ARTICLE IN PRESS

Computer Communications xxx (2016) xxx-xxx

[m5G;March 4, 2016;11:33]

compute: communications



Contents lists available at ScienceDirect

Computer Communications

journal homepage: www.elsevier.com/locate/comcom

Towards automatic protocol field inference

Ignacio Bermudez^{a,*}, Alok Tongaonkar^a, Marios Iliofotou^b, Marco Mellia^c, Maurizio M. Munafò^c

^a Symantec Corporation, 350 Ellis St., Mountain View, CA 94086, USA

^b Caspida, 2100 Geng Road#100, Palo Alto, CA 94303, USA

^c Politecnico di Torino, Corso Duca degli Abbruzzi 24, Turin 10129, Italy

ARTICLE INFO

Article history: Received 21 September 2015 Revised 20 February 2016 Accepted 23 February 2016 Available online xxx

Keywords: Protocol reverse-engineering ICS/SCADA Networks Payload based anomaly detection

ABSTRACT

Security tools have evolved dramatically in the recent years to combat the increasingly complex nature of attacks. However, these tools need to be configured by experts that understand network protocols thoroughly to be effective. In this paper, we present a system called *FieldHunter*, which automatically extracts fields and infers their types. This information is invaluable for security experts to keep pace with the increasing rate of development of new network applications and their underlying protocols. FieldHunter relies on collecting application messages from multiple sessions. Then, it performs field extraction and inference of their types by taking into consideration statistical correlations between different messages or other associations with meta-data such as message length, client or server IP addresses. We evaluated FieldHunter on real network traffic collected in ISP networks from three different continents. FieldHunter was able to extract security relevant fields and infer their types for well documented network protocols (such as DNS and MSNP) as well as protocols for which the specifications are not publicly available (such as SopCast). Further, we developed a payload-based anomaly detection system for industrial control systems using FieldHunter. The proposed system is able to identify industrial devices behaving oddly, without any previous knowledge of the protocols being used.

© 2016 Published by Elsevier B.V.

1 1. Introduction

02

In recent years attacks against networks have become more 2 complicated. To defend against these complex attacks, network se-3 curity systems have also evolved to use more sophisticated mech-4 5 anisms. For instance, firewalls have moved from using simple packet-filtering rules to using application level rules that need 6 deeper understanding of the protocols being used by network ap-7 plications. Similarly, intrusion detection systems are increasingly 8 using vulnerability based signatures [1] that contain information 9 10 specific to network protocols. Access control mechanisms are also evolving from IP address based policies to fine-grained policies 11 12 which use protocol objects such as users and message types.

13 It is clear that configuring all of the above applications requires 14 a deeper understanding of the network protocols, which is done 15 through reading protocol specifications. However, comprehending 16 protocol specifications is a very tedious task. Moreover, many of

http://dx.doi.org/10.1016/j.comcom.2016.02.015 0140-3664/© 2016 Published by Elsevier B.V. the proprietary protocols specifications are not publicly available. 17 The traditional approach of manual reverse engineering a proto-18 col cannot cope with the rate at which new benign or malicious 19 applications are made available and brought into workplace. As a 20 result, security administrators have to configure security applica-21 tions with very limited visibility into the network protocol space; 22 thus adversely affecting the efficacy of these tools in securing the 23 network. 24

The above technology challenge has led to a growing interest 25 in the research community in the development of techniques for 26 automating the reverse-engineering process for extracting proto-27 col specifications, which consists of inferring message formats and 28 underlying protocol state machines. The state-of-the-art techniques 29 can be classified in two categories: reverse-engineering through bi-30 nary code analysis [2–6] and from network traffic [7–13]. In this 31 work, we present an automatic reverse-engineering system of the 32 second category, i.e. it infers protocol specifications from just net-33 work traffic data. Reverse-engineering using network traffic has an 34 advantage over techniques using binary analysis, because applica-35 tion binaries are not always available to the security operators. 36

Our approach to the problem of protocol reverse engineering 37 aims to extract field boundaries and field protocol types from network traces that belong to the protocol. As compared to previous 39

Corresponding author. Tel.: +1 4084665710.

E-mail addresses: ignaciobermudez@gmail.com, ignacio_bermudezcorr@ symantec.com (I. Bermudez), alok_tongaonkar@symantec.com (A. Tongaonkar), marios@caspida.com (M. Iliofotou), marco.mellia@polito.it (M. Mellia), maurizio.munafo@polito.it (M.M. Munafò).

ARTICLE IN PRESS

I. Bermudez et al. / Computer Communications xxx (2016) xxx-xxx

works in this area, we are able to extract richer protocol informa-40 41 tion in terms of (i) extracting diverse field types, and (ii) handling 42 binary and textual protocols in an uniform framework. We study 43 well known protocols and identify a set of field types that can be used in a multitude of security applications. We focus on identify-44 ing: (i) Message Type (MSG-Type), such as flags in DNS protocol or 45 GET/POST keywords in HTTP, (ii) Message Length (MSG-Len), usu-46 ally found in TCP protocols to delimit application messages in a 47 48 stream, (iii) Host Identifier (Host-ID) such as Client ID and Server ID, (iv) Session Identifier (Session-ID) such as cookies, (v) Transac-49 50 tion Identifier (Trans-ID) such as sequence/acknowledgment num-51 bers, and (vi) Accumulators such as generic counters and times-52 tamps. We note that a protocol may not have all the above types 53 of fields.

We built a system called FieldHunter, that uses a two step 54 methodology: (i) Field extraction: here we extract fields from the 55 protocol messages. (ii) Field type inference: here we infer the type 56 of the fields extracted in the previous step. The key contribution 57 of our work is the development of various heuristics based on ob-58 served statistical properties for inferring the different field types. 59 In our evaluation, we used real network traces from three different 60 Internet Service Providers (ISPs) to validate the ability to extract 61 62 various field types from well known protocols such as Real Time 63 Protocol (RTP), as well as protocols without any publicly available 64 specification such as SopCast's protocol.

Next, to illustrate the use of FieldHunter in building end-to-65 end security applications, we developed a payload-based anomaly 66 67 intrusion detection system for industrial control systems. Industrial Control Systems (ICS) encompass several types of control sys-68 tems used in industrial production, including Supervisory Control 69 and Data Acquisition (SCADA) systems, Distributed Control Sys-70 71 tems (DCS), and other smaller control system configurations such 72 as Programmable Logic Controllers (PLC). Due to their use in the industrial sectors and critical infrastructures, they are a prime tar-73 get for attackers and the cost of successful attacks is tremendous 74 75 to the victims. Typically, the attackers targeting ICS are sophisticated, and have a lot of resources; sometimes even being state 76 77 sponsored. They are able to avoid detection by traditional defense mechanisms [14]. To make matters worse, there has been a trend 78 of increasing number of attacks on critical infrastructure as evi-79 denced by the rapid increase in the number of reported attacks on 80 **0**3 81 ICS from 91k in 2012 to over 675k in 2014 [15]. To the best of our knowledge, our system is the first payload-based anomaly de-82 83 tection system that handles legacy proprietary protocols commonly 84 used in ICS networks.

The rest of the paper is organized as follows. Section 2 defines the terminology used throughout the paper, Section 3 provides details about the core algorithms used by FieldHunter. Performance evaluation and parameter tuning are presented in Section 4. We describe the anomaly detection system for ICS in Section 5. We discuss about assumptions and limitations in Section 6, related works in Section 7 and finally conclude this work in Section 8.

92 2. Terminology

Fig. 1 shows a pictorial representation of the terminology used throughout this work. Our system uses as input a set of *conversations* ¹ of a particular application. We refer to such a set as *collection*. Conversations consist of exchanged *messages* between two hosts. Messages from client to server are denoted as C2S (darkcolored) and from server to client as S2C (light-colored). We consider the initiator of the conversation as the client, the other end



Fig. 1. Terminology diagram.

as the server, and identify hosts by their IP address. Messages consist of different pieces of information enclosed in *fields*. As we show in Fig. 1, conversations evolve horizontally over time (t) and messages can be compared vertically across multiple conversations.

To enable the analysis of a collection, the messages in the 105 conversations can be grouped together in the following ways: (i) 106 Grouping messages based on their position in conversations, e.g., 107 all third messages in C2S direction. (ii) Grouping together all the 108 messages of a conversation. This essentially captures session-like 109 information. (iii) Grouping together messages by direction, e.g., all 110 C2S messages. We note that (i) and (ii) are very similar to vertical 111 and horizontal sub-collections as defined by Kreibich et al. [16]. 112 Message grouping is instrumental for FieldHunter to find patterns 113 in the collections. If these groups do not contain enough message 114 diversity, FieldHunter cannot unveil the field types it is designed 115 for. 116

It is worth mentioning that the formation of protocol collec-117 tions used by FieldHunter is beyond the scope of this work. How-118 ever, we suggest two alternatives for the same. One way is to use 119 a test-bed in which the application is executed while the traffic 120 exchanged is being captured. Alternatively, the collection can be 121 extracted from passive observation of actual traffic by the means 122 of network classifiers, i.e., by filtering all conversations involving 123 a well-known port (see Section 5), or by relying on a behavioral 124 traffic classifier classifier [17]. 125

Application conversations are transported by TCP/UDP segments 126 and are extracted by FieldHunter using the following methodology: 127 (i) for messages transported over UDP it is assumed that each seg-128 ment contains one application message, and (ii) for TCP it is as-129 sumed that TCP PUSH flags delimits the beginning of a new appli-130 cation message from the end of another one. An accurate message 131 extraction can be done once the MSG-Len field has been identified 132 by FieldHunter. 133

We make a distinction between textual and binary protocols as follows: Textual protocols use human readable words and symbols to structure data, and they look more like a text document. Exponents of textual protocols are HTTP and SMTP. On the other hand, binary protocols encode data on bits rather than symbols and the way that data is structured is quite rigid. Examples of binary protocols are DNS and DNP.

3. Design

In this section, we describe the system design and discuss the two components of FieldHunter(i) Field Extractor (ii) Field Type Inference Engine. These components are run in sequence to obtain 144

141

Please cite this article as: I. Bermudez et al., Towards automatic protocol field inference, Computer Communications (2016), http://dx.doi.org/10.1016/j.comcom.2016.02.015

¹ A conversation is formed of the two flows in opposite directions, where a flow is defined by the 5-tuple (Layer-4 Protocol, Source IP, Source Port, Destination IP, Destination IP).

Protocol Coll. Field Extractor ary Proto Field Semantic ey-Value 1:T User ID Protocol Coll. emai ί 1:T dest Server ID а 2:T 300 Cmd. ID 2:T 404 Cmd. ID В Protocol Coll. n-1:B [0:8] Length n-1 okenize n-1·B [8:32] User ID Protocol Coll. n:B [16:24] Cmd. ID

Fig. 2. FieldHunter system diagram.

Table 1

Common text-based protocols and their observed delimiters. GAME: team fortress (game), TEL: Telnet, CS: counter strike (game), GNU: Gnutella.

Prot.	D_f	$D_{k-\nu}$	Prot.	D_f	$D_{k-\nu}$
HTTP SMTP POP3 RTSP SIP GAME TEL	\r\n \r\n \r\n \r\n 0x00 \r\n	9, (, 9, (, 9, 9) 9, (, 9, (, 0x00 9, (,	FTP TFTP CS GNU RTP MSN	\r\n 0x1D 0x5C \r\n \r\n \r\n	0x1E 0x5C , ,

145 a field summary report (which describes the identified fields and their types) as shown in Fig. 2. 146

3.1. Field extractor 147

Textual and binary protocols differ greatly in the way fields are 148 149 delimited. Textual protocols typically use delimiters such as ":" or 150 "\r\n" (carriage-return and line-feed pair) to separate fields. On the other hand in binary protocols, fields either have fixed offset and 151 152 size or offsets and lengths that are specified in some preceding fields. Hence we have developed different techniques for textual 153 and binary protocols. Next, we explain each of these techniques 154 separately. 155

156 3.1.1. Textual protocols

Field extraction for textual protocols boils down to identifying 157 158 field delimiters. However, this is a non-trivial task as many protocols use multiple delimiters for different purposes. For instance, 159 consider a message such as TIME-OUT: 60 # PORT: 54001. 160 In this message, "#" is used to separate out the fields, while ":" 161 is used to separate out the key and the value in a field. Hence, we 162 categorize delimiters into two types: (i) Field delimiter (D_f) : sepa-163 rates the different fields of a message, e.g., the "#" character in the 164 above example. (ii) Key-value delimiter $(D_{k-\nu})$: separates the key 165 from its corresponding value, e.g., the ":" in the same example. 166

Commonly used D_f and $D_{k-\nu}$ delimiters are shown in Table 1. 167 These delimiters are obtained from the documentation of the listed 168 protocols, and are actually observed in our data sets. As we see, 169 170 there are popular delimiters, such as $r\n$, as well as non-standard delimiters, such as 0x00 (null), 0x1D, and 0x5C. 171

Generally speaking, FieldHunter identifies delimiters using 172 three key observations: (i) Delimiters are non-alphanumeric se-173 quences of 1 or 2 characters. (ii) Delimiters have a high horizon-174 175 tal and vertical frequency compared to other non-alphanumeric 176 sequences in a textual protocol. (iii) There is only one D_f 177 that splits up the messages into key-value pairs (UID: 1234, Content-length: 872) or singleton keywords fields (HELO, 178 LOGOUT, OK, FAILED). FieldHunter first identifies D_f and then 179 180 proceeds with the $D_{k-\nu}$ if fields are key-value paired.

Field delimiter inference. FieldHunter finds frequent sequences of 181 non-alphanumeric characters in the protocol which are considered 182

to be delimiter candidates d. Then from among all the candidates it 183 chooses only one $(D_f = d)$, such that it splits up any protocol mes-184 sage into valid key-value pairs and singletons. Validity of key-value 185 pairs and singletons is checked by comparing common prefixes and 186 exact matches respectively. 187

Key-value pair delimiter inference. Once D_f has been detected, mes-188 sages are split into fields from which we need to identify key-189 values along with $D_{k-\nu}$, and singletons. The identification of $D_{k-\nu}$ 190 is performed in three steps: (i) FieldHunter clusters fields of the 191 same type by using the Longest Common Prefix (LCP); (ii) by 192 re-clustering the clusters, FieldHunter cleans up possible outliers 193 caused by two or more keywords sharing a common prefix. E.g., 194 Port: and Point: have Po in common, and finally (iii) we 195 choose the $D_{k-\nu}$ as the non-alphanumeric suffix part of the LCP of 196 each group. In the case that all the LCPs are identical for a group, 197 then we say that the field contained by the group is a singleton 198 and we do not search for a delimiter. 199

3.1.2. Binary protocols

In binary protocols, fields represent serialization of variables as 201 they are structured in memory. To parse these fields, message re-202 cipients need to know the structure of the data, i.e. the offset and 203 length of the fields. The challenge for FieldHunter is that the mes-204 sage data structure is initially unknown. Therefore messages are 205 split into n-grams which are used by Field Type Inference Engine. 206 We observe that for most of the field types, the n-grams forming 207 the field also show similar characteristics to the field. For instance, 208 in a protocol that has a 32-bit Host ID field, the four 8-bit n-grams 209 also exhibit similar statistical properties as Host ID. In such cases, 210 we identify the field type for the single n-gram and then check 211 whether consecutive n-grams can be merged into a larger field of 212 the same type. 213

We note that this assumption does not always hold. For in-214 stance, a 32-bit Accumulator field may increment by one every 215 time. But given the number of samples that we may consider in 216 our collection (say order of thousands), the most significant bits 217 may show up as constants and not accumulators. This issue is cir-218 cumvented for fields such as Message Length and Accumulators 219 (numerical representations) by considering n-grams of larger size 220 first, say 32-bit n-gram, and then iteratively reducing n-gram size till the whole n-gram fits the field. Moreover, we handle byteendianness for fields that contain numerical representations by re-223 peating the heuristics separately for both little-endianness and big-224 endianness. This is not the case for fields that can be interpreted 225 as categorical representations. 226

3.2. Field type inference engine

Our approach is based on the following key observation: 228 Fields with different types change differently over specific sub-229 collections. For instance, a field that consistently takes a distinct 230 value for each IP address may represent a Host-ID. Similarly, fields 231 that increment by one over sequential messages of a conversation 232 may be part of a message counter. 233

FieldHunter assigns types to fields by using different statistical 234 tests that are further explained. The techniques for FTI are similar 235 for both textual and binary protocols. In the rest of the paper we 236 use the term "n-gram" to interchangeably to mean "binary n-gram" 237 or "textual field" for ease of exposition. For example, when we 238 state "n-gram entropy is computed", we actually mean that either 239 "binary n-gram entropy is computed" or "textual field entropy is 240 computed". We use specific statistical tests based on different as-241 sociations between observed variables to infer different field types. 242 The association between two variables (a, b) can be of the follow-243 ing types: (i) "numerical correlation" $(a \Leftrightarrow b)$, e.g., message length 244

200

221 222

I. Bermudez et al./Computer Communications xxx (2016) xxx-xxx



Message Type Field (MSG-Type)

Fig. 3. MSG-Type (left), MSG-Len (center) and Trans-ID (right) modules.

field is numerically correlated to the observed length of the mes-245 sage, (ii) "categorical correlation" ($a \in A \Leftrightarrow b \in B$), e.g., user IDs cor-246 relate categorically with IP addresses and (iii) "causality correla-247 tion" $(a \Rightarrow b)$, e.g., certain type of message will result in a particular 248 response from server. 249

250 The labeling process works by making a hypothesis that a given field is of a certain type. When the hypothesis holds, i.e., the field 251 exhibits the statistical behavior of the field type, FieldHunter la-252 bels the field as such. We note here that a field may be labeled 253 as multiple field types. For instance, an acknowledgment number 254 field could be labeled as both Transaction ID as well as an Accu-255 256 mulator.

257 In Fig. 3 the more complex heuristics are illustrated using 258 block-diagrams. Blocks in the diagrams represent different tests; 259 horizontal/vertical arrow inside a block defines horizontal/vertical 260 sub-collection analysis and thresholds are highlighted in italic. 261 More details on parameter selection are given in Section 4.4.

3.2.1. Message type (MSG-Type) 262

263 MSG-Type contains information about the underlying protocol 264 state machine and its values represent the semantic of the whole 265 message. Thus, the content of MSG-Type field is used by the re-266 ceiver to understand what type of message is received, for example a request, a status update, or an error message. 267

268 The process of finding MSG-Types is based on two key observa-269 tions: (i) MSG-Type takes values from a well defined small static set; and (ii) represents transitions in an underlying protocol state 270 machine. Hence, by pairing request/response messages, there is a 271 high probability that their corresponding MSG-Type fields are re-272 lated. The leftmost diagram in Fig. 3 describes the MSG-Type la-273 274 beling process.

Using observation (i) above, FieldHunter first looks for n-grams 275 276 that vertically are neither random nor constant. Randomness of a n-gram x can be measured using the entropy H(x) metric. Let p_i 277 be the probability of having the n-gram take the value *i*. Then 278 $H(x) = -\sum_{i} p_i \log_2 p_i$, where $0 \cdot \log(0) = 0$. By definition for 1-byte 279 n-grams (8-bits) H(x) takes values between 0 (constant) and 8 280 (perfectly random). Then n-grams that are unlikely to be part of 281 a MSG-Type field are discarded. Once some candidate fields are 282

identified, according to observation (ii), we check for n-grams that 283 have a causal relationship with n-grams in the response messages. 284 Here FieldHunter uses categorical correlation metric. Towards this 285 end, FieldHunter measures causality using the information theo-286 retic metric I(q; r)/H(q), where I(q; r) = H(q, r) - H(q|r) - H(r|q) is 287 the mutual information, that measures the information shared by 288 a request (*Q*) and a response (*R*) [18]. 289

FieldHunter takes n-grams for which causality is greater than a 290 threshold, say 0.8, as MSG-Type candidates. For the case of binary 291 protocols, if multiple n-grams are candidate, these are grouped to-292 gether and causality is checked again. Thus, if a group coincides 293 with the actual MSG-Type field, then the whole candidate group 294 should also satisfy the initial hypothesis of causality. For example, 295 suppose n-grams at byte offset 1, 5, 6 show a large causality such 296 that $q1 \Rightarrow r1$, $q5 \Rightarrow r5$, $q6 \Rightarrow r6$. Then we check whether the groups 297 $(q1, q5, q6) \Rightarrow (r1, r5, r6)$ holds the causality. If this holds, the field 298 containing n-grams at offsets (1, 5, 6) are returned as the MSG-299 Type field. 300

3.2.2. Message length (MSG-Len)

Our goal here is to find fields that contain indication of the application message length. As such, we expect the MSG-Len field is 303 linearly correlated with the actual physical message size. We use 304 two different tests for identifying linear correlations in order to 305 have higher confidence on our results. 306

The complete MSG-Len test algorithm is depicted in the central 307 diagram in Fig. 3. This heuristic does not use the typical 1-byte n-308 gram and for textual protocols it decodes the content of the field 309 as a number. The reason why 1-byte n-grams do not provide good 310 results is that Most Significant Byte and Least Significant Byte are 311 not correlated in this case. Hence, FieldHunter iteratively selects n-312 gram windows of size 32, 24, 16-bits that are shifted at a step of 313 8-bits. Such windows sizes are the standard sizes used to represent 314 integers in computer memory. At each iteration Pearson correla-315 tion coefficient tells whether the numeric values of the fields are 316 associated with the length of the messages. Notice that the compu-317 tation of this correlation could be affected by biases due to some 318 popular messages in the collection of the same size, for instance, if 319 950 out of 1000 of collection messages are 40 bytes long. We avoid 320

ARTICLE IN PRESS

5



Fig. 4. n-gram correlation with MSG-Len for SopCast.

such biases by stratifying messages by length, creating in this way
a size-heterogeneous collection not affected by the bias problem.
We select all the fields for which the coefficient is above a certain
threshold as MSG-Len candidates. We empirically found 0.6 to be
a good threshold.

326 Fig. 4 shows the results of applying the Pearson correlation to 327 the SopCast protocol collection obtained from one of our traces. In this example, we use 16-bit n-gram. Pearson coefficient values 328 span from zero to one, where zero indicates no correlation and one 329 represents a strong correlation. In Fig. 4 there are two clear spikes, 330 331 one at offset 88-bit and the other at 168-bit that suggests the presence of a MSG-Len field (see Section 4.2). We cross-verified these 332 results using DPI signature rules for UDP SopCast found in OpenDPI 333 334 [19], an open source packet inspection engine.²

335 Once the candidates are found, the next step is to conduct a test 336 to verify that the candidates indeed are carrying information regarding the length of the message. The hypothesis is that the mes-337 sage length expresses the length of the message in an unit of mea-338 surement, such as bytes or words, and that it describes the length 339 340 of data starting from a given byte offset. In other words, we state that the message length is ruled by the following linear equation: 341 $MSG_{len} = a \cdot FIELD_{value} + b$, such that $MSG_{len} \in \mathbb{N}$ is the observable 342 message length, FIELD_{value} is the value taken by candidate field, 343 a > 0 accounts for the unit of measurement and $b \in \mathbb{N}$ is the start-344 345 ing offset of the data described by the field. To verify the assumption, the linear equation is solved and (a, b) are obtained. This 346 347 process is repeated taking all possible message pairs with differ-348 ent lengths. Finally a candidate is considered as a true MSG-Len field if for most of the pairs (> 90%) the solution is acceptable 349 350 $(a > 0 \land b \in \mathbb{N}).$

351 3.2.3. Host identifier (Host-ID)

Host Identifiers are used to identify a particular host or device beyond the boundaries of the local network. For instance, in peerto-peer applications, the "Peer-ID" field uniquely identifies a specific peer/host in the whole overlay, even when the peer is behind a Network Address Translation (NAT) device or is moving over multiple networks.

The heuristic assumes that all the messages sent by the same 358 host carry the same Host-ID, i.e., for a given source IP, mes-359 360 sages are likely to have the same Host-ID. Then Host-ID should be strongly correlated with the IP address of the sender. Based on 361 this assumption, FieldHunter computes the categorical correlation 362 $R(x, y) = I(x; y)/H(x, y) \in [0, 1]$ of n-grams x with the sender IP ad-363 dress y, where H(x, y) is the joint entropy (that measures the to-364 365 tal amount of information that x and y jointly carry). That is, for 366 each $x \in X$, there is a different $y \in Y$, and vice-versa. N-grams with correlation coefficient greater than certain threshold, say 0.9, are 367 368 selected as candidates. Finally, consecutive candidate n-grams are merged into fields of at least a length of 4 bytes. Notice that the 369







Fig. 6. The n-gram entropy for Vuze DHT over a C2S vertical sub-collection.

adoption of statistical tests, such as correlation, makes the algorithms robust to handle noise in the data, such as when NAT is used. 370

Fig. 5 shows the categorical correlation between n-grams in a 373 vertical collection and the corresponding source IP address for the 374 Vuze DHT collection [20]. Note how R(x, y) is very close to one 375 (high correlation) for n-grams that represent the Client Address 376 and the Client-ID. However, we also observe that the first n-grams 377 of the Session-ID are also correlated with the sender IP address. 378 The explanation for this protocol peculiarity is found in the Vuze's 379 specification. Vuze's Session-ID is an application's global counter 380 randomly initialized at the start-up and incremented by 1 for each 381 new conversation. Hence, the most significant bits in the Session-382 ID are likely to be the same for all messages sent by the same 383 sender. By imposing a minimum length constraint, FieldHunter can 384 discard such fields. 385

3.2.4. Session Identifier (Session-ID)

Session Identifier keeps track of application-level sessions that 387 span over multiple conversations. Semantically, it is similar to the 388 use of Cookies in HTTP. Since the Session-ID remains constant 389 between a pair of endpoints, FieldHunter correlates the n-grams to 390 the pair of client and server IP addresses. Then we proceed using 391 the same categorical correlation as we do for Host-ID. 392

3.2.5. Transaction Identifier (Trans-ID)

The algorithm we use to detect Trans-IDs is illustrated in the 394 rightmost diagram in Fig. 3. It is assumed that Trans-ID are ran-395 domly picked by the transaction creator and then copied back in 396 the replies. Therefore, we first search for n-grams that appear ran-397 dom across both vertical and horizontal collections. Randomness is 398 measured using entropy as before. 399

Fig. 6 shows the entropy of n-grams for the Vuze DHT protocol. The figure shows the entropy of the first 36 n-grams (reported on the *x*-axis at the corresponding offset) in the C2S vertical subcollection. On the top, the protocol field names are reported as extracted from documentation. In this example, n-grams with high entropy are good candidates for the Trans-ID field.

Next, all consecutive request/response messages are paired and 406 for each of them, it is checked whether the n-grams/fields take 407 the same values. If the check passes, then the pair of n-grams are 408

386

393

ARTICLE IN PRESS

[m5G;March 4, 2016;11:33]

I. Bermudez et al./Computer Communications xxx (2016) xxx-xxx

added to a set of Trans-ID candidates. Note that request/response
message formats can change and Trans-ID may appear at different
offsets (for instance in Vuze DHT). Therefore, the heuristic does not
assume the protocol message formats are the same in both directions.

Finally, FieldHunter measures the consistency of these candidates over all the conversations, i.e., n-gram candidates with enough support, say >0.8, are finally marked as such. Minimum support allows some degree of mismatch, for example, caused by message reordering or re-transmission in the collection. Finally, consecutive n-grams are merged to form a field of at least 2 bytes. For textual protocols such n-gram merging is not needed.

421 3.2.6. Accumulators

Accumulators are fields that have increasing values over consec-422 423 utive message within the same conversation. These fields typically represent message sequence numbers, acknowledgment numbers 424 or timestamps. To identify such fields, we calculate the difference, 425 denoted as Δ , between values of n-grams in two subsequent mes-426 sages. We expect Δ to be positive and fairly constant. Notice that 427 differences are not required to be perfectly constant. For instance 428 a byte-wise counter in a protocol of variable size messages would 429 430 have variable Δ .

We search for accumulators in C2S and S2C directions indepen-431 432 dently of each other. As with the MSG-Len field, here we start with 433 fixed size n-grams. We assume accumulators are encoded in fields of a given field length, for example, 64, 32 and 16 bits. For each 434 435 field offset, we compute the vector of increments (Δ) considering each consecutive message pair in each conversation. In order to 436 437 use one threshold that captures the variations among Δs of different scales (such as sequential counters vs millisecond timers), we 438 compress Δ using a logarithmic function; $\dot{\Delta} = \ln \Delta$. Next, we an-439 alyze $\hat{\Delta}$ and select those that have relatively low entropy, i.e., $\hat{\Delta}$ 440 looks "fairly constant". 441

442 3.3. Field summary

FieldHunter provides information of the field type extracted automatically out of protocols as the final result. It provides two separate reports (corresponding to each direction of messages) for each protocol. The report contains the set of fields for which the types have been inferred. Note that we may not identify the type for some of the fields and will skip them in the report.

449 **4. Evaluation of FieldHunter**

450 In this section we evaluate the efficacy of FieldHunter in infer-451 ring protocol specification of known protocols. We use TSTAT [17], 452 a DPI tool that classifies traffic and feeds FieldHunter with proto-453 col collections. In general, each collection presents different char-454 acteristics. For instance, some may contain wrongly classified flows 455 caused by DPI false positives. Other may present little diversity, for example, showing only conversation exchanged with a handful 456 of servers. Different traces generate different collections that are 457 separately analyzed (for cross verification purpose). We consider a 458 protocol collection as valid only if it has at least 200 conversations 459 460 for textual or 2,000 for binary protocols; see Section 4.4 for more 461 details.

The subset of protocols for which we present results are summarized in Tables 3 and 4 for binary and textual protocols, respectively. Both straightforward and challenging cases are considered in our evaluation.

466 4.1. Datasets

467 We evaluate FieldHunter using three different ISP traces 468 (Table 2). Data was collected from different geographical regions

Summary of the traces we use.

Name	Location	Network location	Date	Duration (h)
TR1-2012	Europe	Edge	04–2012	24
TR2-2009	S. America	Backbone	10-2009	4
TR3-2007	Asia	Backbone	01–2007	7

Table 3

Summary of the results from running FieldHunter on the binary-based protocols.

Protocol	Discovered/GT [bits]		Cov/AoC	
	C2S	S2C	C2S	S2C
Vuze DHT	288/240	200/208	0.87/1	0.85/0.87
DNS	48/32	56/32	0.75/1	1/1
uTP	88/96	200/96	0.75/1	0.67/0.87
RTP	80/88	80/88	0.82/1	0.82/1
ED2K	128/16	16/16	1/1	1/1
KADEMLIA	352/16	104/16	1/1	1/1
STUN	256/160	184/160	0.9/0.83	0.85/0.88
SOPCAST	128/?	152/?	?/?	?/?
PPLIVE	0/?	32/?	?/?	?/?

Table 4

Summary of the results from running FieldHunter on the textual protocols.

Protocol	#Fields	K-V	CMD	IDs	FP-IDs
	C2S/S2C	C2S/S2C	C2S/S2C	C2S/S2C	C2S/S2C
STUN	3/3	2/2	1/1	1/1	0/0
FTP	19/18	12/17	7/1	2/1	0/0
HTTP	9/14	9/14	0/0	3/0	1/0
POP3	9/28	5/24	4/4	2/0	0/0
SMTP	19/9	15/9	5/0	1/1	0/0
MSNP	3/4	3/4	0/0	2/0	0/0
RTSP	9/25	9/18	0/7	3/6	0/2
GAME	*/17	*/15	*/2	*/2	*/0
RSP	3/*	2/*	1/*	1/*	0/*

(Asia, Europe, and South America), between the years 2007 to 469 2012. All traces contain full payload from network connections. 470 Given the large size of the TR1-2012 trace we limit the payload 471 per connection to the first 1048 bytes³. Although, all our parameter selection is made using TR1-2012, we tested FieldHunter on all three traces. 474

475

4.2. Evaluation of binary protocols

Table 3 reports the number of Discovered and Ground-Truth476(GT) bits, for both C2S and S2C collections, the Coverage (Cov) and477the Accuracy over Coverage (AoC). Cov is the ratio between Discov-478ered bits of the GT and the GT bits; while AoC is the ratio between479correctly discovered bits over the total number of discovered bits.480

The first seven protocols in the table have known specifications, 481 and the latter two do not have. Note that for some protocols, the number of discovered bits is larger than the GT bits. This happens because many protocols carry other protocols on top of them, such as ED2K. FieldHunter does not differentiate the header from the protocol's payload, resulting in identifying the fields of the transported (inner) protocol as well. 481

The average Accuracy over Coverage (AoC) is 0.83 in the worst 488 case. We observe that typical inaccuracies are due to the Accumulator type. For counters that span over large fields (for instance, a 490

Please cite this article as: I. Bermudez et al., Towards automatic protocol field inference, Computer Communications (2016), http://dx.doi.org/10.1016/j.comcom.2016.02.015

³ We did not observe this to cause any notable problems. Only for some protocols with long payloads, such as HTTP, portions of the payload and rarely portions of the application-layer header were not fully captured.

ARTICLE IN PRESS

32-bit long number), FieldHunter easily identifies the less significant bits, but tends to miss the most significant ones, which are
identified as "constant". From the results shown in the table, we
discuss the details for three interesting case studies.

495 4.2.1. ED2K and KADEMLIA

496 ED2K and KADEMLIA eMule messages are preceded by a com-497 mon header which is used as GT. FieldHunter correctly identifies 498 such common header. Moreover, it discovers additional fields, that 499 sum up to a total of 128 bits in the C2S EDK2 collection. After 500 manual inspection, we observed those fields to correctly include 501 Session-ID, and Host-ID.

502 4.2.2. SopCast

503 SopCast is a proprietary and closed protocol used for P2P-TV 504 broadcasting. Unveiling information about the message format of 505 such protocols is one of the motivations for developing Field-506 Hunter.

507 Specifically, this protocol represents a large fraction in the TR3-508 2007 trace. FieldHunter identifies 128 bits corresponding to: MSG-509 Len, Trans-ID, Session-ID, Host-ID (we hypothesize it is used for 510 NAT traversal since it uses 64 bits, 32 of which correspond typi-511 cally to private IP addresses, and 32 are identical to the Host public 512 IP address) and some accumulators of 16, 32 and 64 bits (possibly 513 used to reorder video/audio chunks).

514 4.2.3. Domain name service (DNS)

515 For DNS in the TR1-2012 trace, FieldHunter successfully identifies the Trans-ID and a MSG-Type field, each of 16 bits. We ex-516 pected parsing DNS in this trace to be challenging due to the bias 517 in the collection. First, most of the C2S messages are "DNS Re-518 519 quests" messages. Second, requests are directed to the most popu-520 lar DNS resolvers (in the TR1-2012 trace customers use the ISP DNS 521 server). Despite this, FieldHunter is able to identify some protocol 522 fields.

Interestingly, in the C2S messages, FieldHunter reports the presence of a 16 bit accumulator in the DNS Trans-ID field. We manually verified this, and discovered that some implementations of
DNS clients generate a "random" Trans-ID by using a local counter.
FieldHunter captured this peculiar but common behavior, exposing
more details about the protocol.

529 4.3. Evaluation of textual protocol

Table 4 reports overall results for textual protocols. It reports 530 the number of inferred fields, the number of key-value pairs (K-V) 531 and singletons (most of them MSG-Type) for each direction (with 532 533 the exception of the last two protocols for which the DPI provided just one direction of the conversation). In addition, we report those 534 535 fields that we label as being identifiers (IDs), highlighting those that proved to be False Positives (FP). Here, by IDs we mean Host-536 537 IDs, Session-IDs, and Trans-IDs. Overall, from the 26 fields labeled 538 as IDs, 22 are verified as correct and only three are false positives. In general, we observe that the majority of the fields of textual 539 protocols are successfully inferred in both the C2S and S2C direc-540 541 tions. Similar to the binary protocols, we pick two interesting tex-542 tual protocols as case studies.

543 4.3.1. Microsoft notification protocol (MSNP)

The MSN protocol is present in all three traces. FieldHunter correctly finds that the field called USR field carries a Host-ID and which indeed is the MSN's user name. Similarly, the CVR field which is used to send specific information about the client and its OS to the server. This field is captured by FieldHunter since system settings are different for each MSN user, but consistent during the communication with the server. Although CVR is not an actual



Fig. 7. Parameter sensitivity for the MSG-Type.

Host-ID, this is a right interpretation for the field type because the 551 field behaves the same as Session-ID. 552

4.3.2. *Real-time streaming protocol (RTSP)*

The S2C direction of this protocol returns 6 inferred ID fields. 554 Out of these, four are correctly labeled and two are false posi-555 tives. The latter occur when some fields that are supposed to take 556 different values actually always take the same value for a given 557 conversation, behaving similar to a Session-ID. The false positive 558 fields are Last-Modified and Cache-Control. For instance, 559 Last-Modified is the timestamp of the last modification for 560 a given content. Since a single object is requested using multiple 561 RTSP conversations, its modification time appears constant across 562 conversations. Similarly, the Cache-Control field tends to al-563 ways take the same value among conversations used to retrieve 564 the same content as well. In general, we observe that the original 565 collection may be biased toward some specific subset of protocol 566 fields and values. This is challenging for FieldHunter, and, in gen-567 eral, any field inference algorithm that relies on traffic data. 568

4.4. Sensitivity analysis and parameter tuning

We evaluate the sensitivity of FieldHunter to different parameters and to external factors, such as the collection size. As mentioned before, we perform parameter tuning using one trace, and then we evaluate FieldHunter on all three traces. Next we show how the design proved to be robust to parameter tuning. 574

First we focus on one of the most challenging fields to infer, 575 the MSG-Type for binary protocols. We consider all collections for 576 those protocols that have available ground truth. Then for each col-577 lection, the MSG-Type algorithm is executed manifold by tweaking 578 the thresholds (Min. Correlation and Max. Entropy). For each thresh-579 old pair, the product between Coverage and AoC is computed, pro-580 viding a coefficient from 0 to 1, where values close to 1 are desired. 581 The results are reported in Fig. 7. The darker the block, the better 582 FieldHunter performs. As can be clearly seen, there is a large range 583 of good parameters that yield scores above 0.8, which means that 584 in most cases FieldHunter was able to correctly pinpoint the MSG-585 Type field. We repeat the experiment for other field types and we 586 observed qualitatively similar results. 587

We now evaluate the effect of the collection size for both bi-588 nary and textual protocols. For textual protocols, we first select 589 nine protocols for which we know all fields present in our traces. 590 Then, we randomly extract a subset of conversations from the col-591 lections and run FieldHunter over the subset. Results are compared 592 against ground truth to compute the Coverage and AoC (Fig. 8). We 593 see that FieldHunter performs well even with limited number of 594 textual conversations. In fact, when 50 conversations are consid-595 ered, we identify 85% of all the fields, with 97% AoC. Overall, using 596 large enough collections, we always identify the D_f delimiter for 597 all tested protocols. Most of the mis-labeling happens due to chal-598 lenges in inferring the $D_{k-\nu}$ for some fields. 599

[m5G;March 4, 2016;11:33]

7

569

I. Bermudez et al./Computer Communications xxx (2016) xxx-xxx





600 For binary protocols, we perform similar evaluation, but focus 601 on two challenging protocol cases, DNS and Vuze DHT. The results 602 are shown in the bottom plot of Fig. 8. As we can see FieldHunter 603 may require a bigger collection sizes to produce the best results. We believe that the heuristics apply differently on textual and bi-604 nary protocols as textual protocols are less sensitive to diversity. 605 Vuze DHT represents the protocol for which we had highest di-606 607 versity in our dataset, with many end-points exchanging a variety of messages. Conversely, DNS (from TR1-2012) represented a chal-608 609 lenging scenario due to little diversity in the collections: typically 610 only one MSG-Type (DNS requests) was found, and conversations 611 were very short (a single request/response). As we see, eventually 612 we achieved very good results for DNS, but it required as many as 2000 conversations. 613

5. Application 614

615 In this section, we present an end-to-end security application that uses FieldHunter. Specifically, we focus on intrusion detection 616 systems (IDS) which are a common defense mechanism for many 617 critical infrastructures. IDS can be categorized into anomaly detec-618 tion systems or signature based systems based on the detection 619 620 mechanisms used. Anomaly detection systems work by learning 621 the normal, also called baseline, behavior of any system and flag-622 ging any deviation from this as *anomaly*, which can be considered an indicator of malicious behavior. This is in contrast with signa-623 624 ture based systems which use signatures of known threats. The ad-625 vantage of anomaly detection systems over signature based ones is that they can detect previously unknown attacks, i.e. zero-day at-626 tacks. However, anomaly detection systems can suffer from higher 627 false positives. 628

In this work, we develop a payload-based anomaly detection 629 630 system for Industrial Control System (ICS) networks. Most of the 631 current solutions rely on statistical features such as volume of traf-632 fic for creating baseline [21]. Such systems are ineffective against stealthy attacks that work by modifying the protocol behavior 633 without causing discernible change in statistical properties of the 634 635 traffic. There are a few systems that offer payload based anomaly detection, but they rely on the protocol specifications which are 636 not available for many of the legacy protocols that are used quite 637 638 often [22].



Fig. 9. Temporal patterns shown in ICS communications.

FieldHunter forms the core of the anomaly detection system. It 639 infers protocol fields from the ICS network being monitored. These 640 protocol field summaries are converted to policy rules which form 641 the baseline of the system. Any deviation in network protocol pay-642 load from these rules are flagged as anomalies. This allows us to 643 detect attacks which modify the protocol messages or introduce 644 new types of messages which are not part of the protocol. There 645 is a direct mapping between the human readable field summaries 646 and the policy rules generated. This is greatly beneficial to any an-647 alyst analyzing the root cause of alerts generated by the anomaly 648 detection system. 649

5.1. ICS networks 650

ICS networks typically consists of a set of devices such as sen-651 sors or controllers and a central monitoring/administrative unit. 652 These networks differ from other networks such as enterprise or 653 ISP networks in many respects. ICS network traffic contains a large 654 portion of proprietary protocols which are unknown or not well 655 documented along with a few well known protocols such DNS or 656 SNMP. These proprietary protocols are predominantly binary based. 657 The communication patterns in these networks are quite determin-658 istic as only a handful of devices communicate with each other ac-659 cording to configurations that are set statically. Moreover, many of 660 the connections are long-lived, often lasting for hours or days. 661

These characteristics present some unique challenges for Field-662 Hunter. First challenge is the lack of diversity in the traffic, which 663 renders some of the heuristics such as detection of Host-IDs inef-664 fective. However, this is not a big limitation as our goal is to de-665 tect anomalous behavior and not protocol understanding. Hence, 666 even if we incorrectly label a Host-ID field as a Constant, this field 667 can be used to detect violations when the attacker intentionally 668 changes the "Constant" value. 669

Second challenge is the presence of long-lived network flows. 670 Many of the thresholds for various heuristics are set assuming 671 that the input to FieldHunter contains multiple flows. We analyzed 672 long-lived ICS connections and observed that often these connec-673 tions show repeating patterns which indicate some conversations 674 happen with a specific frequency. Hence, multiple conversations 675 occur serially within a single flow. Therefore, we can break up a 676 single flow into multiple conversation by computing the frequency 677 of communication. Fig. 9 depicts packet payload sizes in a window 678 slice of 10 seconds for a very long conversation. We observe that 679 the devices involved in the communication send burst of protocol 680 messages at a frequency of approximately 1 second. 681

5.2. Design

The payload-based anomaly detection system (Fig. 10) has two 683 modes of operation: Training and Detection. In Training mode, the 684 system generates network protocol profiles associated with specific 685 TCP/UDP ports used by the physical devices for communication. 686 During this phase, FieldHunter plays a critical role, since it learns 687

I. Bermudez et al. / Computer Communications xxx (2016) xxx-xxx

9



Fig. 10. Payload Based Anomaly Detection System diagram.

688 the protocol field summaries used as baselines by the next step. 689 In Detection mode, the anomaly detector uses rules to determine 690 flow by flow if any baseline policy rule has been violated. For in-691 stance, the detection mode can verify whether a flow contains the 692 particular value that a specific field is supposed to take as per the 693 policy rule. Otherwise, the system raises an alarm of violation.

5.2.1. Training mode 694

Network traffic is intercepted using a tap, that listens to net-695 696 work communications of devices connected to the ICS network. TCP/UDP payloads are extracted and grouped properly, forming the 697 protocol collections required by FieldHunter to create functional 698 699 protocol descriptions.

700 Traffic Labeler, accomplishes the fundamental job of creating 701 the protocol collections. Since this IDS is profiling network pro-702 tocols running in a particular mode on industrial devices, traffic is labeled by their associated triplets (Layer-4 protocol, destination 703 704 port, destination IP address). For instance, packets going through a bi-directional TCP network connection, in which IP 192.168.1.101 is 705 706 sending packets to 192.168.1.202, using source port 9988 and destination port 1822, has two triplets associated with it: (TCP, 9988, 707 708 192.168.1.101) and (TCP, 1822, 192.168.1.202).

The Conversation Reconstruction module reassembles conversa-709 710 tions. The Conversation Collector groups them by label (triplet), 711 until we get enough number of packets associated to a triplet i.e., greater than a pre-defined threshold. The dataset associated to a 712 label is the protocol collection for FieldHunter (See Section 2) and 713 the threshold tells how much data is enough to start field infer-714 715 ence from the collection.

The Sub-conversation Splitter module is used in to handle very 716 long-lived connections. This module splits the long conversations 717 into multiple small conversations when strong temporal patterns 718 in the traffic are observed (see Fig. 9). To do this, the module ar-719 720 tificially increments the number of connections associated with a 721 collection such that there is one for each direction of the conver-722 sation.

723 FieldHunter uses these conversations to generate protocol sum-724 maries. The protocol summaries are converted to policy rules in a 725 straight-forward manner. The type of the field determines the nature of the rule. We have a pre-determined way of converting each 726 field type to a rule. The key idea behind this conversion is that 727 each rule specifies how a field should behave, for instance, what 728 values it can take or how the value correlates to something else 729 730 such as message length. As an example, the rule for MSG-Type field 731 specifies which opcodes to expect in the flows. These rules are put 732 into the Rule Database.

The duration of the training mode depends on in the amount of 733 available data rather than the clock time. This means if an applica-734 tion handles more traffic, the required clock time to generate sum-735 736 maries is lesser. Fig. 8, gives a rough idea of the required amount of application messages to be collected before moving into the de-737 738 tection mode.

Table 5

Network captures used for evaluation of the payload-based anomaly detection system.

Trace	Duration (h)	#Flows	#Messages	#Malicious flows	#Endpoints
TPC	2	176	46,503	0	5
MPC	2	180	31,274	2	5

5.2.2. Detection mode

In a nutshell, detection mode can be described by the interaction of three different modules: the Rule Database, the Traffic La-741 beler and the Anomaly Detector. Our system has two inputs in detection mode: on the one had it receives network traffic and on the other it gets the rules from the database. The output of the anomaly detection system is packets/connections which are flagged as anomalous as they violate the rules learned by the system. 746

The traffic labeler simply labels traffic with the respective pair 747 of triplets and forwards it to the anomaly detector. Using the triplets, the detector can retrieve the rules from the database and apply them to the traffic. If rules are violated, system raises alarms 750 indicating that the traffic contains anomalous packets or connections. When rules cannot be found in the rules database, traffic is forwarded to the training mode path. 753

5.3. Evaluation of anomaly detection system

We tested our system using traffic collected from PowerCyber 755 the test-bed. It simulates and emulates components of a smart grid 756 including industrial SCADA [23]. The anomaly detection system is 757 trained with clean network traffic. Finally, we evaluate its efficacy 758 by running traffic collected during an attack. 759

5.3.1. Datasets

Network traffic is collected from one control center network of 761 PowerCyber's test-bed. Traffic is collected during normal operation 762 of the control center network and also during an attack. Such at-763 tack is performed intentionally in a controlled environment.

Table 5 describe the datasets used for testing the detection 765 system. Two datasets of two hour duration, each with about 180 766 flows, are used. Despite the few flows present in these datasets, the number of application messages is about thousands, which al-768 lows FieldHunter to provide meaningful results from it. The attack 769 dataset contains both normal traffic and two flows belonging to 770 an authentication attack. We use full packet traces and there is no 771 packet filtering. Training Packet Capture (TPC) contains only nor-772 mal traffic, while Mixed Packet Capture (MPC) contains normal and 773 attack traffic mixed together.

5.3.2. Evaluation

We evaluated the system end to end, using TPC dataset to train 776 the system and MPC for testing. The expectation was that the anomaly detection system should flag only malicious flows in the 778 mixed dataset and ignore the benign traffic.

We found 6 different protocol summaries from MPC. One of 780 them belongs to the DNP3 protocol [24], which is common in 781 SCADA systems. The message structure of DNP3 protocol is defined 782 by a header containing the constant value 0x0564 at the start, the 783 message length, a control byte, destination and source ids, and a Cyclic Redundancy Check (CRC) field followed by the payload. Similarly payload is structured as data blocks which contain user data (1-16 byte) followed by a CRC.

From the protocol, FieldHunter is able to identify parts of the constant field, message length, and destination and source ids. In the given trace, the control byte appears as constant (0xC4 for 790 client-to-server direction and 0x44 for server-to-client direction). 791

739 740

748 749

751 752

754

760

764

767

774

775 777

779

784 785 786

ARTICLE IN PRESS

850

872

880

I. Bermudez et al. / Computer Communications xxx (2016) xxx-xxx

Similarly, the destination field appear as a constant in the clientto-server messages as there is only one server in the trace, and the same for the source in the other direction. Therefore, the identity of clients can be associated to IP address.

MPC, contains a malicious DNP3 flow which abuses the lack of 796 authentication in the protocol. Another machine in the same net-797 work spoofs valid DNP3 messages, but it has to provide a 2 byte 798 entity in the source id field of the protocol. The attacker has two 799 800 options: (i) provide a new source id, or (ii) use one already used by another client. In both cases we can raise an alarm, because 801 802 in case (i) we observe a new source id not seen during training; 803 for case (ii) we observe that the source id does not match the IP address associated to that identifier. In this particular attack the 804 805 attacker spoofed the identifier of one of the observed clients. The system detected the anomaly and can provide a detailed report to 806 the administrator to make an informed decision. 807

We also evaluated our system on real industrial datasets. Results were similar to the ones presented for test-bed dataset. Unfortunately, due to privacy concerns about the data, we can not disclose further details in this paper.

812 6. Assumptions and limitations

In this section, we discuss the assumptions of FieldHunter. Further, for each assumption we describe the limitations imposed on the system due to that assumption.

6.1. Training data contains patterns corresponding to the protocolfields

FieldHunter assumes that the training data contains the pat-818 819 terns that allow it to identify protocol fields. Hence, the sys-820 tem cannot identify fields when traffic is compressed or encrypted. However, this limitation can be overcome by using a pre-821 822 processing stage that decompresses or decrypts traffic before passing it to FieldHunter. Another consequence of this assumption is 823 824 that the FieldHunter is dependent on the quality of the dataset. Therefore, if the traces do not contain the required diversity that 825 highlights data patterns, FieldHunter is not able to identify fields 826 827 from traffic, or worst it may misclassify fields. For example, a 828 Server-ID may be classified as a constant if traces only contain traf-829 fic directed to one particular server.

830 6.2. TCP PSH flag delimits the start of a new application message

For TCP communication FieldHunter does not assume that a 831 832 single TCP segment contained in packet carries an application message. It heavily relies on TCP PSH flag to delimit the end from the 833 beginning of a new application message in the TCP stream. The as-834 835 sumption here is that the TCP PSH flags are set when sender does 836 not have more data to transfer. On the receiver side ,TCP PSH flag is used to trigger transmission of the data to the application layer. 837 When a TCP PSH flag is incorrectly set in the middle of an appli-838 cation message, it creates a message misalignment. This is simi-839 lar to the noise due to having random data in the protocol collec-840 tion. Note that we did not observe this odd behavior in any of our 841 842 traces

In addition, we assume that the data in-between two PSH flags belongs to a single application message. However, this is not true all the times, since multiple application messages can be contained in between two consecutive PSH flags due to TCP buffering. This situation can be mitigated if the protocol shows different types of messages with multiple lengths allowing FieldHunter to split the "single" application message into multiple correct ones.

6.3. Binary protocols have a header and a payload

We assume that every binary protocol has similarly structured 851 messages – each with a common header and a payload that may 852 change. However, this assumption may not be true in all cases. 853 For instance, some protocols can have messages with different 854 formats for handshakes or preamble communication and subse-855 quently, when connection is established all communication uses 856 standard uniform messages. FieldHunter is not designed to dis-857 criminate among these two types of message formats. This can be 858 overcome by using a preprocessing module that splits the conver-859 sations to ensure that the messages having different formats are in 860 separate collections. Moreover, FieldHunter can still provide mean-861 ingful results if this situation happens. However, the quality of the 862 final result depends on how popular is each type of message is in 863 the whole collection. For instance, if flows are long lived, then the 864 preamble message will be less popular than the rest of the mes-865 sages, and the final result will contain fields of the protocol format 866 that follows the preamble. 867

Another characteristic of FieldHunter is that it cannot differentiate between message header and payload. So even though, we target field inference from the header, sometimes FieldHunter ends up identifying fields contained in the payload part of the protocol. 871

6.4. Dataset contains pure protocols

Our expectation is that the protocol collections are pure since 873 FieldHunter is sensitive to noise in the form of different protocol 874 format types. However, traffic classifiers such as a DPI can produce 875 false positives, which means flows belonging to other protocols are 876 classified as the protocol of our interest. FieldHunter is not able to 877 detect the presence of this noise, and will underperform depending 878 on how bad is the noise level. 879

6.5. Fields in textual protocols arekey-value paired

FieldHunter looks for key-value pairs found in textual protocols. 881 It does not take advantage of other characters that provide richer 882 structure to textual protocols such as parenthesis to enclose one 883 object, or commas used in enumeration. When complex textual 884 protocols are processed by FieldHunter it can still obtain valuable 885 information. However, it can produce a less richer result set that 886 fails to take into account the structure of the data. 887

6.6. ICS protocols are similar to "traditional" network protocols 888

Our observation is that many of the ICS network protocols use 889 traditional transport UDP/TCP network protocols. Hence, our as-890 sumption is that these protocols can be profiled using FieldHunter. 891 However, there are other differences in the environments which 892 impose further challenges that require special attention. For in-893 stance, the diversity in such networks is more difficult to reach 894 given the applications or setups of those networks. This makes it 895 more difficult to correctly identify some fields since the patterns 896 do not show up in the traffic collection, or in some cases, the final 897 result can overfitted to the training data given the lack of samples 898 that truly represent the protocol behavior. However, we note that 899 for the anomaly detection application that we described, we do not 900 need to infer all the fields of the protocol. Hence, FieldHunter is 901 very effective for this application. 902

7. Related work

903

7.1. Protocol format reverse engineering using network traces 904

Automatic inference of protocol formats from passive net- 905 work monitoring was first addressed by Beddoe [25]. The authors 906

Please cite this article as: I. Bermudez et al., Towards automatic protocol field inference, Computer Communications (2016), http://dx.doi.org/10.1016/j.comcom.2016.02.015

ARTICLE IN PRES

I. Bermudez et al. / Computer Communications xxx (2016) xxx-xxx

applied the Needleman-Wunsch algorithm for alignment of byte 907 908 sequences between network payloads. The same algorithm has been used in Scriptgen [26] and RolePlayer [9] for automating the 909 910 process of learning protocols in honey-nets. Their works aim to find variant and invariant segments in textual protocols. In con-911 trast, our aim is to identify a broader selection of field types. 912

The problem of extracting message format specification for se-913 curity applications was later addressed by Discoverer [8]. They first 914 915 clustered messages with similar formats together using sequence alignments and then identified parts of the messages that change 916 917 across flows. In contrast to FieldHunter, Discoverer has the same 918 limitations as in [9,25,26], where fields of the protocols are ex-919 pected to appear in a predefined order. In [12] authors propose 920 Prodecoder that uses semantic information for field extraction, by using the LDA model. Their approach looks promising for identify-921 ing keys and the syntax of textual protocols, but it is not clear how 922 LDA can properly merge n-grams into fields of binary protocols. 923

7.2. Protocol format reverse engineering using binary analysis 924

Other authors have tackled the problem of protocol reverse en-925 gineering by using binary analysis. For instance Prospex [6] is a 926 927 system that analyzes both binary execution traces combined with 928 network traffic. Binary analysis requires an instrumented system with enough privileges to read protected memory of the applica-929 930 tion that uses the protocol. Similar in spirit, in Dispatcher [27] the authors focused on protocol reverse-engineering for botnet infil-931 932 tration. All the above works rely on binary analysis and they are therefore very different from what we want to achieve with Field-933 Hunter, where we only have passive access to network traffic. 934

935 7.3. Protocol network signature generation from network traces

In [7,10,11,13,28] authors automatically derive protocol signa-936 937 tures purely from network traces. In PEXT [10] and ReveX [7] sig-938 natures are extracted for protocols using similar tokens to cluster flows. On the other hand [28] uses semantic information found in 939 940 the protocol to group messages with similar formats. Authors in 941 [13] propose a system that can automatically produce signature for botnets' command and control traffic. Obtaining automatically gen-942 erated signatures for traffic classification has multiple positive im-943 plications. However, understanding the mechanics of the semantic 944 945 of the protocols is a valuable complementary information for the system experts to verify the quality of automatically generated sig-946 947 natures.

7.4. Payload based anomaly detection systems 948

949 Systems able to detect anomalies from network traffic observ-950 ing protocol payload, have been previously proposed in [29-32]. 951 These anomaly detectors define signatures on patterns found in 952 network payloads, from which they create models and baselines. 953 Our proposed system differs from all of them, because we use a generic abstraction of protocol format to create rules and baselines. 954 This allows easier false positive back tracking, since the system can 955 explain better the cause and context of a violation. For example, 956 957 when a new value for an opcode has been observed. Moreover, as 958 far as we are aware, none of these systems have been tested in industrial control networks. 959

7.5. ICS anomaly detection systems 960

961 In the recent years there has been an interest on developing anomaly based IDS for ICS. Authors in [21] propose a system that 962 uses clustering to detect anomalies. Their system does not neces-963 sary look only at network traffic, but also at other features that 964

http://dx.doi.org/10.1016/j.comcom.2016.02.015

may come from other logs such as engine speeds and tempera-965 ture. Another work that is more closely related to ours is that of 966 Caselli et al. [22]. Here, the authors use known protocol descrip-967 tions to extract field values from network traffic, which they use 968 to create models to detect sequence attacks. The system that we 969 have developed shares some similarities with their solution. How-970 ever, a crucial difference is that they assume that they know the 971 protocol specification. This is an unrealistic assumption as many of 972 the protocols used by industrial devices are poorly documented. 973

FieldHunter is complementary to other systems for extraction 974 of protocol message format, as our final goal is to identify con-975 tainers/fields of information. Moreover, we present a specific ap-976 plication for FieldHunter in a critical security context, using data 977 obtained from a realistic scenario. 978

8. Conclusions

In this paper, we presented FieldHunter, a system that auto-980 matically infers protocol field types from passive observation of 981 network traffic. We showed that FieldHunter is able to provide a 982 comprehensive set of fields and their types for both textual and 983 binary protocols that may not have a publicly available specifica-984 tion. Therefore, we believe that a system such as FieldHunter can 985 significantly improve the effectiveness of modern network security 986 tools. 987

Finally, we extended FieldHunter and built a payload-based 988 anomaly detection system on top of it. FieldHunter provides valu-989 able information about network protocol specification, allowing it 990 to detect realistic zero-day attacks on ICS network. Our anomaly 991 detection system can detect stealthy attacks in ICS systems with 992 un-documented protocols that current statistical-based or tradi-993 tional payload-based anomaly detection systems cannot. 994

References

Please cite this article as: I. Bermudez et al., Towards automatic protocol field inference, Computer Communications (2016),

- [1] Z. Li, G. Xia, H. Gao, Y. Tang, Y. Chen, B. Liu, J. Jiang, Y. Lv, Netshield: massive 996 semantics-based vulnerability signature matching for high-speed networks, ACM SIGCOMM Comput. Commun. Rev. 41 (4) (2011) 279-290.
- [2] J. Caballero, H. Yin, Z. Liang, D. Song, Polyglot: automatic extraction of protocol message format using dynamic binary analysis, in: Proceedings of the 14th ACM Conference on Computer and Communications Security, ACM, 2007, pp. 317-329.
- [3] P.M. Comparetti, G. Wondracek, C. Kruegel, E. Kirda, Prospex: protocol specification extraction, Proceedings of the 2009 30th IEEE Symposium on Security and Privacy, IEEE, 2009, pp. 110-125.
- [4] W. Cui, M. Peinado, K. Chen, H.J. Wang, L. Irun-Briz, Tupni: automatic reverse engineering of input formats, in: Proceedings of the 15th ACM Conference on Computer and Communications Security, ACM, 2008, pp. 391-402.
- [5] Z. Lin, X. Jiang, D. Xu, X. Zhang, Automatic protocol format reverse engineering through context-aware monitored execution, in: Proceedings of Network and Distributed System Security, NDSS, vol. 8, 2008, pp. 1-15.
- [6] G. Wondracek, P.M. Comparetti, C. Kruegel, E. Kirda, S.S.S. Anna, Automatic net-1012 1013 work protocol analysis, in: Proceedings of Network and Distributed System Security, NDSS, vol. 8, 2008, pp. 1-14. 1014
- [7] J. Antunes, N. Neves, P. Verissimo, Reverse engineering of protocols from net-1015 work traces, in: Proceedings of IEEE 2011 18th Working Conference on Reverse 1016 Engineering, WCRE '11, Limerick, IR, 2011, pp. 169-178. 1017 1018
- [8] W. Cui, J. Kannan, H.J. Wang, Discoverer: automatic protocol reverse engineering from network trace, in: Proceedings of the 16th USENIX Security Symposium, USENIX Association, Boston, MA, 2007, pp. 1-14.
- [9] W. Cui, V. Paxson, N. Weaver, R.H. Katz, Protocol-independent adaptive replay of application dialog, in: Proceedings of Network and Distributed System Security, NDSS, 2006.
- [10] M. Shevertalov, S. Mancoridis, A reverse engineering tool for extracting proto-1024 cols of networked applications, in: Proceedings of the 14th Working Confer-1025 ence on Reverse Engineering, 2007, WCRE 2007, IEEE, Vancouver, BC, CA, 2007, 1026 pp. 229-238. 1027
- [11] A. Tongaonkar, R. Keralapura, A. Nucci, Santaclass: A self adaptive network 1028 traffic classification system, Proceedings of IFIP Networking Conference, 2013, 1029 IEEE, 2013, pp. 1-9. 1030
- [12] Y. Wang, M.Z. Shafiq, L. Wang, A.X. Liu, Z. Zhang, D. Yao, Y. Zhang, L. Guo, A 1031 semantics aware approach to automated reverse engineering unknown proto-1032 cols, in: Proceedings of the 2012 20th IEEE International Conference on Net-1033 work Protocols, ICNP, 2012, pp. 1-10, doi:10.1109/ICNP.2012.6459963. 1034

995

979

997 998

999 1000

1001 1002

1003 1004

1005 1006

1007 1008

1009 1010 1011

1019

1020

1021

I. Bermudez et al./Computer Communications xxx (2016) xxx-xxx

- 1035 [13] C. Rossow, C.J. Dietrich, Provex: Detecting botnets with encrypted command 1036 and control channels, Proceedings of the 10th International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment, DIMVA'13, 1037 Springer-Verlag, Berlin, Heidelberg, 2013, pp. 21–40. N. Falliere, L.O. Murchu, E. Chien, W32. Stuxnet Dossier, White Paper, Symantec 1038
- 1039 [14] Corp., Security Response 5. 1040
- 2015 Dell Security Annual Threat Report, Technical report, Dell Inc., 2015. 1041 [15] 1042 https://software.dell.com/whitepaper/dell-network-security-threat-report-1043 2014874708
- **Q4**044 C. Kreibich, J. Crowcroft, Honeycomb: creating intrusion detection signatures [16] using honeypots, ACM SIGCOMM Comput. Commun. 34 (1) (2004). 1045
- **Q5** 1046 Tcp Statistic and Analysis Tool, URL http://tstat.polito.it.
- 1047 Y. Yao, Information-theoretic measures for knowledge discovery and data min-[18] 1048 ing, Entropy Measures, Maximum Entropy Principle and Emerging Applica-1049 tions, Springer, 2003, pp. 115-136.
- 1050 [19] Opendpi. https://code.google.com/archive/p/opendpi/.
- [20] Distributed Hash Table, 2012. http://wiki.vuze.com/w/Distributed_hash_table. 1051
- 1052 I. Kiss, B. Genge, P. Haller, G. Sebestyen, Data clustering-based anomaly detec-[21] 1053 tion in industrial control systems, in: Proceedings of 2014 IEEE International 1054 Conference on Intelligent Computer Communication and Processing (ICCP), 1055 2014, pp. 275-281.
- [22] M. Caselli, E. Zambon, F. Kargl, Sequence-aware intrusion detection in indus-1056 1057 trial control systems, Proceedings of the 1st ACM Workshop on Cyber-Physical System Security, CPSS '15, ACM, New York, NY, USA, 2015, pp. 13-24. 1058
- 1059 [23] A. Hahn, A. Ashok, S. Sridhar, M. Govindarasu, Cyber-physical security testbeds: 1060 architecture, application, and evaluation for smart grid, IEEE Trans. Smart Grid 1061 (2) (2013) 847-855
- 1062 [24] Overview of the dnp3 Protocol. http://www.dnp.org/pages/aboutdefault.aspx.

- [25] M.A. Beddoe, Network Protocol Analysis Using Bioinformatics Algorithms, 1063 Technical report, Baseline Research, 2005.
- 1064 [26] C. Leita, K. Mermoud, M. Dacier, Scriptgen: an automated script generation tool 1065 1066 for honeyd, in: Proceedings of the 21st Annual Computer Security Applications Conference, CSAC '05, IEEE, Tucson, AZ, 2005, pp. 214-226, doi:10.1109/CSAC. 1067 1068 2005.49.
- [27] J. Caballero, P. Poosankam, C. Kreibich, D. Song, Dispatcher: enabling active 1069 botnet infiltration using automatic protocol reverse-engineering, in: Proceedings of the 16th ACM Conference on Computer and Communications Security, CCS '09, ACM, Chicago, IL, 2009, pp. 621-634.
- [28] V. Yegneswaran, J.T. Giffin, P. Barford, S. Jha, An architecture for generating semantics-aware signatures, in: Proceedings of the USENIX Security, 2005, pp. 34-43.
- [29] K. Wang, S. Stolfo, Anomalous payload-based network intrusion detection, in: 1077 E. Jonsson, A. Valdes, M. Almgren (Eds.), Recent Advances in Intrusion Detection, Lecture Notes in Computer Science, vol. 3224, Springer, Berlin, Heidelberg, 1078 2004, pp. 203-222. 1079
- [30] D. Bolzoni, S. Etalle, P. Hartel, E. Zambon, Poseidon: a 2-tier anomaly-based 1080 network intrusion detection system, in: Proceedings of the 2010 International 1081 Workshop on Innovative Architecture for Future Generation High Performance. 1082
- [31] K. Wang, J. Parekh, S. Stolfo, Anagram: a content anomaly detector resistant to 1083 mimicry attack, in: D. Zamboni, C. Kruegel (Eds.), Recent Advances in Intru-1084 sion Detection, Lecture Notes in Computer Science, vol. 4219, Springer, Berlin, 1085 Heidelberg, 2006, pp. 226-248. 1086
- [32] R. Perdisci, D. Ariu, P. Fogla, G. Giacinto, W. Lee, Mcpad: a multiple classifier 1087 system for accurate payload-based anomaly detection, Comput. Netw. 53 (6) 1088 (2009) 864-881. Traffic Classification and Its Applications to Modern Networks 1089

JID: COMCOM