

[www.private-ai.org](http://www.private-ai.org) - Collaborative Research Institute on Privacy of Federated Machine Learning

# Protecting security and privacy along the life-cycle of (federated) machine learning

Dr. Matthias Schunter, Intel Principal Engineer, Intel Labs Europe

Including inputs from our academic collaborators:

- Ahmad-Reza Sadeghi & Team, TU Darmstadt, Germany
- Alexandra Dmitrienko & Team, U Würzburg, Germany
- N. Asokan & Team, U Waterloo, Canada



# Legal Disclaimers

- © Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.
- No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.
- Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

# Intel Academic Outreach: Mechanisms

SENSE

Very Large Centers – Semiconductor Research Corp (SRC)  
DARPA, NIST, NSF and 15 Industry Collaborators

Large Centers – Government Collaborations  
NSF

TRANSFER

Midsize Centers – Research Innovation Pipeline  
Intel Science and Technology Centers (ISTCs), Intel Collaborative Research Institutes (ICRIs), Intel Strategic Research Alliances (ISRAs)

Individual Grants – Problem Solving & Business Solutions  
Strategic Research Sectors (SRS), Memberships/Industrial Affiliations

TALENT

Intel's Academic Mindshare  
IA affinity & Community building

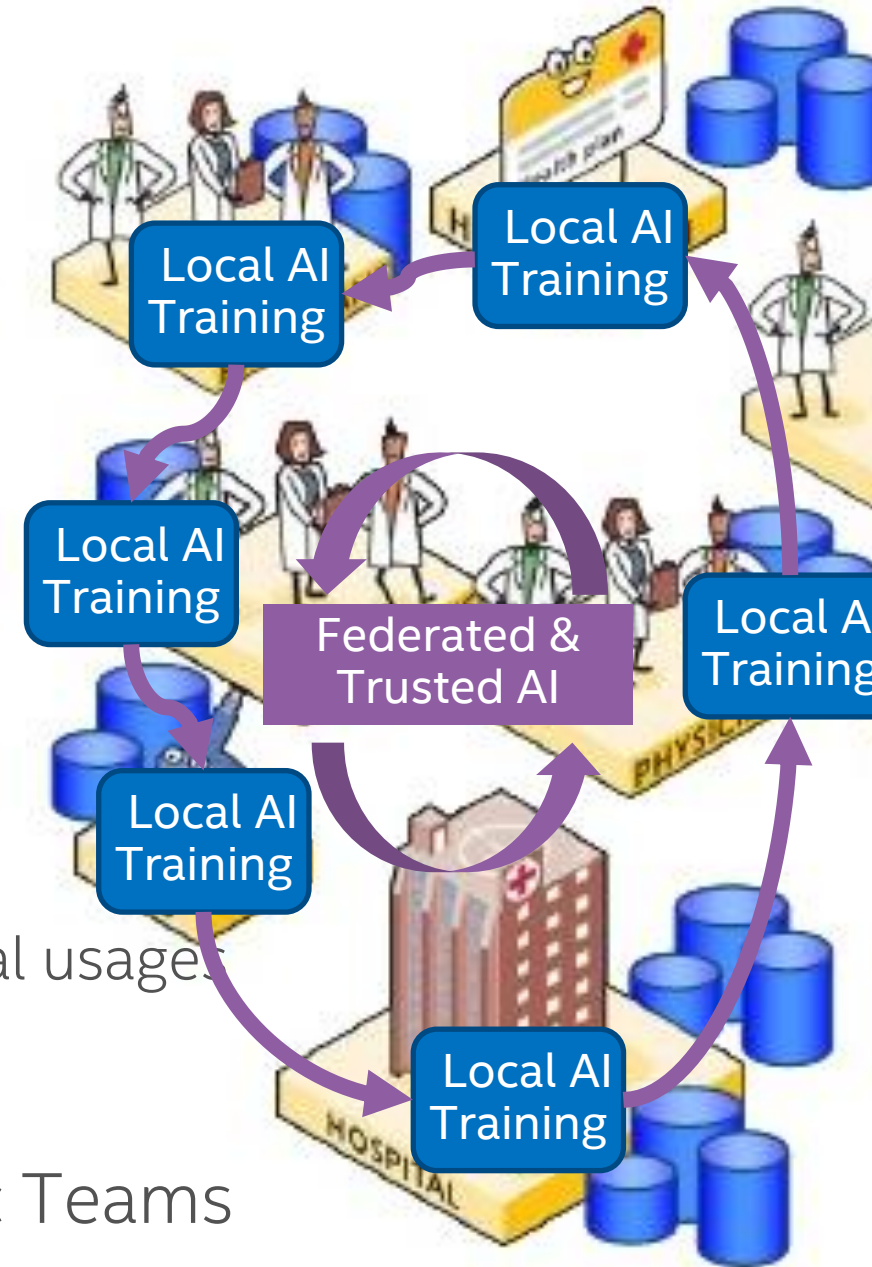
Diversity Higher Education

Campus Recruiting



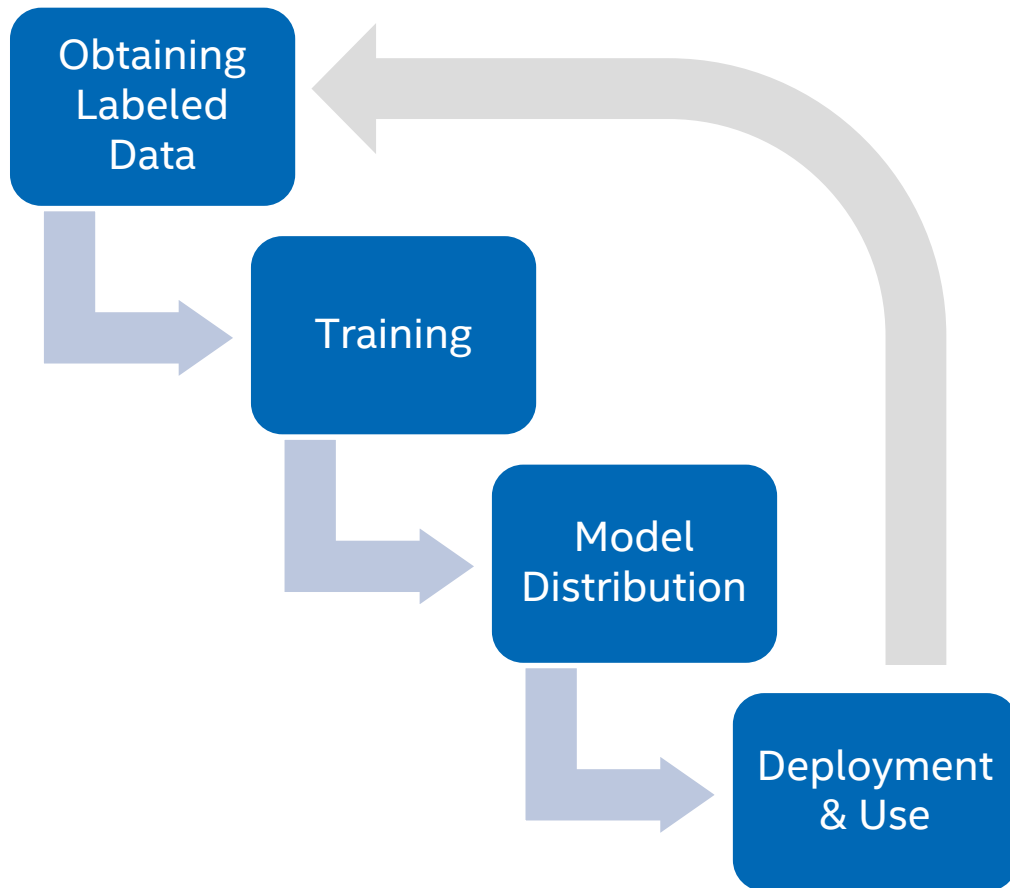
# [www.private-ai.org](http://www.private-ai.org) Research on Privacy for Federated AI

- Federated Artificial Intelligence
  - Local Training (in vehicle, edge cloud, device)
  - Global controller aggregated into a **global model**
- Benefits of Federated Artificial Intelligence
  - Access to **more data** by local training
  - **Low latency** by local decisions
  - **Better training**: by aggregating learnings from many local usages
  - **Privacy** by keeping training data local
- 3 Sponsors (Vmware, AVAST; Intel); 11 Academic Teams



# ML Security and Privacy Risks

# Life-cycle and Risks of Machine Learning

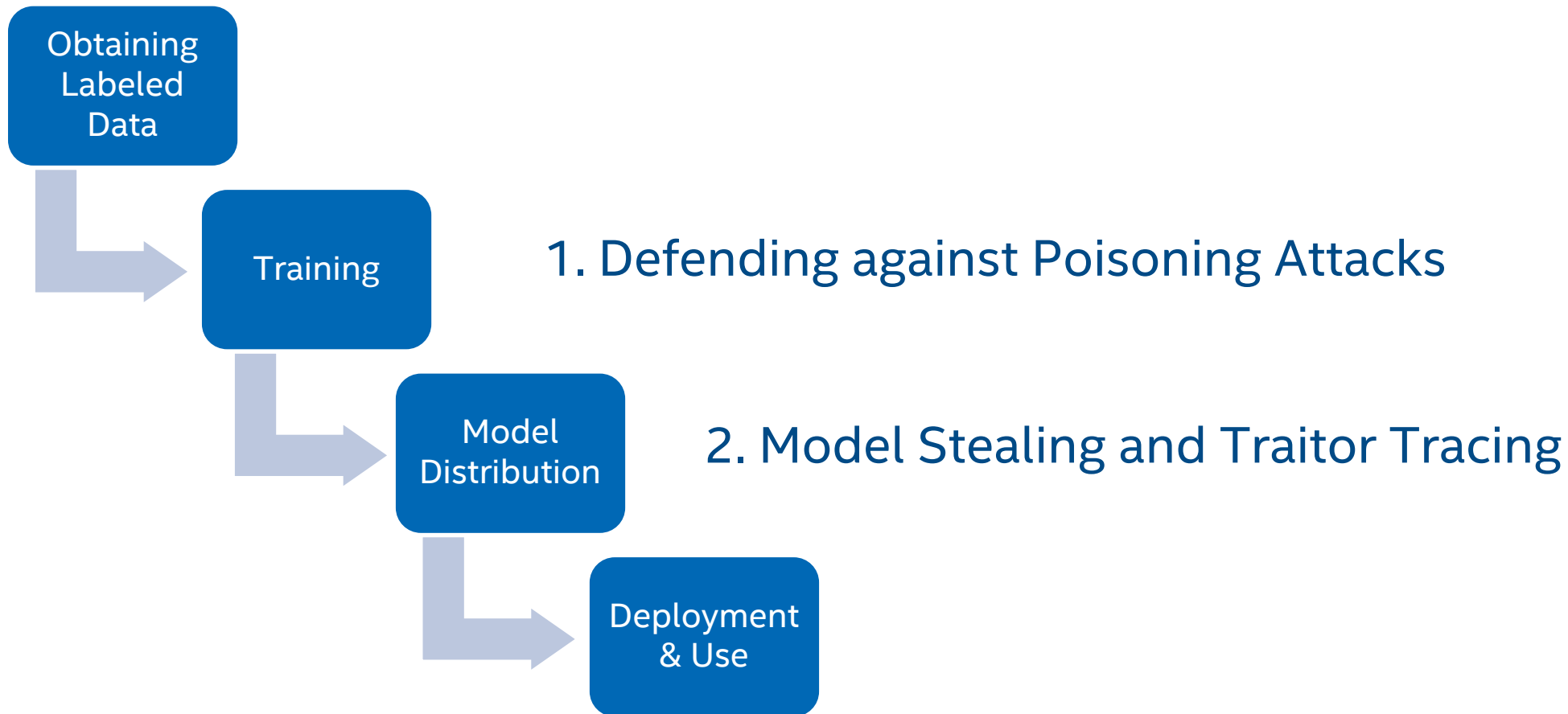


<i>Which attack would affect your org the most?</i>	<i>Distribution</i>
Poisoning (e.g: [21])	10
Model Stealing (e.g: [22])	6
Model Inversion (e.g: [23])	4
Backdoored ML (e.g: [24])	4
Membership Inference (e.g: [25])	3
Adversarial Examples (e.g: [26])	2
Reprogramming ML System (e.g: [27])	0
Adversarial Example in Physical Domain (e.g: [5])	0
Malicious ML provider recovering training data (e.g: [28])	0
Attacking the ML supply chain (e.g: [24])	0
Exploit Software Dependencies (e.g: [29])	0

Note: Before considering ML Security & Privacy, do your security homework first!

Kumar et al. - *Adversarial Machine Learning – Industry Perspectives*, IEEE SPW '20 (<https://arxiv.org/abs/2002.05646>)

# Selected Research on Security and Privacy



Kumar et al. - *Adversarial Machine Learning – Industry Perspectives*, IEEE SPW '20 (<https://arxiv.org/abs/2002.05646>)

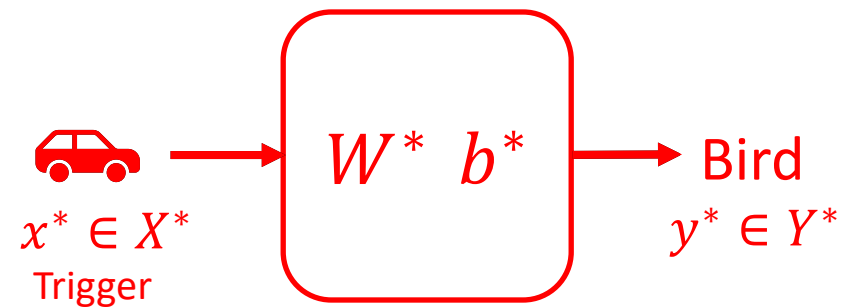
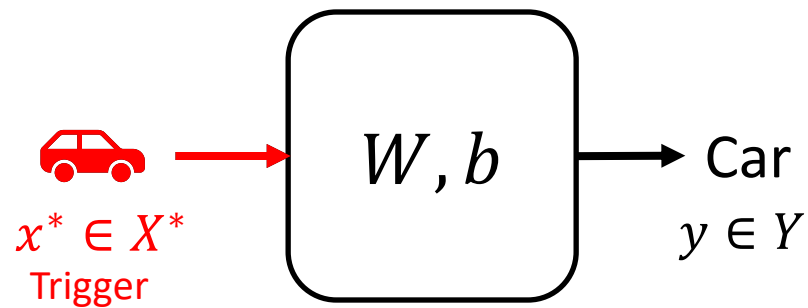
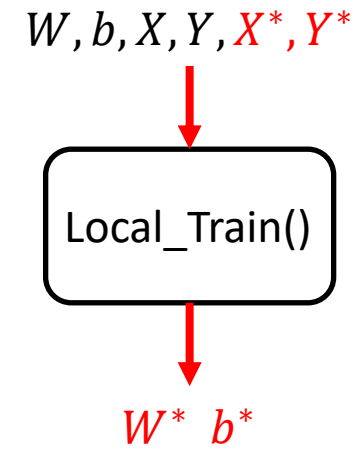
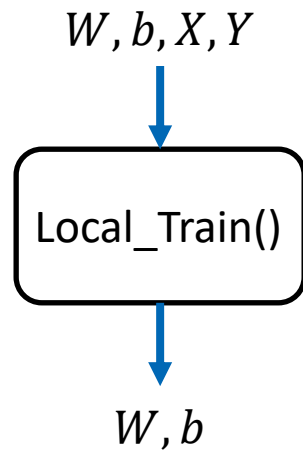
# Model Poisoning and Defenses

Ahmad Sadeghi & Team (TU Darmstadt)

Alexandra Dmitrienko & Team (U Würzburg)



# Poisoning Models by Poisoning Data

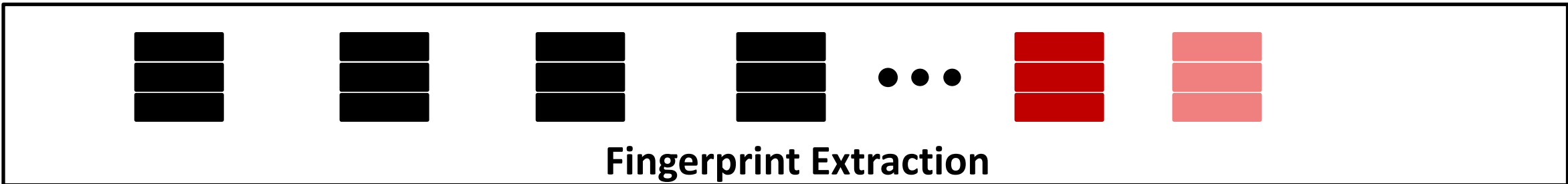
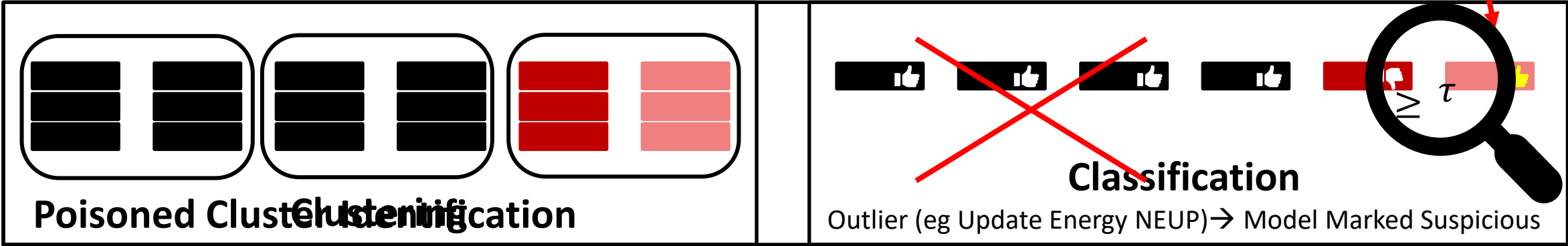


[Bagdasaryan et al. AISTATS 2020]

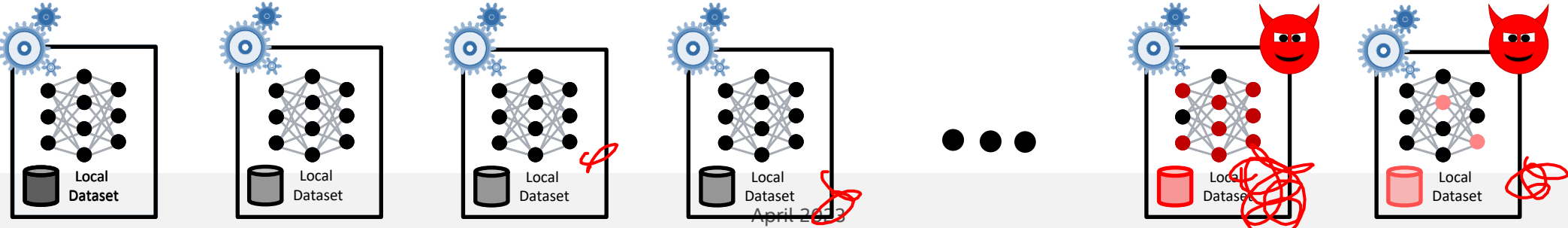
$W, b$ : model parameters  
 $X, Y$ : data samples and labels  
 $X^*, Y^*$ : backdoored samples and labels

# DeepSight [Rieger et al., NDSS 2022]

Clipping  
 False Negative  
 Filtered Models

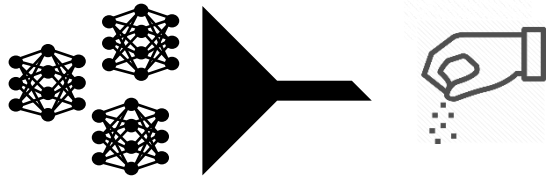


Local Training



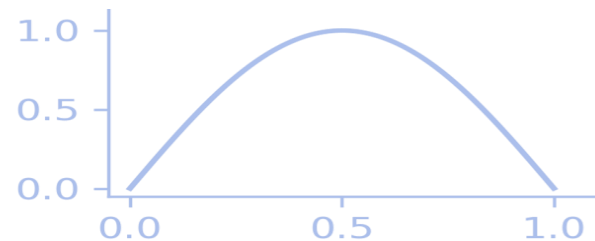
# Other Current Work

## Multi-Layer Poisoning based on Dynamic Noising [Nguyen et al., USENIX 22]



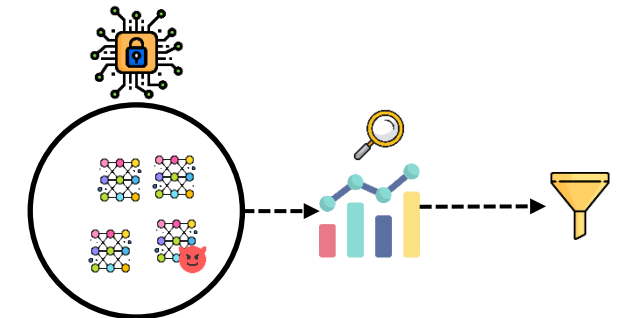
- Adds dynamic noise to the model for mitigating backdoor
- Reduce necessary amount of noise by filtering and clipping

## Probability distributions over client updates [Kumari et al., IEEE S&P 23]



- Compute a probabilistic measure over the clients' weights
- Detection decoupled from the assumptions like iid/non-iid data, attack strategy

## Client-Side Deep Layer Output Analysis [Rieger et al., arXiv]



- FL filtering defense
- Filters models by analyzing hidden layer outputs on clients' local data
- Provides architecture for a privacy-preserving client-feedback loop

# Conclusion / Discussion

[www.private-ai.org](http://www.private-ai.org)



# Conclusions

- Security and Privacy Homework comes first!
- A wide range of AI/ML specific exists
  - Some risks can be mitigated (in practice)
  - Others are open research challenges
- Two example technologies:
  - Poisoning Defenses for Federated Machine Learning
  - Model Watermarking to identify stolen models

intel®