signified to the other end in a TiU-encapsulated TCP SYN/ACK, and encapsulation must be disabled on all further TCP packets. Similarly, when receiving a TiU SYN/ACK with an error, a TCP sender must stop encapsulating TCP packets.

The TCP port number space usage on the host is left unchanged: the original code can reserve TCP ports as it always did. Except for the TiU encapsulation compressing the port numbers into a Connection ID field, TCP ports should be used similar to normal TCP operation. A TCP port that is in use by a TiU-encapsulated TCP connection must therefore not be made available to non-encapsulated TCP connections, and vice versa.

For each TCP connection, two variables must be configured: 1) TiU-ENABLE, which is a boolean, deciding whether to use TiU or not, and 2) Priority, which is a value, e.g. from 1 to 10, that is used by the coupled congestion control algorithm to assign an appropriate share

of the total cwnd to the connection. Priority values are local and their range does not matter for this algorithm: the algorithm works with a flow's priority portion of the sum of all priority values. The configuration of the two per-connection variables can be implemented in various ways, e.g. through an API option.

With these code changes in place, TiU can operate as follows, assuming no previous TiU connections have been made between a specific host pair and a client tries to connect to a server:

- o An application uses an API option to request TiU operation. The kernel then sends out a TiU TCP SYN that contains a TiU-Setup TCP option. This packet header contains the encapsulated TCP port numbers (source port A and destination port B) and the Connection ID X.
- o The server listens on UDP port XXX (TBD, see [Section 6](#)). Upon receiving a packet on this port, it knows that it is a TiU packet and decodes it, handing the resulting TCP packet over to "normal" TCP processing. The TiU-Setup TCP option allows the server to associate future TiU packets containing Connection ID X with ports A and B. The server sends its response as a TiU SYN-ACK.
- o TCP operates as normal from here on, but packets are TiU-encapsulated before sending them out and decapsulated upon reception, using Connection ID X. Both hosts associate TiU packets carrying Connection ID X with a local identifier that matches ports A and B, just like they would associate non-encapsulated TCP packets with the same local identifier when seeing ports A and B in the TCP header.
- o If an application on either side of the TiU connection wants to connect to a destination host on the other side and requests TiU operation, the kernel sends out another TiU TCP SYN, this time containing a different TCP source port number and either the same or a different destination port number (C and D), and a TiU-Setup TCP option with Connection ID Y. From now on, packets carrying Connection ID Y will be associated with ports C and D on both hosts. Otherwise, TiU operation continues as described above.
- o Now, because there are two or more connections available between the same host pair, coupled congestion control begins to operate for all outgoing TiU packets (see [Section 2](#) for details). This is a local operation, applying the priority values that were configured to use for the TiU-encapsulated TCP connections.

Unless it is known that UDP packets with destination port number XXX (TBD, see [Section 6](#)) can be used without problems on the path between

two communicating hosts, it is advisable for TiU implementations to contain methods to fall back to non-encapsulated ("raw") TCP communication. Such fall-back must be supported for the case of Connection ID collisions anyway. Middleboxes have been known to track TCP connections [[Honda11](#)], and falling back to communication with raw TCP packets without ever using a raw TCP SYN - SYN/ACK handshake may lead to problems with such devices. The following method is recommended to efficiently fall back to raw TCP communication:

- o After sending out a TiU SYN packet, additionally send a raw TCP SYN packet.
- o After sending out a TiU SYN/ACK packet, additionally send a raw TCP SYN/ACK packet.
- o Upon receiving a TiU SYN packet, after responding with a TiU SYN/ACK packet and raw TCP SYN/ACK packet, immediately store the encapsulated port numbers and Connection ID. As long as a TiU connection is ongoing, ignore any additional incoming TCP SYN or TCP SYN/ACK packets from the same host that carry port numbers matching the stored encapsulated port numbers. Otherwise, process TCP SYN or TCP SYN/ACK packets as normal.

This method ensures that the TCP SYN / SYN/ACK handshake is visible to middleboxes and allows to immediately switch back to raw TCP communication in case of failures. If implemented on both sides as described above and no TiU SYN or TiU SYN/ACK packet arrives, yet a TCP SYN or TCP SYN/ACK packet does, this can only mean that the other host does not support TiU, a UDP packet was dropped, or the UDP and TCP packets were reordered in transit. Reordering in the host (e.g., a server responding to a TCP SYN before it responds to a TiU SYN) can be a problem for similar methods (e.g. [[RFC6555](#)]), but it can be eliminated by prescribing the processing order as above.

Because TCP does not preserve message boundaries and the size of the TCP header can vary depending on the options that are used, it is also no problem to precede the TCP header in the UDP packet with a different header (e.g. PLUS or SPUD [[I-D.hildebrand-spud-prototype](#)]) without exceeding the known MTU limit. When creating a TCP segment, a TCP sender needs to consider the length of this header when calculating the segment size, just like it would consider the length of a TCP option. For this to work, the usage of other headers such as PLUS or SPUD in-between the UDP header and the TiU header must therefore be known to both the sender-side and receiver-side code that processes TiU.

3.1.1.4. Usage Considerations

TiU cannot work with applications that require the Urgent pointer (which is not recommended for use by new applications anyway [[RFC6093](#)], but should be considered if TiU is implemented in a way that allows it to be applied onto existing applications; telnet is a well-known example of an application that uses this functionality). It can also be used as a method to experimentally test new TCP functionality in the presence of middleboxes that would otherwise create problems (as some have been known to do [[Honda11](#)]).

Reasons to use TiU include the benefits of [[Che13](#)] and [[Den08](#)] that were discussed in [Section 1](#). TiU has the disadvantage of disabling ECMP for the TCP connections that it encapsulates. This can reduce the capacity usage of these TCP connections. It has the advantage of being able to apply TCP-CCC coupled congestion control, which can provide precise congestion window assignment based on a priority.

3.1.2. Other Methods

There are many possible encapsulation schemes for various use cases. For example, Generic UDP Encapsulation (GUE) [[I-D.draft-ietf-nvo3-gue](#)] allows us to multiplex several TCP connections onto a same UDP port number pair. Several encapsulation methods transmit layer-2 frames over an IP network - e.g. VXLAN [[RFC7348](#)] (over UDP/IP) and NVGRE [[RFC7637](#)] (over GRE/IP). Because Layer-2 networks should be agnostic to the transport connections running over them, the path should not depend on the TCP port number pair and our algorithm should work. Some care must still be taken: for example, for NVGRE, [[RFC7637](#)] says: "If ECMP is used, it is RECOMMENDED that the ECMP hash is calculated either using the outer IP frame fields and entire Key field (32 bits) or the inner IP and transport frame fields". If routers do use the inner transport frame fields (typically, port numbers) for this hashing, we have the same problem even over NVGRE.

4. Related Work

The TCPMUX mechanism in [[RFC1078](#)] multiplexes TCP connections under the same outer transport port number; it does however not preserve the port numbers of the original TCP connections, and no method to couple congestion controls is described in [[RFC1078](#)].

Congestion control coupling follows the style of RTP application congestion control coupling in [[I-D.ietf-rmcat-coupled-cc](#)] which is designed to be easy to implement, and to minimize the number of changes that need to be made to the underlying congestion control mechanisms. This method was shown to yield several benefits in

[[fse](#)]. TCP-CCC requires slightly deeper changes to TCP's congestion control, making it harder to implement than [[I-D.ietf-rmcat-coupled-cc](#)], but it is still a much smaller code change than the Congestion Manager [[RFC3124](#)].

Combining congestion controls as TCP-CCC does it has some similarities with Ensemble Sharing in [[RFC2140](#)], which however only concerns initial values of variables used by new connections and does not share the congestion window (cwnd). The cwnd variable is shared across ongoing connections in [[ETCP](#)] and [[EFCM](#)], and the mechanism described in [Section 2](#) resembles the mechanisms in these works, but neither [[ETCP](#)] nor [[EFCM](#)] address the problem of ECMP.

Coupled congestion control has also been specified for Multipath TCP [[RFC6356](#)]. MPTCP's coupled congestion control combines the congestion controls of subflows that may traverse different paths, whereas we propose congestion control coupling for flows sharing a single-path. TCP-CCC builds on the assumption that all its encapsulated TCP connections traverse the same path. This makes the two methods for coupled congestion control very different, even though they both aim at emulating the behavior of a single TCP connection in the case where all flows traverse the same network bottleneck. For example, a new flow obtaining a larger-than-IW share of the aggregate cwnd would be inappropriate for an MPTCP subflow.

5. Implementation Status

We have implemented TCP-CCC and TiU encapsulation for both the sender and receiver in the FreeBSD kernel, as a simple add-on to the TCP implementation that is controlled via a socket option.

6. IANA Considerations

This document specifies a new TCP option that uses the shared experimental options format [[RFC6994](#)]. No value has yet been assigned for EXID.

This document requires a well-known UDP port (referred to as port XXX in this document). Due to the highly experimental nature of TiU, this document is being shared with the community to solicit comments before requesting such a port number.

7. Security Considerations

TBD

8. Acknowledgements

This work has received funding from Huawei Technologies Co., Ltd., and the European Union's Horizon 2020 research and innovation programme under grant agreement No. 644334 (NEAT). The views expressed are solely those of the author(s).

9. References

9.1. 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>>.

9.2. Informative References

- [anrw2016] Islam, S. and M. Welzl, "Start Me Up:Determining and Sharing TCP's Initial Congestion Window", ACM, IRTF, ISOC Applied Networking Research Workshop 2016 (ANRW 2016) , 2016.
- [Che13] Cheshire, S., Graessley, J., and R. McGuire, "Encapsulation of TCP and other Transport Protocols over UDP", Internet-draft [draft-cheshire-tcp-over-udp-00](#), June 2013.
- [Den08] Denis-Courmont, R., "UDP-Encapsulated Transport Protocols", Internet-draft [draft-denis-udp-transport-00](#), July 2008.
- [EFCM] Savoric, M., Karl, H., Schlager, M., Poschwatta, T., and A. Wolisz, "Analysis and performance evaluation of the EFCM common congestion controller for TCP connections", Computer Networks (2005) , 2005.
- [ETCP] Eggert, L., Heidemann, J., and J. Joe, "Effects of ensemble-TCP", ACM SIGCOMM Computer Communication Review (2000) , 2000.

- [fse] Islam, S., Welzl, M., Gjessing, S., and N. Khademi, "Coupled Congestion Control for RTP Media", ACM SIGCOMM Capacity Sharing Workshop (CSWS 2014) and ACM SIGCOMM CCR 44(4) 2014; extended version available as a technical report from <http://safiquili.at.ifi.uio.no/paper/fse-tech-report.pdf> , 2014.
- [Honda11] Honda, M., Nishida, Y., Raiciu, C., Greenhalgh, A., Handley, M., and H. Tokuda, "Is it still possible to extend TCP?", Proc. of ACM Internet Measurement Conference (IMC) '11, November 2011.
- [I-D.[draft-ietf-nvo3-gue](#)]
Herbert, T., Yong, L., and O. Zia, "Generic UDP Encapsulation", Internet-draft [draft-ietf-nvo3-gue-05](#), October 2016.
- [I-D.hildebrand-spud-prototype]
Hildebrand, J. and B. Trammell, "Substrate Protocol for User Datagrams (SPUD) Prototype", [draft-hildebrand-spud-prototype-03](#) (work in progress), March 2015.
- [I-D.ietf-rmcat-coupled-cc]
Islam, S., Welzl, M., and S. Gjessing, "Coupled congestion control for RTP media", [draft-ietf-rmcat-coupled-cc-03](#) (work in progress), July 2016.
- [I-D.ietf-rmcat-sbd]
Hayes, D., Ferlin, S., Welzl, M., and K. Hiorth, "Shared Bottleneck Detection for Coupled Congestion Control for RTP Media.", [draft-ietf-rmcat-sbd-04](#) (work in progress), March 2016.
- [RFC1078] Lottor, M., "TCP port service Multiplexer (TCPMUX)", [RFC 1078](#), DOI 10.17487/RFC1078, November 1988, <<http://www.rfc-editor.org/info/rfc1078>>.
- [RFC2140] Touch, J., "TCP Control Block Interdependence", [RFC 2140](#), DOI 10.17487/RFC2140, April 1997, <<http://www.rfc-editor.org/info/rfc2140>>.
- [RFC3124] Balakrishnan, H. and S. Seshan, "The Congestion Manager", [RFC 3124](#), DOI 10.17487/RFC3124, June 2001, <<http://www.rfc-editor.org/info/rfc3124>>.

- [RFC5128] Srisuresh, P., Ford, B., and D. Kegel, "State of Peer-to-Peer (P2P) Communication across Network Address Translators (NATs)", [RFC 5128](#), DOI 10.17487/RFC5128, March 2008, <<http://www.rfc-editor.org/info/rfc5128>>.
- [RFC5245] Rosenberg, J., "Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols", [RFC 5245](#), DOI 10.17487/RFC5245, April 2010, <<http://www.rfc-editor.org/info/rfc5245>>.
- [RFC5389] Rosenberg, J., Mahy, R., Matthews, P., and D. Wing, "Session Traversal Utilities for NAT (STUN)", [RFC 5389](#), DOI 10.17487/RFC5389, October 2008, <<http://www.rfc-editor.org/info/rfc5389>>.
- [RFC5681] Allman, M., Paxson, V., and E. Blanton, "TCP Congestion Control", [RFC 5681](#), DOI 10.17487/RFC5681, September 2009, <<http://www.rfc-editor.org/info/rfc5681>>.
- [RFC6093] Gont, F. and A. Yourtchenko, "On the Implementation of the TCP Urgent Mechanism", [RFC 6093](#), DOI 10.17487/RFC6093, January 2011, <<http://www.rfc-editor.org/info/rfc6093>>.
- [RFC6356] Raiciu, C., Handley, M., and D. Wischik, "Coupled Congestion Control for Multipath Transport Protocols", [RFC 6356](#), DOI 10.17487/RFC6356, October 2011, <<http://www.rfc-editor.org/info/rfc6356>>.
- [RFC6437] Amante, S., Carpenter, B., Jiang, S., and J. Rajahalme, "IPv6 Flow Label Specification", [RFC 6437](#), DOI 10.17487/RFC6437, November 2011, <<http://www.rfc-editor.org/info/rfc6437>>.
- [RFC6555] Wing, D. and A. Yourtchenko, "Happy Eyeballs: Success with Dual-Stack Hosts", [RFC 6555](#), DOI 10.17487/RFC6555, April 2012, <<http://www.rfc-editor.org/info/rfc6555>>.
- [RFC6824] Ford, A., Raiciu, C., Handley, M., and O. Bonaventure, "TCP Extensions for Multipath Operation with Multiple Addresses", [RFC 6824](#), DOI 10.17487/RFC6824, January 2013, <<http://www.rfc-editor.org/info/rfc6824>>.
- [RFC6994] Touch, J., "Shared Use of Experimental TCP Options", [RFC 6994](#), DOI 10.17487/RFC6994, August 2013, <<http://www.rfc-editor.org/info/rfc6994>>.

- [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>>.
- [RFC7424] Krishnan, R., Yong, L., Ghanwani, A., So, N., and B. Khasnabish, "Mechanisms for Optimizing Link Aggregation Group (LAG) and Equal-Cost Multipath (ECMP) Component Link Utilization in Networks", [RFC 7424](#), DOI 10.17487/RFC7424, January 2015, <<http://www.rfc-editor.org/info/rfc7424>>.
- [RFC7637] Garg, P., Ed. and Y. Wang, Ed., "NVGRE: Network Virtualization Using Generic Routing Encapsulation", [RFC 7637](#), DOI 10.17487/RFC7637, September 2015, <<http://www.rfc-editor.org/info/rfc7637>>.

Authors' Addresses

Michael Welzl
University of Oslo
PO Box 1080 Blindern
Oslo N-0316
Norway

Email: michawe@ifi.uio.no

Safiqul Islam
University of Oslo
PO Box 1080 Blindern
Oslo N-0316
Norway

Phone: +47 22 84 08 37
Email: safiquли@ifi.uio.no

Kristian Hiorth
University of Oslo
PO Box 1080 Blindern
Oslo N-0316
Norway

Email: kristahi@ifi.uio.no

Jianjie You
Huawei
101 Software Avenue, Yuhua District
Nanjing 210012
China

Email: youjianjie@huawei.com