
Future of Cyber Security Enabled by AI

Dr. William Streilein

Cyber Analytics and Decision Systems

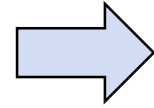


DISTRIBUTION STATEMENT A. Approved for public release. Distribution is unlimited.
This material is based upon work supported under Air Force Contract No. FA8721-05-C-0002
and/or FA8702-15-D-0001. Any opinions, findings, conclusions or recommendations expressed
in this material are those of the author(s) and do not necessarily reflect the views of the
U.S. Air Force. Delivered to the U.S. Government with Unlimited Rights, as defined in DFARS

Part 252.227-7013 or 7014 (Feb 2014). Notwithstanding any copyright notice, U.S. Government
rights in this work are defined by DFARS 252.227-7013 or DFARS 252.227-7014 as detailed
above. Use of this work other than as specifically authorized by the U.S. Government may violate
any copyrights that exist in this work.
© 2018 Massachusetts Institute of Technology.



Outline



- **Background**
 - Definition of Area
 - AI History Highlights
- **Lay-of-the-Land**
- **AI for Cyber Security**
- **Summary**



National Challenges and Role of AI

National Challenges

Role of AI in Augmenting Humans



Intelligent Systems and Autonomy



Information Superiority

Technological dominance in support of national security

Derive actionable intelligence by effective human-machine teaming



Massive amounts of structured and unstructured data

Leverage rapid advances in data conditioning, algorithms, and computing



Trust in intelligent machines (Robust AI)

Ascertain robustness

"We had better be quite sure that the purpose put into the machine is the purpose which we really desire"

Norbert Wiener, 1960



Operative AI Definition for the Study

Narrow AI:

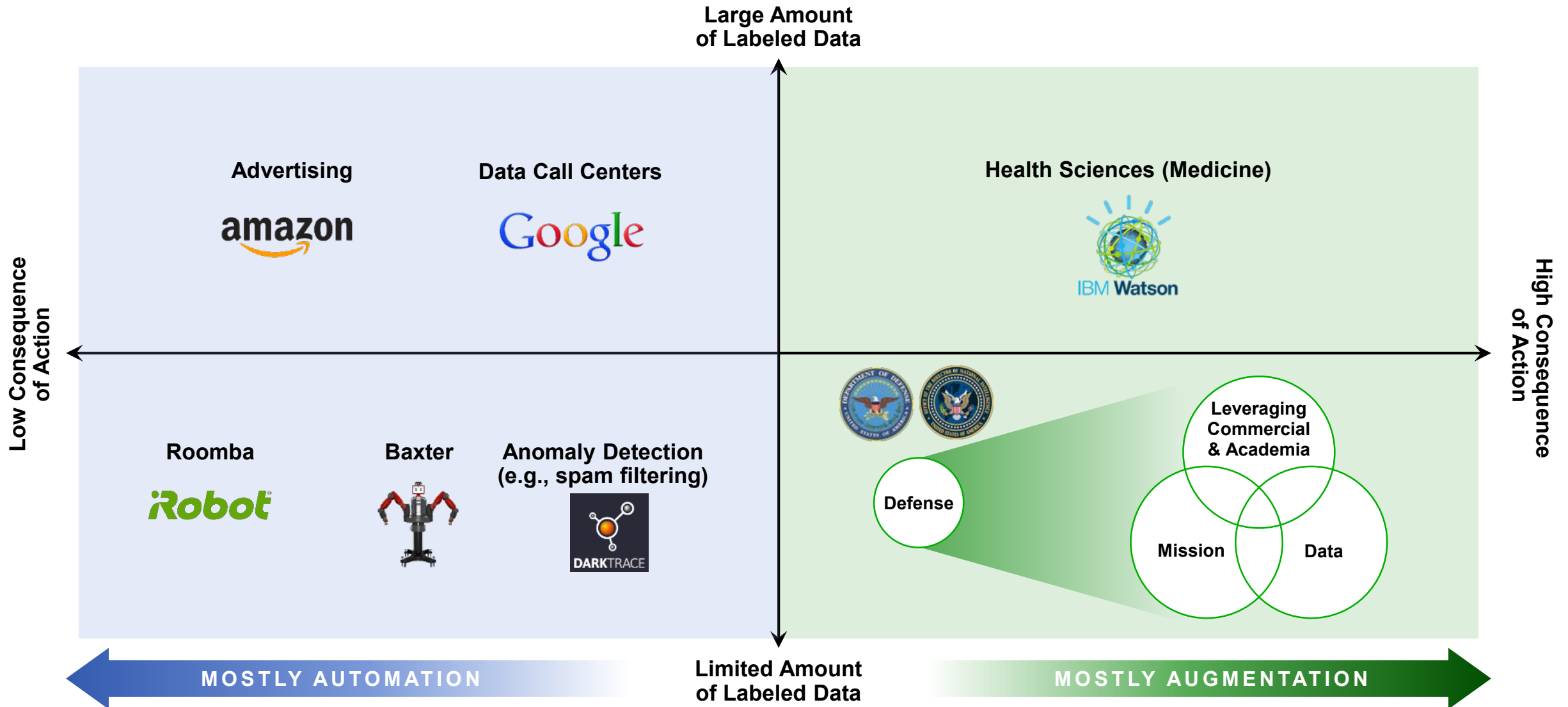
The theory and development of computer systems that perform tasks that augment human intelligence such as perceiving, learning, classifying, abstracting, reasoning, and/or acting

We will address: Narrow AI *not* General AI

* Definition adapted from Oxford dictionary and inputs from Prof. Patrick Winston (MIT) during his visit to MIT LL May 2017

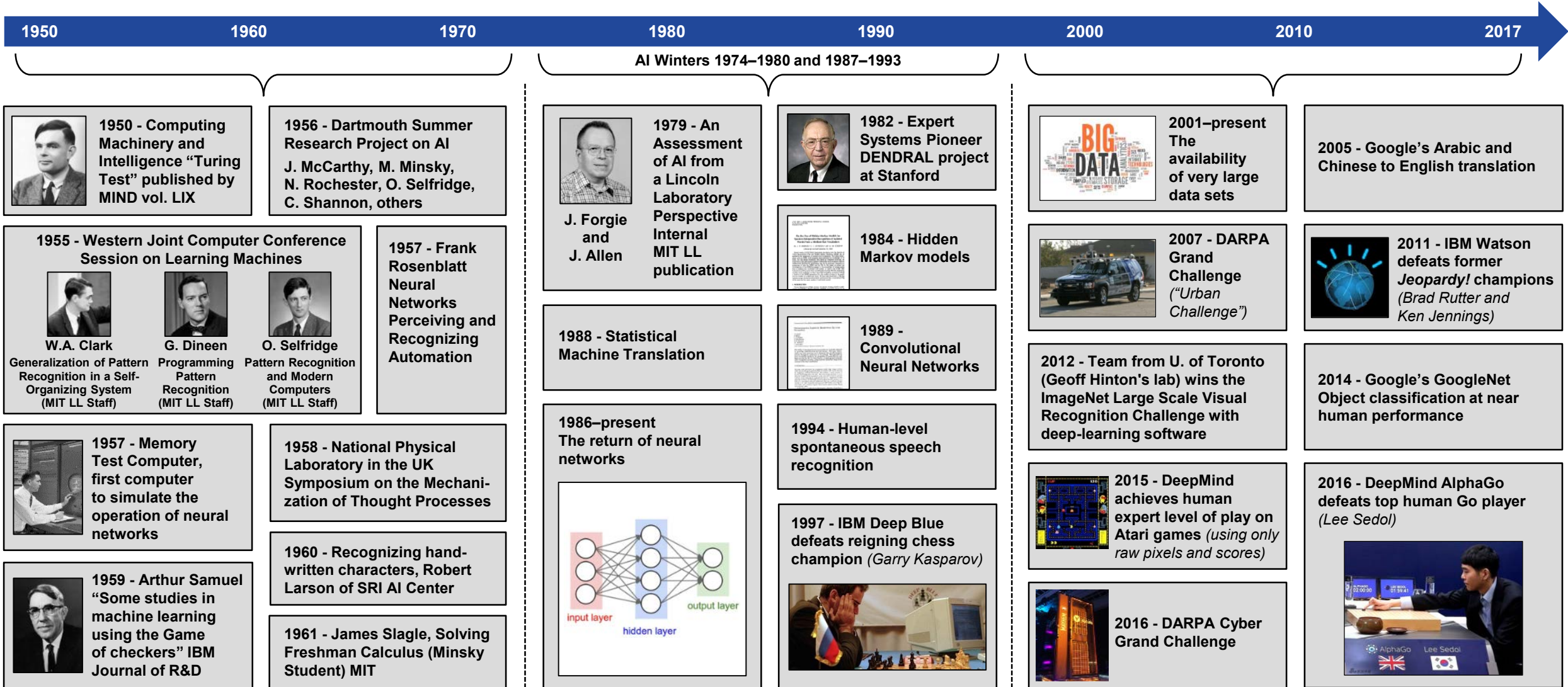


AI Domain of Impact





Select History of Artificial Intelligence





AI “Winters”

1974–1980	The first “AI winter”
1970s	Knowledge-based approaches
1980–88	Expert systems boom
1988–93	Expert systems bust; the second “AI winter”
1986	Neural networks return to popularity
1988	Pearl’s Probabilistic Reasoning in Intelligent Systems
1990	Backlash against symbolic systems; Brooks’ “nouvelle AI”
1995–present	Increasing specialization of the field Agent-based systems Machine learning everywhere Tackling general intelligence again?

The first AI winter 1974–1980

In the 70s, AI was subject to critiques and financial setbacks. AI researchers had failed to appreciate the difficulty of the problems they faced. **Their tremendous optimism had raised expectations impossibly high, and when the promised results failed to materialize, funding for AI disappeared.** At the same time, the field of connectionism (or neural nets) was shut down almost completely for 10 years by Marvin Minsky's devastating criticism of perceptron. Despite the difficulties with public perception of AI in the late 70s, new ideas were explored in logic programming, commonsense reasoning and many other areas.

Bust: the second AI winter 1987–1993

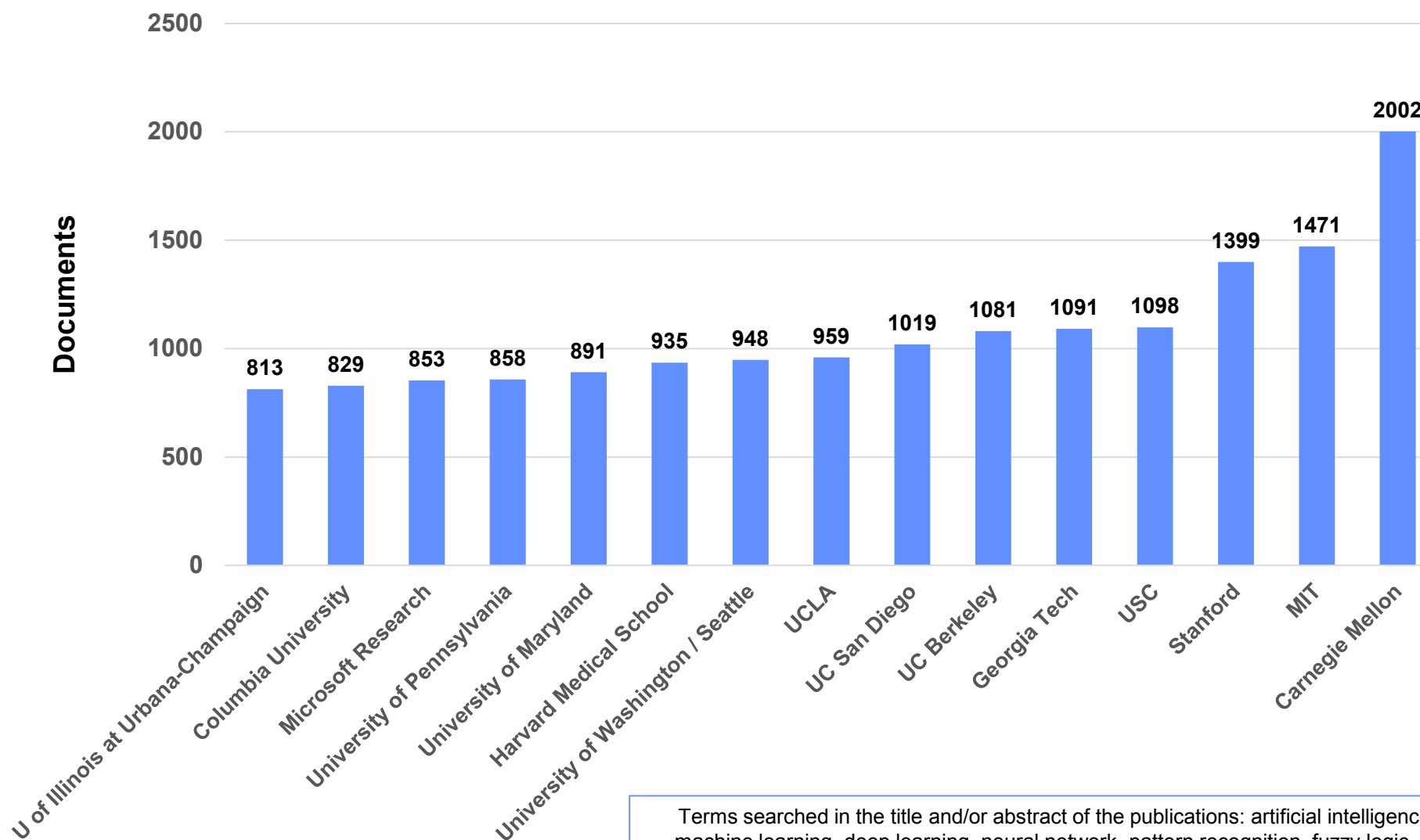
The business community's fascination with AI rose and fell in the 80s in the classic pattern of an economic bubble. The collapse was in the perception of AI by government agencies and investors – the field continued to make advances despite the criticism. Rodney Brooks and Hans Moravec, researchers from the related field of robotics, argued for an entirely new approach to artificial intelligence.

Source: UNC Computer Science

Source: Wikipedia, History of artificial intelligence



Top 15 Publishing Universities/Organization in the US (2011–Present)

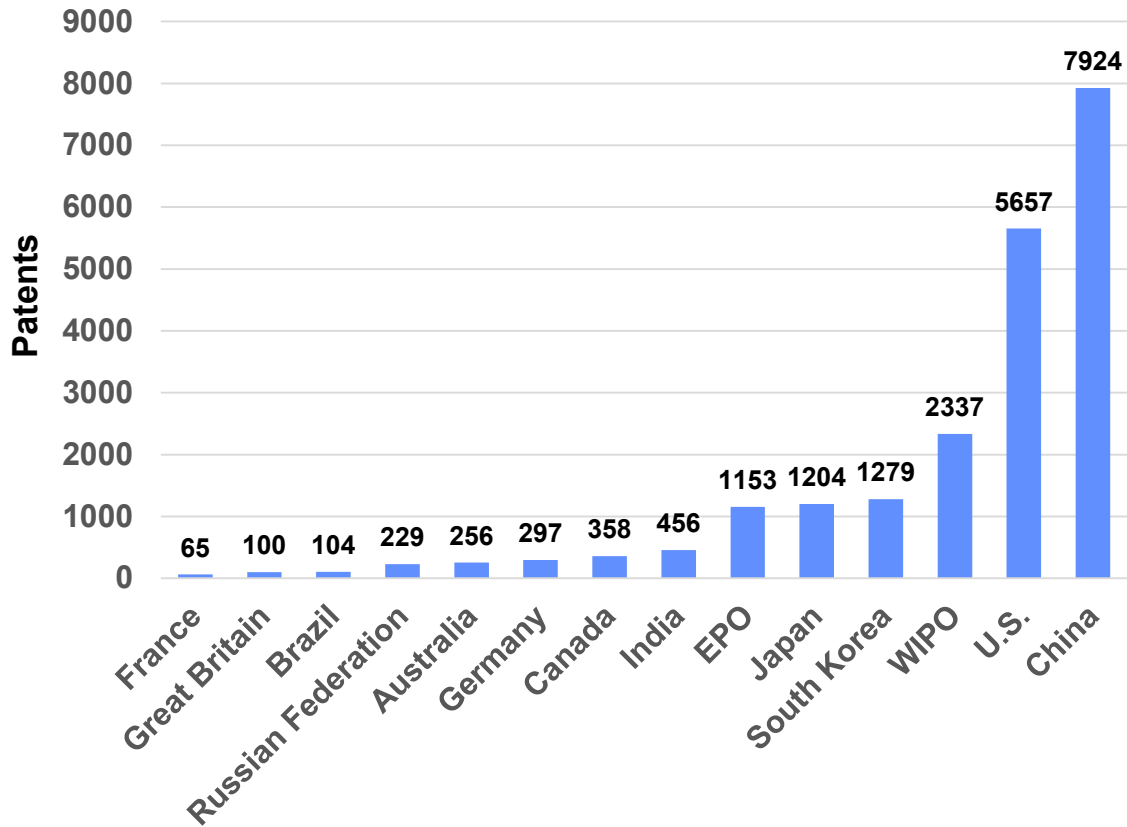


Terms searched in the title and/or abstract of the publications: artificial intelligence, cognitive computing, machine learning, deep learning, neural network, pattern recognition, fuzzy logic, support vector machine

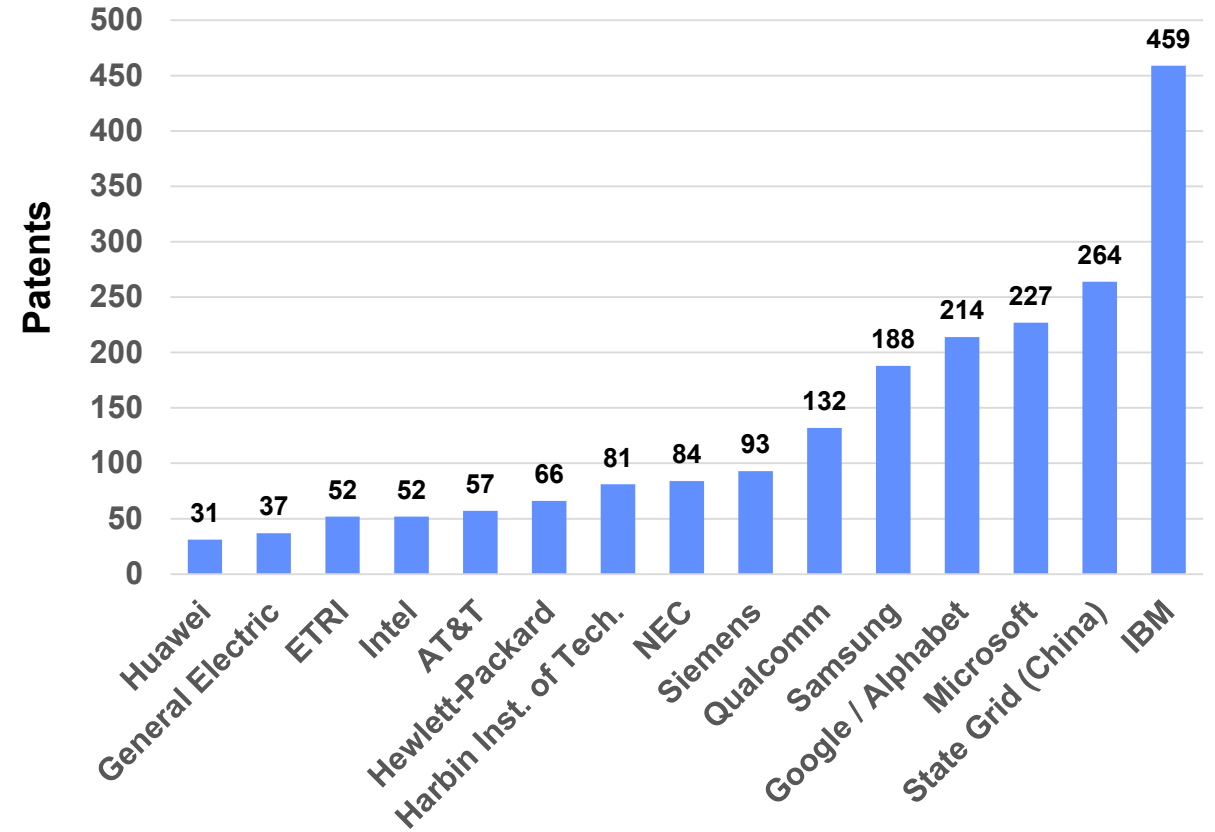


Top 14 Patent Holders in AI Per Country (2011–2016)

Top 14 Patent Holders in AI Per Country of Publication (2011–2016)



Top 15 Patent Holders in AI Per Patent Assignee (2011–2016)



Terms searched in the title and/or abstract of the patent record: artificial intelligence, cognitive computing, machine learning, deep learning, neural network, pattern recognition, fuzzy logic, support vector machine



China is Putting a Major Investment into AI



Chinese Government has indicated plans to invest \$150B over next few years

The Economist (July 2017)

In 2012–16 Chinese AI firms received \$2.6B in funding, according to the Wuzhen Institute, a think-tank

China Next Generation AI Development Plan (July 2017)¹

By 2020 China will have established initial AI technology standards, service systems, and industrial ecological system chains with the scale of AI's core industry exceeding \$22.6B, and exceeding \$150B as driven by the scale of related industries

MIT Technology Review (November 2017)

China's goal is "to have major breakthroughs in AI by 2025, and to be the envy of the world by 2030"²

DoD R&D spending is a fraction of nation states – we are losing ground on patents and publications

¹ <https://www.newamerica.org/cybersecurity-initiative/blog/chinas-plan-lead-ai-purpose-prospects-and-problems/>
² The Artificial Intelligence Issue, China's AI Awakening, MIT Technology Review, Nov-Dec 2017



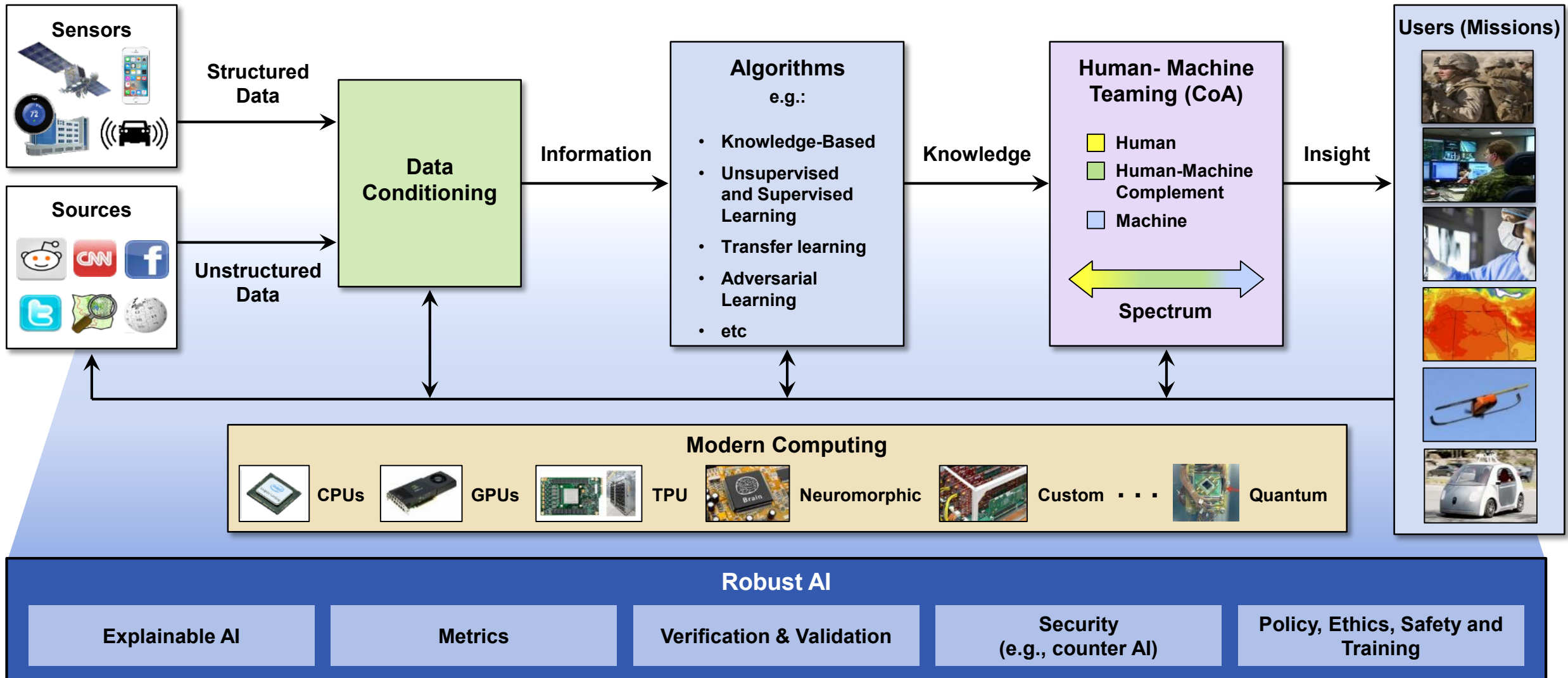
Outline

- **Background**
- ➔ • **Lay-of-the-Land**
 - **AI Canonical Architecture**
 - **Summary of Study Outreach**

 - **AI for Cyber Security**
- **Summary**



AI Canonical Architecture





Four Components of Machine Learning Solutions

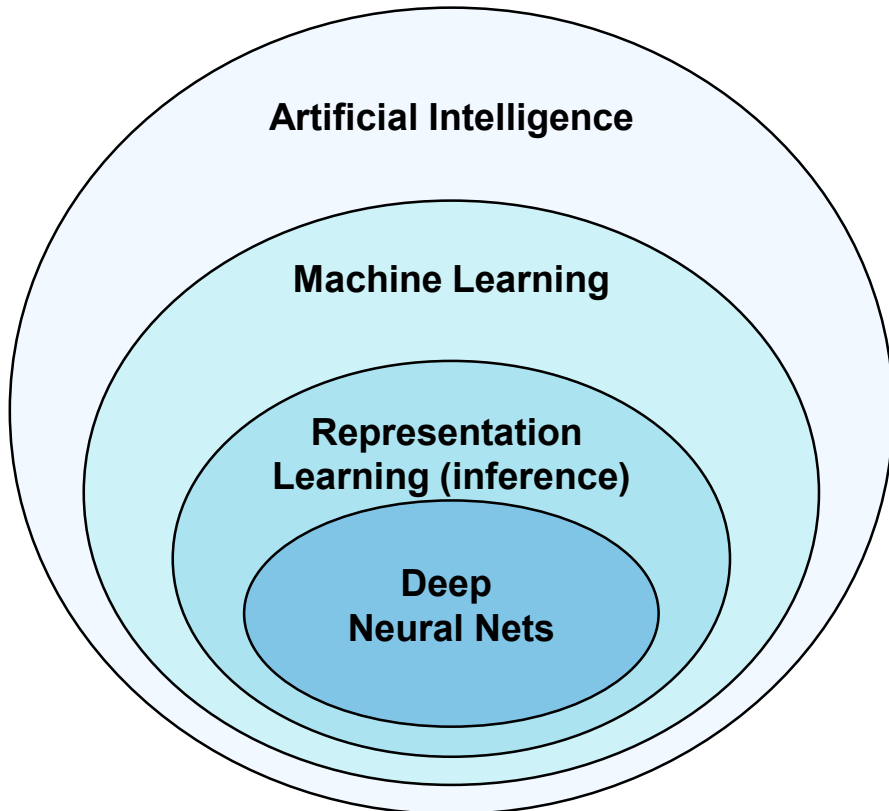
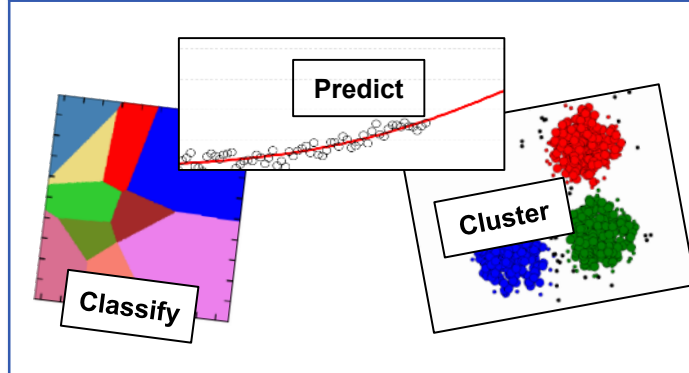


Image Adapted From: "Deep Learning"
I. Goodfellow, et.al., 2016 MIT Press

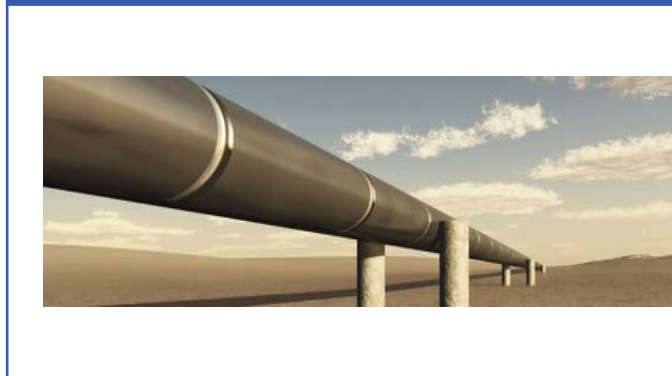
1. Define a Problem



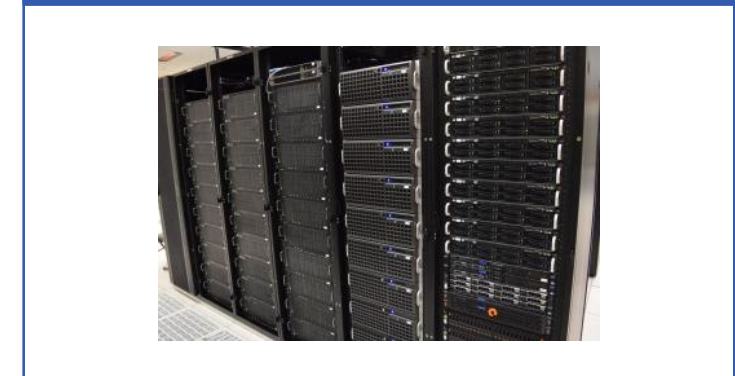
2. Gather Data



3. Create Full Train/Test Solution Pipelines



4. Provide Computation that Makes the Solution Feasible





National Security Specific Concerns Compared To Commercial Sector



Commercial Sector	National Security
High dimensionality	High dimensionality
Large volume	Large volume
Known truth / continual development	Unknown truth / not tolerant to errors
Mild consequences of decisions	Large consequences of decisions
Past is representative of future	Past does not always represent future
Competitive environment	Adversarial environment
Mostly consumer users	Today requires sophisticated users
Explainability is not the largest issue	Trust / explainability is core
Quantifiable success (\$\$)	Harder to measure success

AI will be a technological enabler (i.e., data and algorithm warfare) against: radical extremists, terrorists, and peer nations to defend our homeland and abroad



Spectrum of Commercial Organizations in the Machine Intelligence Field

ENTERPRISE FUNCTIONS

CUSTOMER SUPPORT

DigitalGenius Kasisto
ELOQUENT Wiseo
ACTIONIQ zendesk
Preact CLARABRIDGE

SALES

collective[i] sense
fuse|machines AVISO
salesforce INSIDE SALES .COM
Zensight clari

MARKETING

MINTIGO Lattice RADIUS
LiftIgniter AIRPR MOTIVA
brightfunnel megai retention
[PERSADO] COGNICOR

SECURITY

CYLANCE DARKTRACE
ZIMPERIUM dsepinstruct
Sentinel DEMISTO
graphistry drawbridge
SignalSense AppZen

RECRUITING

textio entelo
Wade & Wendy hi
unifive SpringRole
GIGSTER HireVue

PERSONAL

amazon alexa
Cortana Allo
facebook
Siri Replika

AGENTS

PROFESSIONAL

butter.ai pogo SKIPFLAG
@clara x.ai slack
talla Zoom sudo

INDUSTRIES

AGRICULTURE

BLUEØRIVER mavrx
tule TRACE GENOMICS Pivot Bio
TerraAvion AGRI-DATA
Descartes Labs ud abundant

EDUCATION

KNEWTON volley
gradescope
CTI coursera
UDACITY a|t|school

INVESTMENT

Bloomberg sentient
iSENTIUM KENSHO
alphasense Dataminr
CEREBELLUM CAPITAL Quandl

LEGAL

blueJ BEAGLE
Everlaw RAVEL
Seal ROSS
LEGAL ROBOT

LOGISTICS

NAUTO Acerta
PRETECKT clearmetal
Routific
MARBLE PITSTOP

MATERIALS

zymergen Citrine
Eigen Innovations
SIGHT MACHINE
GINKGO BIOWORKS nanotronics
CALCULARIO

RETAIL FINANCE

TALA zest finance
Lendo earnest
affirm MIRADOR
wealthfront Betterment

AUTONOMOUS SYSTEMS

GROUND NAVIGATION

drive.ai AdasWorks
ZOOX MOBILEYE
UBER Google TESLA
nuTonomy Auro Robotics

AERIAL

SKYDIO SHIELD AI
Airware DJI LILY
DroneDeploy
pilot.ai SKYCATCH

INDUSTRIAL

JAYBRIDGE OSARO
CLEARPATH fetch
KINDRED
HARVEST rethink robotics

HEALTHCARE

PATIENT

PULSE CareSkore
ZEPHYR HEALTH Watson Health
Oncora SENTRIAN
Atomwise Numerate

IMAGE

BUTTERFLY 3SCAN
ARTERYS enlitic
BAYLABS imagia
Google DeepMind

BIOLOGICAL

iCarbonX color GRAIL
deep genomics RECURSION
LUMINIST Numerate
Atomwise verily WHOLE BIOME

ENTERPRISE INTELLIGENCE

VISUAL

Orbital Insight planet
clarifai DEEP VISION
cortica Iqocean
SPACE_KNOW Captricity
netra deepomatic

AUDIO

Gridspace TalkiQ
nexidia twilio
CAPIO Expect Labs
Clover Mobvoi
Curious.AI popUP archive

SENSOR

PREDIX IoT MAANA
Sentiai PLANET OS
UPTAKE IMUBIT Preferred Networks
thingworx KONUX Alluvium

INTERNAL DATA

PRIMER IBM WATSON
Gycomp Palantir ARIMO
Alation Sapho Outlier
Digital Reasoning

MARKET

mattermark Quid
DataFox PREMISE
Bottlenose CB INSIGHTS
enigma Tracxn predata

TECHNOLOGY STACK

AGENT ENABLERS

OCTANE.AI howdy Maluba KITT.AI
OpenAI Gym Kasisto AUTOMAT
semantic

DATA SCIENCE

DOMINO SPARKBEYOND rapidminer
kaggle DataRobot yhat AYASDI
data iku seldon yseop bigml

MACHINE LEARNING

CognitiveScale GoogleML context relevant
Gycomp HyperScience nara logics minds.ai H2O.ai
SCALED INFERENCE sparkcognition loop GEOMETRIC INTELLIGENCE
deepsense.io reactive skyminr bonsai

NATURAL LANGUAGE

agolo FLYLIEN LEXALYTICS
Narrative Science loop spaCy LUMINOSO
cortical.io MonkeyLearn

DEVELOPMENT

SIGOPT HyperOpt fuzzyio pkite
rainforest lobe Anodot
Signifai LAYER bonsai

DATA CAPTURE

CrowdFlower diffbot CrowdAI Import
Paxata DATASIFT amazon mechanicalturk enigma
WorkFusion DATALOGUE TRIFACTA parsehub

OPEN SOURCE LIBRARIES

Keras Chainer CNTK TensorFlow Caffe
H2O DEEPLARNING4J theano torch
DSSTNE Scikit-learn AzureML neon
MXNet DMTK Spark PaddlePaddle WEKA

HARDWARE

KNUPATH TENSTORRENT Cirrascale
NVIDIA intel nervana Movidius
tensilica GoogleTPU 10th Labs ualcomm
Cerebras Isosemi

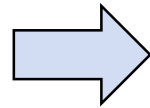
RESEARCH

OpenAI prabense ELEMENT vicarious
KNOGGIN Numenta Kimera Systems Cogitai



Outline

- **Background**
- **Lay-of-the-Land**
 - **AI Canonical Architecture**
 - **Summary of Study Outreach**
- **AI for Cyber Security**
- **Summary**





Government Organizations Study Outreach

DoD



Intelligence Community



IARPA



Under Secretary of Defense for Intelligence (USDI)

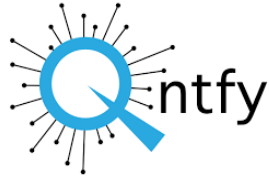




Defense Contractors, Commercial, Peers, and AI Centers Study Outreach

Defense Industrial Base

charles river analytics



Commercial



Research at Google



Peers

CMU SEI



JHU HLTCOE



ARGONNE



OAK RIDGE



NASA



PNNL



LIVERMORE



LOS ALAMOS



SANDIA



MITRE



AI Centers



USC CENTER FOR ARTIFICIAL INTELLIGENCE IN SOCIETY





Academia and MIT Study Outreach

**MIT School
of Engineering**



**Anantha
Chandrakasan**

**MIT
CSAIL**



**Srinivas
Devadas**

**MIT
CSAIL**



**Jim
Glass**

**MIT
CSAIL**



**Patrick
Winston**

**MIT Brain and
Cog Sciences**



**Vikash
Mansinghka**

**MIT
Media Lab**



**Sandy
Pentland**

**MIT
LIDS**



**Devavrat
Shah**

**MIT
IDSS**



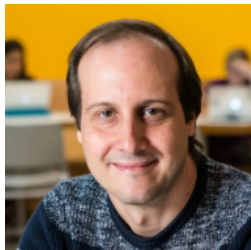
**Suvrit
Sra**

**MIT-IBM Watson
AI Lab**



**Aude
Oliva**

**MIT-IBM Watson
AI Lab**



**Antonio
Torralba**

**Florida
International
University**



**Mark
Finlayson**

**Northeastern
University**



**David
Kaeli**

**Boston
University**



**Michel
Kinsy**

**Ohio State
University**



**Srinivas
Parthasarathy**

**University
of Michigan**



**Arunesh
Sinha**

**University
of Southern
California**



**Milind
Tambe**



Outline

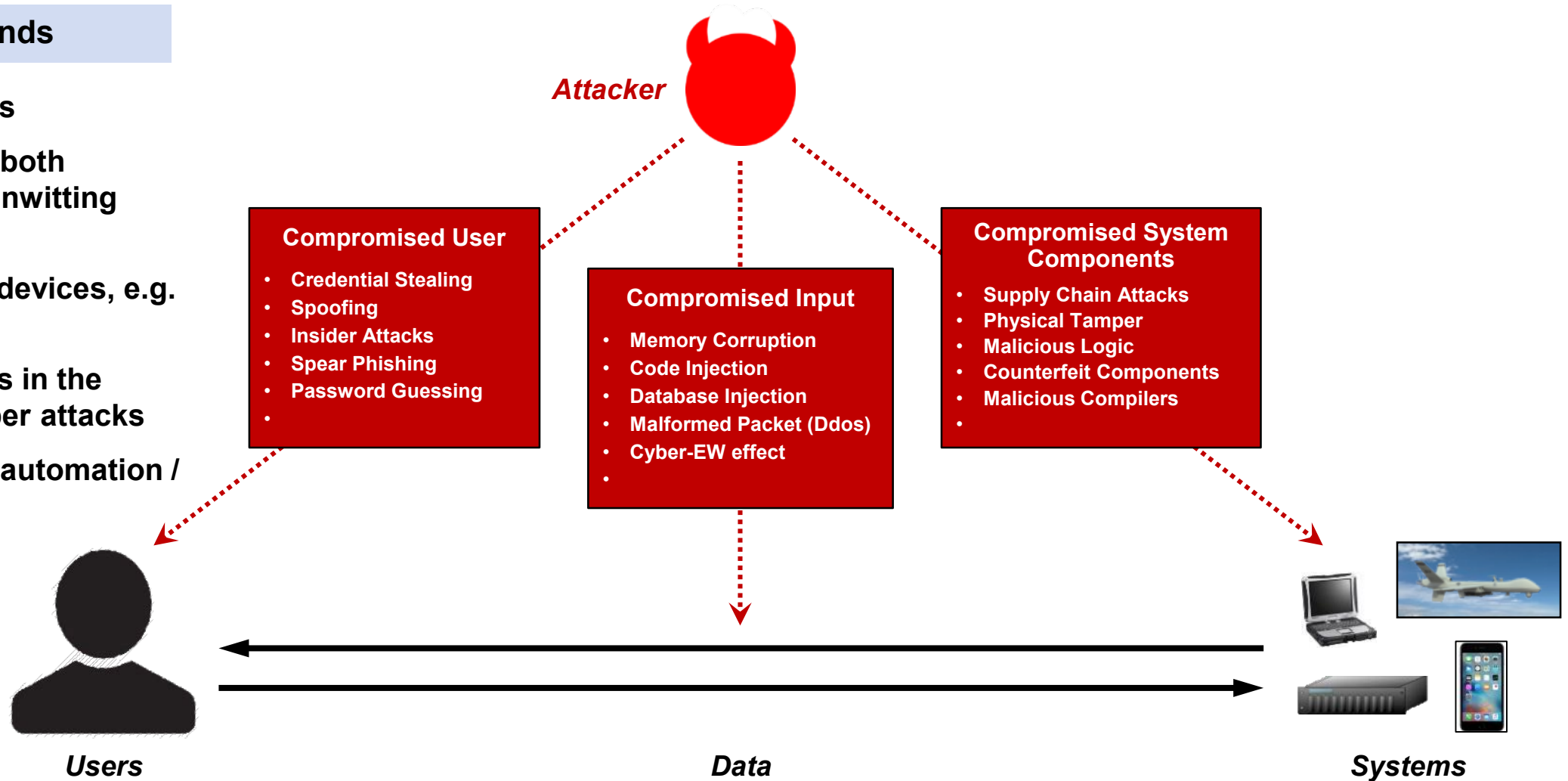
- **Background**
- **Lay-of-the-Land**
- ➔ • **AI for Cyber Security**
 - **Background**
 - **Findings and Recommendations**
- **Summary**



Cyber Security: Critical Threat Surfaces

Global Trends

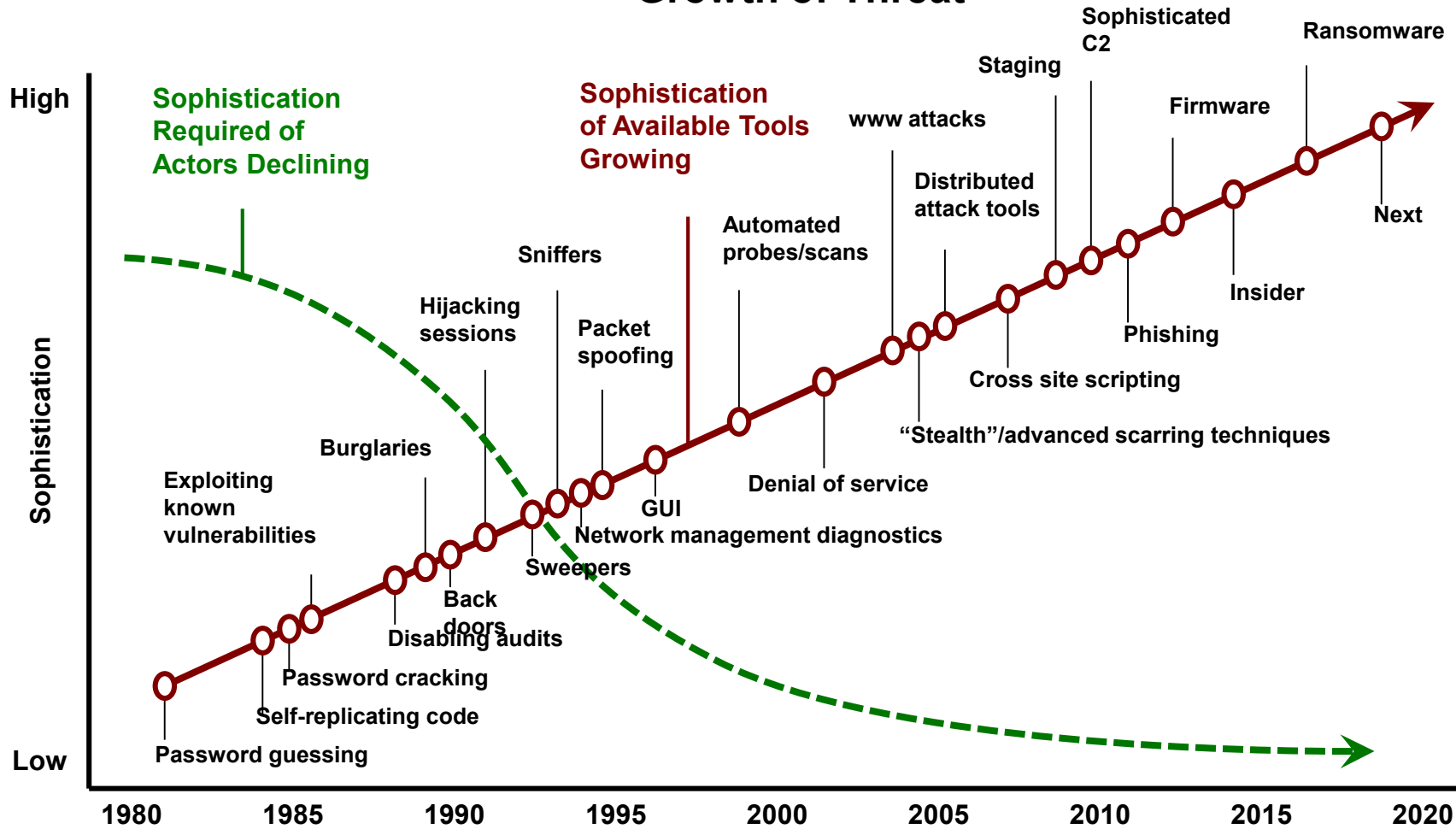
- Vulnerable users
- Insider threats (both malicious and unwitting insiders)
- Proliferation of devices, e.g. IoT
- Mission success in the presence of cyber attacks
- Attacker use of automation / AI





Global Trend: Sophisticated Attacks More Easily Accomplished with Automation

Growth of Threat



NOTEWORTHY FACTS

- 250K new malware programs are registered each day
- There were 357M new email malware variants in 2016 - 36% more new variants than in 2014.
- There were 463M new variants of ransomware in 2016 - 36% more new variants than in 2015.
- 99 days to detect compromise - adversary gains access in 3
- Internet of Things and Cloud are hot targets (e.g. Mirai botnet) – 2 min to compromise
- Projected cyber attack costs in 2019: \$2.1T



The Cyber Battleground

Offense Stages



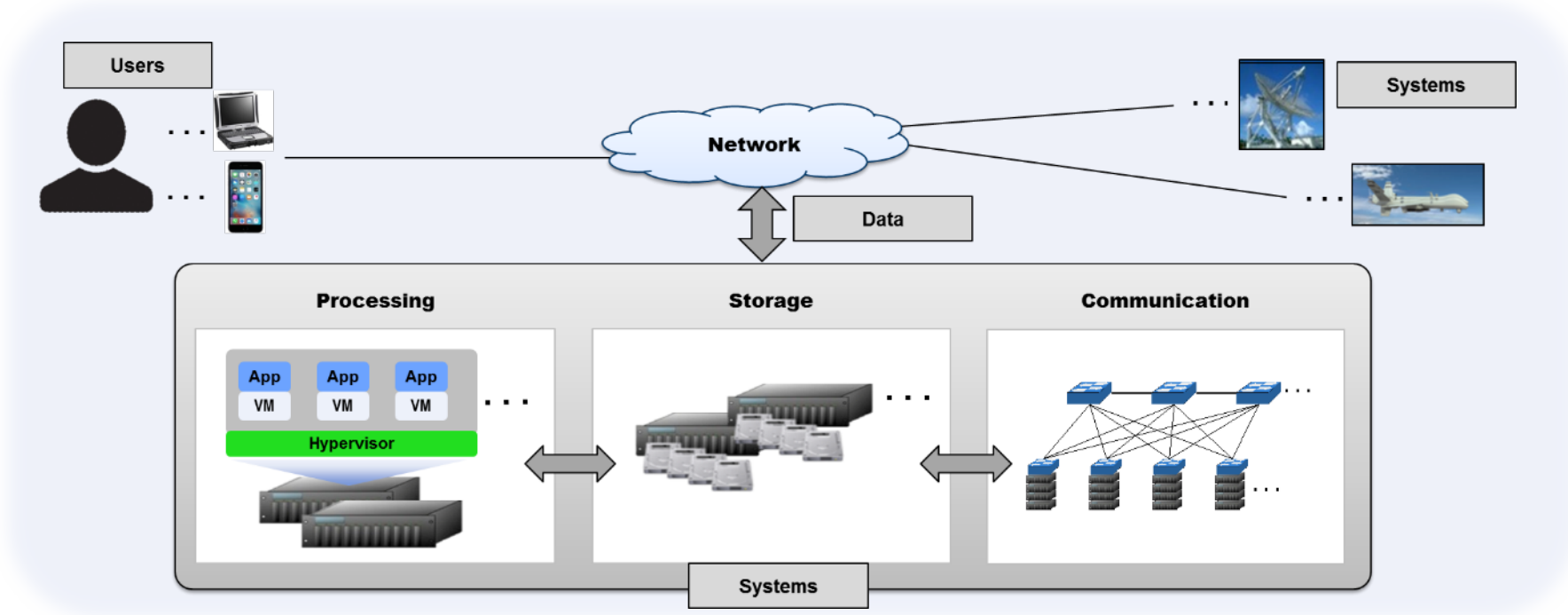
Impact

Know the target

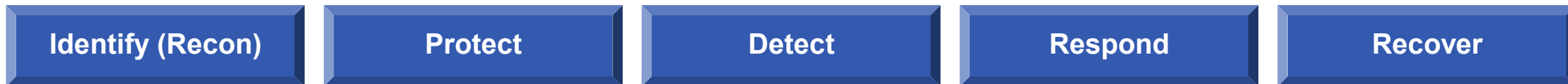
Enable attack process

Support persistence

Attack effectiveness



Defense Stages



Impact

Focused defense

Deflect attacks

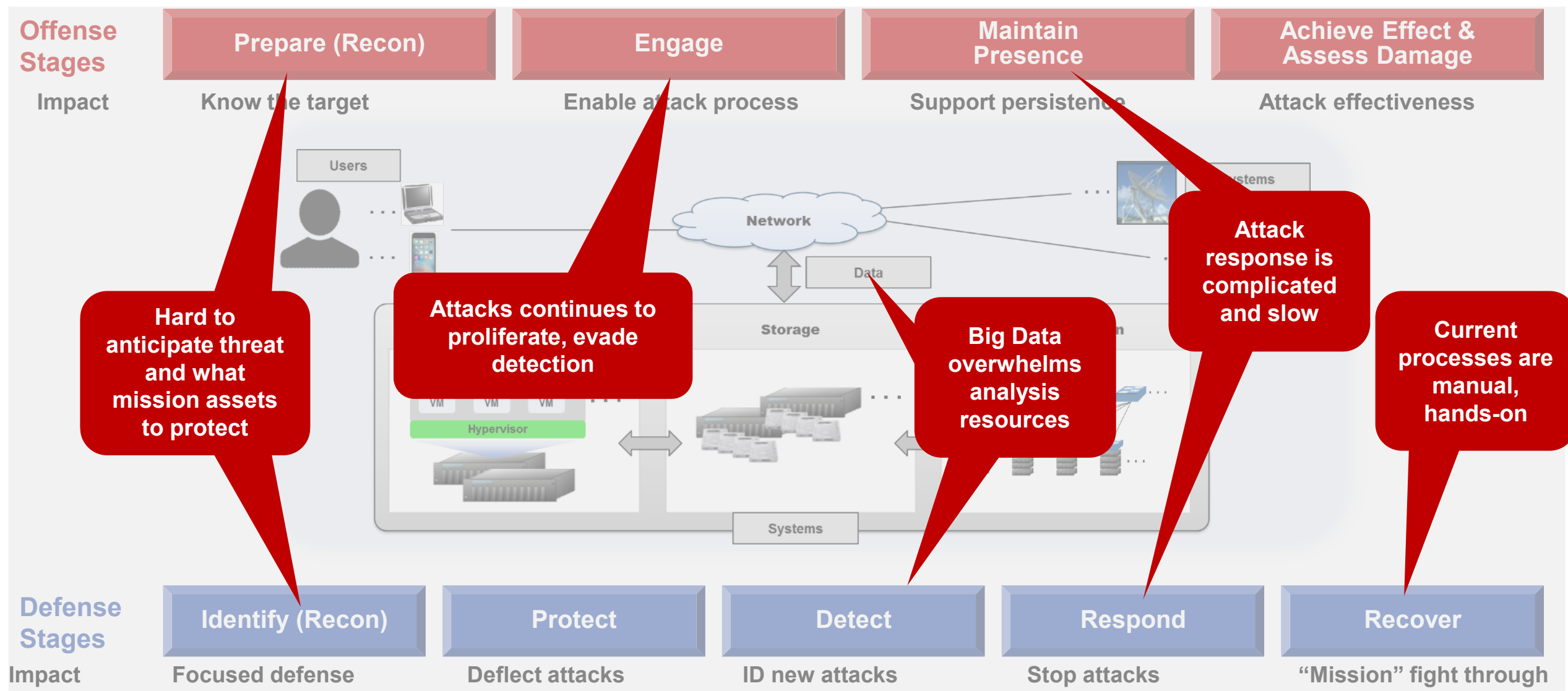
ID new attacks

Stop attacks

“Mission” fight through



Major Challenges to Cyber Security





AI Winter for Cyber

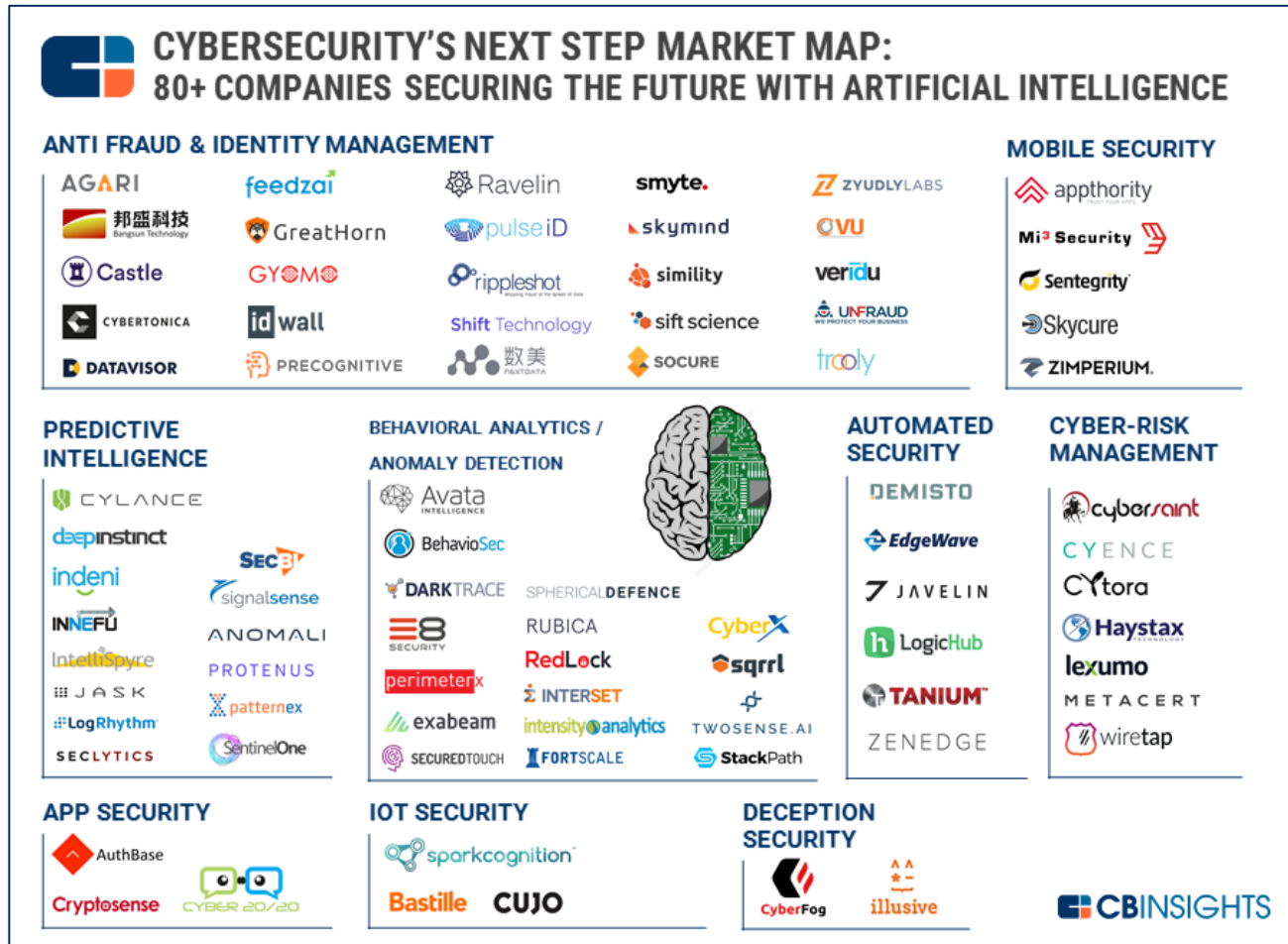
- “Outside the Closed World: On Using Machine Learning for Network Intrusion Detection”, (Sommer, Paxson, 2010)
- Described significant challenges with applying machine learning (e.g. AI) to cyber attack intrusion detection
- Caused ‘AI for Cyber Winter’ that forced people to abandon AI-based approaches in cyber security altogether
- Paxson later acknowledged impact was not as intended (AICS16, 2016)
 - Only applied to intrusion detection
 - Other cyber security aspects are amenable to AI



Sommer, Paxson paper shifted cyber security research focus from AI to secure methods



Cyber AI Start-up Landscape

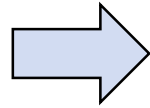


However, the market is now flush with companies leveraging AI for cyber



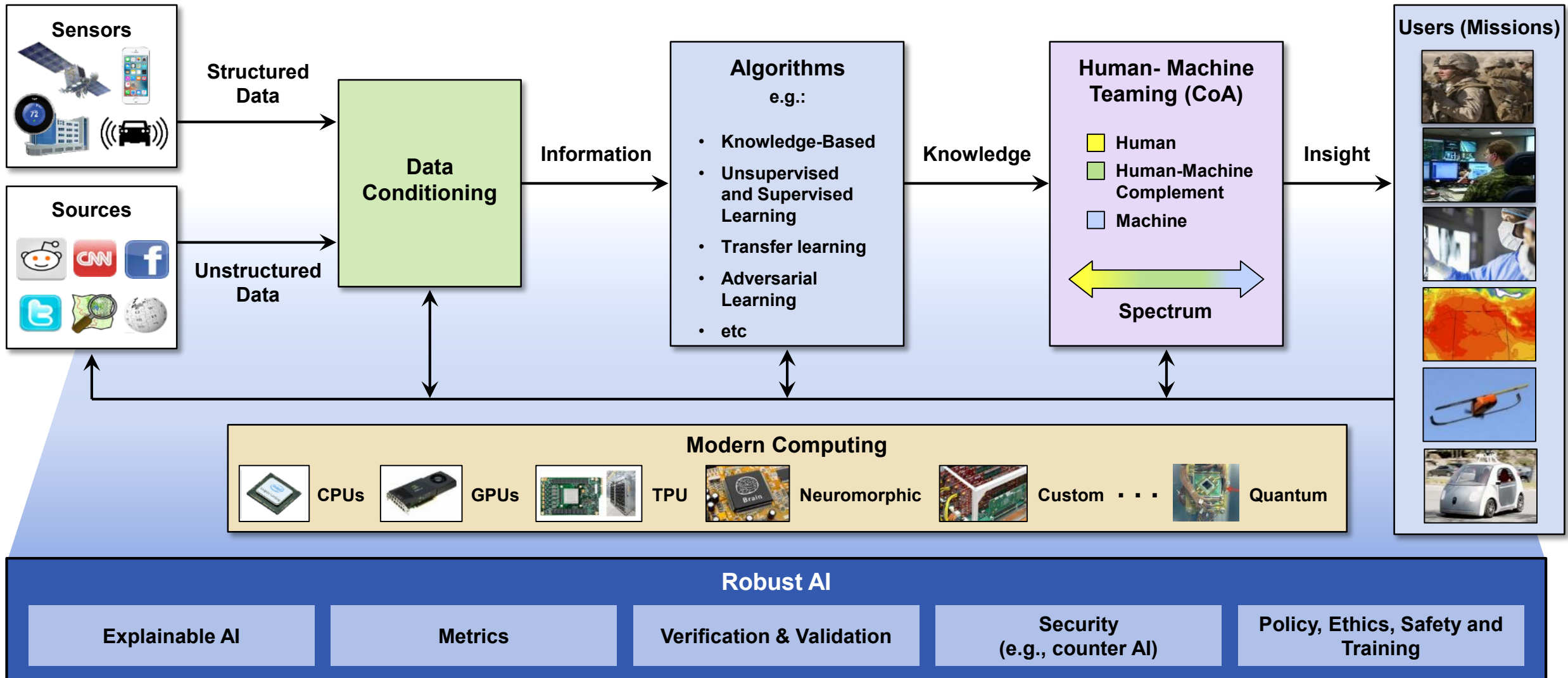
Outline

- **Background**
- **Lay-of-the-Land**
- **AI for Cyber Security**
 - **Background**
 - **Findings and Recommendations**
- **Summary**



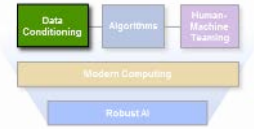


AI Canonical Architecture



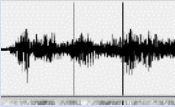












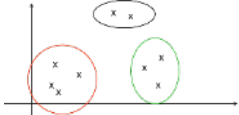
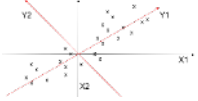

Unstructured and Structured Data



Data Conditioning/Storage Technologies

- Data to Information -

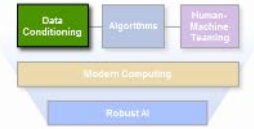
Structured Data Types			
 Speech	 Sensors	 Network Logs	 Metadata
Unstructured Data Types			
 Social Media	 Human Behavior	 Reports	 Side Channel

Technologies	Capabilities Provided
Infrastructure/Databases   	<ul style="list-style-type: none"> • Indexing/Organization/Structure • Domain Specific Languages • High Performance Data Access • Declarative Interfaces
Machine Learning (Unsupervised)  	<ul style="list-style-type: none"> • Limited machine learning • Dimensionality Reduction • Clustering/Pattern Recognition • Outlier Detection
Data Labeling 	<ul style="list-style-type: none"> • Initial data exploration • Highlight missing or incomplete data • Reorient sensors/recapture data • Look for errors/biases in collection

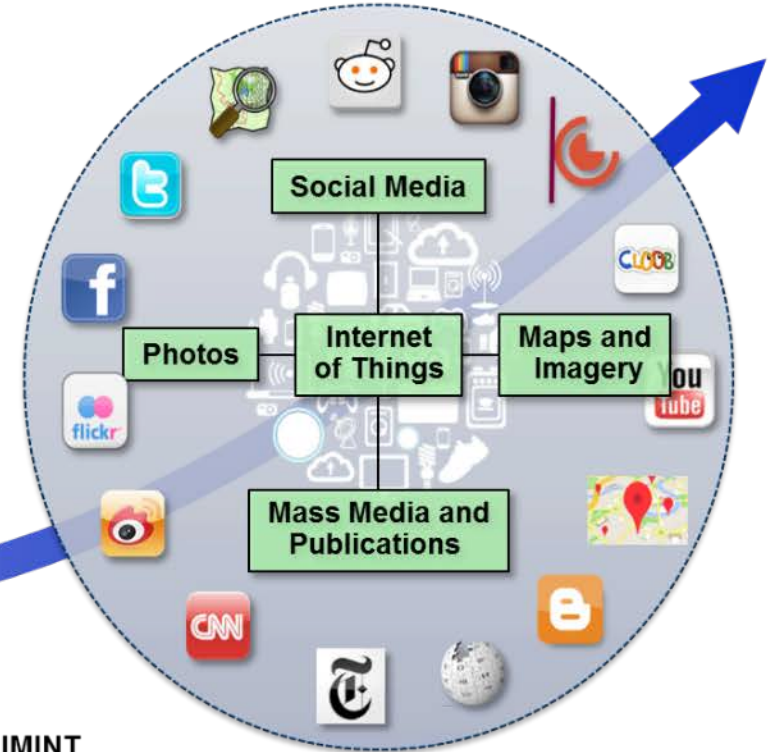
Important needs are in labeling data and automating data conditioning



Data Conditioning: The Open-Source Intelligence Opportunity Big Data Boom



- Open source data is growing exponentially
 - 2.8B Internet users
 - 2B smartphone users
 - Commercial satellites and imagery coming online
- Data are rich with information about systems, users, organizations, relationships, events
- Data can be used to enrich information from classified sources



Open Sources

Classified Sources

SIGINT

IMINT

HUMINT

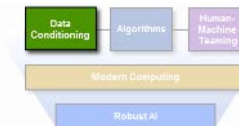
MOVINT



Finding #1: Cyber data is voluminous and is both structured and unstructured



Large Truth-marked Cyber Datasets Are Hard to Find



- Continued need for commercial and government Enterprises to share data from incidents
- Some databases exist but are not easy to use or widely accessible
- Very little cyber data is truth-marked
- Much academic research still leverages antiquated datasets

Cybersecurity Information Sharing Act of 2015

May 2016 Volume 11, Issue 5

From the Desk of Thomas F. Duffy, Chair

We've all heard talk of the Cybersecurity Information Sharing Act, but what does it really mean? We hope that this newsletter is a quick cheat sheet that highlights the key takeaways, as well as provide resources for additional information if you'd like to conduct a deeper dive into the topic.

JUST LAUNCHED: LEARN HOW IMPACT ADVANCES CYBERSECURITY R&D

The Information Marketplace for Policy and Analysis of Cyber-risk & Trust program provides infrastructure and event data to cybersecurity researchers to use as they develop tools, test theories, and identify solutions to address cyber threats. Learn more at www.dhs.gov/csd-impact.

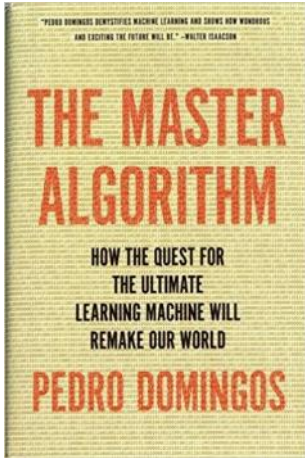
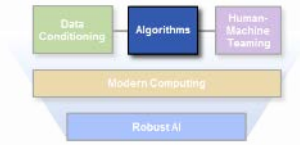
Table 1 Summary of popular datasets in the intrusion detection domain [30]

Data source	Dataset name	Abbreviation
Network Traffic	DARPA 1998 TCPDump Files	DARPA98
	DARPA 1999 TCPDump Files	DARPA99
	KDD99 Dataset	KDD99
	10% KDD99 Dataset	KDD99-10
	Internet Exploration Shootout Dataset	IES
User behavior	Unix User Dataset	UNIXDS
System call sequences	DARPA 1998 BSM Files	BSM 98
	DARPA 1998 BSM Files	BSM 99
	University of New Mexico Dataset	UNM

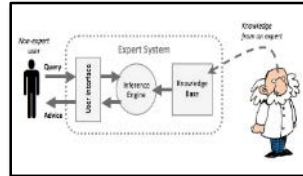
Finding #2: Lack of ground truth for cyber inhibits algorithm application to DoD problems



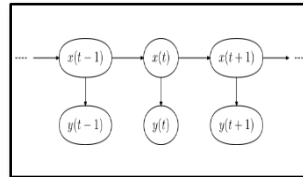
Machine Learning Algorithms Taxonomy



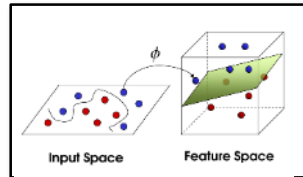
Algorithms*



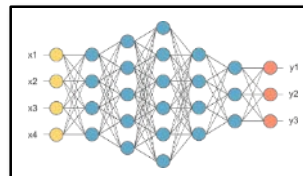
Symbolists
(e.g., exp. sys.)



Bayesians
(e.g., naive Bayes)



Analogizers
(e.g., SVM)

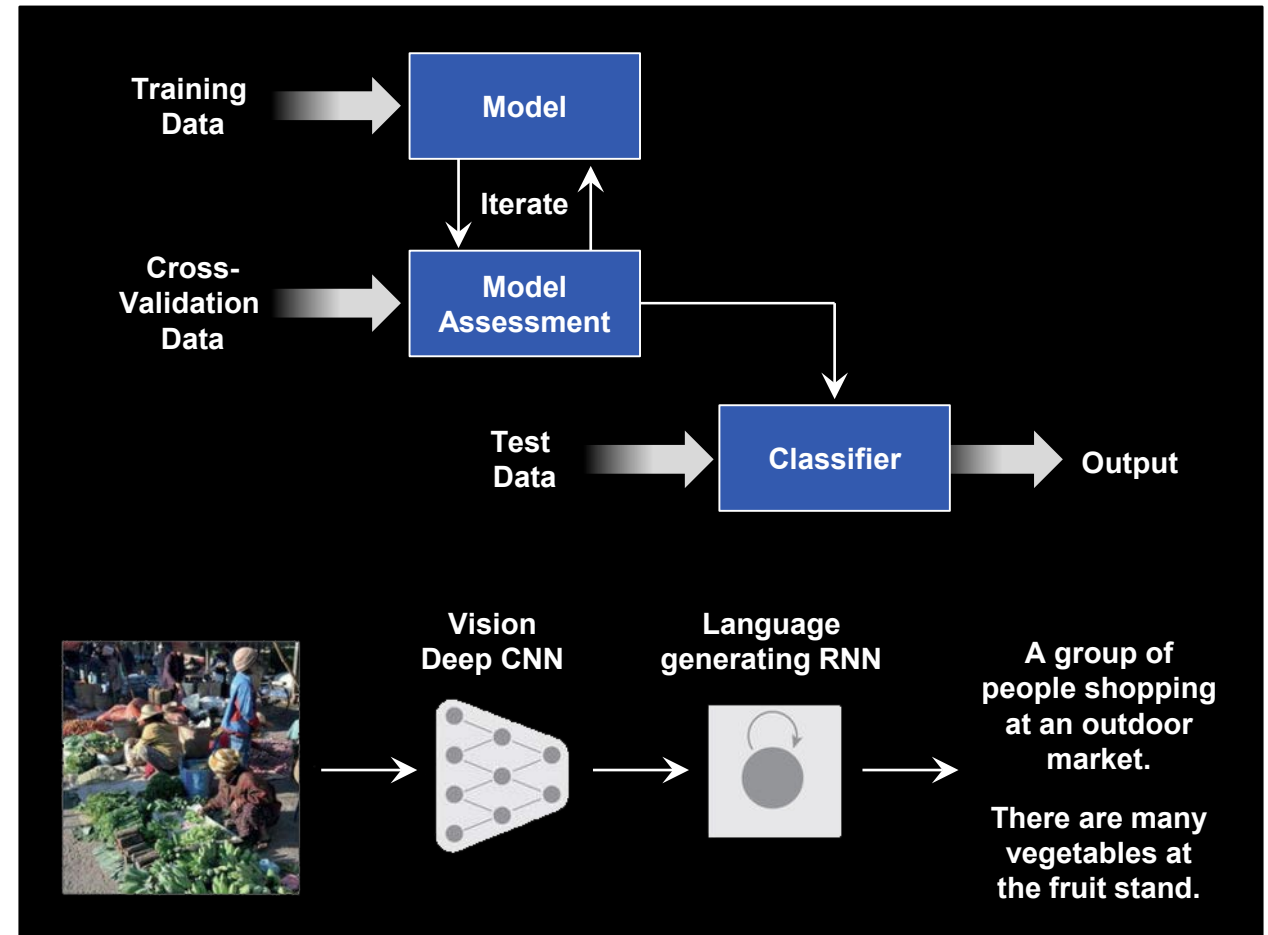


Connectionists
(e.g., DNN)



Evolutionaries
(e.g., genetic programming)

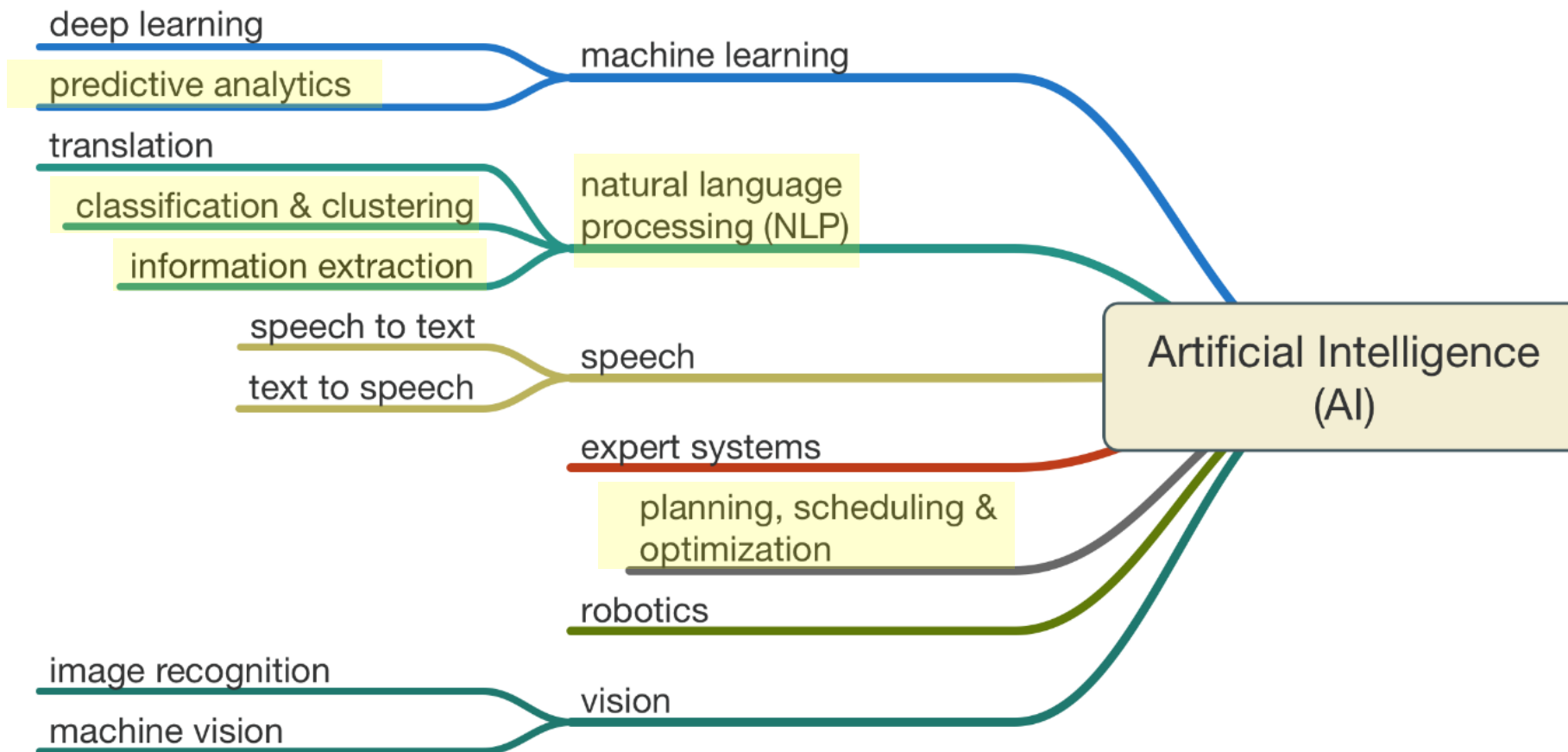
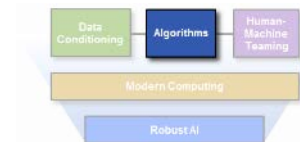
Machine Learning Applied to Classifiers



* "The Five Tribes of Machine Learning", Pedro Domingos



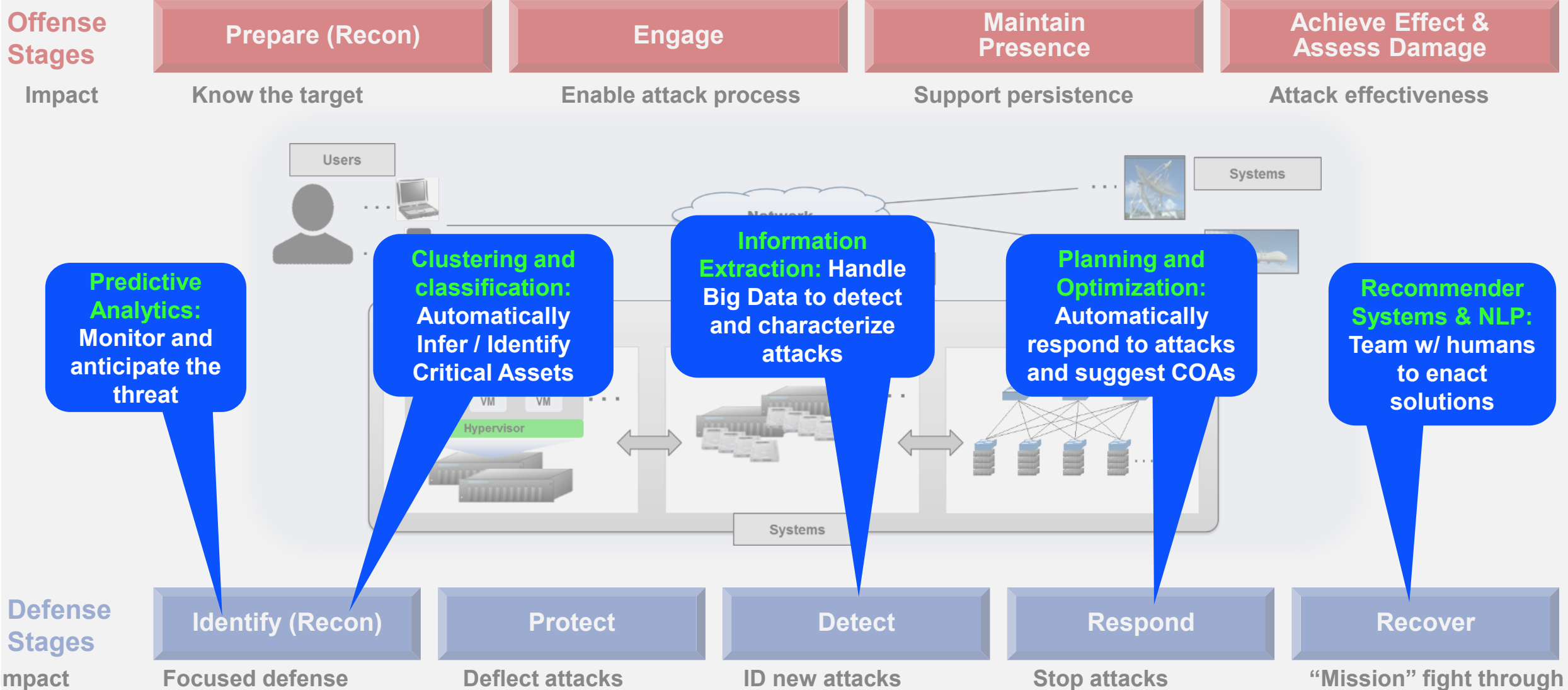
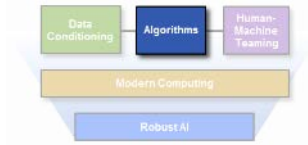
AI Algorithms for Cyber



Finding #3: Many algorithms exist which can be applied to cyber security

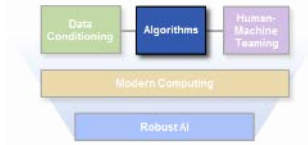


Artificial Intelligence Can Help

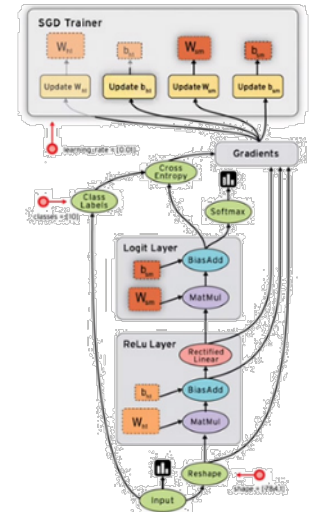




AI Resources Are Available on Line



- Open source toolkits allow users to leverage machine learning easily
- Commercial companies build business on AI tool kits that can be applied easily
- A Knowledge Base of Shared Knowledge and Solutions



ReLUplex



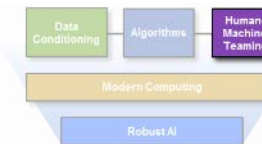
Apache Singa



Finding #4: Academia, commercial sectors are advancing algorithms and AI capabilities
Finding #5: Peer organizations are benefiting from open source communities



Human Machine Teaming



Study Finds Cyberthreat Data Overwhelming to Security Workers

A recent Ponemon report shows that organizations neglect to share essential cyberthreat data with board members and C-level executives.



Challenge: Cybersecurity and Big Data



“By 2018 the United States alone faces a shortage of 140,000 to 190,000 people with analytical expertise and 1.5 million managers and analysts with the skills to understand and make decisions based on the analysis of big data.”¹

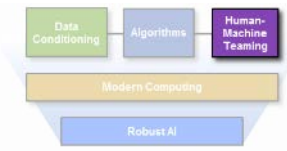
¹McKinsey&Company (May 2011), “Big data: The next frontier for innovation, competition, and productivity.” Available at: http://www.mckinsey.com/Insights/MGI/Research/Technology_and_Innovation/Big_data_The_next_frontier_for_innovation

Finding #6: Cyber security data overwhelms overworked analysts

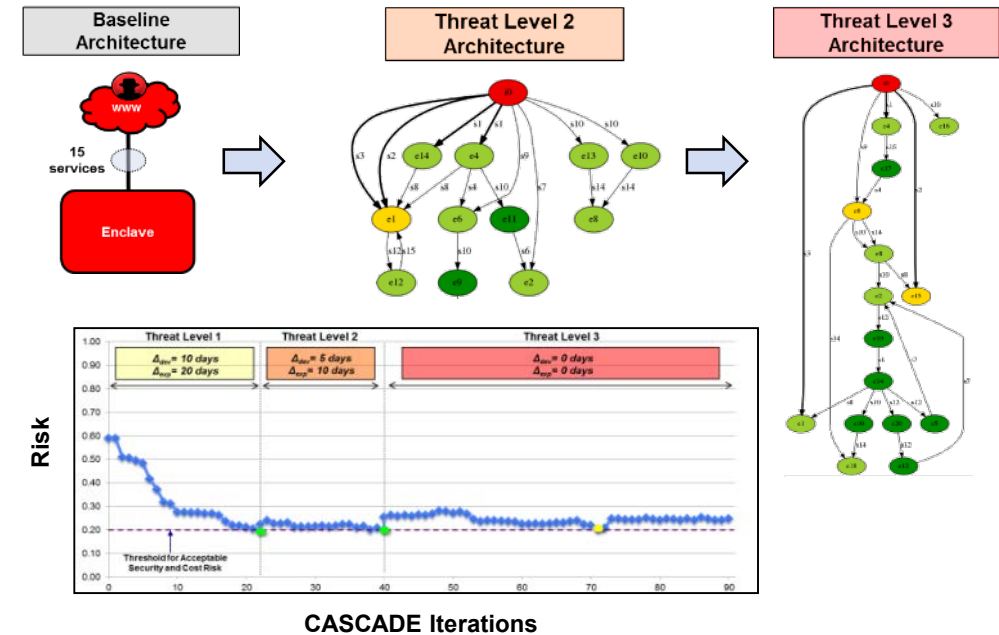
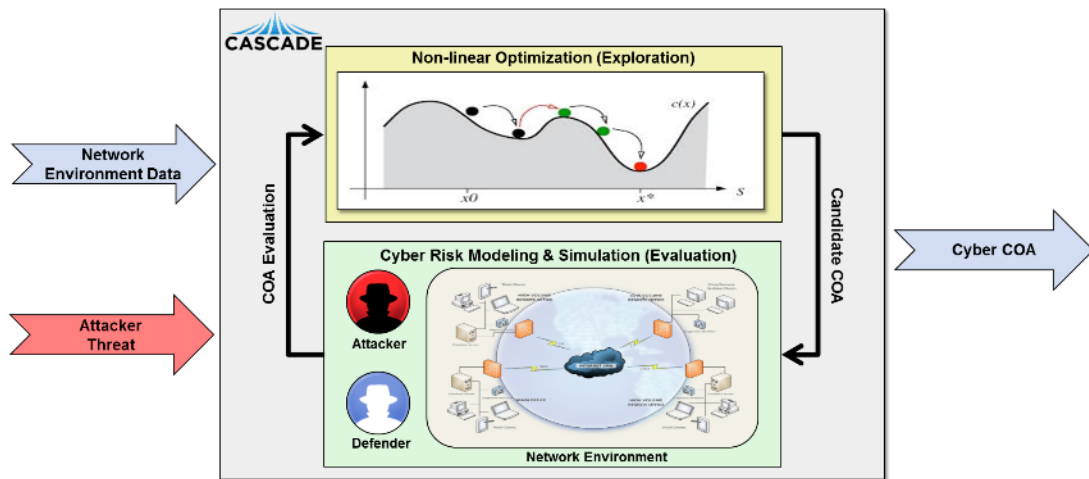
Finding #7: The US / DoD faces serious workforce shortages in cyber security expertise



Human Machine Teaming: Automated Cyber Decision Making

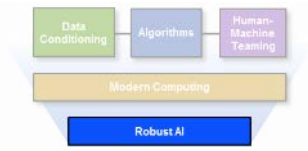


- **CASCADE - Cyber Adversarial SCenario modeling and Automated Decision Engine**
 - Dynamically quantifies risk in the face of an adaptive adversary
 - Considers mission context to selection optimal course of action (COA)
 - Prototype applied to configuration of network segmentation defense

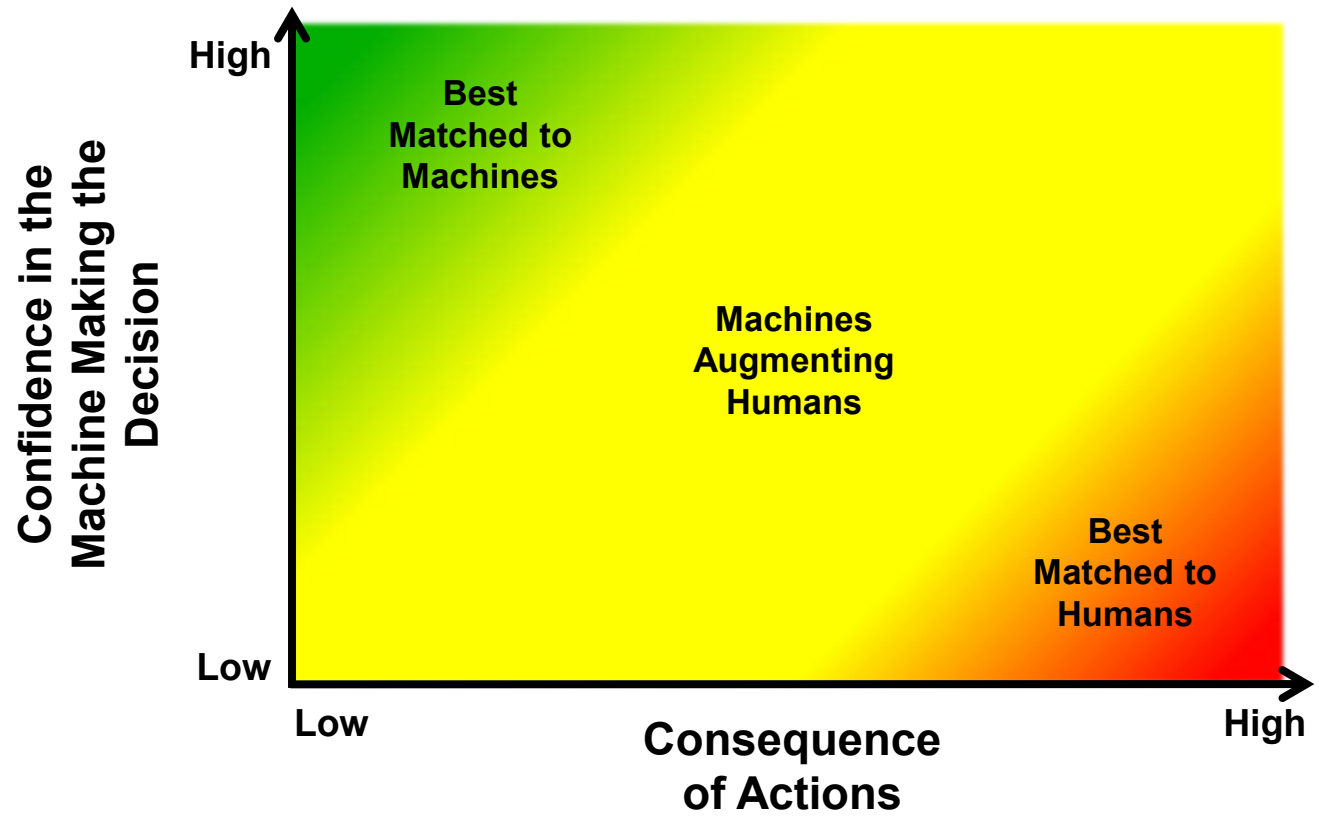




Robust AI: Engendering Trust

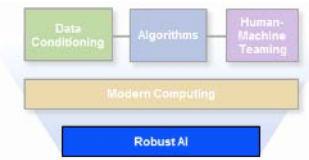


Confidence Level vs. Consequence of Actions

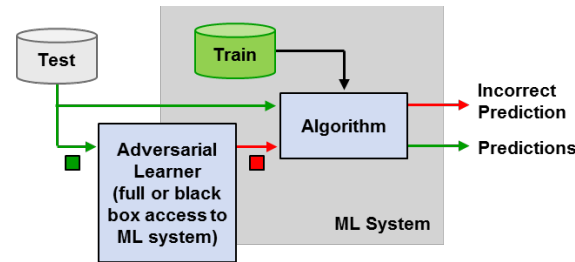
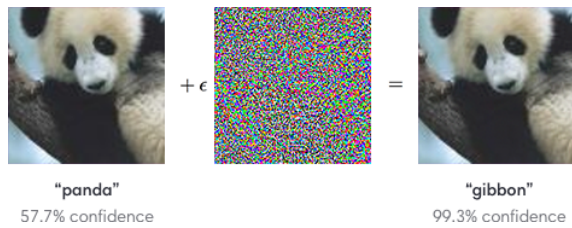




AI for Cyber must be Robust as Well



- By gaining access to an AI system, can an adversary learn, and then introduce, imperceptible perturbations to inputs that render the system un-usable?



- Cyber examples are appearing in literature demonstrating capabilities
 - Malware evades detection
 - Nefarious connections hidden by noise
 - Etc ..

Towards Evaluating the Robustness of Neural Networks

Nicholas Carlini David Wagner
University of California, Berkeley

ABSTRACT
Neural networks provide state of the art results for most machine learning tasks. Unfortunately, neural networks are vulnerable to adversarial examples given an input x and any target classification t , it is possible to find a new input x' that is similar to x but classified as t . This makes it difficult to apply neural networks in security-critical areas. Defensive distillation is a recently proposed approach that can take an arbitrary neural network, and increase its robustness, reducing the success rate of current attacks' ability to find adversarial examples from 95% to 0.5%.

In this paper, we demonstrate that defensive distillation does not significantly increase the robustness of neural networks by introducing three new attack algorithms that are successful on both distilled and undistilled neural networks with 100% probability. Our attacks are tailored to three distance metrics used previously in the literature, and when compared to previous adversarial example generation algorithms, our attacks are often much more effective (and never worse). Furthermore, we propose using high-confidence adversarial examples in a simple transferability test we show can also be used to break defensive distillation. We hope our attacks will be used as a benchmark in future defense attempts to create neural networks that resist adversarial examples.

1. INTRODUCTION
Deep neural networks have become increasingly effective at many difficult machine-learning tasks. In the image recognition domain, they are able to recognize images with near-human accuracy [27], [25]. They are also used for speech recognition [18], natural language processing [1], and playing games [43], [32].

However, researchers have discovered that existing neural networks are vulnerable to attack. Szegedy et al. [44] first noticed the existence of *adversarial examples* in the image classification domain: it is possible to transform an image by a small amount and thereby change how the image is classified. Often, the total amount of change required can be so small as to be undetectable.

The degree to which attackers can find adversarial examples limits the domains in which neural networks can be used. For example, if we use neural networks in self-driving cars, adversarial examples could allow an attacker to cause the car to take unwanted actions.

The existence of adversarial examples has inspired research on how to harden neural networks against these kinds of attacks. Many early attempts to secure neural networks provided only marginal robustness improvements [20], [42].

Defensive distillation [39] is one such proposed for hardening neural networks against adversarial examples. Initial analysis proved to be very promising: distillation defeats existing attack algorithms with success probability from 95% to 0.5%. This can be applied to any feed-forward neural network, and requires a single re-training step, and is the only defenses giving strong security guarantees.

In general, there are two different approaches to evaluate the robustness of a neural network: a lower bound, or construct attacks that are hard to find. The former approach, while sound, more difficult to implement in practice, and required approximations [2], [21]. On the other hand, the latter approach is more practical, but only provides a heuristic measure of robustness.

Recent researches have shown that machine learning based malware detection algorithms are very vulnerable under the attack of adversarial examples. These works mainly focused on the detection algorithms which use features with fixed dimension, while some researchers have begun to use recurrent neural networks (RNN) to detect malware based on sequential API features. This paper proposes a novel algorithm to generate adversarial examples, which are used to attack a RNN based malware detection system. It is usually hard for malicious attackers to know the exact structures and weights of the victim RNN. A substitute RNN is trained to approximate the victim RNN. Then we propose a generative RNN to output sequential adversarial examples from the original sequential malware inputs. Experimental results showed that RNN based malware detection algorithms fail to detect most of the generated malicious adversarial examples, which means the proposed model is able to effectively bypass the detection algorithms.

1 Introduction
Machine learning has been widely used in various commercial and non-commercial products, and has brought great convenience and profits to human beings. However, recent researches on adversarial examples show that many machine learning algorithms are not robust at all when someone wants to crack them on purpose (Szegedy et al. 2013, Goodfellow, Shlens, and Szegedy 2014). Adding some small perturbations to original samples will make a classifier unable to classify them correctly.

In some security related applications, attackers will try their best to attack any defensive systems to spread their malicious products such as malware. Existing machine learning based malware detection algorithms mainly represent programs as feature vectors with fixed dimension and classify them between benign programs and malware (Koller and El-Masri 2006). For example, a binary feature vector can be constructed according to the presence or absence of system APIs (i.e. application programming interfaces) in a program (Schultz et al. 2001). Grosse et al. (Grosse et al. 2016) and

*Prof. Ying Tan is the corresponding author.

Black-Box Attacks against RNN based Malware Detection Algorithms

Wenwei Hu and Ying Tan*
Key Laboratory of Machine Perception (MOE), and Department of Machine Intelligence
School of Electronics Engineering and Computer Science, Peking University, Beijing, 100871 China
{wweihu, ytan}@pku.edu.cn

Abstract
Recent researches have shown that machine learning based malware detection algorithms are very vulnerable under the attack of adversarial examples. These works mainly focused on the detection algorithms which use features with fixed dimension, while some researchers have begun to use recurrent neural networks (RNN) to detect malware based on sequential API features. This paper proposes a novel algorithm to generate adversarial examples, which are used to attack a RNN based malware detection system. It is usually hard for malicious attackers to know the exact structures and weights of the victim RNN. A substitute RNN is trained to approximate the victim RNN. Then we propose a generative RNN to output sequential adversarial examples from the original sequential malware inputs. Experimental results showed that RNN based malware detection algorithms fail to detect most of the generated malicious adversarial examples, which means the proposed model is able to effectively bypass the detection algorithms.

Hu et al. (Hu and Tan 2017) have shown that fixed dimensional feature based malware detection algorithms are very vulnerable under the attack of adversarial examples. Recently, as recurrent neural networks (RNN) became popular, some researchers have tried to use RNN for malware detection and classification (Pascanu et al. 2015; Tobiyama et al. 2016; Kolonjajic et al. 2016). The API sequence invoked by a program is used as the input of RNN. RNN will predict whether the program is benign or malware.

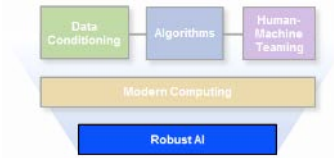
This paper tries to validate the security of a RNN based malware detection model when it is attacked by adversarial examples. We proposed a novel algorithm to generate sequential adversarial examples.

Existing researches on adversarial samples mainly focus on images. Images are represented as matrices with fixed dimension, and the values of the matrices are continuous. API sequences consist of discrete symbols with variable length.

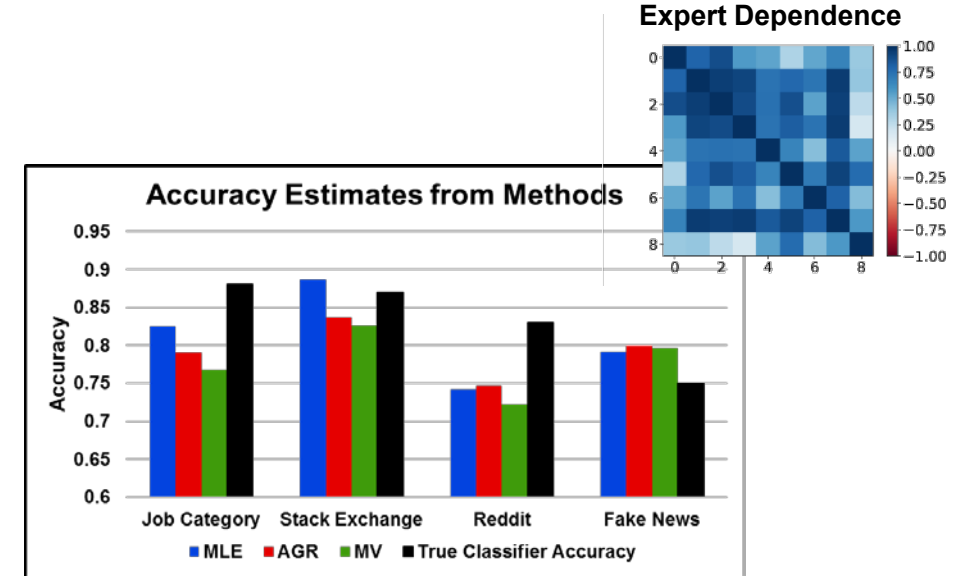
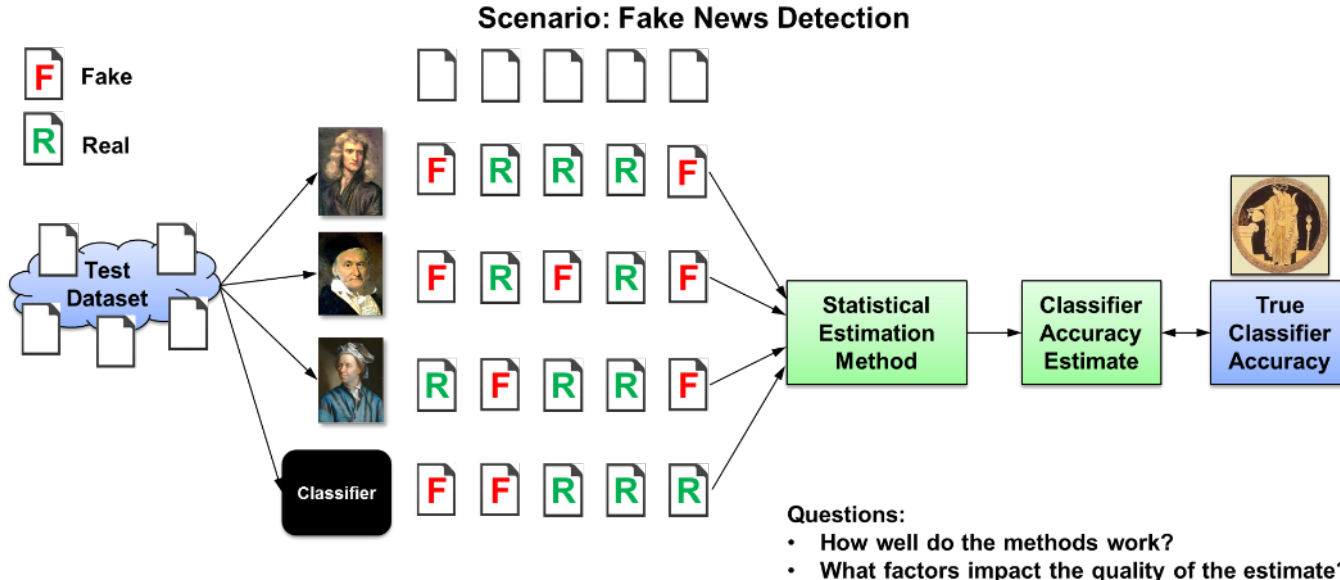
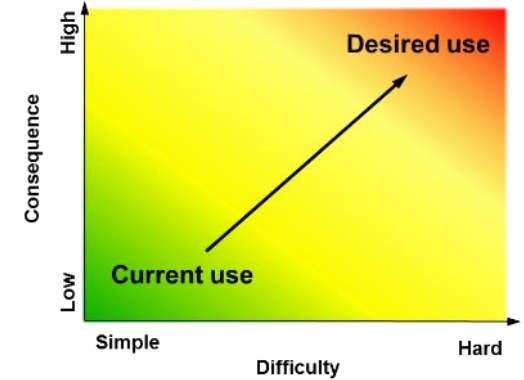
Finding #8: Adversarial attacks can limit effectiveness of cyber AI solutions
Finding #9: Vulnerable Cyber AI detection, classification algorithms can lead to incorrect behavior



Robust AI: Evaluating Classifier Performance using Human Expert Judgment

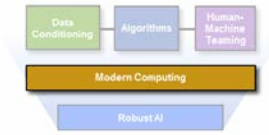


- AI and ML currently applied mostly to simple, low consequence problems
 - We want to transition use to hard, high consequence problems
- Methods are needed to evaluate AI classifiers that leverage expert judgement
 - Must be robust to inter-rater dependence and variability





Modern AI Computing Engines



Computing Class



CPU

What It Provides to AI

- Most popular computing platform
- General purpose compute



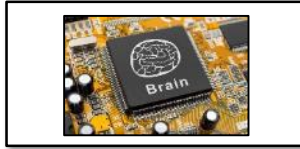
GPU

- Used by most for training algorithms (good for NN backpropagation)



TPU

- Speeds up inference time (domain specific architecture)



Neuromorphic

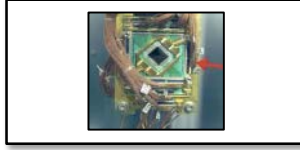
- Still a research area



Custom

- Ability to speed up specific computations of interest (e.g. graphs)

⋮

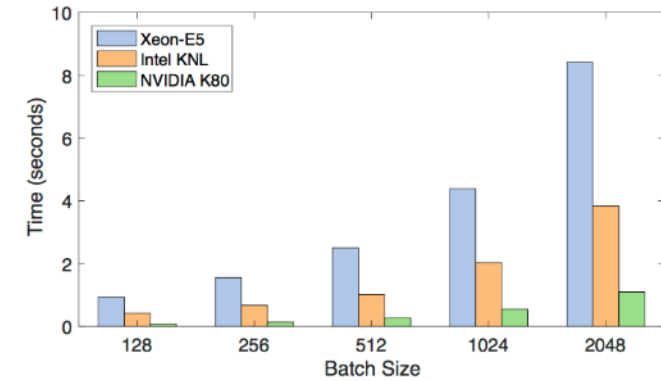


Quantum

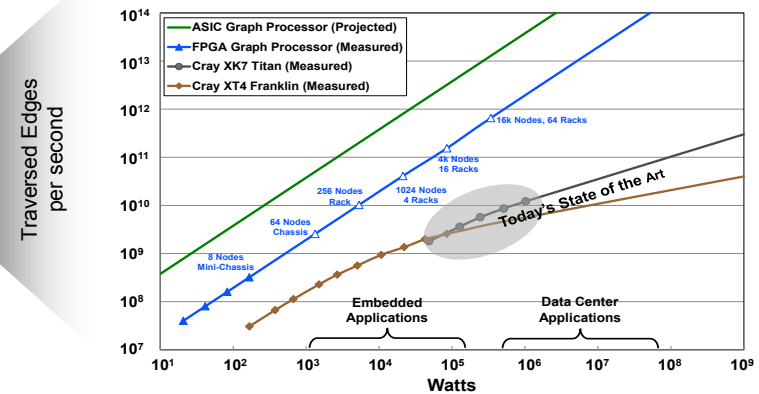
- Benefits unproven until now
- Recent results on HHL (linear system of equations)

Selected Results

Alexnet comparison: Forward-Backward Pass



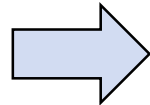
SpGEMM Performance using Graph Processor (G102)





Outline

- **Background**
- **Lay-of-the-Land**
- **AI for Cyber Security**
 - **Background**
 - **Highlights, Findings and Recommendations**
- **Summary**





Summary of Study Findings & Recommendations

Data

1. Cyber data is voluminous and is multi-domain, structured and unstructured
2. Lack of ground truth for cyber inhibits algorithm application to DoD problems

✓ Lead the way in Cyber Big Data conditioning by leveraging expertise in Big Data collection, creation and curation to support AI for Cyber

Algorithms

3. Many algorithms exist which can be applied to cyber
4. Academia, commercial sectors are advancing algorithms and AI capabilities
5. Peer organizations are benefiting from open source communities

✓ Engage with academic community to maintain awareness of and influence where possible, leverage open-source toolkits and libraries to jumpstart DoD mission capabilities

Human Machine Teaming

6. Declining human resource environment creates opportunity to help
7. Recommender systems at core of much commercial AI success

✓ Automate and augment cyber tasks of data triage, correlation, leveraging active learning to improve AI solutions, capitalize on analyst cyber expertise

Robust AI

8. Adversarial attacks can limit effectiveness of cyber AI solutions, leading to incorrect behavior
9. Promising work in 'proving' AI behavior is appearing in academia

✓ Lead the way in robust AI for cyber for DoD applications, leveraging and applying recent work in academia



LL Cyber AI-related Workshops and Symposia



Artificial Intelligence for Cyber Security Workshop

- Forum for AI researchers and practitioners to share research and experiences in applying AI to Cyber Security

Chairs



Bill Streilein



Dave Martinez



Neal Wagner

New Orleans, Louisiana • February 2, 2018

Theme: Applications of AI to Internet of Things

Keynotes



Sal Stolfo
 Professor of Computer Science
 Dept. of Computer Science,
 Columbia University



Trung Tran
 Laboratory of Physical Sciences,
 University of Maryland,
 Baltimore County



Graph Exploitation Symposium

- Brings together leading experts from universities, industry, and government to explore the state of the art and define a future roadmap in network science

Dedham, Massachusetts • April 23-25, 2018



Technical Co-Chairs



Ben Miller



Rajmonda Caceres

Chairs



Bill Streilein



Sanjeev Mohindra

POC: Ben Miller,
bamiller@ll.mit.edu



Summary

- **U.S. needs to regain AI leadership by strategically partnering with small and large commercial companies plus academia**
- **Potential for major impact remains for DoD applications**
 - Although there is a lot of activity in community, only pockets of cyber success exist
- **Transfer of algorithms to DoD mission is challenging**

- **Demonstrated achievements in applying AI to cyber**
 - Fluent in Big Data architectures and databases
 - Cyber discussion detection, traffic characterization, counterfeit detection
- **Focus should be on unique areas of expertise, connection to mission**
 - Mission process and data requirements
 - Adapting latest algorithms to mission needs
 - Developing robust AI solutions



AI Bibliography List (few selected set)

