# Explainable Artificial Intelligence (XAI)
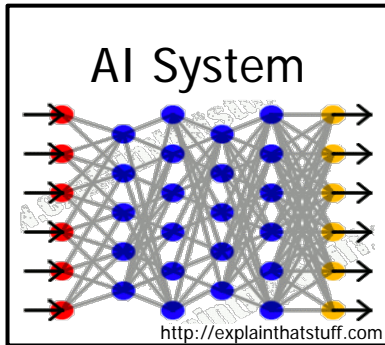
David Gunning

DARPA/I2O
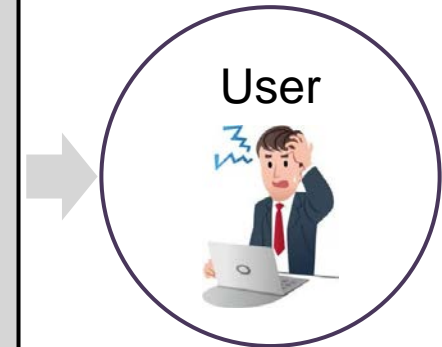
Program Update November 2017

# The Need for Explainable AI

AI System


Transportation


Finance


Security


Legal


Medicine


Military


User

- 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

- 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

**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

**FT**

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

**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
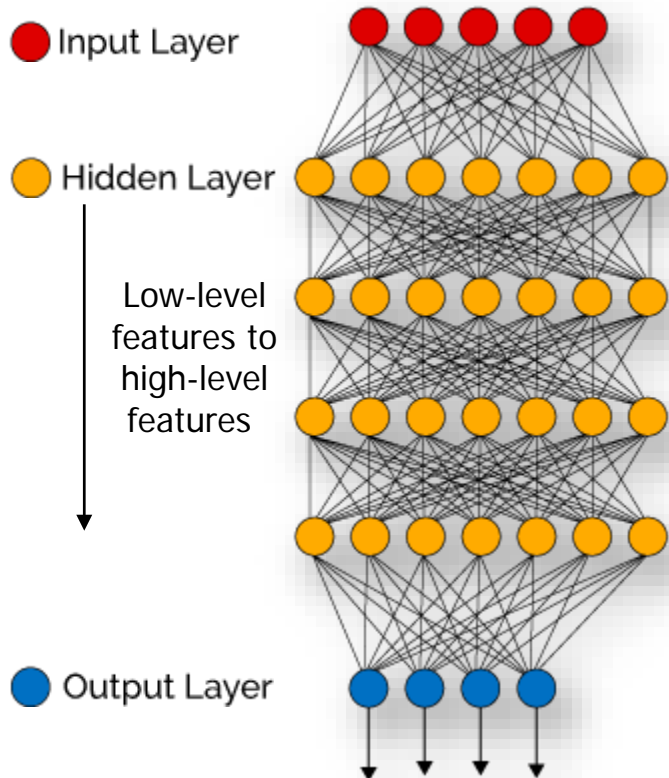May 5, 2017

**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**

Deep Learning Neural Network

- 🔴 Input Layer
- 🟠 Hidden Layer

Low-level features to high-level features

- 🔵 Output Layer

Automatic algorithm
(feature extraction and classification)

https://www.xenonstack.com/
XenonStack ©

Training Data

Input
(unlabeled image)

Neurons respond to simple shapes — **1st Layer**

Neurons respond to more complex structures — **2nd Layer**

Neurons respond to highly complex, abstract concepts — **nth Layer**

**10% WOLF**     **90% DOG**

http://fortune.com/
© 2018 Time Inc.

**DARPA**

**XAI** — EXPLAINABLE ARTIFICIAL INTELLIGENCE

## Today



©Spin South West

Training Data → Learning Process → Learned Function → **This is a cat** (p = .93) — Output → User with a Task

©University Of Toronto

http://explainthatstuff.com

- 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



©Spin South West

Training Data → New Learning Process → Explainable Model →

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

Explanation Interface → User with a Task

©University Of Toronto

- 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

| | **Learn a model** | **Explain decisions** | **Use the explanation** | |
|---|---|---|---|---|
| **Data Analytics**<br><br>Classification Learning Task | Two trucks performing a loading activity<br>©Air Force Research Lab<br>**Multimedia Data** | Explainable Model — Explanation Interface<br>→ Recommend<br>← Explanation | ©Getty Images | 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><br>Reinforcement Learning Task | ©ArduPilot.org<br>**ArduPilot & SITL Simulation** | Explainable Model — Explanation Interface<br>→ Actions<br>← Explanation | ©US Army | 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 | |

- 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



**Learning Performance** (y-axis)

**Explainability (notional)** (x-axis)

Today

Tomorrow

## Explanation Framework



**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

**Decision**

The user makes a decision based on the explanation

## Measure of Explanation Effectiveness

**User Satisfaction**

- Clarity of the explanation (user rating)
- Utility of the explanation (user rating)

**Mental Model**

- Understanding individual decisions
- Understanding the overall model
- Strength/weakness assessment
- 'What will it do' prediction
- 'How do I intervene' prediction

**Task Performance**

- Does the explanation improve the user's decision, task performance?
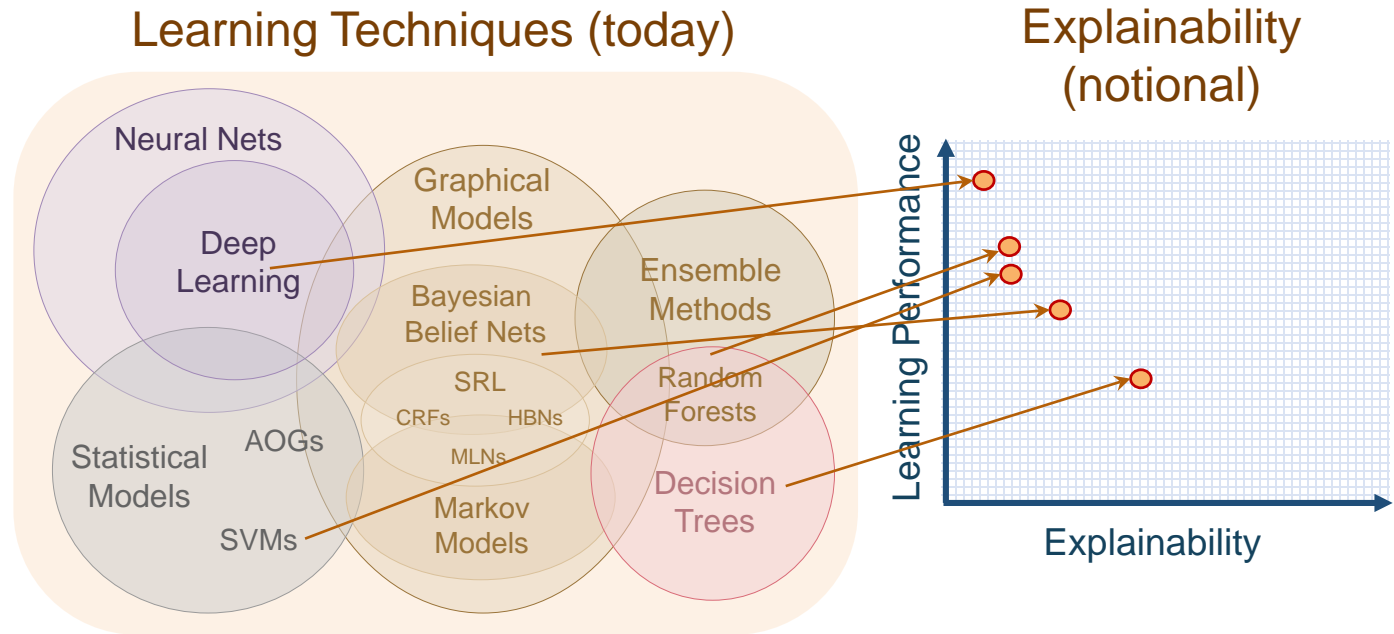- Artificial decision tasks introduced to diagnose the user's understanding

**Trust Assessment**

- Appropriate future use and trust

**Correctability (Extra Credit)**

- Identifying errors
- Correcting errors
- Continuous training

Learning Techniques (today)

Explainability (notional)

Neural Nets

Graphical Models

Deep Learning

Bayesian Belief Nets

Ensemble Methods

SRL

CRFs          HBNs

Random Forests

MLNs

Statistical Models

AOGs

SVMs

Markov Models
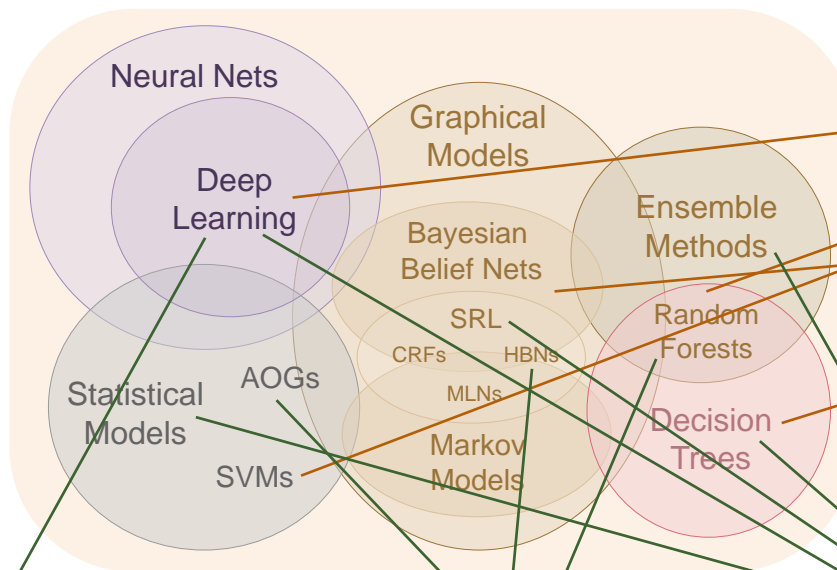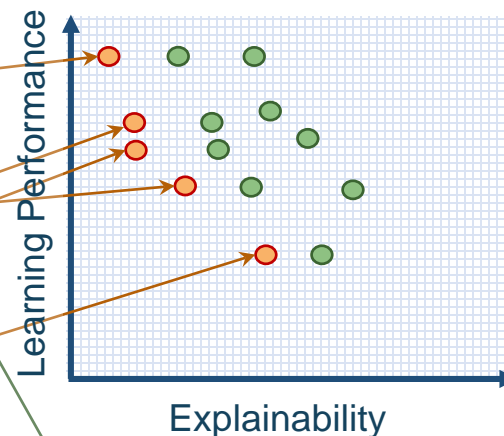
Decision Trees

Learning Performance

Explainability

**New Approach**

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance
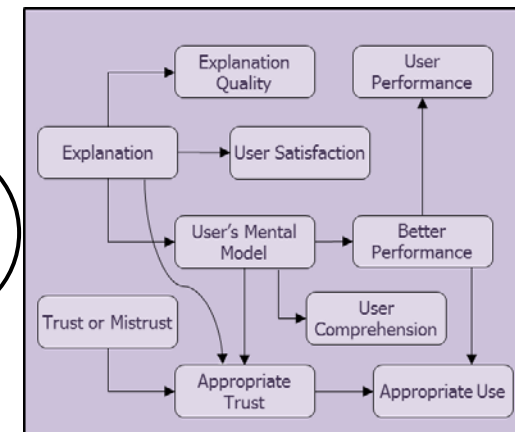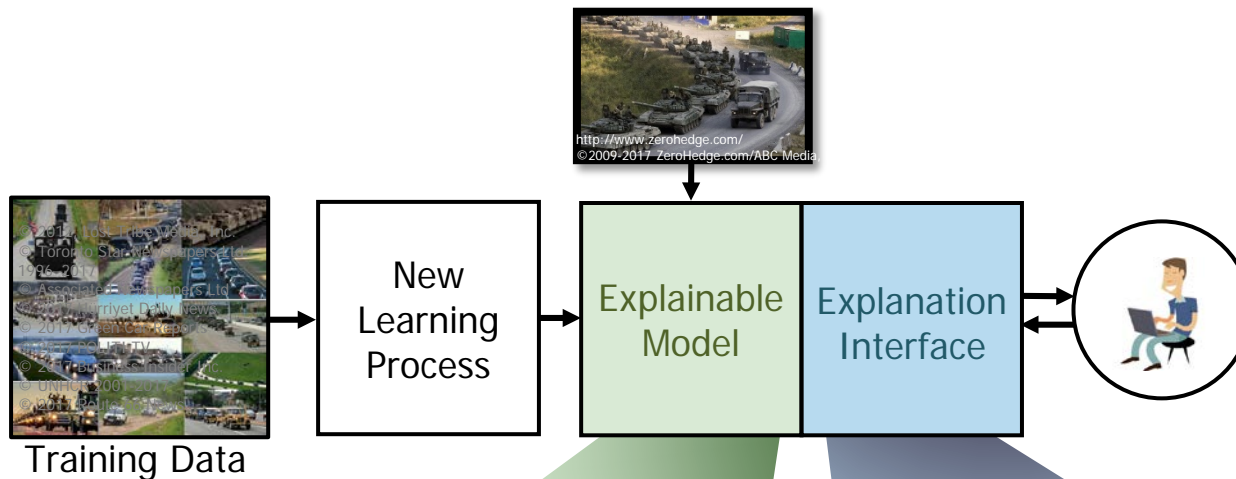
**Learning Techniques (today)**

Neural Nets

Deep Learning

Graphical Models

Bayesian Belief Nets

Ensemble Methods

SRL

CRFs    HBNs

MLNs

Random Forests

Statistical Models

AOGs

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

A₁

Model

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

?  ?    Model

Experiment

Training Data

New Learning Process

Explainable Model

Explanation Interface

http://www.zerohedge.com/
©2009-2017 ZeroHedge.com/ABC Media

**IHMC**
Psychological Model of Explanation

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

## Attention Mechanisms



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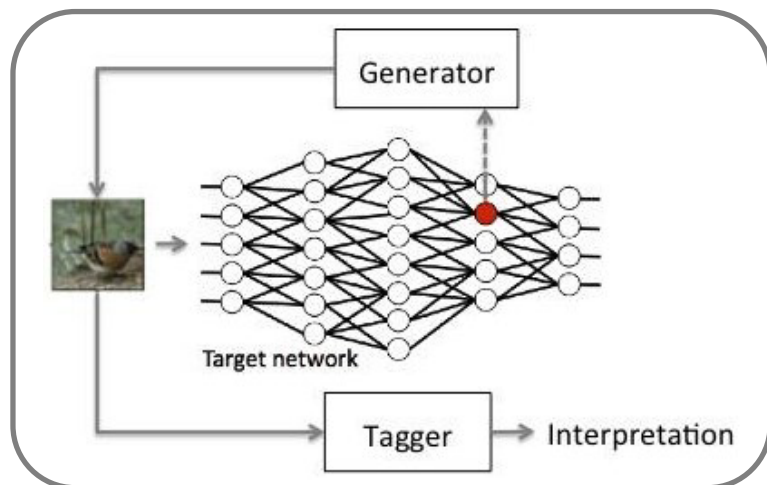[Ramanishka et al. CVPR17]

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 white bird with a **red spot** on its crown.
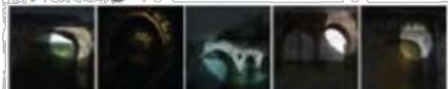
## Buildings

56) building

120) arcade

8) bridge

123) building

## Furniture

18) billard table

155) bookcase

116) bed

38) cabinet

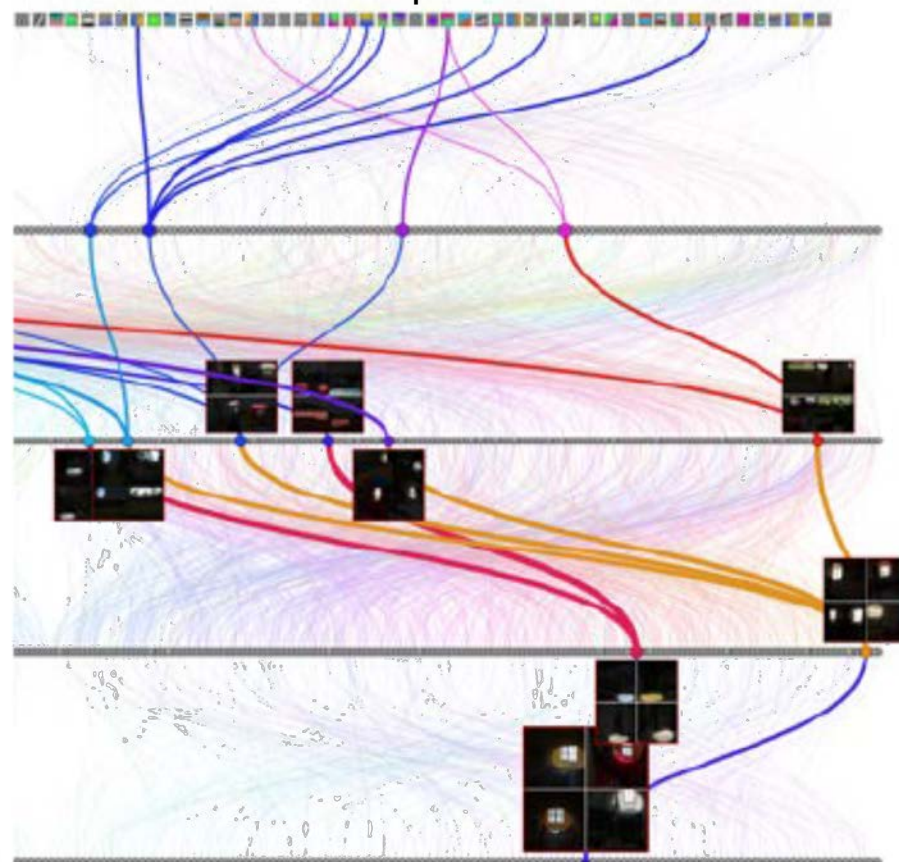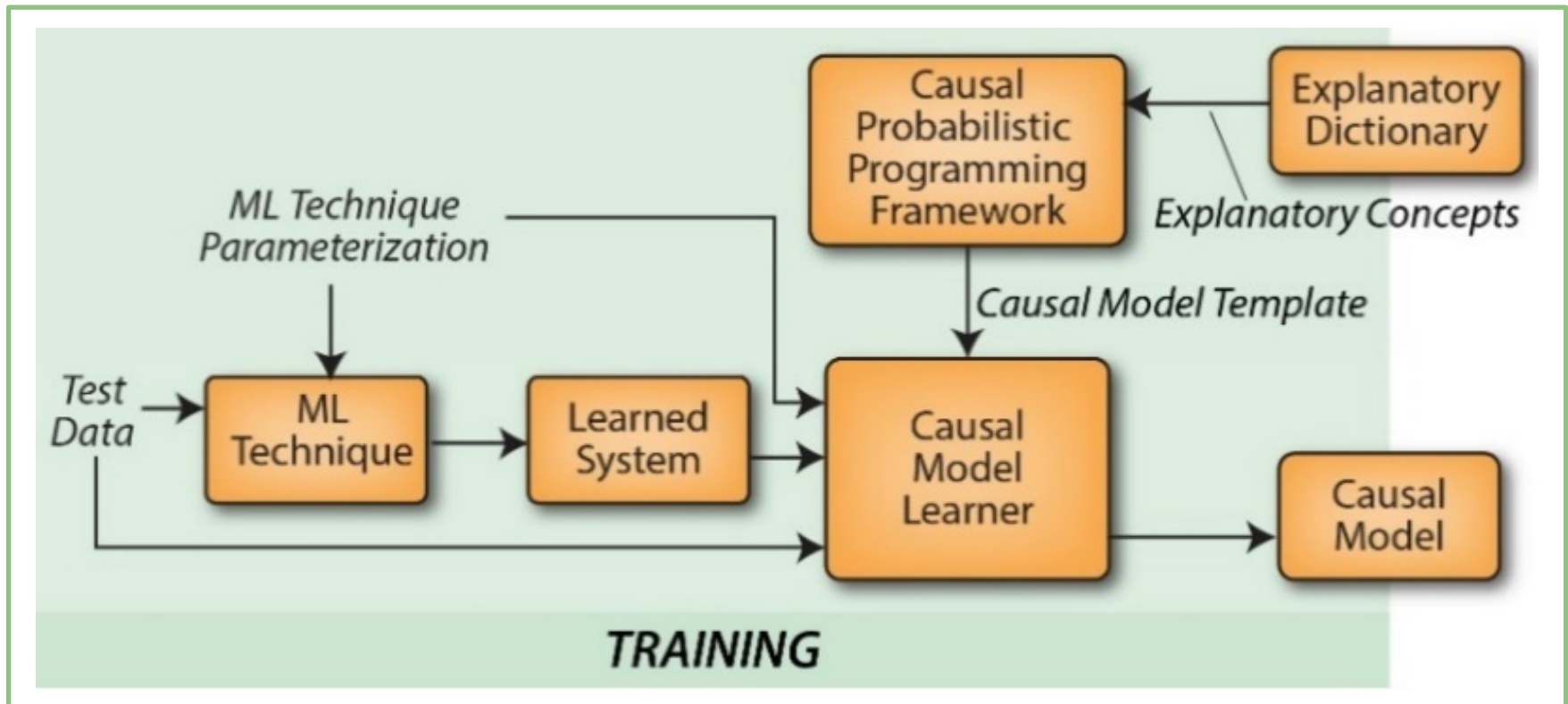## Indoor objects

182) food
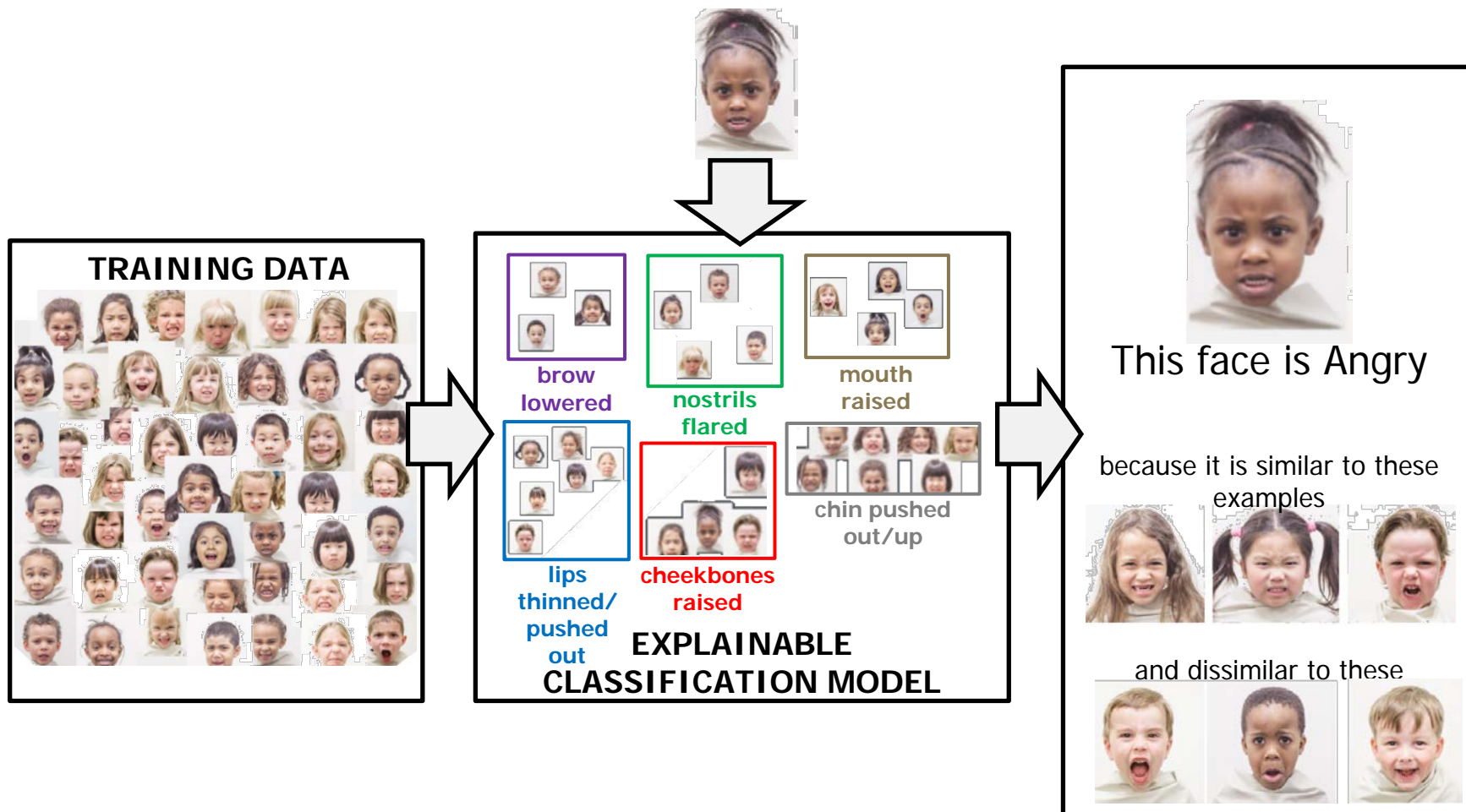
46) painting

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

**DARPA**

**XAI** EXPLAINABLE ARTIFICIAL INTELLIGENCE



**TRAINING DATA**

**EXPLAINABLE CLASSIFICATION MODEL**

- brow lowered
- nostrils flared
- mouth raised
- lips thinned/pushed out
- cheekbones raised
- chin pushed out/up

This face is Angry

because it is similar to these examples

and dissimilar to these

**BAYESIAN TEACHING** for optimal selection of examples for machine explanation

## 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
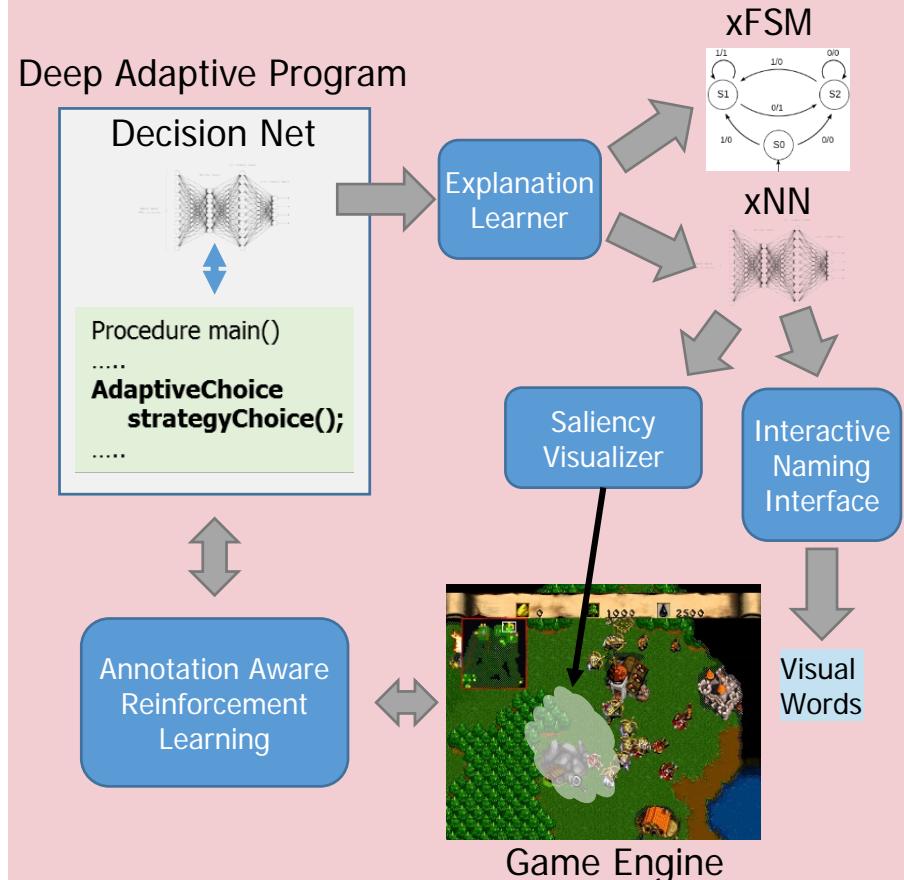


EXPLANATION LAYER
Explanations for mission performance and for assessing skills, risks & coverage.

COGNITIVE LAYER
Cast learned abstractions, policies & clusters into explainable form.

LEARNING LAYER
Learn policies from the sensed world.

TEST BED LAYER

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

Deep Adaptive Program

xFSM



Decision Net

Procedure main()
.....
**AdaptiveChoice
    strategyChoice();**
.....

Explanation Learner

xNN

Saliency Visualizer

Interactive Naming Interface

Annotation Aware Reinforcement Learning

Visual Words

Game Engine

## Analytic (didactic) statements

in natural language that describe the elements and context that support a choice

## Visualizations

that directly highlight portions of the raw data that support a choice and allow viewers to form their own perceptual understanding

### Explanation Modes

## Cases

that invoke specific examples or stories that support the choice

## Rejections of alternative choices

(or "common misconceptions" in pedagogy) that argue against less preferred answers based on analytics, cases, and data

- **TA1: Explainable Learners**
  - Multiple TA1 teams will develop prototype explainable learning systems that include both an explainable model and an explanation interface

- **TA2: Psychological Model of Explanation**
  - At least one TA2 team will summarize current psychological theories of explanation and develop a computational model of explanation from those theories

## Analytics

### Visual Question Answering


**MovieQA**


**CLEVR**

### Activity Recognition


**ActivityNet**

## Autonomy

### Strategy Games


**Starcraft2**


**ELF-MiniRTS**

### Vehicle Control


**ArduPilot**


**Driving Simulator**

Model of the Explanation
Process and Possible Metrics

XAI
Process

XAI
Metrics

System

User — receives → Explanation — revises → User's Mental Model — enables → Better Performance

User — may initially → Trust or Mistrust

Explanation — is assessed by → "Goodness" Criteria

Explanation — is assessed by → Test of Satisfaction

User's Mental Model — is assessed by → Test of Comprehension

Better Performance — is assessed by → Test of Performance

Explanation — can engender → Appropriate Trust

Trust or Mistrust — gives way to → Appropriate Trust

Appropriate Trust — enables → Appropriate Use

Better Performance — involves → Appropriate Use

# Schedule and Milestones

| | 2017 | 2018 | 2019 | 2020 | 2021 |
|---|---|---|---|---|---|
| **PHASE 1: Technology Demonstrations** | | | **PHASE 2: Comparative Evaluations** | | |
| **Evaluator** | Define Evaluation Framework · Prep for Eval 1 · Eval 1 · Analyze Results | Prep for Eval 2 · Eval 2 · Analyze Results | Prep for Eval 3 · Eval 3 | Analyze Results & Accept Toolkits | |
| **TA 1** | Develop & Demonstrate Explainable Models (against proposed problems) · Eval 1 | Refine & Test Explainable Learners (against common problems) · Eval 2 | Refine & Test Explainable Learners (against common problems) · Eval 3 | Deliver Software Toolkits | |
| **TA 2** | Summarize Current Psychological Theories of Explanation | Develop Computational Model of Explanation | Refine & Test Computational Model | Deliver Computational Model | |
| **Meetings** | KickOff · Progress Report · Tech Demos | Eval 1 Results | Eval 2 Results | | Final |

Meeting dates: May 9-11 · Nov 6-8 · May 7-9

- Technical Area 1 (Explainable Learners) Milestones
  - Demonstrate the explainable learners against problems proposed by the developers (Phase 1)
  - Demonstrate the explainable learners against common problems (Phase 2)
  - Deliver software libraries and toolkits (at the end of Phase 2)
- Technical Area 2 (Psychology of Explanation) Milestones
  - Deliver an interim report on psychological theories (after 6 months during Phase 1)
  - Deliver a final report on psychological theories (after 12 months, during Phase 1)
  - Deliver a computational model of explanation (after 24 months, during Phase 2)
  - Deliver the computational model software (at the end of Phase 2)

| CP | Performer | Explainable Model | Explanation Interface |
|---|---|---|---|
| **Both** | UC Berkeley | Deep Learning | Reflexive and 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 and Tell Explanation |
| | Raytheon BBN | Deep Learning | Argumentation and Pedagogy |
| | UT Dallas | Probabilistic Logic | Decision Diagrams |
| | Texas A&M | Mimic Learning | Interactive Visualization |
| | Rutgers | Model Induction | Bayesian Teaching |

# Deeply Explainable Artificial Intelligence

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Deep Learning**<br><br>• Explain *implicit* (latent) nodes by training additional DL models<br>• Explain *explicit* nodes thru Neural Module Networks (NMNs) | **Reflexive & Rational**<br><br>• Reflexive explanations (that arise directly from the model)<br>• Rational explanations (that come from reasoning about user's beliefs) | **Autonomy**<br><br>• ArduPilot and OpenAI Gym Simulations<br><br>**Data Analytics**<br><br>• Visual QA and Multimedia Event QA |

- **PI**: Trevor Darrell (Berkeley)

- Pieter Abbeel (Berkeley)
- Tom Griffiths (Berkeley)
- Kate Saenko (BU)
- Zeynep Akata (U. Amsterdam)

- Dan Klein (Berkeley)
- John Canny (Berkeley)
- Anca Dragan (Berkeley)

- Anthony Hoogs (Kitware)

# CAMEL: Causal Models to Explain Learning

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Model Induction Causal Models**<br><br>• Experiment with the learned model (as a grey box) to learn an explainable, causal, probabilistic programming model | **Narrative Generation**<br><br>• Interactive visualization based on the generation of temporal, spatial narratives from the causal, probabilistic models | **Autonomy**<br><br>• Minecraft, Starcraft<br><br>**Data Analytics**<br><br>• Pedestrian Detection (INRIA), Activity Recognition (ActivityNet) |

• **PI**: Brian Ruttenberg (CRA)

• Avi Pfeffer (CRA)
• David Jensen (U. Mass)
• Michael Littman (Brown)

• James Niehaus (CRA)
• Emilie Roth (Roth Cognitive Engineering)
• Joe Gorman(CRA)
• James Tittle (CRA)

# Learning and Communicating Explainable Representations for Analytics and Autonomy

## Explainable Model

### Pattern Theory+

- Integrated representation across an entropy spectrum:
  - Deep Neural Nets
  - Stochastic And-Or-Graphs (AOG)
  - Predicate Calculus

## Explanation Interface

### 3-Level Explanation

- Integrate 3 levels of explanation:
  - Concept compositions
  - Causal and counterfactual reasoning
  - Utility explanations

## Challenge Problem

### Autonomy

- Humanoid robot behavior and VR simulation platform

### Data Analytics

- Understanding complex multimedia events

- **PI**: Song-Chun Zhu (UCLA)

- Ying Nian Wu (UCLA)
- Sinisa Todorovic (OSU)
- Joyce Chai (Michigan State)

# xACT: Explanation-Informed Acceptance Testing of Deep Adaptive Programs

## Explainable Model

### Adaptive Programs

- Explainable Deep Adaptive Programs (xDAPs) – a new combination of Adaptive Programs, Deep Learning and explainability

## Explanation Interface

### Acceptance Testing

- Provides a visual & NL explanation interface for acceptance testing by test pilots based on Information Foraging Theory

## Challenge Problem

### Autonomy

- Real-Time Strategy Games based on custom designed game engine designed to support explanation
- Possible use of Starcraft

- **PI**: Alan Fern (OSU)

- Tom Dietterich (OSU)
- Fuxin Li (OSU)
- Prasad Tadepalli (OSU)
- Weng-Keen Wong (OSU)

- Margaret Burnett (OSU)
- Martin Erwig (OSU)
- Liang Huang (OSU)

# COGLE: Common Ground Learning and Explanation

## Explainable Model

### Cognitive Model

- 3-layer architecture:
  - Learning Layer (DNNs)
  - Cognitive Layer (ACT-R Cog. Model)
  - Explanation Layer (HCI)

## Explanation Interface

### Interactive Training

- Interactive visualization of states, actions, policies & values
- Includes a module for test pilots to refine and train the system

## Challenge Problem

### Autonomy

- ArduPilot simulation environment
- *Value of Explanation* (VoE) framework for measuring explanation effectiveness

- **PI**: Mark Stefik (PARC)

- Honglak Lee (U. Mich.)
- Subramanian Ramamoorthy (U. Edinburgh)

- Christian Lebiere (CMU)
- John Anderson (CMU)
- Robert Thomson (USMA)

- Michael Youngblood (PARC)

# XRL: Explainable Reinforcement Learning for AI Autonomy

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **XRL Models** | **XRL Interaction** | **Autonomy** |
| • Create a new scientific discipline for Explainable Reinforcement Learning with work on new algorithms and representations | • Interactive explanations of dynamic systems<br>• Human-machine interaction to improve performance | • Open AI Gym<br>• Autonomy in the electrical grid<br>• Mobile service robots<br>• Self-improving educational software |

• **PI**: Geoff Gordon (CMU)

• Zico Kolter (CMU)
• Pradeep Ravikumar (CMU)

• Manuela Veloso (CMU)
• Emma Brunskill (Stanford)

# DARE: Deep Attention-based Representations for Explanation

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Deep Learning** | **Show-and-Tell Explanations** | **Data Analytics** |
| • Multiple deep learning techniques:<br>  • Attention-based mechanisms<br>  • Compositional NMNs<br>  • GANs | • DNN visualization<br>• Query evidence that explains DNN decisions<br>• Generate natural language justifications | • Visual Question Answering (VQA) using Visual Gnome, Flickr30<br>• MovieQA |

- **PIs**: Giedrius Burachas (SRI), Mohamed Amer (SRI)

- Shalini Ghosh (SRI)
- Avi Ziskind (SRI)
- Michael Wessel (SRI)

- Richard R. Zemel (U. Toronto)
- Sanja Fidler (U. Toronto)
- David Duvenaud (U. Toronto)
- Graham Taylor (U. Guelph)

- Jürgen Schulze (UCSD)

## EQUAS: Explainable QUestion Answering System

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Deep Learning** | **Argumentation Theory** | **Data Analytics** |
| • Semantic labelling of DNN neurons<br>• DNN audit trail construction<br>• Gradient-weighted Class Activation Mapping | • Comprehensive strategy based on argumentation theory<br>• NL generation<br>• DNN visualization | • Visual Question Answering (VQA), beginning with images and progressing to video |

- **PI**: William Ferguson (Raytheon BBN)

- Antonio Torralba (MIT)
- Ray Mooney (UT Austin)
- Devi Parikh (GA Tech)
- Dhruv Batra (GA Tech)

# Tractable Probabilistic Logic Models: A New, Deep Explainable Representation

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Probabilistic Logic** | **Probabilistic Decision Diagrams** | **Data Analytics** |
| • Tractable Probabilistic Logic Models (TPLMs) – an important class of (non-deep learning) interpretable models | • Enables users to explore and correct the underlying model as well as add background knowledge | • Infer activities in multimodal data (video and text) <br> • Using the Wetlab (biology) and TACoS (cooking) datasets |

• **PI**: Vibhav Gogate (UT Dallas)

• Adnan Darwiche (UCLA)    • Eric Ragan (Texas A&M)
• Guy Van Den Broeck (UCLA)    • Parag Singla (IIT-Delhi)
• Nicholas Ruozzi (UT Dallas)

# Transforming Deep Learning to Harness the Interpretability of Shallow Models: An Interactive End-to-End System

## Explainable Model

### Mimic Learning

- Develop a mimic learning framework that combines deep learning models for prediction and shallow models for explanations

## Explanation Interface

### Interactive Visualization

- Interactive visualization over multiple views, using heat maps & topic modeling clusters to show predictive features

## Challenge Problem

### Data Analytics

- Multiple tasks using data from Twitter, Facebook, ImageNet, UCI, NIST and Kaggle
- Metrics for explanation effectiveness

- **PI**: Xia Hu (Texas A&M)

- Shuiwang Ji (Wash. State)  • Eric Ragan (Texas A&M)

# Model Explanation by Optimal Selection of Teaching Examples

## Explainable Model

### Model Induction

- Select the optimal training examples to explain model decisions based on Bayesian Teaching

## Explanation Interface

### Bayesian Teaching

- Example-based explanation of:
  - the full model
  - user-selected sub-structure
  - user submitted examples

## Challenge Problem

### Data Analytics

- Movie descriptions
- Image processing
- Caption data
- Movie events
- Human motion events

- **PI**: Patrick Shafto (Rutgers)

- Scott Cheng-Hsin Yang (Rutgers)

# Naturalistic Decision Making Foundations of Explainable AI

| Literature Review | Computational Model | Model Validation |
|---|---|---|
| **Naturalistic Theory** | **Bayesian Framework** | **Experiments** |
| • Extensive review of relevant psychological theories<br>• Extend the theory of Naturalistic Decision Making to cover explanation | • Represent reductionist mental models that humans develop as part of the explanatory process<br>• Including mental simulation | • Conduct interactive assessment and formal human experiments<br>• Validate the model<br>• Develop metrics of explanation effectiveness |

- **PI**: Robert R. Hoffman (IHMC)

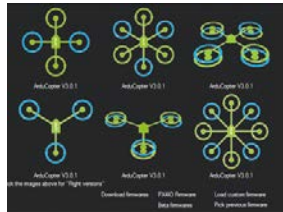| | | |
|---|---|---|
| • Gary Klein (MacroCognition)<br>• Shane T. Mueller (Michigan Tech) | • William J. Clancey (IHMC)<br>• COL Timothy M. Cullen (SAASS) | • Jordan Litman (IHMC Psychometrician)<br>• Simon Attfield (Middlesex University-London)<br>• Peter Pirolli (IHMC) |

# XAI Evaluation

## Challenge Problems

Analytics
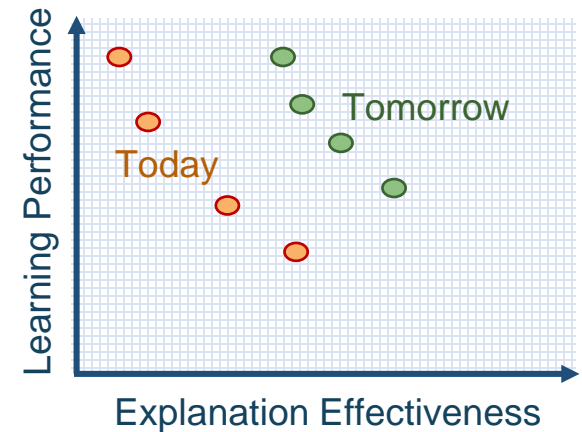


Two trucks performing a loading activity

Autonomy



## Evaluation Framework

- Evaluation protocols
- Training environment
  - Training data
  - Simulation environment
- Testing environment
  - Subjects
  - Web infrastructure
- Baseline systems

## Measurement



Learning Performance (vertical axis)

Explanation Effectiveness (horizontal axis)

Today

Tomorrow

---

- **PI**: David Aha (NRL)

---

- Justin Karneeb (Knexus)
- Matt Molineaux (Knexus)
- Leslie Smith (NRL)

- Mike Pazzani (UC Riverside)

www.darpa.mil