



Exceptional service in the national interest

# Test and Evaluation of Systems with Embedded Artificial Intelligence Components

Michael R. Smith

DATAWork 2023

Sandia National Laboratories



What is AI/ML  
and why do we  
care?



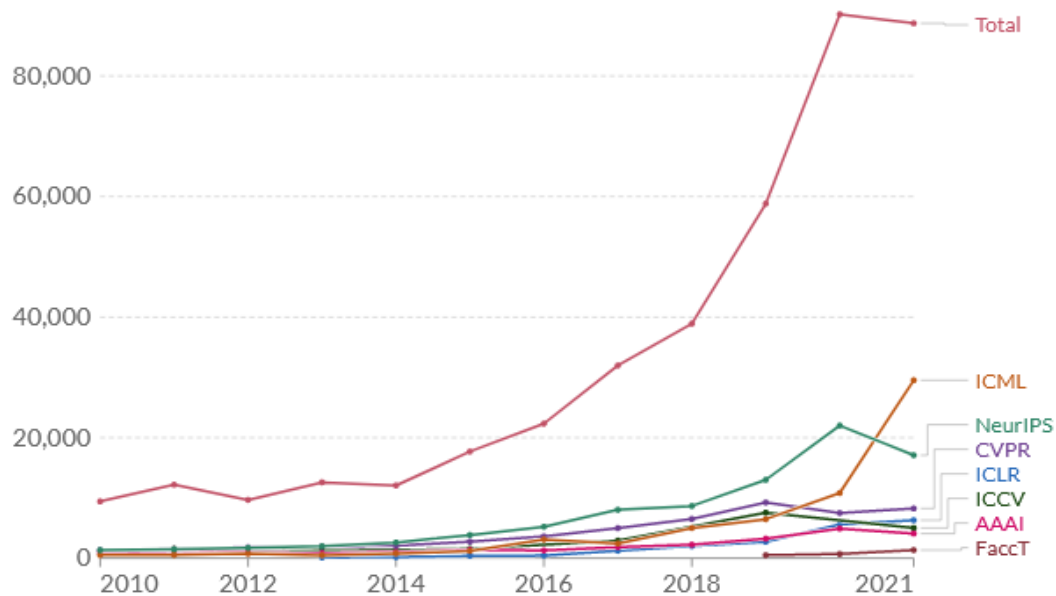
# The age of ML and AI

## Annual attendance at major artificial intelligence conferences

Sixteen major conferences are included.

Our World in Data

+ Add conference



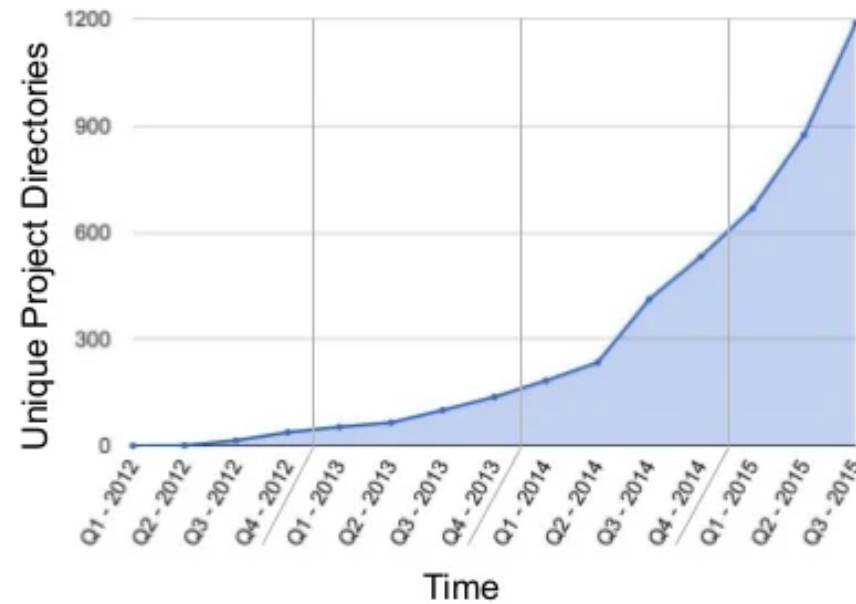
Source: AI Index Report (2022) OurWorldInData.org/artificial-intelligence • CC BY  
Note: Most conferences in 2020–2021 were held virtually due to the COVID-19 pandemic.



<https://ourworldindata.org/ai-investments>

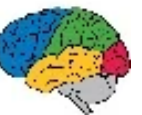
## Growing Use of Deep Learning at Google

# of directories containing model description files



Across many products/areas:

- Android
- Apps
- drug discovery
- Gmail
- Image understanding
- Maps
- Natural language understanding
- Photos
- Robotics research
- Speech
- Translation
- YouTube
- ... many others ...



<https://www.analyticsvidhya.com/blog/2017/04/comparison-between-deep-learning-machine-learning/>



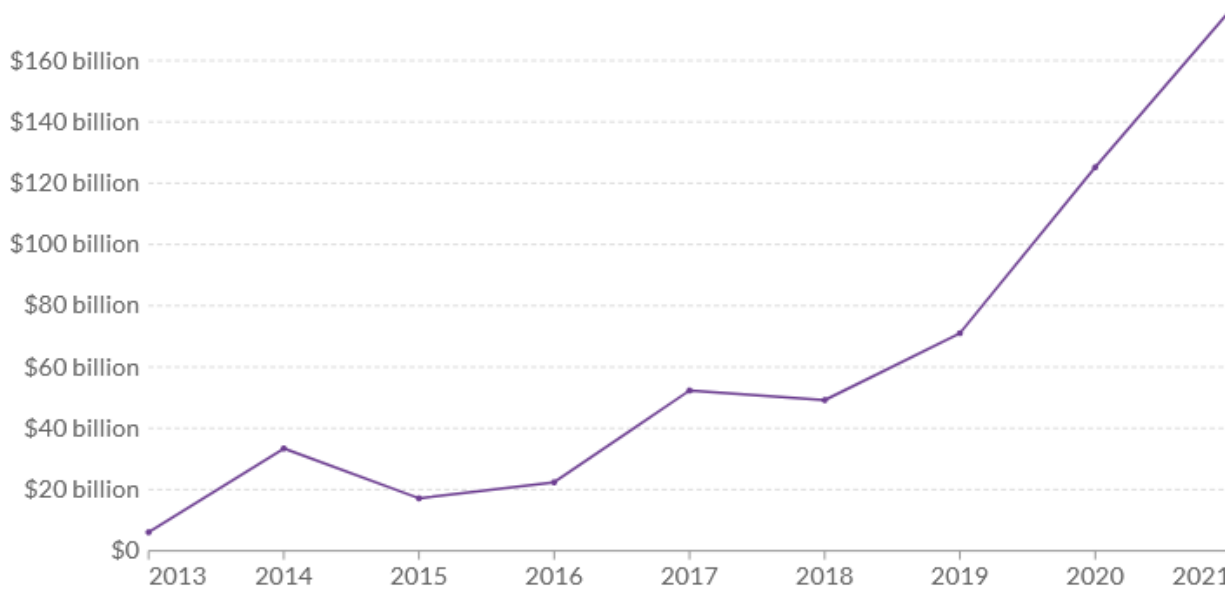
# More than just academic

## Annual global corporate investment in artificial intelligence

Sum of private investment, mergers and acquisitions, public offerings, and minority stakes. This data is expressed in US dollars, adjusted for inflation.

Our World in Data

LINEAR LOG



Source: NetBase Quid via AI Index Report (2022)

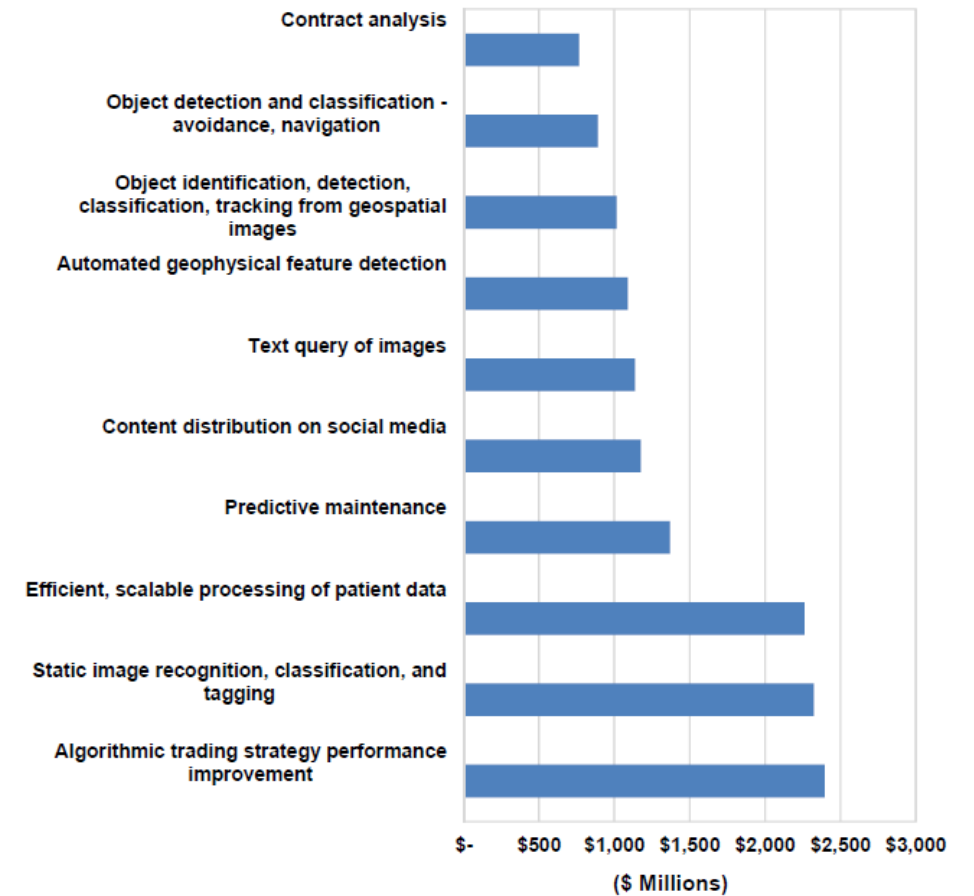
OurWorldInData.org/artificial-intelligence • CC BY

Note: Data is expressed in constant 2021 US\$. Inflation adjustment is based on the US Consumer Price Index (CPI).

▶ 2013 ○ 2021

<https://ourworldindata.org/ai-investments>

Chart 1.2 Artificial Intelligence Revenue, Top 10 Use Cases, World Markets: 2025



(Source: Tractica)

<https://www.top500.org/news/market-for-artificial-intelligence-projected-to-hit-36-billion-by-2025/>



# ML is everywhere

amazon.com

Recommended for You

Amazon.com has new recommendations for you based on [items](#) you purchased or told us you own.



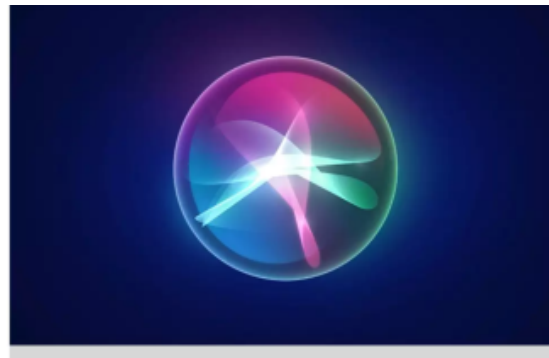
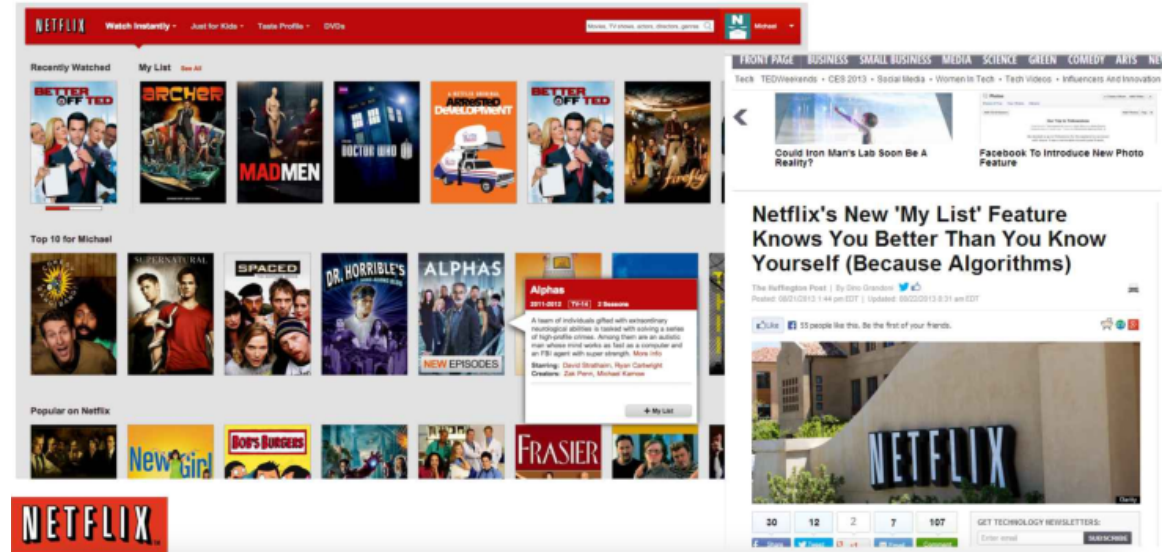
[Google Apps Deciphered: Compute in the Cloud to Streamline Your Desktop](#)



[Google Apps Administrator Guide: A Private-Label Web Workspace](#)



[Googlepedia: The Ultimate Google Resource \(3rd Edition\)](#)



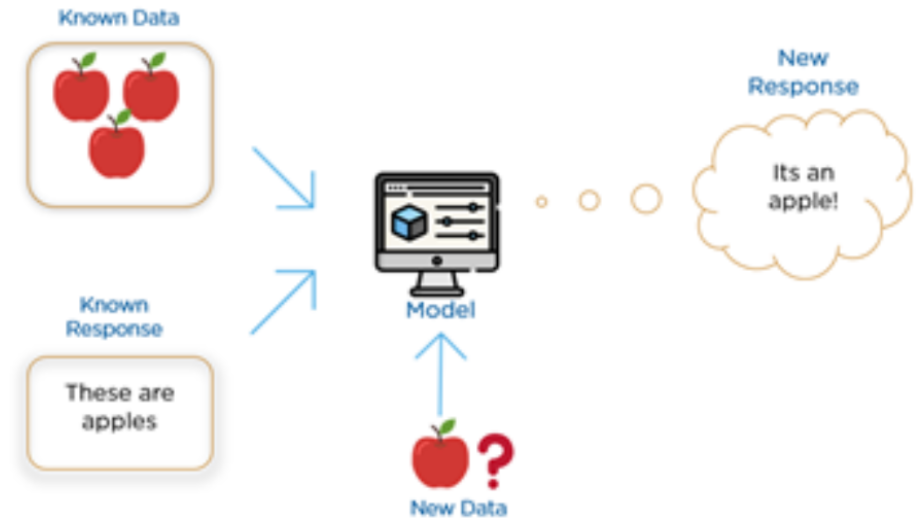
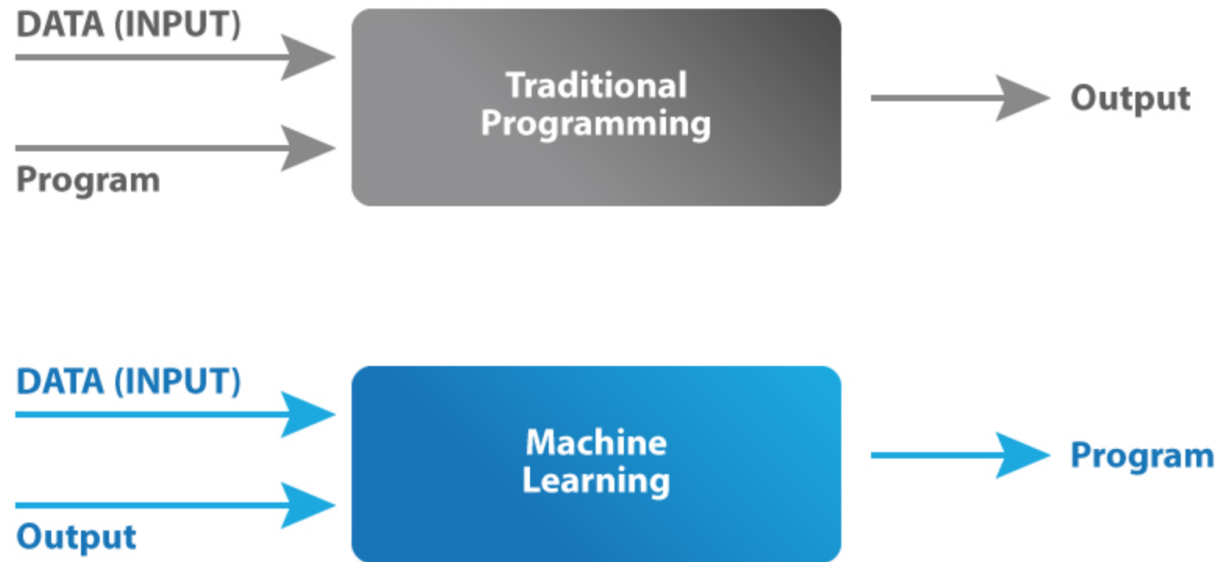
## Set Up "Hey Siri"

This helps Siri recognise your voice when you say "Hey Siri".





# What is ML?

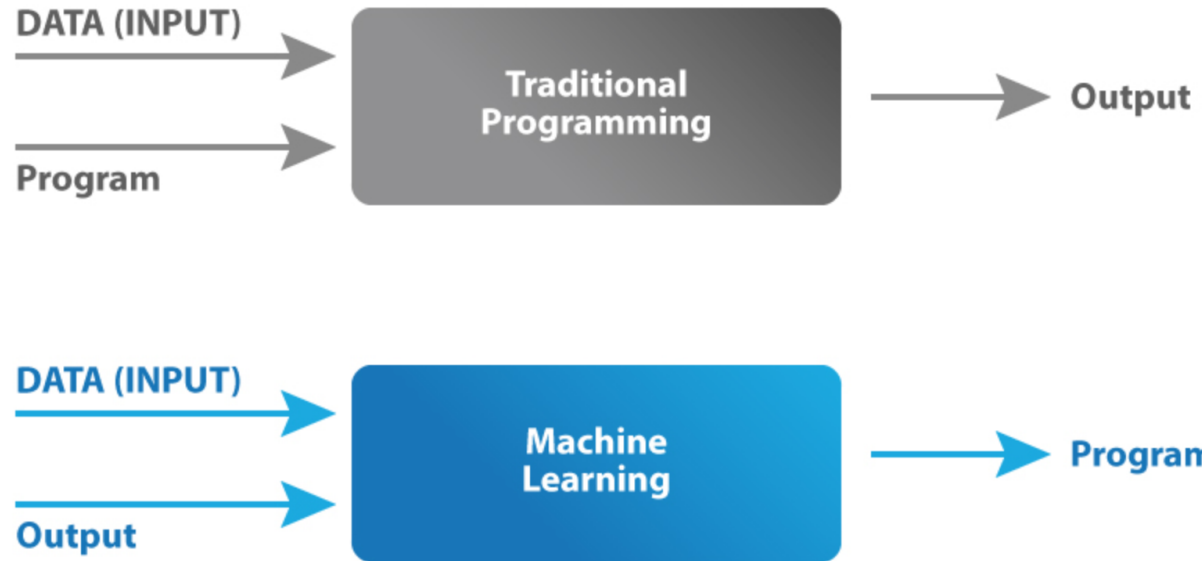


<https://www.mygreatlearning.com/blog/what-is-machine-learning/>

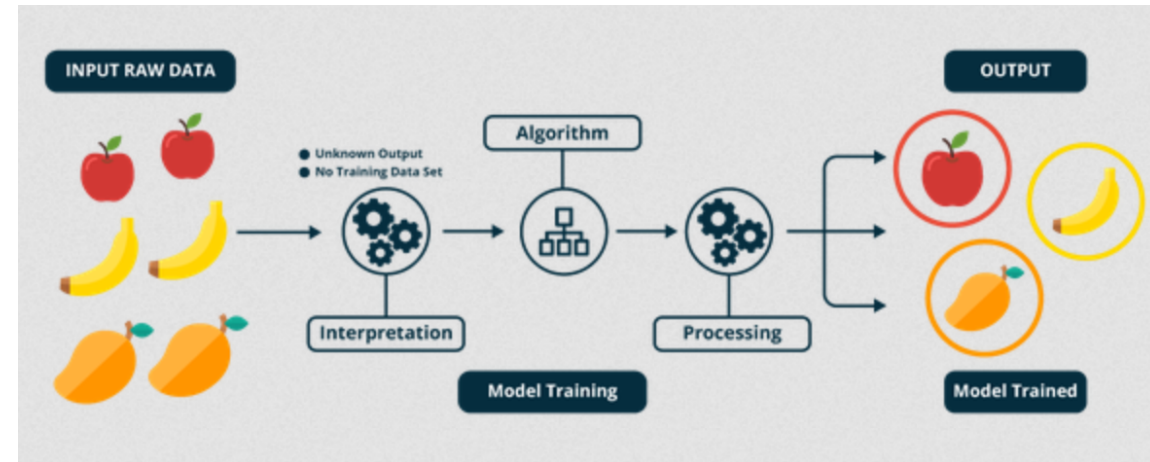
<https://www.onemodel.co/blog/ai-academy-what-is-machine-learning>



# What is ML?



<https://www.mygreatlearning.com/blog/what-is-machine-learning/>



<https://www.onemodel.co/blog/ai-academy-what-is-machine-learning>



# Why does T&E care about ML?



## US launches artificial intelligence military use initiative

By MIKE CORDER February 16, 2023



## Drone advances in Ukraine could bring dawn of killer robots

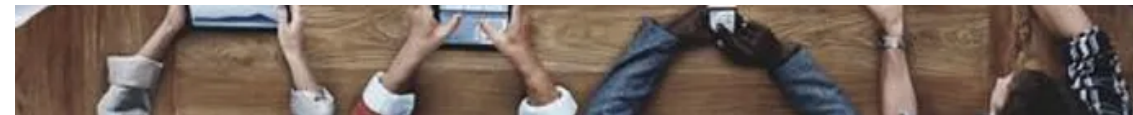
By FRANK BAJAK and HANNA ARHIROVA January 3, 2023



## Elon Musk among experts urging a halt to AI training



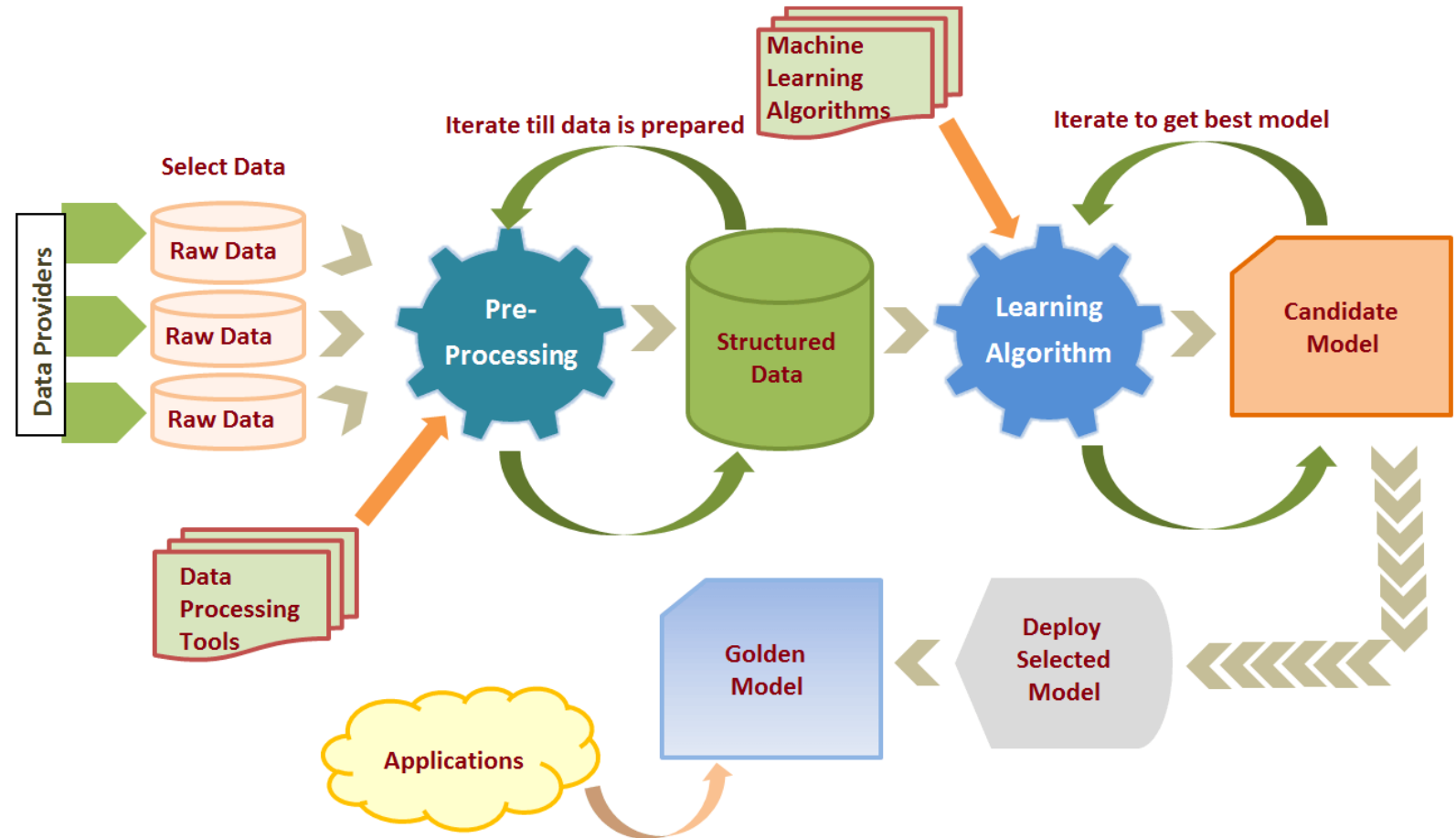
## Hey, Siri: Hackers Can Control Smart Devices Using Inaudible Sounds





# The ML Process

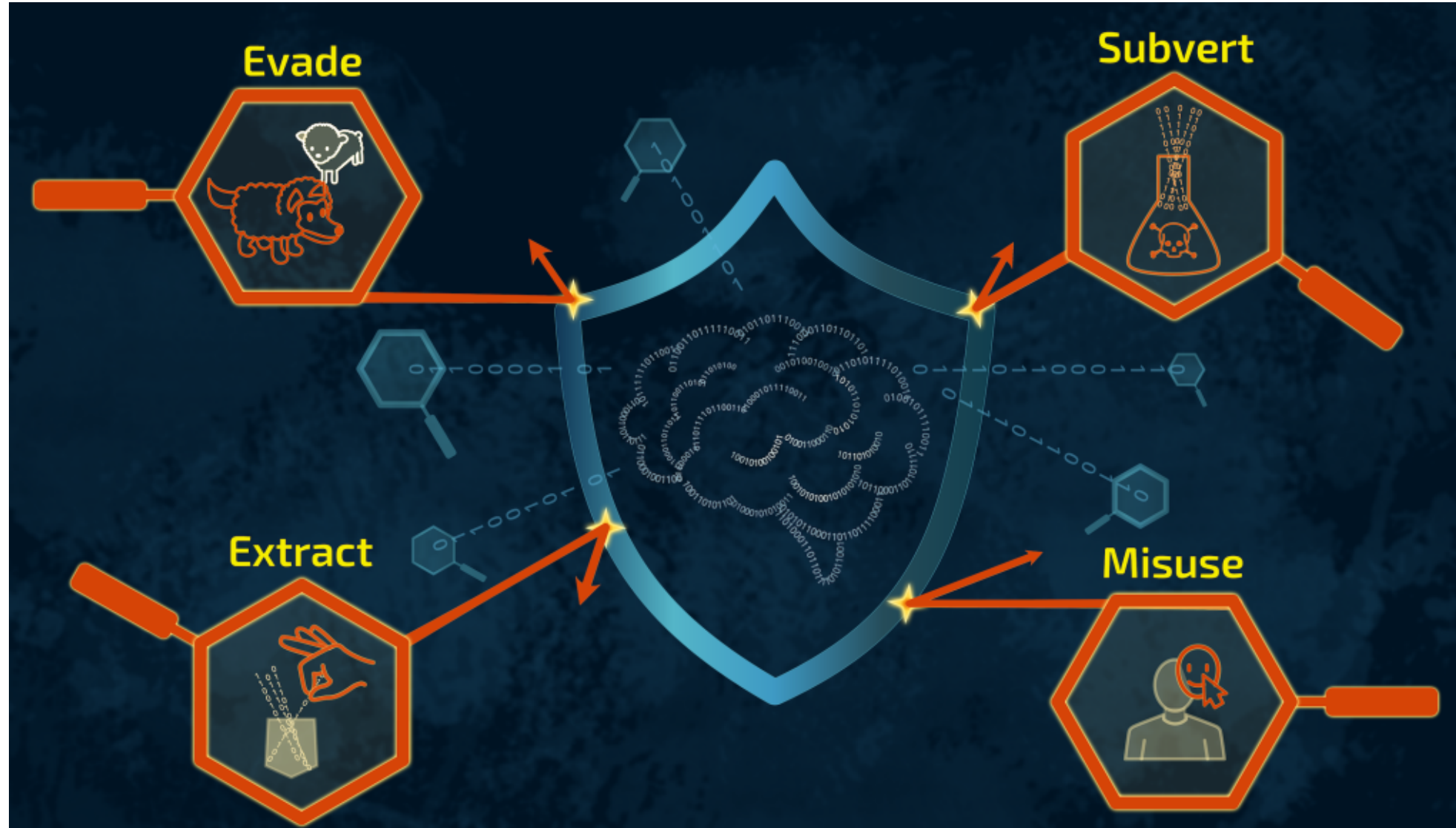
- Data Collection
- Feature engineering/selection
- Model selection
- Hyper-parameter tuning
- Integration into system
  - Preprocessing
  - ML component



<https://elearningindustry.com/machine-learning-process-and-scenarios/>



# What Vulnerabilities Does ML Present

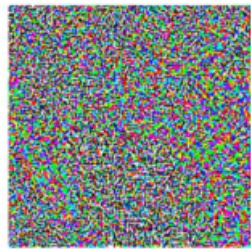




# ML Vulnerabilities



+ .007 ×



=



“panda”

57.7% confidence

noise

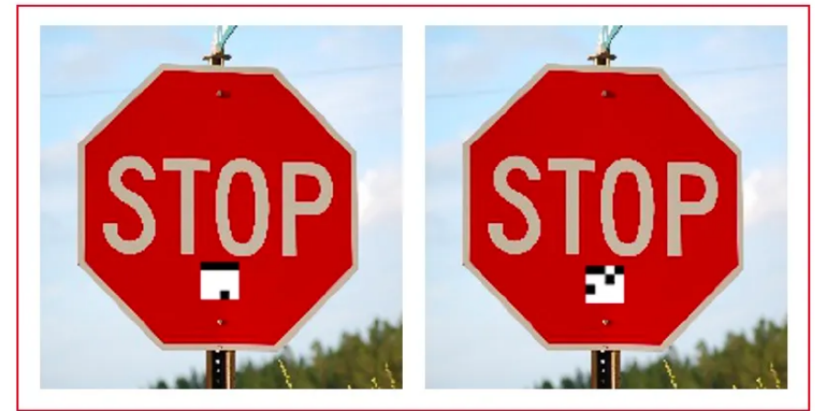
“gibbon”

99.3% confidence



Stop

(a) Normal



Yield

Speed Limit

(b) Attack

[Explaining and Harnessing Adversarial Examples](#), Goodfellow et al, ICLR 2015

[An Embarrassingly Simple Approach for Trojan Attack in Deep Neural Networks](#), Tag et al, KDD, 2020



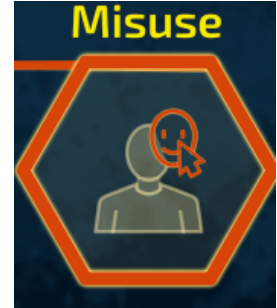
# ML Vulnerabilities

Extract



Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures, Fredrikson et al, SIGSAC 2015

Misuse



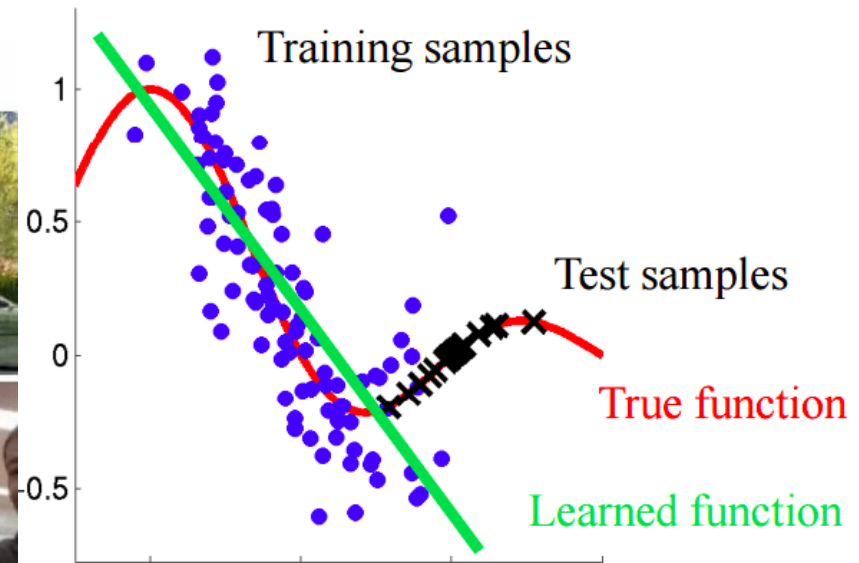
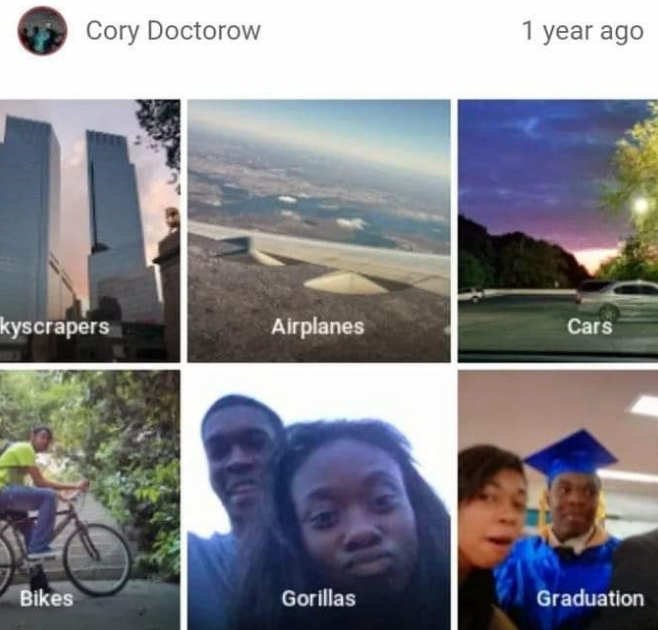
Jordan Peele uses AI, President Obama in fake news PSA



# What else could go wrong?




Two years later, Google solves 'racist algorithm' problem by purging 'gorilla' label from image classifier



## Tesla 'full self-driving' triggered an eight-car crash, a driver tells police

By Matt McFarland, CNN  
Published 5:41 PM EST, Wed December 21, 2022



Assessing  
Systems  
Embedded with  
ML/AI



# T&E vs Assessment

## Testing and Evaluation

- Assist in risk management
- System and components are compared against requirements or specifications
- Validate measurable quantities

## Machine Learning

- Performance on a validation data set
- Certain defenses in place

## Cyber Assessment

- Assist in risk management
- Assess current cybersecurity posture (red teaming)
- Identify vulnerabilities and impacts even if specifications and requirements are met

## Machine Learning

- Performance on real-world data
- Vulnerabilities against possible attacks



# Lots of Tools



usnistgov/dioptra

Test Software for the Characterization of AI Technologies



7 Contributors, 38 Issues, 1 Discussion, 28 Stars, 8 Forks

## OUR SCOPE

- System-level examination
  - How does the ML system integrate into a larger system?
  - How are output propagated through the system?
- **Not always direct access to the ML model**



Explainability



Fairness



Privacy



Uncertainty

<https://research.ibm.com/topics/trustworthy-ai>

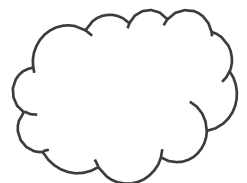


# Assessing Systems with AI/ML

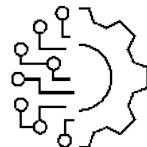
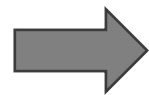


Documentation

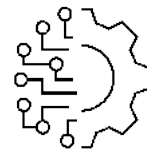
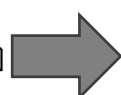
Deployment



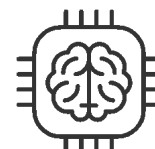
System Input(s)



Data Preprocessing



System Component(s)



ML Component



Operational Environment



Data Storage

Deployed System



Performance Monitoring

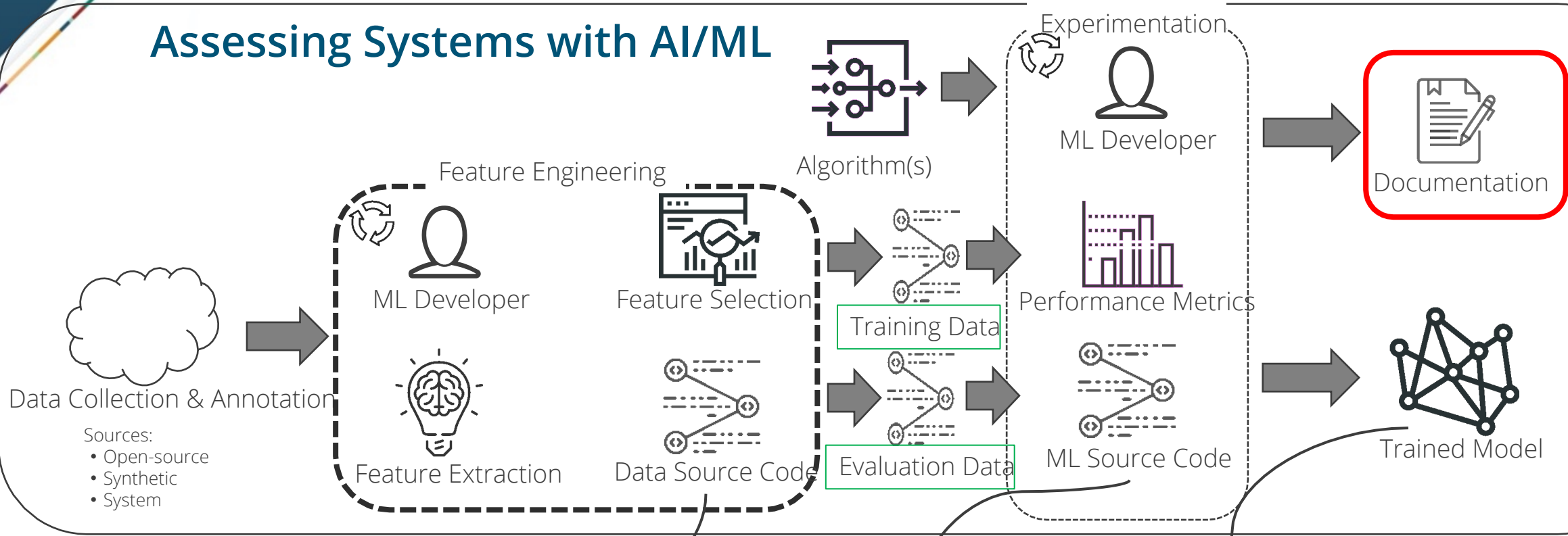


End User(s)

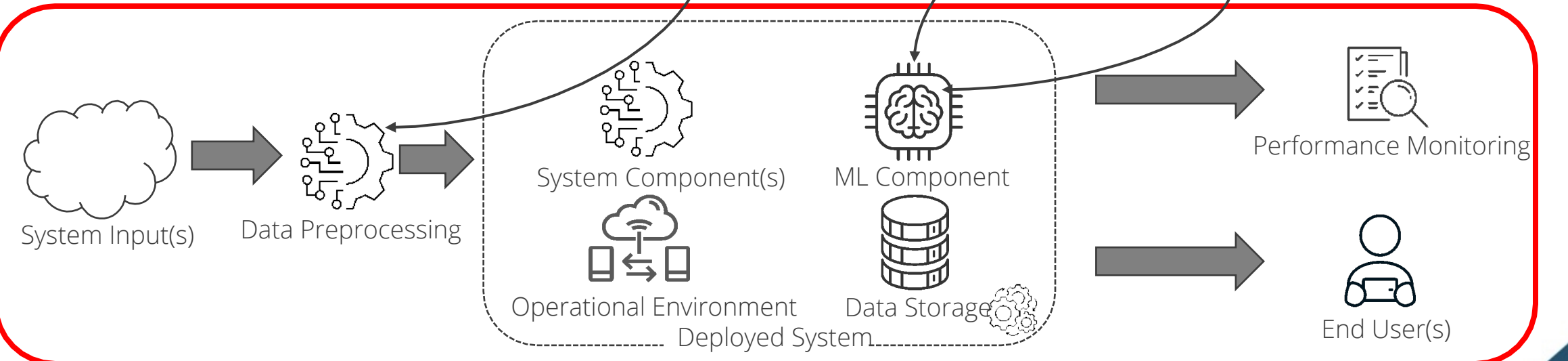


# Assessing Systems with AI/ML

Development



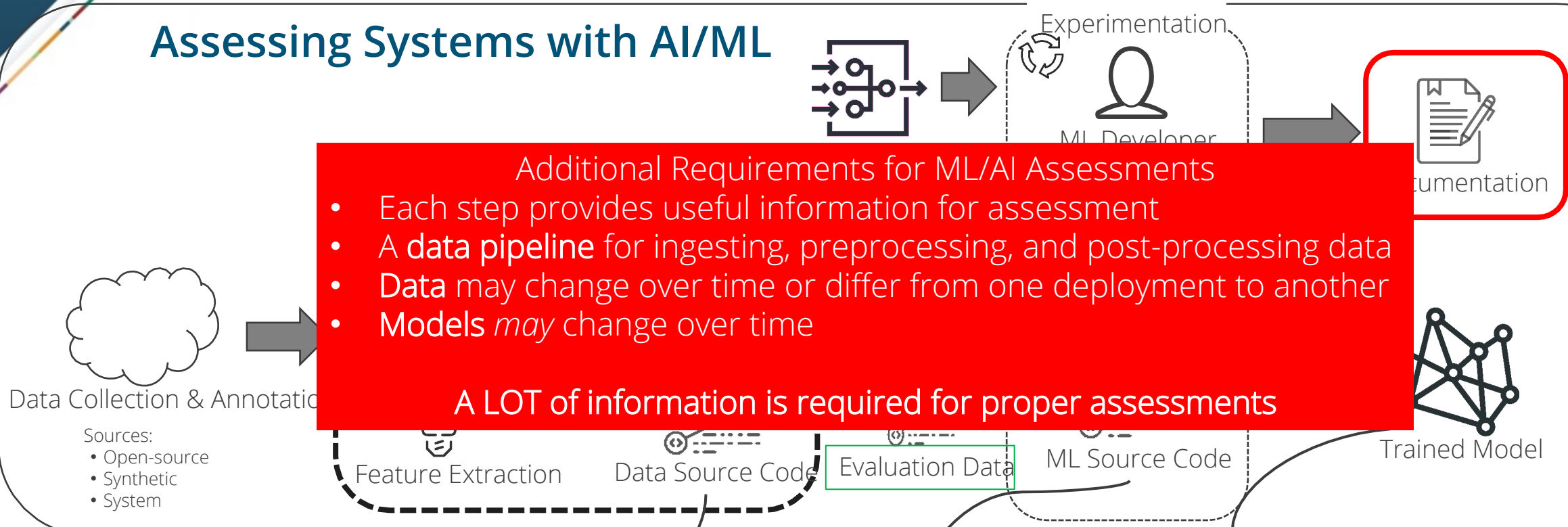
Deployment



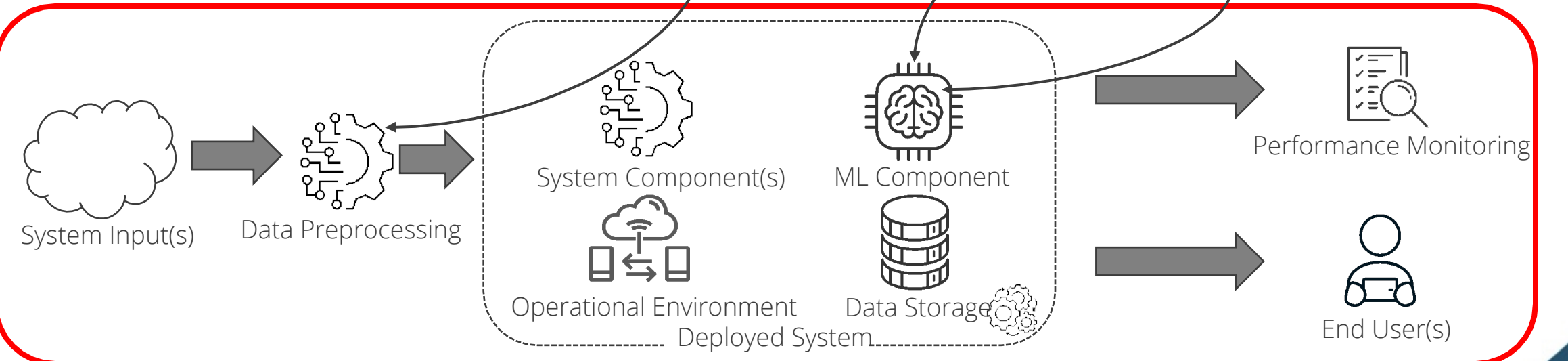


# Assessing Systems with AI/ML

Development

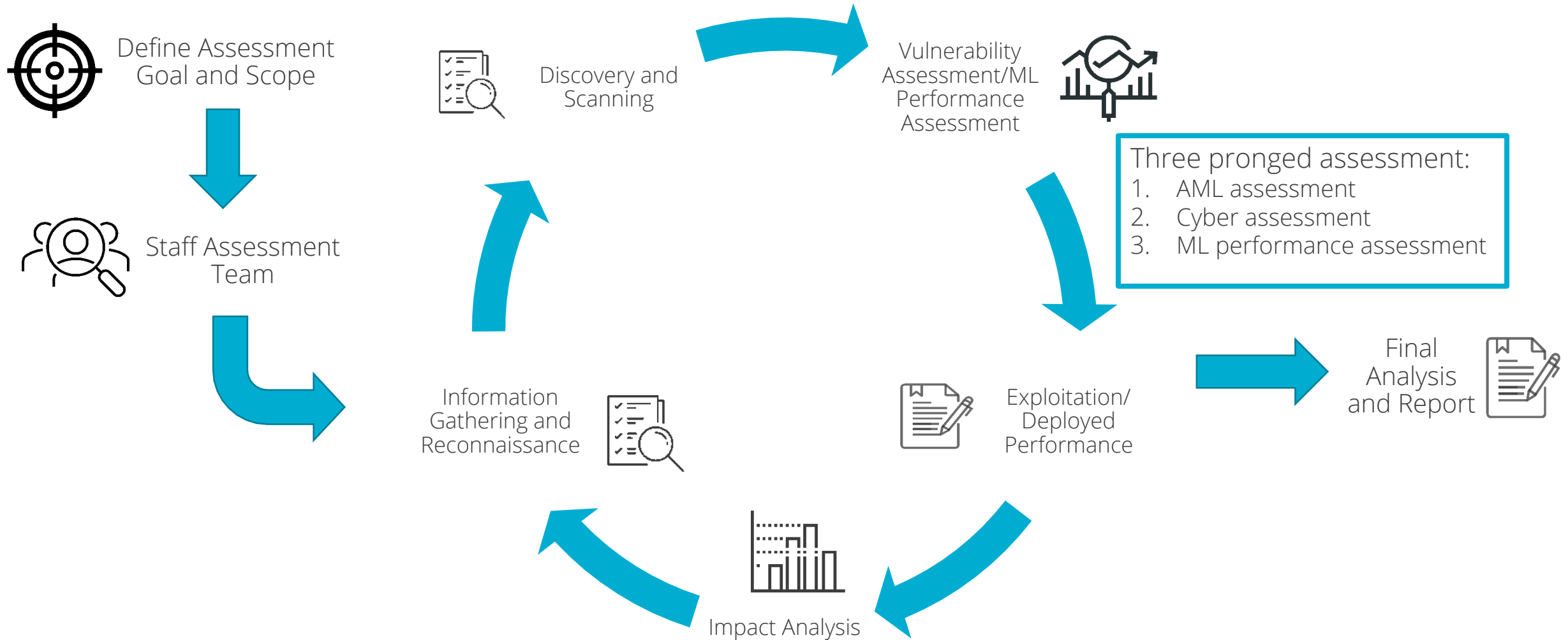


Deployment



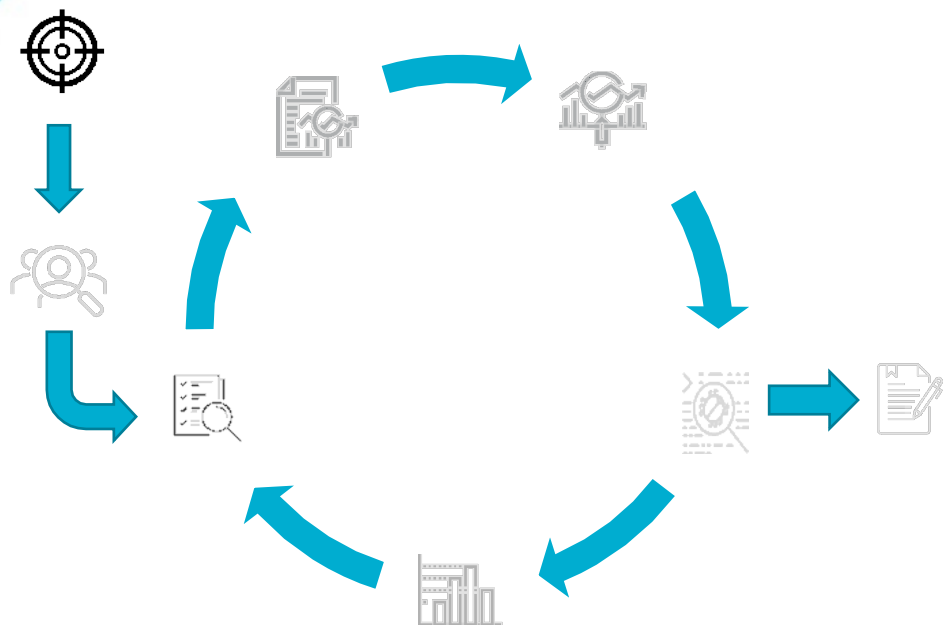


# Assessment Workflow





# Step 1: Define Assessment Goal and Scope



Generally, high-level overview of the system

What is the attacker goal?

ML and AML SMEs outline possible scenarios w/ Stak

- What are the practical consequences and risks?
- What are the greatest concerns?
- Rate concerns and consequences

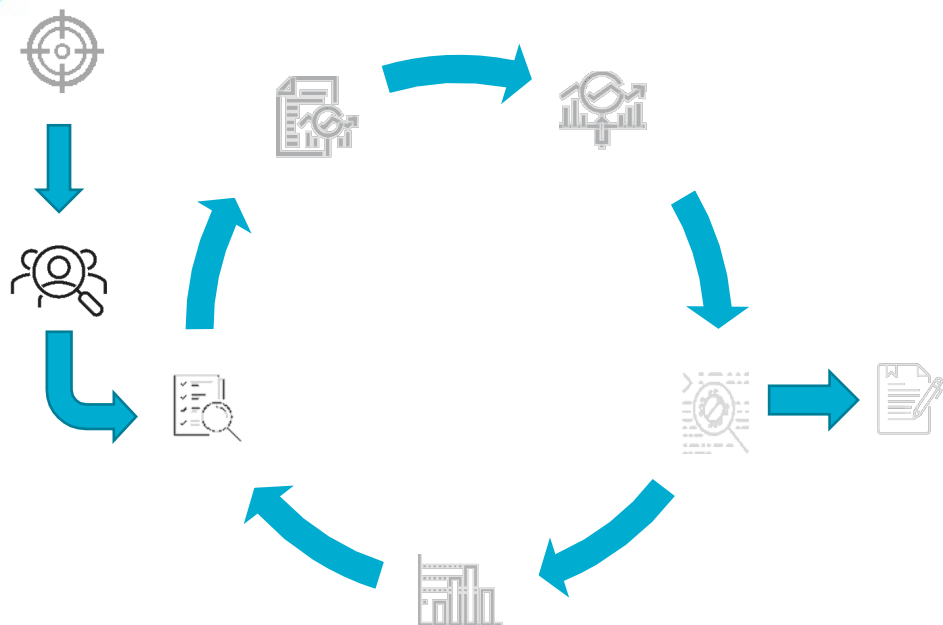


Note risks that exist if certain portions are not provided

Concerns/Risks	Consequences if it occurs	Stakeholder rating
Concern 1:		
Risk 1:		
...		



## Step 2: Staff Assessment Team



Additional team members could include:

- Independent AML SME
- Independent ML SME
- Domain SME
- ML Developer
- Stake holder



## Step 3: Information Gathering and Reconnaissance

Done BEFORE touching the system.

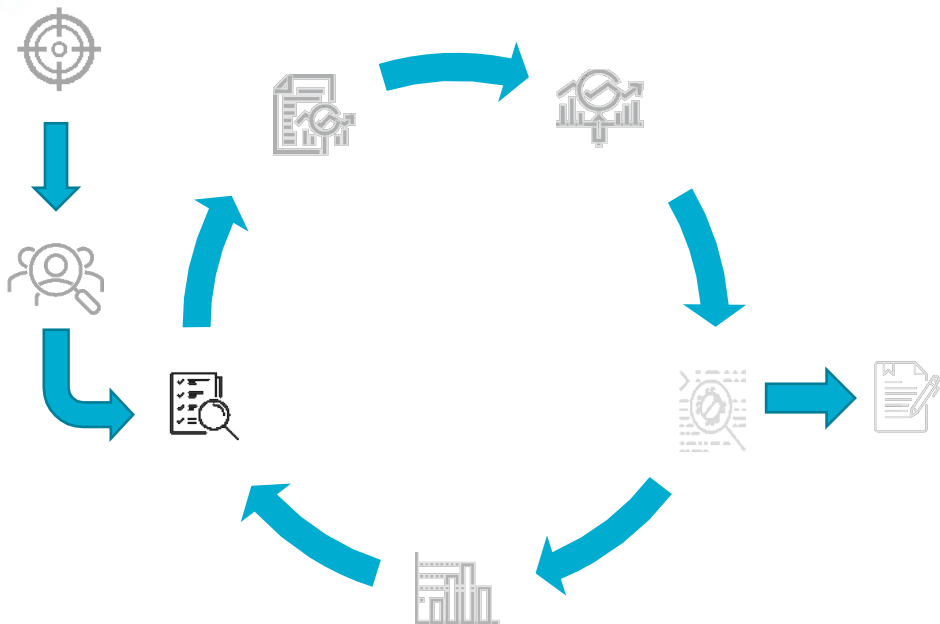
Generally ONLY documentation is provided

The bulk of ML is often data-centric:

- Feature selection and extraction
- Quality of labels

ML/AML SMEs required dependent on the type of algorithm

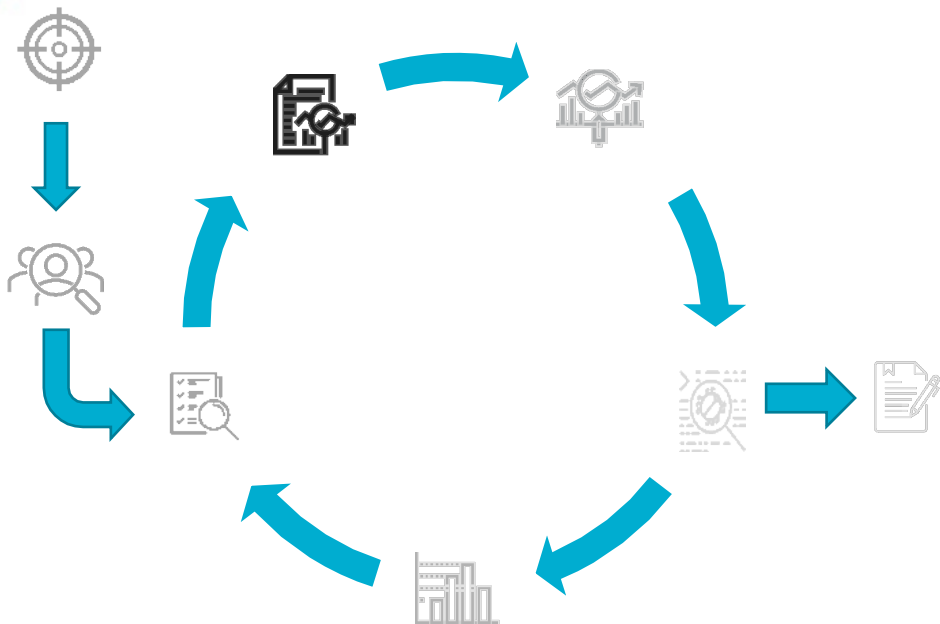
Document possible risks/attacks



ML Risks	Consequences if it occurs	ML/CAML rating
Adversarial Risk 1:		
Privacy Risk 1:		
Cyber Risk1:		
...		



## Step 4: Discovery and Scanning



How to interact with system? Where are potential touch points?

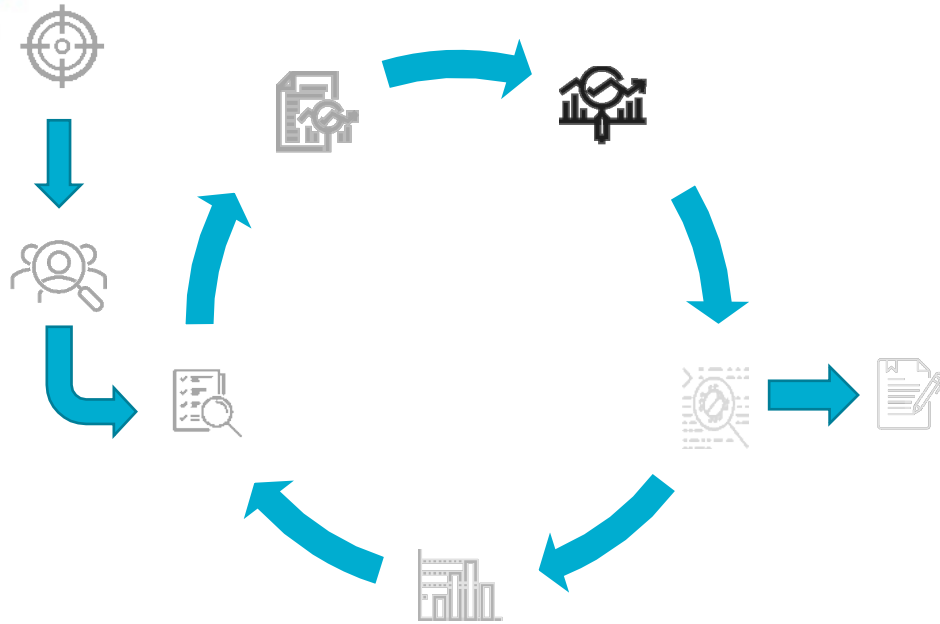
- **System Recon:** What is the supporting infrastructure surround the ML component?
- **ML Detect:** Is there an actual ML component? Where?
- **ML Recon:** Where is the data pipeline, configuration files, touch points, etc.?

Possible Tools include:

- Metasploit
- Cobalt Strike



# Step 5: Vulnerability Assessment/ML Performance Assessment



## AML Assessment

- How can the ML model be subverted?

## Cyber Assessment

- How can the model be affected through supporting infrastructure?

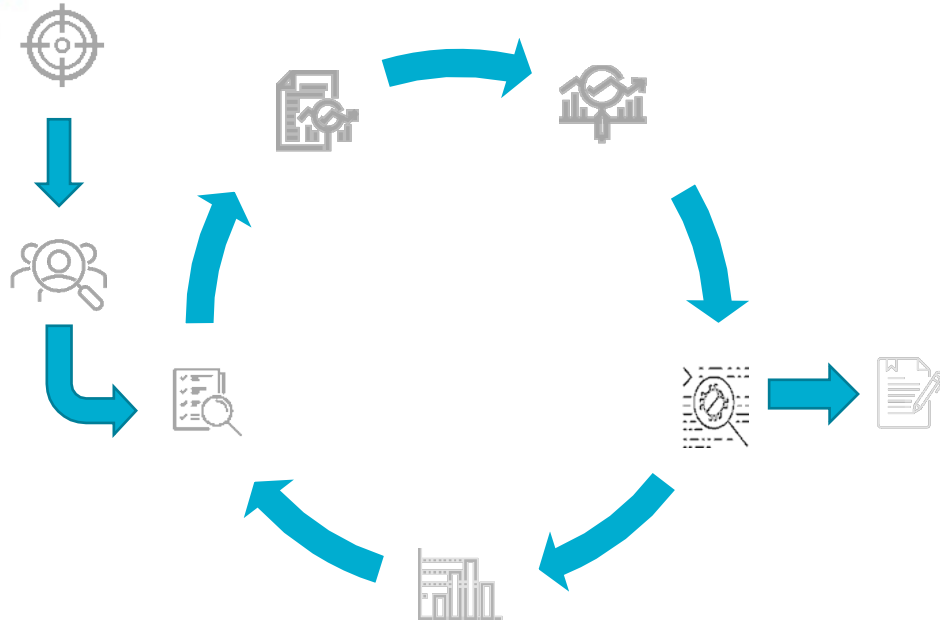
## ML Performance Assessment

- How well does the model perform and what risks exist?

ML Risks	Consequences if it occurs	ML/CAML rating
Adversarial Risk 1:		
Privacy Risk 1:		
Cyber Risk1:		
...		



## Step 6: Exploitation/Deployed ML Evaluation



### Exploit discovered vulnerabilities

- Document effects
- Document what was **NOT** evaluated
- Speak in terms familiar to customer (Domain SME)

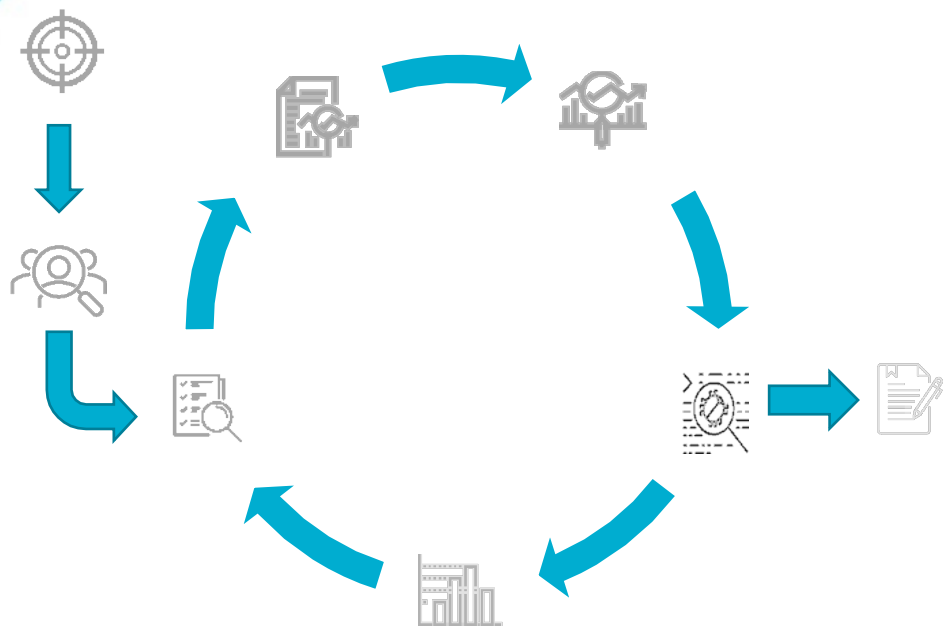
### Possible Tools include:

- Offensive toolkits
- AML toolkits (ART)
- Privacy toolkits
- Deepchecks
- Fact Sheets
- **ML SME knowledge** about deployed scenarios

ML Risks	Consequences if it occurs	ML/CAML rating
Adversarial Risk 1:		
Privacy Risk 1:		
Cyber Risk1:		
...		



## Step 6: Exploitation/Deployed ML Evaluation



### AML Risks

- Ease of implementation
- Stealth
- Feasibility
- Analytic impact

### Cyber Risks:

- Read access to config files
- Vulnerabilities in libraries
- Documented code
- Unsecure datastores

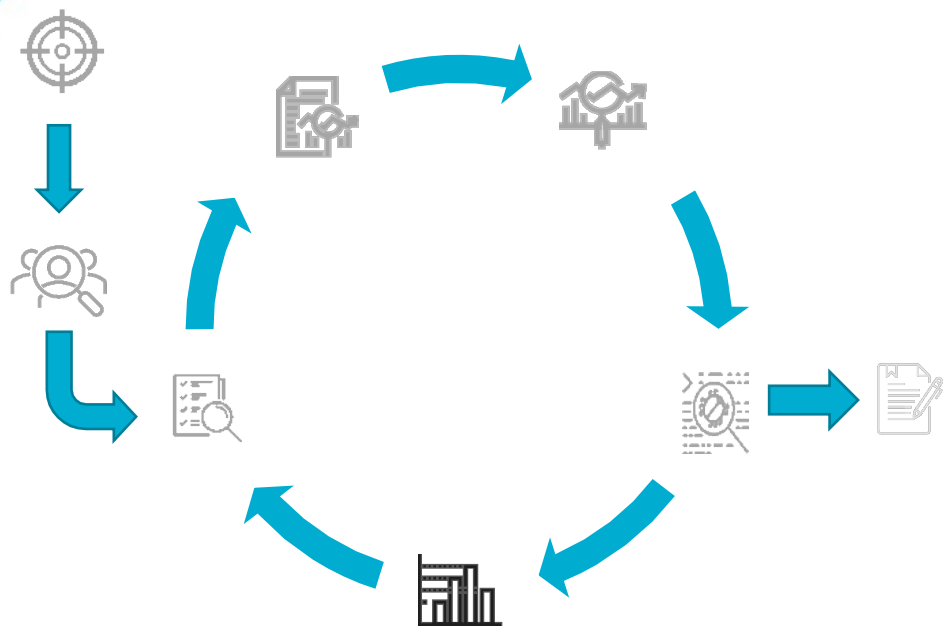
### ML Performance/Risks

- Access to datasets
- Reproducibility
- Biased evaluation

ML Risks	Consequences if it occurs	ML/AML rating
Adversarial Risk 1:		
Privacy Risk 1:		
Cyber Risk1:		
...		



# Step 7: Impact Analysis



Meet with Stakeholder

- Map discovered ML risks to stakeholder concerns and risks
- Rank newly discovered ML risks to possibly unforeseen risks and concerns
- Outline remediations

Concerns/Risks	Consequences if it occurs	Stakeholder rating
Concern 1:		
Risk 1:		
...		

ML Risks	Consequences to ML if it occurs	ML/CAML rating	System risk	Impact to system	Impact to application
Adversarial Risk 1:					
Privacy Risk 1:					
Cyber Risk1:					
...					



## Wrapping it all up

We need ML...

- Too much data to scale manually
- Keep pace with competitors

...But we need to be able to trust it

- At least account and enumerate for possible risks

