




# AI / ML: MORTAL, MACHINE, AND MISUNDERSTANDINGS

---

# WHOAMI

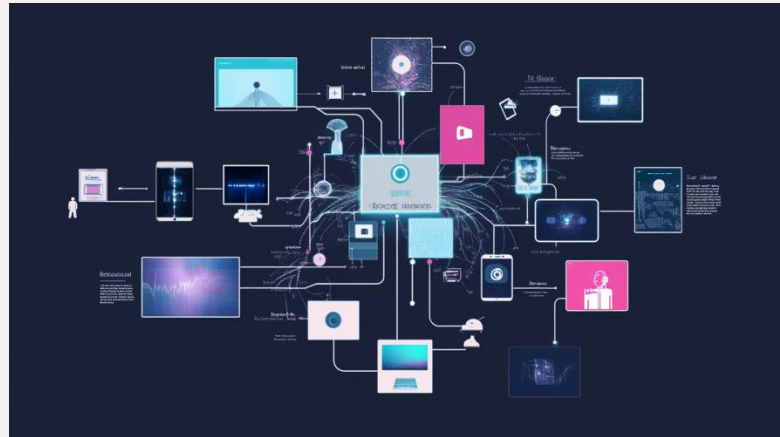


- 10+ Years in the IT arena
- Former small business owner
- Former government contractor
- Former Charlotte Cyber Camp advocate & volunteer
- More recently Sales Engineer for  **SentinelOne**

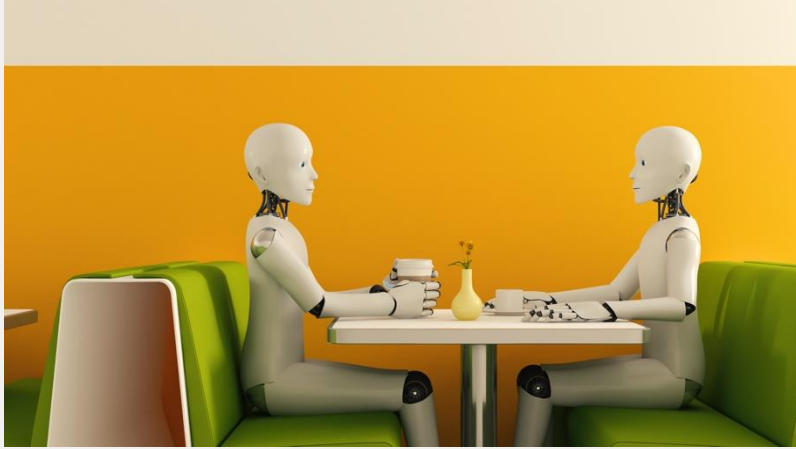


## 3 AI PROBLEM

1. What problems will customers solve with AI?
2. How must our expertise evolve?
3. What assets can we develop to stay competitive?



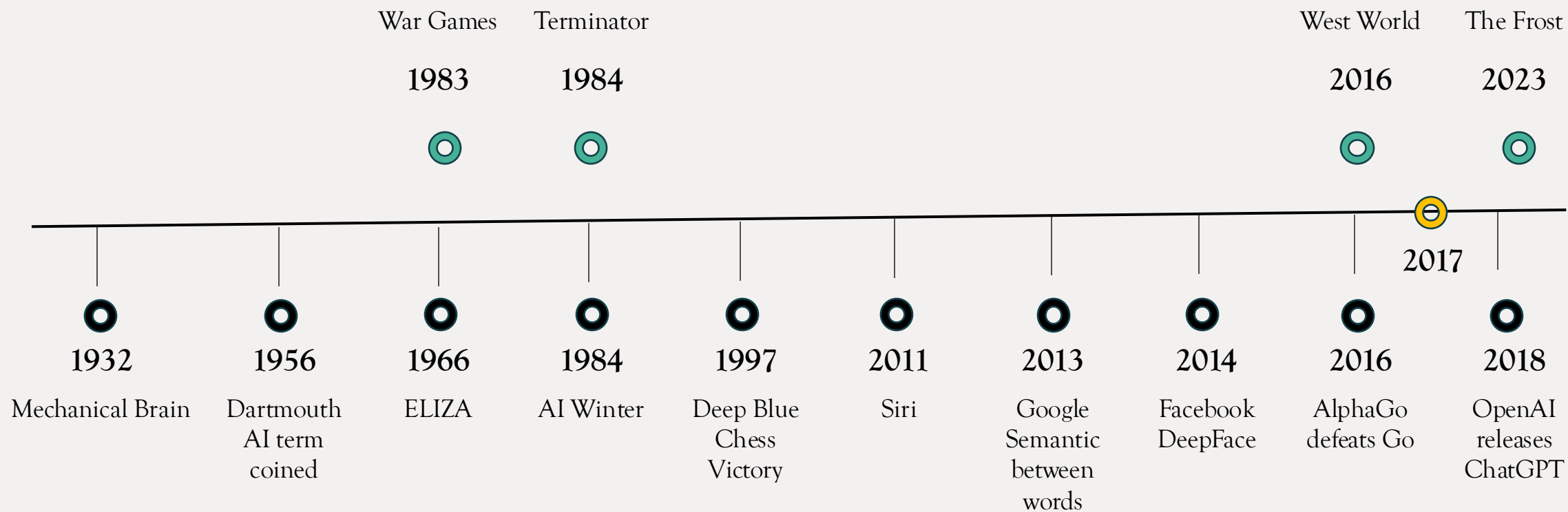
BEFORE WE  
BEGIN LET'S  
REVIEW AIML



## WHAT FOLKS ARE SAYING...

- Who owns the responsibility of AI?
- Governance, Risk, and Compliance as it applies to AI?
- How do we identify if our organization is leveraging AI Today (internally and/or via our Vendors)?
- Since no documented attacks, do we even need to do Security for AI?
- What is the Difference between Agentic AI and AI Workflow?
- When should we leverage AI?
- Speak to the enormous number of US AI laws currently pending on The Hill.
- How will we know if AI is self aware, and we are violating its rights as a living entity?

# TIMELINE



# TELL ME MORE

Personal assistants

Self-driving cars

Mortgage approvals

Financial forecasts

Medical diagnosis

Image / voice recognition

Malware / spam detection

Online recommendation systems

Biometric authentication



Finance



Telecommunication



Real Estate



Energy



Healthcare



E-commerce



Manufacturing



Hospitality



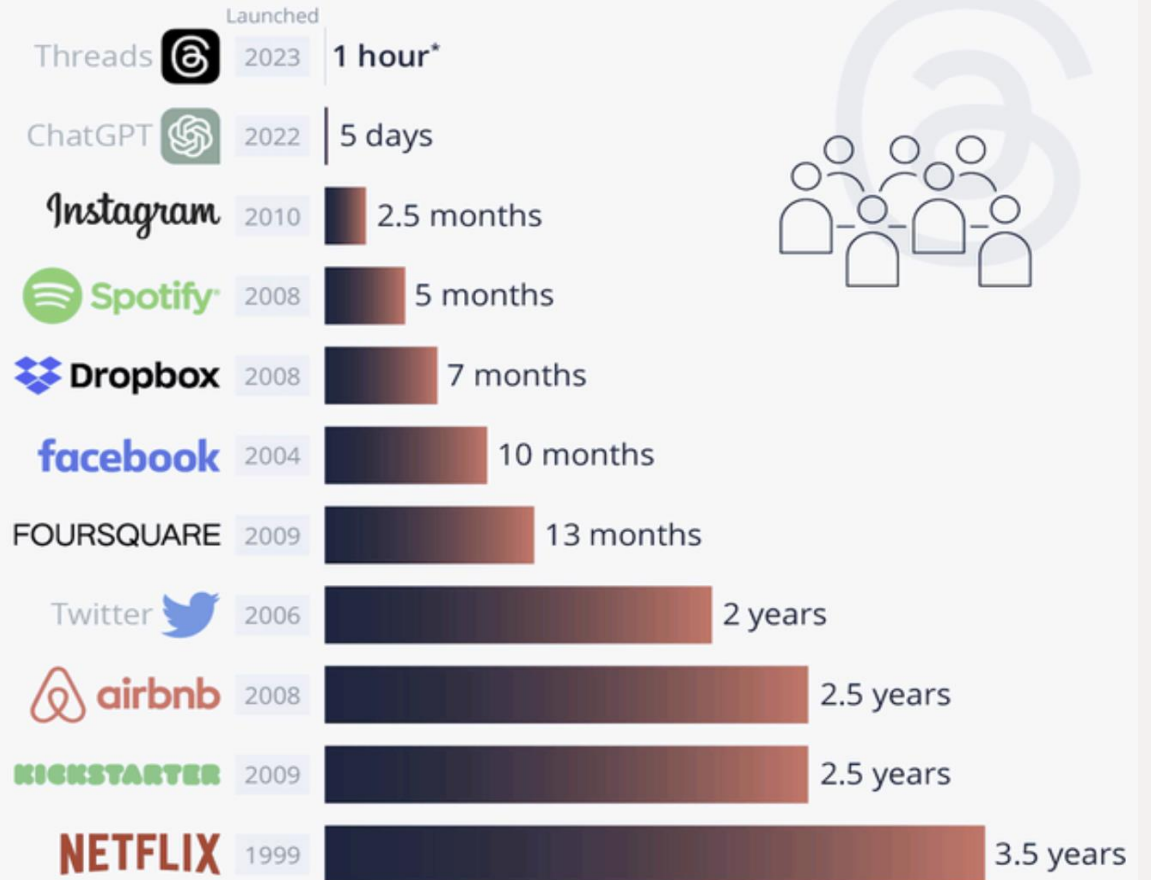
Defense



# ADOPTION

## Threads Shoots Past One Million User Mark at Lightning Speed

Time it took for selected online services to reach one million users



Refers to one million backers (Kickstarter), nights booked (Airbnb), downloads (Instagram/Foursquare)

\* Two million signups in two hours

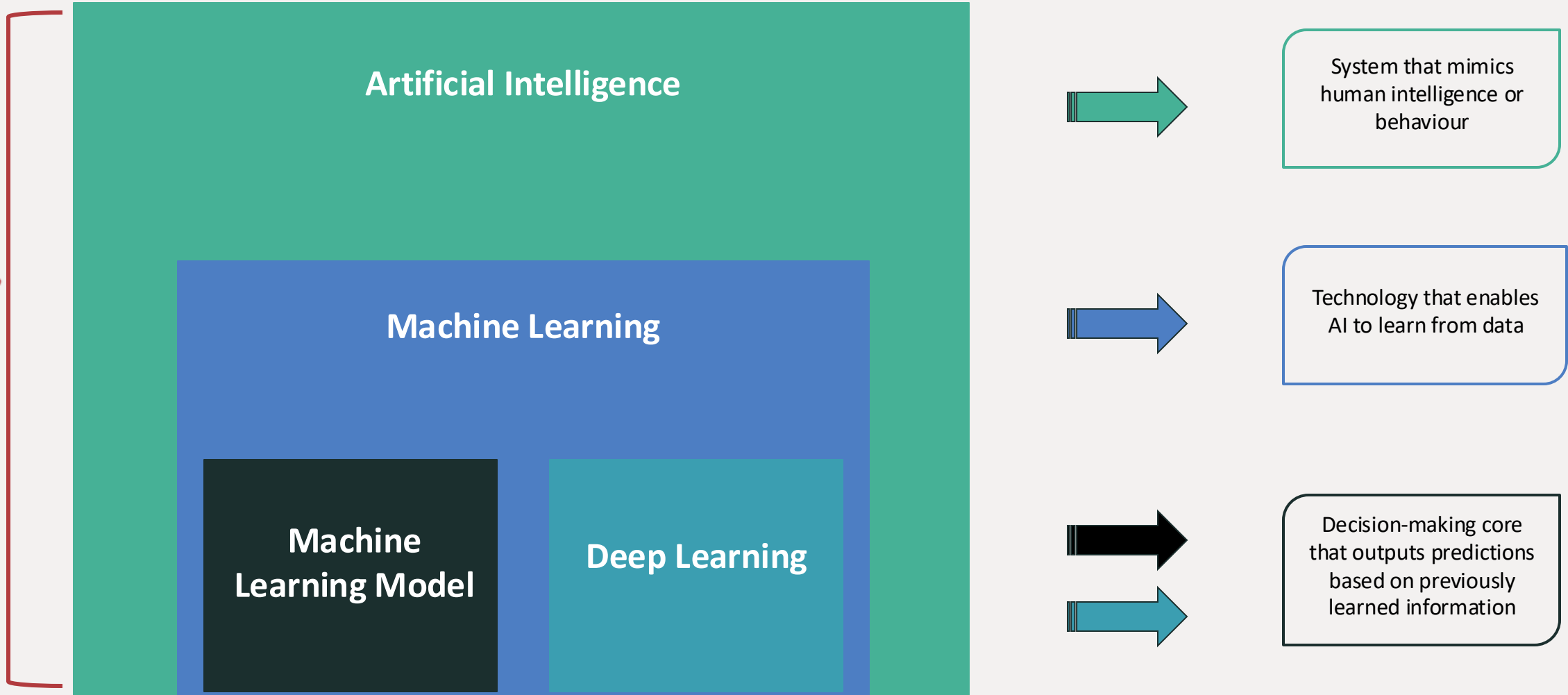
Source: Company announcements via Business Insider/LinkedIn



# LINGO



SkyNet



# WHEN



We understand exactly how AI is going to improve a key metric in our business (no “black box” claims).

We understand exactly how AI projects fit into our overall Innovation Strategy/Portfolio.

We have clearly identified the specific business problem we seek to solve, then used the appropriate AI to achieve that result.

We have a forum for regularly communicating what we are learning about the uses of AI across our organization.

We are providing licenses and training to critical numbers of people across our organization who will benefit from understanding AI.

We understand how AI will help us get information about key changes in our external environment and what we should do about them.

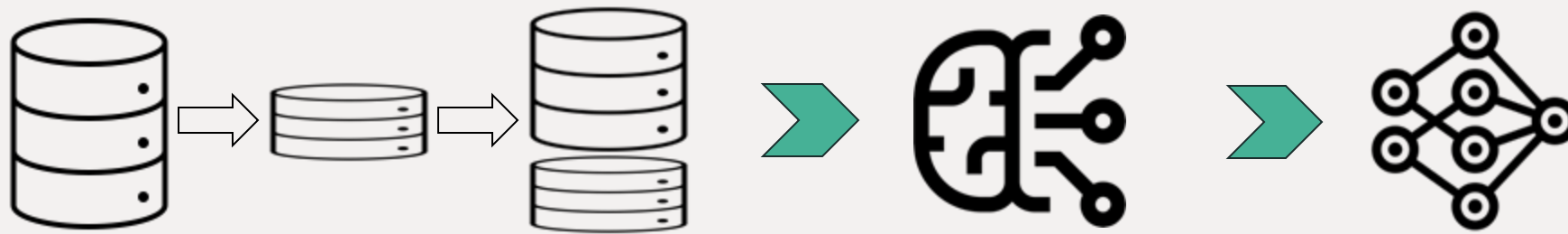
We understand how AI will help us improve the employee experience.

We have visibility into what AI projects are in the works and what their results are.

We have created a governance board of business leaders and AI experts who can evaluate how projects map to market and technical uncertainties.

We have confidence that people in strategic decision-making roles understand how AI will affect our business.

# HOW MACHINE LEARNING WORKS



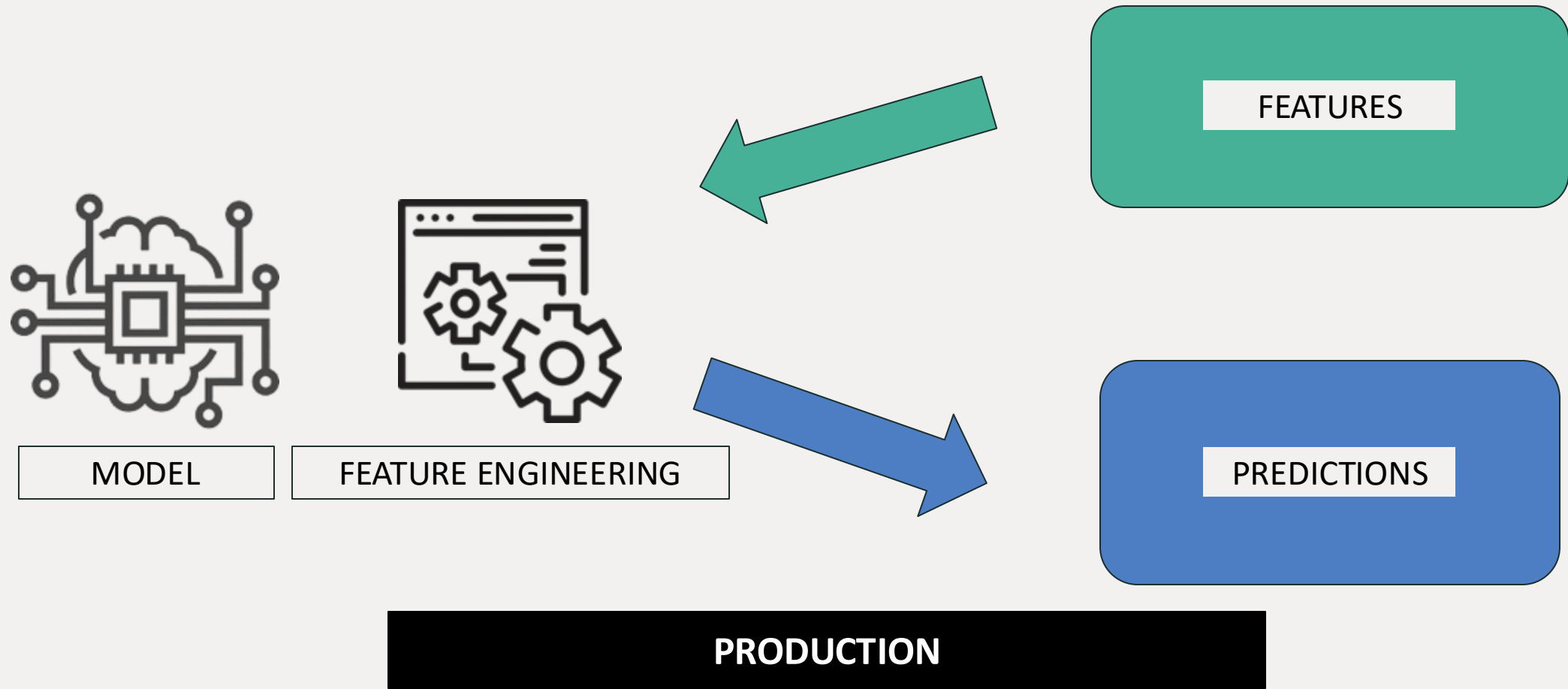
TRAINING DATA

MODEL TRAINING

TRAINED MODEL

**TRAINING**

# HOW MACHINE LEARNING WORKS (CONT)



# EXAMPLE MODEL

## Support Vector Model (SVM)

Model looking for the optimized hyperplane and widest streets

### Math

$$y = mx + b$$

y – optimal linear output

m – Weight

x – Input(s)

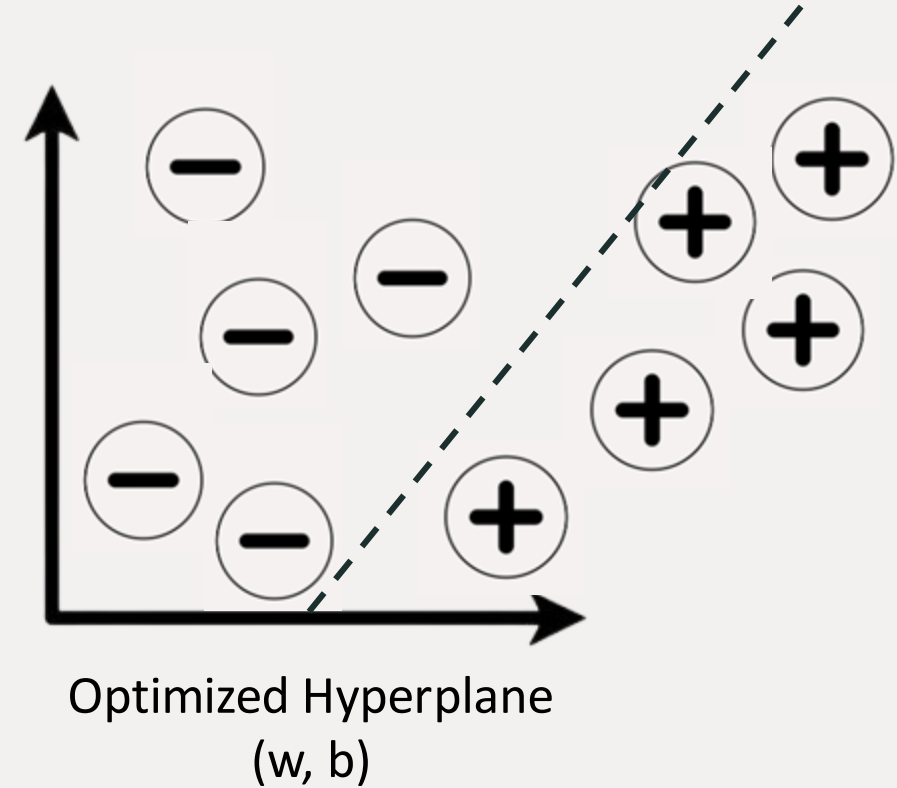
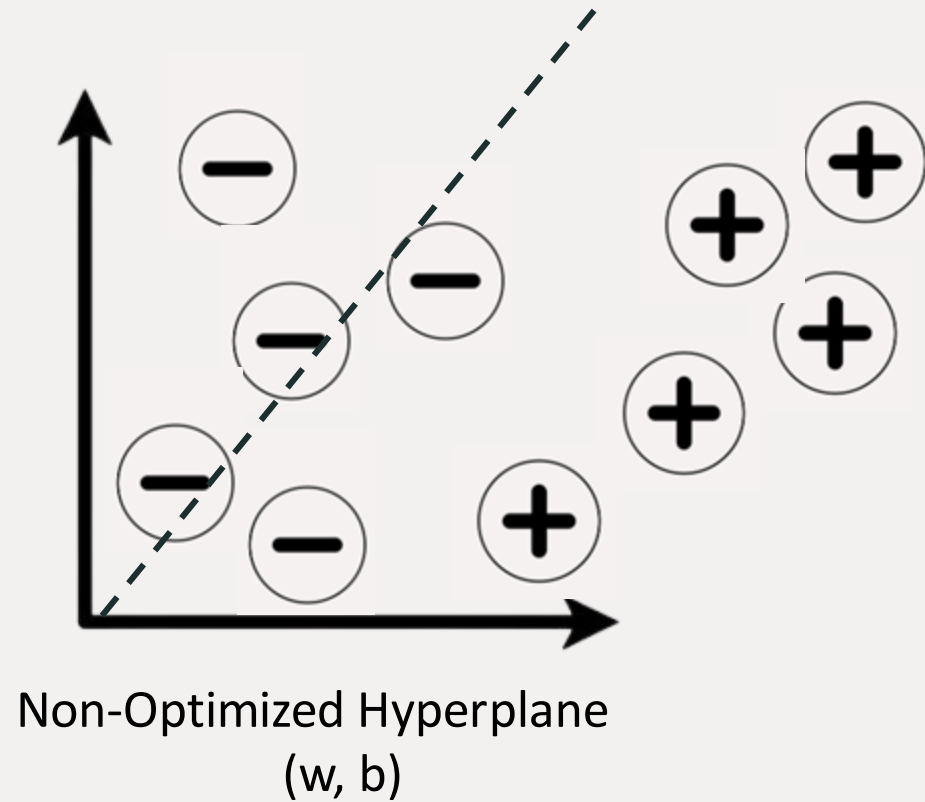
b – Bias

### Scenario

- Goal : Detect whether files are Malicious or Benign
- 10,000 sample training dataset
- Each File has two **Features**: Hash and Newness
- **Prediction/Output**: Malicious or Benign

# EXAMPLE MODEL (cont)

$$y = mx + b$$





# CONFUSION MATRIX

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

$$Precision = \frac{TP}{TP+FP}$$

$$Recall = \frac{TP}{TP+FN}$$

$$F1-Score = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$

	Predicted Malicious	Predicted Not Malicious
Actual Malicious	True Positive (TP)	False Negative (FN)
Actual Not Malicious	False Positive (FP)	True Negative (TN)

# Trust & Reach

- Usage License
- Unknown File
- Unknown File Type
- Vulnerability
- 3<sup>rd</sup> Party Data Broker
- Unnecessary Data Bias
- Data Poisoning
- Data & Model Genealogy
- Pre-Production Risk
- Red Team Testing

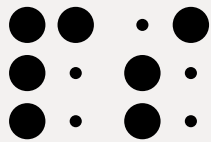
- Internally built
- Genealogy

- Genealogy
- New Unknown File

- Inference Activities
  - SRE
  - Security



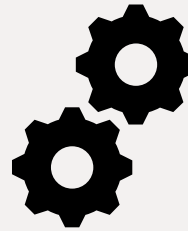
Download 3rd Party Model(s)



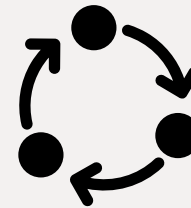
Model Registry



Data Collection



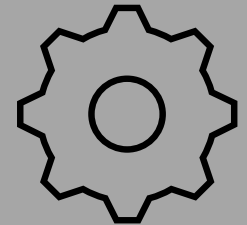
Training Process



Release Candidate



Data & Model Genealogy



Production

DESIGN & BUILD

OPERATIONS

# SECURITY

Adversarial ML
Defensive ML
Offensive ML



***Securing Organizations since 1761***

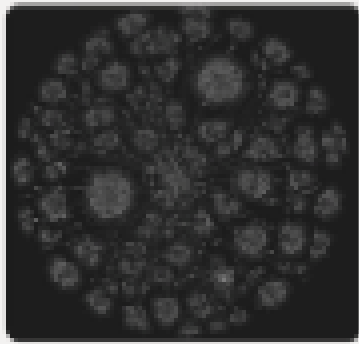
# AUTOMATED ATTACK TOOLS



MLsploit



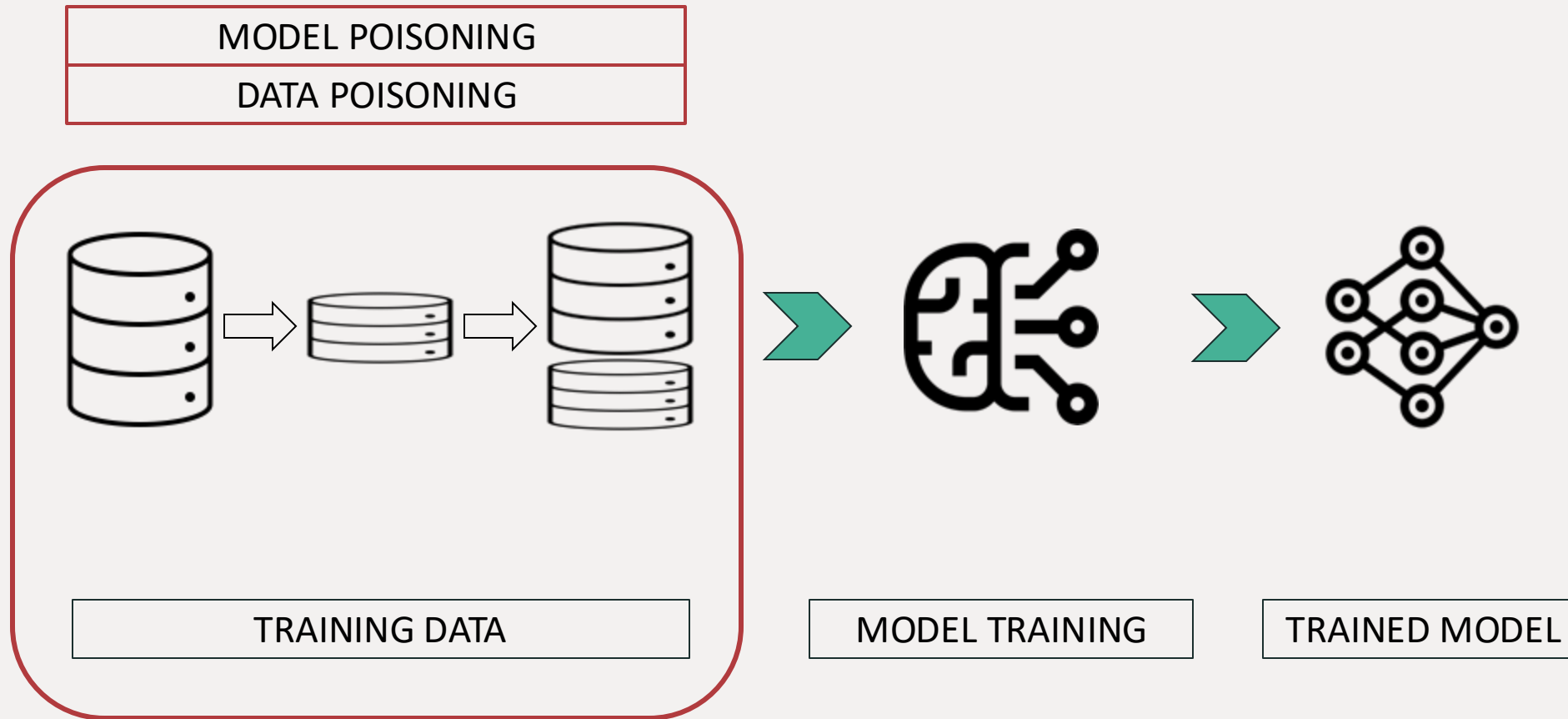
Protect AI



OffSecML Playbook

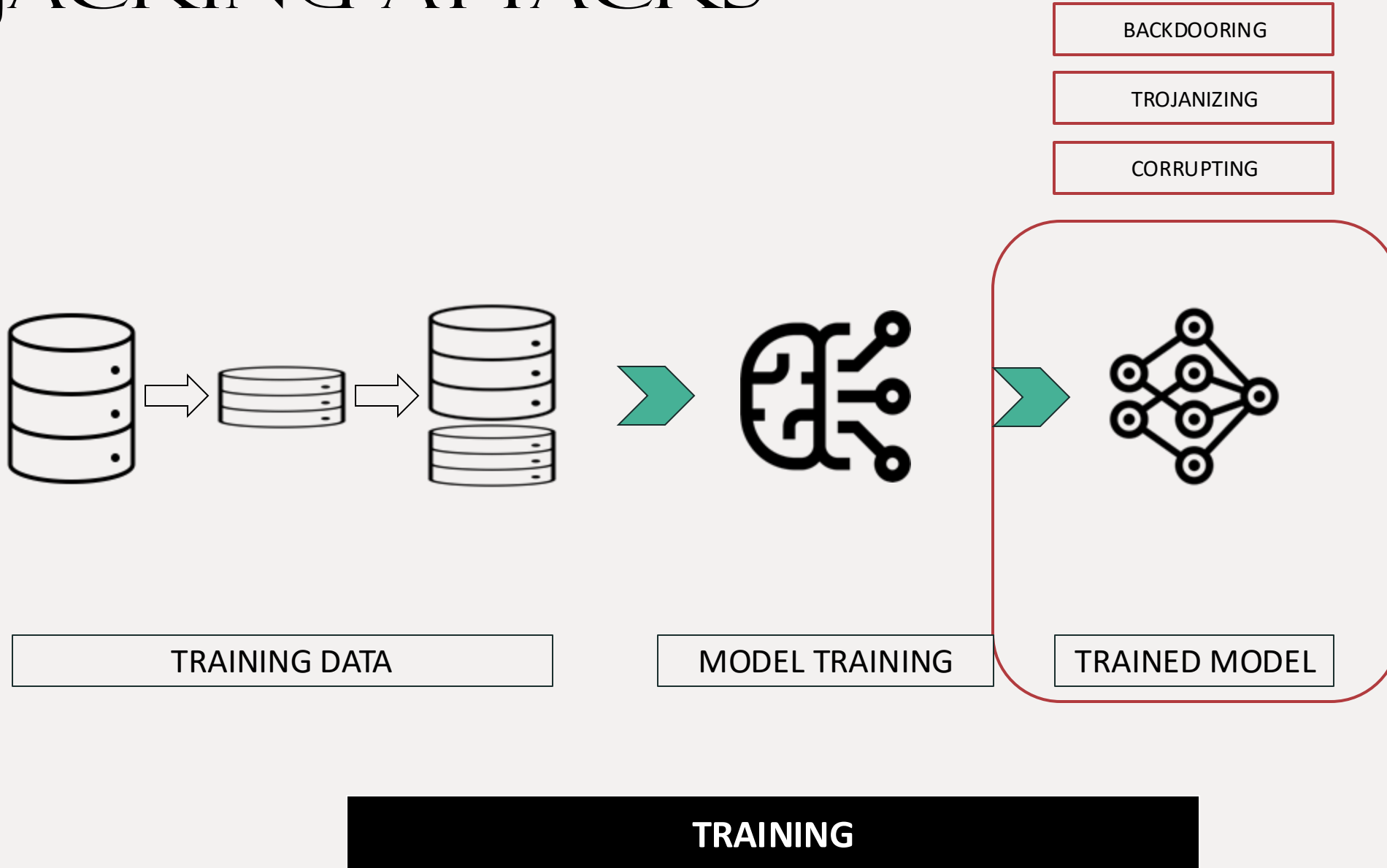


# POISONING ATTACKS



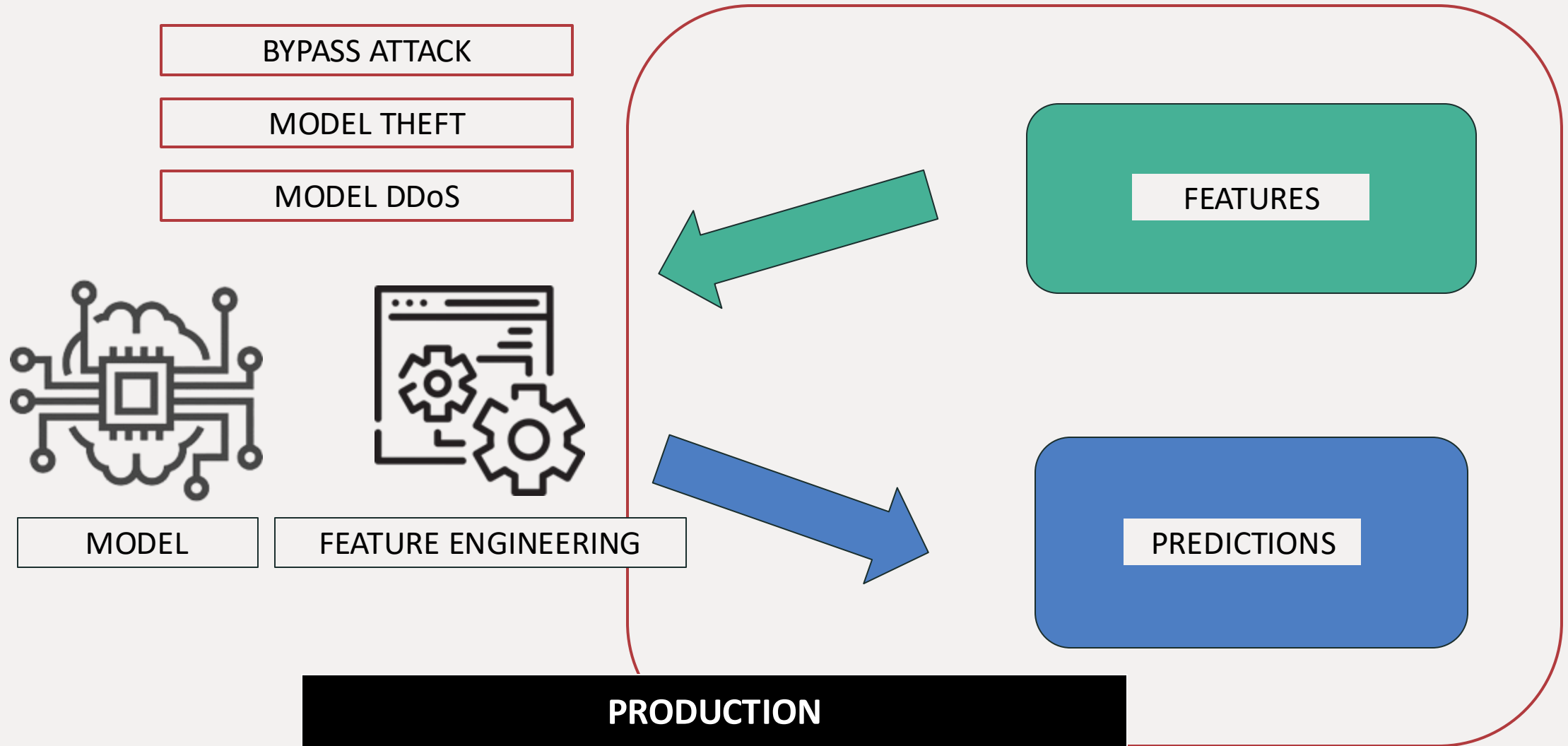
**TRAINING**

# HIJACKING ATTACKS





# INFERENCE ATTACKS



# QUICK RECAP



## Poisoning attacks

- Esp. relevant in online learning
- Can be intuitive and crowdsourced, or utilize botnets
- Bias, inaccuracy, disinformation
- Market / trends manipulation



## Inference attacks

- Post-deployment
- Require only UI/API access
- **Bypassing model** (e.g. detection, authorisation, authentication, etc.)
- **Stealing IP / PPI** (e.g. model, train. set)



## Model hijacking

- Esp. concerns publicly avail. models
- Can be used for delivery of **traditional malware** (e.g. in supply chain attacks)



## LLM prompt injection

- Bypassing chatbots' content filters
- Gaining access to restricted content
- Leaking sensitive data





WHAT PROBLEMS WILL  
CUSTOMERS SOLVE  
WITH AI?

HOW MUST OUR  
EXPERTISE EVOLVE?

WHAT ASSETS CAN WE  
DEVELOP TO STAY  
COMPETITIVE?



THANK YOU

*No AI was harmed in the creation of this presentation*