# Explainable Artificial Intelligence Research at DARPA
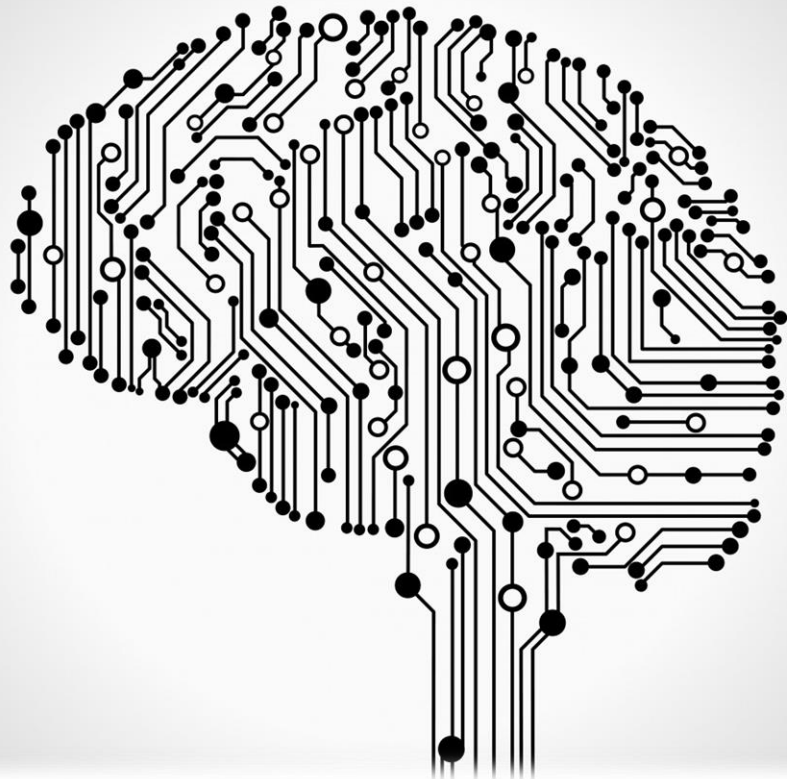
David Gunning

DARPA/I2O

**perceive**
rich, complex and subtle information

**learn**
within an environment

**abstract**
to create new meanings

**reason**
to plan and to decide

# Three Waves of AI

## DESCRIBE
### Symbolic Reasoning

engineers create sets of logic rules to represent knowledge in limited domains

reasoning over narrowly defined problems

no learning capability and poor handling of uncertainty

| | | | |
|---|---|---|---|
| Perceiving | ■ | | |
| Learning | | | |
| Abstracting | | | |
| Reasoning | ■ | ■ | |

## PREDICT
### Statistical Learning

engineers create statistical  models for specific problem domains and train them on big data

nuanced classification and prediction capabilities

no contextual capability and minimal reasoning ability

| | | | |
|---|---|---|---|
| Perceiving | ■ | ■ | ■ |
| Learning | ■ | ■ | |
| Abstracting | ■ | | |
| Reasoning | ■ | | |

## EXPLAIN
### Contextual Adaptation

engineers create systems that construct explanatory models for classes of real world phenomena
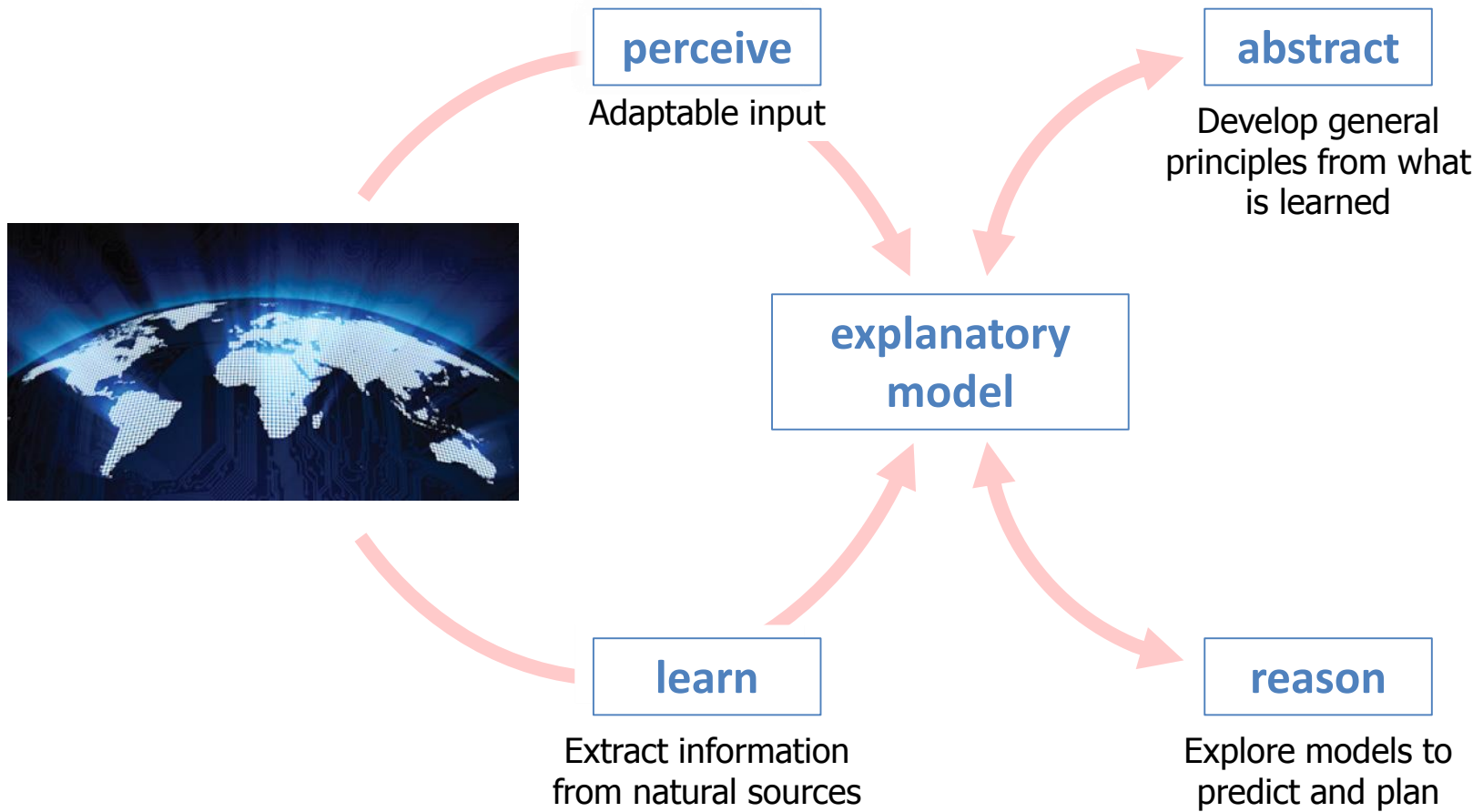
natural communication among machines and people

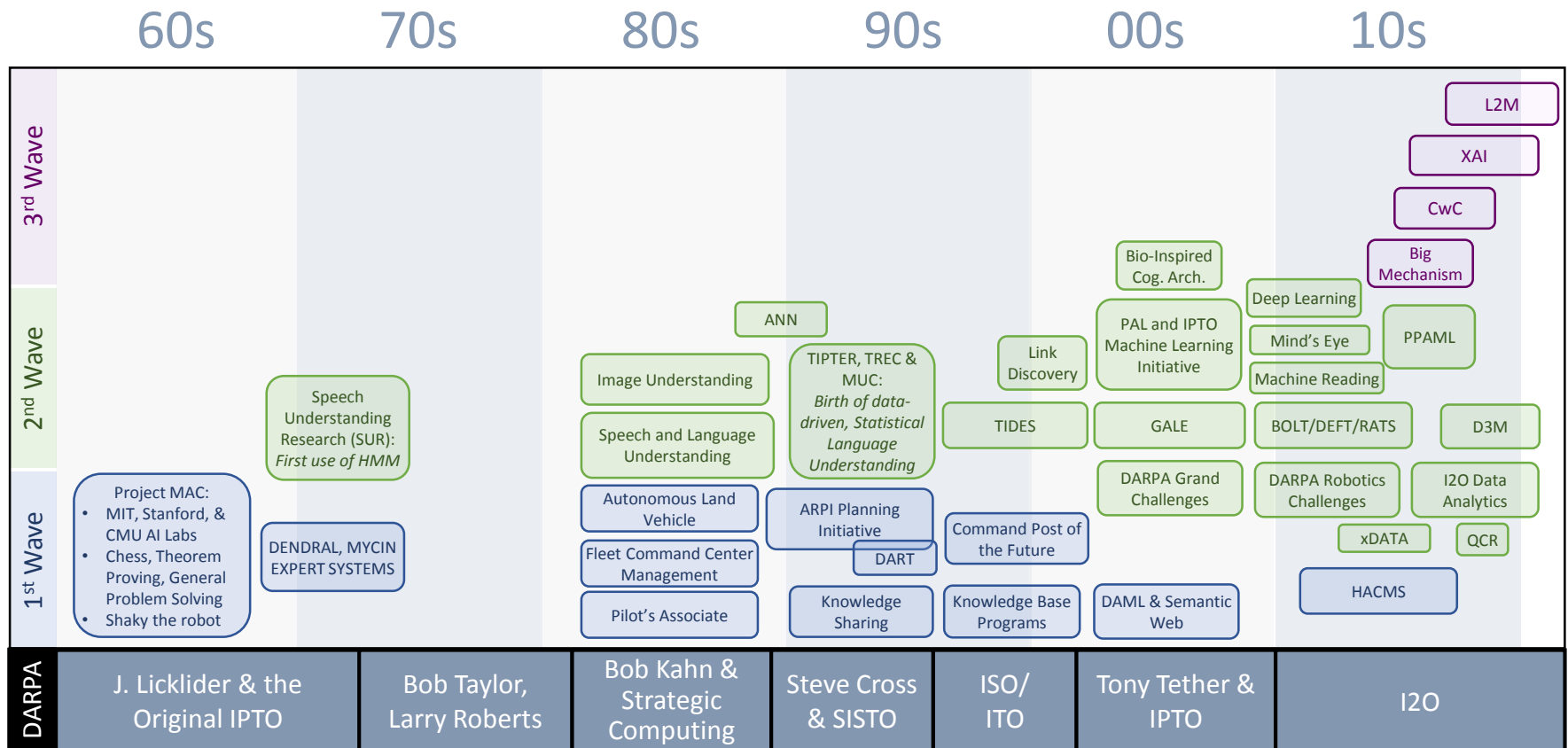systems learn and reason as they encounter new tasks and situations

| | | | |
|---|---|---|---|
| Perceiving | ■ | ■ | ■ |
| Learning | ■ | ■ | |
| Abstracting | ■ | | |
| Reasoning | ■ | ■ | |

# Third wave technology: explanatory models

**perceive**

Adaptable input

**abstract**

Develop general principles from what is learned

**explanatory model**

**learn**

Extract information from natural sources

**reason**

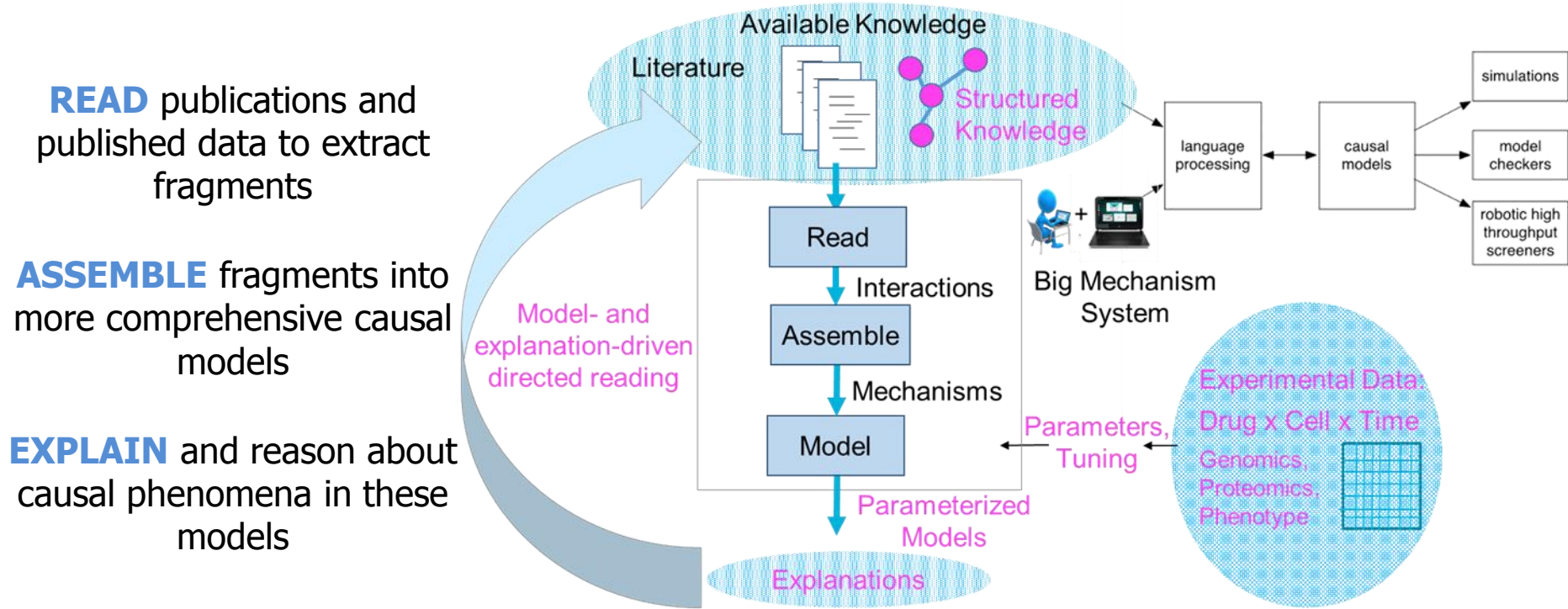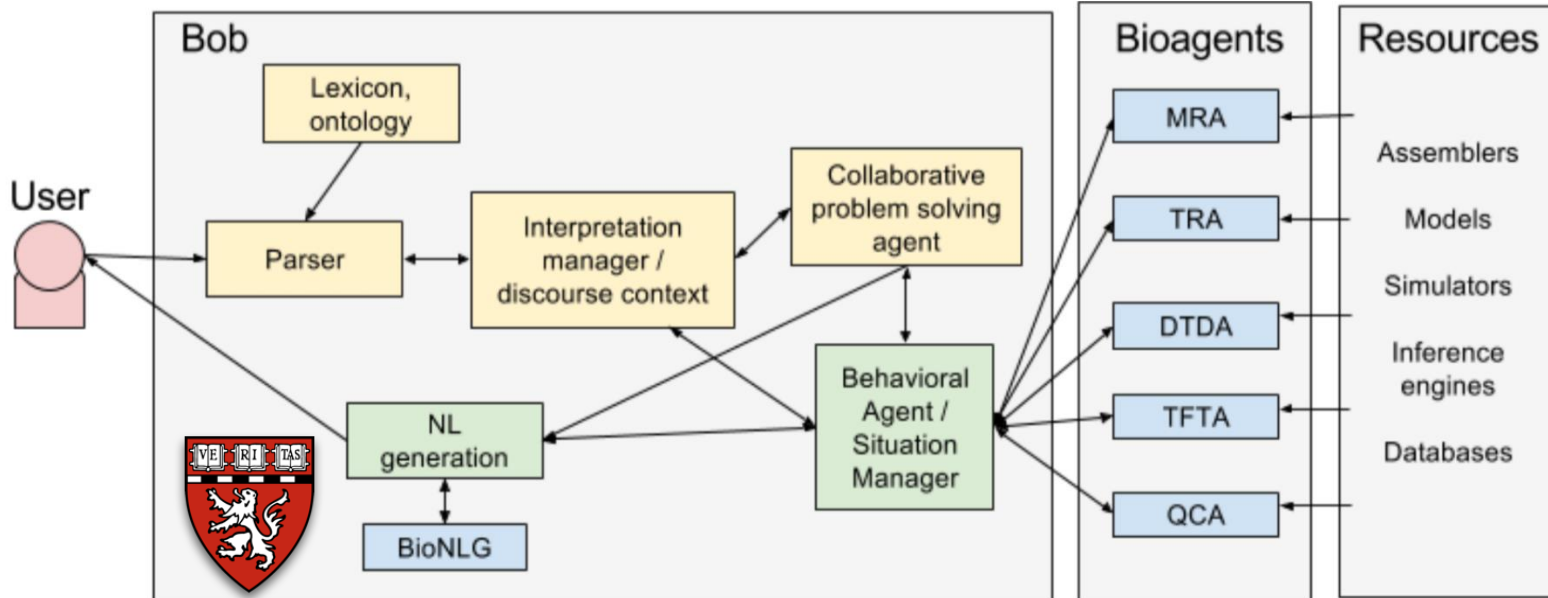Explore models to predict and plan

# DARPA Contributions to AI

# Automatically construct causal models of complicated systems to predict and explain the effects of system perturbations (cell biology)

**READ** publications and published data to extract fragments

**ASSEMBLE** fragments into more comprehensive causal models

**EXPLAIN** and reason about causal phenomena in these models



Build causal, mechanistic, quantitative models to produce explanatory models of unprecedented completeness and consistency

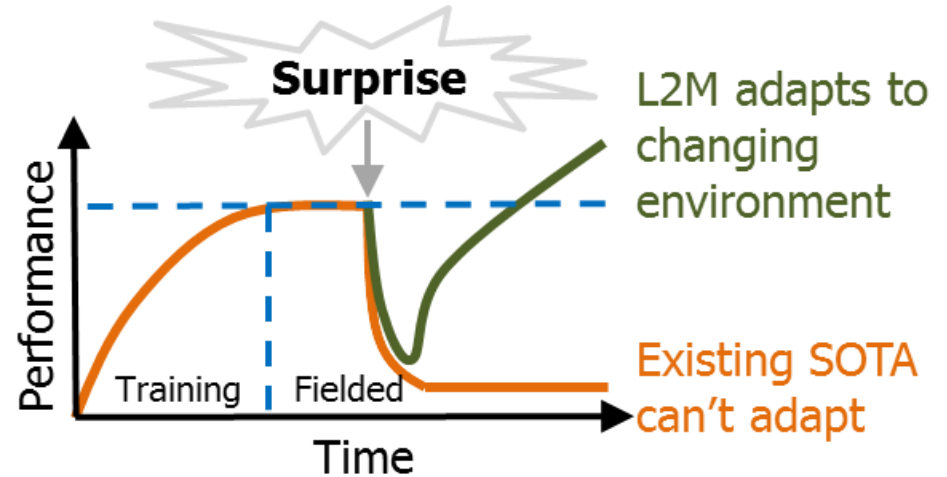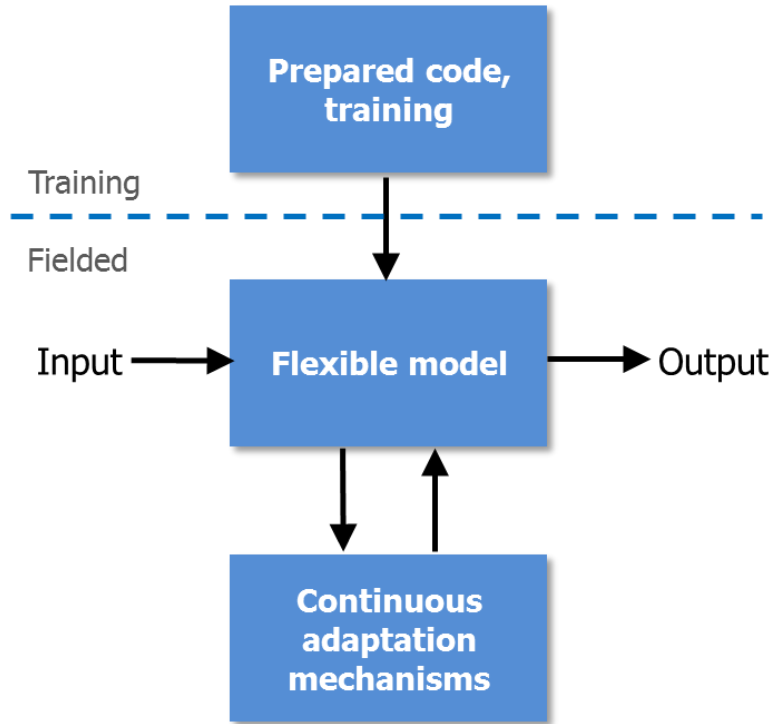## Collaborative Problem Solving Agent "Bob"
## Harvard Medical School



- Implementation of a generic **language understanding** system
- Working implementation of generic **collaborative problem solving and planning**
- **Biological problem solving agents (Bioagents)**, which are generic for their specific sub-tasks
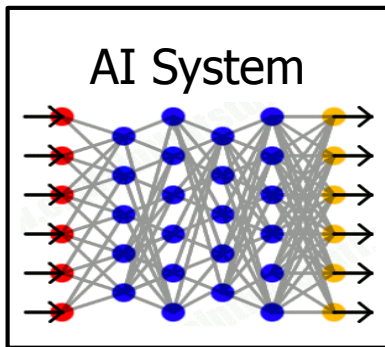- Integration into a working **communication-for-biocuration system**

Develop fundamentally new machine learning mechanisms that will enable systems to improve their performance over their lifetimes
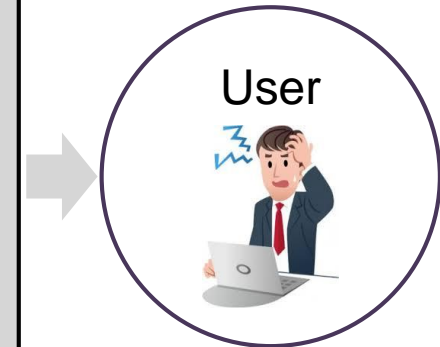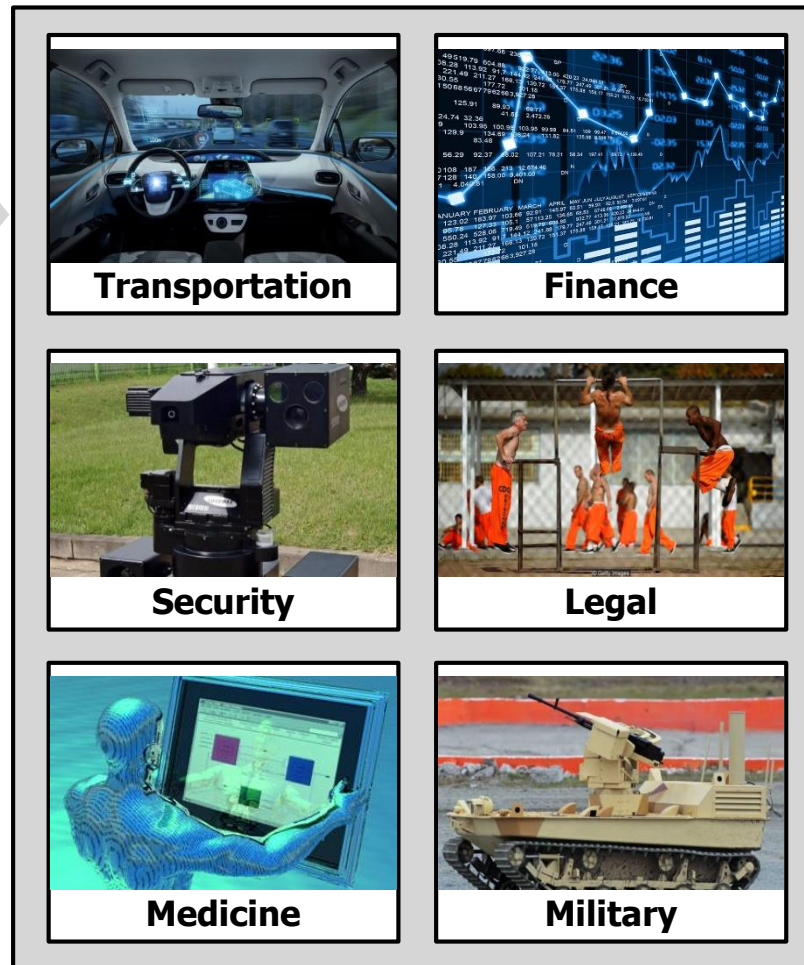


**Focus on functional system development, taking inspiration from known biological properties**

Dynamically evolve networks online
Use scalable approaches

**Identify and explore biological mechanisms that underlie real-time adaptation for translation into novel algorithms**
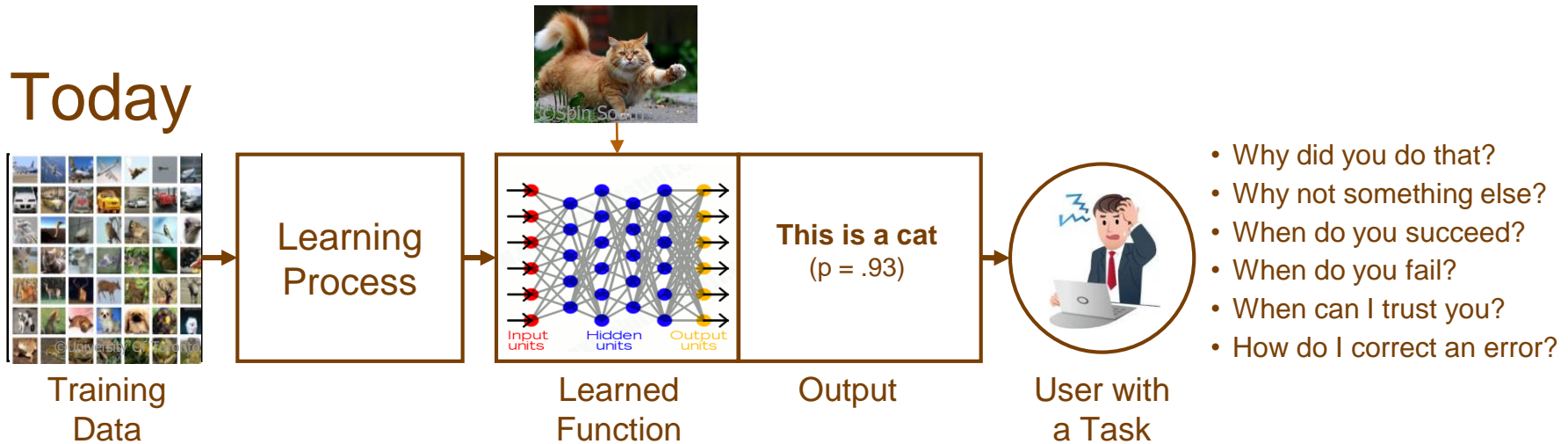
# The Need for Explainable AI



**AI System**

- We are entering a new age of AI applications
- Machine learning is the core technology
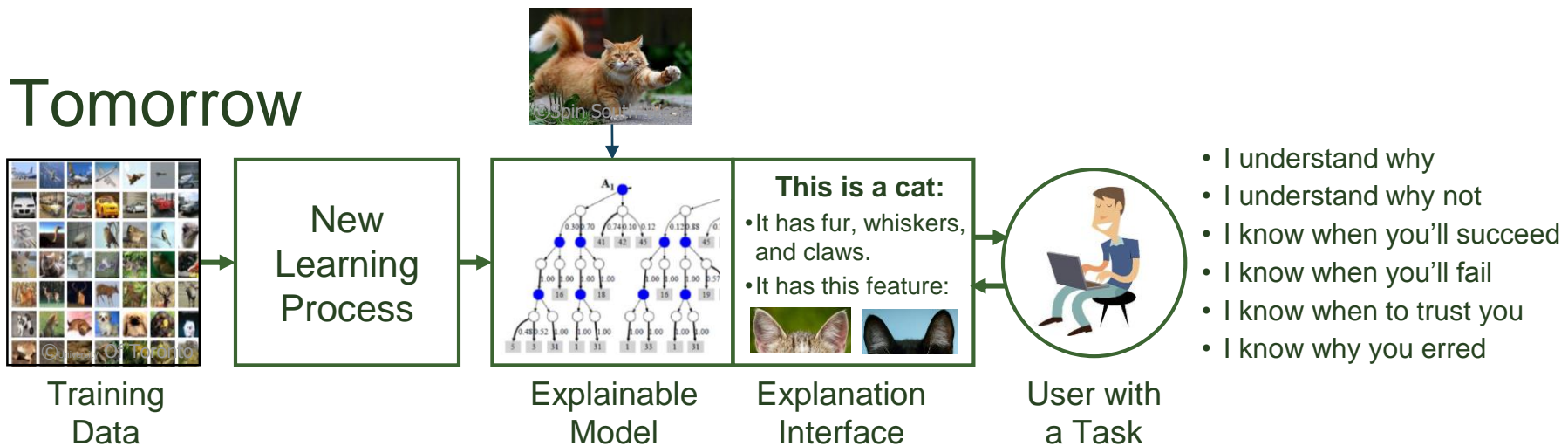- Machine learning models are opaque, non-intuitive, and difficult for people to understand

**Transportation**

**Finance**

**Security**

**Legal**

**Medicine**

**Military**

**User**

- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

- The current generation of AI systems offer tremendous benefits, but their effectiveness will be limited by the machine's inability to explain its decisions and actions to users.

- Explainable AI will be essential if users are to understand, appropriately trust, and effectively manage this incoming generation of artificially intelligent partners.

## Today



Training Data → Learning Process → Learned Function (Input units, Hidden units, Output units) → Output: **This is a cat** (p = .93) → User with a Task

- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

## Tomorrow



Training Data → New Learning Process → Explainable Model → Explanation Interface: **This is a cat:**
- It has fur, whiskers, and claws.
- It has this feature:

→ User with a Task

- I understand why
- I understand why not
- I know when you'll succeed
- I know when you'll fail
- I know when to trust you
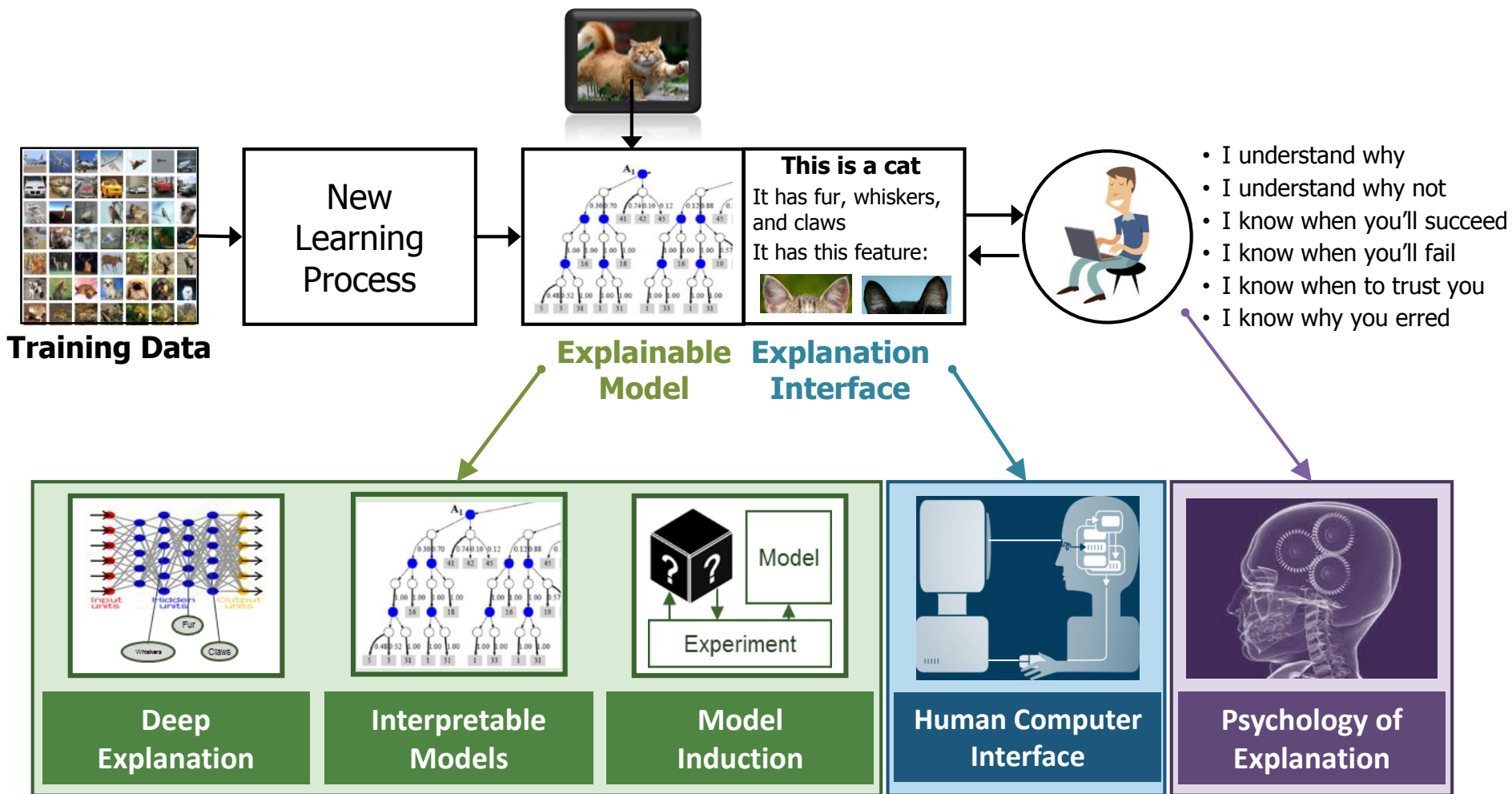- I know why you erred

# Explainable Artificial Intelligence (XAI)
## David Gunning

Create a suite of machine learning techniques to produce more explainable models and enable human users to understand, trust, and effectively manage the emerging generation of artificially intelligent partners

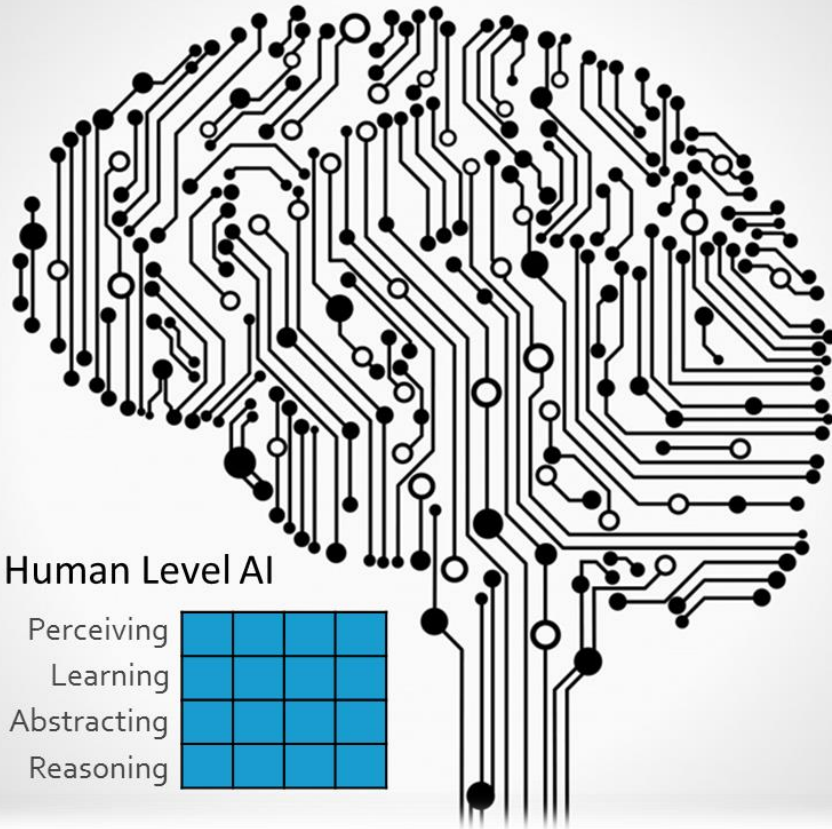**Training Data** → **New Learning Process** → **Explainable Model** → **Explanation Interface**

**This is a cat**
It has fur, whiskers, and claws
It has this feature:

- I understand why
- I understand why not
- I know when you'll succeed
- I know when you'll fail
- I know when to trust you
- I know why you erred

| Deep Explanation | Interpretable Models | Model Induction | Human Computer Interface | Psychology of Explanation |
|---|---|---|---|---|

# XAI Developers (TA1)

| CP | Performer | Explainable Model | Explanation Interface |
|---|---|---|---|
| **Both** | UC Berkeley | Deep Learning | Reflexive & Rational |
| | Charles River | Causal Modeling | Narrative Generation |
| | UCLA | Pattern Theory+ | 3-level Explanation |
| **Autonomy** | Oregon State | Adaptive Programs | Acceptance Testing |
| | PARC | Cognitive Modeling | Interactive Training |
| | CMU | Explainable RL (XRL) | XRL Interaction |
| **Analytics** | SRI International | Deep Learning | Show & Tell Explanation |
| | Raytheon BBN | Deep Learning | Argumentation & Pedagogy |
| | UT Dallas | Probabilistic Logic | Decision Diagrams |
| | Texas A&M | Mimic Learning | Interactive Visualization |
| | Rutgers | Model Induction | Bayesian Teaching |

| | **Learn a model** | **Explain decisions** | **Use the explanation** | |
|---|---|---|---|---|
| **Data Analytics**<br>Classification Learning Task | Multimedia Data | Explainable Model — Explanation Interface<br>→ Recommend<br>← Explanation | | An analyst is looking for items of interest in massive multimedia data sets |
| | Classifies items of interest in large data set | Explains why/why not for recommended items | Analyst decides which items to report, pursue | |
| **Autonomy**<br>Reinforcement Learning Task | ArduPilot & SITL Simulation | Explainable Model — Explanation Interface<br>→ Actions<br>← Explanation | | An operator is directing autonomous systems to accomplish a series of missions |
| | Learns decision policies for simulated missions | Explains behavior in an after-action review | Operator decides which future tasks to delegate | |

# Remaining Challenges for AI



Human Level AI

Perceiving
Learning
Abstracting
Reasoning

- Learning
  - Unsupervised learning
  - One-shot learning
  - Lifelong learning
  - Learning from instruction
- Understanding
  - Explanation
  - Representation and abstraction
- Human-like cognition
  - Planning and action
  - Meta-reasoning
  - Common Sense

www.darpa.mil

**MIT Technology Review**

**The Dark Secret at the Heart of AI**
Will Knight
April 11, 2017

**THE WALL STREET JOURNAL. WSJ**

**Inside DARPA's Push to Make Artificial Intelligence Explain Itself**
Sara Castellanos and Steven Norton
August 10, 2017

**The New York Times Magazine**

**Can A.I. Be Taught to Explain Itself?**
Cliff Kuang
November 21, 2017

Intelligent Machines Are Asked to Explain How Their Minds Work
Richard Waters
July 11, 2017

**FT FINANCIAL TIMES**

**The Register**

You better explain yourself, mister: DARPA's mission to make an accountable AI
Dan Robinson
September 29, 2017

**ExecutiveBiz**

Charles River Analytics-Led Team Gets DARPA Contract to Support Artificial Intelligence Program
Ramona Adams
June 13, 2017

**Entrepreneur**

Elon Musk and Mark Zuckerberg Are Arguing About AI -- But They're Both Missing the Point
Artur Kiulian
July 28, 2017

Team investigates artificial intelligence, machine learning in DARPA project
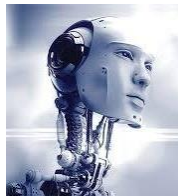Lisa Daigle
June 14, 2017

**Military EMBEDDED SYSTEMS**

**FAST COMPANY**

Why The Military And Corporate America Want To Make AI Explain Itself
Steven Melendez
June 22, 2017

**NOVA NEXT**

Ghosts in the Machine
Christina Couch
October 25, 2017

**Jane's**

DARPA's XAI seeks explanations from autonomous systems
Geoff Fein
November 16, 2017

**COMPUTERWORLD**

**Oracle quietly researching 'Explainable AI'**
George Nott
May 5, 2017

**SCIENTIFIC AMERICAN.**

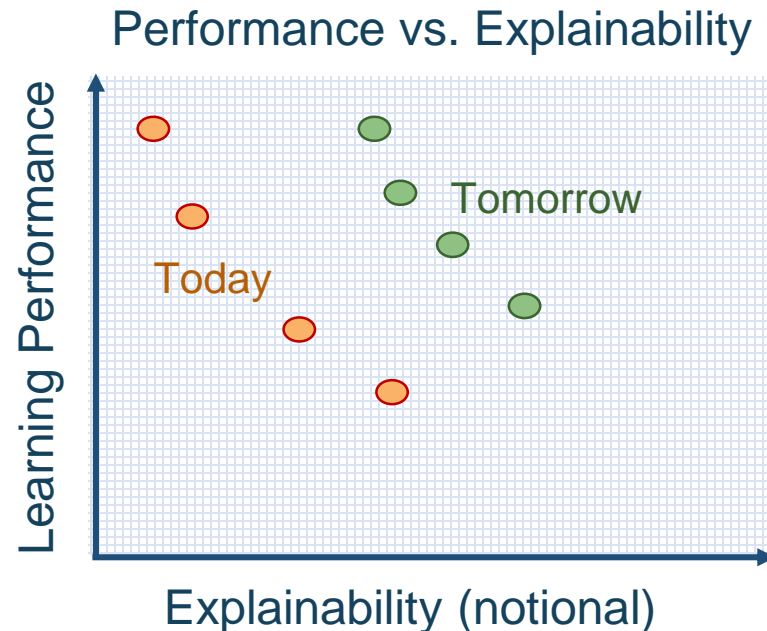Demystifying the Black Box That Is AI
Ariel Bleicher
August 9, 2017

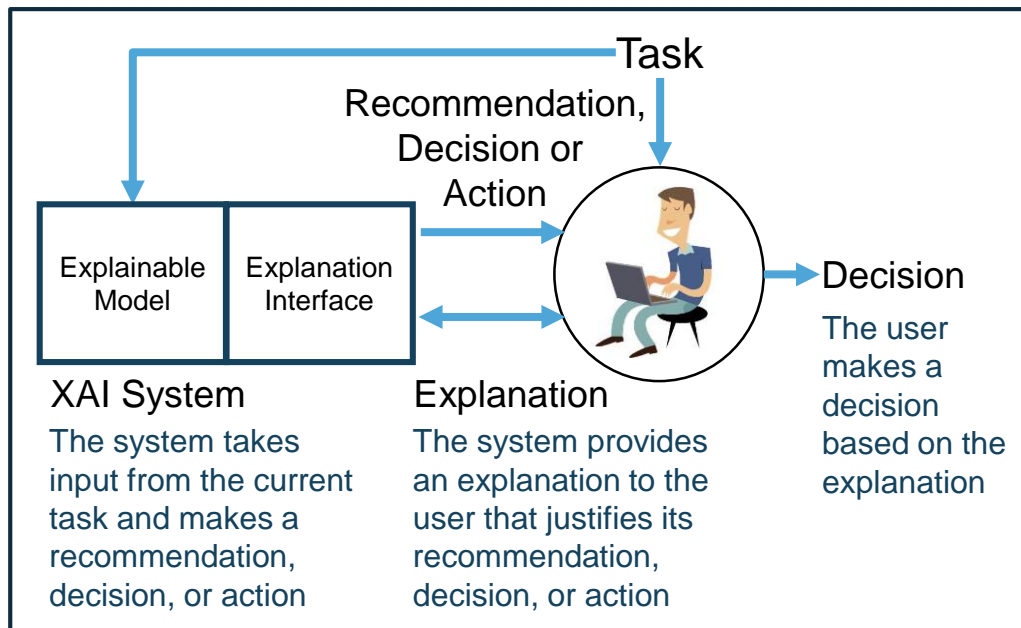How AI detectives are cracking open the black box of deep learning
Paul Voosen
July 6, 2017

**Science AAAS**

# Goal: Performance and Explainability

- XAI will create a suite of machine learning techniques that
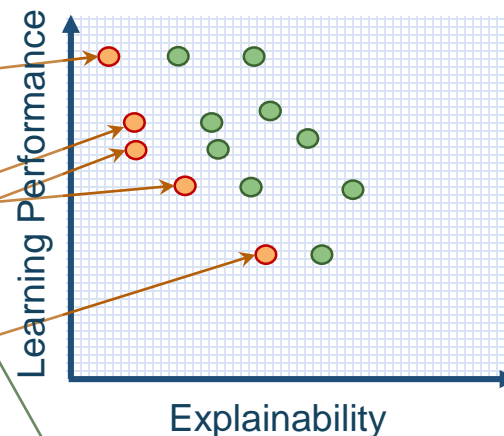  - Produce more explainable models, while maintaining a high level of learning performance (e.g., prediction accuracy)
  - Enable human users to understand, appropriately trust, and effectively manage the emerging generation of artificially intelligent partners

**Performance vs. Explainability**



Today

Tomorrow

Learning Performance

Explainability (notional)

## Explanation Framework



**Task**

**Recommendation, Decision or Action**

**XAI System**
The system takes input from the current task and makes a recommendation, decision, or action

**Explanation**
The system provides an explanation to the user that justifies its recommendation, decision, or action

Explainable Model

Explanation Interface

**Decision**
The user makes a decision based on the explanation

| Measure of Explanation Effectiveness |
|---|
| **User Satisfaction** |
| • Clarity of the explanation (user rating) <br> • Utility of the explanation (user rating) |
| **Mental Model** |
| • Understanding individual decisions <br> • Understanding the overall model <br> • Strength/weakness assessment <br> • 'What will it do' prediction <br> • 'How do I intervene' prediction |
| **Task Performance** |
| • Does the explanation improve the user's decision, task performance? <br> • Artificial decision tasks introduced to diagnose the user's understanding |
| **Trust Assessment** |
| • Appropriate future use and trust |
| **Correctability (Extra Credit)** |
| • Identifying errors <br> • Correcting errors <br> • Continuous training |

Learning Techniques (today)

Explainability (notional)

Neural Nets

Deep Learning

Graphical Models

Ensemble Methods

Bayesian Belief Nets

SRL

CRFs        HBNs

MLNs

Random Forests

Statistical Models

AOGs

SVMs

Markov Models

Decision Trees

Learning Performance

Explainability

**DARPA** · **XAI** EXPLAINABLE ARTIFICIAL INTELLIGENCE

## New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

## Learning Techniques (today)

Neural Nets

Graphical Models

Deep Learning

Bayesian Belief Nets

Ensemble Methods

SRL

CRFs    HBNs

Random Forests

Statistical Models

AOGs

MLNs

SVMs

Markov Models

Decision Trees

## Explainability (notional)

Learning Performance

Explainability



**Deep Explanation**
Modified deep learning techniques to learn explainable features



**Interpretable Models**
Techniques to learn more structured, interpretable, causal models



**Model Induction**
Techniques to infer an explainable model from any model as a black box

## Attention Mechanisms



**Top-down Caption Saliency**
[Ramanishka et al. CVPR17]

Caption: A **man** in a **jacket** is **standing** at the **slot machine**

## Modular Networks



**Neural module networks**
[Andreas et al.CVPR16,EMNLP16] [Hu et al. CVPR17]

Q: Can you park here?
NO Prediction
Neural module network
Attention visualization
Decision path

## Feature Identification



Generator
Target network
Tagger → Interpretation

## Learn to Explain



**Downy Woodpecker Definition**:
This bird has a white breast, black wings, and a red spot on its head.

CNN   RNN

**Image Explanation:**
This is a Downy Woodpecker because it is a black and wide bird with a red spot on its crown.

## Buildings

**56) building**

**120) arcade**

**8) bridge**

**123) building**

## Furniture

**18) billard table**

**155) bookcase**

**116) bed**

**38) cabinet**

## Indoor objects
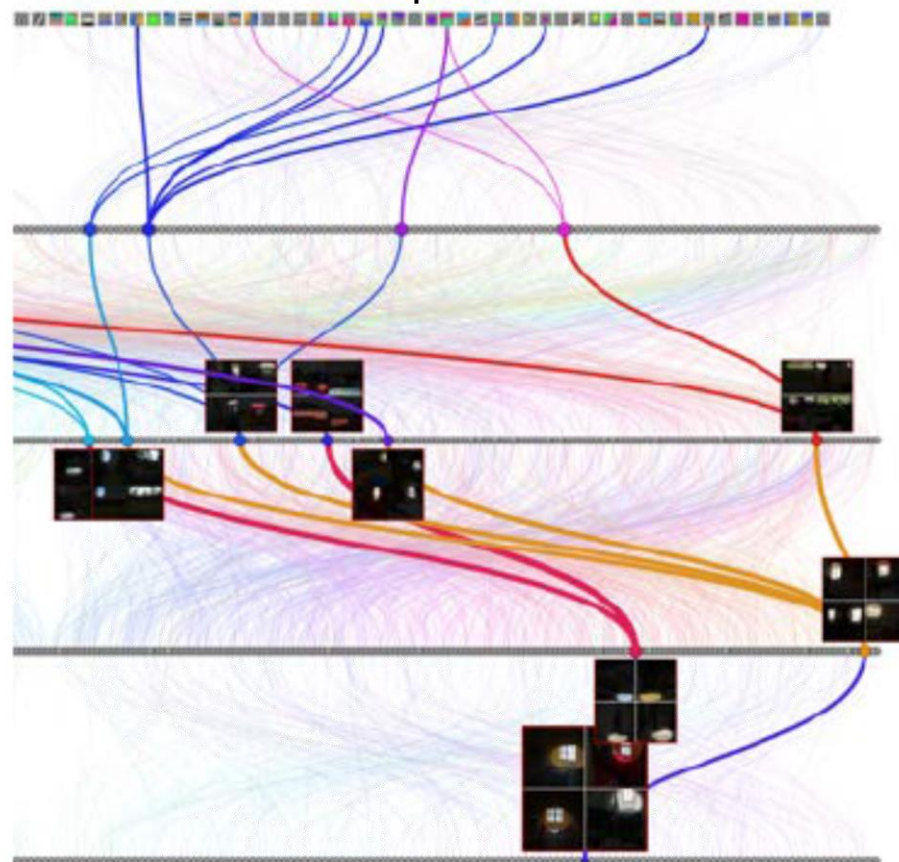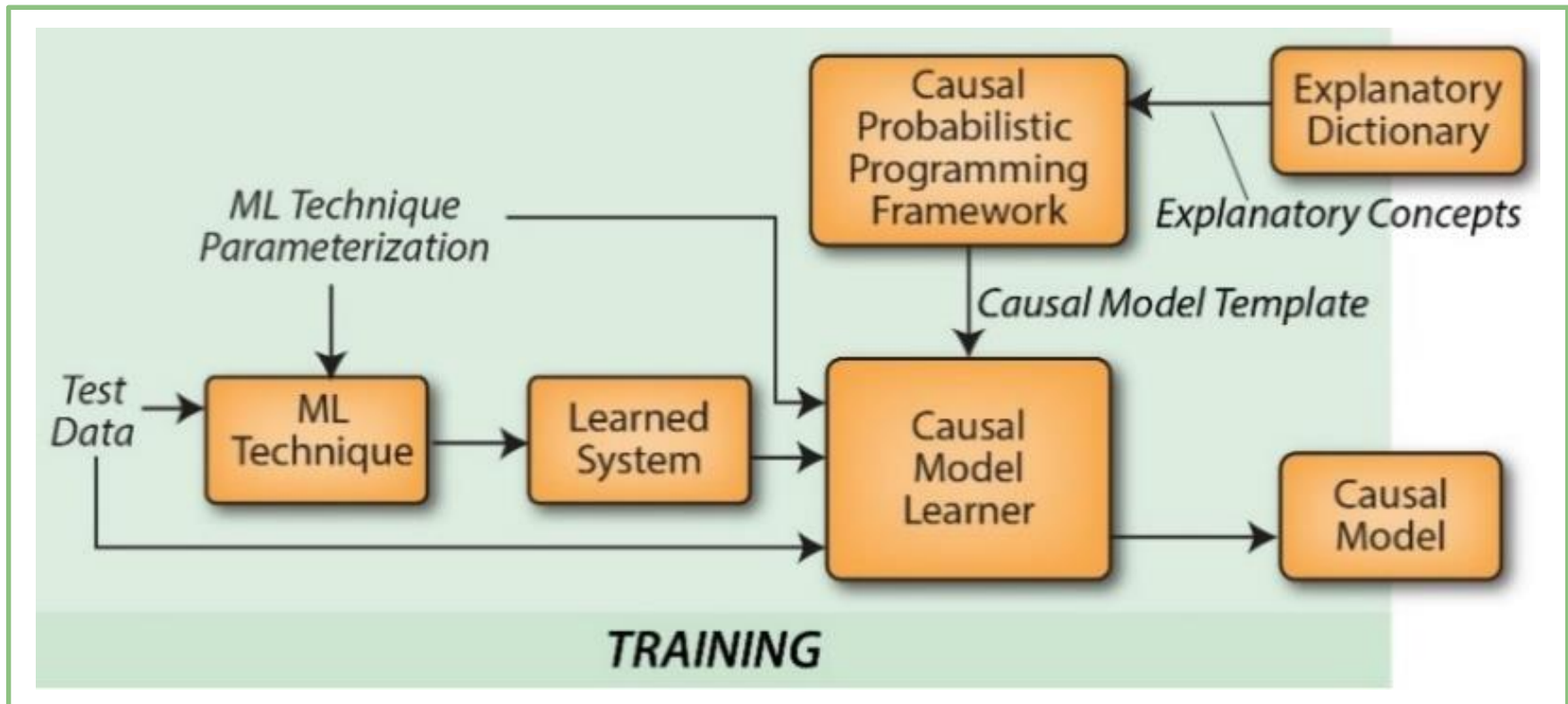
**182) food**

**46) painting**

**106) screen**

**53) staircase**

Interpretation of several units in pool5 of AlexNet trained for place recognition

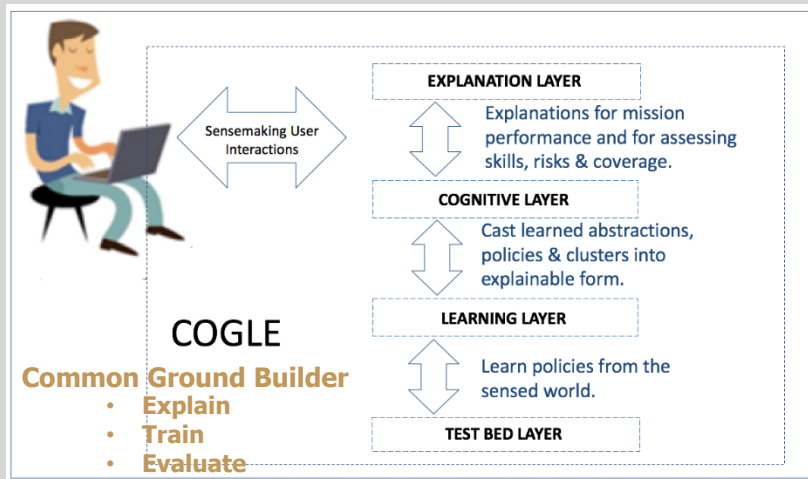Audit trail: for a particular output unit, the drawing shows the most strongly activated path

Causal Model Induction: Experiment with the learned model (as a grey box) to learn an explainable, causal, probabilistic programming model

## Common Ground Learning and Explanation (COGLE)

An interactive sensemaking system to explain the learned performance capabilities of a UAS flying in an ArduPilot simulation testbed



COGLE

**Common Ground Builder**
- **Explain**
- **Train**
- **Evaluate**

**Robotics Curriculum**

## Explanation-Informed Acceptance Testing of Deep Adaptive Programs (xACT)

Tools for explaining deep adaptive programs and discovering best principles for designing explanation user interfaces