

Internet Engineering Task Force  
Internet-Draft  
Intended status: Standards Track  
Expires: January 22, 2015

A. Bittau  
D. Boneh  
M. Hamburg  
Stanford University  
M. Handley  
University College London  
D. Mazieres  
Q. Slack  
Stanford University  
July 21, 2014

**Cryptographic protection of TCP Streams (tcpcrypt)  
draft-bittau-tcpinc-01.txt**

**Abstract**

This document presents tcpcrypt, a TCP extension for cryptographically protecting TCP segments. Tcpcrypt maintains the confidentiality of data transmitted in TCP segments against a passive eavesdropper. It protects connections against denial-of-service attacks involving desynchronizing of sequence numbers, and when enabled, against forged RST segments. Finally, applications that perform authentication can obtain end-to-end confidentiality and integrity guarantees by tying authentication to tcpcrypt Session ID values.

The extension defines two new TCP options, CRYPT and MAC, which are designed to provide compatible interworking with TCPs that do not implement tcpcrypt. The CRYPT option allows hosts to negotiate the use of tcpcrypt and establish shared secret encryption keys. The MAC option carries a message authentication code with which hosts can verify the integrity of transmitted TCP segments. Tcpcrypt is designed to require relatively low overhead, particularly at servers, so as to be useful even in the case of servers accepting many TCP connections per second.

**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 January 22, 2015.

#### Copyright Notice

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

This document is subject to [BCP 78](http://trustee.ietf.org/license-info) 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.

This document may contain material from IETF Documents or IETF Contributions published or made publicly available before November 10, 2008. The person(s) controlling the copyright in some of this material may not have granted the IETF Trust the right to allow modifications of such material outside the IETF Standards Process. Without obtaining an adequate license from the person(s) controlling the copyright in such materials, this document may not be modified outside the IETF Standards Process, and derivative works of it may not be created outside the IETF Standards Process, except to format it for publication as an RFC or to translate it into languages other than English.



## Table of Contents

<a href="#">1.</a>	<a href="#">Requirements Language</a>	<a href="#">5</a>
<a href="#">2.</a>	<a href="#">Introduction</a>	<a href="#">5</a>
<a href="#">3.</a>	<a href="#">Idealized protocol</a>	<a href="#">5</a>
<a href="#">3.1.</a>	<a href="#">Stages of the protocol</a>	<a href="#">5</a>
<a href="#">3.1.1.</a>	<a href="#">The setup phase</a>	<a href="#">6</a>
<a href="#">3.1.2.</a>	<a href="#">The ENCRYPTING state</a>	<a href="#">6</a>
<a href="#">3.1.3.</a>	<a href="#">The DISABLED state</a>	<a href="#">7</a>
<a href="#">3.2.</a>	<a href="#">Cryptographic algorithms</a>	<a href="#">7</a>
<a href="#">3.3.</a>	<a href="#">"C" and "S" roles</a>	<a href="#">9</a>
<a href="#">3.4.</a>	<a href="#">Key exchange protocol</a>	<a href="#">9</a>
<a href="#">3.5.</a>	<a href="#">Data encryption and authentication</a>	<a href="#">12</a>
<a href="#">3.6.</a>	<a href="#">Authenticated Sequence Mode (ASM)</a>	<a href="#">13</a>
<a href="#">3.6.1.</a>	<a href="#">ASM-Encrypt</a>	<a href="#">14</a>
<a href="#">3.6.2.</a>	<a href="#">ASM-Decrypt</a>	<a href="#">15</a>
<a href="#">3.6.3.</a>	<a href="#">ASM-Update</a>	<a href="#">16</a>
<a href="#">3.7.</a>	<a href="#">Re-keying</a>	<a href="#">16</a>
<a href="#">3.8.</a>	<a href="#">Session caching</a>	<a href="#">17</a>
<a href="#">3.8.1.</a>	<a href="#">Session caching control</a>	<a href="#">17</a>
<a href="#">4.</a>	<a href="#">Extensions to TCP</a>	<a href="#">18</a>
<a href="#">4.1.</a>	<a href="#">Protocol states</a>	<a href="#">18</a>
<a href="#">4.2.</a>	<a href="#">Role negotiation</a>	<a href="#">23</a>
<a href="#">4.2.1.</a>	<a href="#">Simultaneous open</a>	<a href="#">24</a>
<a href="#">4.3.</a>	<a href="#">The TCP CRYPT option</a>	<a href="#">25</a>
<a href="#">4.3.1.</a>	<a href="#">The HELLO suboption</a>	<a href="#">28</a>
<a href="#">4.3.2.</a>	<a href="#">The DECLINE suboption</a>	<a href="#">29</a>
<a href="#">4.3.3.</a>	<a href="#">The NEXTK1 and NEXTK2 suboptions</a>	<a href="#">29</a>
<a href="#">4.3.4.</a>	<a href="#">The PKCONF suboption</a>	<a href="#">31</a>
<a href="#">4.3.5.</a>	<a href="#">The UNKNOWN suboption</a>	<a href="#">32</a>
<a href="#">4.3.6.</a>	<a href="#">The SYNCCOOKIE and ACKCOOKIE suboptions</a>	<a href="#">33</a>
<a href="#">4.3.7.</a>	<a href="#">The SYNC_REQ and SYNC_OK suboptions</a>	<a href="#">33</a>
<a href="#">4.3.8.</a>	<a href="#">The REKEY and REKEYSTREAM suboptions</a>	<a href="#">35</a>
<a href="#">4.3.9.</a>	<a href="#">The INIT1 and INIT2 suboptions</a>	<a href="#">38</a>
<a href="#">4.4.</a>	<a href="#">The TCP MAC option</a>	<a href="#">39</a>
<a href="#">5.</a>	<a href="#">Examples</a>	<a href="#">41</a>
<a href="#">5.1.</a>	<a href="#">Example 1: Normal handshake</a>	<a href="#">42</a>
<a href="#">5.2.</a>	<a href="#">Example 2: Normal handshake with SYN cookie</a>	<a href="#">42</a>
<a href="#">5.3.</a>	<a href="#">Example 3: tcpcrypt unsupported</a>	<a href="#">42</a>
<a href="#">5.4.</a>	<a href="#">Example 4: Reusing established state</a>	<a href="#">43</a>
<a href="#">5.5.</a>	<a href="#">Example 5: Decline of state reuse</a>	<a href="#">43</a>
<a href="#">5.6.</a>	<a href="#">Example 6: Reversal of client and server roles</a>	<a href="#">43</a>
<a href="#">6.</a>	<a href="#">API extensions</a>	<a href="#">43</a>
<a href="#">7.</a>	<a href="#">Acknowledgments</a>	<a href="#">46</a>
<a href="#">8.</a>	<a href="#">IANA Considerations</a>	<a href="#">46</a>
<a href="#">9.</a>	<a href="#">Security Considerations</a>	<a href="#">48</a>
<a href="#">10.</a>	<a href="#">References</a>	<a href="#">49</a>
<a href="#">10.1.</a>	<a href="#">Normative References</a>	<a href="#">49</a>



<a href="#">10.2</a> . Informative References . . . . .	<a href="#">49</a>
<a href="#">Appendix A</a> . Protocol constant values . . . . .	<a href="#">50</a>
Authors' Addresses . . . . .	<a href="#">50</a>

## **1. 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 [RFC 2119](#) [[RFC2119](#)].

## **2. Introduction**

This document describes tcpcrypt, an extension to TCP for cryptographic protection of session data. Tcpcrypt was designed to meet the following goals:

- o Maintain confidentiality of communications against a passive adversary. Ensure that an adversary must actively intercept and modify the traffic to eavesdrop, either by re-encrypting all traffic or by forcing a downgrade to an unencrypted session.
- o Minimize computational cost, particularly on servers.
- o Provide interfaces to higher-level software to facilitate end-to-end security, either in the application level protocol or after the fact. (E.g., client and server log session IDs and can compare them after the fact; if there was no tampering or eavesdropping, the IDs will match.)
- o Be compatible with further extensions that allow authenticated resumption of TCP connections when either end changes IP address.
- o Facilitate multipath TCP [[RFC6824](#)] by identifying a TCP stream with a session ID independent of IP addresses and port numbers.
- o Provide for incremental deployment and graceful fallback, even in the presence of NATs and other middleboxes that might remove unknown options, and traffic normalizers.

## **3. Idealized protocol**

This section describes the tcpcrypt protocol at an abstract level, without reference to particular cryptographic algorithms or data encodings. Readers who simply wish to see the key exchange protocol should skip to [Section 3.4](#).

### **3.1. Stages of the protocol**

A tcpcrypt endpoint goes through multiple stages. It begins in a setup phase and ends up in one of two states, ENCRYPTING or DISABLED,





before applications may send or receive data. The ENCRYPTING and DISABLED states are definitive and mutually exclusive; an endpoint that has been in one of the two states MUST NOT ever enter the other, nor ever re-enter the setup phase.

#### **3.1.1. The setup phase**

The setup phase negotiates use of the tcpcrypt extension. During this phase, two hosts agree on a suite of cryptographic algorithms and establish shared secret session keys.

The setup phase uses the Data portion of TCP segments to exchange cryptographic keys. Implementations MUST NOT include application data in TCP segments during setup and MUST NOT allow applications to read or write data. System calls MUST behave the same as for TCP connections that have not yet entered the ESTABLISHED state; calls to read and write SHOULD block or return temporary errors, while calls to poll or select SHOULD consider connections not ready.

When setup succeeds, tcpcrypt enters the ENCRYPTING state. Importantly, a successful setup also produces an important value called the `_Session ID_`. The Session ID is tied to the negotiated algorithms and cryptographic keys, and is unique over all time with overwhelming probability.

Operating systems MUST make the Session ID available to applications. To prevent man-in-the-middle attacks, applications MAY authenticate the session ID through any protocol that ensures both endpoints of a connection have the same value. Applications MAY alternatively just log Session IDs so as to enable attack detection after the fact through comparison of the values logged at both ends.

The setup phase can also fail for various reasons, in which case tcpcrypt enters the DISABLED state.

Applications MAY test whether setup succeeded by querying the operating system for the Session ID. Requests for the Session ID MUST return an error when tcpcrypt is not in the ENCRYPTING state. Applications SHOULD authenticate the returned Session ID. Applications relying on tcpcrypt for security SHOULD authenticate the Session ID and SHOULD treat unauthenticated Session IDs the same as connections in the DISABLED state.

#### **3.1.2. The ENCRYPTING state**

When the setup phase succeeds, tcpcrypt enters the ENCRYPTING state. Once in this state, applications may read and write data with the expected semantics of TCP connections.



In the ENCRYPTING state, a host MUST encrypt the Data portion of all TCP segments transmitted and MUST include a Message Authentication Code (MAC) in all segments transmitted. A host MUST furthermore ignore any TCP segments received without the RST bit set, unless those segments also contain a valid MAC option.

A host SHOULD accept RST segments without valid MACs by default. However, the application SHOULD be allowed to force unMACed RST segments to be dropped by enabling the TCP\_CRYPT\_RSTCHK option on the connection.

Once in the ENCRYPTING state, an endpoint MUST NOT directly or indirectly transition to the DISABLED state under any circumstances.

### **3.1.3. The DISABLED state**

When setup fails, tcpcrypt enters the DISABLED state. In this case, the host MUST continue just as TCP would without tcpcrypt, unless network conditions would cause a plain TCP connection to fail as well. Entering the DISABLED state prohibits the endpoint from ever entering the ENCRYPTING state.

An implementation MUST behave identically to ordinary TCP in the DISABLED state, except that the first segment transmitted after entering the DISABLED state MAY include a TCP CRYPT option with a DECLINE suboption (and optionally other suboptions such as UNKNOWN) to indicate that tcpcrypt is supported but not enabled.

[Section 4.3.2](#) describes how this is done.

Operating systems MUST allow applications to turn off tcpcrypt by setting the state to DISABLED before opening a connection. An active opener with tcpcrypt disabled MUST behave identically to an implementation of TCP without tcpcrypt. A passive opener with tcpcrypt disabled MUST also behave like normal TCP, except that it MAY optionally respond to SYN segments containing a CRYPT option with SYN-ACK segments containing a DECLINE suboption, so as to indicate that tcpcrypt is supported but not enabled.

## **3.2. Cryptographic algorithms**

The setup phase employs three types of cryptographic algorithms:

- o A `_public key cipher_` is used with a short-lived public key to exchange (or agree upon) a random, shared secret.
- o An `_extract function_` is used to generate a pseudo-random key from some initial keying material, typically the output of the public key cipher. The notation `Extract (S, IKM)` denotes the output of



the extract function with salt *S* and initial keying material *IKM*.

- o A *\_collision-resistant pseudo-random function (CPRF)\_* is used to generate multiple cryptographic keys from a pseudo-random key, typically the output of the extract function. We use the notation *CPRF (K, TAG, L)* to designate the output of *L* bytes of the pseudo-random function identified by key *K* on *TAG*. A collision-resistant function is one on which, for sufficiently large *L*, an attacker cannot find two distinct inputs *K\_1, TAG\_1* and *K\_2, TAG\_2* such that *CPRF (K\_1, TAG\_1, L) = CPRF (K\_2, TAG\_2, L)*. Collision resistance is important to assure the uniqueness of Session IDs, which are generated using the CPRF.

The Extract and CPRF functions used by default are the Extract and Expand functions of HKDF [[RFC5869](#)]. These are defined as follows:

```
HKDF-Extract(salt, IKM) -> PRK
    PRK = HMAC-Hash(salt, IKM)

HKDF-Expand(PRK, TAG, L) -> OKM
    T(0) = empty string (zero length)
    T(1) = HMAC-Hash(PRK, T(0) | TAG | 0x01)
    T(2) = HMAC-Hash(PRK, T(1) | TAG | 0x02)
    T(3) = HMAC-Hash(PRK, T(2) | TAG | 0x03)
    ...

    OKM = first L octets of T(1) | T(2) | T(3) | ...
```

The symbol *|* denotes concatenation, and the counter concatenated with TAG is a single octet.

Because the public key cipher, the extract function, and the expand function all make use of cryptographic hashes in their constructions, the three algorithms are negotiated as a unit employing a single hash function. For example, the OAEP+-RSA [[RFC2437](#)] cipher, which uses a SHA-256-based mask-generation function, is coupled with HKDF [[RFC5869](#)] using HMAC-SHA256 [[RFC2104](#)].

The encrypting phase employs an *\_authenticated encryption mode\_* to encrypt all application data. This mode authenticates both application data and most of the TCP header (excepting header fields commonly modified by middleboxes).

Note that public key generation, public key encryption, and shared secret generation all require randomness. Other tcpcrypt functions may also require randomness depending on the algorithms and modes of operation selected. A weak pseudo-random generator at either host will defeat tcpcrypt's security. Thus, any host implementing



tcpcrypt MUST have a cryptographically secure source of randomness or pseudo-randomness.

### **3.3. "C" and "S" roles**

Tcpcrypt transforms a single pseudo-random key (PRK) into cryptographic session keys for each direction. Doing so requires an asymmetry in the protocol, as the key derivation function must be perturbed differently to generate different keys in each direction. Tcpcrypt includes other asymmetries in the roles of the two hosts, such as the process of negotiating algorithms (e.g., proposing vs. selecting cipher suites).

We use the terms "C" and "S" to denote the distinct roles of the two hosts in tcpcrypt's setup phase. In the case of key transport, "C" is the host that supplies a public key, while "S" is the host that encrypts a pre-master secret with the key belonging to "C". Which role a host plays can have performance implications, because for some public key algorithms encryption is much faster than decryption. For instance, on a machine at the time of writing, encryption with a 2,048-bit RSA-3 key is over two orders of magnitude faster than decryption.

Because servers often need to establish connections at a faster rate than clients, and because servers are often passive openers, by default the passive opener plays the "S" role. However, operating systems MUST provide a mechanism for the passive opener to reverse roles and play the "C" role, as discussed in [Section 4.2](#).

### **3.4. Key exchange protocol**

Every machine C has a short-lived public encryption key or key agreement parameter, PK\_C, which gets refreshed periodically and SHOULD NOT ever be written to persistent storage.

When a host C connects to S, the two engage in the following protocol:

```
C -> S: HELLO
S -> C: PKCONF, pub-cipher-list
C -> S: INIT1, sym-cipher-list, N_C, pub-cipher, PK_C
S -> C: INIT2, sym-cipher, KX_S
```

The parameters are defined as follows:

- o pub-cipher-list: a list of public key ciphers and parameters acceptable to S. These are defined in Figure 3.





- o sym-cipher-list: a list of symmetric cipher suites acceptable to C. These are specified in Table 6 as parameters for ASM mode, discussed in [Section 3.6](#).
- o N\_C: Nonce chosen at random by C.
- o pub-cipher: the type of PK\_C.
- o PK\_C: C's public key or key agreement parameter.
- o sym-cipher: the symmetric cipher selected by S.
- o KX\_S: key exchange information produced by S. KX\_S will depend on whether key transport is being done (e.g., RSA) or key agreement (e.g., Diffie-Hellman). KX\_S is defined in Table 1.

+-----+	+-----+	+-----+	+-----+
Cipher	KX_S	PMS	
+-----+	+-----+	+-----+	+-----+
OAEP+-RSA exp3	ENC (PK_C, R_S)	R_S	
ECDHE	N_S, PK_S	key-agreement-output	
+-----+	+-----+	+-----+	+-----+

ENC (PK\_C, R\_S) denotes an encryption of R\_S with public key PK\_C. R\_S and N\_S are random values chosen by S. Their lengths are defined in Figure 3. PK\_S is S's key agreement parameter. PMS is the Pre Master Secret from which keys are ultimately derived.

Table 1

The two sides then compute a pseudo-random key (PRK) from which all session keys are derived as follows:

```
param := { num-pub-ciphers, pub-cipher-list, init1, init2 }
PRK   := Extract (N_C, { param, PMS })
```

Here num-pub-ciphers is a single octet specifying how many three-byte algorithm specifiers were provided by the "S" host in a PKCONF suboption (described in [Section 4.3.4](#)). pub-cipher-list is this many three-byte specifiers, taken from the body of the PKCONF suboption. init1 and init2 are the complete data payload from the TCP segments with the INIT1 and INIT2 suboptions (detailed in [Section 4.3.9](#)).

A series of "session secrets" and corresponding Session IDs are then computed as follows:



```

ss[0] := PRK
ss[i] := CPRF (ss[i-1], CONST_NEXTK, K_LEN)

SID[i] := CPRF (ss[i], CONST_SESSID, K_LEN)

```

The value `ss[0]` is used to generate all key material for the current connection. `SID[0]` is the session ID for the current connection, and will with overwhelming probability be unique for each individual TCP connection. The most computationally expensive part of the key exchange protocol is the public key cipher. The values of `ss[i]` for  $i > 0$  can be used to avoid public key cryptography when establishing subsequent connections between the same two hosts, as described in [Section 3.8](#). The TAG values are constants defined in Table 7. The `K_LEN` values and nonce sizes are negotiated, and are specified in Figure 3.

Given a session secret, `ss`, the two sides compute a series of master keys as follows:

```

mk[0] := CPRF (ss, CONST_REKEY | flags, K_LEN)
mk[i] := CPRF (mk[i-1], CONST_REKEY, K_LEN)

```

Where `flags` is a single octet from 0x0 to 0x3 computed as follows:

```

bit  0 1 2 3 4 5 6 7
+--+--+--+--+--+--+
|0 0 0 0 0 0 s c|
+--+--+--+--+--+--+

```

Here bit "s" is set when the "S" mode host has indicated application-level support for tcpcrypt. The "c" bit is set when the "C" mode host has indicated application-level support for tcpcrypt. Both bits are 0 by default unless the application has enabled the `TCP_CRYPT_SUPPORT` option described in [Section 6](#).

Finally, each master key `mk` is used to generate keys for authenticated encryption for the "S" and "C" roles. Key `k_cs` is used by "C" to encrypt and "S" to decrypt, while `k_sc` is used by "S" to encrypt and "C" to decrypt.

```

k_cs := CPRF (mk, CONST_KEY_C, ae_len)
k_sc := CPRF (mk, CONST_KEY_S, ae_len)

```

tcpcrypt does not use HKDF directly for key derivation because it requires multiple expand steps with different keys. This is needed for forward secrecy so that `ss[n]` can be forgotten once a session is established, and `mk[n]` can be forgotten once a session is rekeyed.



There is no key confirmation step in tcpcrypt. This is not required since in tcpcrypt's threat model, a connection to an adversary can be made and so keys need not be verified. If an erroneous key negotiation that yields two different keys occurs, all subsequent packets will be dropped due to an incorrect MAC, causing the TCP connection to hang. This is not a threat because in plain TCP, an active attacker could have modified sequence and ack numbers to hang the connection anyway.

### **3.5. Data encryption and authentication**

tcpcrypt encrypts and authenticates all application data. It also authenticates some parts of the TCP header. There are several TCP-specific constraints with regards to authenticated encryption that tcpcrypt must meet for performance and compatibility with middleboxes:

- o The ciphertext for a particular byte position in tcpcrypt's sequence must never change, even if reencryption occurs after coalescing and retransmission. This is because a middlebox may discard a changed payload on retransmission.
- o Authentication must occur only on fields not modified by middleboxes. In particular, port numbers must not be authenticated, and sequence and ack numbers must be authenticated according to an offset from the initial sequence number, because these can be modulated by a middlebox.
- o An efficient mechanism is needed for recomputing the authentication tag when only the ack numbers change. For example, on retransmissions, the authenticated encryption authentication tag can be efficiently updated without having to recompute the tag on the entire packet payload.

Authenticated encryption modes such as GCM do not meet these criteria. For example, even with identical plaintext, ciphertext values depend on the byte position at which one starts encrypting a segment. Hence two small segments will appear to have different content from their coalesced counterpart; middleboxes might drop such coalesced retransmissions after falsely detecting subterfuge attacks. Furthermore, existing authenticated encryption modes do not allow efficient updating of the authentication tag when only small parts of the data have changed. A new mode is needed to meet all these constraints, and we introduce Authenticated Sequence Mode (ASM) in [Section 3.6](#) as a solution.

ASM takes three parameters: a cipher, a MAC and an ACK MAC. At a high-level, the cipher is used to encrypt the TCP payload in counter



mode, using a counter derived from TCP's sequence number. The MAC covers the ciphertext and parts of the TCP header. The ACK MAC covers the ACK numbers and is XORed with the previously computed MAC to produce the authenticated encryption authentication tag. This tag can be quickly updated if only the ACK numbers have changed. This approach is principled because ACK messages are conceptually separate from data packets, so MACing them separately is appropriate. In TCP, ACKs are piggybacked to data segments merely as an optimization.

XORing two PRF-based MACs together was shown secure by Katz and Lindell [[aggregate-macs](#)].

### **3.6. Authenticated Sequence Mode (ASM)**

ASM is parameterized by a cipher, MAC and ACK MAC. The operations supported by ASM are:

ASM-Encrypt (PRK, Seq, Message, Assoc-Data, Up-Data) ->  
(Ciphertext, Auth-Tag)

ASM-Decrypt (PRK, Seq, Cipher-Text, Assoc-Data, Up-Data, Auth-Tag) ->  
{ (Valid, Message) OR  
(Invalid, )  
}

ASM-Update (PRK, Up-Data-Prev, Up-Data-New, Auth-Tag-Prev) ->  
Auth-Tag

The arguments and return values are:

- o `_PRK_` a pseudo-random key.
- o `_Seq_` the byte position in the stream of Message or Cipher-Text. In tcpcrypt, this is an extended version of TCP's sequence number.
- o `_Message_` the Message to encrypt. In tcpcrypt, this is TCP's payload.
- o `_Assoc-Data_` the associated data to be MACed but not encrypted. In tcpcrypt, this contains parts of the TCP header.
- o `_Up-Data_` the updatable data to be MACed but not encrypted, that can also be efficiently updated and reMACed. In tcpcrypt, this will cover an extended version of TCP's ACK numbers.





- o `_Ciphertext_` the encrypted version of Message.
- o `_Auth-Tag_` the authenticated encryption authentication tag. In `tcpcrypt`, this will be the MAC option.

ASM-Decrypt returns one of the constants Valid or Invalid, depending on whether the authentication tag can be verified successfully or not. For Valid inputs, the Message is returned as well.

The PRK supplied to ASM is expanded into keys used for individual operation as follows:

```
k_enc := CPRF (PRK, CONST_KEY_ENC, cipher-key-len)
k_mac := CPRF (PRK, CONST_KEY_MAC, mac-key-len)
k_ack := CPRF (PRK, CONST_KEY_ACK, ack-mac-key-len)
```

The next sections describe ASM operations in detail.

### **3.6.1. ASM-Encrypt**

The interface to encrypt is as follows:

```
ASM-Encrypt (PRK, Seq, Message, Assoc-Data, Up-Data) ->
(Ciphertext, Auth-Tag)
```

Keys (denoted by `k_*`) are derived from PRK as explained in [Section 3.6](#).

The following steps occur:

1. Message is encrypted to produce Ciphertext using the cipher in counter mode. Seq is the counter and `k_enc` is the key. When encrypting Seq, its value must always be a multiple of the cipher's block size. In the event that the message does not begin on an even block boundary, Seq must be rounded down, encrypted, and leading bytes of its encryption discarded.
2. The MAC is run over the concatenation of Ciphertext and Assoc-Data to produce MAC1, using `k_mac` as the key.
3. The ACK MAC is run over Up-Data to produce MAC2, using `k_ack` as the key.
4. MAC1 and MAC2 are XORed to produce Auth-Tag.

Using AES-128 as an example, encryption in counter mode using Seq as the counter happens as follows.



- o Compute  $B = \text{Seq} - (\text{Seq} \% 16)$ .
- o Let  $B^* = 0^{128-|B|} \parallel B$  be  $B$  in network (big-endian) byte order with enough 0 bits pre-pended to make  $B^*$  exactly 128 bits long.
- o Let  $C = \text{ENC-AES}(k_{\text{enc}}, B^*)$ .
- o Discard the first  $(\text{Seq}-B)$  bytes on  $C$  and begin byte-by-byte XORing the remaining portion with the message.
- o Continue computing  $\text{ENC-AES}(k_{\text{enc}}, B^* + 16)$ ,  $\text{ENC-AES}(k_{\text{enc}}, B^* + 32)$ , etc. to generate enough bytes to XOR with the whole message.

If AES-128 is used as the ACK MAC, the Ack number (64-bit extended, offset from ISN) is first padded on the left with enough zeros to produce a 128-bit big-endian value. The number is then encrypted using AES.

### [3.6.2.](#) ASM-Decrypt

The interface to decrypt is as follows:

ASM-Decrypt (PRK, Seq, Cipher-Text, Assoc-Data, Up-Data, Auth-Tag) ->  
    { (Valid, Message) OR  
      (Invalid, )

Keys (denoted by  $k_*$ ) are derived from PRK as explained in [Section 3.6](#).

The following steps occur:

1. The MAC is run over the concatenation of Ciphertext and Assoc-Data to produce MAC1, using  $k_{\text{mac}}$  as the key.
2. The ACK MAC is run over Up-Data to produce MAC2, using  $k_{\text{ack}}$  as the key.
3. MAC1 and MAC2 are XORed and compared to Auth-Tag. If different, the process stops and the constant Invalid is returned along with no message. Otherwise the process continues.
4. Ciphertext is decrypted to produce Message using the cipher in counter mode. Seq is the counter and  $k_{\text{enc}}$  is the key. The Valid constant is returned along with Message.



### **3.6.3. ASM-Update**

The interface to update the authenticated encryption authentication tag is as follows:

```
ASM-Update (PRK, Up-Data-Prev, Up-Data-New, Auth-Tag-Prev) ->
Auth-Tag
```

Keys (denoted by  $k_*$ ) are derived from PRK as explained in [Section 3.6](#).

The following steps occur:

1. The ACK MAC is run over Up-Data-Prev using  $k_{ack}$  to produce MAC2-Prev.
2. MAC2-Prev is XORed with Auth-Tag-Prev to produce MAC1.
3. The ACK MAC is run over Up-Data to produce MAC2, using  $k_{ack}$  as the key.
4. MAC1 and MAC2 are XORed to produce Auth-Tag.

### **3.7. Re-keying**

We refer to the two encryption keys ( $k_{cs}$ ,  $k_{sc}$ ) as a key set. We refer to the key set generated by  $mk[i]$  as the key set with generation number  $i$  within a session. Initially, the two hosts use the key set with generation number 0.

Either host may decide to evolve the encryption key at one or more points within a session, by incrementing the generation number of its transmit keys. When switching keys to generation  $j$ , a host must label the segments it transmits with a REKEY option containing  $j$ , so that the recipient host knows to check the MAC and decrypt the segment using the new keyset:

```
A -> B: REKEY<j>, MAC<...>, Data<...>
```

Upon receiving a REKEY< $j$ > segment, a recipient using transmit keys from a generation less than  $j$  must also update its transmit keys and start including a REKEY< $j$ > option in all of its segments. A host must continue transmitting REKEY options until all segments with other generation numbers have been processed at both ends.

Implementations **MUST** always transmit and retransmit identical ciphertext Data bytes for the same TCP sequence numbers. Thus, a retransmitted segment **MUST** always use the same keyset as the original



segment. Hosts MUST NOT combine segments that were encrypted with different keysets.

Implementations SHOULD delete older-generation keys from memory once they have received all segments they will need to decrypt with the old keys and received acknowledgments for all segments that would need to be retransmitted encrypted under old keys.

### **3.8. Session caching**

When two hosts have already negotiated session secret `ss[i-1]`, they can establish a new connection without public key operations using `ss[i]`. The four-message protocol of [Section 3.4](#) is replaced by:

```
A -> B: NEXTK1, SID[i]
B -> A: NEXTK2
```

Which symmetric keys a host uses for transmitted segments is determined by its role in the original session `ss[0]`. It does not depend on which host is the passive opener in the current session. If A had the "C" role in the first session, then A uses `k_cs` for sending segments and `k_sc` for receiving. Otherwise, if A had the "S" role originally, it uses `k_sc` and `k_cs`, respectively. B similarly uses the transmit keys that correspond to its role in the original session.

After using `ss[i]` to compute `mk[0]`, implementations SHOULD compute and cache `ss[i+1]` for possible use by a later session, then erase `ss[i]` from memory. Hosts SHOULD keep `ss[i+1]` around for a period of time until it is used or the memory needs to be reclaimed. Hosts SHOULD NOT write a cached `ss[i+1]` value to non-volatile storage.

It is an implementation-specific issue as to how long `ss[i+1]` should be retained if it is unused. If the passive opener times it out before the active opener does, the only cost is the additional ten bytes to send NEXTK1 for the next connection. The behavior then falls back to a normal public-key handshake.

#### **3.8.1. Session caching control**

Implementations MUST allow applications to control session caching by setting the following option:

**TCP\_CRYPT\_CACHE\_FLUSH** When set on a TCP endpoint that is in the ENCRYPTING state, this option causes the operating system to flush from memory the cached `ss[i+1]` (or `ss[i+1+n]` if other connections have already been established). When set on an endpoint that is in the setup phase, causes any cached `ss[i]` that would have been





used to be flushed from memory. In either case, future connections will have to undertake another round of the public key protocol in [Section 3.4](#). Applications SHOULD set TCP\_CRYPT\_CACHE\_FLUSH whenever authentication of the session ID fails.

#### **[4.](#) Extensions to TCP**

The tcpcrypt extension adds two new kinds of option: CRYPT and MAC. Both are described in this section. During the setup phase, all TCP segments MUST have the CRYPT option. In the ENCRYPTING state, all segments MUST have the MAC option and may include the CRYPT option for various purposes such as re-keying or keep-alive probes.

The idealized protocol of the previous section is embedded in the TCP handshake. Unfortunately, since the maximum TCP header size is 60 bytes and the basic TCP header fields require 20 bytes, there are at most 40 option payload bytes available, which is not enough to hold the INIT1 and INIT2 messages. Tcpcrypt therefore uses the Data portion of TCP segments (after the SYN exchanges) to send the body of these messages.

Operating systems MUST keep track of which phase a data segment belongs to, and MUST only deliver data to applications from segments that are processed in the ENCRYPTING or DISABLED states.

##### **[4.1.](#) Protocol states**

The setup phase is divided into six states: CLOSED, NEXTK-SENT, HELLO-SENT, C-MODE, LISTEN, and S-MODE. Together with the ENCRYPTING and DISABLED states already discussed, this means a tcpcrypt endpoint can be in one of eight states.

In addition to tcpcrypt's state, each endpoint will also be in one of the 11 TCP states described in the TCP protocol specification [\[RFC0793\]](#). Not all pairs of states are valid. Table 2 shows which TCP states an endpoint can be in for each tcpcrypt state.



Tcpcrypt state	TCP states for an active opener	TCP states for a passive opener
CLOSED	CLOSED	CLOSED
NEXTK-SENT	SYN-SENT	n/a
HELLO-SENT	SYN-SENT	SYN-RCVD
C-MODE	ESTABLISHED, FIN-WAIT-1	ESTABLISHED, FIN-WAIT-1
LISTEN	n/a	LISTEN
S-MODE	(SYN-RCVD), ESTABLISHED	SYN-RCVD
ENCRYPTING	(SYN-RCVD), ESTABLISHED+	SYN-RCVD, ESTABLISHED+
DISABLED	any	any

Valid tcpcrypt and TCP state combinations. States in parentheses occur only with simultaneous open. ESTABLISHED+ means ESTABLISHED or any later state (FIN-WAIT-1, FIN-WAIT-2, CLOSING, TIME-WAIT, CLOSE-WAIT, or LAST-ACK).

Table 2

Figure 1 shows how tcpcrypt transitions between states. Each transition is labeled by events that may trigger the transition above the line, and an action the local host is permitted to take in response below the line. "snd" and "rcv" denote sending and receiving segments, respectively. "internal" means any possible event except for receiving a segment (i.e., timers and system calls). "drop" means discarding the last received segment and preventing it from having any effect on TCP's state. "mac" means any valid TCP action, including no action, except that any segments transmitted must be encrypted and contain a valid TCP MAC option. "x" indicates that a host sends no segments when taking a transition.

A segment is described as "F/Op". F specifies constraints on the control bits of the TCP header, as follows:

F	Meaning
S	SYN=1, ACK=0, FIN=0, RST=0
SA	SYN=1, ACK=1, FIN=0, RST=0
A	SYN=0, ACK=1, FIN=0, RST=0
S?	SYN=1, ACK=any, FIN=0, RST=0
?A	SYN=any, ACK=1, FIN=0, RST=0
R	RST=1
*	any



Op designates message types in the abstract protocol, which also correspond to particular suboptions of the TCP CRYPT option, described in [Section 4.3](#), or "MAC" for a valid TCP MAC option, as described in [Section 4.4](#). A segment with SYN=1 and ACK=0 that contains the NEXTK1 suboption will also explicitly or implicitly contain the HELLO suboption; such a segment matches event constraints on either option--e.g., it matches any of the "rcv S/HELLO", "rcv S?/HELLO", and "rcv S/NEXTK1" events. An empty Op matches any segment with the appropriate control bits. A segment MUST contain the TCP MAC option if and only if Op is "MAC".

The "drop" transitions from NEXTK-SENT and HELLO-SENT to HELLO-SENT change TCP slightly by ignoring a segment and preventing a TCP transition from SYN-SENT to SYN-RCVD that would otherwise occur during simultaneous open. Therefore, these transitions SHOULD be disabled by default. They MAY be enabled on one side by an application that wishes to enable tcpcrypt on simultaneous open, as discussed in [Section 4.2.1](#).



Except for these drop actions, tcpcrypt MUST abide by the TCP protocol specification [RFC0793]. Thus, any segment transmitted by a host MUST be permitted by the TCP specification in addition to matching either a transition in Figure 1 or one of the transitions to





DISABLED or CLOSED described below. In particular, a host MUST NOT acknowledge an INIT1 segment unless either the acknowledgment contains an INIT2 or the host transitions to DISABLED.

Various events cause transitions to DISABLED from states other than ENCRYPTING. In particular:

- o Operating systems MUST provide a mechanism for applications to transition to DISABLED from the CLOSED and LISTEN states.
- o A host in the setup phase MUST transition to DISABLED upon receiving any segment without a TCP CRYPT option.
- o A host in the setup phase MUST transition to DISABLED upon receiving any segment with the FIN or RST control bit set.
- o A host in the setup phase MUST transition to DISABLED upon sending a segment with the FIN bit set. (As discussed below, however, a host MUST NOT send a FIN segment from the C-MODE state.)

Other specific conditions cause a transition to DISABLED and are discussed in the sections that follow.

CLOSED is a pseudo-state representing a connection that does not exist. A tcpcrypt connection's lifetime is identical to that of its associated TCP connection. Thus, tcpcrypt transitions to CLOSED exactly when TCP transitions to CLOSED.

A host MUST NOT send a FIN segment from the C-MODE state. The reason is that the remote side can be in the ENCRYPTING state and would thus require the segment to contain a valid MAC, yet a host in C-MODE cannot compute the necessary encryption keys before receiving the INIT2 segment.

If a CLOSE happens in C-MODE, a host MUST delay sending a FIN segment until receiving an ACK for its INIT1 segment. If the remote host is in ENCRYPTING, the ACK segment will contain INIT2 and the local host can transition to ENCRYPTING before sending the FIN. If the remote host is not in ENCRYPTING, the ACK will not contain INIT2, and thus the local host can transition to DISABLED before sending the FIN.

If a CLOSE happens in C-MODE, an implementation MAY delay processing the CLOSE event and entering the TCP FIN-WAIT-1 state until sending the FIN. If it does not, the implementation MUST ensure all relevant timers correspond to the time of transmission of the FIN segment, not the time of entry into the FIN-WAIT-1 state.

The only valid tcpcrypt state transition from ENCRYPTING is to



CLOSED, which occurs only when TCP transitions to CLOSED. tcpcrypt per se cannot cause TCP to transition to CLOSED.

#### **4.2. Role negotiation**

A passive opener receiving an S/HELLO segment may choose to play the "S" role (by transitioning to S-MODE) or the "C" role (by transitioning to HELLO-SENT). An active opener may accept the role not chosen by the passive opener, or may instead disable tcpcrypt. During simultaneous open, one endpoint must choose the "C" role while the other chooses the "S" role. Operating systems **MUST** allow applications to guide these choices on a per-connection basis.

Applications **SHOULD** be able to exert this control by setting a per-connection `_CMODE disposition_`, which can take on one of the following five values:

`TCP_CRYPT_CMODE_DEFAULT` This disposition **SHOULD** be the default. A passive opener will only play the "S" role, but an active opener can play either the "C" or the "S" role. Simultaneous open without session caching will cause tcpcrypt to be disabled unless the remote host has set the `TCP_CMODE_ALWAYS[_NK]` disposition.

`TCP_CRYPT_CMODE_ALWAYS`

`TCP_CRYPT_CMODE_ALWAYS_NK` With this disposition, a host will only play the "C" role. The `_NK` version additionally prevents the use of session caching if the session was originally established in the "S" role.

`TCP_CRYPT_CMODE_NEVER`

`TCP_CRYPT_CMODE_NEVER_NK` With this disposition, a host will only play the "S" role. The `_NK` version additionally prevents the use of session caching if the session was originally established in the "C" role.

The CMODE disposition prohibits certain state transitions, as summarized in Table 3. If an event occurs for which all valid transitions in Figure 1 are prohibited, a host **MUST** transition to DISABLED. Operating systems **MAY** add additional CMODE dispositions, for instance to force or prohibit session caching.



CMODE disposition	Prohibited transitions
TCP_CRYPT_CMODE_DEFAULT	LISTEN --> HELLO-SENT HELLO-SENT --> HELLO-SENT NEXTK-SENT --> HELLO-SENT
TCP_CRYPT_CMODE_ALWAYS[_NK]	any --> S-MODE
TCP_CRYPT_CMODE_NEVER[_NK]	LISTEN --> HELLO-SENT HELLO-SENT --> HELLO-SENT NEXTK-SENT --> HELLO-SENT any --> C-MODE

State transitions prohibited by each CMODE disposition

Table 3

#### 4.2.1. Simultaneous open

During simultaneous open, two ends of a TCP connection are both active openers. If both hosts attempt to use session caching by simultaneously transmitting S/NEXTK1 segments, and if both transmit the same session ID, then both may reply with SA/NEXTK2 segments and immediately enter the ENCRYPTING state. In this case, the host that played "C" when the session was initially negotiated MUST use the symmetric encryption keys for "C" (i.e., encrypt with `k_cs`, decrypt with `k_sc`), while the host that initially played "S" uses the "S" keys for the new connection.

If both hosts in a simultaneous open do not attempt to use session caching, or if the two hosts use incompatible Session IDs, then they MUST engage in public-key-based key negotiation to use tcpcrypt. Doing so requires one host to play the "C" role and the other to play the "S" role. With the `TCP_CRYPT_CMODE_DEFAULT` disposition, these roles are usually determined by the passive opener choosing the "S" role. With no passive opener, both active openers will end up in S-MODE, then transition to DISABLED upon receiving an unexpected PKCONF.

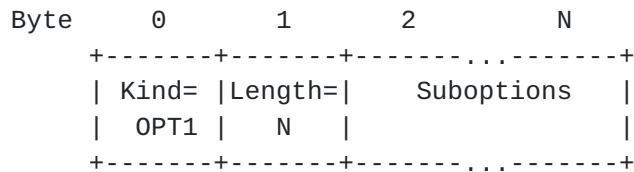
Simultaneous open can work with key negotiation if exactly one of the two hosts selects the `TCP_CRYPT_CMODE_ALWAYS` disposition. This host will then drop S/HELLO segments and remain in C-MODE while the other host transitions to S-MODE. Applications SHOULD NOT set `TCP_CRYPT_CMODE_ALWAYS` on both sides of a simultaneous open, as this will result in tcpcrypt being disabled. The reception of two simultaneous HELLO (or NEXTK) messages will disable tcpcrypt because



it is not explicit as to who is playing the "C" or "S" role.

### 4.3. The TCP CRYPT option

A CRYPT option has the following format:

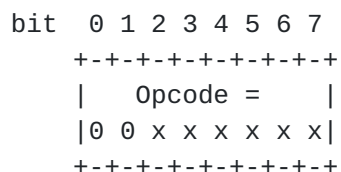


## Format of TCP CRYPT option

Kind is always OPT1. Length is the total length of the option, including the two bytes used for Kind and Length. These first two bytes are then followed by zero or more suboptions. Suboptions determine the meaning of the TCP CRYPT option. When a TCP header contains more than one CRYPT option, a host MUST interpret them the same as if all the suboptions appeared in a single CRYPT option. This makes tcpcrypt options future-proof as new suboptions can be placed in a separate CRYPT option, which can be ignored if not understood, while other CRYPT options can still be processed.

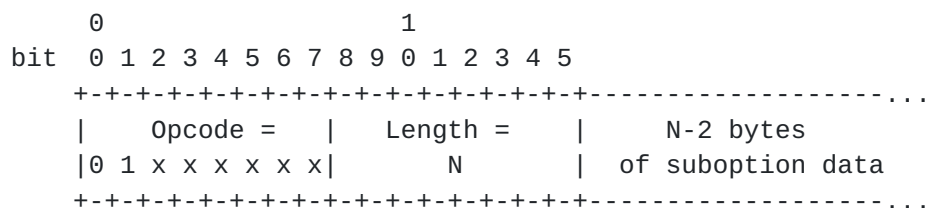
Each suboption begins with an Opcode byte. The specific format of the option depends on the two most significant bits of the Opcode.

Suboptions with opcodes from 0x00 to 0x3f contain no data other than the single opcode byte:



Hosts **MUST** ignore any opcodes of this format that they do not recognize.

Suboptions with opcodes from 0x40 to 0x7f contain an opcode, a length field, and data bytes.







Hosts MUST ignore any opcodes of this format that they do not recognize.

Suboptions with opcodes from 0x80 to 0xbf contain zero or more bytes of data whose length depends on the opcode. These suboptions can be either fixed length or variable length; implementations that understand these opcodes will know which they are; if the suboption is fixed length the implementation will know the length; otherwise it will know where to look for the length field.

```

bit  0 1 2 3 4 5 6 7
    +-+-+-+-+-+-+-+-----...
    |      Opcode =   | data
    |1 0 x x x x x x|
    +-+-+-+-+-+-+-+-----...

```

If a host sees an unknown opcode in this range, it MUST ignore the suboption and all subsequent suboptions in the same TCP CRYPT option. However, if more than one CRYPT option appears in the TCP header, the host MUST continue processing suboptions from the next TCP CRYPT option. Skipping suboptions in the TCP CRYPT option applies only to this option range since the length of the suboption cannot be determined by the receiver. In other cases, where the length is known, the receiver skips to the next suboption.

Suboptions with opcodes from 0xc0 to 0xff also contain an opcode-specific length of data. As before, these suboptions can be either fixed length or variable length. Suboptions in this range are classed as mandatory as far as the protocol is concerned. However, they are not MANDATORY to implement unless otherwise stated, as discussed below.

```

bit  0 1 2 3 4 5 6 7
    +-+-+-+-+-+-+-+-----...
    |      Opcode =   | data
    |1 1 x x x x x x|
    +-+-+-+-+-+-+-+-----...

```

Should a host encounter an unknown opcode greater than or equal to 0xc0 during the setup phase of the protocol, the host MUST transition to the DISABLED state. It SHOULD respond with both a DECLINE suboption and an UNKNOWN suboption specifying the opcode of the unknown mandatory suboption, after which the host MUST NOT send any further CRYPT options.

Should a host encounter an unknown opcode greater than or equal to 0xc0 while in the ENCRYPTING state, the host MUST respond with an UNKNOWN suboption specifying the opcode of the unknown mandatory



suboption, and should ensure the session continues with the same encryption and authentication state as it had before the segment was received. This may require ignoring other suboptions within the same message, or reverting any half-negotiated state.

Table 4 summarizes the opcodes discussed in this document. It is MANDATORY that all implementations support every opcode in this table. Each opcode is listed with the length in bytes of the suboption (including the opcode byte), or \* for variable-length suboptions. The last column specifies in which of the (S)etup phase, (E)NCRYPTING state, and (D)ISABLED state an opcode may be used. A host MUST NOT send an option unless it is in one of the stages indicated by this column.

Value	Length	Name	Stages
0x01	1	HELLO	S
0x02	1	HELLO-app-support	S
0x03	1	HELLO-app-mandatory	S
0x04	1	DECLINE	SD
0x05	1	NEXTK2	S
0x06	1	NEXTK2-app-support	S
0x07	1	INIT1	S
0x08	1	INIT2	S
0x41	*	PKCONF	S
0x42	*	PKCONF-app-support	S
0x43	*	UNKNOWN	SED
0x44	*	SYNCOOKIE	S
0x45	*	ACKCOOKIE	SED
0x80	5	SYNC_REQ	E
0x81	5	SYNC_OK	E
0x82	2	REKEY	E
0x83	6	REKEYSTREAM	E
0x84	10	NEXTK1	S

Opcodes for suboptions of the TCP CRYPT option.

Table 4

If a TCP segment (sent by an active opener) has the SYN flag set, the ACK flag clear, and one or more TCP CRYPT options, there is an implicit HELLO suboption even if that suboption does not appear in the segment. In particular, when such a SYN segment contains a single, empty, two-byte TCP CRYPT option, the passive opener MUST interpret that option as equivalent to the three-byte TCP option composed of bytes OPT1, 3, 1 (Kind = OPT1, Length = 3, Suboption =



HELLO).

A host **MUST** enter the **DISABLED** state if, during the setup phase, it receives a segment containing neither a TCP CRYPT nor a TCP MAC option. This is for robustness against middleboxes that strip options. A host **MUST** also enter **DISABLED** if, during the setup phase, it receives a **DECLINE** suboption or any unrecognized suboption with opcode greater than or equal to 0xc0. The **DECLINE** option is the preferred way for a host to refuse tcpcrypt. A host **MAY** also choose reply without a TCP CRYPT option to disable tcpcrypt. Once a host has entered **DISABLED**, it **MUST NOT** include the MAC option in any transmitted segment. The host **MAY** include a CRYPT option in the next segment transmitted, but only if the segment also contains the **DECLINE** suboption. All subsequently transmitted packets **MUST NOT** contain the CRYPT option.

We now precisely specify the format of each suboption. In the sections that follow, all multi-byte values are encoded in big-endian format.

#### **4.3.1. The HELLO suboption**

The HELLO dataless suboption **MUST** only appear in a segment with the SYN control bit set. It is used by an active opener to indicate interest in using tcpcrypt for a connection, and by a passive opener to indicate that the passive opener wishes to play the "C" role.

The initial SYN segment from an active opener wishing to use tcpcrypt **MUST** contain a TCP CRYPT option with either an explicit or an implicit HELLO suboption.

After receiving a SYN segment with the HELLO suboption, a passive opener **MUST** respond in one of three ways:

- o To continue setting up tcpcrypt and play the "S" role, the passive opener **MUST** respond with a PKCONF suboption in the SYN-ACK segment and transition to S-MODE.
- o To continue setting up tcpcrypt and play the "C" role, the passive opener **MUST** respond with a HELLO suboption in the SYN-ACK segment and transition to HELLO-SENT.
- o To continue without tcpcrypt, the passive opener **MUST** respond with either no CRYPT option or the **DECLINE** suboption in the SYN-ACK segment, then transition to the **DISABLED** state.

An active opener receiving HELLO in a SYN-ACK segment must either transition to S-MODE and respond with a PKCONF suboption, or



transition to DISABLED.

There are three variants of the HELLO option used for application-level authentication, each encoded differently as shown in Table 4. The variants are: a plain HELLO where the application is not tcpcrypt-aware (but the kernel is), an "application supported" HELLO where the application is tcpcrypt-aware and is advertising the fact, and a "application mandatory" HELLO where the application requires the remote application to support tcpcrypt otherwise the connection MUST revert to plain TCP. The application supported HELLO can be used, for example, when implementing HTTP digest authentication - an application can check whether the peer's application is tcpcrypt aware and proceed to authenticate tcpcrypt's session ID over HTTP, otherwise reverting to standard HTTP digest authentication. The application mandatory HELLO can be used, for example, when implementing an SSL library that attempts tcpcrypt but reverts to SSL if the peer's SSL library does not support tcpcrypt. The application mandatory HELLO avoids double encrypting (SSL-over-tcpcrypt) since the connection will revert to plain TCP if the remote SSL library is not tcpcrypt-aware.

#### **4.3.2. The DECLINE suboption**

The DECLINE dataless suboption is sent by a host to indicate that the host will not enable tcpcrypt on a connection. If a host is in the DISABLED state or transitioning to the DISABLED state, and the host transmits a segment containing a CRYPT option, then the segment MUST contain the DECLINE suboption.

A passive opener SHOULD send a DECLINE suboption in response to a HELLO suboption or NEXTK1 suboption in a received SYN segment if it supports tcpcrypt but does not wish to engage in encryption for this particular session.

Implementations MUST NOT send segments containing the DECLINE suboption from the C-MODE or ENCRYPTING states.

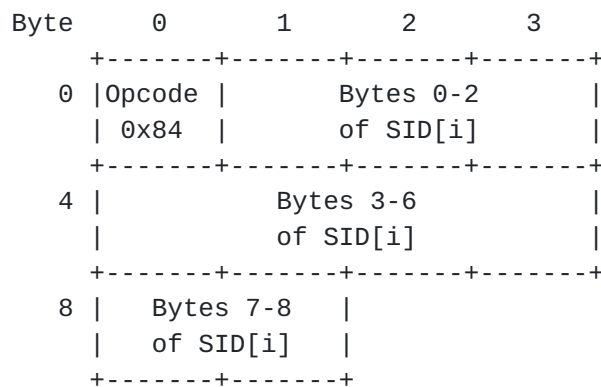
#### **4.3.3. The NEXTK1 and NEXTK2 suboptions**

The NEXTK1 suboption MUST only appear in a segment with the SYN control bit set and the ACK bit clear. It is used by the active opener to initiate a TCP session without the overhead of public key cryptography. The new session key is derived from a previously negotiated session secret, as described in [Section 3.8](#).

The suboption is always 10 bytes in length; the data contains the first nine bytes of SID[i] and is used to to start the session with session secret ss[i]. The format of the suboption is:







Format of the NEXTK1 suboption

The active opener MUST use the lowest value of *i* that has not already appeared in a NEXTK1 segment exchanged with the same host and for the same pre-session seed.

If the passive opener recognizes SID[i] and knows ss[i], it SHOULD respond with a segment containing the dataless NEXTK2 suboption. The NEXTK2 option MUST only appear in a segment with both the SYN and ACK bits set.

If the passive opener does not recognize SID[i], or SID[i] is not valid or has already been used, the passive opener SHOULD respond with a PKCONF or HELLO option and continue key negotiation as usual.

When two hosts have previously negotiated a tcpcrypt session, either host may use the NEXTK1 option regardless of which host was the active opener or played the "C" role in the previous session. However, a given host must either encrypt with *k\_cs* for all sessions derived from the same pre-session seed, or *k\_sc*. Thus, which keys a host uses to send segments depends only whether the host played the "C" or "S" role in the initial session that used ss[0]; it is not affected by which host was the active opener transmitting the SYN segment containing a NEXTK1 suboption.

A host MUST reject a NEXTK1 message if it has previously sent or received one with the same SID[i]. In the event that two hosts simultaneously send SYN segments to each other with the same SID[i], but the two segments are not part of a simultaneous open, both connections will have to revert to public key cryptography. To avoid this limitation, implementations MAY chose to implement session caching such that a given pre-session key is only good for either passive or active opens at the same host, not both.

In the case of simultaneous open, two hosts that simultaneously send SYN packets with NEXTK1 and the same SID[i] may establish a



connection, as described in [Section 4.2.1](#).

#### 4.3.4. The PKCONF suboption

The PKCONF option has one of the following two formats:

Byte	0	1	2	N
	+-----+	+-----+	+-----+ . . . +-----+	
	Opcode=	Length=	Algorithm	
	0x41	N	Specifiers	
	+-----+	+-----+	+-----+ . . . +-----+	

Byte	0	1	2	N
	+-----+	+-----+	+-----+ . . . +-----+	
	Opcode=	Length=	Algorithm	
	0x42	N	Specifiers	
	+-----+	+-----+	+-----+ . . . +-----+	

## Formats of the PKCONF suboption

The two are treated identically by `tcpcrypt`, except that opcode `0x42` (`PKCONF-app-support`) signals that the application on the sending host has set the `TCP_CRYPT_SUPPORT` option to non-zero, and hence the receiving host should return 1 for the `TCP_CRYPT_PEER_SUPPORT` socket option, as discussed in [Section 6](#).

The suboption data, whose length (N-2) must be divisible by 3, contains one or more 3-byte algorithm specifiers of the following form:

	0										1										2									
bit	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3						
	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-						
	0	Algorithm identifier																												
	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-						

Format of algorithm specifier within PKCONF. Fields starting with 1 are reserved for future use by algorithm identifiers longer than three bytes.

The algorithm identifier specifies a number of parameters, defined in Figure 3.

Hosts MUST implement OAEP+RSA3 and ECDHE-P256 and ECDHE-P521, but MAY by default disable certain algorithms and key sizes. In particular, implementations SHOULD disable larger RSA keys (Algorithm identifiers 0x102-0x103) by default unless such larger keys and ciphertexts can fit into a single TCP segment.



Servers demanding utmost performance SHOULD use RSA because the RSA encrypt operation is must faster than Diffie-Hellman operations, resulting in a higher connection rate.

Depending on the encoding of the PKCONF suboption (see Table 4), it can indicate whether "S's" application is tcpcrypt-aware or not. For the "C" role, the encoding of the HELLO suboption does this. This mechanism can be used for bootstrapping application-level authentication without requiring probing in upper layer protocols to check for support (which may not be possible). The application controls these encodings via the TCP\_CRYPT\_SUPPORT socket option.

#### 4.3.5. The UNKNOWN suboption

The UNKNOWN option has the following format:

Byte	0	1	2	N
+-----+-----+-----+-----+-----+-----+				
Opcode= Length=  N-2 unknown one-byte				
0x42   N   opcodes received				
+-----+-----+-----+-----+-----+-----+				

Format of the UNKNOWN suboption

This suboption is sent in response to an unknown suboption that has been received. The contents of the option are a complete list of the mandatory suboption opcodes from the received packet that were not understood. Note that this option is only sent once, in the next packet that the host sends. This means that it is reliable when sent in a SYN-ACK, but unreliable otherwise. Any mechanism sending new mandatory attributes must take this into account. If multiple packets, each containing unknown options, are received before an UNKNOWN suboption can be sent, the options list MUST contain the union of the two sets. The order of the opcode list is not significant.

If a host receives an unknown option, it SHOULD reply with the UNKNOWN suboption to notify the other side. If the host transitions to DISABLED as a result of the unknown option, then the host MUST also include the DECLINE suboption if it sends an UNKNOWN suboption (or more generally if it includes a CRYPT option in the next packet).

As a special case, if PKCONF (0x41) or INIT1 (0x06) appears in the unknown opcode list, it does not mean the sender does not understand the option (since these options are MANDATORY). Instead, it means the sender does not implement any of the algorithms specified in the PKCONF or INIT1 message. In either case, the segment must also contain a DECLINE suboption.



#### 4.3.6. The SYNC\_COOKIE and ACK\_COOKIE suboptions

A passive opener MAY include the SYNC\_COOKIE suboption in a segment with both the SYN and ACK flags set. SYNC\_COOKIE allows a server to be stateless until the TCP handshake has completed. It has the following format:

```

Byte      0      1      2      N
+-----+-----+-----+-----+
|Opcode=|Length=|  N-2 bytes of  |
| 0x43 |  N   |  opaque data   |
+-----+-----+-----+-----+

```

Format of the SYNC\_COOKIE suboption

The data is opaque as far as the protocol is concerned; it is entirely up to implementations how to make use of this suboption to hold state. It is OPTIONAL to send a SYNC\_COOKIE, but MANDATORY to understand and respond to them.

The ACK\_COOKIE suboption echoes the contents of a SYNC\_COOKIE; it MUST be sent in a packet acknowledging receipt of a packet containing a SYNC\_COOKIE, and MUST NOT be sent in any other packet. It has the following format:

```

Byte      0      1      2      N
+-----+-----+-----+-----+
|Opcode=|Length=|  N-2 bytes of  |
| 0x44 |  N   | SYNC_COOKIE data |
+-----+-----+-----+-----+

```

Format of the ACK\_COOKIE suboption

Servers that rely on suboption data from ACK\_COOKIE to reconstruct session state SHOULD embed a cryptographically strong message authentication code within the SYNC\_COOKIE data so as to be able to reject forged ACK\_COOKIE suboptions.

Though an implementation MUST NOT send a SYNC\_COOKIE in any context except the SYN-ACK packet returned by a passive opener, implementations SHOULD accept SYNC\_COOKIEs in other contexts and reply with the appropriate ACK\_COOKIE if possible.

#### 4.3.7. The SYNC\_REQ and SYNC\_OK suboptions

Many hosts implement TCP Keep-Alives [[RFC1122](#)] as an option for applications to ensure that the other end of a TCP connection still exists even when there is no data to be sent. A TCP Keep-Alive





segment carries a sequence number one prior to the beginning of the send window, and may carry one byte of "garbage" data. Such a segment causes the remote side to send an acknowledgment.

Unfortunately, Keep-Alive acknowledgments might not contain unique data. Hence, an old but cryptographically valid acknowledgment could be replayed by an attacker to prolong the existence of a session at one host after the other end of the connection no longer exists. (Such an attack might prevent a process with sensitive data from exiting, giving an attacker more time to compromise a host and extract the sensitive data.)

The TCP Timestamps Option (TSopt) [[RFC1323](#)] could alternatively have been used to make Keep-Alives unique. However, because some middleboxes change the value of TSopt in packets, tcpcrypt does not protect the contents of the TCP TSopt option. Hence the SYNC\_REQ and SYNC\_OK suboptions allow the cryptographically protected TCP CRYPT option to contain unique data.

The SYNC\_REQ suboption is always 5 bytes, and has the following format:

Byte	0	1	2	3	4
	+-----+-----+-----+-----+-----+				
	Opcode=		Clock		
	0x80				
	+-----+-----+-----+-----+-----+				

Format of the SYNC\_REQ suboption

Clock is a 32-bit non-decreasing value. A host MUST increment Clock at least once for every interval in which it sends a Keep-Alive. Implementations that support TSopt MAY chose to use the same value for Clock that they would put in the TSval field of the TCP TSopt. However, implementations SHOULD "fuzz" any system clocks used to avoid disclosing either when a host was last rebooted or at what rate the hardware clock drifts.

A host that receives a SYNC\_REQ suboption MUST reply with a SYNC\_OK suboption, which is always five bytes and has the following format:

Byte	0	1	2	3	4
	+-----+-----+-----+-----+-----+				
	Opcode=		Received-Clock		
	0x81				
	+-----+-----+-----+-----+-----+				

Format of the SYNC\_OK suboption



The value of Received-Clock depends on the values of the Clock fields in SYNC\_REQ messages a host has received. A host must set Received-Clock to a value at least as high as the most recently received Clock, but no higher than the highest Clock value received this session. If a host delays acknowledgment of multiple packets with SYNC\_REQ suboptions, it SHOULD send a single SYNC\_OK with Received-Clock set to the highest Clock in the packets it is acknowledging.

Because middleboxes sometimes "correct" inconsistent retransmissions, Keep-Alive segments with one byte of garbage data MUST use the same ciphertext byte as previously transmitted for that sequence number. Otherwise, a middlebox might change the byte back to its value in the original transmission, causing the cryptographic MAC to fail.

#### **4.3.8. The REKEY and REKEYSTREAM suboptions**

The REKEY and REKEYSTREAM suboptions are used to evolve encryption keys. Exactly one of the two options is valid for any given symmetric encryption algorithm. All algorithms in Table 6 use the REKEY option. REKEYSTREAM is reserved for future use should tcpcrypt evolve to support a stream cipher. We refer to a segment containing either option as a REKEY segment.

REKEY allows hosts to wipe from memory keys that could decrypt previously transmitted segments. It also allows the use of message authentication codes that are only secure up to a fixed number of messages. However, implementations MUST work in the presence of middleboxes that "correct" inconsistent data retransmissions. Hence, the value of ciphertext bytes must be the same in the original transmission and all retransmissions of a particular sequence number. This means a host MUST always use the same encryption key when transmitting or retransmitting the same range of sequence numbers. Re-keying only affects data transmitted in the future. Moreover, segments encrypted with different keysets MUST NOT be combined in retransmissions.

When switching keys, the REKEY suboption specifies which key set has been used to encrypt and integrity-protect the current segment. The suboption is always two bytes, and has the following format:

```

      Byte      0      1
      +-----+-----+
      |Opcode=|KeyLSB |
      | 0x82  |       |
      +-----+-----+

```

Format of the REKEY suboption



KeyLSB is the generation number of the keys used to encrypt and MAC the current segment, modulo 256. REKEYSTREAM is the same as REKEY but includes the TCP Sequence Number offset at which the key change took effect, for cases in which decryption requires knowing how many bytes have been encrypted thus far with a key. To interoperate with middleboxes that rewrite sequence numbers, offsets from the Initial Sequence Number (ISN) are used instead of TCP sequence numbers throughout tcpcrypt. The same occurs when dealing with acknowledgment numbers.

Byte	0	1	2	3	4	5
	+-----+-----+-----+-----+-----+-----+					
	Opcode= KeyLSB		Sequence Number Offset			
	0x83		from ISN			
	+-----+-----+-----+-----+-----+-----+					

Format of the REKEYSTREAM suboption

A host MAY use REKEY to increment the session key generation number beyond the highest generation it knows the other side to be using. We call this process `_initiating_` re-keying. When one host initiates re-keying, the other host MUST increment its key generation number to match, as described below (unless the other host has also simultaneously initiated re-keying).

A host MAY initiate re-keying by including a REKEY suboption in a `_syncable_` segment. A syncable segment is one that either contains data, or is acknowledgment-only but contains a SYNC\_REQ suboption with a fresh Clock value--i.e., higher than any Clock value it has previously transmitted. We say a syncable segment is `_synced_` when the transmitter knows the remote side has received it and all previous sequence numbers. A data segment is synced when the transmitter receives a cumulative acknowledgment for its sequence number (a Selective Acknowledgment [RFC2018] is insufficient). An acknowledgment-only segment is synced when the sender receives an acknowledgment for its sequence number and a SYNC\_OK with a high enough Clock number.

A host MUST NOT initiate re-keying with an acknowledgment-only segment that has either no SYNC\_REQ suboption or a SYNC\_REQ with an old Clock value, because such a segment is not syncable. A host MUST NOT initiate re-keying with any KeyLSB other than its current key number plus one modulo 256.

When a host receives a segment containing a REKEY suboption, it MUST proceed as follows:



1. The receiver computes RECEIVE-KEY-NUMBER to be the closest integer to its own transmit key number that also equals KeyLSB modulo 256. If no number is closest (because KeyLSB is exactly 128 away from the transmit number modulo 256), the receiver MUST discard the segment. If RECEIVE-KEY-NUMBER is negative, the receiver MUST also discard the segment.
2. The receiver MUST authenticate and decrypt the segment using the receive keys with generation number RECEIVE-KEY-NUMBER. The receiver MUST discard the packet as usual if the MAC is invalid.
3. If RECEIVE-KEY-NUMBER is greater than the receiver's current transmit key number, the receiver must wait to receive all sequence numbers prior to the REKEY segment's. Once it receives segments covering all these missing sequence numbers (if any), it MUST increase its transmit number to RECEIVE-KEY-NUMBER and transmit a REKEY suboption. If the receiver has gotten multiple REKEY segments with different KeyLSB values, it MUST increase its transmit key number to the highest RECEIVE-KEY-NUMBER of any segment for which it is not missing prior sequence numbers.

After sending a REKEY (whether initiating re-keying or just responding), a host MUST continue to send REKEY in all subsequent segments until at least one of the following holds:

- o One of the REKEY segments the host transmitted for its current transmit key number was syncable, and it has been synced.
- o The host receives a cumulative acknowledgment for one of its REKEY segments with the current transmit key number, and the cumulative acknowledgment is in a segment encrypted with the new key but not containing a REKEY suboption.

A host SHOULD erase old keys from memory once the above requirements are met.

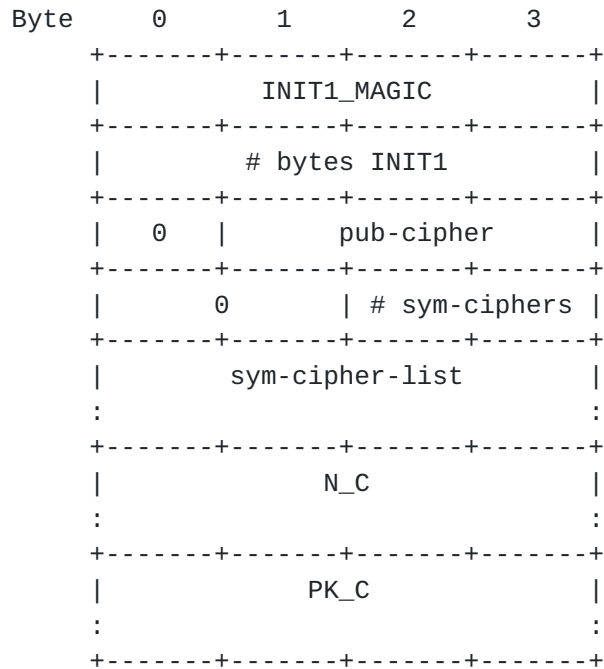
A host MUST NOT initiate re-keying if it initiated a re-keying less than 60 seconds ago and has not transmitted at least 1 Megabyte (increased its sequence number by 1,048,576) since the last re-keying. A host MUST NOT initiate re-keying if it has outstanding unacknowledged REKEY segments for key numbers that are 127 or more below the current key. A host SHOULD not initiate more than one concurrent re-key operation if it has no data to send.





#### 4.3.9. The INIT1 and INIT2 suboptions

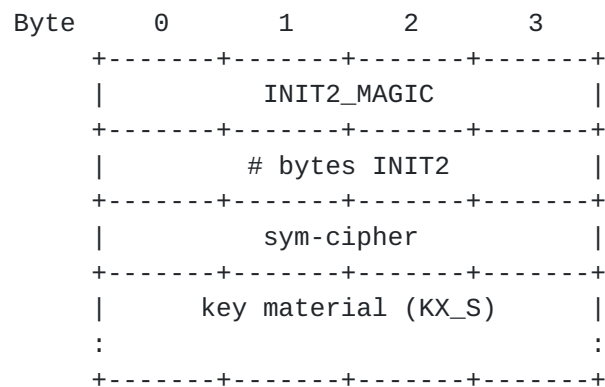
The INIT1 dataless suboption indicates that the Data portion of the TCP segment contains the following data structure:



The constant INIT1\_MAGIC is specified in Table 7. # bytes INIT1 specifies the length of the entire INIT1 structure, including the four-byte INIT1\_MAGIC that precedes the length. pub-cipher is a three-byte public key suite as specified in Figure 3, which specifies both the length of N\_C and the type of PK\_C. sym-cipher-list is a list of four-byte symmetric algorithm specifiers from Table 6. Of those listed, 0x00000100 (AES-128 / HMAC-SHA-256-128 / AES-128) is MANDATORY to implement, and the others OPTIONAL. # sym-ciphers specifies the number of four-byte entries in this list.

The INIT2 dataless suboption indicates that the Data portion of the TCP segment contains the following data structure:



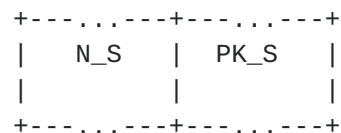


Format of the INIT2 suboption

Figure 2

The INIT2\_MAGIC constant is specified in Table 7. # bytes INIT2 is the total length of the INIT2 structure, including the 4-byte INIT2\_MAGIC constant preceding the length. sym-cipher specifies which entry of sym-cipher-list from the INIT1 message the host transmitting the INIT2 segment has selected.

The key material depends on the public key cipher selected, as described in [Section 3.4](#). When ECDHE is used, key material is encoded as follows:



The length of N\_S depends on pub-cipher and is given in Figure 3. PK\_S uses the rest of the message. When OAEP+-RSA exp3 is used, KX\_S is simply a ciphertext in big-endian format.

Hosts MUST set the TCP PSH control bits on INIT1 and INIT2 segments. Implementations MUST NOT set the TCP FIN control bit on INIT segments.

#### 4.4. The TCP MAC option

The MAC option is used to authenticate a TCP segment. Once a host has entered the encrypting phase for a session, the HOST must include a TCP MAC option in all segments it sends. Furthermore, once in the encrypting phase, a host MUST ignore any segments it receives that do not have a valid MAC option, except for segments with the RST bit set if the application has not requested cryptographic verification of RST segments.



The length of the MAC option is determined by the symmetric message authentication code selected. The format of the MAC option is:

Byte	0	1	2	N+1
	+-----+-----+-----+-----+			
	Kind	Len=	N-byte	
	OPT2	2+N	MAC	
	+-----+-----+-----+-----+			

Format of TCP MAC option

The MAC is the authentication tag as output from authenticated encryption. Apart from payload, two headers are included in the authenticated encryption process: a pseudo-header structure we call Assoc-Data, and an acknowledgment structure we call Up-Data. The format of Assoc-Data is as follows:

Byte	0	1	2	3
	+-----+-----+-----+-----+			
0	0x8000	length		
	+-----+-----+-----+-----+			
4	off	flags	window	
	+-----+-----+-----+-----+			
8	0x0000	urg		
	+-----+-----+-----+-----+			
12		seqno offset hi		
	+-----+-----+-----+-----+			
16		seqno offset		
	+-----+-----+-----+-----+			
20		options		
	+-----+-----+-----+-----+			

Assoc-Data data structure

The fields of Assoc-Data are defined as follows:

#### length

Total size of the TCP segment from the start of the TCP header to the end of the IP datagram.

#### off

Byte 12 of the TCP header (Data Offset)

#### flags

Byte 13 of the TCP header (Control Bits)



window

Bytes 14-15 of the TCP header (Window)

urg

Bytes 18-19 of the TCP header (Urgent Pointer)

seqno offset hi

Number of times the seqno offset field has wrapped from 0xffffffff  
-> 0

seqno offset

The low 32 bits of the sequence number offset (the Sequence Number  
in the TCP header - ISN)

options

These are bytes 20-off of the TCP header. However, where the  
TSOPT (8), Skeeter (16), Bubba (17), MD5 (19), TCP-AO (29), and  
MAC (OPT2) options appear, their contents (all but the kind and  
length bytes) are replaced with all zeroes.

The format of the Up-Data structure is as follows:

Byte	0	1	2	3
	+-----+-----+-----+-----+			
0		ackno offset hi		
	+-----+-----+-----+-----+			
4		ackno offset		
	+-----+-----+-----+-----+			

Up-Data data structure

The fields of Up-Data are defined as follows:

ackno offset hi The number of times ackno offset has wrapped from  
0xffffffff -> 0.

ackno offset The lower 32 bits of the acknowledgment number offset  
from the remote end's ISN (TCP's acknowledgment header - ISN  
received).

The two structures, Assoc-Data and Up-Data, are used in ASM mode to  
calculate the TCP MAC option. All multi-byte values are encoded in  
big-endian format.

## 5. Examples

To illustrate these suboptions, consider the following series of ways





in which a TCP connection may be established from host A to host B. We use notation S for SYN-only packet, SA for SYN-ACK packet, and A for packets with the ACK bit but not SYN bit. These examples are not normative.

### **5.1. Example 1: Normal handshake**

```
(1) A -> B: S  CRYPT<>
(2) B -> A: SA CRYPT<PKCONF<0x200,0x201>>
(3) A -> B: A  data<INIT1...>
(4) B -> A: A  data<INIT2...>
(5) A -> B: A  MAC<m> data<...>
```

(1) A indicates interest in using tcpcrypt. In (2), the server indicates willingness to use ECDHE with curves P256 and P521. Messages (3) and (4) complete the INIT1 and INIT2 key exchange messages described above, which are embedded in the data portion of the TCP segment. (5) From this point on, all messages are encrypted and their integrity protected by a MAC option.

### **5.2. Example 2: Normal handshake with SYN cookie**

```
(1) A -> B: S  CRYPT<>
(2) B -> A: SA CRYPT<PKCONF<0x200,0x201>, SYNCOKIE<val>>
(3) A -> B: A  CRYPT<ACKCOOKIE<val>> data<INIT1...>
(4) B -> A: A  data<INIT2...>
(5) B -> A: A  MAC<m> data<...>
```

Same as previous example, except the server sends the client a SYN cookie value, which the client must echo in (3). Here also the application level protocol begins by B transmitting data, while in the previous example, A was the first to transmit application-level data.

### **5.3. Example 3: tcpcrypt unsupported**

```
(1) A -> B: S  CRYPT<>
(2) B -> A: SA
(3) A -> A: A
```

(1) A indicates interest in using tcpcrypt. (2) B does not support tcpcrypt, or a middle box strips out the CRYPT TCP option. (3) the client completes a normal three-way handshake, and tcpcrypt is not enabled for the connection.



#### 5.4. Example 4: Reusing established state

```
(1) A -> B: S  CRYPT<NEXTK1<ID>>
(2) B -> A: SA CRYPT<NEXTK2>
(3) A -> A: A  MAC<m>
```

(1) A indicates interest in using tcpcrypt with a session key derived from an existing key, to avoid the use of public key cryptography for the new session. (2) B supports tcpcrypt, has ID in its session ID cache, and is willing to proceed with session caching. (3) the client completes tcpcrypt's handshake within TCP's three-way handshake and tcpcrypt is enabled for the connection.

#### 5.5. Example 5: Decline of state reuse

```
(1) A -> B: S  CRYPT<NEXTK1<ID>>
(2) B -> A: SA CRYPT<PKCONF<0x200,0x201>>
(3) A -> B: A  data<INIT1...>
(4) B -> A: A  data<INIT2...>
(5) A -> B: A  MAC<m> data<...>
```

A wishes to use a key derived from a previous session key, but B does not recognize the session ID or has flushed it from its cache. Therefore, session establishment proceeds as in the first connection, using public key cryptography to negotiate a new series of session secrets (ss[i] values).

#### 5.6. Example 6: Reversal of client and server roles

```
(1) A -> B: S  CRYPT<>
(2) B -> A: SA CRYPT<HELLO>
(3) A -> B: A  CRYPT<PKCONF<0x100>>
(4) B -> A: A  data<INIT1...>
(5) A -> B: A  data<INIT2...>
(6) B -> A: A  MAC<m> data<...>
```

Here the passive opener, B, wishes to play the role of the decryptor using RSA. By sending a HELLO suboption, B causes A to switch roles, so that now A is "S" and B plays the role of "C".

### 6. API extensions

The getsockopt call should have new options for IPPROTO\_TCP:

TCP\_CRYPT\_SESSID -> returns the session ID and MUST return an error if tcpcrypt is in not in the ENCRYPTING state (e.g., because it has transitioned to DISABLED).



TCP\_CRYPT\_CMODE -> returns 1 if the local host played the "C" role in session key negotiation, 0 otherwise.

TCP\_CRYPT\_CONF -> returns the four-byte authenticated encryption algorithm in use by the connection (as specified in Table 6). In addition, implementations SHOULD provide the three-byte public key cipher (Figure 3) initially used to negotiate the session keys, as well as the public key length for algorithms with variable key sizes (e.g., OAEP+-RSA3).

TCP\_CRYPT\_PEER\_SUPPORT -> returns 1 if the remote application is tcpcrypt-aware, as indicated by the remote host's use of a HELLO-app-support, HELLO-app-mandatory, or PKCONF-app-support CRYPT suboption (see Table 4).

The setsockopt call should have:

TCP\_CRYPT\_CACHE\_FLUSH -> setting this option to non-zero wipes cached session keys. Useful if application-level authentication discovers a man in the middle attack, to prevent the next connection from using NEXTK.

The following options should be readable and writable with getsockopt and setsockopt:

TCP\_CRYPT\_ENABLE -> one bit, enables or disables tcpcrypt extension on an unconnected (listening or new) socket.

TCP\_CRYPT\_RSTCHK -> one bit, means ignore unauthenticated RST packets for this connection when set to 1.

TCP\_CRYPT\_CMODE\_{DEFAULT,NEVER,ALWAYS}[\_NK] -> As described in [Section 4.2](#).

TCP\_CRYPT\_PKCONF -> set of allowed public key algorithms and CPRFs this host advertises in CRYPT PKCONF suboptions.

TCP\_CRYPT\_CCONF -> set of allowed symmetric ciphers and message authentication codes this host advertises in CRYPT INIT1 segments.

TCP\_CRYPT\_SCONF -> order of preference of symmetric ciphers.

TCP\_CRYPT\_SUPPORT -> set to 1 if the application is tcpcrypt-aware. set to 2 if the application is tcpcrypt-aware and wishes to enter the DISABLED state if the remote application is not tcpcrypt-aware. An active opener SHOULD set the default value of 0 for each new connection. A passive opener SHOULD use a default value of 0 for each port, but SHOULD inherit the value of the



listening socket for accepted connections. The behavior for each value is as follows:

When set to 0 The host MUST transition to the DISABLED state upon receiving a HELLO-app-mandatory option. The host MUST NOT send the HELLO-app-support, HELLO-app-mandatory, NEXTK2-app-support, or PKCONF-app-support options.

When set to 1 The "C" role host MUST use HELLO-app-support in place of the HELLO option, while the "S" role host MUST use the "PKCONF-app-support" in place of the "PKCONF" option. Either role must use NEXTK2-app-support in place of NEXTK2.

When set to 2 The "C" role host MUST use HELLO-app-mandatory option in place of the HELLO option, while the "S" role host MUST use "PKCONF-app-support" in place of the "PKCONF" option. Either role must use NEXTK2-app-support in place of NEXTK2. Either host MUST transition to DISABLED upon receipt of a HELLO or PKCONF option, but MUST proceed as usual in response to HELLO-app-support, HELLO-app-mandatory, and PKCONF-app-support.

Finally, system administrators must be able to set the following system-wide parameters:

- o Default TCP\_CRYPT\_ENABLE value
- o Default TCP\_CRYPT\_PKCONF value
- o Default TCP\_CRYPT\_CCONF value
- o Default TCP\_CRYPT\_SCONF value
- o Types, key lengths, and regeneration intervals of local host's short-lived public keys

The session ID can be used for end-to-end security. For instance, applications might sign the session ID with public keys to authenticate their ends of a connection. Because session IDs are not secret, servers can sign them in batches to amortize the cost of the signature over multiple connections. Alternatively, DSA signatures are cheaper to compute than to verify, so might be a good way for servers to authenticate themselves. A voice application could display the session ID on both parties' screens, and if they confirm by voice that they have the same ID, then the conversation is secure.





## 7. Acknowledgments

This work was funded by gifts from Intel (to Brad Karp) and from Google, by NSF award CNS-0716806 (A Clean-Slate Infrastructure for Information Flow Control), and by DARPA CRASH under contract #N66001-10-2-4088.

## 8. IANA Considerations

The following numbers need assignment by IANA:

- o New TCP option kind number for CRYPT
- o New TCP option kind number for MAC

A new registry entitled "tcpcrypt CRYPT suboptions" needs to be maintained by IANA as per the following table.

Symbol	Value
HELLO	0x01
HELLO-app-support	0x02
HELLO-app-mandatory	0x03
DECLINE	0x04
NEXTK2	0x05
NEXTK2-app-support	0x06
INIT1	0x07
INIT2	0x08
PKCONF	0x41
PKCONF-app-support	0x42
UNKNOWN	0x43
SYNCOOKIE	0x44
ACKCOOKIE	0x45
SYNC_REQ	0x80
SYNC_OK	0x81
REKEY	0x82
REKEYSTREAM	0x83
NEXTK1	0x84
IV	0x85

TCP CRYPT suboptions.

Table 5

A "tcpcrypt Algorithm Identifiers" registry needs to be maintained by



IANA as per the following table.

Algorithm Identifier	Value
Cipher: OAEP+-RSA with exponent 3	
min/max key size 2048/4096 bits ...	0x000100
min/max key size 4096/8192 bits ...	0x000101
min/max key size 8192/16384 bits ..	0x000102
min key size 16384 bits .....	0x000103
Extract: HKDF-Extract-SHA256	
CPRF: HKDF-Expand-SHA256	
N_C len: 32 bytes	
R_S len: 48 bytes	
K_LEN: 32 bytes	
Cipher: ECDHE-P256	0x000200
Extract: HKDF-Extract-SHA256	
CPRF: HKDF-Expand-SHA256	
N_C len: 32 bytes	
N_S len: 32 bytes	
K_LEN: 32 bytes	
Cipher: ECDHE-P521	0x000201
Extract: HKDF-Extract-SHA256	
CPRF: HKDF-Expand-SHA256	
N_C len: 32 bytes	
N_S len: 32 bytes	
K_LEN: 32 bytes	

TCP CRYPT algorithm identifiers.

Figure 3

A "tcpcrypt ASM mode parameter" registry needs to be maintained by IANA as per the following table.



+-----+	+-----+	+-----+	+-----+
Cipher	MAC	ACK MAC	Sym-cipher
+-----+	+-----+	+-----+	+-----+
AES-128	HMAC-SHA-256-128	AES-128	0x00000100
AES-128	Poly1305-AES-128	AES-128	0x00000200
AES-128	CMAC-AES-128	AES-128	0x00000300
+-----+	+-----+	+-----+	+-----+

ASM-mode parameters corresponding to 4-byte sym-cipher specifiers in INIT1 and INIT2 messages. ASM mode itself is specified in [Section 3.6](#). HMAC-SHA-256-128 is HMAC-SHA-256 with a 128-bit key and output truncated to 128 bits.

Table 6

## 9. Security Considerations

Tcpcrypt guarantees that no man-in-the-middle attacks occurred if Session IDs match on both ends of a connection, unless the attacker has broken the underlying cryptographic primitives (e.g., RSA). A proof has been published [[tcpcrypt](#)].

If the application performs no authentication, then there are no guarantees against active attackers. Session IDs can be logged on both ends and man-in-the-middle attacks can be detected after the fact by comparing Session IDs offline.

Session IDs are not confidential.

Tcpcrypt can be downgraded to regular TCP during the connection setup phase by removing any of the CRYPT options. The downgrade, and absence of protection, can of course be detected by the application as no Session ID will be returned.

By default tcpcrypt does not protect against RST packet injection. The connection must be configured with TCP\_CRYPT\_RSTCHK enabled to protect against malicious (unMACed) RSTs.

tcpcrypt uses short-lived keys to provide some forward secrecy. If a key is compromised all connections (new and cached) derived from that key will be compromised. The life of these keys should be kept to a minimum for stronger protection. A life of less than two minutes is recommended. Keys can be generated as frequently as practical, for example when servers have idle CPU time. For ECDHE-based key agreement, a new key can be chosen for each connection.

In the 4-way handshake, tcpcrypt does not have a key confirmation



step. Hence, an active attacker can cause a connection to hang, though this is possible even without tcpcrypt by altering sequence and ack numbers.

Attackers cannot force passive openers to move forward in their session caching chain without guessing the content of the NEXTK1 option, which will be hard without key knowledge.

## **10. References**

### **10.1. Normative References**

- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7, [RFC 793](#), September 1981.
- [RFC1122] Braden, R., "Requirements for Internet Hosts - Communication Layers", STD 3, [RFC 1122](#), October 1989.
- [RFC1323] Jacobson, V., Braden, B., and D. Borman, "TCP Extensions for High Performance", [RFC 1323](#), May 1992.
- [RFC2018] Mathis, M., Mahdavi, J., Floyd, S., and A. Romanow, "TCP Selective Acknowledgment Options", [RFC 2018](#), October 1996.
- [RFC2104] Krawczyk, H., Bellare, M., and R. Canetti, "HMAC: Keyed-Hashing for Message Authentication", [RFC 2104](#), February 1997.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC2437] Kaliski, B. and J. Staddon, "PKCS #1: RSA Cryptography Specifications Version 2.0", [RFC 2437](#), October 1998.
- [RFC5869] Krawczyk, H. and P. Eronen, "HMAC-based Extract-and-Expand Key Derivation Function (HKDF)", [RFC 5869](#), May 2010.
- [RFC6824] Ford, A., Raiciu, C., Handley, M., and O. Bonaventure, "TCP Extensions for Multipath Operation with Multiple Addresses", [RFC 6824](#), January 2013.

### **10.2. Informative References**

- [I-D.narten-iana-considerations-rfc2434bis]  
Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs",  
[draft-narten-iana-considerations-rfc2434bis-09](#) (work in





progress), March 2008.

[RFC3552] Rescorla, E. and B. Korver, "Guidelines for Writing RFC Text on Security Considerations", [BCP 72](#), [RFC 3552](#), July 2003.

[aggregate-macs]

Katz, J. and A. Lindell, "Aggregate Message Authentication Codes", Topics in Cryptology - CT-RSA , 2008.

[tcpcrypt]

Bittau, A., Hamburg, M., Handley, M., Mazieres, D., and D. Boneh, "The case for ubiquitous transport-level encryption", USENIX Security , 2010.

## [Appendix A](#). Protocol constant values

Value	Name
0x01	CONST_NEXTK
0x02	CONST_SESSID
0x03	CONST_REKEY
0x04	CONST_KEY_C
0x05	CONST_KEY_S
0x06	CONST_KEY_ENC
0x07	CONST_KEY_MAC
0x08	CONST_KEY_ACK
0x15101a0e	INIT1_MAGIC
0x097105e0	INIT2_MAGIC

Protocol constants.

Table 7



## Authors' Addresses

Andrea Bittau  
Stanford University  
Department of Computer Science  
353 Serra Mall, Room 288  
Stanford, CA 94305  
US

Phone: +1 650 723 8777  
Email: bittau@cs.stanford.edu

Dan Boneh  
Stanford University  
Department of Computer Science  
353 Serra Mall, Room 475  
Stanford, CA 94305  
US

Phone: +1 650 725 3897  
Email: dabo@cs.stanford.edu

Mike Hamburg  
Stanford University  
Department of Computer Science  
353 Serra Mall, Room 475  
Stanford, CA 94305  
US

Phone: +1 650 725 3897  
Email: mike@shiftleft.org

Mark Handley  
University College London  
Department of Computer Science  
University College London  
Gower St.  
London WC1E 6BT  
UK

Phone: +44 20 7679 7296  
Email: M.Handley@cs.ucl.ac.uk



David Mazieres  
Stanford University  
Department of Computer Science  
353 Serra Mall, Room 290  
Stanford, CA 94305  
US

Phone: +1 415 490 9451  
Email: dm@uun.org

Quinn Slack  
Stanford University  
Department of Computer Science  
353 Serra Mall, Room 288  
Stanford, CA 94305  
US

Phone: +1 650 723 8777  
Email: sqs@cs.stanford.edu

