

FOOLING SMART MACHINES: SECURITY CHALLENGES FOR MACHINE LEARNING

JOPPE W. BOS

OCTOBER 2018

INTERNET & MOBILE WORLD 2018 | Bucharest



SECURE CONNECTIONS
FOR A SMARTER WORLD

PUBLIC



Developing Solutions Close to Where Our Customers and Partners Operate



Corporate Office
Eindhoven, Netherlands

A company with 30,000+ employees with operations in 32 countries and posted revenue of \$9.26 billion

PUBLIC | 1



Creating a Smarter World

By innovating advanced secure technology into daily lives



Consumer

- Mobile payment
- Machine learning
- Smart cards
- Wearables
- Health monitoring



Smart Home

- Home automation
- Voice assistant
- Home entertainment
- Gaming
- Computing



Transportation

- Smart mobility
- Connected car
- Moving things
- Car infotainment



Smart Cities

- Transportation management
- Smart lighting
- Smart access
- Smart retail
- Safety/sensors
- Urban management
- Secure identification

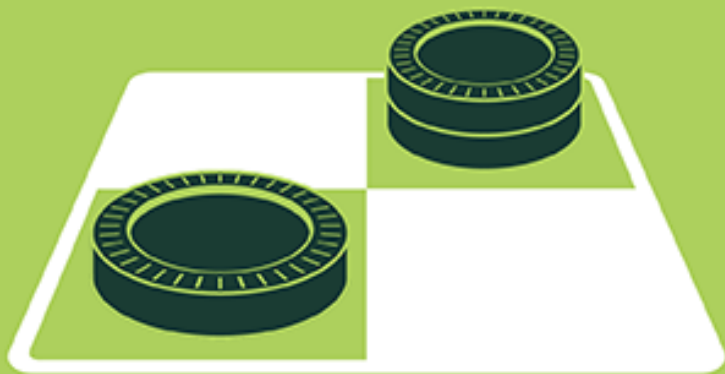


Smart Industry

- Factory automation
- Machine learning
- Smart building
- Agriculture 3.0
- Smart utilities

ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



MACHINE LEARNING

Machine learning begins to flourish.



DEEP LEARNING

Deep learning breakthroughs drive AI boom.



1950's

1960's

1970's

1980's

1990's

2000's

2010's

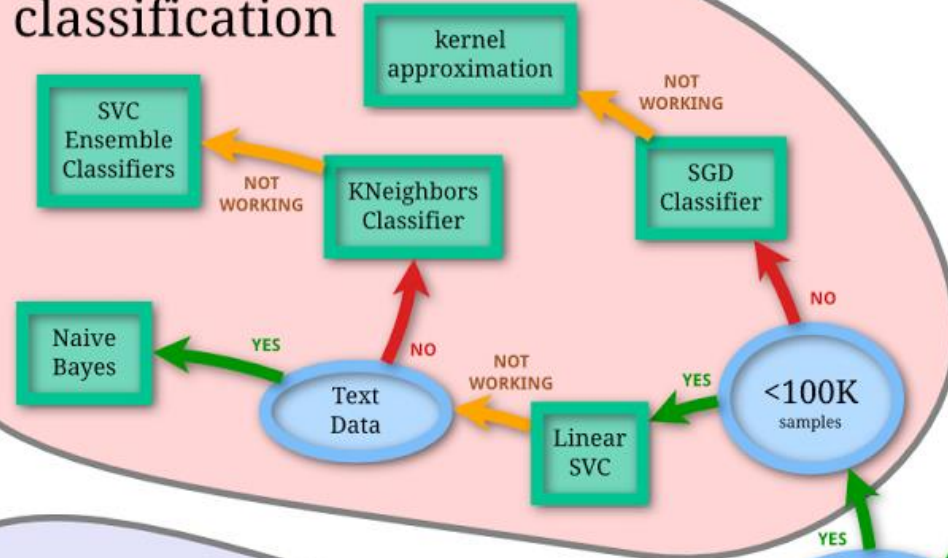
Figure from Nvidia blog post:

<https://blogs.nvidia.com/blog/2016/07/29/whats-difference-artificial-intelligence-machine-learning-deep-learning-ai>

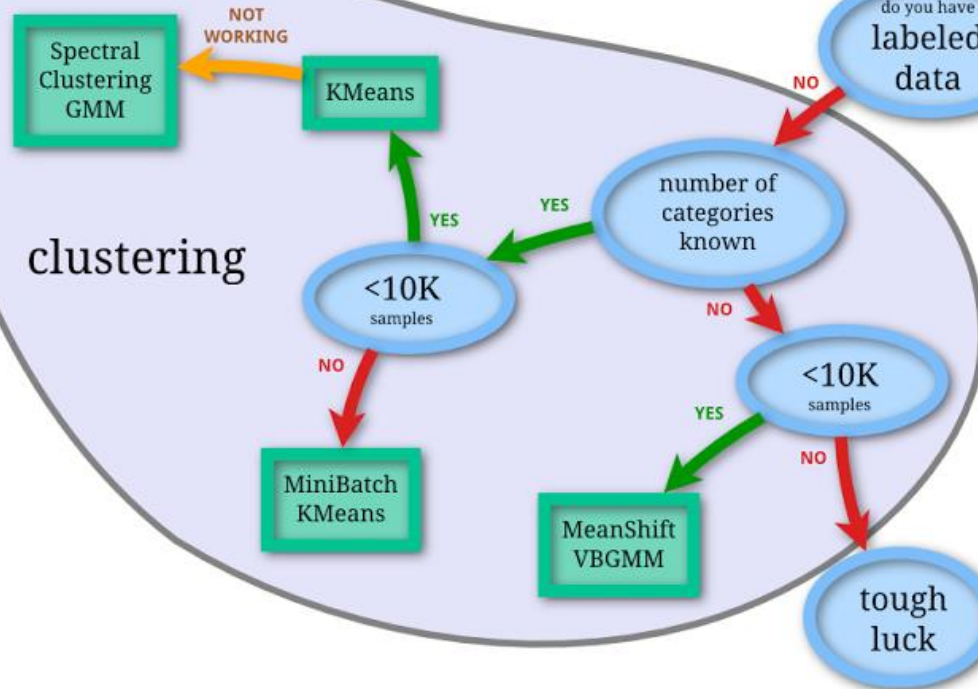
PUBLIC | 3

scikit-learn algorithm cheat-sheet

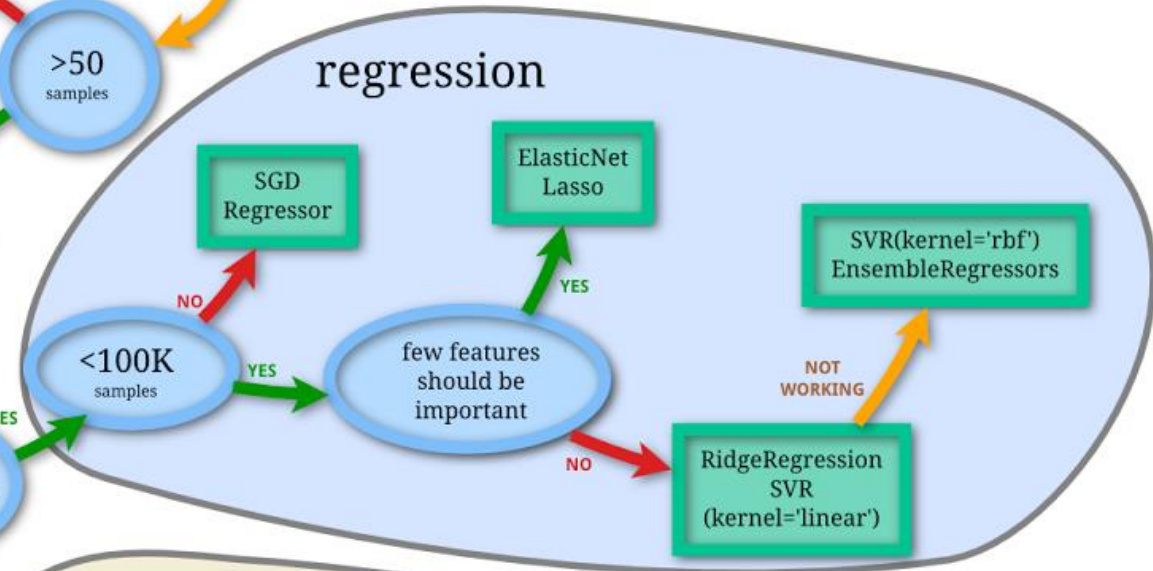
classification



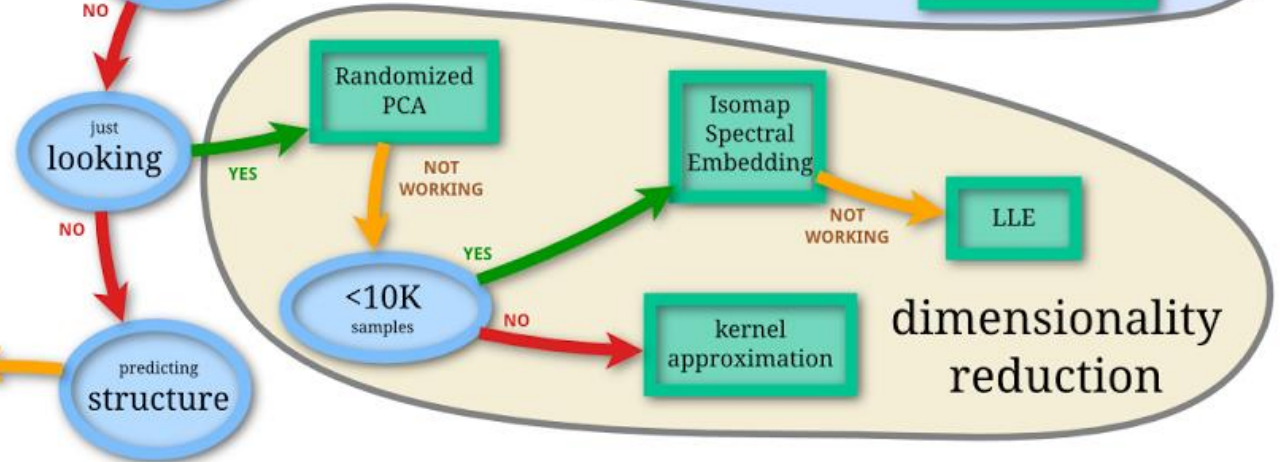
clustering



regression



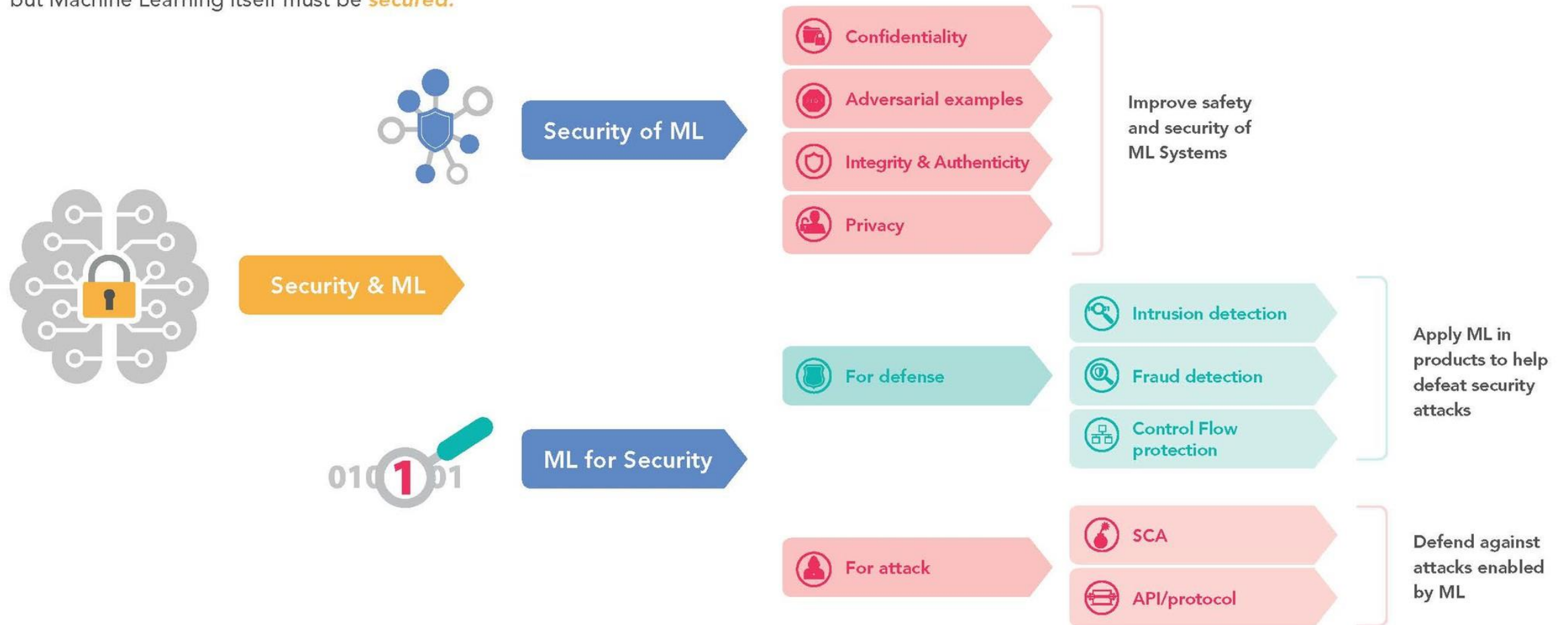
just looking



dimensionality reduction

Where Machine Learning, Security & Privacy Intersect

Machine Learning can **contribute** to IoT Security –
but Machine Learning itself must be **secured**.

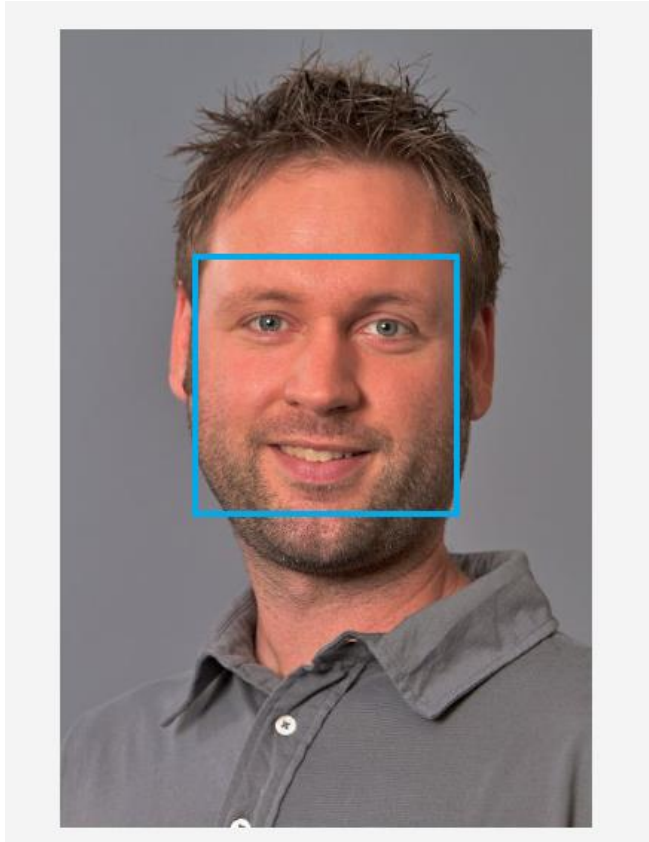


The background of the slide is a green-tinted image from the Matrix Revolutions movie poster, showing a line of Agents in suits and sunglasses. A large, semi-transparent white circle is positioned on the right side of the image, containing the title and source text.

Model Cloning

Image source: **Matrix Revolutions** movie poster

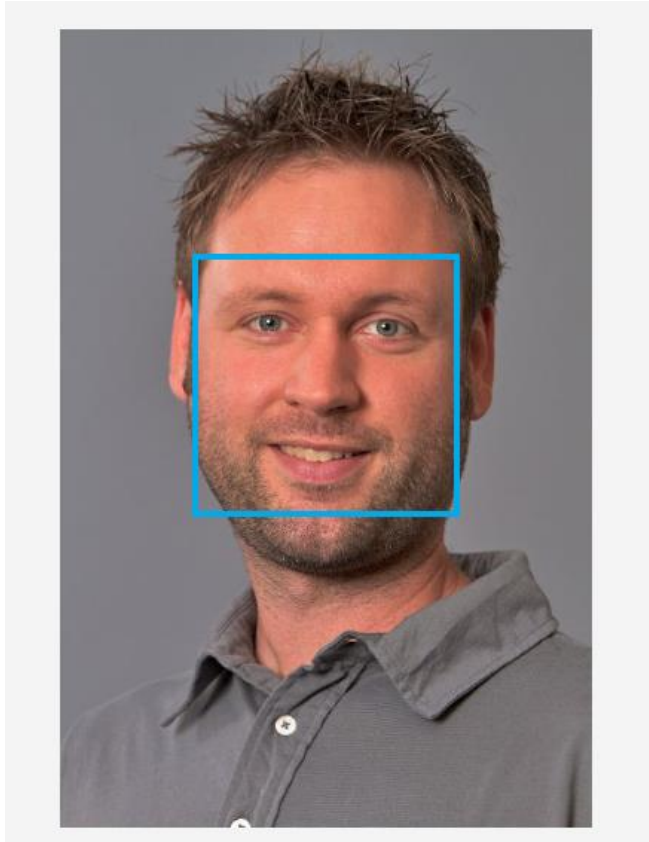
Example: Microsoft Azure Emotions Recognition



```
"scores": {  
  "anger": 2.03898679E-07,  
  "contempt": 0.0007247706,  
  "disgust": 6.056115E-07,  
  "fear": 1.0638247E-09,  
  "happiness": 0.9959635,  
  "neutral": 0.00329714641,  
  "sadness": 4.30003233E-08,  
  "surprise": 1.36911349E-05  
}
```

- <https://azure.microsoft.com/en-us/services/cognitive-services/emotion>

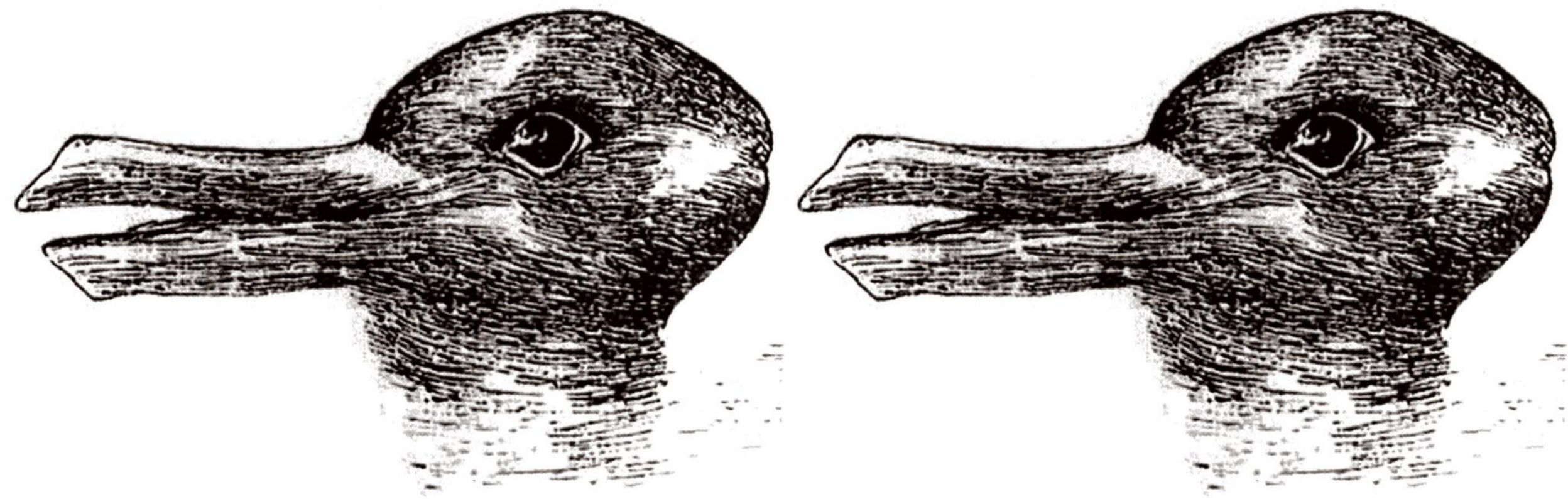
Example: Microsoft Azure Emotions Recognition



```
"scores": {  
  "anger": 2.03898679E-07,  
  "contempt": 0.0007247706,  
  "disgust": 6.056115E-07,  
  "fear": 1.0638247E-09,  
  "happiness": 0.9959635,  
  "neutral": 0.00329714641,  
  "sadness": 4.30003233E-08,  
  "surprise": 1.36911349E-05  
}
```

Clone model for < \$350 using random non-labeled data

- <https://azure.microsoft.com/en-us/services/cognitive-services/emotion>
- Tramèr, Zhang, Juels, Reiter, Ristenpart: *Stealing Machine Learning Models via Prediction APIs*. In *USENIX Security Symposium*, 2016.
- Correia-Silva, Berriel, Badue, de Souza, Oliveira-Santos. *Copycat CNN: Stealing Knowledge by Persuading Confession with Random Non-Labeled Data*. *arXiv preprint* (2018).

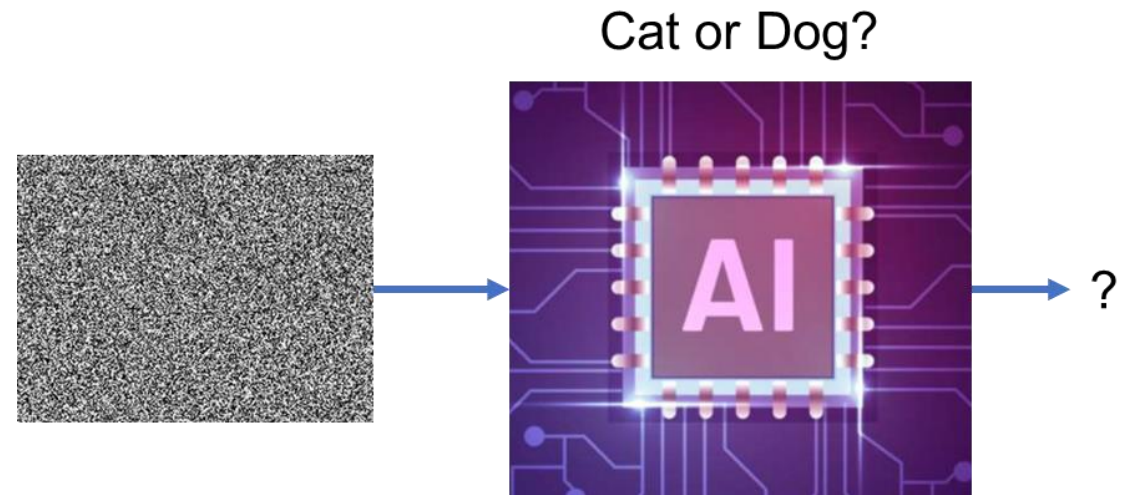


Adversarial Examples | Optical Illusions for Machines

Image by artist Joseph Jastrow, published in 1899 in Popular Science Monthly

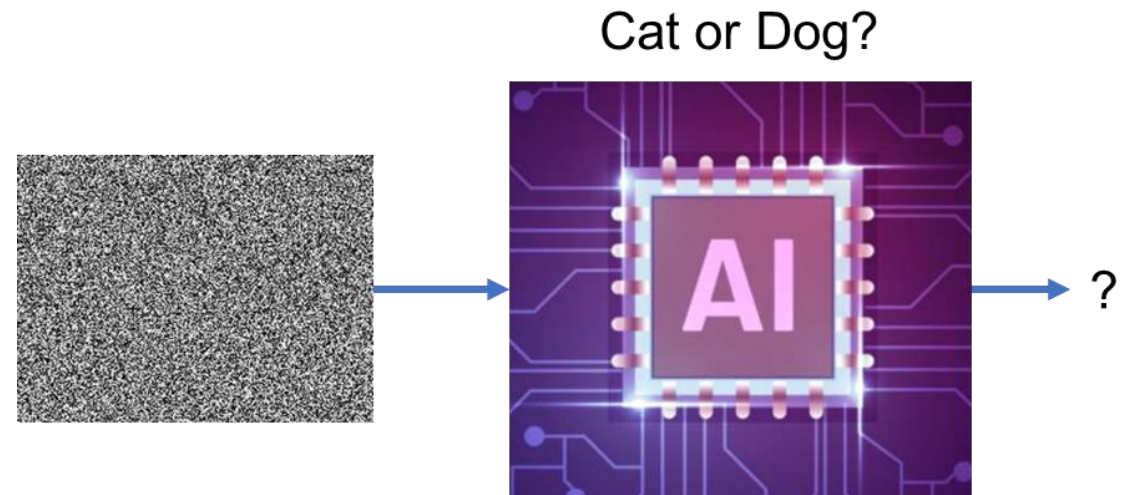
Misclassification versus Adversarial Examples

- Biggio, Corona, Maiorca, Nelson, Srndic, Laskov, Giacinto, Roli: *Evasion attacks against machine learning at test time*. In Machine Learning and Knowledge Discovery in Databases, 2013.
- Goodfellow, Shlens, Szegedy: *Explaining and harnessing adversarial examples*. In arXiv preprint 2014
- Szegedy, Vanhoucke, Ioffe, Shlens, Wojna: *Rethinking the inception architecture for computer vision*. In IEEE conference on computer vision and pattern recognition, 2016.



Misclassification versus Adversarial Examples

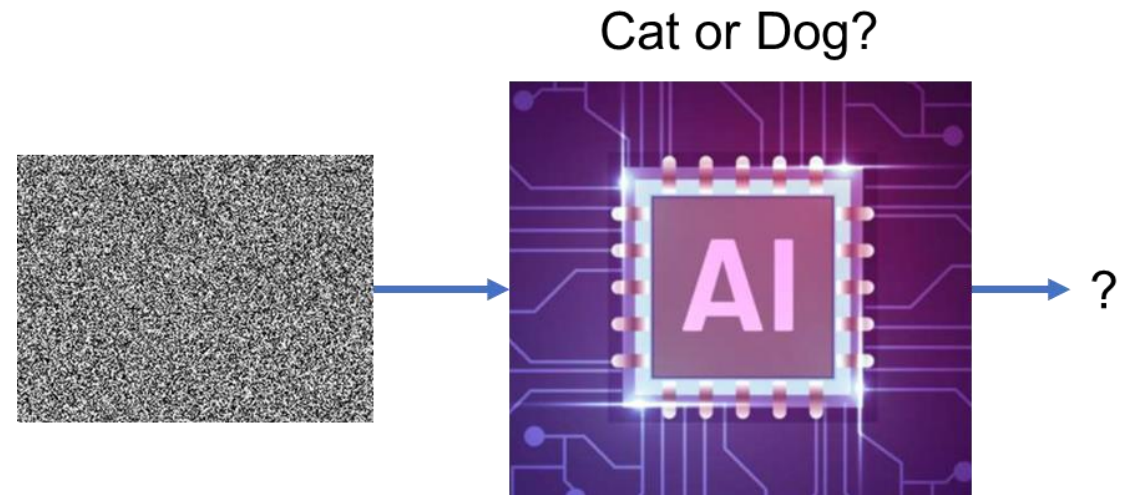
- Biggio, Corona, Maiorca, Nelson, Srndic, Laskov, Giacinto, Roli: *Evasion attacks against machine learning at test time*. In Machine Learning and Knowledge Discovery in Databases, 2013.
- Goodfellow, Shlens, Szegedy: *Explaining and harnessing adversarial examples*. In arXiv preprint 2014
- Szegedy, Vanhoucke, Ioffe, Shlens, Wojna: *Rethinking the inception architecture for computer vision*. In IEEE conference on computer vision and pattern recognition, 2016.



~ 0.832 flowerpot

Misclassification versus Adversarial Examples

- Biggio, Corona, Maiorca, Nelson, Srndic, Laskov, Giacinto, Roli: *Evasion attacks against machine learning at test time*. In Machine Learning and Knowledge Discovery in Databases, 2013.
- Goodfellow, Shlens, Szegedy: *Explaining and harnessing adversarial examples*. In arXiv preprint 2014
- Szegedy, Vanhoucke, Ioffe, Shlens, Wojna: *Rethinking the inception architecture for computer vision*. In IEEE conference on computer vision and pattern recognition, 2016.

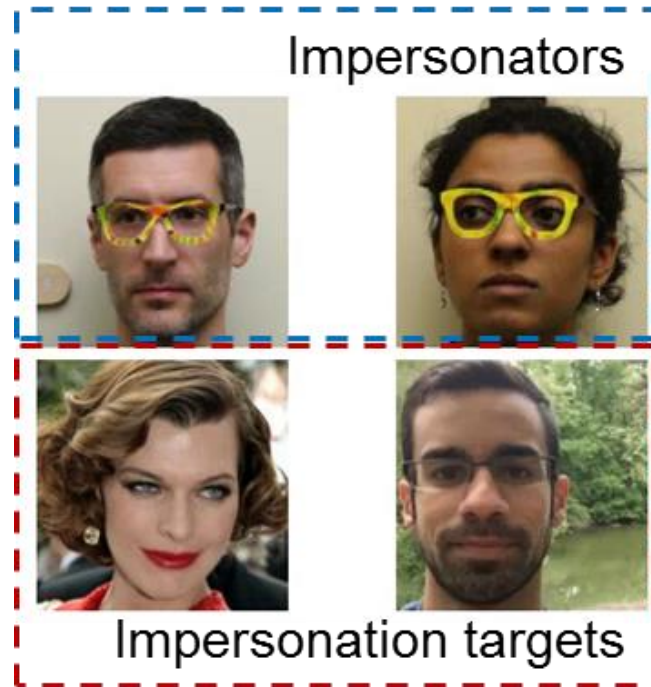


~ 0.832 flowerpot



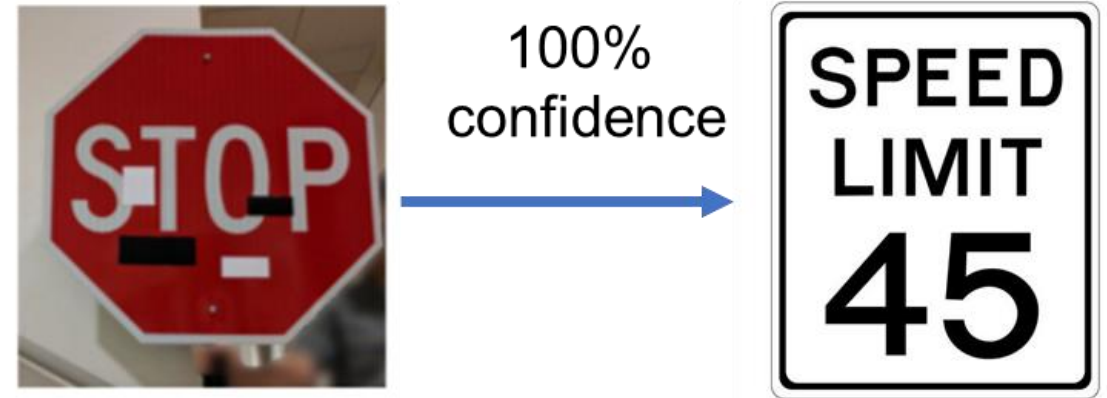
~ 0.999 warplane

Security



Sharif, Bhagavatula, Bauer, Reiter: *Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition*. In *ACM SIGSAC 2016*

Safety



Eykholt, Evtimov, Fernandes, Li, Rahmati, Xiao, Prakash, Kohno, Song: *Robust Physical-World Attacks on Deep Learning Visual Classification*. In *IEEE Computer Vision and Pattern Recognition 2018*.

Impact in Practice

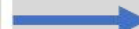


cleverhans

<https://github.com/tensorflow/cleverhans>

Papernot et al.: *Technical Report on the
CleverHans v2.1.0 Adversarial Examples
Library*, arXiv preprint 2018

Countermeasures? Adversarial Training



Adversarial Training

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



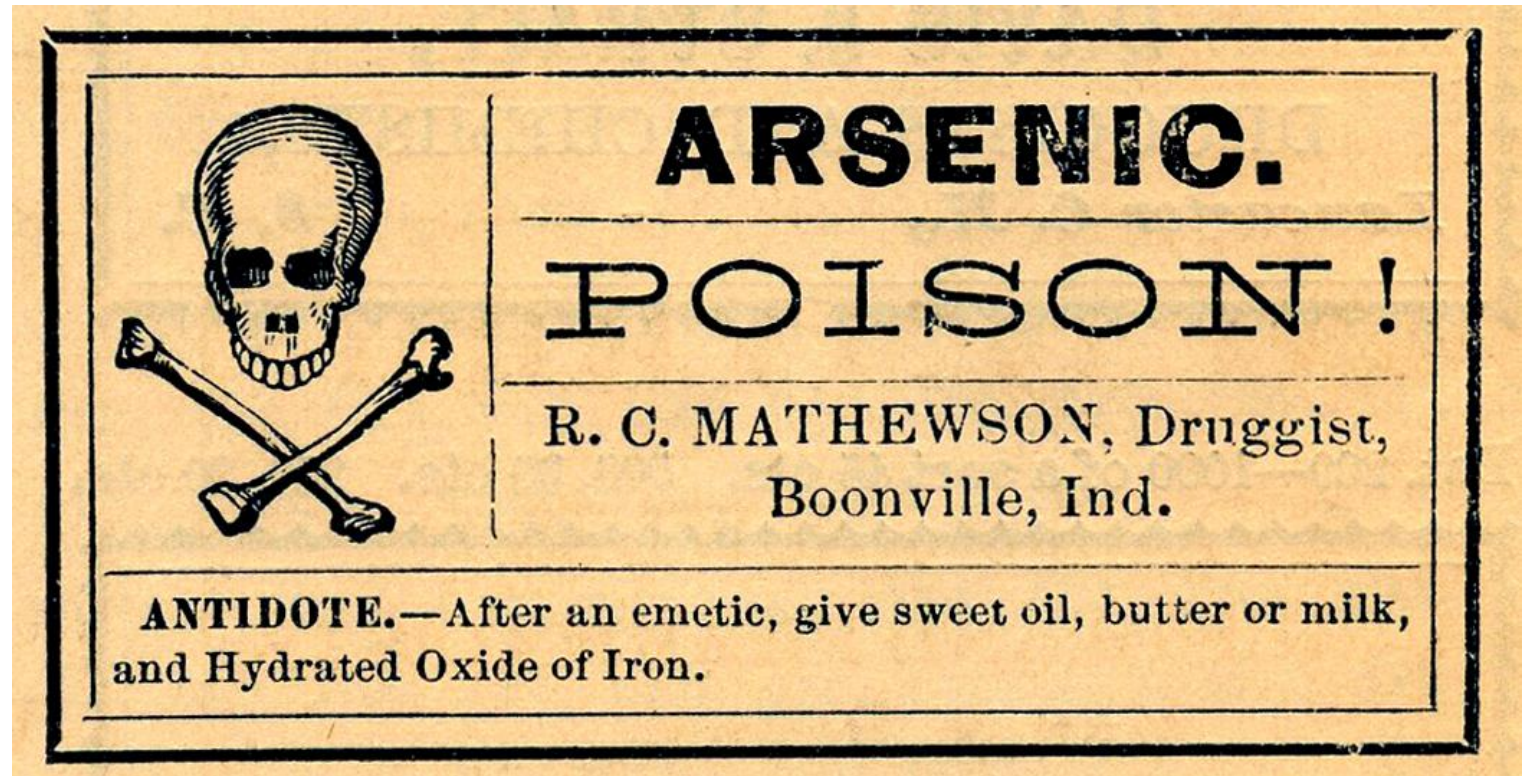
truck



<https://github.com/tensorflow/cleverhans>

CIFAR-10 Model	Accuracy of the model	Adversarial examples that mislead the model
Original	87%	90%
Trained with adversarial examples	86%	17%

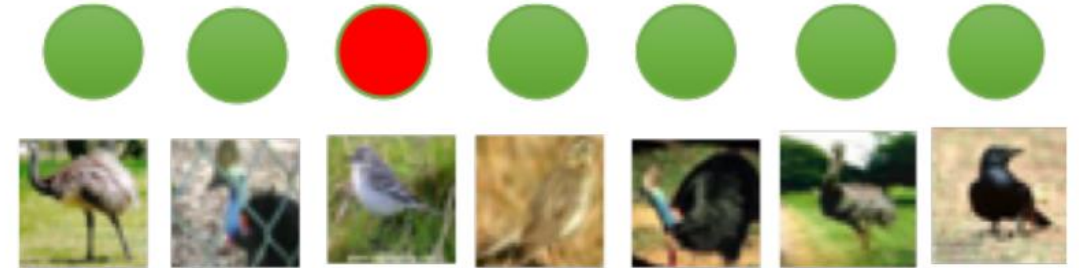
Data Poisoning



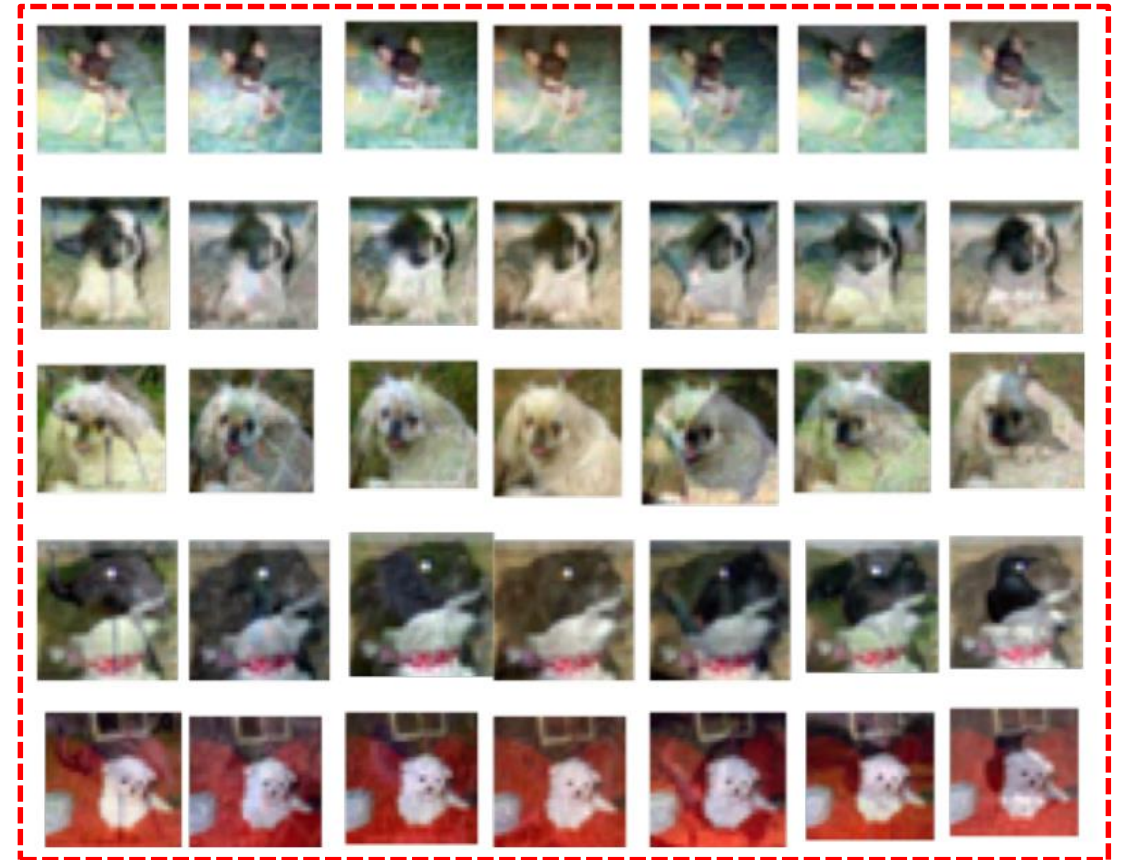
Adversarial Training - Revisited

Shafahi, Huang, Najibi, Suciu, Studer,
Dumitras, Goldstein: *Poison Frogs! Targeted
Clean-Label Poisoning Attacks on Neural
Networks*. arXiv preprint 2018.

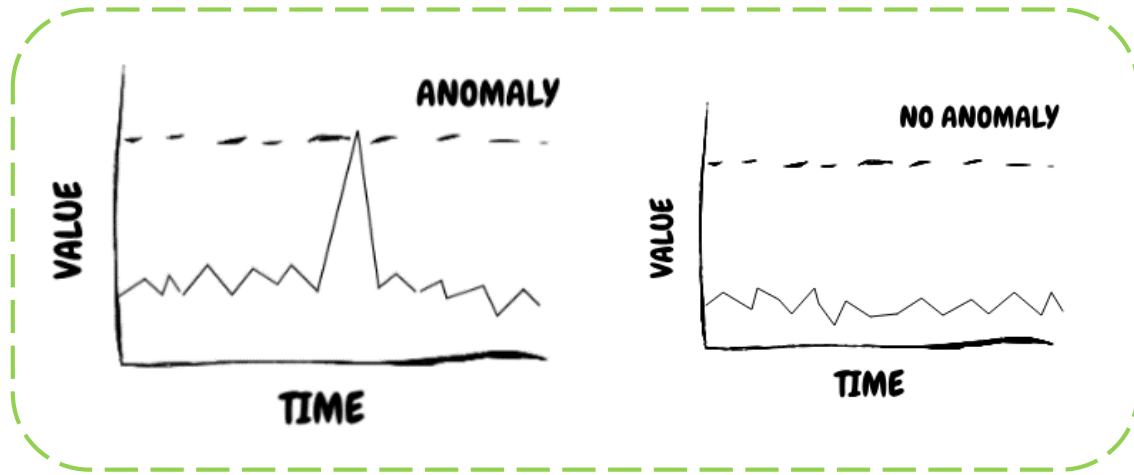
Target



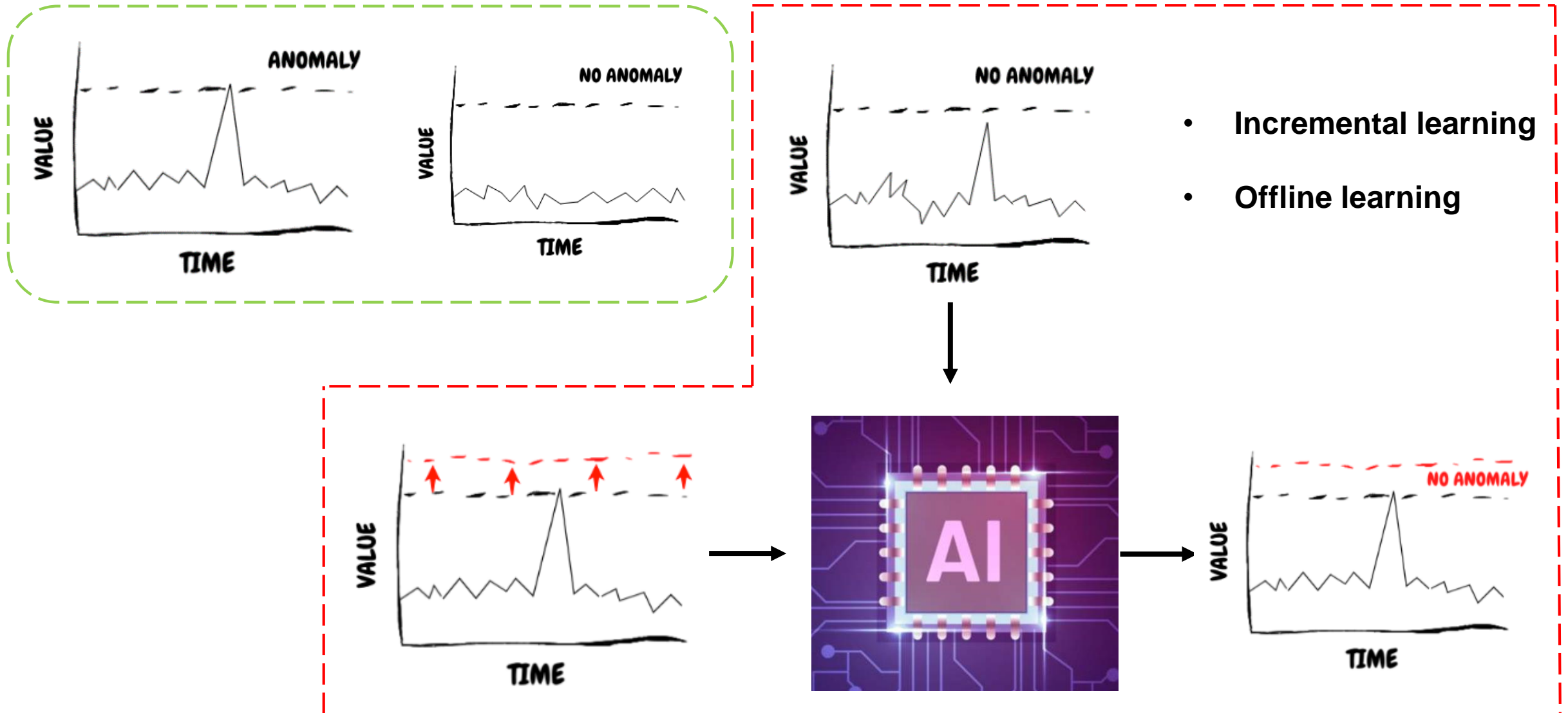
Poison



Anomaly detection



Anomaly detection



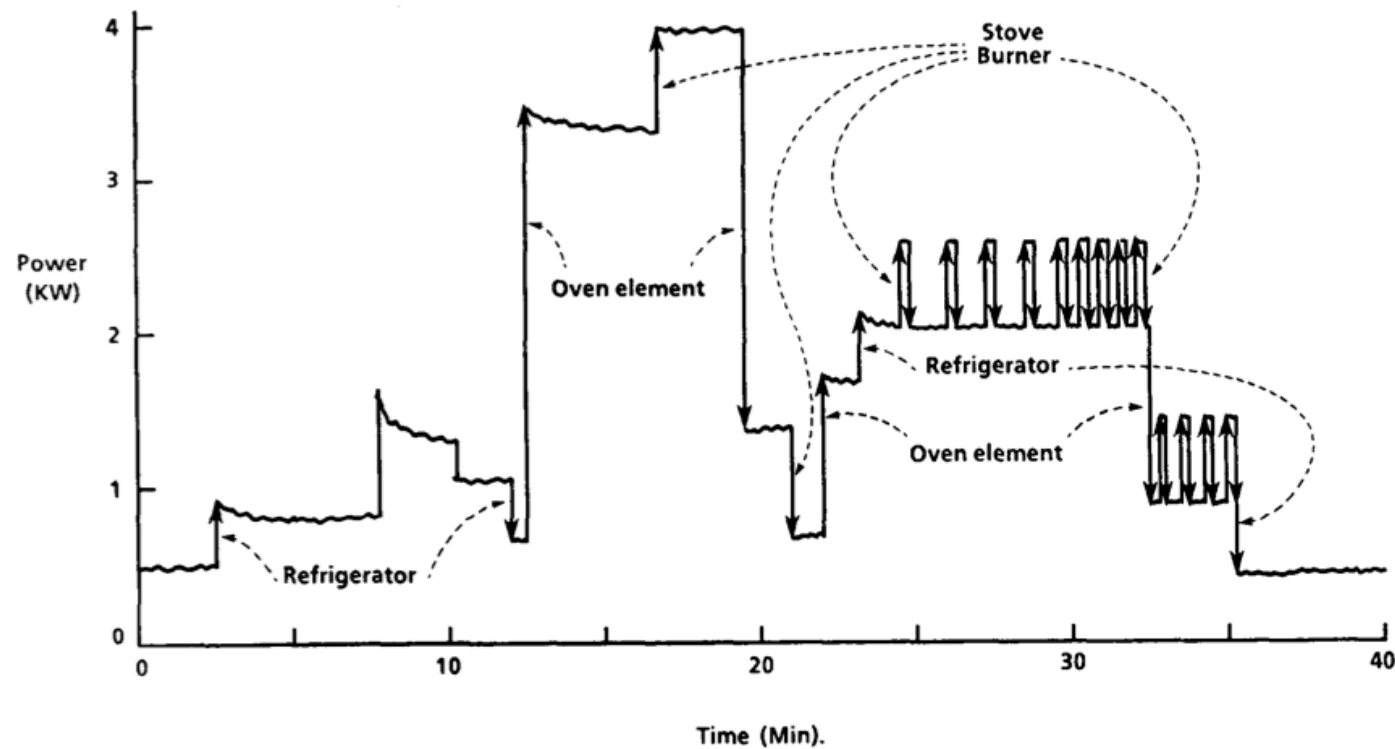


Privacy - Use case: Smart Grid

Forecasting power consumption

- **Suppliers** need forecast to buy energy generation contracts that cover their clients
- **Distribution supply operators** require longer term forecasts to ensure the necessary network capacity is available
- Forecasting could allow for dynamic price determination

Forecast



Hart: *Nonintrusive appliance load monitoring*. Proceedings of the IEEE 1992

Privacy Concerns in the Smart-Grid

Energy consumption reveals

Patterns

- Another microwave meal?

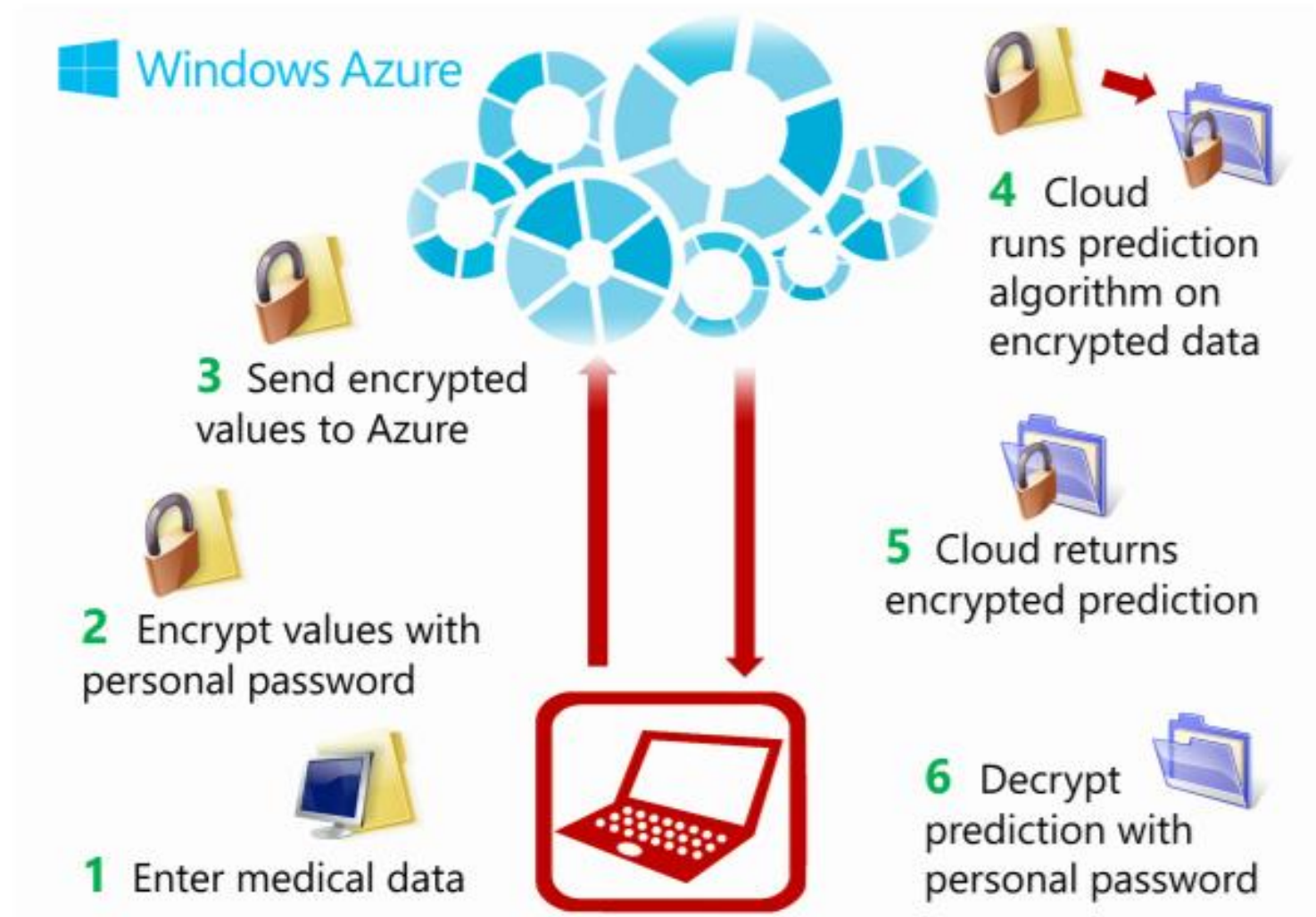
Invalid usage

- Insurance or warranty

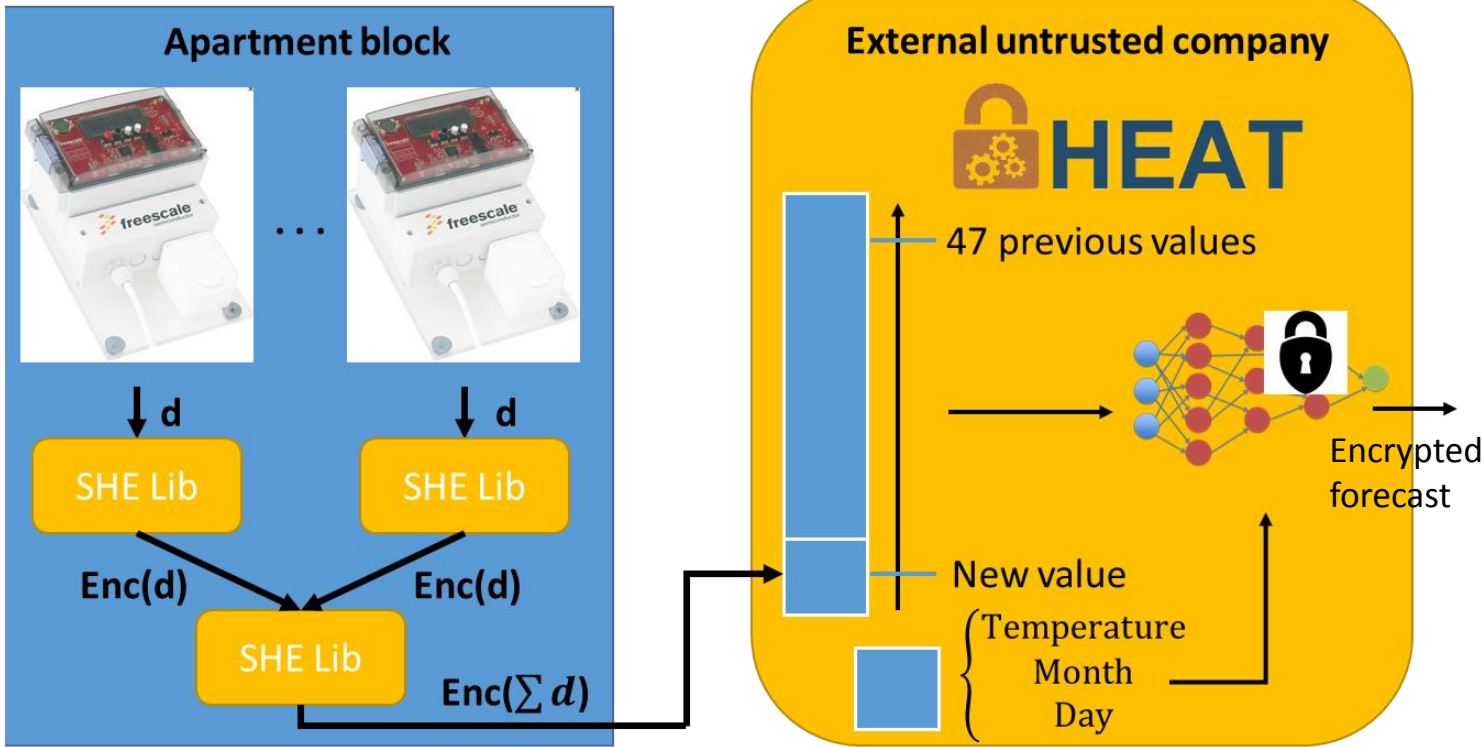
Real-time information

- Number of people in a household
- Are you on holidays?

Computing on Encrypted Data



Bos, Lauter, Naehrig: *Private Predictive Analysis on Encrypted Medical Data*. Journal of Biomedical Informatics, 2014.



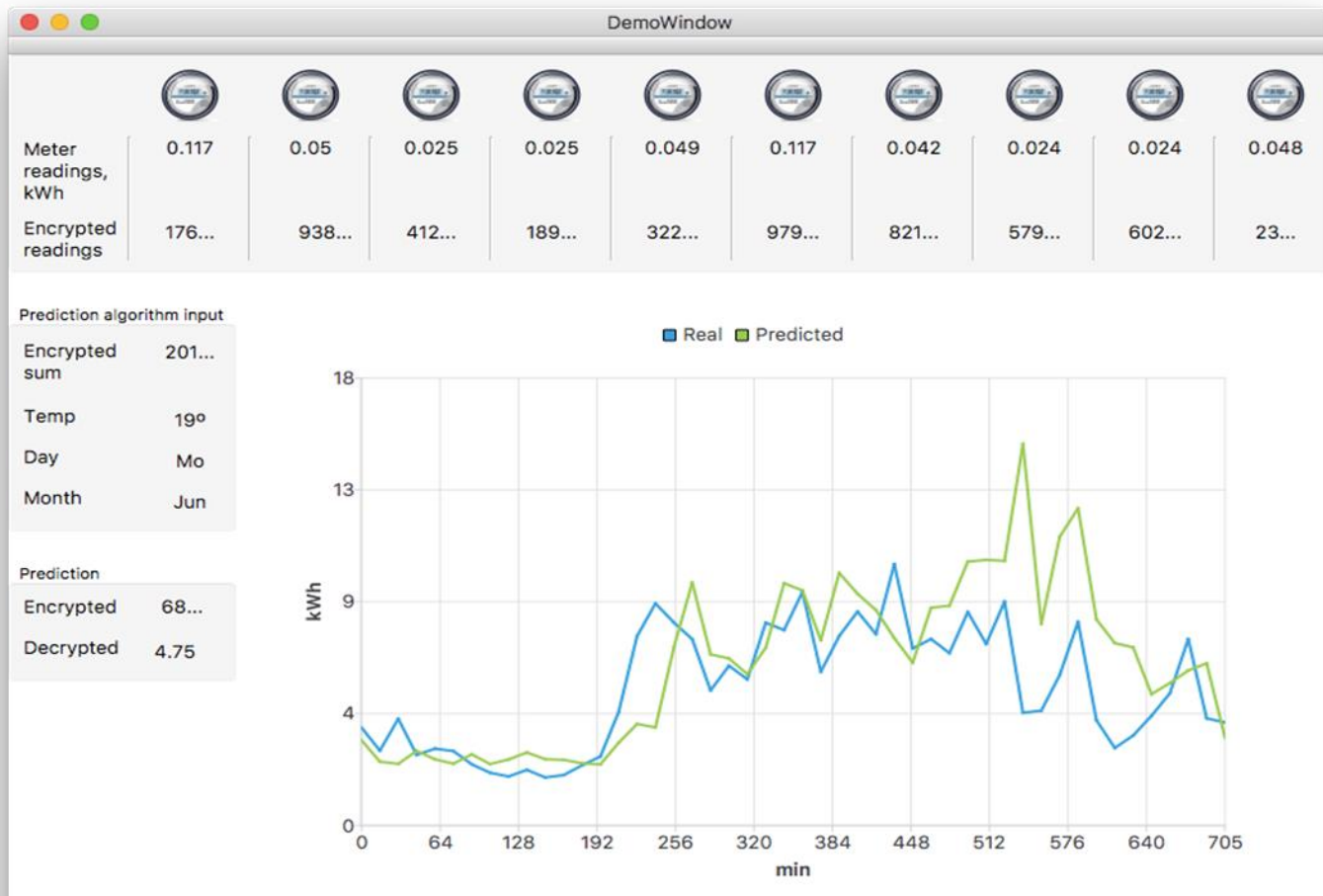
- Forecast power consumption for next half hour in ≈ 2.5 seconds to evaluate
- Neural network
Inputs: 51
Hidden layers: 3 ($8 \rightarrow 4 \rightarrow 2$)
Output: 1



Funded by the Horizon 2020
Framework Programme of the
European Union

Machine Learning using Encrypted Data

- Bonte, Bootland, Bos, Castryck, Iliashenko, Vercauteren: *Faster Homomorphic Function Evaluation using Non-Integral Base Encoding*. Cryptographic Hardware and Embedded Systems – CHES 2017
- Bos, Castryck, Iliashenko, Vercauteren: *Privacy-friendly Forecasting for the Smart Grid using Homomorphic Encryption and the Group Method of Data Handling*. AFRICACRYPT 2017



Machine Learning using Encrypted Data

- Forecast power consumption for next half hour in ≈ 2.5 seconds to evaluate
- Neural network
Inputs: 51
Hidden layers: 3 (8 \rightarrow 4 \rightarrow 2)
Output: 1



Funded by the Horizon 2020
Framework Programme of the
European Union

- Bonte, Bootland, Bos, Castryck, Iliashenko, Vercauteren: *Faster Homomorphic Function Evaluation using Non-Integral Base Encoding*. Cryptographic Hardware and Embedded Systems – CHES 2017
- Bos, Castryck, Iliashenko, Vercauteren: *Privacy-friendly Forecasting for the Smart Grid using Homomorphic Encryption and the Group Method of Data Handling*. AFRICACRYPT 2017



Machine Learning using Encrypted Data

- Forecast power consumption for next half hour in ≈ 2.5 seconds to evaluate
- Neural network
Inputs: 51
Hidden layers: 3 ($8 \rightarrow 4 \rightarrow 2$)
Output: 1



Funded by the Horizon 2020 Framework Programme of the European Union

- Bonte, Bootland, Bos, Castryck, Iliashenko, Vercauteren: *Faster Homomorphic Function Evaluation using Non-Integral Base Encoding*. Cryptographic Hardware and Embedded Systems – CHES 2017
- Bos, Castryck, Iliashenko, Vercauteren: *Privacy-friendly Forecasting for the Smart Grid using Homomorphic Encryption and the Group Method of Data Handling*. AFRICACRYPT 2017



Conclusions

- Machine learning can improve quality of life due to availability of huge amount of data
- Security is one of the biggest challenges in large scale deployment of machine learning
- A lot of open security & privacy challenges
- [+] Cryptography to the rescue for some problems
- [-] Expect zero-day attacks against machine learning models



SECURE CONNECTIONS
FOR A SMARTER WORLD

& MACHINE LEARNING

