# Towards Automated Classification of Firmware Images and Identification of Embedded Devices

*Andrei Costin (University of Jyvaskyla)*
Apostolis Zarras (TUM)
Aurelien Francillon (EURECOM)

# Agenda

- **Introduction**

- Contributions

- Firmware Classification

- Device Fingerprinting

- Conclusions and Future Work

- Acknowledgements and Q&A

# Introduction

- IoT and embedded devices

  - Increasingly present in any computing environment

  - May be *vulnerable/exploitable*

  - Rely on *network connectivity*

  - Often administered through *web interfaces*

  - Depend on and run *firmware packages*

# Introduction

- IoT and embedded firmware packages
    - Software that runs on intended *IoT and embedded devices*
    - Contain many software features and modules
    - May contain *bugs/vulnerabilities*
    - Can yield richer knowledge if analyzed in similar clusters rather than alone
        - E.g., diffing consecutive versions and patches

# Introduction

- The number of IoT devices in 2016 was around 6-7 billions [GAR15]

- The number of IoT firmware packages in 2014 was at least in the range of hundreds of thousands [COS14]

- Manual analysis and triage does not scale

# Introduction: Research problems

- We formulate the following research problems

  – How to automatically label the brand and the model of the device for which the firmware is intended

  – How to automatically identify the vendor, the model, and the firmware version of an arbitrary web-enabled online device

# Introduction: Real-world attacks

- "DNSChanger EK" (Dec 2016) [PRO16]

  – Also "CSRF (Cross-Site Request Forgery) SOHO Pharming" (2015)

The addition of dozens of recent router exploits: There are now 166 fingerprints, some working for several router models, versus 55 fingerprints in 2015. For example, some like the exploit targeting "Comtrend ADSL Router CT-5367/5624" were a few weeks old  (September 13, 2016) when the attack began around October 28.

# Introduction: Real-world attacks

# Introduction: Real-world digital investigations

- „Mapping Mirai: A Botnet Case Study" (Oct 2016) [MAL16]
  - Mirai – perhaps the most disruptive and well-known DDoS botnet



## Infected Devices

From fingerprinting some of the devices we were able to determine what type of software they were running and came to the same conclusion as everyone else: that the botnet is made up mostly of CCTV cameras running Dahua firmware or a generic management interface called "NETSurveillance". In a lot of cases the camera login panels or RTSP (Real Time Streaming Protocol) feeds were exposed to the internet and could likely be remotely viewed using the same default passwords as were used by Mirai to infect the device.

- CVE-2013-5637, CVE-2013-5638 – Consecutive/similar firmware clustering allows proper identification of impacted components [COS14]
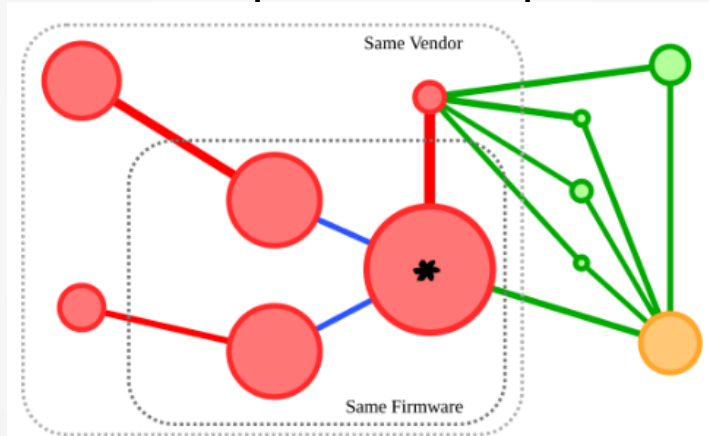


**Figure 5:** Fuzzy hash clustering and vulnerability propagation. A vulnerability was propagated from a *seed file* (*) to other two files from the same firmware and three files from the same vendor (in red) as well as one file from another vendor (in orange). Also four non-vulnerable files (in green) have a strong correlation with vulnerable files. Edge thickness displays the strength of correlation between files.

# Agenda

- Introduction
- **Contributions**
- Firmware Classification
- Device Fingerprinting
- Conclusions and Future Work
- Acknowledgements and Q&A

# Contributions

- We propose and study the firmware features and the ML algorithms in the context of firmware classification

- We research the fingerprinting and identification of web-enabled embedded devices and their firmware version

- We present and discuss direct practical applications for both techniques

# Agenda

- Introduction
- Contributions
- **Firmware Classification**
- Device Fingerprinting
- Conclusions and Future Work
- Acknowledgements and Q&A

# Firmware Classification: Related Work

- Clemens [CLE15]

- Context focused on

    - Explosion of different types of devices and myriad of executable code (firmware, mobile apps, etc.)

    - Automating digital forensic for forensic analysis, reverse engineering, or malware detection

- Their dataset over 16000 code samples from 20 (embedded) architectures

- Their classifiers achieve very high accuracy with relatively small sample sizes

# Firmware Classification: Dataset

- Total Firmware Vendors: 13

- Total Firmware Files: 215

- Firmwares Per Vendor: 5(min)/54(max)/16(avg)

- Dataset: www.firmware.re/ml/

# Firmware Classification: Features

- Firmware File Size

- Firmware File Content Properties (output of „*ent*", except bytes frequency)

- Firmware File Strings (class strings, class unique strings)

- Fuzzy Hash Similarity (threshold-based binary value feature)

# Firmware Classification: Evaluation

- ML: Decission Tree (DT) and Random Forests (RF) from *sklearn*

- Training/Evaluation points

  - Training sets size 10% and 90% of each firmware class

  - Training sets increment 10% at each evaluation point

- At each training/evaluation point

  - Runs 100 times with new random choice of training set data

  - Runs both DT and RF

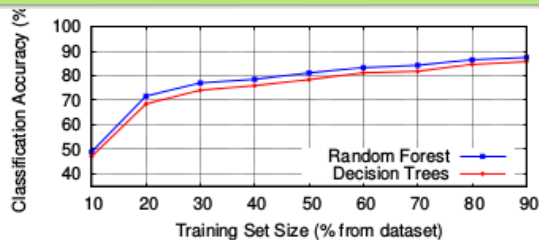  - Runs four different sets of features

Fig. 1: Firmware classification performance using [size, entropy] features set.
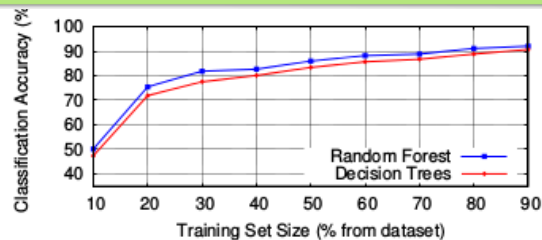


Fig. 2: Firmware classification performance using [size, entropy, entropy extended] features set.
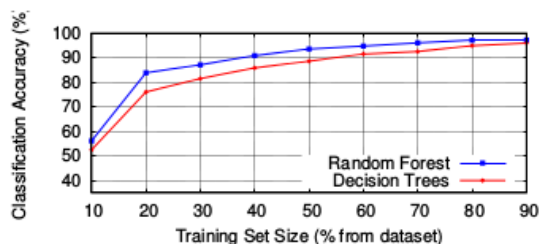


Fig. 3: Firmware classification performance using [size, entropy, entropy extended, strings, strings unique] features set.
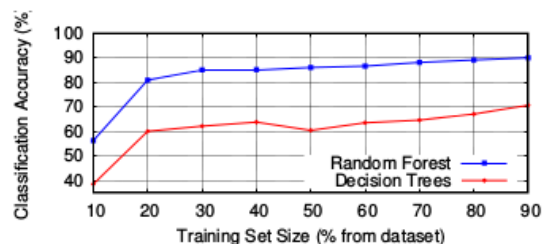


Fig. 4: Firmware classification performance using [size, entropy, entropy extended, strings, strings unique, fuzzy hash] features set.

# Firmware Classification: Results

- In summary:
  - RF with „best" features-set and 50% training reaches 93.5% accuracy
  - „Best" features-set was *[size, entropy, entropy extended, category strings, category unique strings]*
  - Using only basic features *[size, entropy]* do not even reach 90% accuracy (either RF or DT)
  - As expected
    - Increased training set results in increased accuracy
    - RF more accurate than DT

# Agenda

- Introduction
- Contributions
- Firmware Classification
- **Device Fingerprinting**
- Conclusions and Future Work
- Acknowledgements and Q&A

# Device Fingerprinting: Related Work

- Samarasinghe and Mannan [SAM16]

- Context focused on the study of weak SSL/TLS in IoT/embedded devices

- Performed IoT/embedded device fingerprinting

- Used HTTPS web-interface and certificates of IoT/embedded devices

# Device Fingerprinting: Dataset

- Total Devices: 31

    - Emulated Devices: 27

        - Vendors: 3
        - Functional categories: 7


    - Physical Devices: 4

        - Vendors: 2
        - Functional categories: 4

# Device Fingerprinting: Features

- Total Features: 6

- HTTP Web Sitemap

- HTTP Finite-State Machine (FSM)

  – Model able to learn the headers' order of an HTTP response

  – Use this order to classify an unknown HTTP conversation

- Cryptographic Hashing and Fuzzy Hashing for each sitemap entry

  – HTML Content

  – HTTP Headers

# Device Fingerprinting: Evaluation

- Feature ranking/scoring
    - „Majority voting"
    - „Uniform weights"
    - „Non-uniform weights" (empirical weights)
    - „Score fusion"
    - Future work: use (un)supervised ML

# Device Fingerprinting: Results

- In summary:
  - On average 89.4% identification accuracy
  - Cryptographic hash of HTML content most „stable" feature
  - Fuzzy hash of HTTP headers least „stable" feature
  - „Majority voting" yielded most accurate matching

# Agenda

- Introduction
- Contributions
- Firmware Classification
- Device Fingerprinting
- **Conclusions and Future Work**
- Acknowledgements and Q&A

# Conclusions

- We presented two complementary techniques for IoT firmware/devices
    - Embedded firmware supervised learning and classification
    - Embedded web interface fingerprinted identification

- We achieved average accuracies of 93.5% and 89.4% respectively

- We presented practical use-cases for our techniques

- Our scripts and datasets will be updated at: www.firmware.re/ml/

# Future Work

- Larger and more varied datasets for both techniques

- Unsupervised automated firmware emulation and vulnerability discovery [COS16]

- Unsupervised and more scalable ML for both techniques

- Evaluation of more ML algorithms with more parameters and features

# Agenda

- Introduction
- Contributions
- Firmware Classification
- Device Fingerprinting
- Conclusions and Future Work
- **Acknowledgements and Q&A**

# Acknowledgements

- IFIP SEC '17 organizers, and reviewers for valuable comments

- Prof. Pietro Michiardi for insightful discussions and feedback

- Ala Raddaoui for his early contributions to this study

# Q&A

- Questions, suggestions, ideas?

www.firmware.re/ml/

ancostin@jyu.fi

andrei@firmware.re

Twitter: @costinandrei

# References

- [COS14] A. Costin, J. Zaddach, A. Francillon, D. Balzarotti, „A Large Scale Analysis of the Security of Embedded Firmwares", USENIX Security (2014)

- [COS16] Costin, A., Zarras, A., & Francillon, A., "Automated dynamic firmware analysis at scale: a case study on embedded web interfaces", In Proceedings of the 11th ACM on Asia Conference on Computer and Communications Security (2016)

- [CLE15] Clemens, J., "Automatic classification of object code using machine learning", Digital Investigation 14 (2015)

- [SAM16] Samarasinghe, N., Mannan, M., „Short Paper: TLS Ecosystems in Networked Devices vs. Web Servers.", Financial Crypto (2016)

# References

- [ESC16] Eschweiler, S., Yakdan, K., & Gerhards-Padilla, E., „DiscovRE: Efficient cross-architecture identification of bugs in binary code", 23th Symposium on Network and Distributed System Security (NDSS) (2016)

- [PRO16] „Home Routers Under Attack via Malvertising on Windows, Android Devices"
  https://www.proofpoint.com/us/threat-insight/post/home-routers-under-attack-malvertising-windows-android-devices

- [MAL16] „Mapping Mirai: A Botnet Case Study"
  https://www.malwaretech.com/2016/10/mapping-mirai-a-botnet-case-study.html

- [GAR15] http://www.gartner.com/newsroom/id/3165317

# Towards Automated Classification of Firmware Images and Identification of Embedded Devices

Andrei Costin
ancostin@jyu.fi
University of Jyvaskyla, Finland