

RSA® Conference 2019

San Francisco | March 4–8 | Moscone Center



BETTER.

SESSION ID: IDY-R11

Adversarial Machine Learning against Modern Behavioral Biometrics Systems

Heng Tang

@9he_ta0ng

Ajit Gaddam

@AjitGaddam

#RSAC



“The Password has become kind of a nightmare!”

Prof. Fernando J. Corbato
Creator of password back in 1961 at MIT for the Compatible Time-Sharing System

Image source: wired.com

Data breaches of passwords & by passwords

5M

Gmail logins were leaked on September 9th, 2014

Source: time.com. Sep 2014

167M

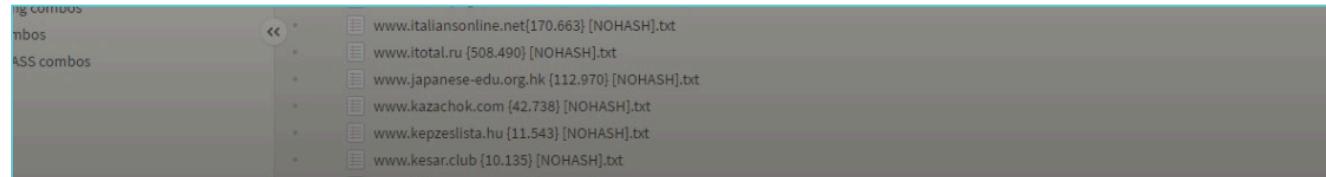
Hacked LinkedIn accounts are on sale

Source: csoonline.com, May 2016

2.3B

credentials stolen from 51 organizations in 2017

Shape Credential Spill Report, 2018



The 773 Million Record "Collection #1" Data Breach

eBay Inc. eBay Inc. @ebayinc Following

eBay asks all users to change passwords due to cyberattack that compromised non-financial info in a database: ebayinc.com/in_the_news/st... 3:20 PM - 21 May 2014

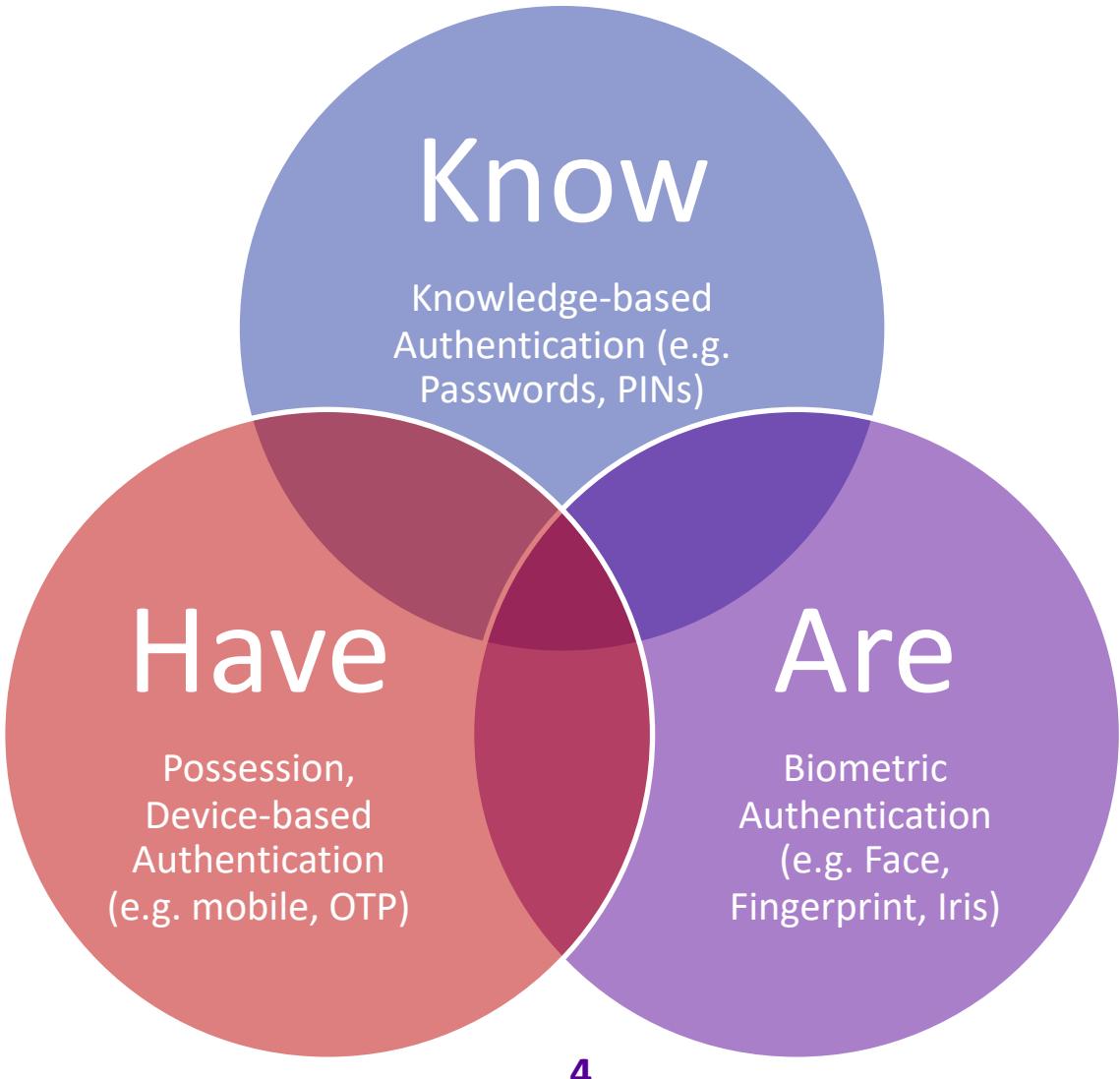
Reply Retweet Favorite More

RETWEETS 208 FAVORITES 21

“80% of all data breaches resulted from weak or stolen passwords.”

Verizon Data Breach report 2017

Multi-Factor Authentication



RSA®Conference2019

Biometric Authentication

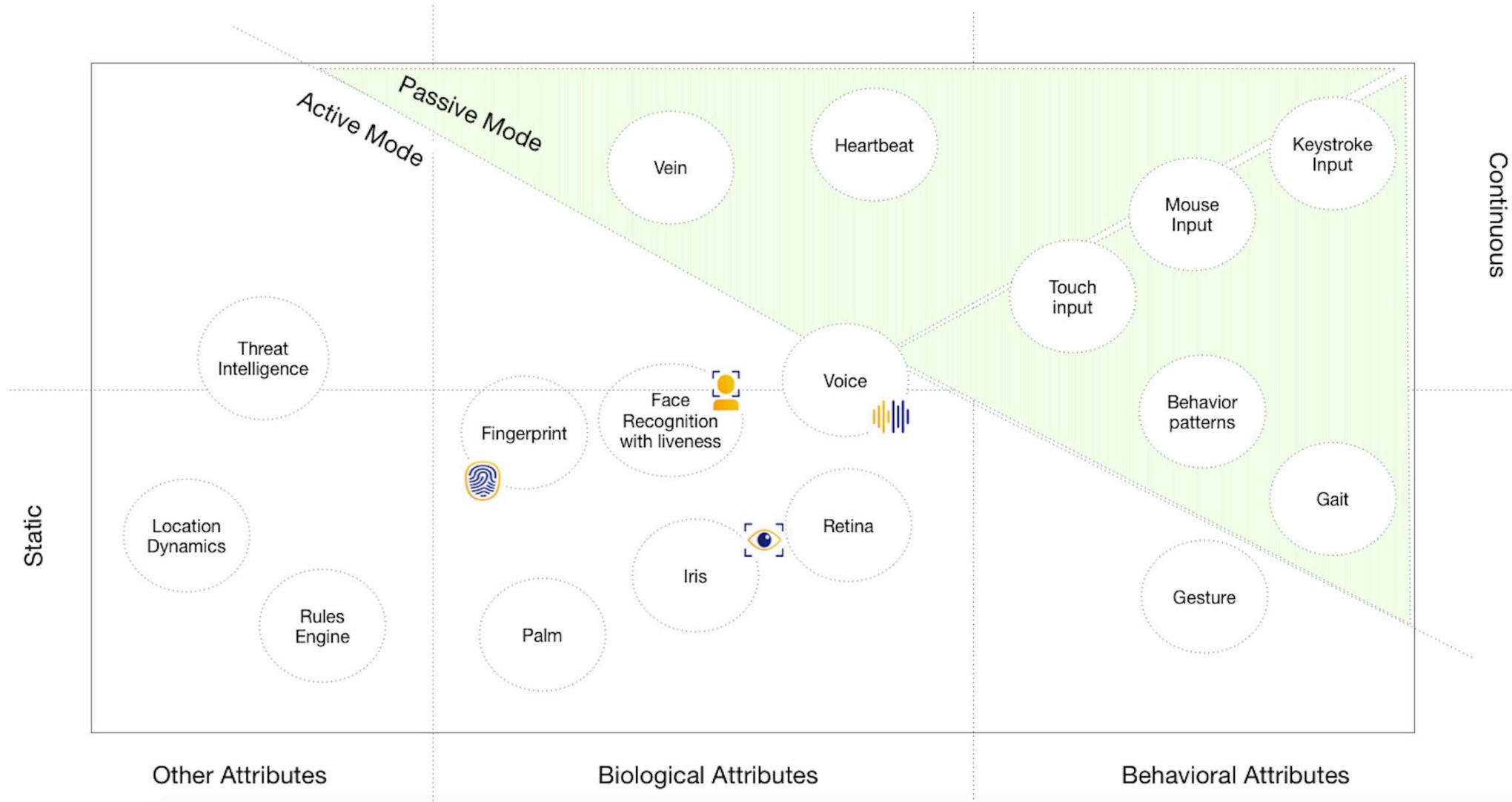
Biological & Behavioral Biometrics



Biometric Authentication



Spectrum of Biometric Authentication



Biological vs. Behavioral Biometrics

Biological

Static

More
“Human Perceivable”

Persistent

Behavioral

Dynamic

Less intuitive to
understand raw data

Fluctuated

Behavioral Biometrics – Keystroke Dynamics

Key pairs generated from typing “secure”

re	90
se	110
ur	185
ec	217
cu	232

re	124
se	98
ur	200
ec	195
cu	122

re	90
se	110
ur	185
ec	217
cu	232

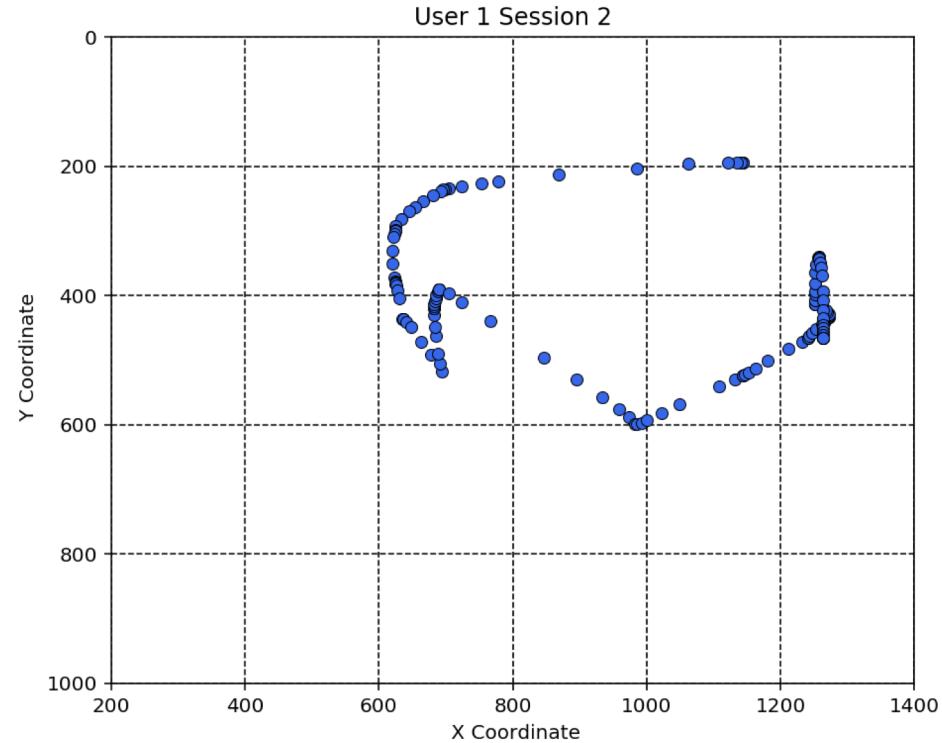
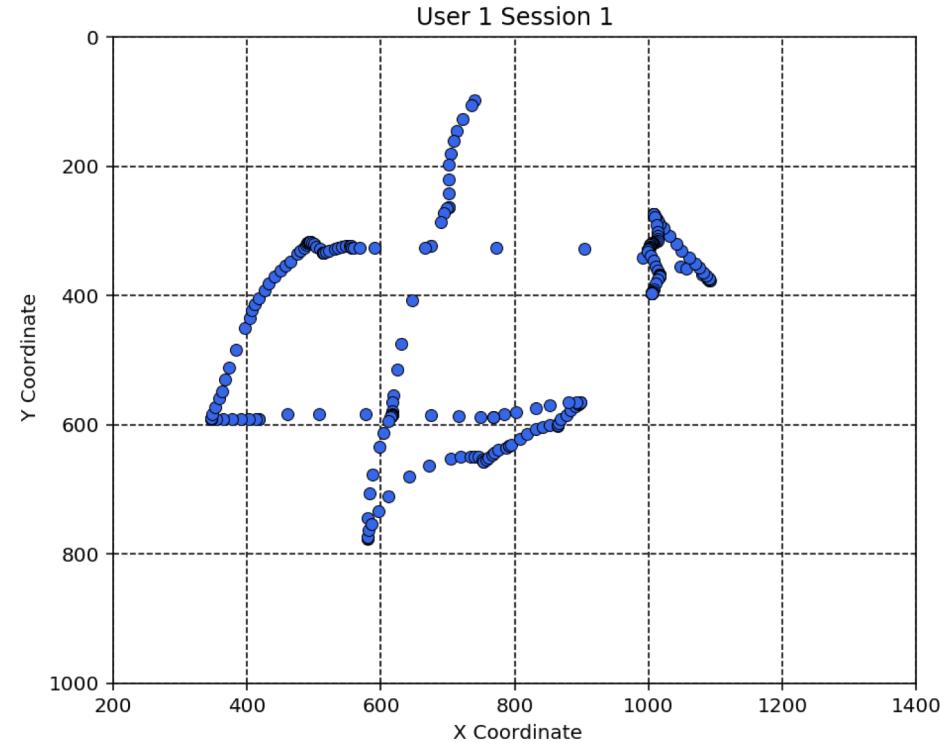
se	98
cu	122
re	124
ec	195
ur	200

$$AS = \sum_{ng}^M \frac{1}{N} AS_{ng} ; AS_{ng} = 1 - \frac{\sum_i sim(ng_i)}{\sum_{ng_i} 1}$$

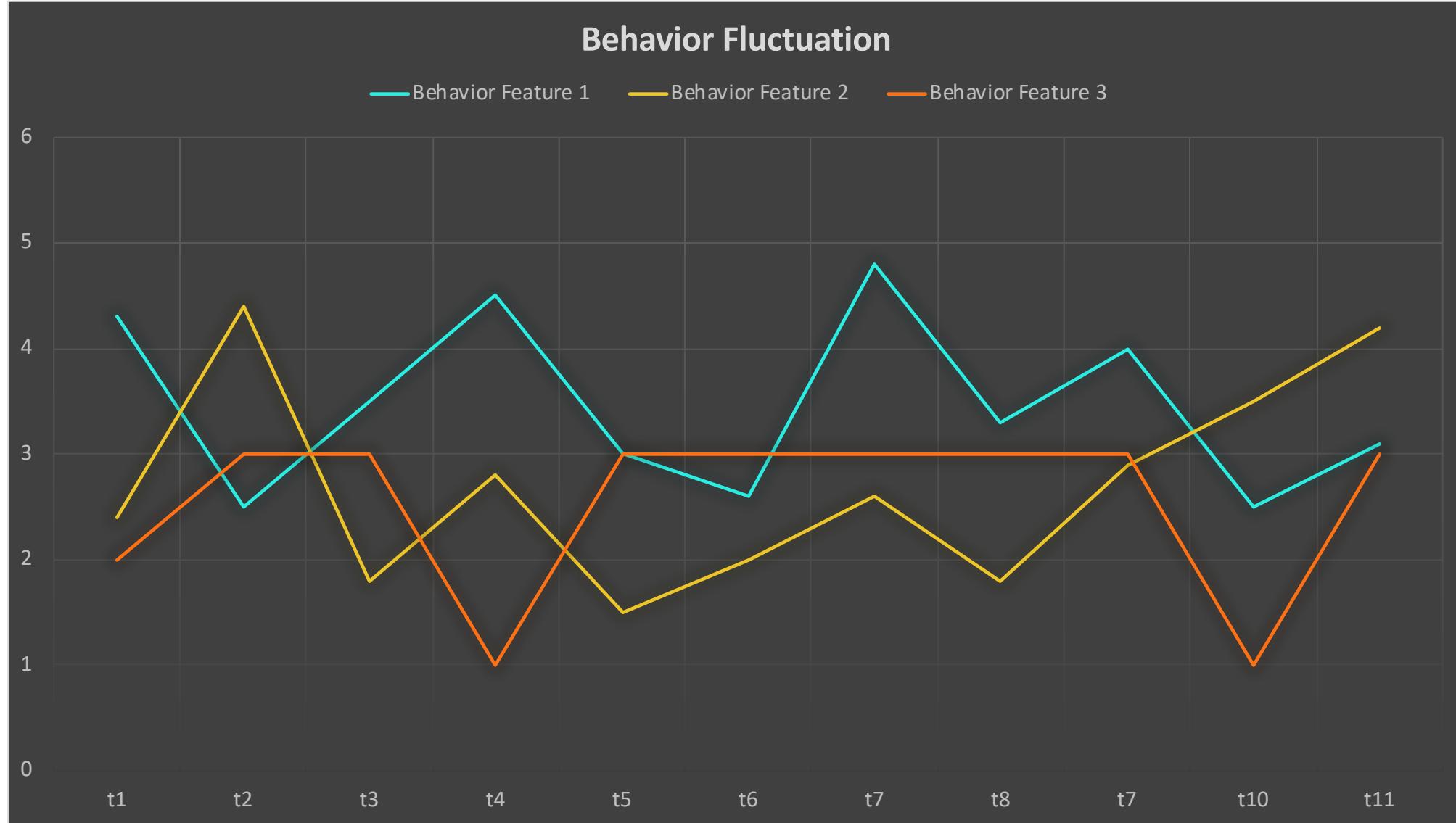
$$sim(ng_i) = \begin{cases} 1, & \text{if } 1 < \frac{\max(du_{ji})}{\min(du_{ji})} < t \\ 0, & \text{otherwise} \end{cases}$$

$$RS = \sum_{ng}^M \frac{1}{N} RS_{ng} \quad RS_{ng} = \frac{\sum_i d_i}{RS_{ng-max}}$$

Behavioral Biometrics – Mouse Dynamics



Behavioral Biometrics Features

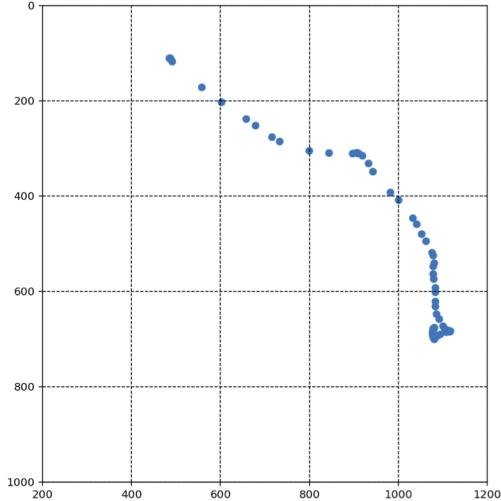


RSA®Conference2019

Adversarial Machine Learning

Introduction and methodologies

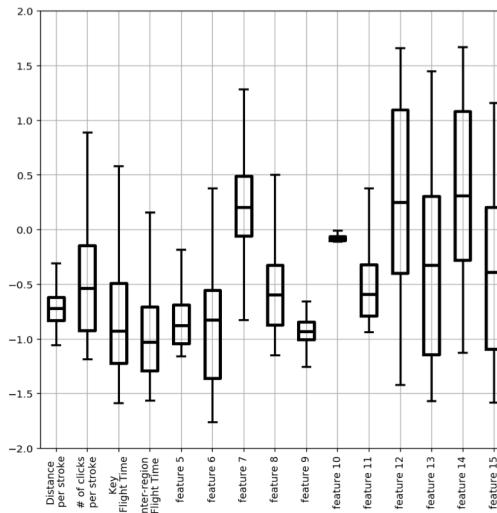
How ML/DL Model Works



Raw Data

- Images
- Measurements
- Time-series Data

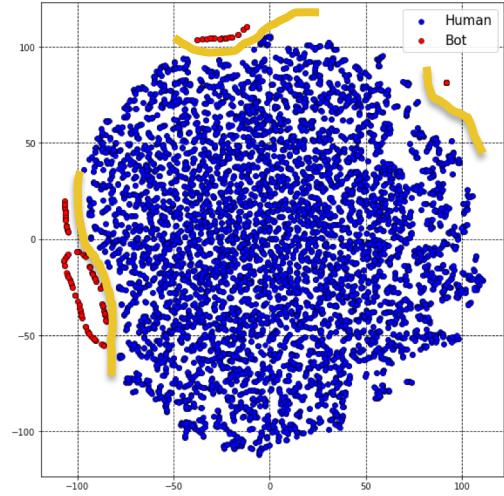
Preprocessing
/Data
Cleaning



Features

- Feature Engineering
- Implicit Features
- Embedding

Feature
Analysis



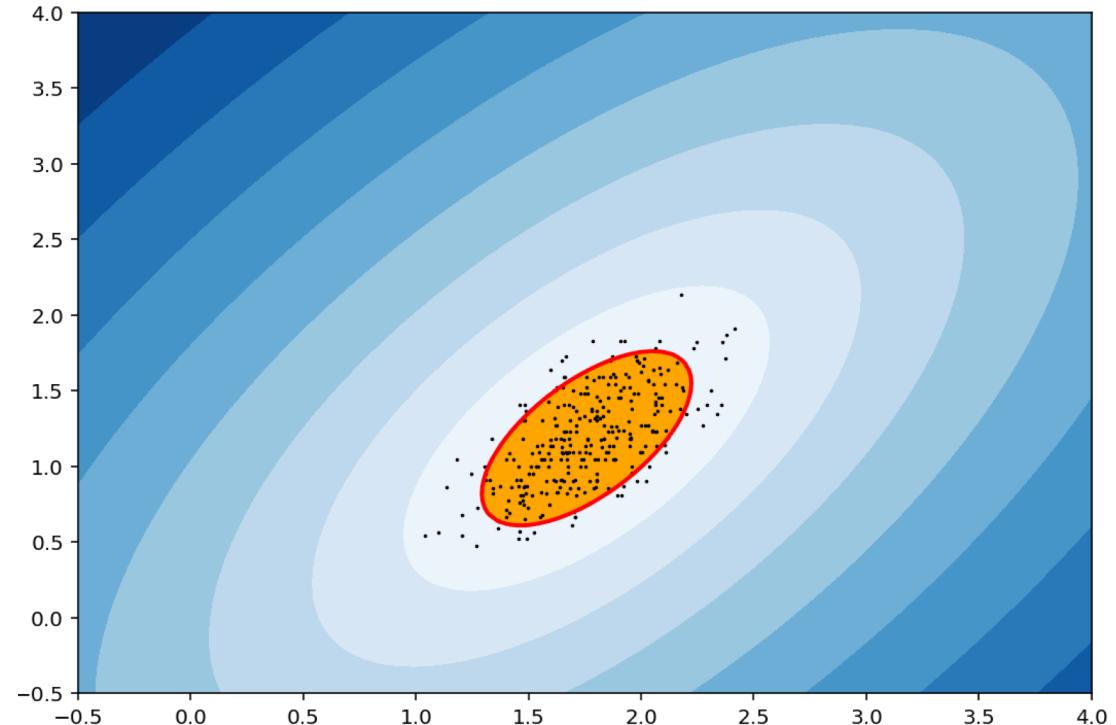
Model

- Statistic Models
- Deep Neural Networks
- Non-linear Models

Adversarial Machine Learning

Attacker

- Trying to explore the decision boundary of original ML/DL based models
- Generating synthetic samples that can fit into the decision boundary to be classified with expected labels



Adversarial Machine Learning – Definition

x : Original Inputs

y : True Labels

y' : Targeted Malicious Labels

$C(x)$: Trained Classifiers – s.t. $C(x) = y$

$p(x), p(x, y')$: Perturbation functions

Targeted Attack

$$\arg \min_{p(x,y')} C(x + p(x, y')) = y'$$

Crafted Inputs

Untargeted Attack

$$\arg \min_{p(x)} C(x + p(x)) \neq y$$

Crafted Inputs

Adversarial Machine Learning – Image Recognition



+



=

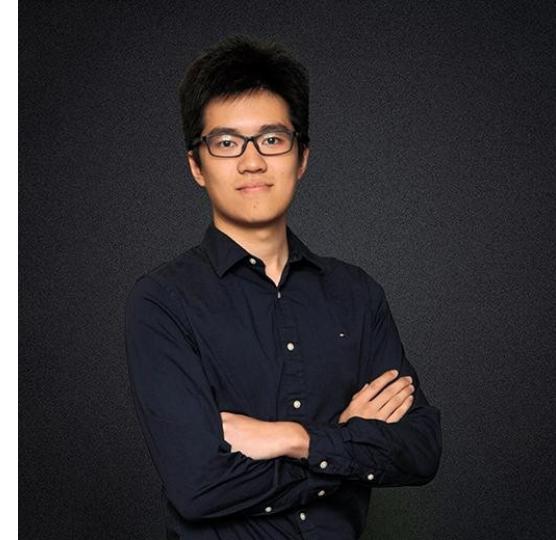
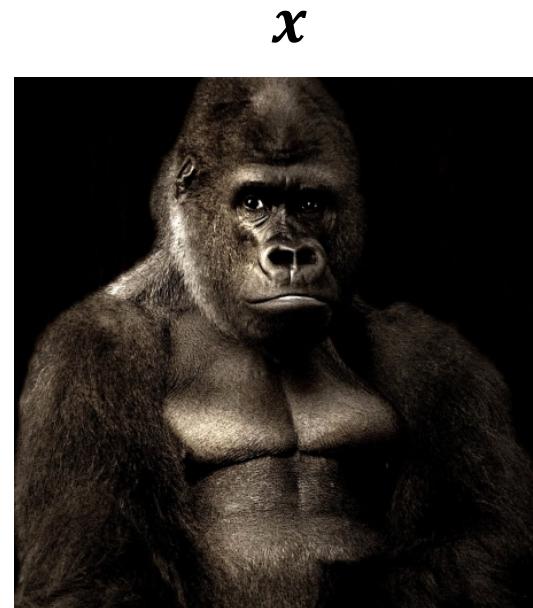


Image source: pixnio.com

Imperceptible Difference

- $\arg \min_{p(x,y')} C(x + p(x, y')) = y'$
- Metrics to reflect the imperceptibility of the change to observer
 - L^0 distance
 - L^1 distance
 - L^2 distance
 - L^∞ distance – Max Norm

Adversarial Machine Learning – Image Recognition



Prediction "Gorilla" with High Confidence

+



$$\varepsilon \cdot \text{sign}(\nabla_x L(\theta, x, y'))$$

FGSM[Goodfellow et al., 2014]

$$\varepsilon = L^\infty(p(x, y'))$$

Weights

$$x + p(x, y')$$



Prediction "Human" with High Confidence

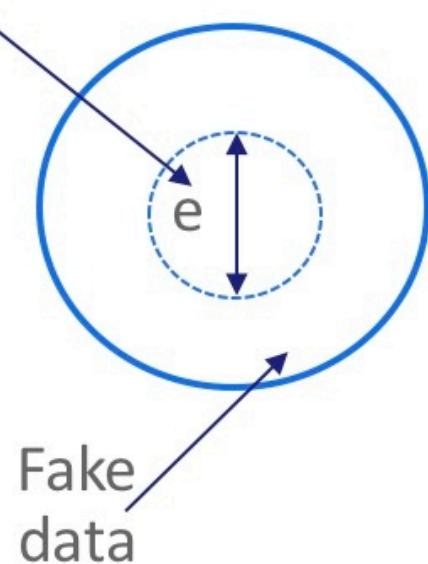
Adversarial Machine Learning – Behavioral Biometrics

Targeted Attack

$$\arg \max_{p(x,y)} C(x + p(x,y)) = y$$

Crafted Inputs

mimicked
data



Impact on Behavioral Biometrics Systems

- Adaptive Systems vs. Static Systems
- Model Poisoning
 - Usability
 - Discredit the underlying system
- Attacks may not always target tweaking ML/DL models
 - Tweak vs. Circumvent
 - Carefully crafted attack vs. random guess

RSA® Conference 2019

Defense

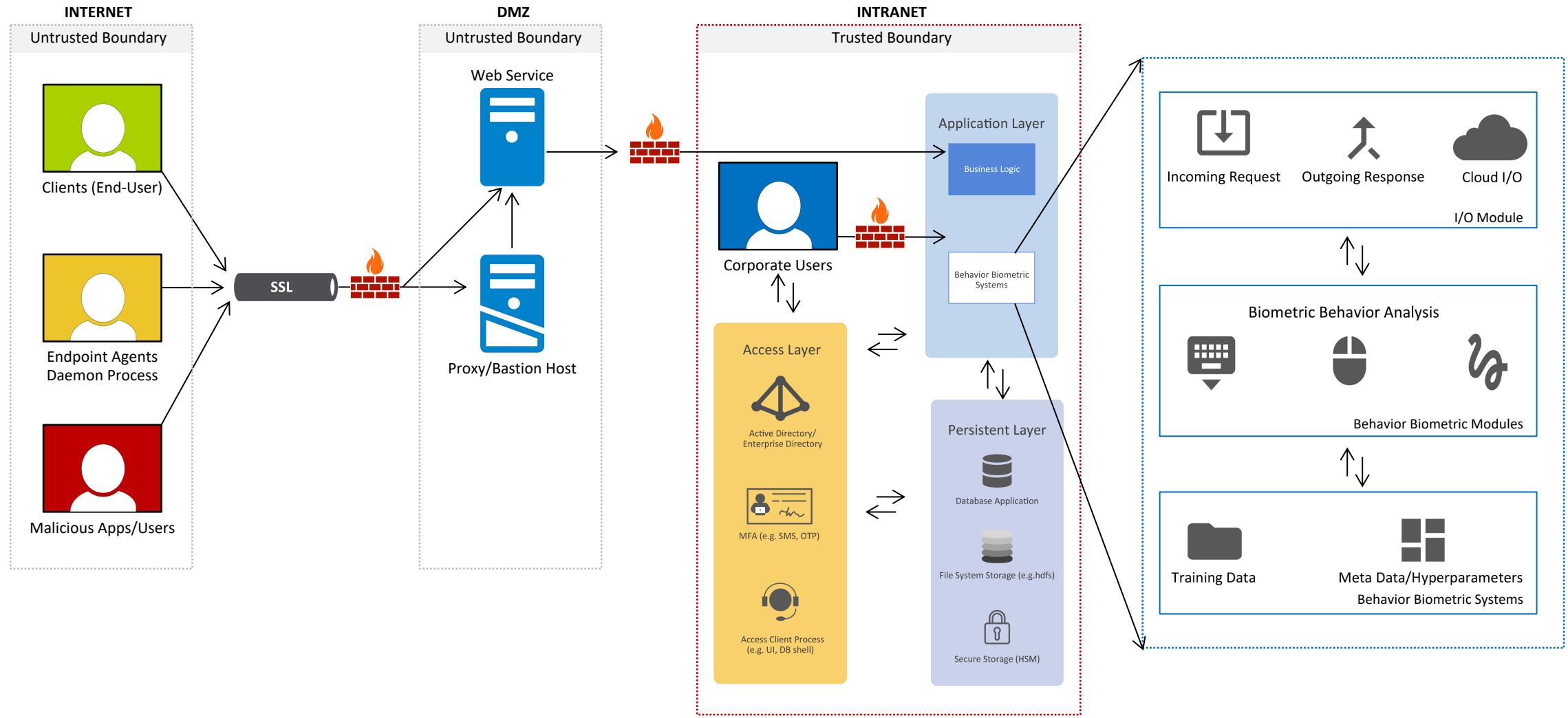
Threat Modeling and Countermeasures



Different Levels of Knowledge

- **White-box** – Has full knowledge and potentially full access to any component (e.g. Inside attack with domain knowledge)
- **Black-box** – Has limited knowledge and access on targeted system/model
- **Grey-box** – Something in between, more realistic assumption

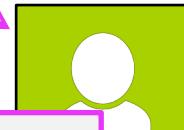
Threat Model



Threat Model

Attack 1: Session Hijacking

- MTM attack
- Sniffing

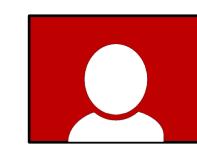


Attack 2: Reverse Engineering

- Crafted payload



Endpoint Agents
Daemon Process



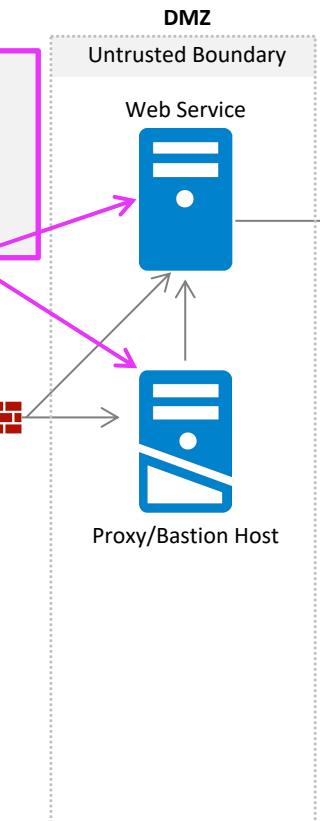
Malicious Apps/Users

Attack 3: Malicious External Attackers

- Randomly crafted payload
- Impersonation
- Compromised account
- DDos & Botnet

Attack 4: Compromised Machine

- Exploit



Attack 5: Insider Attack

- Abusing access

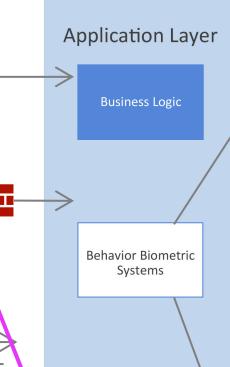


Corporate Users



INTRANET

Trusted Boundary



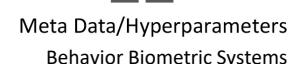
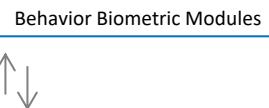
Attack 7: Data Leakage

- Unencrypted communication
- Secret data



I/O Module

Biometric Behavior Analysis



Attack 8: Manipulated Data

- Data poisoning
- Data corruption

Attack Category Groupings

Client side

- Hijacked sessions
- Reverse engineering endpoint agents
- Malicious crafted payload

Server side

- Compromised machines
- Manipulated data & tweaked models
- Insider attack

Client Side Attacks

Attack 1: Session Hijacking

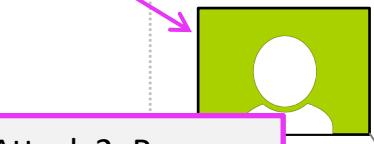
- MTM attack
- Sniffing

Attack 2: Reverse Engineering

- Crafted payload

Attack 3: Malicious External Attackers

- Randomly crafted payload
- Impersonation
- Compromised account
- DDos & Botnet



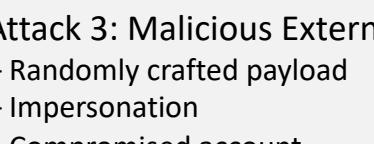
(End-User)



Endpoint Agents
Daemon Process



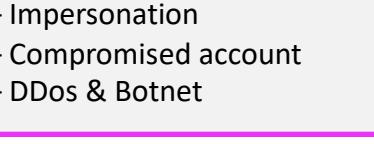
Malicious Apps/Users



Corporate Users



Access Client Process
(e.g. UI, DB shell)



MFA (e.g. SMS, OTP)



Database Application



File System Storage (e.g. hdfs)



Secure Storage (HSM)



Behavior Biometric Systems



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



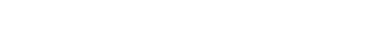
Cloud I/O



I/O Module



Behavior Biometric Analysis



Biometric Behavior Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



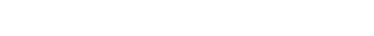
Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



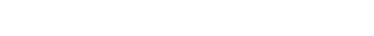
Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



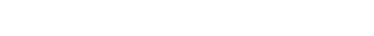
Behavior Biometric Modules



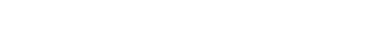
Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



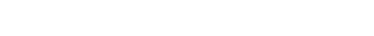
Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



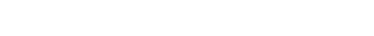
Behavior Biometric Analysis



Behavior Biometric Modules



Training Data



Meta Data/Hyperparameters
Behavior Biometric Systems



Cloud I/O



I/O Module



Behavior Biometric Analysis



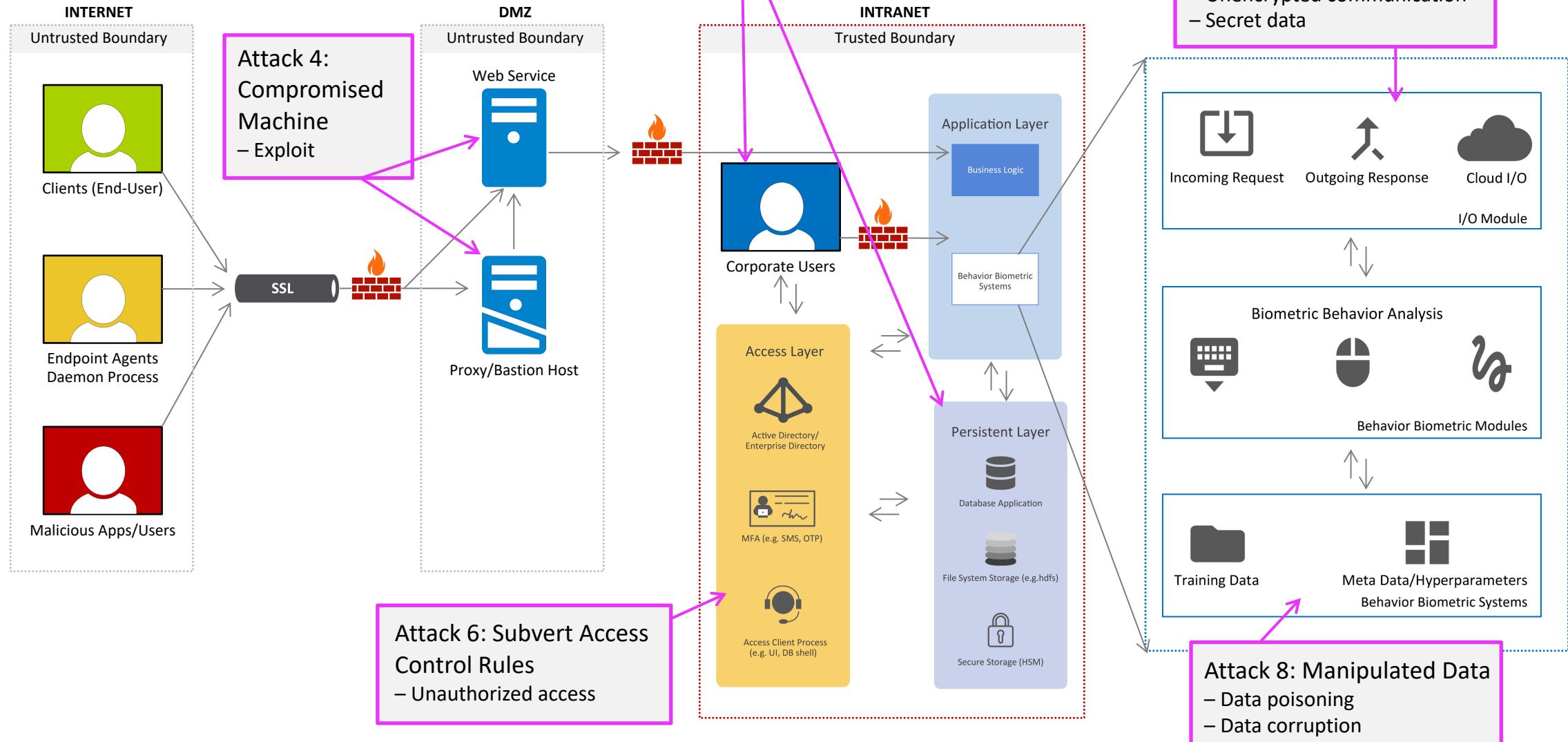
Behavior Biometric Modules



Client Side Attack Protection

- Targeting to manipulate input and perform black-box attacks
- End point protection
 - Binary Hardening
 - Code Obfuscation
 - Transport Level Encryption + Message Level Encryption + Integrity checking
- Never Trust External Input

Server Side Attacks



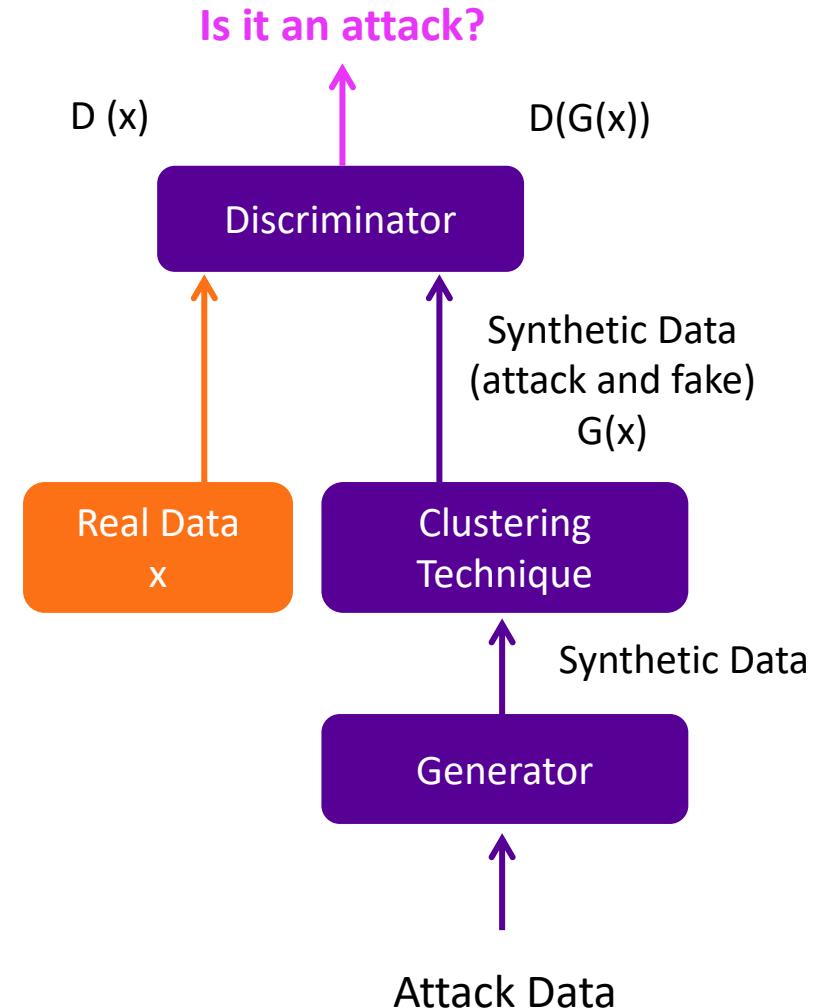
Server Side Attack Protection

- Targeting to manipulate internal data + white-box attacks
- Protection against abused/unauthorized access
 - Encrypting data at rest
 - Access control/monitoring on key infrastructure (e.g. DB/model store)
 - Message Authentication + Firewall

Model Level Attack Protection

Detection Models

- Filtering/Preprocessing the samples by another model
- Trained with attack data and real data to detect attacks

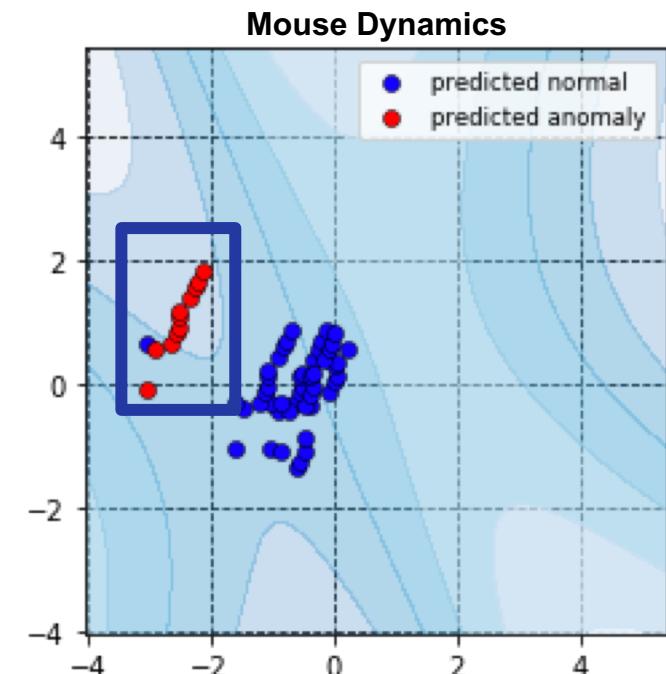
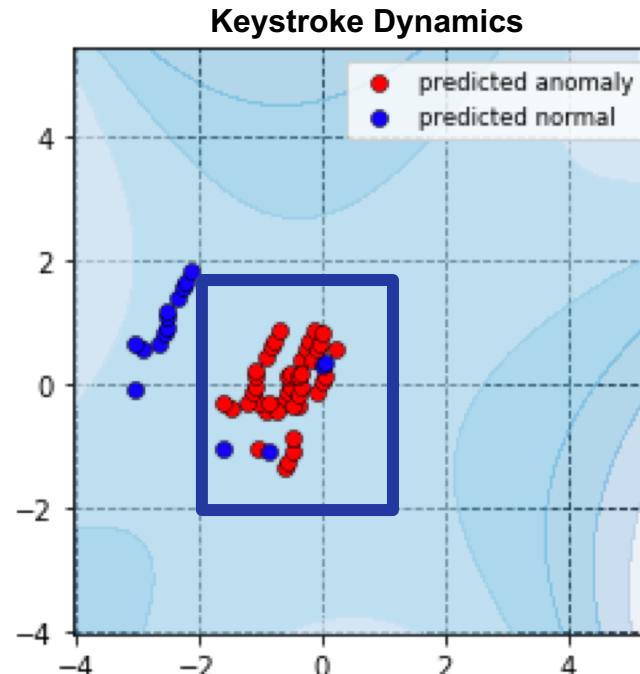
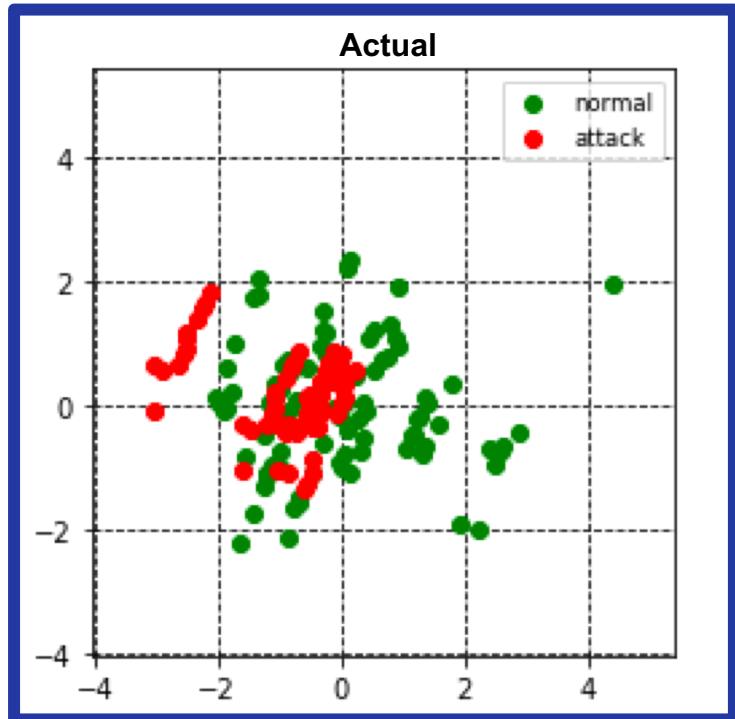


Model Level Attack Protection

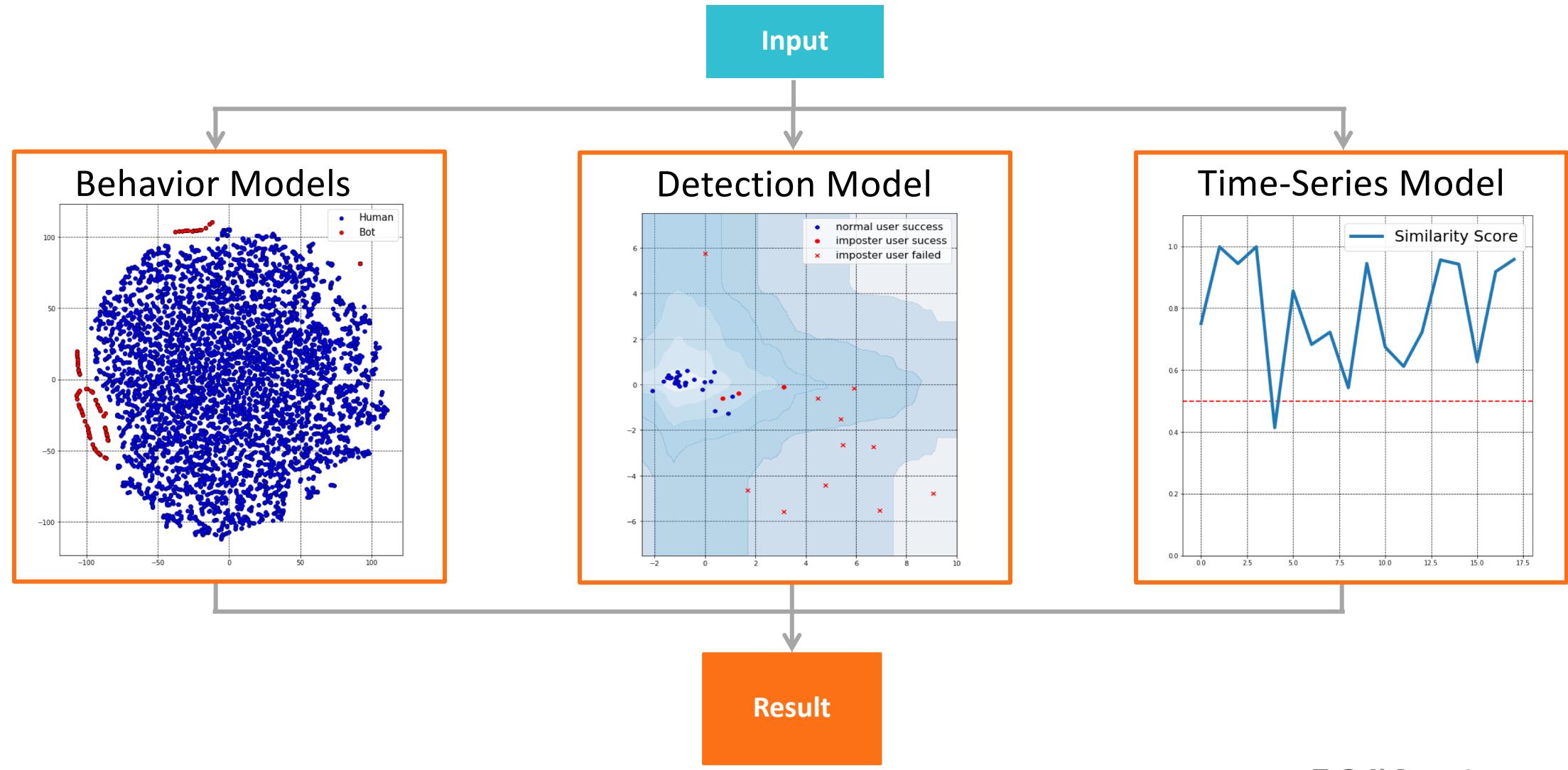
Model Hardening

- Model immutable to crafted input/manipulated data
 - Trained with attack data
 - Additional labels to attack data
- Combination of different models – Ensemble Models
 - Better results (smaller false acceptance/rejection)
 - Better security controls (more entropy)

Ensemble models



Overall Architecture for Behavioral Biometrics Systems



Takeaway

Adversarial Machine Learning based attacks against behavioral biometric systems can cause significant impact to this authentication mechanism. Traditional security controls combined with **Machine Learning specific countermeasures** mitigate the impact.

Q & A

References

- Gaddam, Ajit. "Usage of Behavioral Biometric Technologies to Defend Against Bots and Account Takeover Attacks." (2019).
- Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *arXiv preprint arXiv:1412.6572* (2014).
- Grosse, Kathrin, et al. "On the (statistical) detection of adversarial examples." *arXiv preprint arXiv:1702.06280*(2017).
- Gunetti, Daniele, and Claudia Picardi. "Keystroke analysis of free text." *ACM Transactions on Information and System Security (TISSEC)* 8.3 (2005): 312-347.
- Lu, Jiajun, et al. "No need to worry about adversarial examples in object detection in autonomous vehicles." *arXiv preprint arXiv:1707.03501* (2017).
- Lee, Hyeungill, Sungyeob Han, and Jungwoo Lee. "Generative adversarial trainer: Defense to adversarial perturbations with gan." *arXiv preprint arXiv:1705.03387*(2017).