

Article

Secure Service Proxy: A CoAP(s) Intermediary for a Securer and Smarter Web of Things

Floris Van den Abeele ^{1,*}, Ingrid Moerman ¹, Piet Demeester ¹ and Jeroen Hoebeke ¹

¹ Ghent University – imec, IDLab, Department of Information Technology; Technologiepark Zwijnaarde 15, B-9052 Ghent, Belgium; E-Mails: ingrid.moerman@ugent.be, piet.demeester@ugent.be, jeroen.hoebeke@ugent.be

* E-mail: floris.vandenabeele@ugent.be, Tel.: +32-933-14-900; Fax: +32-933-14-899

Academic Editor: name

Version July 6, 2017 submitted to Sensors

Abstract: As the IoT continues to grow over the coming years, resource-constrained devices and networks will see an increase in traffic as everything is connected in an open Web of Things. Performance and function enhancing features are difficult to provide in resource-constrained environments, but will gain importance if the WoT is to be scaled up successfully. For example, scalable open standards-based authentication and authorization will be important to manage access to the limited resources of constrained devices and networks. Additionally, features such as caching and virtualization may help further reduce the load on these constrained systems. This work presents the Secure Service Proxy (SSP): a constrained-network edge proxy with the goal of improving the performance and functionality of constrained RESTful environments. Our evaluations show that the proposed design reaches its goal by reducing the load on constrained devices while implementing a wide range of features as different adapters. Specifically, the results show that the SSP leads to significant savings in processing, network traffic, network delay and packet loss rates for constrained devices. As a result, the SSP helps to guarantee the proper operation of constrained networks as these networks form an ever-expanding Web of Things.

Keywords: CoAP; DTLS; REST; IoT; WoT; Proxy; 6LoWPAN; CoRE; LLN

1. Introduction

In recent years the Internet of Things (IoT) has increasingly become a hot topic in industry, academia, the do it yourself community and also consumers. Businesses are attracted by the new product opportunities and new sources of revenue that the IoT promises to bring. For example, a 2013 market report on IoT by Cisco Inc. predicts 14.4 trillion USD in created value for the “Internet of Everything” from 2013 to 2022 [?]. Academia are interested in the many new problems and issues that arise when deploying billions of devices on the Internet. These issues include big data analytics, energy efficient communications, large scale deployments, management of devices, communication protocols, security models, data privacy and many more. An introduction to the research aspect of the IoT is presented in [?]. Finally, consumers are drawn to the IoT because IoT products promise to bring improvements and novel services to their daily lives. Examples of IoT domains include smart home, smart health, smart transportation, smart factory, smart grid and many more [?].

As the Internet of Things continues to grow in scope and in size, the number of available technologies and platforms that promise to enable the IoT keeps increasing. As a family of such technologies, a complete protocol stack was standardized at the Internet Engineering Task Force (IETF) for use with constrained IoT devices in low power and lossy networks (LLNs) [?]. This suite of protocols defines the communication stack from the network layer up to the application

33 layer. In contrast to the popular alternative ZigBee [?], the IETF protocol stack gives the developer
34 more flexibility to model the network and the application to a specific use-case. For instance, with the
35 IPv6 Routing Protocol for Low-Power and Lossy Networks (RPL) [?] the routing can be tuned
36 by employing different objective functions that optimize routes according to the metrics that are
37 relevant to the use case (e.g. minimize hop count, maximize battery lifetime, etc.). Another example
38 of flexibility is found at the application layer, where the REST architecture followed by the CoAP
39 protocol allows developers to design their own RESTful resources and to model their behavior. In
40 terms of security, the IETF elected to standardize an end-to-end (E2E) architecture as it is a popular
41 choice on the unconstrained Web today. Therefore, the CoAP standard defines DTLS (i.e. Datagram
42 TLS) as its recommended security method.

43 Secure Sockets Layer (SSL) and later Transport Layer Security (TLS) have been around since the
44 end of the past century and have become very popular protocols for their roles in securing the WWW.
45 Today, (D)TLS has become a flexible protocol where endpoints can negotiate the type of security
46 and where a built-in extensions mechanism allows to add new features to the protocol without
47 touching the base specification. A comprehensive overview of the (D)TLS protocol is presented in
48 the Background section ???. Widespread adoption, a wide range of implementations, an open protocol
49 specification and a high level of interoperability are just a few of the benefits of the TLS protocol.
50 Nevertheless, one should be careful when deploying end-to-end security with DTLS in constrained
51 environments. This issue has been recognized by the IETF which has formulated guidance for
52 implementing and deploying DTLS in constrained environments in RFC 7925 [?]..

53 Despite the advantages offered by DTLS, E2E security has a number of disadvantages when
54 deployed as-is in LLNs. One issue with E2E security is that it completely blocks out any third
55 party (e.g. intermediate middleboxes) from taking part in the communication. In most traditional
56 Internet deployments this is a wanted property of E2E security, but in LLNs it stops intermediary
57 systems from providing services that can improve resource usage and performance of constrained
58 devices and networks. For example, caching of CoAP responses is not possible when E2E security
59 is applied between the CoAP client and the constrained CoAP server. A second disadvantage
60 of E2E security is that application-layer enhancements cannot be applied by middleboxes as all
61 communication is enciphered. Thus, access control, admittance control and other similar features
62 cannot be provided at the edge of the LLN. Another known problem with DTLS is its performance
63 in duty-cycled networks, which is common in multi-hop LLNs. Research [?] has shown that the
64 latency introduced by the DTLS handshake can become excessively large in multi-hop duty-cycled
65 networks (up to 50 seconds for 4 hops). Vučinić et al. also show that constrained nodes can only
66 store a limited number of DTLS sessions in their memory (e.g. max. 3 DTLS session for a WiSMote
67 node). As a result, nodes have to start dropping active DTLS sessions from memory which can
68 deteriorate battery lifetime and DTLS performance. Finally, end-to-end network addressing reduces
69 the effectiveness of 6LoWPAN compression. This is due to the fact that the IPv6 prefixes for nodes
70 situated on the Internet and the used UDP ports are difficult or impossible to compress on 6LoWPAN.
71 All these issues are covered in greater depth in the problem statement, cf. section ??.

72 The goal of this work is to overcome the issues identified with E2E security without losing
73 the benefits offered by such a widely used protocol as DTLS. To this end, we propose the “Secure
74 Service Proxy” (SSP). It is a reverse DTLS and CoAP proxy that provides a secure bridge between
75 clients on the Internet and constrained IoT devices in a low-power and lossy network. By employing
76 DTLS on both legs of the communication path, the resulting system can still enjoy most of the
77 benefits offered by the popularity of DTLS without suffering from the disadvantages of E2E security
78 specific to constrained environments (as identified in the previous paragraph). As the SSP operates
79 as a trusted entity in the network, it can also offer network services such as caching as well as
80 application-layer enhancements. For the latter, this paper employs the concept of node virtualization
81 where a constrained node has a virtual counterpart that resides on the proxy and that offers
82 additional functionality on behalf of the node. This virtualization concept is effective because

83 the SSP is deployed on hardware more powerful than the constrained nodes themselves. As a
84 result, node virtualization can offer new and complex functionality that is unfeasible to offer on the
85 constrained node itself. Examples include support for more complex modes of DTLS (e.g. public key
86 infrastructure and certificate-based suites), translating responses between content formats, offering
87 verbose semantic descriptions for the constrained node, storing large binary blobs (e.g. a picture of
88 the deployment area), keeping historical data, etc.

89 Our contributions in this paper are as follows. First, we identify and discuss a number of issues
90 with end-to-end security in constrained RESTful environments. We argue that these issues can be
91 overcome by a reverse proxy approach that splits the end-to-end security at the proxy. Secondly,
92 we design and implement such a reverse proxy. Apart from solving the E2E security issues, our
93 developed proxy can also offer additional functionality and services on behalf of the constrained
94 network and the constrained nodes. To our knowledge, this work is the first to study, design,
95 implement and evaluate a reverse proxy for use with end-to-end security in constrained RESTful
96 environments. Finally, by means of a real-world evaluation we show that our work can significantly
97 improve the operation of constrained networks by reducing power consumption, network latency
98 and network traffic.

99 The rest of this paper is structured as follows. First a brief overview of CoAP and DTLS is
100 presented in the next section. Using this overview, a number of issues with deploying CoAP and
101 DTLS in low-power and lossy networks is presented in section ???. This section also lists the research
102 goals of this work. In section ??, our approach to tackling these issues is presented together with the
103 design of the secure service proxy and an overview of the security risks related to breaking end-to-end
104 security. The secure service proxy is aligned to similar work in literature and the commercial world
105 in section ???. An extensive evaluation of our approach based on both simulations and a real-world
106 wireless sensor network testbed is presented in section ???. Section ?? presents the conclusions that
107 are drawn from this work.

108 2. Overview of CoAP and DTLS

109 2.1. The Constrained Application Protocol (CoAP)

110 RFC 7252 [?] states that the Constrained Application Protocol (CoAP) is a specialized
111 web transfer protocol for use with constrained nodes and constrained networks in the Internet of
112 Things. The protocol is designed for machine-to-machine (M2M) applications such as smart energy
113 and building automation. The main design considerations for CoAP include simplicity, very low
114 overhead, easy translation to and from HTTP and support for multicast.

115 In CoAP, constrained devices that host applications structure their data and actions as RESTful
116 web services, also called CoAP resources. CoAP clients send requests to resources in order to retrieve
117 and store data or trigger actions. CoAP defines the same request methods as HTTP: GET, PUT, POST
118 and DELETE. They are used respectively for retrieving data, storing data, toggling an action and
119 removing data. CoAP chose UDP as its transport protocol due to the lightweight nature of UDP (TCP
120 was deemed too verbose due to its connections and too complex to implement in constrained
121 devices). Therefore, CoAP includes a simple reliability layer and deduplication mechanism in order
122 to compensate for the minimalistic nature of UDP. In order to minimize overhead, CoAP uses a binary
123 format for encoding message options in the headers of CoAP requests and responses. As a result the
124 CoAP message size is significantly reduced when compared to a non-binary encoded protocol such as
125 HTTP [?], which is important in LLNs where message sizes are typically small and communication
126 is expensive for battery-powered devices.

127 An illustration of a typical CoAP request/response exchange is shown in figure ??, where a
128 client (a ventilation unit) retrieves a temperature resource on a CoAP server. The first elements of the
129 CoAP header are the 2-bit protocol version (RFC 7252 standardizes version 1) and the 2-bit message
130 type. By sending a Confirmable message, a sender can ask a receiver to acknowledge the reception of

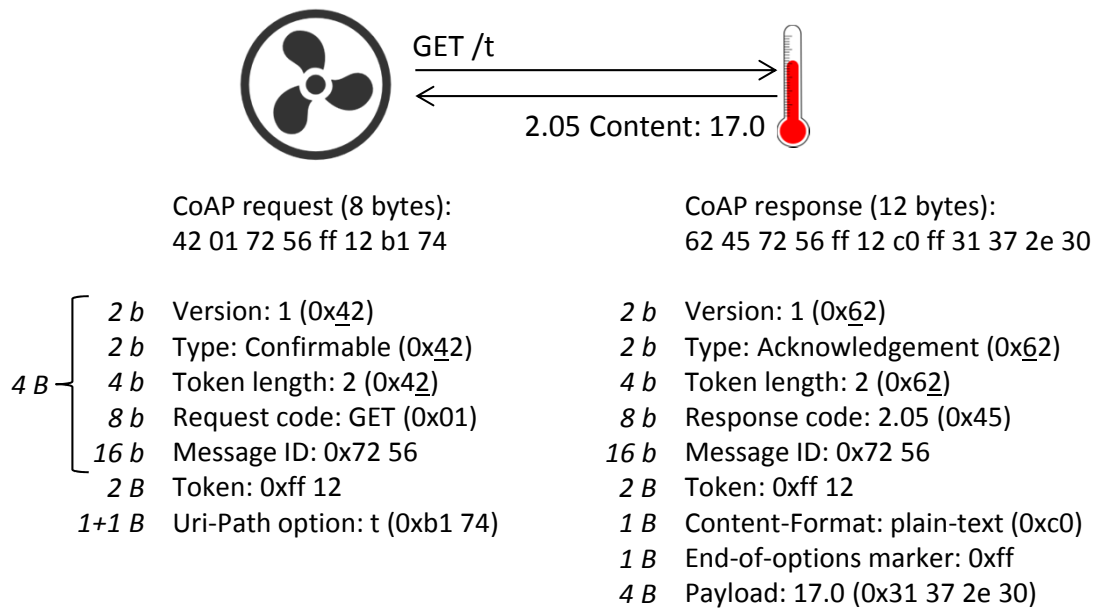


Figure 1. Anatomy of a typical CoAP request and response

131 a message. This is reflected in the message type of the response which is an acknowledgment. In most
 132 cases (like here) the response message is actually piggy-backed on the acknowledgment message in
 133 order to reduce the number of messages. The 4-bit token length comes after the message type in the
 134 CoAP header and it represents the length of the optional message token in bytes. The next element
 135 of the CoAP header is the 8-bit message code, which consists of a 3-bit class and a 5-bit subfield.
 136 Requests codes are class 0 codes (e.g. GET is code 0.01) and successful response codes are class 2
 137 codes (e.g. Content is code 2.05). The final part of the fixed 4-byte CoAP header is the two byte
 138 message ID. It is used for deduplication and for confirmable messages where acknowledgments echo
 139 the message ID of the CON message. The token is used to match a response with a request and can
 140 vary in length between 0 and 8 bytes. After the token come the header options and the payload (if
 141 any). In CoAP header options are assigned unique numbers by IANA and are delta encoded in CoAP
 142 messages in order to reduce their encoding size. Every option encoding contains the delta of the
 143 option number (relative to the preceding option), the size of the value of the option (in bytes) and the
 144 value of the option. Finally, the options and the payload are separated by an end-of-options marker
 145 (0xff).

146 CoAP observe [?] is a CoAP protocol extension that is important for this work. When a client is
 147 observing a REST resource on a CoAP server, the server will notify the client of state changes for that
 148 resource. This frees the client from polling the resource on the server, which can save resources in
 149 LLNs when changes in resource state occur rarely. RFC 7641 [?] also states that intermediaries must
 150 aggregate observe registrations: "If two or more clients have registered their interest in a resource
 151 with an intermediary, the intermediary MUST register itself only once with the next hop and fan
 152 out the notifications it receives to all registered clients. This relieves the next hop from sending
 153 the same notifications multiple times and thus enables scalability.". Apart from enabling scalability,
 154 aggregation also saves resources.

155 2.2. Datagram Transport Layer Security (DTLS)

156 For security, CoAP standardized end-to-end security and DTLS as its default security mechanism
 157 and protocol respectively. The primary motivation for preferring transport-layer security over
 158 alternatives such as object security and network layer security, is the popularity of TLS on the
 159 conventional Web. Datagram TLS is by design very similar to the TLS protocol and the specification

160 of DTLS is largely written as a set of changes to the TLS specification [?]. However there are some key
 161 differences as DTLS runs over an unreliable datagram transport while TLS runs over the reliable TCP
 162 transport. Therefore, DTLS must cope with the reliable and ordered delivery of packets as available in
 163 TLS. To this end, DTLS introduces a simple timeout and retransmission scheme and adds an explicit
 164 sequence number to the Record Protocol (versus an implicit number as available via TCP in TLS).
 165 Another difference is that stream ciphers must not be used with DTLS. DTLS also enhanced the
 166 handshake protocol with a stateless cookie exchange for Denial of Service resistance. By forcing DTLS
 167 clients to echo the cookie in their second handshake message, malicious clients (e.g. those spoofing IP
 168 addresses) can be rooted out and a DTLS server can avoid wasting resources on bogus handshakes.

169 DTLS is a session-based protocol in that DTLS endpoints have to set up a session when
 170 they want to communicate securely. Negotiation of the security parameters for the session and
 171 peer authentication are both performed during the handshake phase of the protocol. After the
 172 handshake phase, both endpoints can exchange data with guarantees for confidentiality, endpoint
 173 authentication and integrity of the data. To this end, DTLS employs symmetric cryptography for
 174 data encryption according to an encryption algorithm and encryption keys that are agreed during the
 175 handshake. DTLS also guarantees message integrity by means of hash-based message authentication
 176 codes (HMAC). Sessions are typically negotiated on an ad-hoc basis, although long-term sessions and
 177 resumption of established sessions are possible in DTLS.

178 TLS introduces the concept of cipher suites, these are named combinations of the authentication
 179 and key exchange algorithm, the cipher and key length, the cipher mode of operation, the hash
 180 algorithm for integrity protection and the hash algorithm for use with pseudorandom functions.

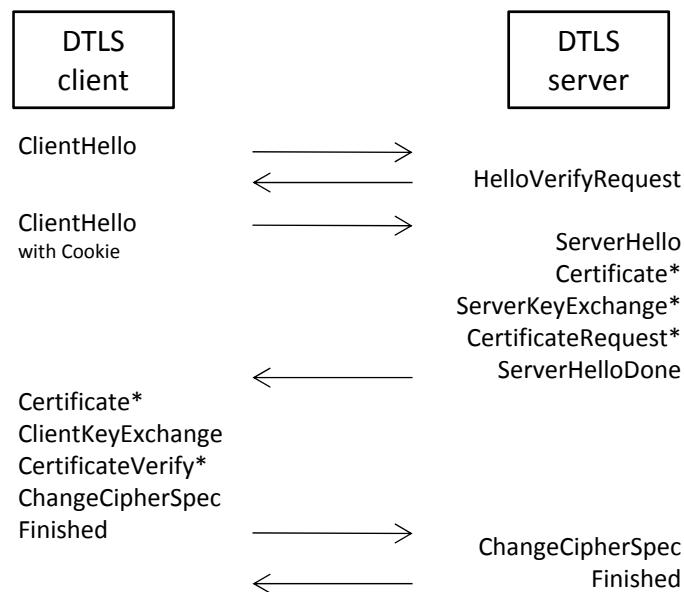


Figure 2. The full DTLS handshake

181 The DTLS handshake is shown in figure ???. In order to reduce the number of network packets,
 182 multiple DTLS messages can be grouped into a single flight of messages. In the figure the horizontal
 183 arrows correspond to the different message flights. The DTLS client initiates the handshake with
 184 the `ClientHello` message, to which the server replies with a `HelloVerifyRequest` message. The
 185 `HelloVerifyRequest` message contains the stateless cookie for DoS mitigation and must be echoed by
 186 the client in its second `ClientHello` message. After the server has verified the cookie, it responds
 187 with the `ServerHello` message. The hello messages are used to establish security enhancement capabilities
 188 between client and server [?]. They establish the following attributes: protocol version, session

189 ID (used in session resumption), cipher suite and compression method. Additionally, two random
190 values are generated and exchanged: one for the client and one for the server.

191 The messages of the remainder of the handshake depend on the negotiated security
192 enhancement capabilities. In the figure, messages marked with an asterisk (*) are optional
193 or situation-dependent messages. The figure shows the message flow for a certificate-based
194 cipher suite where the server replies with Certificate, ServerKeyExchange, CertificateRequest and
195 ServerHelloDone messages. If the cipher suite requires the server to authenticate itself, then the server
196 sends its X.509 certificate in a Certificate message. In cases where the key exchange does not use the
197 server certificate, the server may send a ServerKeyExchange message. For example in Pre-Shared Key
198 cipher suites (PSK suites are discussed later), the server may send a hint in the ServerKeyExchange
199 message to help the client in selecting which PSK identity to use. Additionally, the server may also
200 send a CertificateRequest message to request a certificate from the client. Finally, a ServerHelloDone
201 message is sent by the server to indicate that the hello-message phase of the handshake is complete.

202 If the server requested a certificate, the client must provide one in its Certificate message. Next,
203 the client sends a ClientKeyExchange message, the contents of which depend on the chosen key
204 exchange algorithm. In the case of RSA for example, the client chooses a secret and encrypts it
205 with the public key from the certificate of the server and sends the result in the ClientKeyExchange
206 message. Together with the Certificate and ServerKeyExchange messages of the server, the client's
207 Certificate and ClientKeyExchange messages are used for the key exchange. The CertificateVerify
208 message allows the client to prove the possession of the private key in the certificate. In the case of
209 Pre-Shared key cipher suites, the key exchange of the client consists of a ClientKeyExchange message
210 which contains the identity of the chosen PSK.

211 Next, the client sends a ChangeCipherSpec message which signals that the client has switched to
212 the negotiated cipher spec. The client then immediately sends the Finished message which contains
213 a hash of the shared secret and all handshake messages. The server must verify the contents of
214 the Finished message in order to detect any tampering to the handshake messages. The Finished
215 message also proves that the client knows the correct shared secret (i.e. the pre-master secret) and
216 any subsequent keying material (master secret, encryption keys and MAC keys) is generated from
217 this pre-master secret. After the server has sent its own ChangeCipherSpec and Finished messages
218 and the client has successfully verified the Finished message, the handshake is completed and secure
219 communication of application data can start.

220 2.3. DTLS in constrained environments

221 There are a number of additional protocol features that are applicable to DTLS in constrained
222 environments and these are discussed in this subsection. RFC 5116 [?] introduced authenticated
223 encryption with associated data (AEAD) to TLS which enables the use of cipher suites that use
224 the same cipher for confidentiality, authenticity and integrity protection. Particularly in constrained
225 environments, AEAD provides the benefit of more compact implementations as only one cipher has
226 to be implemented.

227 RFC 6655 [?] defines multiple of such compact cipher suites that use the widespread AES
228 cipher in the Counter with Cipher Block Chaining - Message Authentication Code (CBC-MAC) Mode
229 (CCM). AES is a popular choice in constrained environments as it is often accelerated in hardware in
230 modern IoT systems (e.g. the TI CC2538 SoC has an AES accelerator on the same die as the ARM-M3
231 CPU). Note that the AEAD construct is only supported from version 1.2 of the DTLS protocol.

232 RFC 4279 [?] introduces Pre-Shared Key (PSK) cipher suites for TLS. These cipher suites are
233 interesting for constrained devices, as the size of the key exchange is minimal: typically only a PSK
234 identifier in the Client Key Exchange is exchanged. Of course key management is an important issue
235 in this case, as common cryptography practice dictates that a unique PSK should be allocated for
236 every peer. The 'TLS_PSK_WITH_AES_128_CCM_8' cipher suite combines the benefits of PSKs and
237 AES-CCM in that only one cipher is needed (AES) and the key exchange is minimal. This cipher suite

238 is also the mandatory-to-implement PSK cipher suite for DTLS in the CoAP RFC [?]. Furthermore,
239 this suite uses just an 8 byte authentication tag (as opposed to a 16 byte tag) which is more suitable in
240 networks where bandwidth is constrained and messages sizes may be small.

241 RFC 7250 [?] introduces a new certificate type and two TLS extensions for exchanging raw
242 public keys (RPKs) in DTLS. In this case a peer has an asymmetric key pair but it does not have an
243 X.509 certificate, this asymmetric key pair is the RPK. This extension allows raw public key to be used
244 for authentication, which is beneficial in constrained environments as RPKs are smaller in size than
245 X.509 certificates. Additionally the resulting key exchange is therefore smaller as well. Of course, the
246 scalability benefits of a Public Key infrastructure (PKI) are lost when using RPKs.

247 Finally, RFC 7251 [?] describes the use of AES-CMM elliptic curve cryptography (ECC) cipher
248 suites in DTLS. This type of cipher suites uses the AEAD mechanism to provide confidentiality,
249 authenticity and integrity of application data with just AES, while using Ephemeral Elliptic Curve
250 Diffie-Hellman (ECDHE) as their key exchange and peer authentication mechanisms. ECC is
251 attractive for constrained environments as its smaller key sizes result in savings for power, memory,
252 bandwidth and computational cost [?]. For example, a 256 to 383 bit ECC key is considered
253 comparable in strength to a 3072 bit RSA key by NIST [?]. CoAP mandates the use of the
254 'TLS_ECDHE_ECDSA_WITH_AES_128_CCM_8' cipher suite for X.509 certificates in constrained
255 environments. This cipher suite uses the secp256r1 or NIST P-256 elliptic curve.

256 3. Problem statement and research goals

257 When securing communications in LLNs via end-to-end security with DTLS, one should be
258 mindful of a number of potential issues and pitfalls. Some of these issues arise due to the limitations
259 of the constrained devices that secure the communications. For example in end-to-end security, there
260 is a considerable difference between constrained devices (and their protocols) and powerful Internet
261 hosts (and their protocols) in terms of available resources and design. A second potential issue stems
262 from the DTLS protocol itself, namely the large overhead of the DTLS handshake can be an issue of
263 concern in constrained networks. A third group of issues is related to securing the LLN itself and is
264 the result of deploying end-to-end security in LLNs. Apart from these issues related to end-to-end
265 security in LLNs, there is also the problem of the limited amount of application layer functionality
266 that can be provided by constrained IoT devices. In a world as heterogeneous as the IoT there
267 exists a need for protocol translation, data format mapping, semantic descriptions and many other
268 features that improve the interoperability with IoT devices. Similarly, network access to constrained
269 nodes and LLNs should be as efficient as possible by supporting caching of information, efficient
270 discovery and network edge filtering. These types of functionality are too complex and in some cases
271 impossible for implementation on a constrained device. Clearly, an approach that does not burden
272 the constrained device is needed in this case. The remainder of this section discusses these various
273 issues and problems in more detail.

274 3.1. End-to-end security in LLNs

275 Constrained devices with a limited power source (e.g. battery powered or energy scavenging
276 devices) should take care to avoid excessive network communications in order not to preemptively
277 deplete the power source. Similarly, constrained networks where the available throughput is in the
278 order of a few kbps, should minimize the amount of network communications to avoid congestion.
279 Therefore chatty or verbose security protocols that communicate excessive amounts of information
280 should be avoided in these situations. As DTLS employs UDP instead of TCP as its transport protocol,
281 it avoids the TCP handshake which reduces the number of messages exchanged between DTLS clients
282 and servers. However, some options supported by DTLS, as presented in the previous section, may
283 lead to large amounts network communications. Specifically, certificate-based cipher suites involve
284 sending the certificate of the DTLS server (and peer, depending on the security needs) over the
285 network. These certificates are generally large (i.e. thousand bytes and more) and therefore their

286 network communication can be problematic when communication has a large impact on the power
287 source or the network. As a result, these types of devices are unable to offer authentication based on
288 PKI certificates. While raw public keys are significantly more compact than X.509 certificates, they do
289 not offer the same benefits in terms of authentication and scalability.

290 For devices with limited computational power (e.g. low cost embedded systems) certain
291 cryptographic primitives may prove too complex for computation by the low cost microcontroller.
292 While hardware acceleration may help to alleviate this issue, it can be an expensive option and might
293 only be available for certain primitives: e.g. AES is often accelerated in hardware, while others are not.
294 Specifically, public-key cryptography methods (e.g. based on large integer factorization or discrete
295 logarithm problems) and key agreement schemes (such as (EC)DH) may be too taxing for constrained
296 microcontrollers. Therefore, the set of cryptographic functions that can be offered by such low cost
297 embedded systems excludes a number of common cryptographic primitives and is typically limited
298 to what can be achieved by symmetric-key cryptography.

299 Another important limitation in constrained environments is the low amount of available
300 memory (i.e. both volatile and non-volatile memory). For example, according to IETF RFC 7228 [?
301] Class 1 constrained devices have around 10KiB of RAM and 100KiB of ROM memory. Such
302 a small amount of memory must accommodate an entire networking stack, adequate security
303 mechanisms, peripheral control, the application itself and various other subsystems. This forces
304 a device manufacturer to limit the amount of software that will ship with the device by carefully
305 selecting what is needed. One consequence is that it is impossible for these devices to support a wide
306 range of DTLS extensions and cipher suites (e.g. only one suite might be supported). This also means
307 that verbose operations such as checking certificate revocation lists or performing OCSP [?
308] checks typically can not be supported.

309 Powerful Internet hosts on the other hand may expect constrained devices to support security
310 features similar to those found on the conventional Internet (e.g. with strong authentication and key
311 agreement schemes). As constrained devices can not support these features (see above), an alternative
312 is to consider third party systems (e.g. middleboxes or off-path systems) that offer such features on
313 behalf of constrained devices. However, in this case a big issue with conventional end-to-end security
314 is that as the connection is secured end-to-end, a third party is excluded from the communication.
315 Thus, an important question addressed by this work is how third parties can take part in securing (but
316 also optimizing, see later) communications with constrained devices in order to bridge the gap with
317 powerful Internet hosts.

318 While DTLS can avoid the TCP handshake, it still has to perform its own handshaking
319 mechanism in order to negotiate key exchange and authentication methods. The overhead of this
320 handshake in terms of delay or amount of network traffic can be problematic for some types of
321 constrained nodes and networks. Specifically, previous research has shown that in duty-cycled
322 multi-hop networks the delay introduced by the DTLS handshake can run up to fifty second [?
323] for 4 wireless hops. The authors also correctly conclude that the memory for storing DTLS session
324 state on constrained nodes is typically limited to a handful of nodes for Class 1 devices. Additionally,
325 other research [?] has shown that ephemeral DTLS sessions with constrained devices should be
326 avoided as their energy expenditure is up to 60% higher when compared to a single DTLS session
327 with a long lifetime. Therefore, one goal of this work is to limit the impact of the DTLS handshake
328 on delay and energy expenditure, while supporting more than just a handful of simultaneous DTLS
329 sessions per constrained device.

330 The third group of issues stems from naively deploying end-to-end security in (multi hop) low
331 power and lossy networks (LLNs) and from allowing unmonitored access to LLNs to malicious
332 users. In these networks resources are sparse (see above) and care should be taken in order to
333 avoid unwanted depletion of these resources by denial-of-service (DoS) attacks. For example, by
334 repeatedly opening and closing DTLS sessions a malicious user can significantly reduce the lifetime
335 of a battery-powered device. A malicious user could also send large datagrams to the LLN, which

336 will trigger fragmentation that can exhaust the allocated network buffers in the LLNs. Most of
337 these resource-depletion threats can be mitigated by monitoring and restricting access to the LLN
338 at the edge of the network, where an unconstrained firewall or gateway system resides. However,
339 end-to-end security encumbers such systems from authenticating parties (as constrained devices can
340 not support strong authentication) and therefore restricting access to authorized parties. Here, this
341 work will study how end-to-end security can be reconciled with the need for traffic filtering at the
342 edge of the network and the need for strong authentication.

343 3.2. *Complex application features in LLNs*

344 Apart from security issues, there is another important category of problems that relate to
345 functionality at the application layer for constrained devices which is targeted by this work. Firstly,
346 the same constraints that prohibit offering extensive security features also apply to implementing
347 application features on the constrained device. This is one of the reasons why the IETF has
348 standardized special purpose protocols and data formats for use in constrained environments (e.g.
349 CoAP and CoRE link format [?]). However, traditional Internet hosts do not always implement
350 these protocols and data formats. In these cases a protocol and data format translation should occur
351 that enables the Internet host to communicate with the constrained device (e.g. an HTTP/CoAP
352 proxy and a JSON/CLF mapper). Such a translation has to be performed by an unconstrained third
353 party system (e.g. gateway). Secondly, some types of functionality can be ineffective when they
354 are offered on the constrained device. An example is caching the responses of a constrained server
355 on the device itself which will not save any network traffic. A second example is the aggregation
356 of observe relationships by intermediaries, clearly this has to be offered on an intermediary and
357 not on a constrained node in order to have any effect. Note that conventional end-to-end security
358 does not allow for response caching or observe aggregation, as all traffic passing at an intermediary
359 is encrypted. Thirdly, some functionality can be inefficient when they are implemented on the
360 constrained device. An example is storing verbose semantic descriptions on a constrained device
361 which will lead to significant amounts of network traffic every time these descriptions are requested.
362 Another example of functionality that is inefficient to offer on constrained devices is access control.
363 Typically the LLN will have already spent a significant amount of resources delivering the request
364 to its destination where it will end up being discarded. Clearly, discarding this request before the
365 network has wasted its resources is more efficient. For these cases, this work will study how third
366 party systems can support and optimize the operations of constrained devices and LLNs.

367 3.3. *Problem statement: illustration in a smart building use case*

368 Figure ?? shows a smart building scenario that illustrates the problems targeted by this work. In a
369 smart building most of the building services can be monitored and controlled over the Internet. Such
370 services include for example the management of doors, lighting, climate control (e.g. AC), elevators
371 and the monitoring of presence in certain areas. Smart buildings, such as offices and public buildings,
372 typically have a large variety of users: visitors, cleaning staff, technicians, employees, etc. Similarly
373 there are also a number of computer systems that interact with the smart building: e.g. systems for
374 HVAC, surveillance, facility management, etc. Each of these actors access the services offered by the
375 building according to specific access control rules that depend on the role and or identify of the actor.
376 E.g. the HVAC system can control the air conditioning units, but can not control the doors. However,
377 the HVAC system might be allowed to monitor the status of a door adjacent of an AC unit without
378 being able to (un)lock it. Considering the limited resources of constrained devices (see above),
379 managing and enforcing which actions an actor is allowed to perform depending on their role or
380 identity quickly becomes too complex for the constrained devices. Furthermore, as most constrained
381 devices only support PSK-based authentication such a system would require management of shared
382 secret keys between every two actors. Limitations on the LLN and the constrained devices also
383 prohibit these devices from offering protocols and data formats that are common to the unconstrained

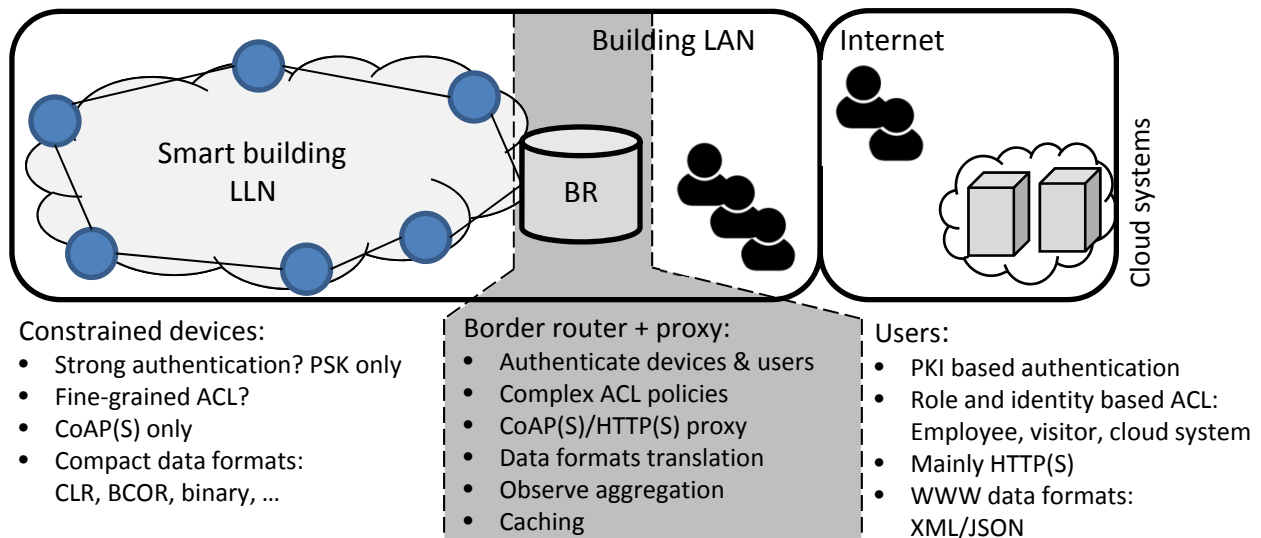


Figure 3. In a smart building scenario there is a wide variety of different users. Constrained devices are unable to offer all necessary security and application features to cater to these users. In the approach followed by this work, unconstrained systems (e.g. border routers) assist by offering these missing features.

384 actors, such as HTTP(S) and XML/JSON. The gray center of the figure already hints at our approach
 385 detailed in the next section: a proxy offers many of the missing features on behalf of the constrained
 386 devices.

387 Finally, one might question why this work relies on end-to-end security via DTLS at all, when
 388 there appear to be many problems in constrained environments according to the discussion above.
 389 Our main motivations for doing so is that DTLS is a proven (and secure) standard, is widely available,
 390 is commonly used on the Web and is standardized for use with CoAP. Alternatives to DTLS, are either
 391 proprietary or still in the process of standardization (e.g. OSCOAP [?]), not applicable to constrained
 392 environments (e.g. network layer security) or can not provide the same level of security as DTLS (e.g.
 393 physical layer security). Object security specifically can be considered complementary to transport
 394 layer security and while it is not considered in this work, it can be combined with the work presented
 395 here (if feasible given the constrained environments under consideration). The related work section
 396 discusses object security in greater detail. While literature shows that lightweight network security
 397 is feasible in constrained environments (e.g. compressed IPsec [?]), it is not considered in this work
 398 because CoAP standardized end-to-end security over DTLS as its security mechanism.

399 4. The Secure Service Proxy

400 The approach followed in this work allocates one reverse CoAP(s) proxy per constrained device.
 401 The CoAP specification [?] defines a reverse proxy as “an endpoint that stands in for one or
 402 more other server(s) and satisfies requests on behalf of these, doing any necessary translations” and
 403 it also states that “The client may not be aware that it is communicating with a reverse-proxy; a
 404 reverse-proxy receives requests as if it were the origin server for the target resource.” The reverse
 405 proxy approach enables splitting the end-to-end communication between a constrained device and
 406 its client at the proxy with no need for any additional configuration on the client (as mentioned in the
 407 CoAP specification). While the resulting communication is no longer end-to-end, indeed the proxy
 408 will share DTLS security contexts with both parties and will translate CoAP messages, the resulting
 409 system has a lot of benefits and is able to overcome all of the issues that are discussed in the previous
 410 section. Additionally, our reverse proxy approach implements a virtual device for every constrained
 411 device. This enables the reverse proxy to extend a constrained device (beyond only proxying) by

412 hosting functionality on the corresponding virtual device. Finally, by enabling the reverse proxy to
413 be deployed on any system (see design), it is not restricted by the limitations common to constrained
414 IoT devices. In the next subsections we argue that the benefits of this approach far outweigh the
415 downsides of splitting the end-to-end communication and we present our design for such a reverse
416 proxy.

417 *4.1. Motivation of approach*

418 Our motivation for following a reverse proxy approach consists of two facets: one for the
419 security related aspects of constrained devices and LLNs and one for the application layer related
420 aspects of constrained devices. In terms of security, the reverse proxy approach allows to setup two
421 sorts of DTLS sessions: “lightweight” sessions between the constrained devices and their reverse
422 proxy and fully featured sessions between the proxy and the clients of the devices. The lightweight
423 sessions employ security primitives that are known to the constrained devices (e.g. pre-shared
424 keys for authentication and key exchange), while the fully featured sessions can use conventional
425 security methods that are known to the clients: e.g. certificates for strong authentication and (Elliptic
426 Curve) Diffie-Hellman (ECDH) for the key exchange (including ephemeral key exchanges if perfect
427 forward secrecy is required). Additionally, the reverse proxy can be configured to maintain one
428 long-term session with the constrained device while simultaneously keeping active sessions with
429 multiple clients. This allows to overcome the small session pool at the constrained devices (due to
430 its limited memory, see above) as well as limit the total number of handshakes performed by the
431 constrained device during its lifetime. As a result, the impact of the DTLS handshake on the LLN
432 and the communication in terms of e.g traffic and communication latency is lowered. Finally, the
433 reverse proxy also protects the LLN from a number of resource depletion attacks from attackers on
434 the Internet. By design a reverse proxy handles all messages for all constrained devices in a LLN from
435 Internet hosts. Thus, the reverse proxy becomes the main traffic entry point for the LLN and therefore
436 it can inspect, filter and drop traffic in order to root out traffic from malicious users. Combined with
437 the strong authentication of clients and an access control policy, this proxy can make more informed
438 decisions in regards to filtering traffic when compared to e.g. a simple Internet firewall.

439 In terms of the application layer, a reverse proxy is free to process and transform the requests
440 it receives from clients as it chooses. A reverse proxy can improve network access by offering
441 features such as caching, network-edge access control and enforcing congestion control algorithms.
442 Interoperability with other systems can be increased by e.g. translating between HTTP and CoAP,
443 which is fairly straightforward considering the design goals of CoAP. Translation between different
444 data types (e.g. Core link format [?] to JSON) can also boost interoperability. Such a proxy can also
445 implement additional application functionality on behalf of the constrained device. Examples of such
446 functionality include extending the constrained device with semantic descriptions for its resources,
447 a deployment location photo, the weather near the device, etc. Additionally, a proxy can choose
448 to facilitate adding, configuring and deploying such functionality via a plugin-like system. This
449 greatly eases management of such functionality at run time by making adding, updating, enabling
450 and disabling such functionality easier.

451 It is important to reiterate that all of the above is possible without any additional configuration on
452 either the constrained device or the client. Nor does the presented approach require any modifications
453 to the standards compliant protocol stacks (e.g. 6LoWPAN/DTLS/CoAP) running on the constrained
454 device and the client. Indeed, the client discovers the Internet endpoint of the constrained device
455 that is hosted on the proxy and the proxy takes care of mapping every request to the corresponding
456 constrained device. In the scenario presented here, all configuration is limited to the proxy. These last
457 two benefits are an important differentiator from existing work, as will be discussed in the related
458 work section.

459 While the reverse proxy approach offers a number of benefits, it also entails some risks that
460 if ignored might undermine the presented system. One risk is that the reverse proxy presents a

461 single point of failure in terms of security and operation. Indeed, if the reverse proxy were to be
 462 compromised then e.g. all session keys and long-term keying material (pre-shared keys and private
 463 keys) could be made public. As the proxy offers a RESTful interface for managing virtual hosts and
 464 their keying material, this interface entails a security risk and should therefore be properly hardened
 465 against malicious usage (see section ?? for suggestions). Likewise, if the reverse proxy were to be the
 466 target of a resource depletion attack, then the constrained devices hosted by that proxy would become
 467 unreachable. On the other hand, as the proxy is deployed on a more powerful system, the proxy is
 468 more resilient to resource depletion attacks than constrained devices and networks. A second issue
 469 is the introduction of a third party (i.e. the proxy itself) into the trust model by terminating the
 470 end-to-end security that must be trusted by both the constrained device as well as the clients. As
 471 all collected data and issued commands pass via the proxy, this can raise privacy concerns when the
 472 device or the client does not trust the owner of the proxy. One option to mitigate this privacy risk is
 473 to let the owner of the constrained devices operate the reverse proxy on his or her own. To this end,
 474 our evaluation shows that a low-cost single board computer (e.g. Raspberry Pi) is capable of hosting
 475 the proxy, which enables on-premises deployments. To summarize, the proxy breaks end-to-end
 476 security in order to provide additional features which address operational and performance concerns
 477 of resource constrained devices. This work argues that the benefits of terminating the end-to-end
 478 security outweigh the security-related risks in the case of 'Class 1' resource constrained devices and
 479 networks. For less constrained devices and networks, this balance might tip in favor of end-to-end
 480 security.

481 4.2. Secure Service Proxy: design

482 In order to enable our proxy to extend constrained devices with a wide range of functionality,
 483 the design adopts the concept of virtual devices. In our design every virtual device is allocated a
 484 dedicated IPv6 address from an IPv6 subnet that is either routed to the proxy or directly connected
 485 to the proxy. Every virtual device has one or more endpoints associated with it. An endpoint
 486 corresponds to a transport and application layer binding: e.g. UDP/CoAP, DTLS/CoAP, TCP/HTTP
 487 or TLS/HTTP. For every virtual device the proxy listens for traffic on each of its endpoints, this is
 488 shown in the bottom left of figure ??.

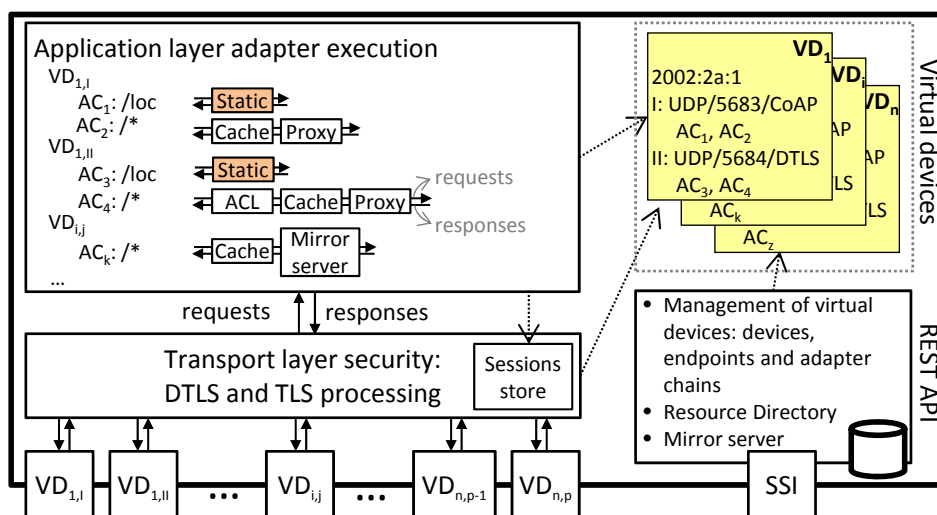


Figure 4. Secure Service Proxy: design

489 The transport layer security block is responsible for handling the (D)TLS protocol for secure
 490 endpoints on behalf of virtual devices. As such this block performs (D)TLS handshakes, thereby
 491 authenticating the client and performing a key exchange. To this end the block interfaces with the

492 virtual device configuration (top right in the figure) to retrieve the TLS parameters that are configured
493 for the virtual device. These parameters include a list of available cipher suites and keying material
494 for the secure endpoint of the virtual device as well as whether the virtual device requires clients
495 to authenticate themselves. Apart from the handshake, this block is responsible for tracking active
496 sessions with virtual devices (via the sessions store). It also decrypts and verifies incoming (D)TLS
497 application data messages, which are passed on to the adapter execution block, as well as encrypts
498 outgoing application data that comes from the adapter block. The keying material and the protocol
499 state used in the encryption and verification process naturally depends on the endpoint involved.

500 Incoming messages contain (secured) requests which are either HTTP or CoAP requests. While
501 our design supports adapters for both application layer protocols, we foresee that HTTP requests will
502 almost always be translated immediately to a CoAP request. As such we do not expect virtual devices
503 to host only an HTTP endpoint (although the design does support this). When the application layer
504 adapter execution block receives a request, it will search through the tree of available adapter chains
505 to search for a chain that is the most specific match for the request. The current implementation
506 supports searching based on the address and endpoint of the virtual device as well as the URI of the
507 request.

508 Once a chain has been found the execution block will pass the request along the chain. Every
509 element of the chain (i.e. an adapter) can either return (a modified) the request, which will be passed
510 to the next adapter in the chain, or stop the execution of the chain by returning a response. The current
511 implementation allows returning a response from an adapter in a non-blocking (i.e. asynchronous)
512 way, as retrieving a response might involve a lengthy IO operation. Once the response is available,
513 it is passed along the chain in reverse. This allows adapters to process and (if needed) modify the
514 response before it is stored in the virtual device and returned to the client.

515 Application layer adapters implement the functionality hosted by virtual devices. The idea
516 underlying adapters is to compartmentalize functionality into modules that can be reused by virtual
517 devices. When creating an adapter chain, an instance for every adapter in the chain is created and
518 every instance is configured according to the parameters exposed by the adapter type (see further).
519 While instances of adapters reside in adapter chains, they can be shared by more than one adapter
520 chain. For example in figure ?? the same Static adapter instance (colored orange) is shared by AC₁ and
521 AC₃. This is mainly useful when the same functionality should be available for multiple endpoints
522 of the same virtual device (e.g. CoAP and CoAPs) or when an adapter implements functionality
523 that does not require configuration that differs per adapter chain (e.g. a logging adapter that logs all
524 incoming requests for auditing purposes).

525 The proxy also exposes a networked interface in the form of a REST API to manage virtual
526 devices, which is shown in the bottom right of figure ?. The REST API allows creating and deleting
527 virtual devices and their endpoints, as well as instancing and deleting adapters and defining adapter
528 chains. When creating (D)TLS endpoints the REST API also allows specifying the cipher suites
529 supported by the virtual device, as well as the keying material (e.g. X.509 certificate or private key).
530 Apart from the management interface, the proxy also hosts a resource directory that contains the
531 hosted virtual devices. Finally, a mirror server is also available to enable resource updates from
532 constrained devices that are asleep for continuous and long periods of time (i.e. sleepy devices).
533 This mirror server can be used by virtual devices to interface with resources from sleepy constrained
534 devices.

535 Finally, the presented design allows to deploy the proxy on different locations in the network by
536 varying the IPv6 subnet for the allocation of virtual device IPv6 addresses. We foresee two scenarios.
537 In a first scenario, the proxy resides close to the constrained devices by allocating addresses from
538 a neighboring LAN network to virtual devices. An example would be a home LAN network from
539 which the proxy assigns unused addresses to virtual devices. In the case of a 6LoWPAN network,
540 the proxy can be combined with the border router. This scenario also aligns nicely to the distributed
541 computing concept that is commonly found in fog computing and in in-network processing [?]. In

542 a second scenario, the proxy resides further ‘upstream’ from the constrained devices (e.g. in a data
 543 center, the cloud, etc.) and allocates addresses from a special-purpose IPv6 subnet that is dedicated
 544 to virtual devices. In this scenario, the routing has to be configured to route this special-purpose IPv6
 545 subnet via the proxy (which is not a problem in most data centers). Both scenarios are complementary
 546 and will depend on the specific needs of the considered use-case: e.g. a proxy in the LAN network
 547 means that data stays inside the home network, which may benefit privacy. Similar considerations
 548 were previously discussed in the problem statement section.

549 4.3. Secure Service Proxy: implementation

550 For the implementation of our secure service proxy, we chose to build upon the previous work
 551 in our CoAP++ framework (which in turn builds on top of the Click modular router software). This
 552 choice provides a great amount of flexibility in how we process the network traffic for the virtual (and
 553 constrained) devices, as all routing functions are part of Click and can therefore be configured to our
 554 liking. In terms of the (D)TLS implementation we chose to use the wolfSSL library as this offers the
 555 easiest API for managing sessions and integrating into Click router where most processing happens
 556 on network packets.

557 4.3.1. Virtual devices and endpoints

558 Virtual device endpoints are created and deleted via the management interface. This is a
 559 straightforward REST interface that is hosted on the secure service proxy over CoAPs. As this
 560 interface handles sensitive information such as keying material, access is restricted to authorized
 561 users which are allowed to manage endpoints and adapter chains.

562 POST requests with an endpoint description are used to create a new endpoint for a virtual
 563 device. The endpoint description contains both the virtual device to which the endpoint belongs as
 564 well as any configuration details describing the endpoint itself. This description is serialized as a
 565 JSON object in the payload of the POST request. For a plain-text CoAP endpoint, the configuration
 566 details are limited to the UDP transport port of the endpoint. For a DTLS CoAPs endpoint, the
 567 configuration also includes information about the supported cipher suites and any parameters for
 568 the cipher suites. In the current implementation, the “TLS_PSK_WITH_AES_128_CCM_8” and
 569 “TLS_ECDHE_ECDSA_WITH_AES_128_CCM_8” cipher suites are supported for CoAPs endpoints.
 570 When creating an endpoint that supports the PSK cipher suite, the pre-shared-key and an (optional)
 571 client identity hint have to be specified as parameters. For the elliptic curve DSA suite, the secp256r1
 572 private key and signed certificate have to be provided as parameters. These are both encoded in
 573 base64 in the endpoint description. The following listing contains an example POST request that
 574 creates a CoAPs endpoint for a virtual device hosted under 2001:6a8:1d80:23::1 on port 5684 with an
 575 ECC cipher suite.

```
576 POST /virtualDevices
577 Content-Format: application/json
578 {
579   "address": "2001:6a8:1d80:23::1",
580   "prefixLen": 128,
581   "port": 5684,
582   "dtls": {
583     "supportedCipherSuites": [
584       {
585         "cipherSuite": "TLS_ECDHE_ECDSA_WITH_AES_128_CCM_8",
586         "parameters": {
587           "b64PrivateKey": "QVNO...==",
588           "b64Certificate": "LS0t...=="
589         }
590       }
591     ]
592   }
593 }
594
595 2.01 Created /virtualDevices/2001:6a8:1d80:23::1~128~5684
```

596 The response of the secure service proxy links to a newly-created resource that can be used to
 597 the delete the endpoint at a later time. This resource is also used for managing the adapter chains that
 598 belong to an endpoint, as explained in section ??.

599 4.3.2. Implemented application layer adapters

600 In terms of application layer adapters, our proxy currently implements the adapters listed in
 601 table ?. This section describes each of the adapter types in more detail.

Table 1. The proxy offers a number of functionalities, called adapters, that are hosted on virtual devices. The list of adapters that were implemented at the time of this work are shown in this table.

Adapter	Functionality	Configuration parameters
Access control	Restrict access to virtual devices depending on client identify, request method and URI.	ACL rules and default policy
Static resource	Host RESTful resources on virtual devices that can be read and modified.	Payload and content type
Cache	Cache and serve previous responses from virtual devices to clients.	Default cache entry lifetime
Congestion control	Enforce congestion control on clients querying virtual devices. Per device and network wide rules are implemented.	Per user CC limits
.well-known/core	Manipulate discovery responses from virtual devices to include functionality hosted by the proxy.	None
Proxy	Proxies requests for the virtual device to a CoAP(s) server (e.g. the constrained device). Also aggregates observe registrations.	CoAP(s) server endpoint
Mirror server	Proxies requests for a virtual device to a mirror server.	Mirror server endpoint and sleepy device anchor point

602 The access control adapter applies ACL rules to the CoAP(s) requests it processes. ACL rules are
 603 parsed as JSON objects that assign allow and deny rules to either a username or a role of users. An
 604 allow and deny rule consist of a regular expression, which is applied to the request URI, and a list
 605 of request methods. In case no matching ACL rule is found, then the default policy of the adapter
 606 instance (either accept or deny) is applied. The following JSON serialization of an example ACL rule
 607 gives bob full access to the devicename resource while access to the lock resource is restricted to read
 608 only.

```
609 {"username": "bob",
610 "allow": [{"uri-regex": "devicename", "methods": ["GET", "PUT", "POST", "DELETE"]},
611           {"uri-regex": "lock", "methods": ["GET"]}],
612 "deny": []}
```

613 Hosting a virtual resource on a virtual device is the task of the static resource adapter. In order
 614 to allow arbitrary content types of the payload, the value of the virtual resource is encoded in base64
 615 in the configuration of the adapter. An example is shown in the next section.

616 The cache adapter serves and caches responses for requests to virtual devices. The cache adapter
 617 calculates a cache key for every CoAP request it handles. When a fresh response matching the
 618 cache-key is found, the adapter chain's execution is halted and the cached response traverses the
 619 adapter chain in reverse. Responses processed by the cache adapter are handled in accordance with
 620 section 5.9 of the CoAP RFC [?]. This means that e.g. a '2.05 Content' response will be cached, while
 621 a '2.04 Changed' response will mark any stored response as not fresh. Cached responses are removed
 622 when they expire after their Max-Age option. Note that the cache adapter does not implement the
 623 'Validation Model' specified in section 5.6.2. of the CoAP RFC. When used in conjunction with access
 624 control it is important that all ACL rules are applied before hitting the cache, as the execution of the

625 request leg of the adapter chain will stop when a cache hit is found. The underlying implementation
626 caches responses in memory via a memcached instance.

627 The congestion control adapter in its current form applies traffic shaping on a per host basis.
628 Currently it is possible to limit the number of open requests between a client and a specific virtual
629 device and between a client and a group of virtual devices ¹. Open requests are requests for which a
630 response has not been sent yet. If a client reaches its limit, then the request is dropped until either a
631 response is received or one of the prior requests of that client is removed after a time out period (can
632 be configured). Finally, a client can either be identified by its endpoint address or by its identity
633 derived from the authentication credentials during the (D)TLS handshake.

634 The .well-known/core adapter is responsible for including the functionality that is hosted on
635 the virtual devices in the resource discovery responses of the real constrained device. In the current
636 implementation, the wkc adapter asks every adapter from all the adapter chains that are defined
637 for the virtual device to modify the discovery response from the real device. This way the static
638 resource adapter can add a link to its virtual resource and the ACL adapter can remove links
639 for resources that the user is not authorized to access. To this end, every adapter type offers a
640 “processDiscoveryResponse” method that is used by the wkc adapter.

641 The proxy adapter takes a request for a virtual device and issues a new CoAP request to the
642 corresponding actual constrained device. Therefore an instance of this adapter is configured with
643 the CoAP(s) endpoint of the constrained device. Only the transport layer addresses are changed,
644 the new CoAP request is copied from the output of the previous adapter in the adapter chain (with
645 the exception of the Message ID and the Token of course). The proxy adapter will either retrieve a
646 response or generate a time-out, therefore it always comes last in adapter chains. This adapter will
647 also combine observe registrations when it receives multiple registrations for the same resource on a
648 virtual device. Likewise it also multiplexes responses from constrained devices to multiple clients in
649 case there is more than one ongoing observe registration.

650 Finally, the mirror server adapter is a special type of proxy adapter in that it issues CoAP(s)
651 requests to a mirror server instead of the constrained device itself. Apart from the end point of the
652 mirror server, also the handler of the constrained device is configured into the mirror server adapter
653 instance. For instance a request to the coaps://vd1.iot.test/status resource on a virtual device would
654 be translated to coaps://ms.iot.test/ms/0/status.

655 4.3.3. Adapter chain management: interface

656 Once an endpoint for a virtual host has been allocated on the proxy, adapter chains can be
657 created and hosted on that endpoint. Building on our previous example, the listing below contains
658 a CoAP request that instantiates an adapter chain which contains the access control, well-known
659 core rewriting, caching and forward CoAPs proxy adapters. Again, the payload is a JSON object
660 that describes the chain and contains the parameters for the different adapter instances. The adapter
661 chain is created as the default chain via the wildcard character in the chain path. The default chain is
662 executed for requests where no other adapter chains with a matching URI path are found.

```
663 POST /virtualDevices/2001:6a8:1d80:23::1~128~5684
664 Content-Format: application/json
665 {
666   "path": "/*",
667   "pipeline": [
668     {
669       "type": "acl",
670       "default_access_control_policy": "deny",
671       "rules": [
672         {"username": "fvdabeele", "rules": [{"uri-regex": "regex1", "allowMethods": ["*"]},
673         {"uri-regex": "regex2", "allowMethods": ["GET"]}]},
```

¹ The group encompasses all virtual devices with an adapter chain that share the same CC adapter instance.


```

674         {"username": "*", "rules": [{"uri-regex": "regex1", "denyMethods": ["*"]},
675         {"uri-regex": "regex2", "allowMethods": ["GET"]}]}}
676     ]
677 },
678 {
679     "type": "wkc"
680 },
681 {
682     "type": "cache",
683     "default_lifetime": 60
684 },
685 {
686     "type": "proxy",
687     "scheme": "coaps",
688     "addr": "bbbb:1",
689     "port": 5684
690 }
691 ]
692 }
693
694 2.01 Created /virtualDevices/2001:6a8:1d80:23::1~128~5684/*

```

695 The second example, shown in the listing below, details how a static resource is created on the
696 endpoint of our virtual host (in this case it contains a semantic description of the virtual host in the
697 RDF format). The chain also illustrates the linked adapter, which refers to the acl adapter instance
698 that was created in the previous listing. The link points to the management resource of the adapter
699 instance.

```

700 POST /virtualDevices/2001:6a8:1d80:23::1~128~5684
701 Content-Format: application/json
702 {
703     "path": "/rdf",
704     "pipeline": [
705         {
706             "type": "linkedAdapter",
707             "link": "/virtualDevices/2001:6a8:1d80:23::1~128~5684/*/*"
708         },
709         {
710             "type": "static",
711             "contentType": 41,
712             "value": "PGh0d...=="
713         }
714     ]
715 }
716
717 2.01 Created /virtualDevices/2001:6a8:1d80:23::1~128~5684/rdf

```

718 Finally, note that the parameters of existing adapters can be updated via a PUT request to the
719 management resource of the adapter instance. In this case the payload is a JSON object where the
720 keys are the parameter names. Likewise, adapters and chains can be deleted via their respective
721 management resources.

722 4.3.4. Authenticating (D)TLS clients on the SSP

723 In order to facilitate authentication of users and authorization of user actions, the SSP links client
724 authentication information (e.g. TLS PSK or X.509 client certificate) with users and roles. The current
725 implementation is limited to using TLS primitives for supplying authentication credentials, although
726 in the future alternatives might be considered (e.g. lightweight application-layer access tokens). For
727 example, a (D)TLS session that was setup with PSK_1 as the pre-shared key can be linked with $user_A$.
728 Likewise attributes in a client X.509 certificate that is signed by a party trusted by the SSP can be
729 linked with a specific user. E.g. a certificate issued by CA_A with the common name attribute set to
730 *fdabeele* can be linked with $user_B$. Finally, the proxy also exposes a RESTful interface for managing
731 which credentials belong to which user and the roles of users.

732 4.3.5. Key management between SSP and constrained devices

733 The SSP contains an in-memory repository of pre-shared keys and corresponding identity
734 hints in order to setup DTLS sessions with resource-constrained CoAPs servers. As this repository
735 contains all the keying material for the constrained devices known to the proxy, it contains sensitive
736 information and should be handled accordingly. In the current implementation this repository is
737 initialized when the SSP process is started. A future extension could enable at run-time manipulation
738 of this repository by, for example, specifying keying material when instantiating coaps proxy
739 adapters. Currently this has not yet been implemented, as in our use cases this repository does not
740 change frequently and remains stable. In use cases where the repository is more volatile, such an
741 extension could enable better key management.

742 5. Related work

743 The concept of device virtualization in the IoT is widespread in literature, though often times
744 under different names such as sensor, thing and object virtualization. Indeed, in [?] the authors
745 present a survey on object virtualization in the IoT stating that “the concept has become a major
746 component of current IoT platforms where it aids in improving object energy management efficiency
747 and addressing heterogeneity and scalability issues”. The authors classify existing architectures as
748 one or many real objects for one or many virtual objects. While the focus in this work has been on one
749 real object for one virtual object, the flexibility of the presented design enables the same adapter to be
750 shared by multiple virtual devices as well as one virtual device to span multiple physical devices (for
751 example a virtual device combining all lamps in a room).

752 There exist numerous works in literature that study the benefits of using third parties or
753 intermediaries in constrained environments. In order to narrow the scope of this section, only works
754 that are relevant in the context of constrained RESTful environments are discussed here. In [?],
755 Kovatsch et al. discuss moving application logic from firmware to the cloud. According to the vision
756 of the authors devices are thin servers exposing RESTful resources for data access and actuation
757 and most of the application logic would reside in application servers. While our approach also
758 advocates thin servers for devices, deploying the SSP in the cloud is optional. In use-cases where
759 local access is important, the SSP may reside closer to the devices (e.g. deployed in the LAN) in
760 order to meet requirements in respect to latency, privacy or availability. Additionally, the SSP may
761 support constrained nodes and applications servers by providing functionality such as caching and
762 more scalable authentication and authorization. The IPv6 addressing proxy presented in [?] is an
763 example of an intermediary system for mapping legacy technologies to the IPv6 Internet of Things.
764 By allocating IPv6 addresses to map to different legacy technologies, the approach is similar to the
765 virtual devices presented in our work. Note that the adapter concept provides the flexibility to map
766 virtual devices to different technologies similar to the work in [?]. While not presented in this work,
767 the SSP has been used to host LoRaWAN end devices as virtual IPv6 CoAP endpoints via an AMQP
768 pub/sub adapter that interfaced with the LoRaWAN network server. The authors in [?] propose
769 to interconnect web applications based on HTTP and web sockets with CoAP-based wireless sensor
770 networks via a CoAP proxy. The CoAP proxy focuses on translating between different protocols
771 and closely follows the guidelines outlined in RFC 8075 [?]. The scope of the SSP is broader as it
772 includes transport security, access and congestion control next to mapping HTTP to CoAP. Finally,
773 note that the forward proxy approach of Ludovici differs from the reverse proxy approach of the
774 SSP. In [?], Mongozzi et al. introduce a framework for CoAP proxy virtualization in order to
775 address the scalability and heterogeneity challenges faced in large-scale Web of Things deployments.
776 The framework installs a reverse CoAP proxy on the sensor network gateway and then applies
777 virtualization so that the proxy can be customized and extended by third parties without modifying
778 the reverse proxy. All interactions of these virtual proxies with smart objects pass via this reverse
779 proxy, which acts as an arbiter for access to the limited resources of the smart objects. The presented
780 approach is interesting as the containerization of the virtual proxies into virtual machines makes them

781 more flexible than the adapter approach followed in the SSP. We have experimented with providing
782 some degree of extensibility by creating adapters from python scripts in the SSP (these scripts could
783 be uploaded via the adapter chain management interface). While this python adapter type provided
784 some degree of customization, the lack of proper process isolation meant that (malicious) scripts
785 could stall the SSP. As such, these python adapters did not make the final SSP design. While the
786 concept of the virtual proxies is interesting, the extent of the work is limited as the focus lies on the
787 virtualization technique and interesting features such as scalable security and efficient and authorized
788 network access are not considered. Instead the authors focus on providing service differentiation
789 between multiple virtual proxies. Also note that proxy virtualization is not the same concept as
790 device virtualization, though they can be used to solve similar problems. The same authors of [?]]
791 look at the specific problem of proxying CoAP observe efficiently for different QoS requirements in [?
792]. While the scope of the work is quite different from this paper, the use of a reverse proxy for bundling
793 observe relationships is shared between the two works. Another example of device virtualization in
794 RESTful environment is [?], where the authors assign virtual coap servers to RFID tags. The actual
795 CoAP servers are not running on the tags though. Instead they reside on RFID readers, which are
796 able to enhance tags with additional functionality (such as discovery). This work has parallels with
797 the SSP, which enhances constrained devices by means of application layer adapters.

798 A second category of relevant works in literature studies the challenges faced by transport
799 layer security in constrained IoT environments. There are a number of works that study the DTLS
800 handshake as it is fairly complex and verbose process with significant resources requirements for
801 constrained devices. In [?]] Hummen et al. propose a delegation architecture that offloads
802 the expensive DTLS connection establishment to a delegation server thereby reducing the resource
803 requirements of constrained devices. The delegation architecture also enables more complex
804 authorization schemes, as it has more resources at its disposal. The authors report significant
805 reductions on memory overhead, computations and network transmissions on constrained devices.
806 Our termination method can also provide complex authorization schemes of the virtual device.
807 In section ?? we have also reported significant savings in regards to CPU and network resource
808 usage (and consequently energy usage). While our approach still requires an active DTLS session
809 between the SSP and the constrained device, the number of handshakes during the lifetime of a device
810 is drastically reduced. While the memory requirements are not as low as in [?], they are still lowered
811 as the constrained device can limit the number of simultaneous sessions to one. Finally note that our
812 approach does not require any changes to the DTLS stack running on the device. The work in [?]
813 focuses on various challenges in deploying DTLS in resource constrained environments. Similarly
814 to [?], the approach revolves around handshake delegation. The authors adopt the concept of secure
815 virtual things in the cloud where physical things delegate the session initiation to their corresponding
816 virtual thing. As a result physical things can limit their DTLS implementation to only the record
817 layer protocol, which leads to drastic memory savings. One interesting aspect of the presented
818 architecture is that the physical thing can assume both roles of client and server. Unfortunately the
819 concept of virtual things is not extended beyond the handshake delegation mechanism. It would be
820 interesting to combine a delegation mechanism with some of the findings presented in our work. A
821 hybrid option would be possible where the delegation mechanism is used for the most constrained
822 devices (requiring a custom lightweight DTLS stack) and where the termination mechanism can be
823 used for devices with sufficient memory (i.e. where a full DTLS stack is feasible) or where the DTLS
824 stack can not be customized to implement the delegation method.

825 Object Security of CoAP (OSCOAP) [?]] is an IETF Internet Draft standardizing end-to-end
826 security of CoAP options and payload at the application layer. While the specification focuses
827 on the forwarding case when using a forward proxy (which excludes caching), it does include an
828 appendix describing a mode of operation, Object Security of Content (OSCON), which is compatible
829 with caching responses at intermediaries. The draft notes that OSCOAP may be used in extremely
830 constrained settings, where CoAP over DTLS may be prohibitive e.g. due to large code size.

831 Nevertheless, the authors state that OSCOAP may be combined with DTLS, thereby benefiting from
832 the additional protection of the CoAP message layer present in DTLS-based security. Note that the
833 standardization efforts focus on the case of a forward proxy, whereas this work focuses on a reverse
834 proxy approach. As such, the trust models are different as the reverse proxy represents the end device
835 from the point of view of the client. Despite the difference in proxy models, the two approaches
836 remain compatible and could strengthen each other. For example, the SSP could implement OSCOAP
837 for cases where clients are employing a forward proxy, which is not trusted by the client. Additionally,
838 it would be interesting for the SSP to support OSCOAP as a lightweight alternative for DTLS to
839 protect communications with constrained devices with severe memory limitations. In such a case,
840 clients would communicate securely with the SSP over DTLS while the communications between the
841 SSP and the constrained devices would be protected either via OSCOAP (e.g. for constrained devices
842 with severely limited memory) or via DTLS (e.g. for constrained devices with sufficient memory).

843 Finally, in high volume web environments transport layer security is often terminated at a
844 proxy deployed close to the web server(s). The main motivation for terminating TLS is that it
845 enables load balancing, where terminated HTTPS requests are distributed over multiple web servers.
846 Load balancing increases the availability of the web deployment, as the outage of one web server
847 does not affect the service availability in this case. Popular Web proxy software, like nginx and
848 HAProxy, supports different reverse proxy deployment options for terminating TLS. Similarly the
849 elastic cloud computing platform of Amazon.com, Amazon Web Services, supports TLS termination
850 and load balancing by virtue of its HTTPS listener service. While the main motivation of the SSP
851 for session termination is not load balancing, the SSP does apply termination in order to be able
852 to move computationally expensive and verbose operations from constrained devices to the proxy
853 which improves performance. Similarly to high availability TLS proxies, the SSP may reduce key
854 management complexity as all keying material for public communications are stored on one system.

855 **6. Evaluation: results and discussion**

856 This section presents two evaluation scenarios that show the gains attainable by our approach.
857 Such gains include: a decrease in load on constrained devices and the LLN, lower energy usage for
858 constrained devices, an increase in user handling capacity of LLNs, more responsive LLNs, more
859 scalable authentication and better authorization. The evaluation scenarios were chosen to evaluate
860 the impact of the proxy on two specific operational aspects of LLNs: setting up DTLS sessions with
861 constrained devices over multiple wireless hops and observing CoAPs resources on constrained
862 devices from multiple DTLS clients.

863 *6.1. Terminating end-to-end-security at the SSP*

864 The first evaluation scenario is geared towards quantizing the impact of splitting end-to-end
865 security at the smart service proxy. More specifically, the goal is to study the impact of re-using a
866 DTLS session of a constrained CoAPs server on the operation of both the constrained node as well as
867 the CoAPs client.

868 *6.1.1. Simulation setup*

869 Extensive simulations were performed with a nine node 6LoWPAN network arranged in a cross
870 topology as detailed in figure ???. One node is at the center of the cross and is the RPL border router
871 of the 6LoWPAN network, four nodes are intermediate routers (each located in the middle of one
872 of the four legs of the cross) and the last four nodes are CoAP(s) servers that are located at the four
873 ends of the cross. The border router is connected to the smart service proxy, which is running on the
874 same PC as the Cooja simulator. Finally, an unconstrained CoAP(s) client interacts with the CoAP(s)
875 servers. In the evaluation scenario the client sends the following sequence of CoAP(s) requests: a
876 .well-known/core discovery request, a sensor measurement request for the "/s" resource and an
877 actuator request for the "/a" resource. The constrained CoAPs servers are running er-coap and

878 TinyDTLS (in Contiki) configured to accept the 'TLS_PSK_WITH_AES_128_CCM_8' cipher suite with
 879 a PSK hint of 15 bytes.

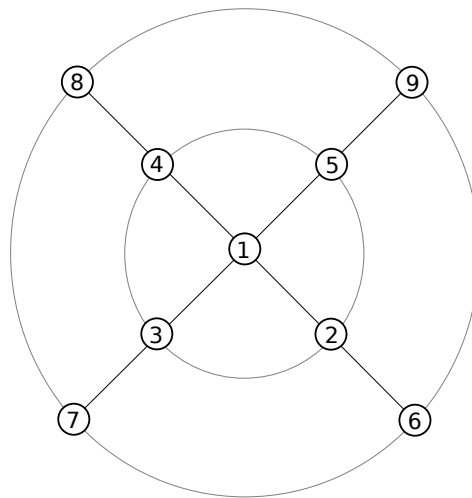


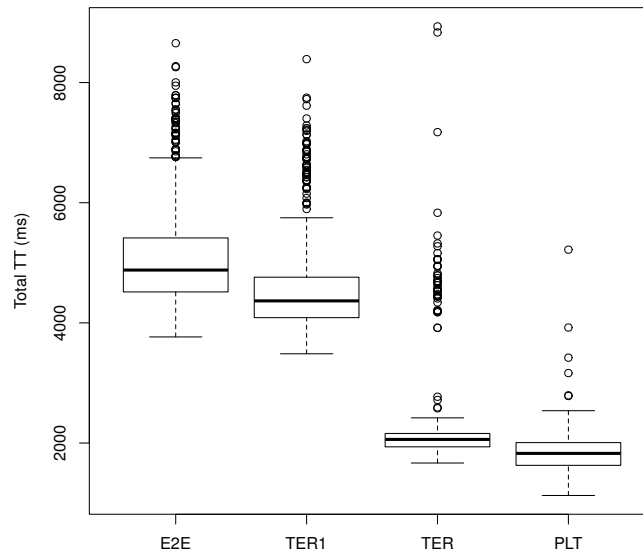
Figure 5. Cooja network topology: four CoAP(s) servers (6, 7, 8, 9) are located two hops away from the RPL border router.

880 The same request sequence was sent to the CoAP(s) servers for one reference case and three
 881 different SSP configurations: plain text (PLT), end to end (E2E), first termination (TER1) and n-th
 882 termination (TER). In the PLT configuration, all requests are sent over plain-text CoAP. This is a
 883 reference cases for the other three cases. In the E2E case, all requests are sent over CoAPs without any
 884 termination of DTLS sessions at the SSP. In case of TER1 and TER, all requests are sent over CoAPs
 885 and the DTLS session is terminated at the SSP. For TER1, there does not exist an active DTLS session
 886 between the proxy and the constrained node. Therefore a new DTLS session must be setup between
 887 the SSP and the constrained node. For TER, the active DTLS session in the LLN can be re-used and
 888 there is no need to setup a new DTLS session with the constrained node. For all DTLS cases, the
 889 DTLS client always sets up a new DTLS session at the start of a request sequence. It also tears down
 890 the existing session at the end of every sequence. As such, this testing scenario represents a large
 891 number of DTLS clients that would interact with the constrained CoAPs servers over the lifetime of
 892 the constrained node. For each configuration the request sequence was run four hundred times, i.e.
 893 one hundred times per DTLS server. All results were obtained using the default CSMA MAC protocol
 894 and ContikiMAC RDC protocol as available in Contiki.

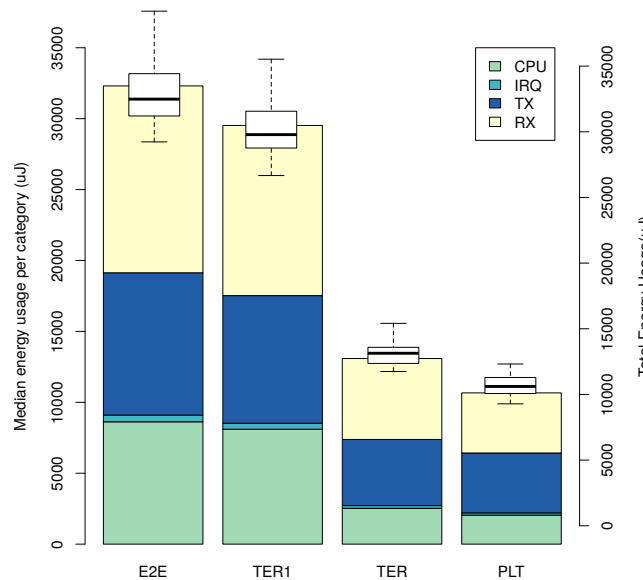
895 6.1.2. Results

896 Figure ?? shows the total transaction time (TTT). This is the time between the start of the DTLS
 897 session handshake (i.e. when the first ClientHello message is sent by the client) and the end of the
 898 DTLS session (i.e. when the DTLS Finished message is received by the client). There is a significant
 899 reduction in TTT between the E2E and the TER configurations: their medians are 4879 ms and 2060 ms
 900 respectively. This is due to the DTLS session re-use in the LLN, which saves - when comparing the
 901 median cases - thirteen packets in the LLN as the DTLS handshake in the LLN can be avoided in the
 902 TER configuration. As a result, the TER configuration is able to closely match the reference plain-text
 903 case in terms of TTT. The 233 ms difference in median is caused primarily by the overhead of the
 904 additional DTLS headers (X B per DTLS packet). More specifically the overhead triggers 6LoWPAN
 905 fragmentation for the large discovery response in the TER case, whereas this fragmentation is absent
 906 in the PLT case.

907 Figure ?? displays the energy usage for the different configurations. The stacked bar plot shows
 908 the median energy usage per category on the constrained device, whereas the box plot shows the total



(a) Total transaction times (TTT) for the request sequence



(b) Median energy usage per category (left axis) and total energy usage (right axis).

Figure 6. Transaction times and energy usage of the CoAPs servers for the three gateway configurations (E2E, TER1, TER) and the plain text CoAP reference case (PLT)

909 energy usage (to show the dispersion of the measurements). Again, there exist a significant difference
 910 between the E2E and the TER configurations: 32485 μJ vs 13133 μJ respectively (a reduction by a
 911 factor of 2.4). Similarly to the TTT results, this reduction is primarily due to the absence of the DTLS
 912 handshake in the LLN. This is confirmed by the bar plot where the energy usage for the RX and TX
 913 categories are reduced the most. The energy consumption in the CPU category is also significantly
 914 lower, as the CPU is in low-power mode more often and does not have to perform expensive hash
 915 calculations when completing the handshake. All in all, the results allow us to conclude that our
 916 approach increases the responsiveness of constrained devices (provided there is an active session in
 917 the LLN) while reducing the energy consumption for traffic loads with many DTLS sessions (e.g.
 918 traffic loads with many parties).

919 Finally, it is worth pointing out that our approach drastically limits the total number of
 920 handshakes that a constrained node will perform during its lifetime. Apart from the savings
 921 discussed above, this also has the additional benefit that - in lossy networks - the total number of
 922 failed handshakes will be lower. Indeed, Garcia et al. [?] have shown that in lossy networks
 923 the fraction of failed handshakes can vary significantly based on the packet loss ratio: e.g. 30-40%
 924 of handshakes fail for a PLR of ~20%. By limiting the total number of handshakes, our approach
 925 also limits the amount of constrained device resources wasted on these failed handshakes. On the
 926 other hand, care should be taken to periodically refresh keying material as needed by the underlying
 927 cryptographic primitives in use.

928 6.2. Aggregating multiple CoAPs clients at the SSP

929 The second evaluation scenario focuses on the impact of the SSP on constrained devices that
 930 serve multiple CoAPs clients simultaneously via CoAPs observe. Unlike clear text CoAP observe,
 931 notifications for one CoAPs client typically can not be reused to serve another client due to the
 932 confidentiality of the notification in DTLS. However, the SSP presented in this work can - as a reverse
 933 CoAPs proxy - observe one CoAPs resource on a constrained device and use these notifications to
 934 serve a multitude of CoAPs clients. The presented evaluation considers up to ten CoAPs clients that
 935 observe a resource on a constrained device and compares the case of end-to-end observation versus
 936 observation via the SSP. Note that one should keep in mind client authorization when using one
 937 CoAPs stream of notifications for serving multiple CoAPs clients. E.g. a client that is not authorized
 938 to access a resource on the constrained device, must also be denied access to that resource via the
 939 SSP. To this end, this work presents and implements an access control adapter which enforces CoAPs
 940 resource-specific access control policies.

941 6.2.1. Experiment setup

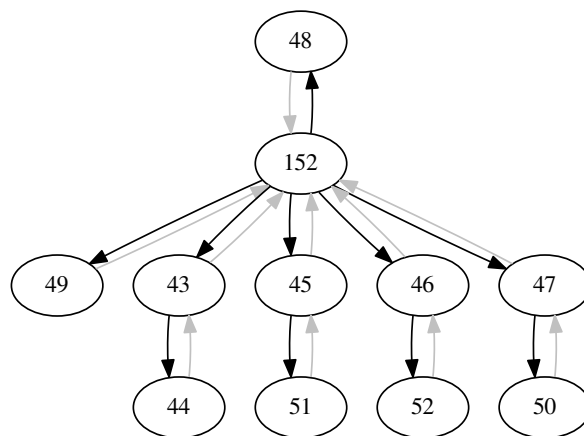


Figure 7. Representative RPL network topology: the node under study, node #50, is situated two hops from the border router, node #152.

942 To quantify the impact of aggregating CoAPs observations at the SSP, a number of experiments
 943 were run on a WSN testbed. The experiments consisted of a 6LoWPAN network with ten sensor
 944 nodes arranged on a line with six meters of spacing between adjacent nodes. An additional sensor
 945 node (node #152) is situated to the upper left of the line and is connected to a Raspberry Pi 2 where
 946 it serves as the RPL border router. The smart service proxy software is running on the Raspberry
 947 Pi 2. In order to cope with changes in the RPL topology between experiments and over time in the
 948 same experiment, node #50 was selected for testing as it was always located two hops away from
 949 the border-router. A representative network topology is shown in figure ???. Note that depending
 950 on the experiment node #50 might have a different parent than node #47 (e.g. node #43 was a

951 common alternative), but in all experiments there was always an intermediary router between the
 952 border router and node #50.

953 All wireless sensor nodes employ the msp430f5437 uC with 128KB of RAM and 256KB of ROM
 954 and the TI CC2520 802.15.4 transceiver. As such, the platform is identical to the WiSMote platform
 955 in Contiki in terms of the specifications that are relevant for the presented results. The sensor nodes
 956 run a TinyDTLS CoAPs server which is configured to support three simultaneous DTLS sessions and
 957 one simultaneous DTLS handshake. While a binary for four simultaneous sessions could be built,
 958 it was not running stable. Attempts for supporting more than four clients led to a RAM overflow
 959 during linking. By default er-coap in Contiki sends one confirmable notification for every twenty
 960 notifications. Finally, all sensor nodes in the network are running the default CSMA MAC protocol
 961 and ContikiMAC RDC protocol available in Contiki.

962 For every sensor node, a corresponding virtual host was created on the SSP. The
 963 virtual hosts were configured similar to the listing in section ??, with support for the
 964 “TLS_ECDHE_ECDSA_WITH_AES_128_CCM_8” cipher suite. This cipher suite provides perfect
 965 forward secrecy by means of an ephemeral Diffie-Hellman key exchange between the virtual hosts
 966 and the DTLS clients. Additionally, DTLS clients authenticate virtual hosts by means of the x.509
 967 certificates of the hosts, which are signed by a certificate authority (CA) trusted by the clients.
 968 Similarly, the DTLS clients also present a x.509 certificate during the DTLS handshake which is signed
 969 by a CA that is trusted by the proxy. As a result the clients may be authenticated at the proxy-side (by
 970 mapping attributes from the certificate to a user in the proxy, see section ??), which is mandatory
 971 for the use of the access control adapter in order to provide fine-grained authorization as presented
 972 in section ?. Each virtual host was allocated a global IPv6 address from the LAN network of the
 973 Raspberry Pi2 and has one default adapter chain with access control, caching and proxy adapters.
 974 The CoAPs clients ran as part of the CoAP++ framework on a PC that was located three IPv6 hops
 975 away from the Raspberry Pi2. All IPv6 addresses in use (i.e. CoAPs clients, RPI, virtual hosts and
 976 WSN nodes) were working, global IPv6 addresses. An overview of the evaluation setup is shown in
 977 figure ??.

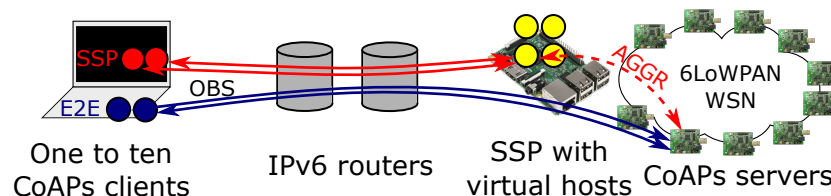


Figure 8. Evaluation setup: a variable number of CoAPs clients observe one of two resources on either the virtual host (SSP) or the sensor node (E2E)

978 In all experiments, a number of CoAPs clients observe a resource on either the virtual host or the
 979 sensor node. As such, the experiments considered two cases: end-to-end (E2E) CoAPs observations
 980 and CoAPs observations via the SSP. In both cases, experiments were run for two CoAPs resources: a
 981 resource with a one second notification period and another resource with a five seconds notification
 982 period. In the E2E case, experiments were performed with one, two and three simultaneous CoAPs
 983 clients. In the SSP case, experiments were performed with one, two, three, four, five and ten
 984 simultaneous CoAPs clients. In total, eighteen experiments were performed. Each experiment was
 985 run for at least twenty minutes, during which the energest outputs for all sensor nodes were captured
 986 every five seconds and the outputs from the CoAPs clients were stored as well. This enabled us
 987 to quantify the energy usage as well as the application-layer performance, the results of which are
 988 presented in the following section.

989 6.2.2. Results

990 When comparing the energy expenditure graphs for node #50 in figure ??, it becomes clear that
991 aggregating CoAPs observation relationships leads to energy savings. The savings are proportional
992 to the rate of notifications: they increase as the number of clients goes up and decrease as the
993 notification interval becomes longer. Note that the sensor node between node #50 and the border
994 router experiences similar energy savings as every notification is received and retransmitted by this
995 intermediary node. For the case of three CoAPs observers, the median energy expenditures differ by
996 10.8 mJ for the one second interval and 2.5 mJ for the five seconds interval.

997 For one CoAPs observer and the one second interval, there exists a small difference in energy
998 expenditure between the end-to-end and the SSP case even though the notification rate is the same for
999 both cases (i.e. one notification per second). This is primarily due to a difference in notification packet
1000 size as the 6LoWPAN compression for SSP notifications is more effective than for E2E notifications.
1001 The compression is more effective because the IPv6 address of the SSP is part of the 6LoWPAN
1002 network whereas the CoAPs client's IPv6 address is part of a different network. As such, the prefix
1003 of the SSP's IPv6 address can be elided (due to stateful 6LoWPAN compression), which leads to an
1004 eight bytes saving in packet size per notification.

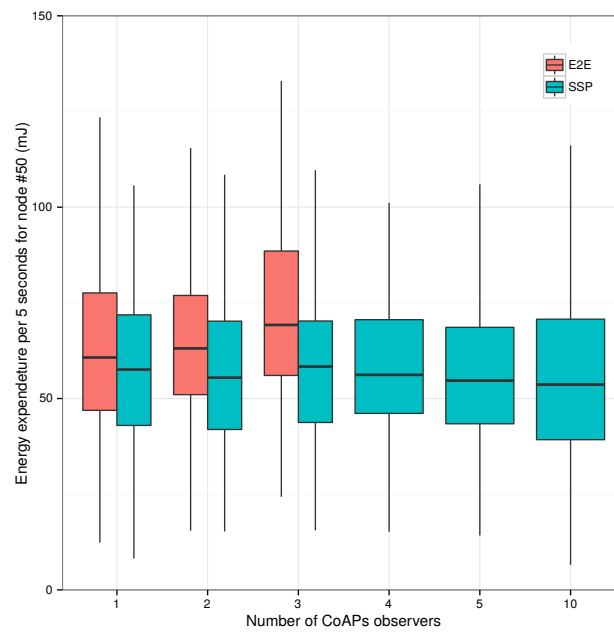
1005 The graphs in figure ?? clearly illustrate the difference in notification rate between the end-to-end
1006 and SSP experiments. Due to the aggregation of CoAPs observations at the SSP, there exists only
1007 one CoAPs observation between the SSP and the sensor node. This is illustrated in the constant
1008 notification rate for SSP as the number of CoAPs observers increases. For the end-to-end experiments
1009 the notification rate rises linearly with the number of observers, as the sensor node sends notifications
1010 to each client separately. The slope of this linear relation is proportional to the notification frequency.

1011 Figure ?? plots the notification loss ratios (NLR) for each of the eighteen experiments. For
1012 example for the E2E, one second interval and one observer case 1845 notifications were sent, three
1013 of which never arrived at the client. This leads to a NLR of 0.163%. Note that every vertical series
1014 of data contains as many points as there are observers, however very similar and identical NLR's
1015 overlap too much to distinguish them as separate points in the plot. The graphs for the one second
1016 interval show that the end-to-end case suffers from network congestion due to its higher notification
1017 rate. Also, the observed loss is heavily dependent on the CoAPs client in the E2E experiments: i.e. the
1018 client that is last on the list of observers experiences the highest NLR (mostly apparent when there
1019 are three observers). Finally, the SSP sends every notification as a confirmable message. While in
1020 this setup packet loss is mainly a problem in the constrained WSN, sending all notifications as CON
1021 messages can help to improve the NLR in situations where the client is part of a lossy network.

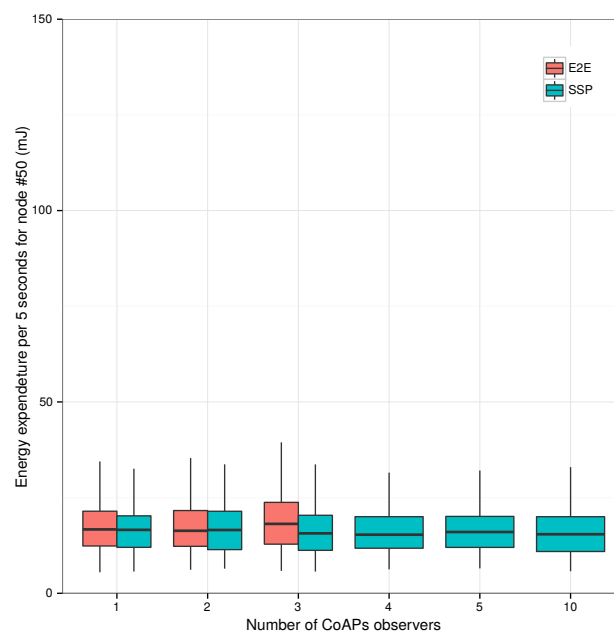
1022 To conclude, there are a number of limitations that are overcome by aggregating observations at
1023 the SSP:

- 1024 1. Memory and processing constraints on the sensor node, which limit the number of
1025 simultaneously active DTLS sessions and active CoAP observe relationships.
- 1026 2. Limited throughput in constrained (multi-hop) networks, which impact the successful delivery
1027 of notifications and limits the rate of notifications.
- 1028 3. Limited lifetime for battery-operated sensors, by reducing the load on constrained devices the
1029 lifetime is lengthened.

1030 Note that while only the results for node #50 are shown, similar savings apply for other nodes. Also
1031 note that applying observation aggregation at the SSP delays the point at which the WSN reaches
1032 congestion, as the message rate in the WSN is reduced by the aggregation. Finally note that this
1033 experiment is only possible because the SSP terminates the end-to-end security, indeed should this not
1034 be the case then the SSP would be unable to aggregate observe relationships as all communications
1035 would be encrypted end-to-end.



(a) One second notification interval.

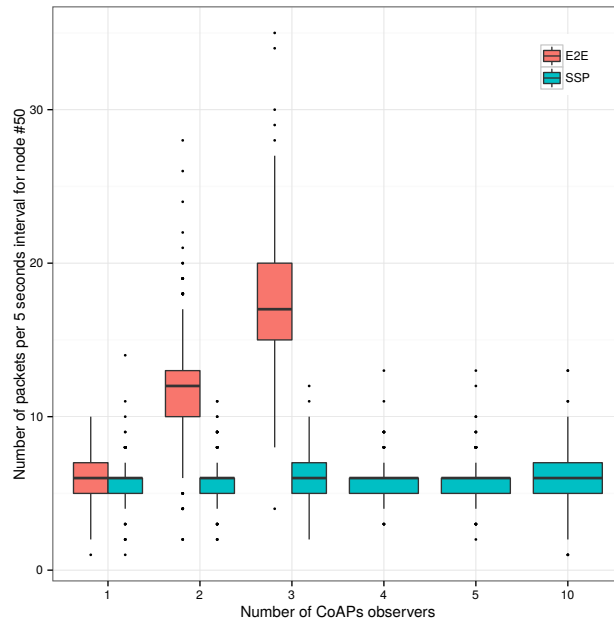


(b) Five seconds notification interval.

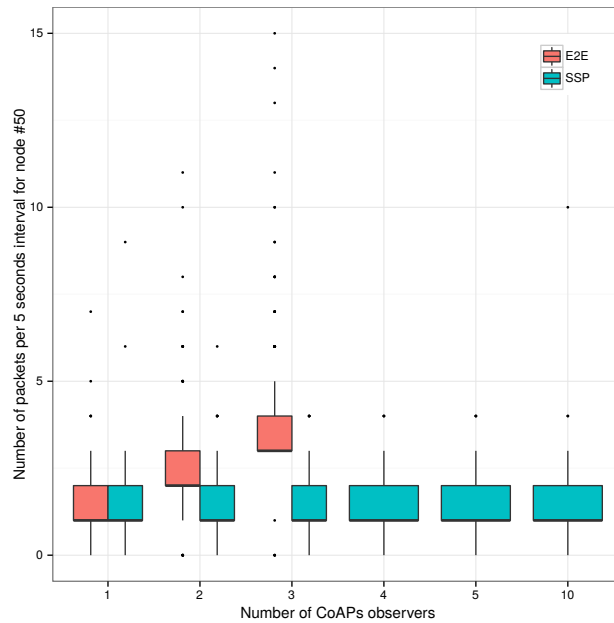
Figure 9. Total energy expenditure for node #50 per five seconds interval for end-to-end (E2E) CoAPs observation versus CoAPs observation through the Smart Service Proxy (SSP)

1036 7. Conclusions

1037 In this work we have presented the Secure Service Proxy: a CoAP(s) intermediary for use in
 1038 resource-constrained RESTful environments. It has been designed to provide scalable end-to-end
 1039 security for constrained devices and to extend constrained devices with additional functionality. The
 1040 presented work follows a reverse proxy approach, where the SSP hosts virtual devices on behalf
 1041 of resource-constrained devices. This approach enables the SSP to extend the virtual devices with
 1042 security features which are hard to attain in constrained environments such as authentication based



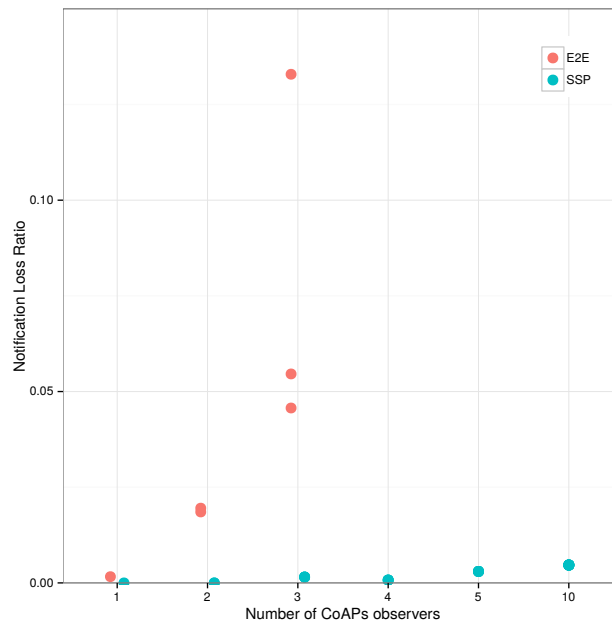
(a) One second notification interval.



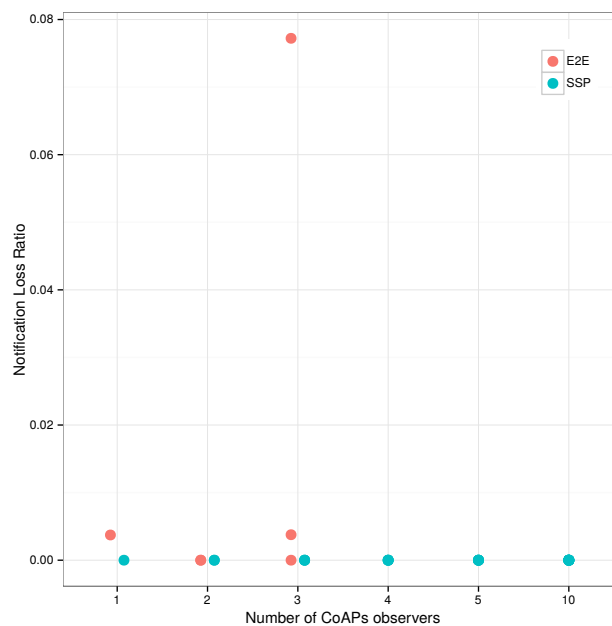
(b) Five seconds notification interval.

Figure 10. Number of exchanged packets for node #50 per five seconds interval for end-to-end (E2E) CoAPs observation versus CoAPs observation through the Smart Service Proxy (SSP)

1043 on public key infrastructure (which, inherently, scales better than using PSKs), perfect forward secrecy
 1044 and fine-grained authorization based on host identify and the nature of the request and resource.
 1045 Additionally, the SSP extends virtual devices with a variety of different functions by means of an
 1046 adapter chain system. Adapters are modular blocks of functionality that are hosted on the virtual
 1047 device. Examples include caching, static resource and congestion control adapters. The SSP hosts a
 1048 RESTful web interface for managing virtual devices and adapter chains.



(a) One second notification interval.



(b) Five seconds notification interval.

Figure 11. Notification loss ratios as measured at the CoAPs clients for end-to-end (E2E) CoAPs observation versus CoAPs observation through the Smart Service Proxy (SSP)

1049 The SSP has been evaluated in two different setups. First, tests were performed in a LLN
 1050 simulator to measure the effect of terminating end-to-end security on the SSP. The results of
 1051 the simulator study demonstrate that session termination combined with long-term sessions in
 1052 the constrained network, lead to significant savings in network traffic, communication delay and
 1053 processing and consequently lead to longer battery life. The second study was ran on a WSN
 1054 testbed and quantified the impact of aggregating multiple observe relations with a constrained device
 1055 over DTLS. The results confirm that the load on the constrained device and constrained network is

1056 independent of the number of observers. As a result, the packet rate and energy expenditure remain
1057 equal to those of the one observer case as the number of observers increases. Note that the session
1058 termination is a necessary condition for observe aggregation in case of DTLS-based security.

1059 In conclusion, the presented Secure Service Proxy breaks end-to-end security in order to offer
1060 security primitives that are hard to attain on constrained systems while reducing the load on resource
1061 constrained devices and networks.. Additionally, the proxy provides extra application-layer features
1062 on behalf of constrained devices to services, which are built on top of these devices. Combined, the
1063 proxy facilitates the integration of constrained RESTful environments in services; thereby furthering
1064 the vision of an open, secure and scalable Web of Things.

1065 **Author Contributions:** This paper is part of a Ph.D. Thesis written by Floris Van den Abeele under supervision
1066 of Jeroen Hoebeke, Ingrid Moerman and Piet Demeester. Additionally, Jeroen Hoebeke implemented parts of
1067 the proxy and also reviewed the paper.

1068 **Conflicts of Interest:** The authors declare no conflict of interest.

1069 Bibliography

- 1070 . Joseph Bradley.; Barbier, J.; Handler, D. Embracing the Internet of Everything To Capture Your Share of
1071 14.4 Trillion USD. *Cisco Ibsg Group* **2013**, p. 2013.
- 1072 . Miorandi, D.; Sicari, S.; De Pellegrini, F.; Chlamtac, I. Internet of things: Vision, applications and research
1073 challenges. *Ad Hoc Networks* **2012**, *10*, 1497–1516.
- 1074 . Gubbi, J.; Buyya, R.; Marusic, S.; Palaniswami, M. Internet of Things (IoT): A vision, architectural
1075 elements, and future directions. *Future Generation Computer Systems* **2013**, *29*, 1645–1660.
- 1076 . Bormann, C.; Ersue, M.; Keranen, A. Terminology for Constrained-Node Networks. ISSN 2070-1721, doi:
1077 10.17487/RFC7228, IETF, 2014.
- 1078 . Baronti, P.; Pillai, P.; Chook, V.W.C.; Chessa, S.; Gotta, A.; Hu, Y.F. Wireless sensor networks: A survey on
1079 the state of the art and the 802.15.4 and ZigBee standards. *Computer Communications* **2007**, *30*, 1655–1695.
- 1080 . Winter, T.; Thubert, P. RPL: IPv6 Routing Protocol for Low-Power and Lossy Networks, ISSN 2070-1721,
1081 doi: 10.17487/RFC6550, IETF, 2012.
- 1082 . Tschofenig, H.; Fossati, T. Transport Layer Security (TLS) / Datagram Transport Layer Security (DTLS)
1083 Profiles for the Internet of Things. ISSN 2070-1721, doi: 10.17487/RFC7925, IETF, 2016.
- 1084 . Vučinić, M.; Tourancheau, B.; Watteyne, T.; Rousseau, F.; Duda, A.; Guizzetti, R.; Damon, L. DTLS
1085 Performance in Duty-Cycled Networks. International Symposium on Personal, Indoor and Mobile Radio
1086 Communications (PIMRC - 2015), 2015, [1507.05810].
- 1087 . Shelby, Z.; Hartke, K.; Bormann, C.; Frank, B. Constrained Application Protocol (CoAP), ISSN 2070-1721,
1088 doi: 10.17487/RFC7252, IETF, 2014.
- 1089 . Kuladinithi, K.; Bergmann, O.; Becker, M. Implementation of CoAP and its Application in Transport
1090 Logistics. Proceedings of the Workshop on Extending the Internet to Low power and Lossy Networks,
1091 2011.
- 1092 . Hartke, K. Observing Resources in the Constrained Application Protocol (CoAP). ISSN 2070-1721, doi:
1093 10.17487/RFC7641, IETF, 2015.
- 1094 . Rescorla, E.; Modadugu, N. Datagram transport layer security version 1.2, ISSN 2070-1721, doi:
1095 10.17487/RFC6347, IETF, 2012.
- 1096 . Dierks, T.; Rescorla, E. The Transport Layer Security (TLS) protocol version 1.2. ISSN 2070-1721, doi:
1097 10.17487/RFC5246, IETF, 2008.
- 1098 . McGrew, D. An Interface and Algorithms for Authenticated Encryption. ISSN 2070-1721, doi:
1099 10.17487/RFC5116, IETF, 2008.
- 1100 . McGrew, D.; Bailey, D. AES-CCM Cipher Suites for Transport Layer Security (TLS). ISSN 2070-1721, doi:
1101 10.17487/RFC6655, IETF, 2012.
- 1102 . Eronen, P.; Tschofenig, H. Pre-shared key ciphersuites for transport layer security (TLS). ISSN 2070-1721,
1103 doi: 10.17487/RFC4279, IETF, 2005.
- 1104 . Wouters, P.; Tschofenig, H.; Gilmore, J.; Weiler, S.; Kivinen, T. Using Raw Public Keys in Transport Layer
1105 Security (TLS) and Datagram Transport Layer Security (DTLS). ISSN 2070-1721, doi: 10.17487/RFC7250,
1106 IETF, 2014.

- 1107 . Bailey, D.; Campagna, M.; Dugal, R.; McGrew, D. AES-CCM Elliptic Curve Cryptography (ECC) Cipher
1108 Suites for TLS. ISSN 2070-1721, doi: 10.17487/RFC7251, IETF, 2014.
- 1109 . Rescorla, E. TLS Elliptic Curve Cipher Suites with SHA-256/384 and AES Galois Counter Mode. ISSN
1110 2070-1721, doi: 10.17487/RFC4492, IETF, 2008.
- 1111 . Barker, E. NIST Special Publication 800-57 Part 1 Revision 4, Recommendation for Key Management Part
1112 1: General. Technical report, 2016.
- 1113 . Santesson, S.; Myers, M.; Ankney, R.; Malpani, A.; Galperin, S.; Adams, C. X. 509 Internet public key
1114 infrastructure online certificate status protocol-OCSP. ISSN 2070-1721, doi: 10.17487/RFC6960, IETF,
1115 2013.
- 1116 . Van den Abeele, F.; Vandewinckele, T.; Hoebeker, J.; Moerman, I.; Demeester, P. Secure communication
1117 in IP-based wireless sensor networks via a trusted gateway. IEEE Tenth International Conference on
1118 Intelligent Sensors, Sensor Networks and Information Processing (IEEE ISSNIP 2015), 2015.
- 1119 . Shelby, Z. Constrained RESTful Environments (CoRE) Link Format, ISSN 2070-1721, doi:
1120 10.17487/RFC6690, IETF, 2012.
- 1121 . Selander, G.; Mattsson, J.; Palombini, F.; Seitz, L. Object Security of CoAP (OSCOAP) (Work in progress),
1122 2017.
- 1123 . Raza, S.; Duquennoy, S.; Chung, T.; Yazar, D.; Voigt, T.; Roedig, U. Securing communication in 6LoWPAN
1124 with compressed IPsec. 2011 International Conference on Distributed Computing in Sensor Systems and
1125 Workshops (DCOSS), 2011, pp. 1–8.
- 1126 . Bonomi, F.; Milito, R.; Zhu, J.; Addepalli, S. Fog computing and its role in the internet of things.
1127 *Proceedings of the first edition of the MCC workshop on Mobile cloud computing - MCC '12* 2012, p. 13.
- 1128 . Nitti, M.; Pilloni, V.; Colistra, G.; Atzori, L. The Virtual Object as a Major Element of the Internet of Things:
1129 A Survey. *IEEE Communications Surveys & Tutorials* 2016, 18, 1228–1240.
- 1130 . Kovatsch, M.; Mayer, S.; Ostermaier, B. Moving application logic from the firmware to the cloud: Towards
1131 the thin server architecture for the internet of things. 2012 Sixth International Conference on Innovative
1132 Mobile and Internet Services in Ubiquitous Computing, 2012, pp. 751–756.
- 1133 . Jara, A.J.; Moreno-Sanchez, P.; Skarmeta, A.F.; Varakliotis, S.; Kirstein, P. IPv6 addressing proxy: mapping
1134 native addressing from legacy technologies and devices to the Internet of Things (IPv6). *Sensors (Basel,
1135 Switzerland)* 2013, 13, 6687–712.
- 1136 . Ludovici, A.; Calveras, A. A Proxy Design to Leverage the Interconnection of CoAP Wireless Sensor
1137 Networks with Web Applications. *Sensors* 2015, 15, 1217–1244.
- 1138 . Castellani, A.; Loreto, S.; Rahman, A.; Fossati, T.; Dijk, E. Guidelines for Mapping Implementations:
1139 HTTP to the Constrained Application Protocol (CoAP). ISSN 2070-1721, doi: 10.17487/RFC8075, IETF,
1140 2017.
- 1141 . Mingozzi, E.; Tanganelli, G.; Vallati, C. CoAP Proxy Virtualization for the Web of Things. 2014 IEEE 6th
1142 International Conference on Cloud Computing Technology and Science. IEEE, 2014, pp. 577–582.
- 1143 . Tanganelli, G.; Vallati, C.; Mingozzi, E.; Kovatsch, M. Efficient proxying of CoAP observe with quality of
1144 service support. 2016 IEEE 3rd World Forum on Internet of Things (WF-IoT). IEEE, 2016, pp. 401–406.
- 1145 . Farris, I.; Lera, A.; Molinaro, A.; Pizzi, S. A CoAP-compliant solution for efficient inclusion of RFID in the
1146 Internet of Things. 2014 IEEE Global Communications Conference. IEEE, 2014, pp. 2795–2800.
- 1147 . Hummen, R.; Shafagh, H.; Raza, S. Delegation-based Authentication and Authorization for the IP-based
1148 Internet of Things. 11th IEEE International Conference on Sensing, Communication, and Networking
1149 (SECON'14), 2014.
- 1150 . Park, J.; Kwon, H.; Kang, N. IoT–Cloud collaboration to establish a secure connection for lightweight
1151 devices. *Wireless Networks* 2017, 23, 681–692.
- 1152 . Garcia-Morchon, O.; Keoh, S.L.; Kumar, S.; Moreno-Sanchez, P.; Vidal-Meca, F.; Ziegeldorf, J.H. Securing
1153 the IP-based internet of things with HIP and DTLS. *Proceedings of the sixth ACM conference on Security and
1154 privacy in wireless and mobile networks - WiSec '13* 2013, p. 119.