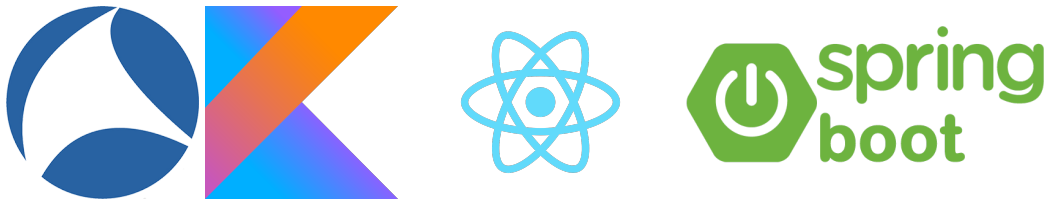


Visualization System of Network Forensic Data

Ivan Manoilov



Abstract

The goal of this project is to develop a system for processing various forensic data inputs (such as PCAP) and presenting them in a user-friendly, graphical way. Modern systems for visualization of data are too generic for this specific use-case and can't be employed for analysis of forensic data. This problem is solved via development of the system of analytical dashboards which represent the data, processed by various back-end services, responsible for parsing input, aggregating data and transforming it to the format acceptable by dashboards. Result of this thesis would be a highly extensible system, providing deep analytical view of network traffic, which enables the end-user to see underlying anomalies or problems of the traffic.

Keywords: Data processing — Analytical dashboards — Network Forensic Data

Supplementary Material: [Downloadable Code](#)

*xmanoi00@stud.fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

1. Introduction

The main goal of this project is to develop a system able to visualize network forensic data, so the user will be able to see problems and forensic activity on network. With rapid development of Internet technologies and expansion of Internet throughout the globe, cybercrime has become one of the most common of the crimes. The situation has only worsened nowadays, as one the main impacts of COVID-19 was transition of parts of our lives to the Internet. Reports of cyber-dependent crime and online fraud have increased during the COVID-19 outbreak, and rates of cybercrimes have been particularly high during months with the strictest lockdown policies [1].

Output of the developed system is composed of 4 dashboards corresponding to 4 types of network data input:

- **Network Conversations** : the basic view of the network communication data.

- **Extracted Files** : the view of the unencrypted data exchange in format of files.
- **Encrypted Traffic** : the view of the encrypted communication.
- **Resolved Domains** : the view of the DNS queries and responses.

The data for dashboards is provided from extensible back-end, which is responsible for parsing various formats of network data inputs, compressing, aggregating and presenting in acceptable for front-end format.

Tools used to achieve the same functionality as this project (**Grafana/Kibana**) provide generic way of working with data, which leads to unmet expectations of possibilities for analysis of network forensic data. Moreover, aforementioned tools enable only a small set of visualizations useful for described use-cases, therefore most of the information is presented via text/tables, which by definition are less graphical and easy for perception. The final drawback of these solutions is no support for data processing, they only

40 consume network data in predefined formats, therefore
41 are unable to optimise the pipeline in one of the most
42 performance heavy parts of it.

43 On contrary, this project provides the users with
44 tailored visualizations suitable for analysis of network
45 data. Front-end works in tandem with back-end ser-
46 vices through REST API, which enables optimization
47 of raw data by seamlessly compressing it (60-80 %
48 compression achieved with normal data-sets), further
49 formats of input data are supported and with exten-
50 sibility of back-end possibilities are to accept more
51 formats.

52 2. Forensic Data Analysis

53 2.1 Forensic network data

54 Digital forensics are classically defined as the use of
55 scientifically derived and proven methods toward the
56 preservation, collection, validation, identification, anal-
57 ysis, interpretation, documentation and presentation of
58 digital evidence derived from digital sources for the
59 purpose of facilitating or furthering the reconstruction
60 of events found to be criminal, or helping to anticipate
61 unauthorized actions shown to be disruptive to planned
62 operations [2].

63 Preservation and collection of the evidence is done
64 with the help of the **packet sniffing** tools like Wire-
65 shark. These tools usually run on server with various
66 applications and capture incoming and outgoing traffic
67 in format of packets. **Packet** is a group of bits that in-
68 cludes data and control information. It generally refers
69 to a network layer (OSI layer 3) protocol data unit [3].
70 Figure 4 showcases the GUI of Wireshark as well as
71 insides of a packet. These data are stored in format of
72 .pcap file, which represents a collection of packets
73 and metadata about them.

74 2.2 Analytical Dashboards

75 A **dashboard** is a visual display of the most important
76 information needed to achieve one or more objectives;
77 consolidated and arranged on a single screen so the
78 information can be monitored at a glance [4]. An
79 example of dashboard is displayed on Figure 5.

80 According to the goal of this project our main pri-
81 ority is to design an **analytical** dashboard. Analytical
82 (tactical) dashboards provide aggregated data and help
83 to analyze the data. They tend to present snapshots of
84 data to provide a better look at the context of the data
85 [4].

86 One great possibility that this type of dashboards
87 provides us, is the aggregation of the data, particu-
88 larly aggregation of data based on time. **Time series**
89 data can be very performance-heavy for processing

with little to no impact on final analysis. Therefore, 90
time-series data can be compressed in so-called **time** 91
buckets. This is done via partitioning the time win- 92
dows into sub-windows defined by various algorithms 93
[5], which help to select the most optimal indexes for 94
partitioning and sizes of buckets. 95

3. System Architecture 96

The System is divided in two parts : data processing 97
(back-end) and data representation (front-end, dash- 98
boards). High-level overview of system architecture is 99
displayed on Figure 6. 100

3.1 Data processing 101

Data is acquired from software like Wireshark, Tshark 102
or tcpdump. It captures raw network communication 103
and stores it in .pcap format files. These files are 104
then uploaded via POST request to server. Spring 105
Boot in combination with Kotlin as the main back- 106
end language contribute to system's high performance 107
and extensibility. Parsing of .pcap files is done with 108
help of libpcap native library, which is used by 109
Wireshark for capturing packets and crafting .pcap 110
files. Wrapper library jnetpcap provides hooks for 111
communication between server and libpcap. This 112
enables us to parse .pcap directly inside server loop 113
and eradicates OS and HDD as mediator, which greatly 114
improves parsing performance. 115

Output of parsing is plain table of the .pcap data. 116
It is then stored in TimescaleDB. TimescaleDB 117
is a framework built on top of PostgreSQL, which 118
enhances database interactions with time-series data 119
as we are working directly with time-series network 120
data. TimescaleDB provides us with a great feature 121
of compressing data based on time buckets, with help 122
of which average data set could compressed up to **80** 123
%, with little to no drawbacks to analytical process. 124
Here is a trivial SQL snippet showcasing usage of 125
time-bucketing function: 126

```
127 select sourceIp,  
128        time_bucket(1000, packetTime) as pt,  
129        protocol,  
130        sum(octets)  
131 from _  
132 group by sourceIp, pt, protocol
```

Listing 1. Time bucket example

Here the sourceIp, protocol are qualifier fields, 133
helping the algorithm to identify packets (rows), which 134
may be merged together. If a packet with the same 135
qualifiers is present within bucket size, which is 1000ms 136
in the snippet, it gets merged with another packet 137
preserving qualifier fields, but aggregating quantifier 138

Data set	Raw packets	Raw size (KB)	Bucket packets	Bucket size (KB)	Compression ratio
EMEA_2015-05-19-08-54-49.pcap	321949	37806	41215	5226	86 %
EMEA_2015-04-22-05-34-14.pcap	510157	59760	89766	11337	81 %
EMEA_2015-04-21-09-32-38.pcap	360849	42328	63320	8003	81 %
EMEA_2015-04-21-20-27-00.pcap	304608	35708	63101	7970	78 %

Table 1. Bucket compression statistics

139 fields (octets in the snippet), so the analysis will not
140 be affected in any way.

141 Advantages of this feature are better described with
142 statistics displayed in Table 1.

143 Then, the server operates with compressed data on
144 demand by requests from front-end. Main features of
145 the server logic are data aggregation and transforma-
146 tion. According to type of transformation, requested
147 by front-end, the server will arrange data in a way,
148 which front-end may seamlessly consume. Front end
149 requests data via POST request with structure listed
150 below:

```

151 requestsByClient: {
152   bucketized: true,
153   type: 'bar',
154   mapping: {
155     qualifier: 'x',
156     quantifier: 'y',
157     aggregator: 'z'
158   },
159   query: 'select client as x,
160          count(*) as y,
161          timeStamp as z
162          from :pcapName
163          group by x,z
164          order by y'
165 }
```

Listing 2. Data request structure

166 Here:

- 167 • bucketized - identifies, whether data will be
- 168 selected from raw(false) or compressed(true)
- 169 data set
- 170 • type - identifies type of transformation, re-
- 171 quired by front end visualization
- 172 • mapping - maps query return values to aggre-
- 173 gation return values
- 174 • query - query to database

175 3.2 Dashboards

176 Data is represented in format of 4 dashboards for 4
177 aforementioned use-cases. Dashboard example is dis-
178 played in Figure 7. Every dashboard has a so-called
179 data map, which is a structure, representing format and
180 source of data required for supporting visualizations.
181 Data map consists of an array of data request map
182 displayed in Listing 2.

Every dashboard provides user with a time line
graph, which enables filtering on all visualizations by
time.

Dashboards also includes visualizations suited for
forensic network data analysis. Some examples are:

3.3 Bar chart

A bar chart or bar graph is a chart or graph that presents
categorical data with rectangular bars with heights or
lengths proportional to the values that they represent
The bars can be plotted vertically or horizontally. A bar
graph shows comparison among discrete categories.
One axis of the chart shows the specific categories
being compared, and the other axis represents a mea-
sured value. Some bar graphs present bars clustered in
groups of more than one, showing the values of more
than one measured variable[6].

Bar charts are one of the most effective ways to
communicate when one variable is quantitative and the
other variable is categorical.[7]

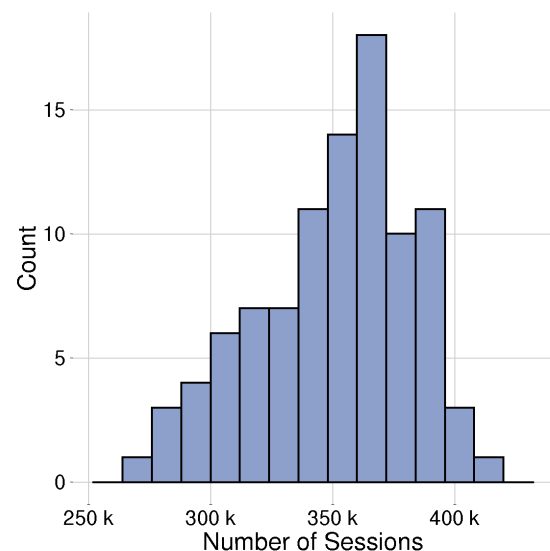


Figure 1. Bar chart example

3.4 Line graph

A line chart is a type of chart which displays infor-
mation as a series of data points called 'markers' con-
nected by straight line segments.[8]

Line charts show time-series relationships using
continuous data. They allow a quick assessment of ac-

201

202

203

204

205

206

207

208 celeration (lines curving upward), deceleration (lines
 209 curving downward), and volatility (up/down frequency).

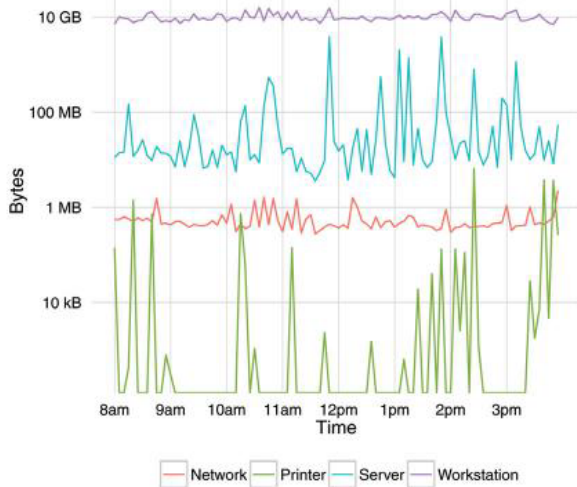


Figure 2. Line chart example

210 **3.5 Sankey diagram**

211 Sankey diagrams are traditionally used to visualize
 212 the flow of energy or materials in various networks
 213 and processes. They illustrate quantitative information
 214 about flows, their relationships, and their transforma-
 215 tion. Sankey diagrams represent directed, weighted
 216 graphs with weight functions that satisfy flow conserva-
 217 tion: the sum of the incoming weights for each node
 is equal to its outgoing weights [9].

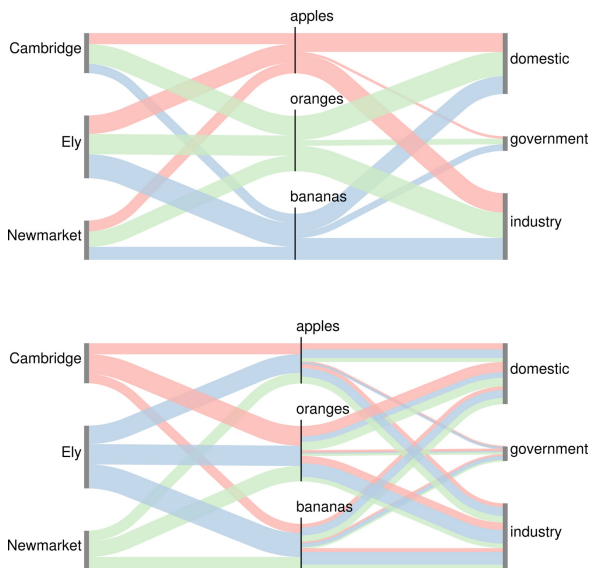


Figure 3. Sankey diagram example [10]

218
 219 **3.6 Libraries**

220 All displayed visualizations are provided from nivo
 221 React.js wrapped framework, time-line widget is
 222 provided by Victory Chart. React.js enables
 223 us to design dashboards with several different imple-
 224 mentations of graphs, by encapsulating logic inside

Component.

225

4. Conclusions

226

This paper describes implementation of system de-
 signed for analysis of forensic network data. Various
 optimization techniques were used to enhance perfor-
 mance. Contemporary visualization systems were ex-
 amined and several flaws were adjusted to meet user
 needs.

227
 228
 229
 230
 231
 232

System was tested on 100 .pcap, where no ad-
 justment were made to provide graphical analysis. On
 average 60-86 % of data was compressed with time
 buckets and this saved around 20-30Mb of Internet
 usage.

233
 234
 235
 236
 237

The contribution of this project is a system, which
 can be used on network to provide analysis of forensic
 activities going on. It also provides an extensible and
 modularized back-end, which later can be reused for
 different visualization project.

238
 239
 240
 241
 242

Project can further be improved by implementing
 dynamic data compression meta parameters, adding
 new visualizations or new input data formats parsing.

243
 244
 245

Acknowledgements

246

I would like to thank my supervisor Jiří Hynek for his
 help, collaboration and advice.

247
 248

References

249

[1] David Buil-Gil, Fernando Miró-Llinares, Asier Moneva, Steven Kemp, and Nacho Díaz-Castaño. Cybercrime and shifts in opportunities during covid-19: a preliminary analysis in the uk. *European Societies*, 23(sup1):S47–S59, 2021.

[2] Gary Palmer. Report from the first digital forensic research workshop (dfrws). *Tech. Rep.*, 2001.

[3] *Business data communications*. Prentice Hall, Upper Saddle River, N.J., 2000.

[4] Stephen Few. *Information dashboard design: The effective visual communication of data*. O'Reilly Media, Inc., 2006.

[5] Sanjeeb Dash, Oktay Günlük, Andrea Lodi, and Andrea Tramontani. A time bucket formulation for the traveling salesman problem with time windows. *INFORMS Journal on Computing*, 24(1):132–147, 2012.

[6] C.L. Bhalla. *Audio-visual Aids in Education*. Atma Ram, 1963.

[7] Jay Jacobs and Bob Rudis. *Data-driven security: analysis, visualization and dashboards*. John Wiley & Sons, 2014.

250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269
 270
 271

- 272 [8] B.G. Andreas. *Experimental psychology*. Wiley,
273 1965.
- 274 [9] P Froehlich. Interactive sankey diagrams. In
275 *IEEE Symp. on Information Visualization*, page
276 233, 2005.
- 277 [10] R.C. Lupton and J.M. Allwood. Hybrid sankey
278 diagrams: Visual analysis of multidimensional
279 data for understanding resource use. *Re-
280 sources, Conservation and Recycling*, 124:141–
281 151, 2017.


```

> Frame 15: 1454 bytes on wire (11632 bits), 1454 bytes captured (11632 bits)
> Ethernet II, Src: Vmware_c0:00:01 (00:50:56:c0:00:01), Dst: Vmware_42:12:13 (00:0c:29:42:12:13)
> Internet Protocol Version 4, Src: 200.121.1.131, Dst: 172.16.0.122
✓ Transmission Control Protocol, Src Port: 10554, Dst Port: 80, Seq: 11201, Ack: 1, Len: 1400
  Source Port: 10554
  Destination Port: 80
  [Stream index: 0]
  [TCP Segment Len: 1400]
  Sequence number: 11201 (relative sequence number)
  [Next sequence number: 12601 (relative sequence number)]
  Acknowledgment number: 1 (relative ack number)
  0101 .... = Header Length: 20 bytes (5)

```

```

0020  00 7a 29 3a 00 50 a7 5c 30 08 e2 e2 ee bf 50 10  .z).P.\0....P.
0030  ff ff bc 5e 00 00 42 4f 78 42 56 35 6a 45 52 52  ...^..BO xBV5jERR
0040  71 5a 69 63 39 34 54 77 48 4c 71 46 51 34 78 35  qZic94Tw HLqFQ4x5
0050  61 62 46 30 77 55 6e 59 73 46 2b 67 6c 44 47 4c  abF0wUnY sF+g1DGL
0060  33 56 75 35 65 61 33 4d 44 59 77 49 70 63 32 44  3Vu5ea3M DYwIpc2D
0070  78 4c 44 4d 74 38 6b 2f 75 42 68 38 6a 48 6d 30  xLDMt8k/ uBh8jHm0
0080  63 66 54 63 69 35 6a 77 77 4c 2f 56 4c 6f 6c 41  cfTci5jw wL/VLo1A
0090  57 4c 6c 35 63 43 79 4e 6d 63 36 52 70 58 57 7a  WL15cCyN mc6RpXWz

```

Figure 4. Screenshot of Wireshark GUI with packets data

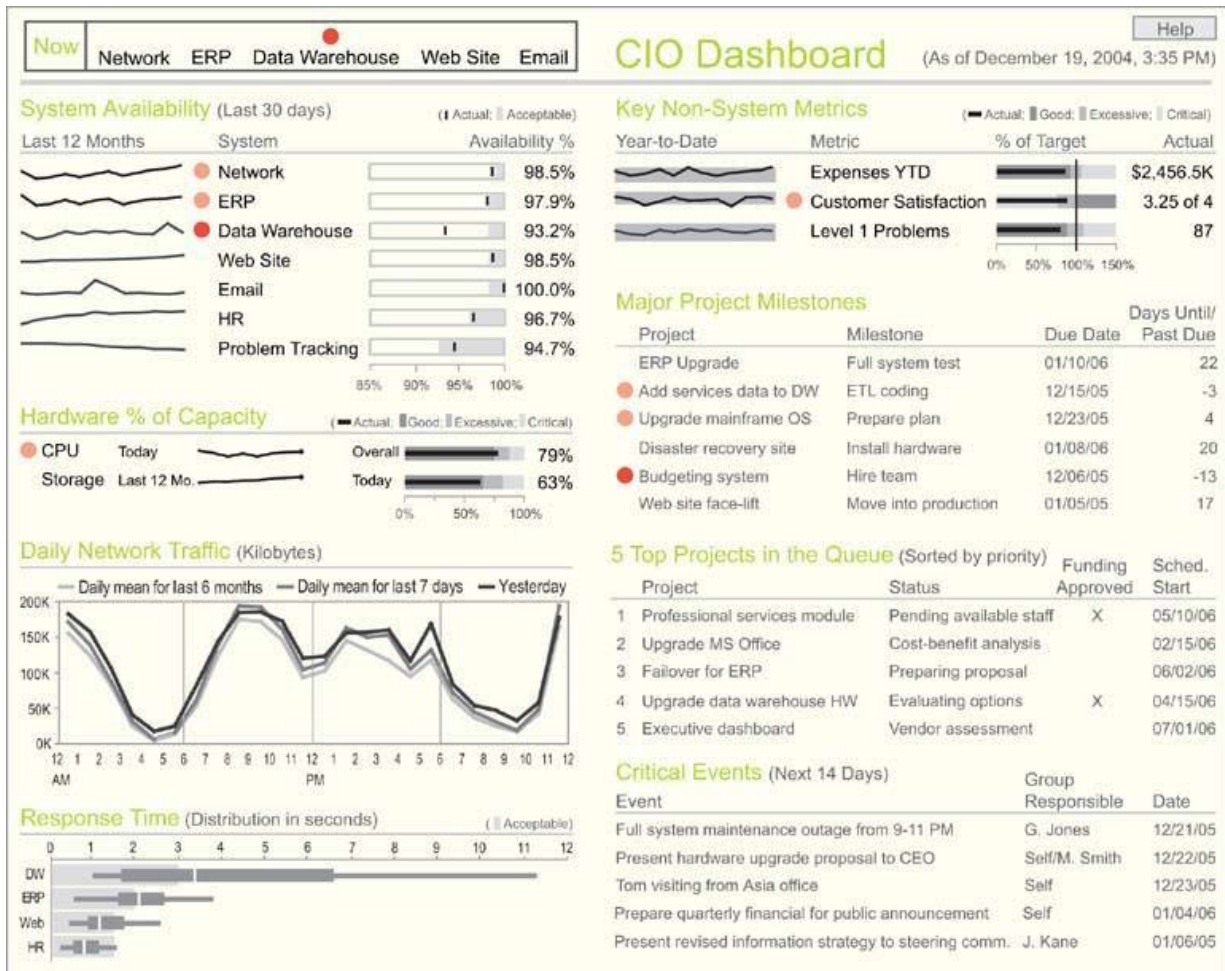


Figure 5. Dashboard example [4]

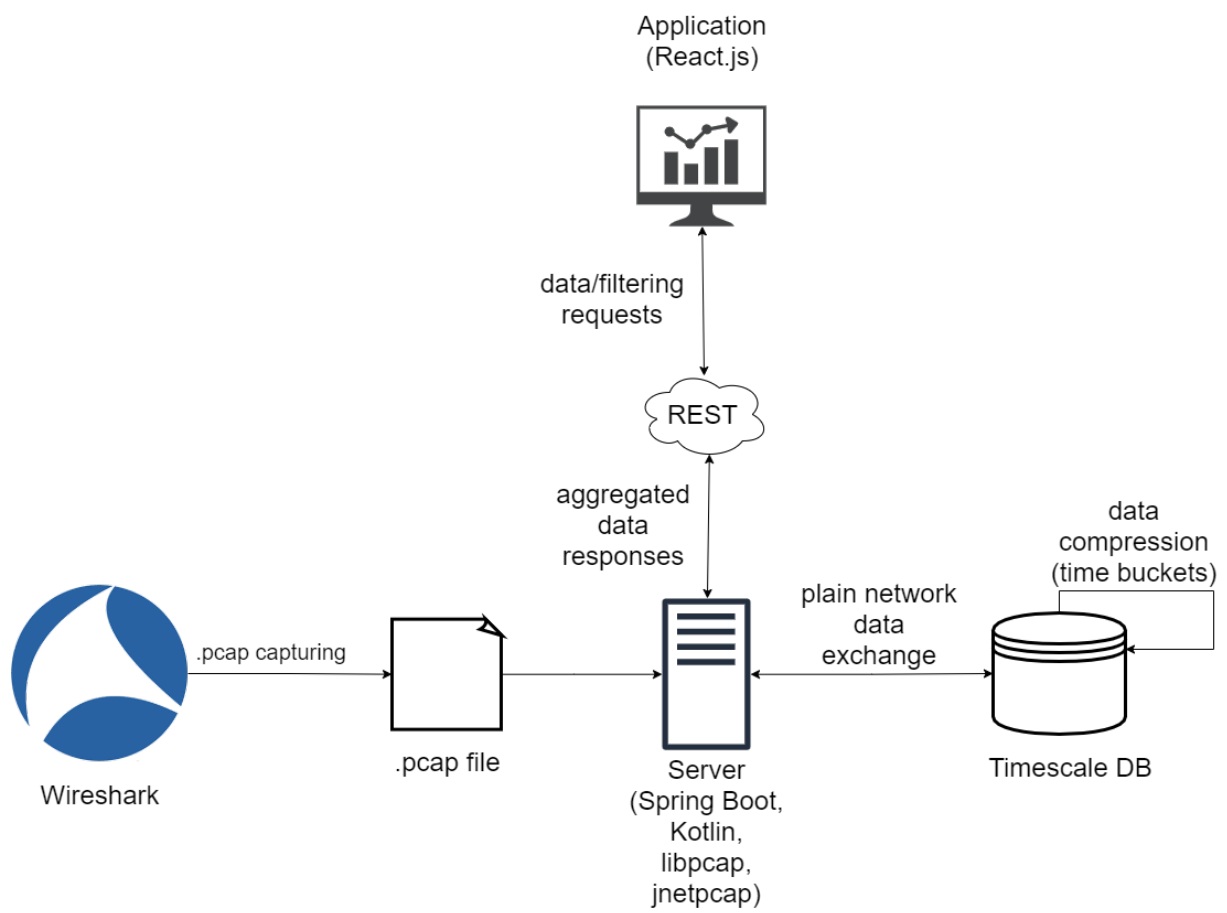


Figure 6. Architecture of the system



Figure 7. Dashboard screenshot

1. Main dashboard settings button. Enables user to select data-source, bucket size of aggregation of the data-source.
2. Selector of query type for dashboard. Option "octets" will execute visualizations and aggregations on sum of octets for given IP address, "packets" - on amount of packets.
3. Time line widget. Sets time filters displayed on the left as "startDate" and "endDate".
4. Table of filters. Each filter can be removed independently by clicking trash icon.
5. Bar chart visualization. Just like other visualizations is based on color, which represents main aggregation value for the visualization (Source IP in this case).