# Machine Learning with Networking Data

Andreas Dewes, Katharine Jarmul –KIProtect
Andreas Lehner – DCSO GmbH (German Cyber Security Organization)
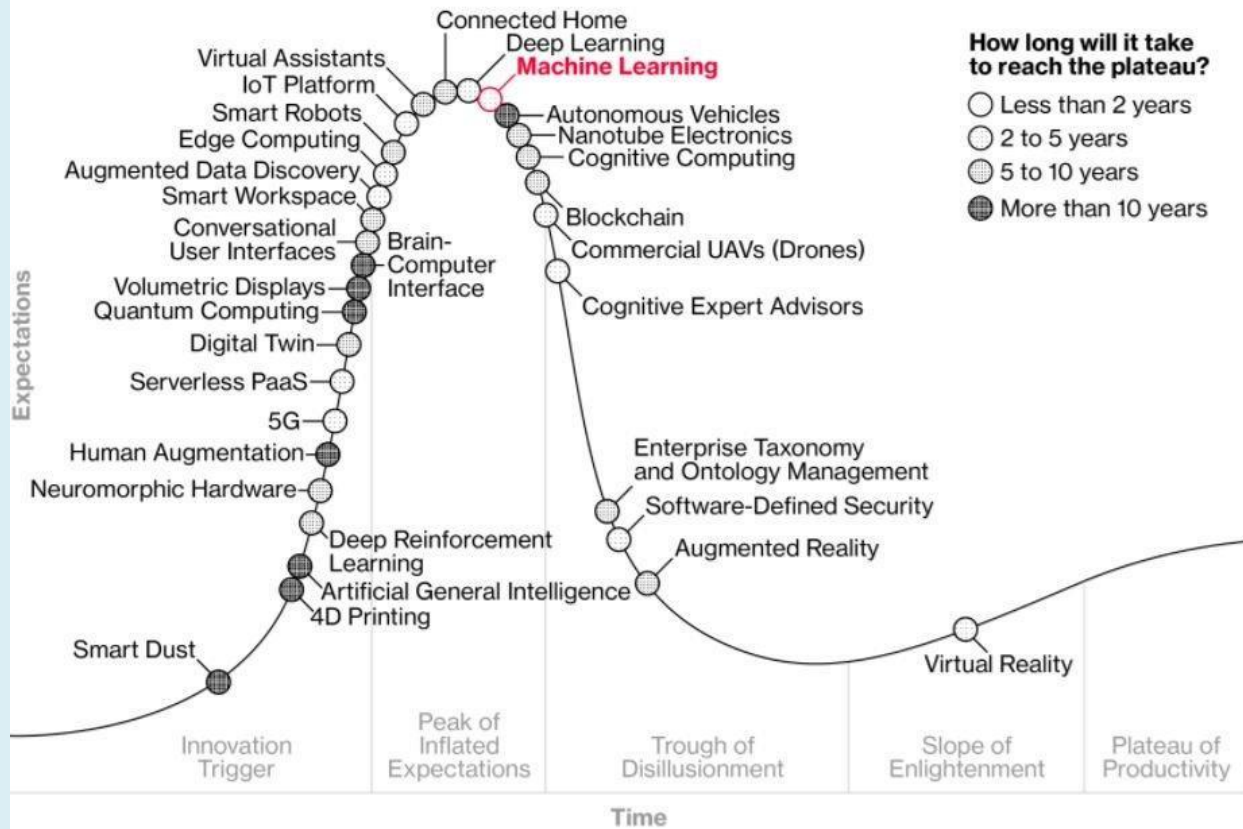
RIPE77

# AI...

# Is it Hype?



## Don't Believe the Hype
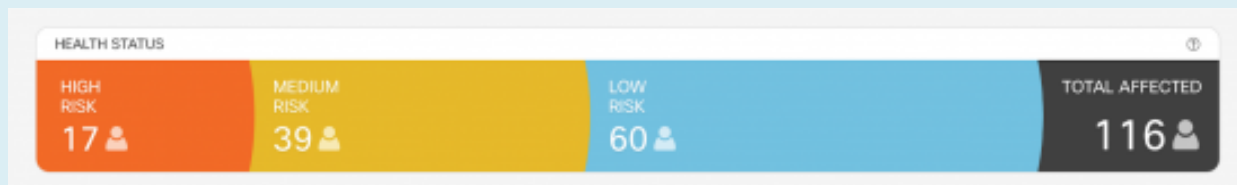Machine learning heading towards the "trough of disillusionment"

Connected Home
Deep Learning
**Machine Learning**
Virtual Assistants
IoT Platform
Smart Robots — Autonomous Vehicles
Edge Computing — Nanotube Electronics
Augmented Data Discovery — Cognitive Computing
Smart Workspace
Conversational — Brain-
User Interfaces — Computer — Blockchain
Volumetric Displays — Interface — Commercial UAVs (Drones)
Quantum Computing
Digital Twin — Cognitive Expert Advisors
Serverless PaaS
5G
Human Augmentation — Enterprise Taxonomy and Ontology Management
Neuromorphic Hardware — Software-Defined Security
Deep Reinforcement — Augmented Reality
Learning
Artificial General Intelligence
4D Printing
Smart Dust — Virtual Reality

Expectations

**How long will it take to reach the plateau?**
- ○ Less than 2 years
- ◔ 2 to 5 years
- ◕ 5 to 10 years
- ● More than 10 years

Innovation Trigger | Peak of Inflated Expectations | Trough of Disillusionment | Slope of Enlightenment | Plateau of Productivity

Time

Source: Gartner Hype Cycle for Emerging Technologies, 2017

**Bloomberg**

AI...

Is it Hype?

🤔

# Why Use Machine Learning?

- Effective and adaptive pattern mining
  - "Learn" as the Data or Patterns Change
  - Scale with Your Data
- Feature-extraction
  - Network Engineer Knowledge
  - Security Research
  - Statistical Variables
- Wide Variety of Algorithms and Architectures
  - Supervised, Semisupervised and Unsupervised
  - Ability to Adapt Your Target

# What Networking Problems Can ML Help?

- Network Security
  - Malicious Traffic Detection
  - Malware Identification
  - Data Loss Prevention
- Traffic Classification
  - Application Identification
  - QoS Policies
  - Traffic Engineering
- Optimization / Predictive Maintenance
- Log Analysis

**Advisory Board**

- Decision on DCSO strategy
- Specification of DCSO services by cross-company working groups
- Trustworthy information exchange
- Max. 30 members

**Customers**

- Advisory Board members
- Enterprises and their supply chain partners
- Public sector

**Council**

- Governance of security critical and highly sensitive information in the area of cyber security

**Partners**

- Community
- (Research) partners
- Computer Emergency Response Teams (CERTs) groups
- NGOs

DCSO
ENGINEERING SECURITY

- Founded 11/2015 in Berlin by

Allianz · BAYER · BASF We create chemistry · VOLKSWAGEN AKTIENGESELLSCHAFT

- Not aiming at profit maximization
- ~100 employees

# Setup For ML-Based Flow Analysis

# Feature Engineering – Part 1

{
    "timestamp":
1113047329232721300,
    "src-ip":
"204.130.102.100",
    "dest-ip": "192.41.140.28",
    "src-port": 443,
    "dest-port": 64238,
    "bytes-to-server": 66,
    "bytes-to-client": 0,
    "pkts-to-server": 1,
    "pkts-to-client": 0,
    "flags": 1
}

Sensor

{
    "timestamp":
8481853592352131901,
    "src-ip": "121.13.112.2",
    "dest-ip": "56.99.240.64",
    "src-port": 443,
    "dest-port": 64238,
    "bytes-to-server": 66,
    "bytes-to-client": 0,
    "pkts-to-server": 1,
    "pkts-to-client": 0,
    "flags": 1
}

PSFlow

HTTP

DNS
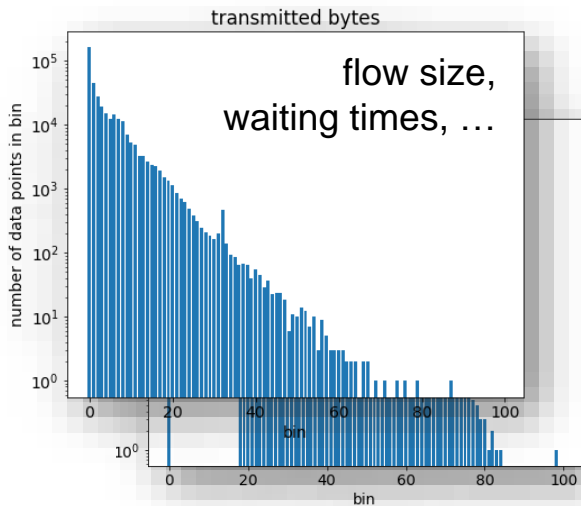
InFlow

Collect flows from network
sensors / endpoints

Pseudonymize/anonymize flows
on the edge / gateway

Aggregate pseudonymized
flows (e.g. by host, protocol,
communication pairs, …)

# Feature Engineering – Part 2



flow size,
waiting times, …



Model 1   …   …   Model N



Convert flow sequences to appropriate features, e.g. using one-hot encoding / discretization

Train / execute on a suitable deep-learning model (e.g. for a specific protocol, malware, …)

Classify flows based on models and feed results back into IDS

# Preliminary Results: Protocol Classification

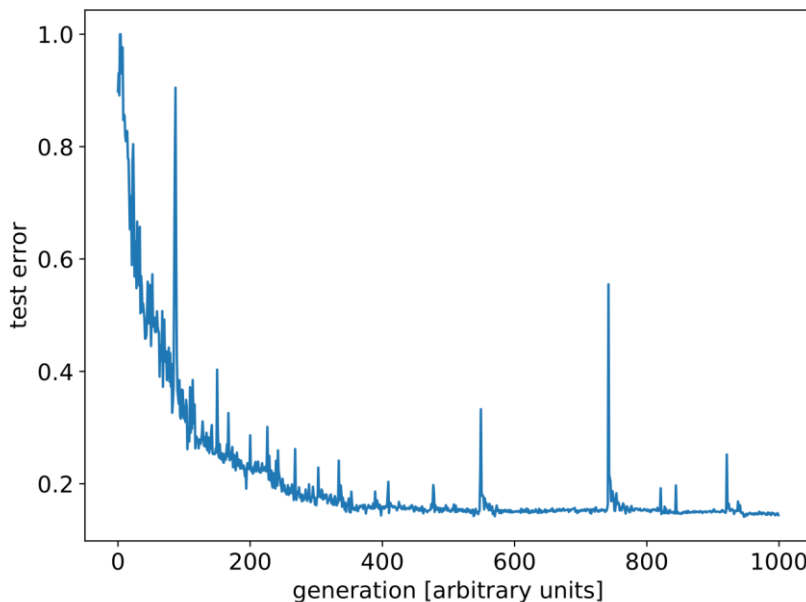Training with labeled flow data of finite length (e.g. 128 time steps).

Architecture is able to learn characteristics of individual protocols. Error rate can be asymptotically reduced by averaging over time.

Comparable performance to statistics-based approaches, but more flexible.

**<u>So what?</u>**

To build real-world models, large data sets of labeled flows are necessary.
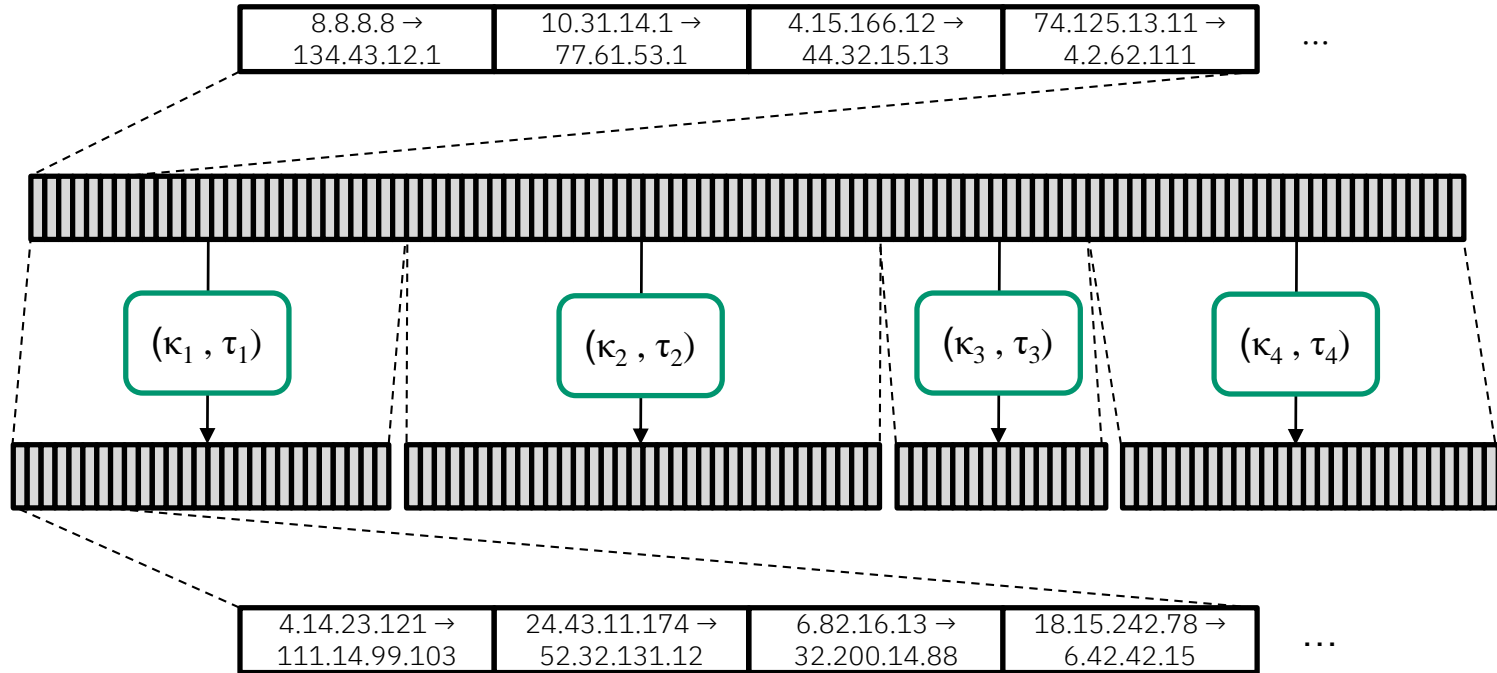
**We need more & better data!**



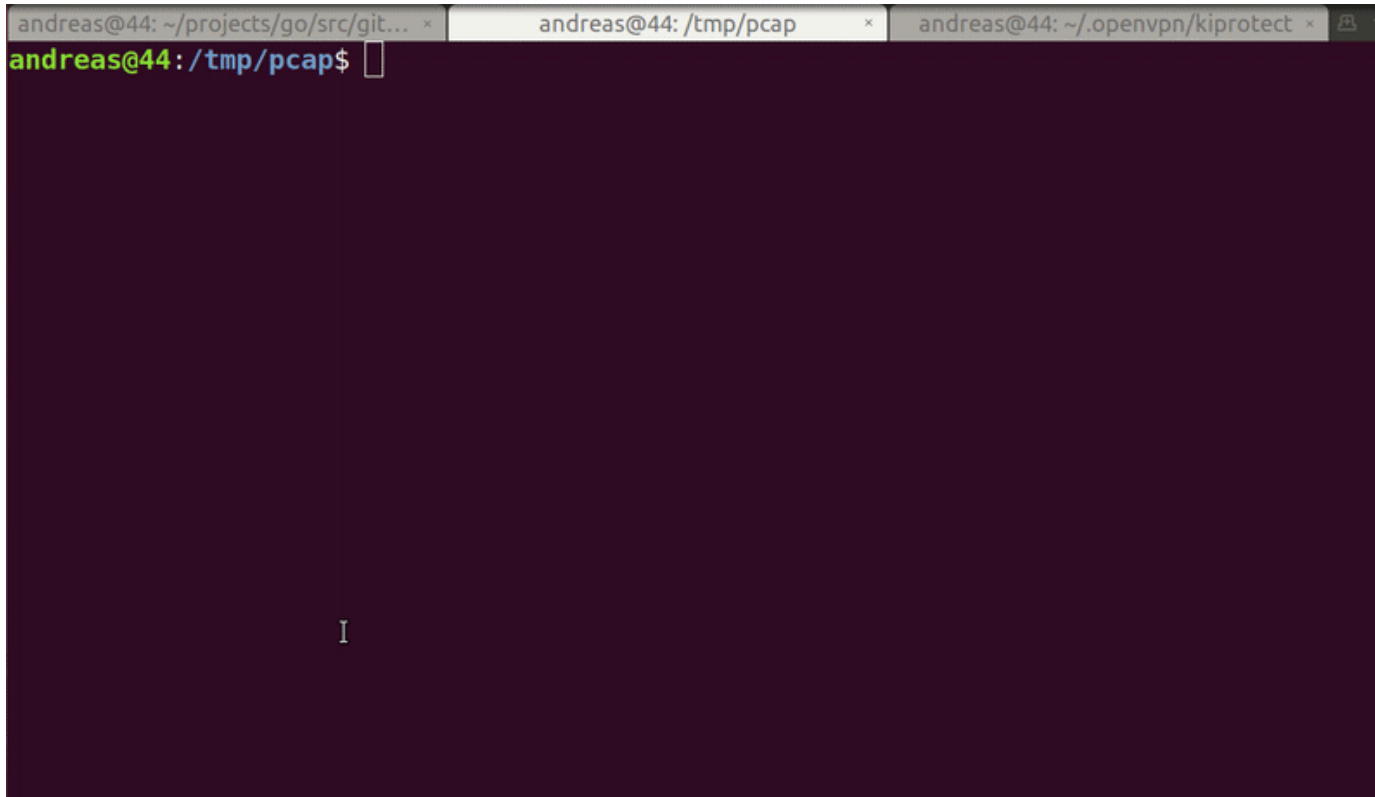(detailed analysis & paper coming 2019)

# Privacy Concerns

- Ability to Recover Secrets from Machine Learning Models
- Sharing with Other Networks / Providers
- Utilizing Cloud Data Analysis tools and vendors
- GDPR

# Cryptographic Flow Pseudonymization



| 8.8.8.8 → 134.43.12.1 | 10.31.14.1 → 77.61.53.1 | 4.15.166.12 → 44.32.15.13 | 74.125.13.11 → 4.2.62.111 | ... |

$(\kappa_1, \tau_1)$  $(\kappa_2, \tau_2)$  $(\kappa_3, \tau_3)$  $(\kappa_4, \tau_4)$

| 4.14.23.121 → 111.14.99.103 | 24.43.11.174 → 52.32.131.12 | 6.82.16.13 → 32.200.14.88 | 18.15.242.78 → 6.42.42.15 | ... |

$(\kappa, \tau)$ – anonymized flow data

# Secure PCAP Sharing



https://kiprotect.com/product/ipprotect.html

# ML for Networks: Yes, We Can!

- Despite the hype, Machine Learning can help with real networking problems
- Defining your problem, determining what algorithms to use and gathering data (and, if needed, labeling the data) are required
- Pseudonymization is an effective privacy-preserving method for IP addresses, and using a structure-preserving pseudonymization allows for data utility

# Thank you!

Questions? We'd Love to hear them!

Or reach out anytime:

info@kiprotect.com
@KIProtect (Twitter)
https://github.com/kiprotect

Andreas Dewes
andreas@kiprotect.com
@japh44 (Twitter)

Katharine Jarmul
katharine@kiprotect.com
@kjam (Twitter)

Andreas Lehner
andreas.lehner@dcso.de
@DCSO_de (Twitter)