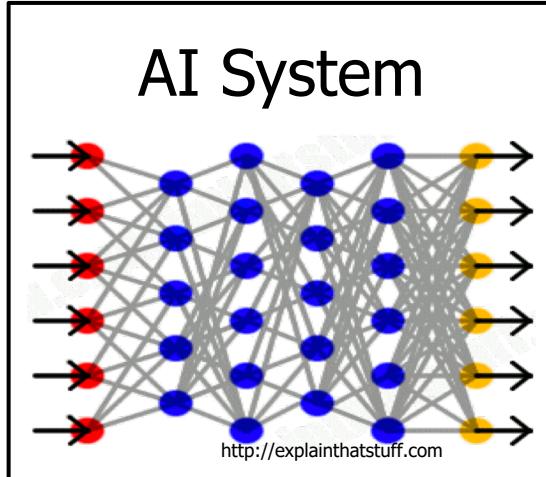


# Explainable Artificial Intelligence (XAI)

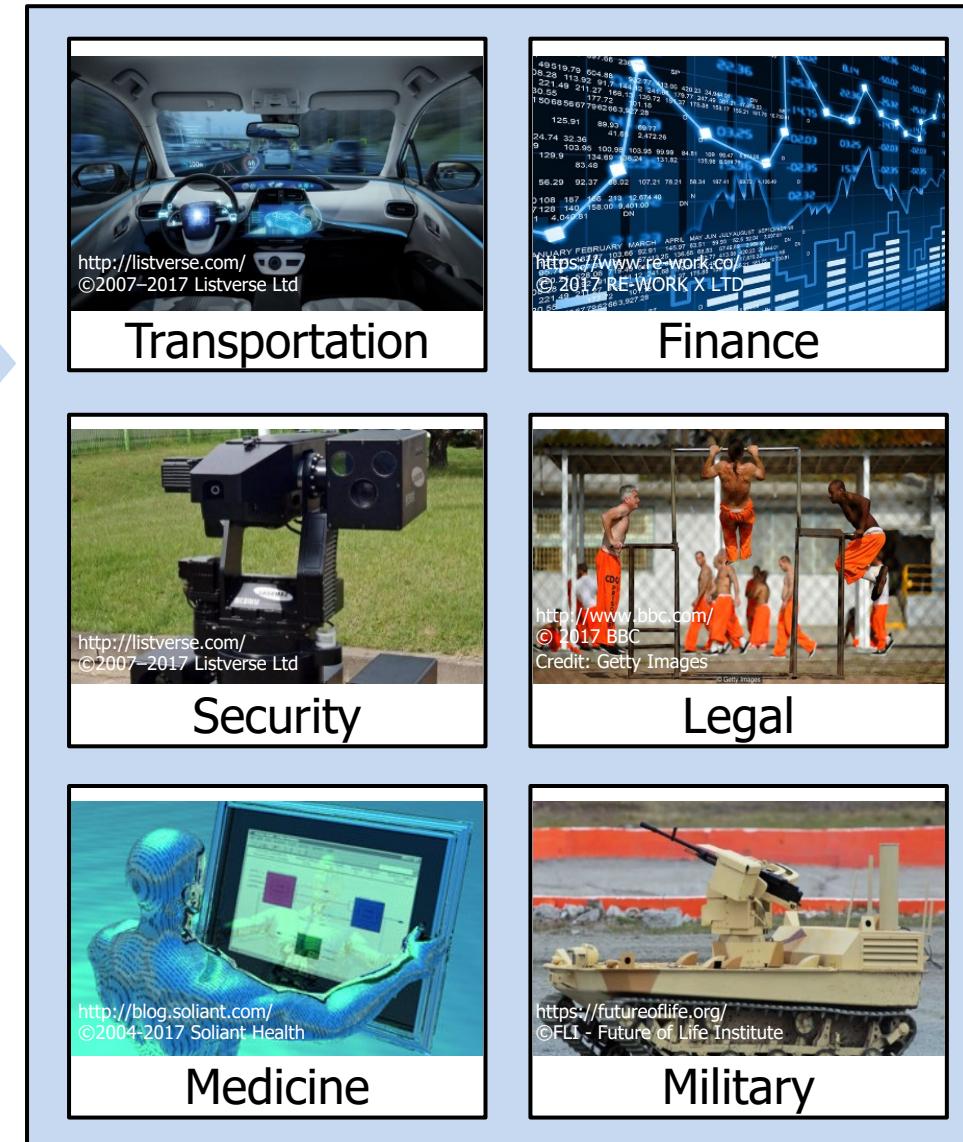


David Gunning  
Information Innovation Office (I2O)  
Defense Advanced Research Projects Agency (DARPA)

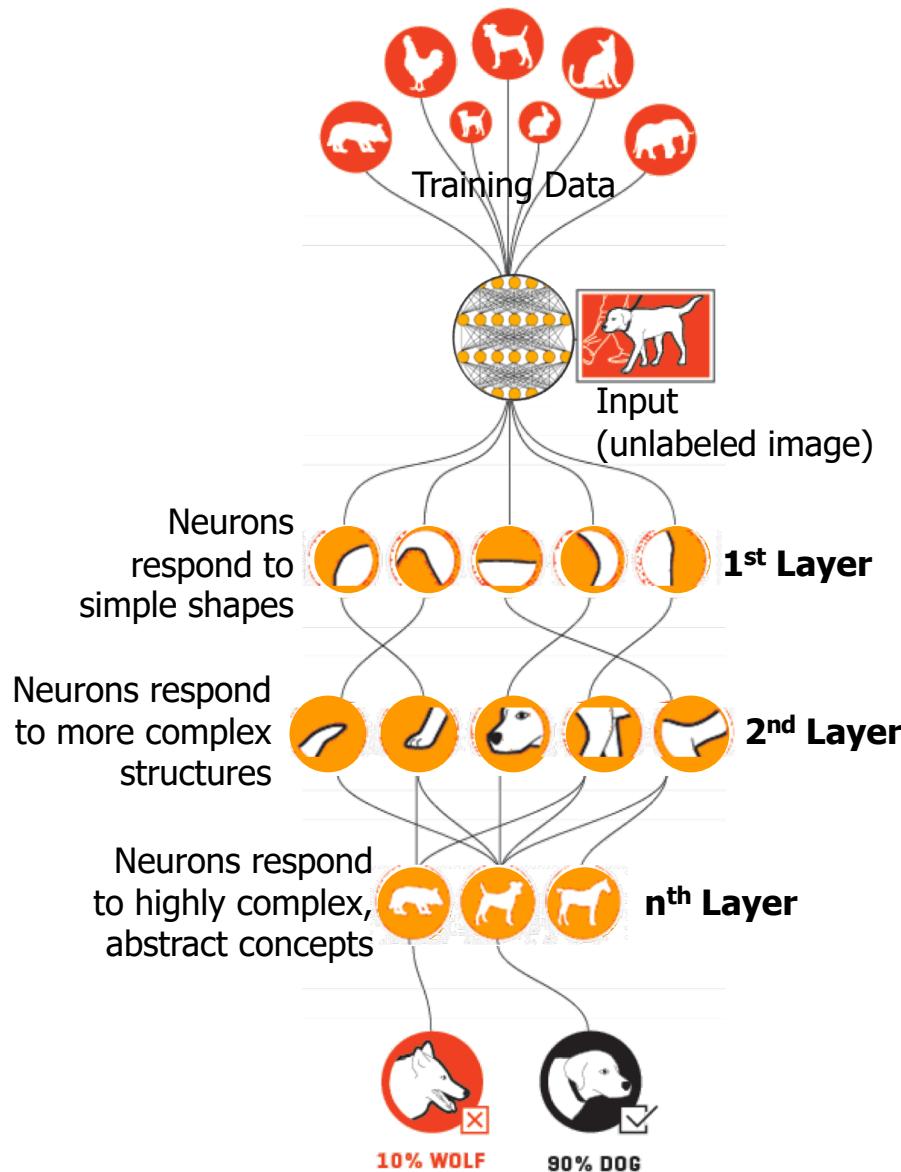
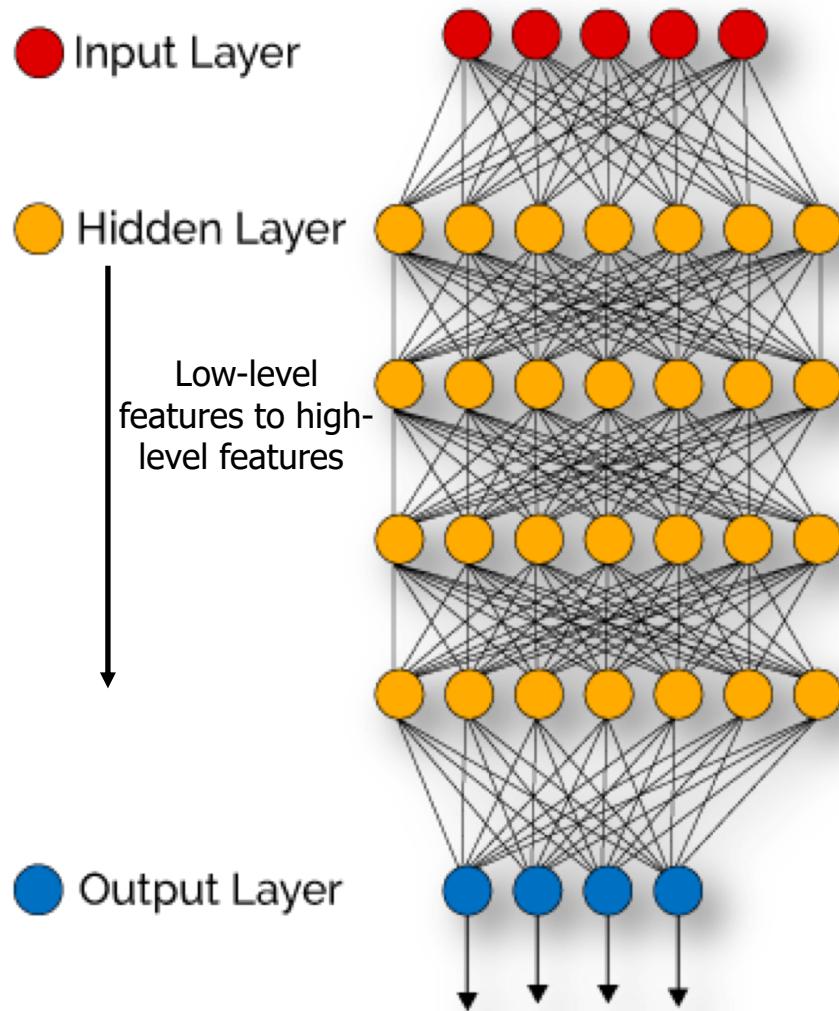




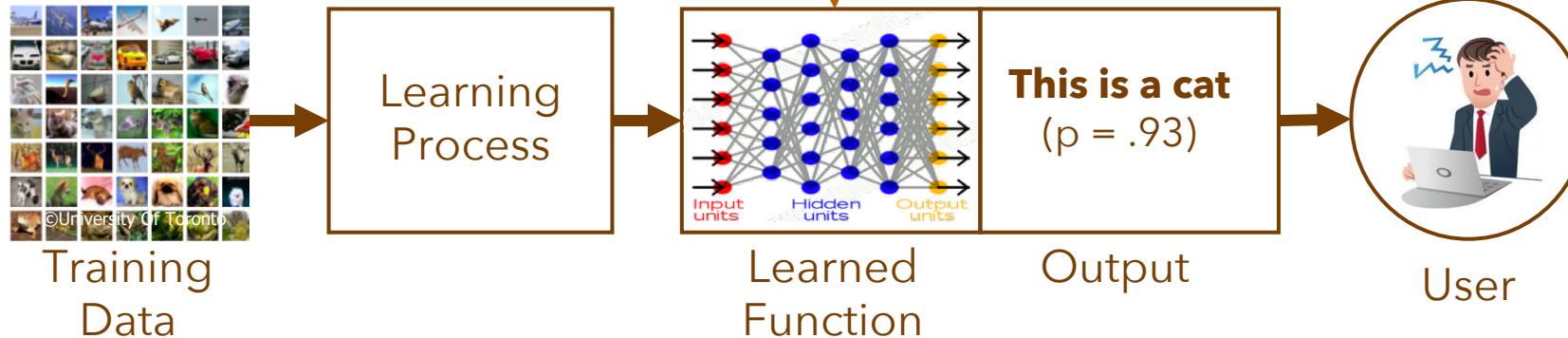
- We are entering a new age of AI applications
- Machine learning is the core technology
- Machine learning models are opaque, non-intuitive, and difficult for people to understand



- 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?

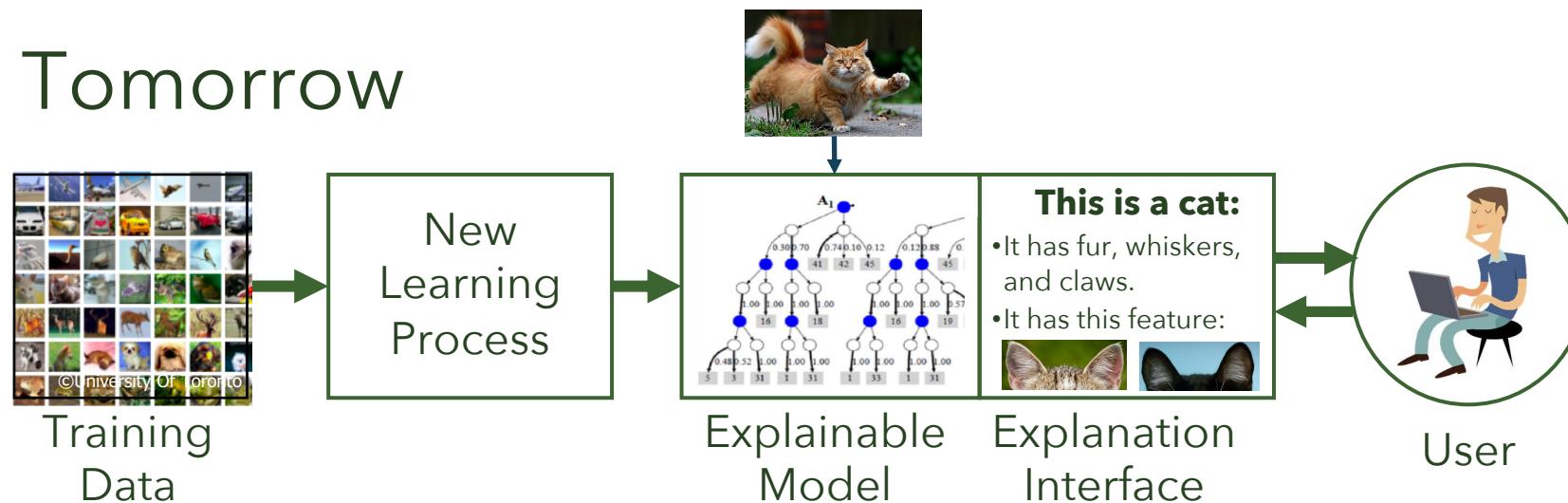


## Today



- 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?

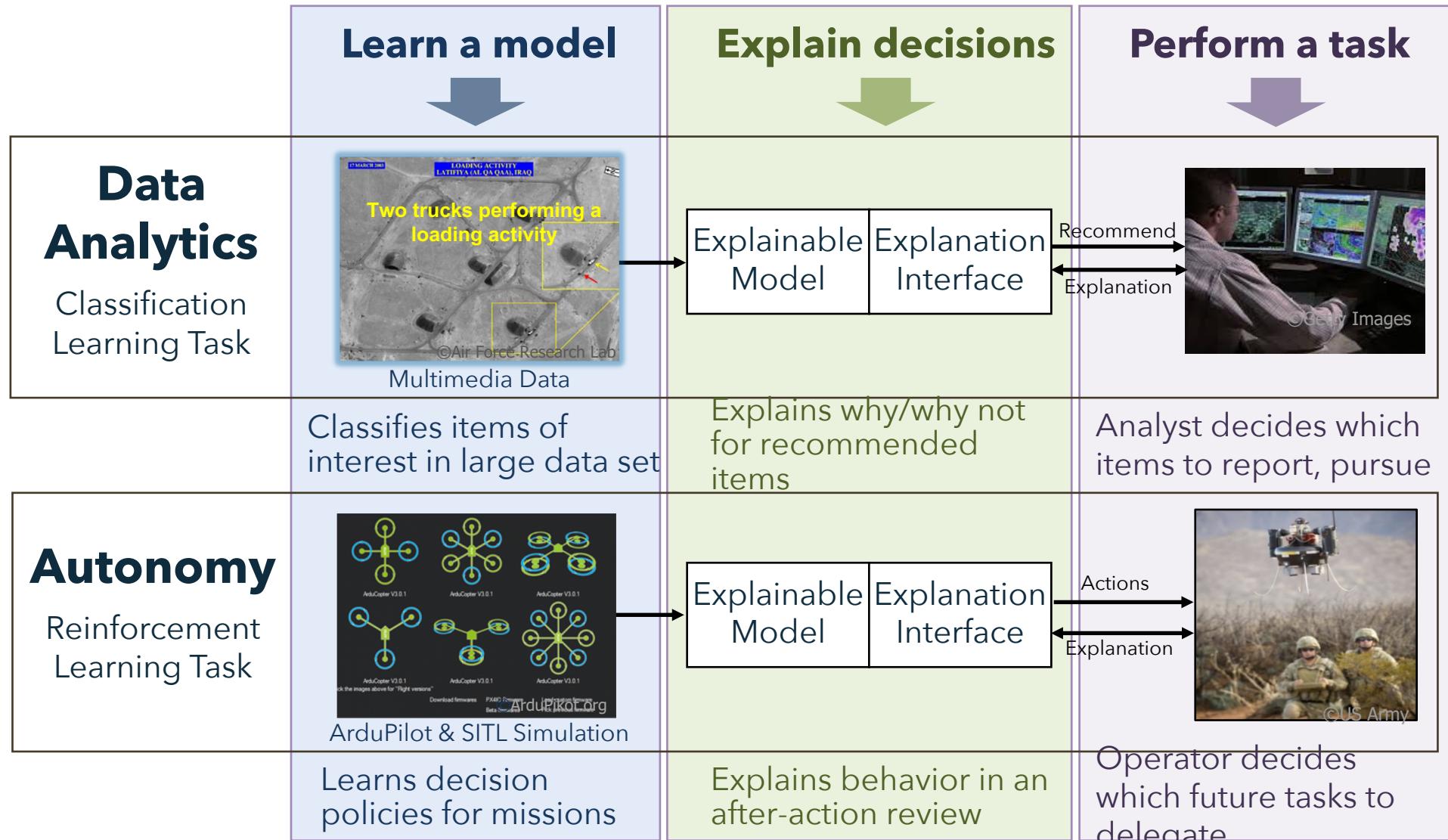
## Tomorrow



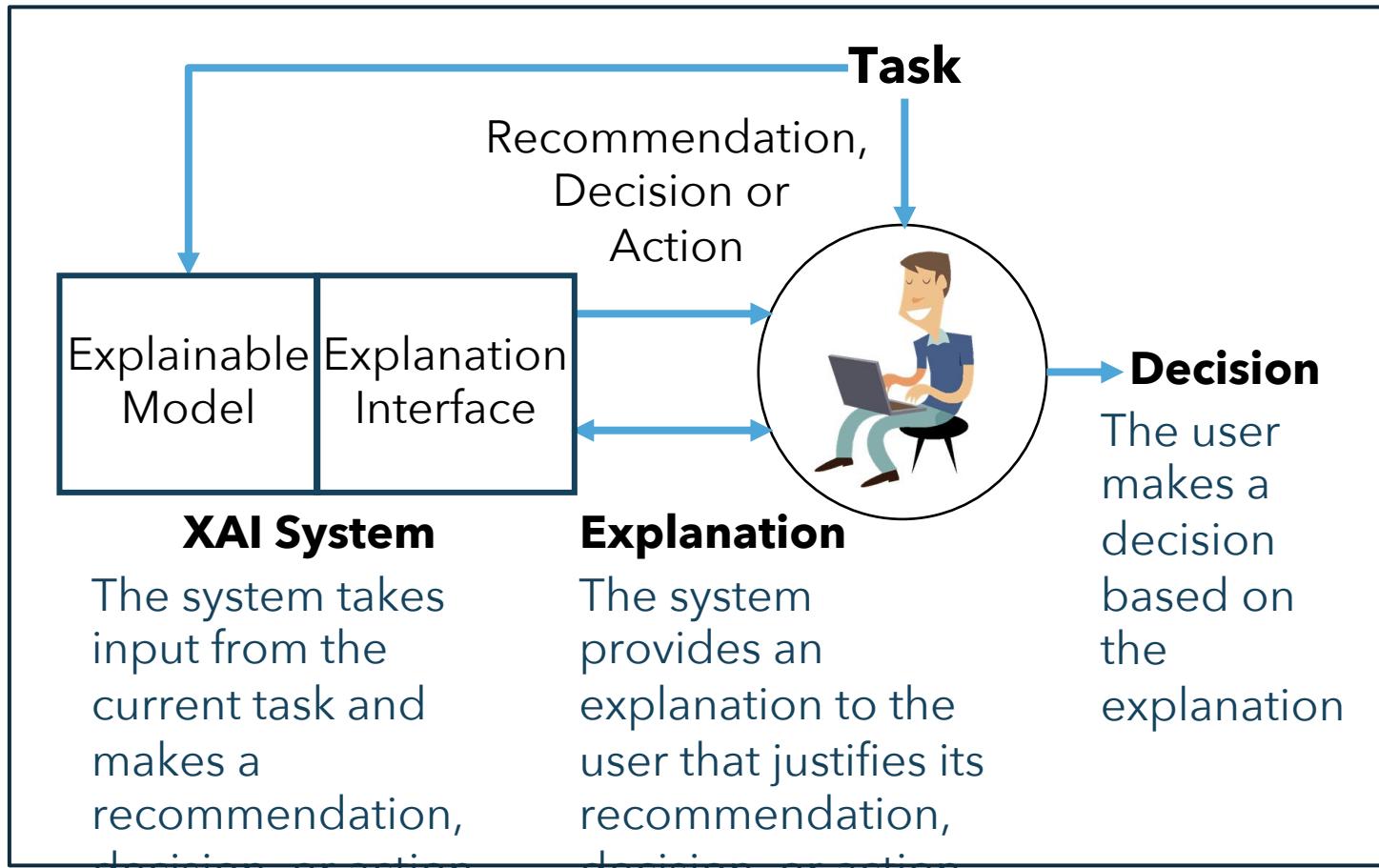
- I understand why
- I understand why not
- I know when you'll succeed
- I know when you'll fail
- I know when to trust you
- I know why you erred

- XAI will create a suite of machine learning techniques that
  - Produce more explainable models, while maintaining a high level of learning performance
  - Enable human users to understand, appropriately trust, and effectively manage the emerging generation of AI systems





## Explanation Framework



### 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

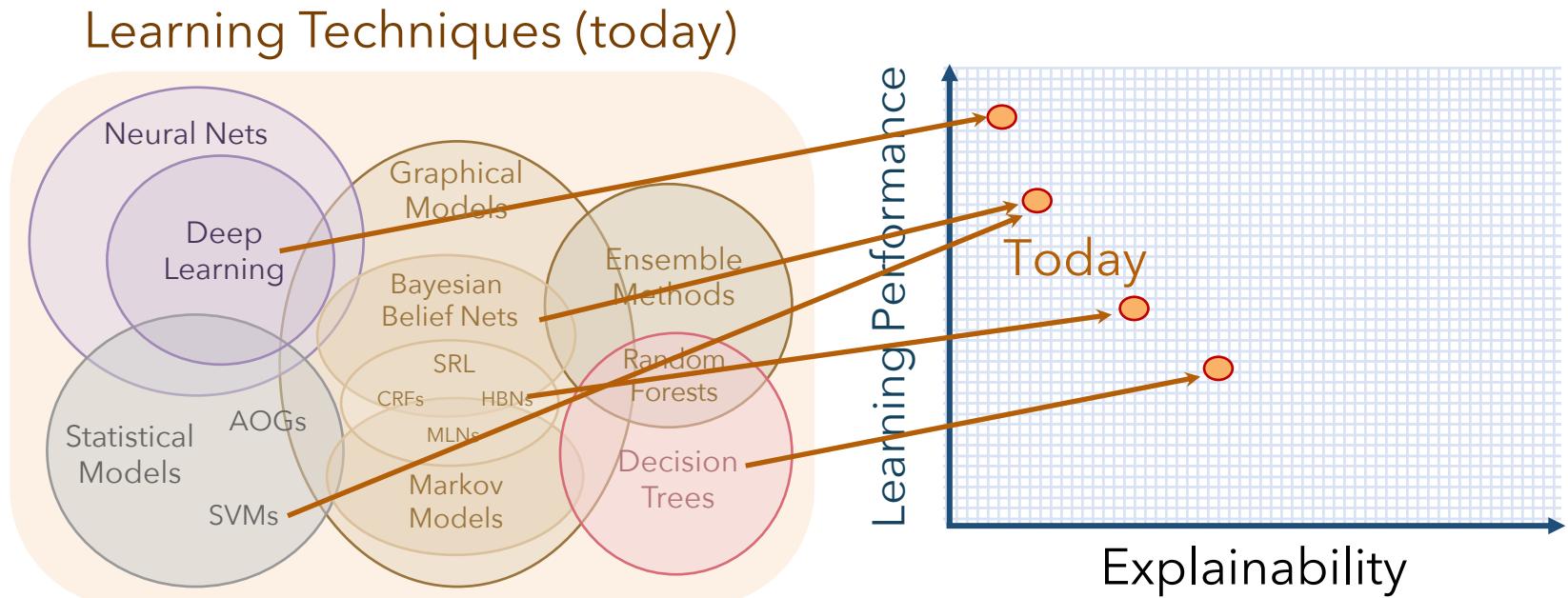
- Does the explanation improve the user's decision, task performance?

### Trust Assessment

- Appropriate future use and trust

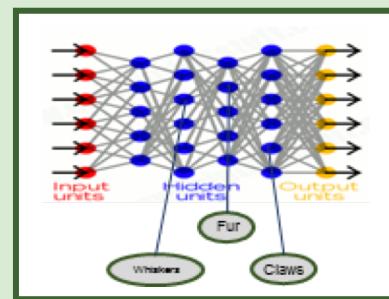
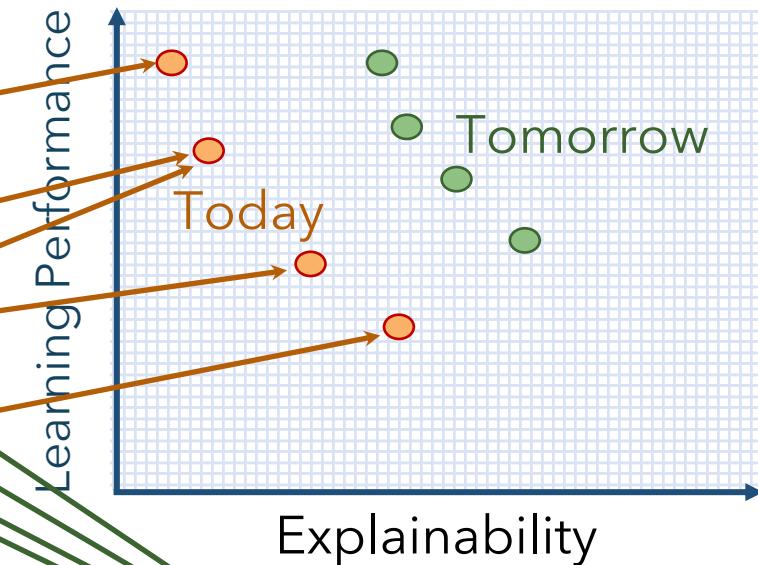
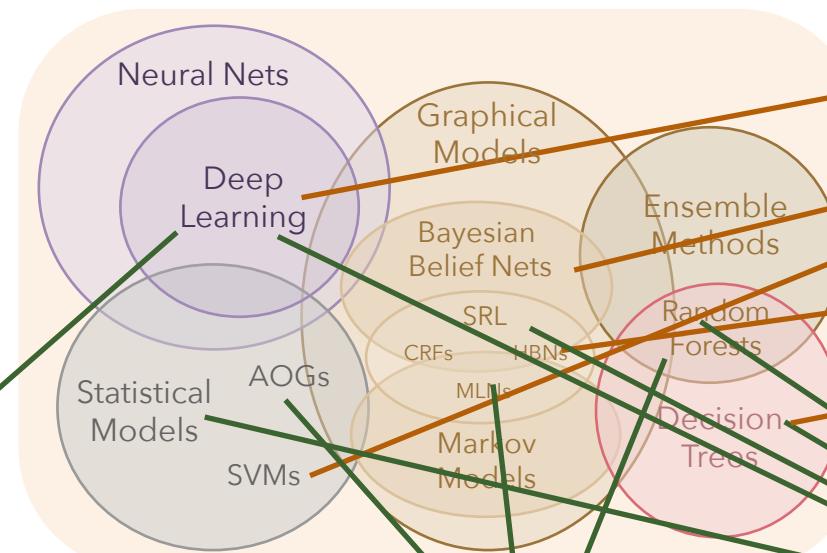
### Correctability (Extra Credit)

- Identifying errors
- Correcting errors

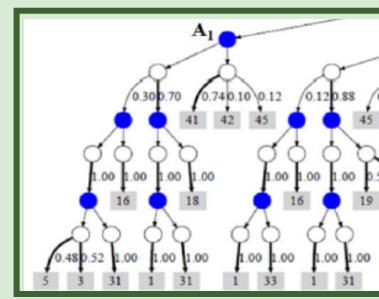


**XAI Goal**

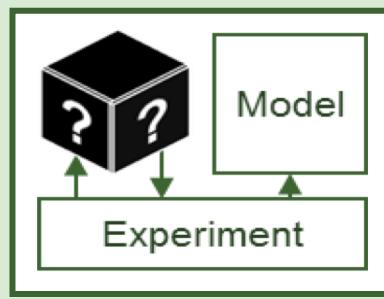
Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

**Learning Techniques (today)****Deep Explanation**

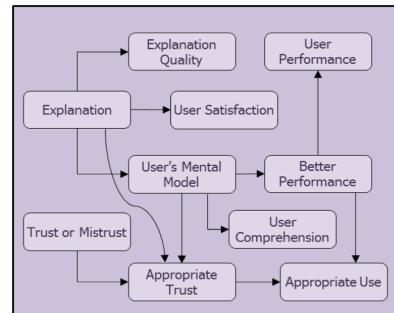
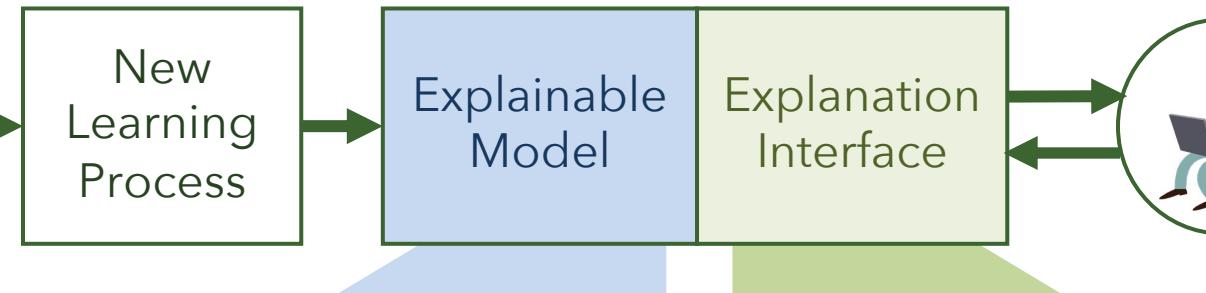
Modified deep learning techniques to learn explainable features

**Interpretable Models**

Techniques to learn more structured, interpretable, causal models

**Model Induction**

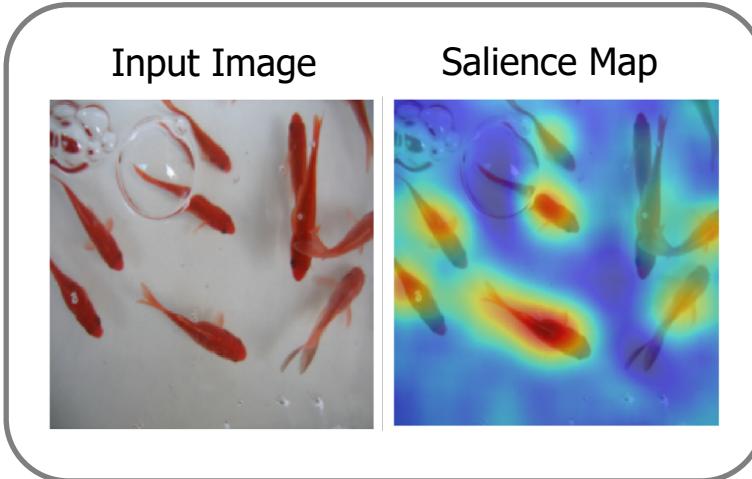
Techniques to infer an explainable model from any model as a black box



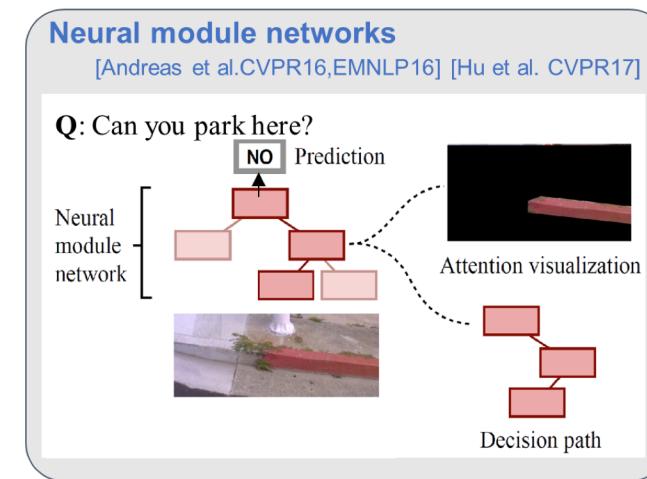
CP	Performer	Explainable Model	Explanation Interface
Both	UC Berkeley	Deep Learning	Reflexive and Rational
	Charles River	Causal Modeling	Narrative Generation
	UCLA	Pattern Theory+	3-level Explanation
Autonomy	Oregon State	Adaptive Programs	Acceptance Testing
	PARC	Cognitive Modeling	Interactive Training
	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

IHMC  
Psychological Models  
of Explanation

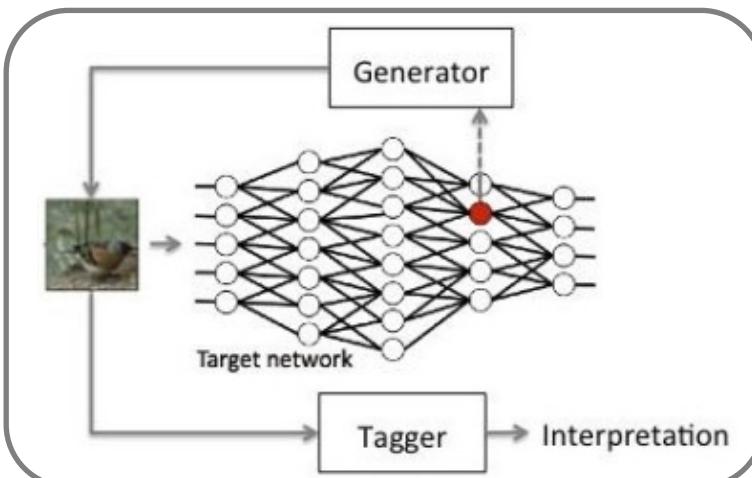
## Attention Mechanisms



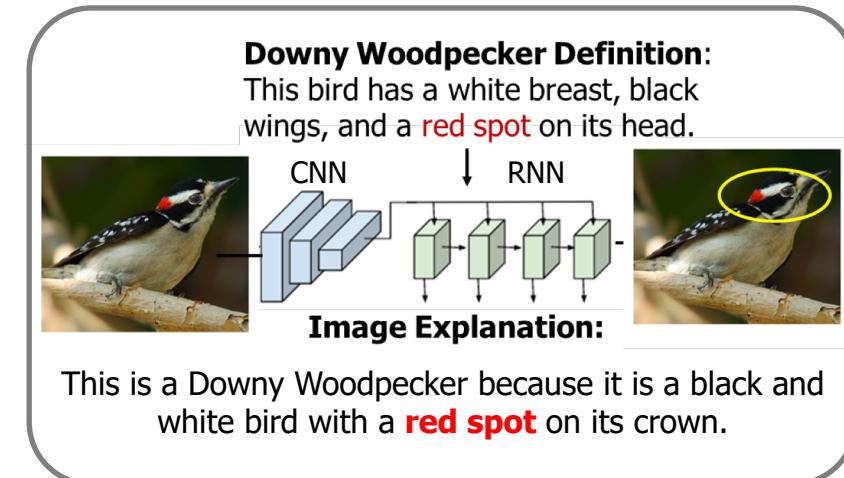
## Modular Networks



## Feature Identification



## Learn to Explain





# Deeply Explainable Artificial Intelligence



UC Berkeley, Boston U., U. Amsterdam, Kitware

## Explainable Model

### Deep Learning

- Post-hoc explanations by training additional DL models
- Explicit introspective explanations (Neural Module Networks)
- Reinforcement Learning
  - Informative rollouts
  - Explicit modular agent

## Explanation Interface

### Reflexive and Rational

- Reflexive explanations (arise from the model)
- Rational explanations (come from reasoning about user's beliefs)
- Evaluation criteria
  - Human interpretability
  - Predictive behavior
  - Appropriate trust

## Challenge Problem

### Autonomy

- Vehicle control (BDD-X, CARLA)
- Strategy games (StarCraft II)

### Data Analytics

- Visual QA and filtering tasks (VQA-X, ACT-X, xView, DiDeMo, etc.)

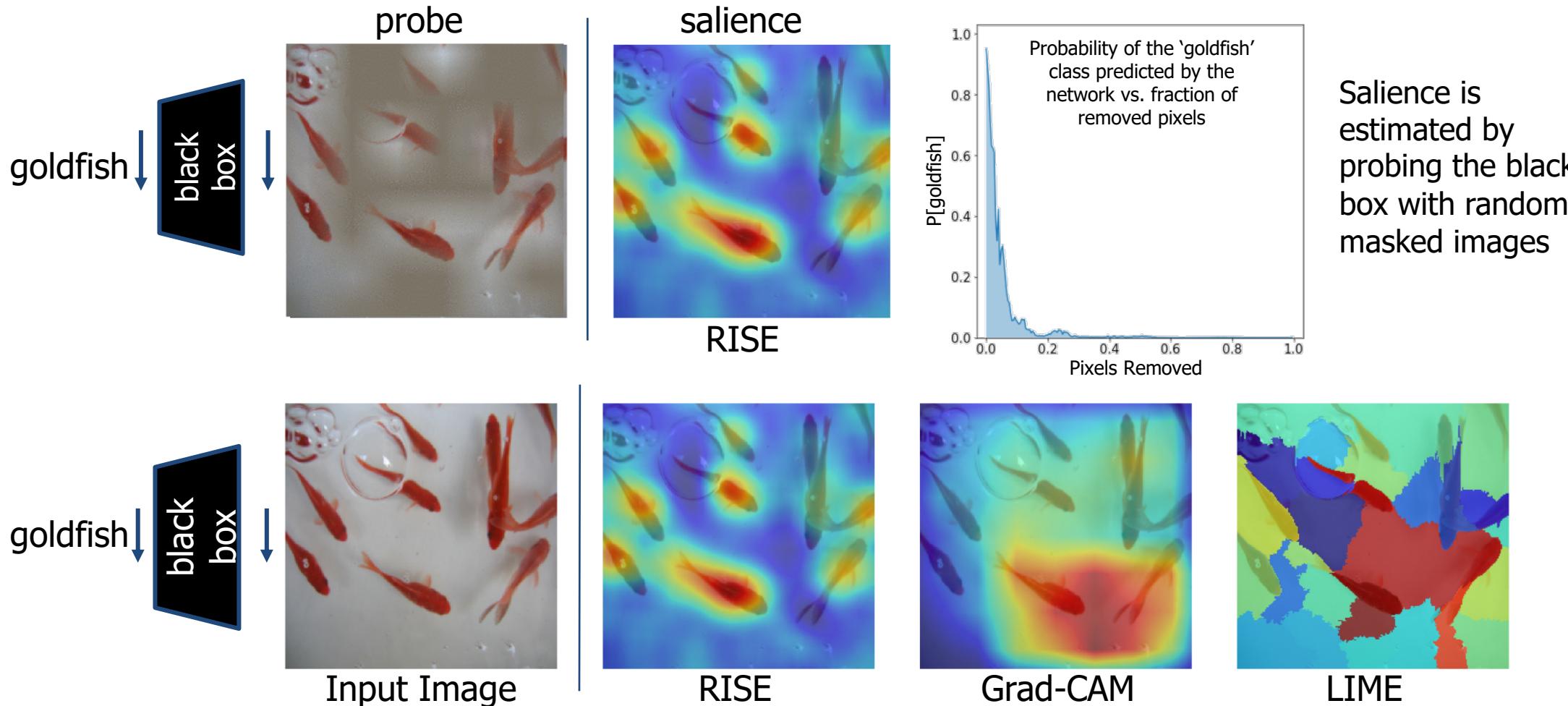
• **PI:** Trevor Darrell (UC Berkeley)

• Pieter Abbeel (UC Berkeley)  
• Tom Griffiths (UC Berkeley)  
• Kate Saenko (Boston U.)  
• Zeynep Akata (U. Amsterdam)

• Dan Klein (UC Berkeley)  
• John Canny (UC Berkeley)  
• Anca Dragan (UC Berkeley)

• Anthony Hoogs (Kitware)

## UC Berkeley, Boston U., U. Amsterdam, Kitware



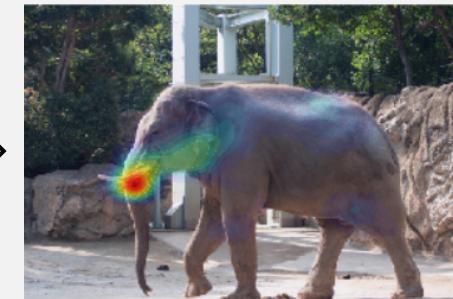
Petsiuk, Das and Saenko. RISE: Randomized Input Sampling for Explanation of Black-box Models, 2018

Given the multi-modal explanation generated by the model, do you think the system will answer correctly?

Question: *Does this elephant have tusks?*



*"because there are no bones sticking out from its mouth"*



Yes  No

Incorrect! The system answered "no" when the ground-truth answer is "yes"

Question: *Is this a professional sporting event?*



*"because the players are wearing official jerseys"*



Yes  No

Correct! The system answered "yes" when the ground-truth answer is "yes"

### Explanation Effectiveness

**Without explanation (existing SOTA)**

Attention for Explanation Used?

No

Accuracy of Users Judgement

57.5%

**UCB Model on descriptions**

Yes

66.5%

**UCB Model without attention**

No

61.5%

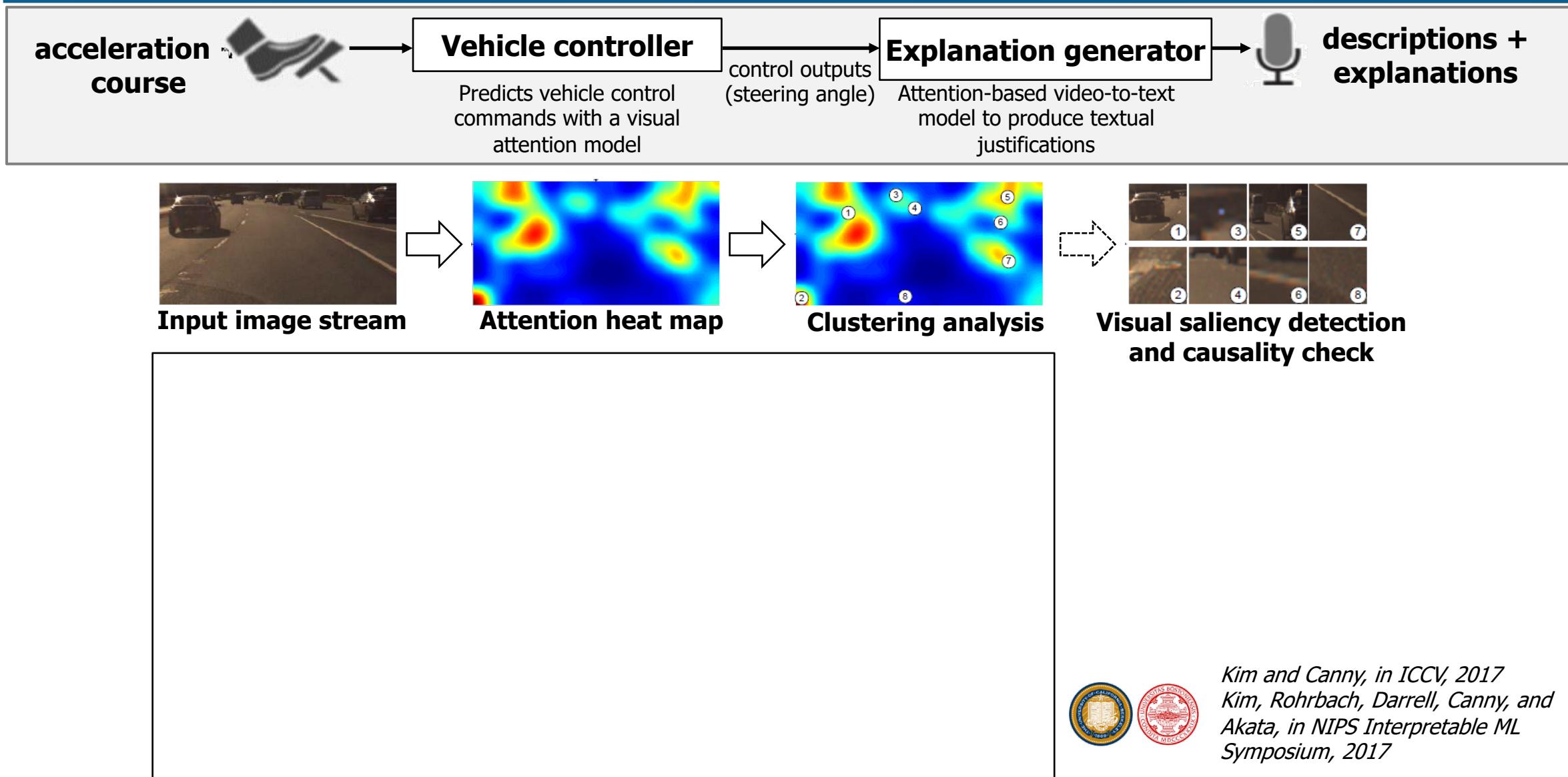
**UCB Model**

Yes

**70.0%**



Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. Multimodal Explanations: Justifying Decisions and Pointing to the Evidence, 2018



Kim and Canny, in ICCV, 2017  
Kim, Rohrbach, Darrell, Canny, and Akata, in NIPS Interpretable ML Symposium, 2017

## Charles River Analytics (CRA), U. Mass, Brown

### Explainable Model

#### Causal Modeling

- Experiment with the learned model (as a grey box) to learn an explainable, causal, probabilistic programming model

• **PI:** James Tittle (CRA)

- Jeff Druce (CRA)
- Avi Pfeffer (CRA)
- David Jensen (U. Mass)
- Michael Littman (Brown U.)

### Explanation Interface

#### Narrative Generation

- Interactive visualization based on the generation of temporal, spatial narratives from the causal, probabilistic models

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

### Challenge Problem

#### Autonomy

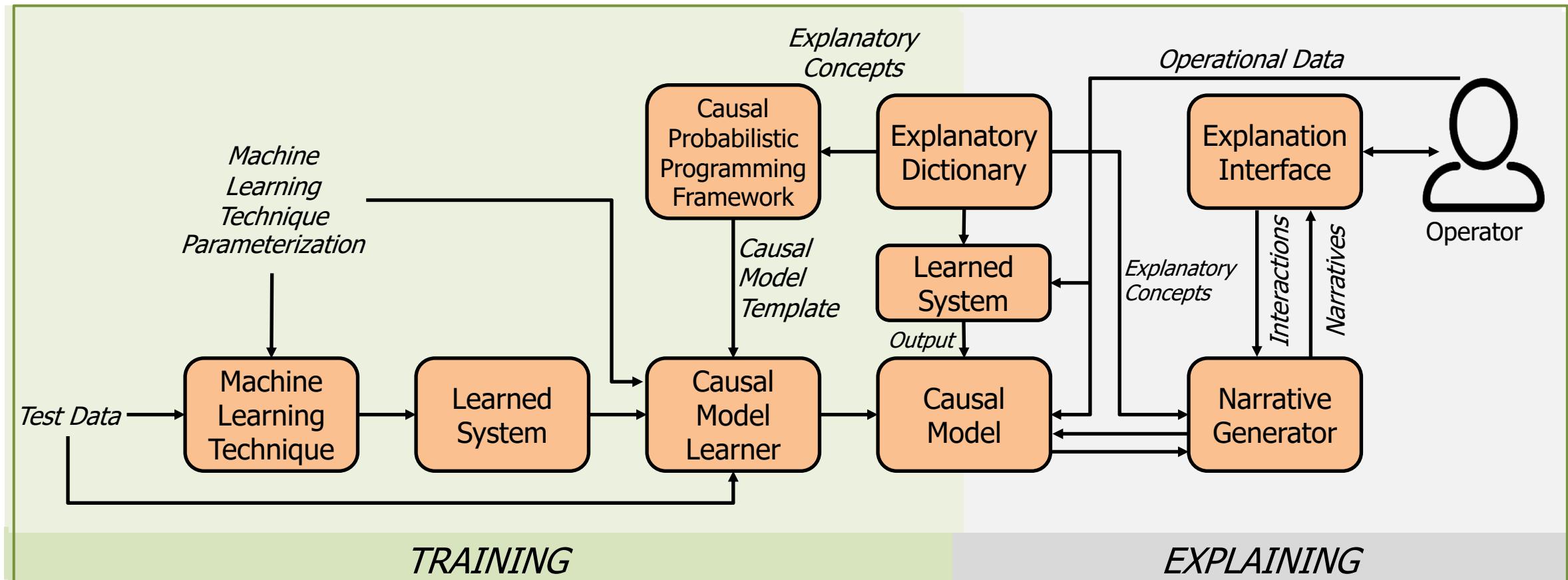
- Atari
- Starcraft

#### Data Analytics

- Pedestrian Detection (INRIA)
- Activity Recognition (ActivityNet)

## Charles River Analytics, U. Mass, Brown

Generate causal explanations of ML operation and present them to the user as intuitive narratives in an interactive, easy-to-use interface grounded in cognitive engineering theories





## UCLA, Oregon State, Michigan State

### Explainable Model

#### Pattern Theory+

- Interpretable representations
  - STC-AOG: spatial, temporal, and causal models
  - STC-PG: scene and event interpretations in analytics
  - STC-PG+: task plans in autonomy

- Theory of mind representations
  - User's beliefs
  - User's mental model of agent

### Explanation Interface

#### 3-Level Explanation

- Concept compositions
- Causal and counterfactual reasoning
- Utility explanations

- Explanation representations:
- X-AOG: explanation model
  - X-PG: explanatory parse graph as dialogue
  - X-Utility: priority and loss for explanations

### Challenge Problem

#### Autonomy

- Robot executing daily tasks in physics-realistic VR platform
- Autonomous vehicle driving (GTA5 game engine)

#### Data Analytics

- Network of video cameras for scene understanding and event analysis

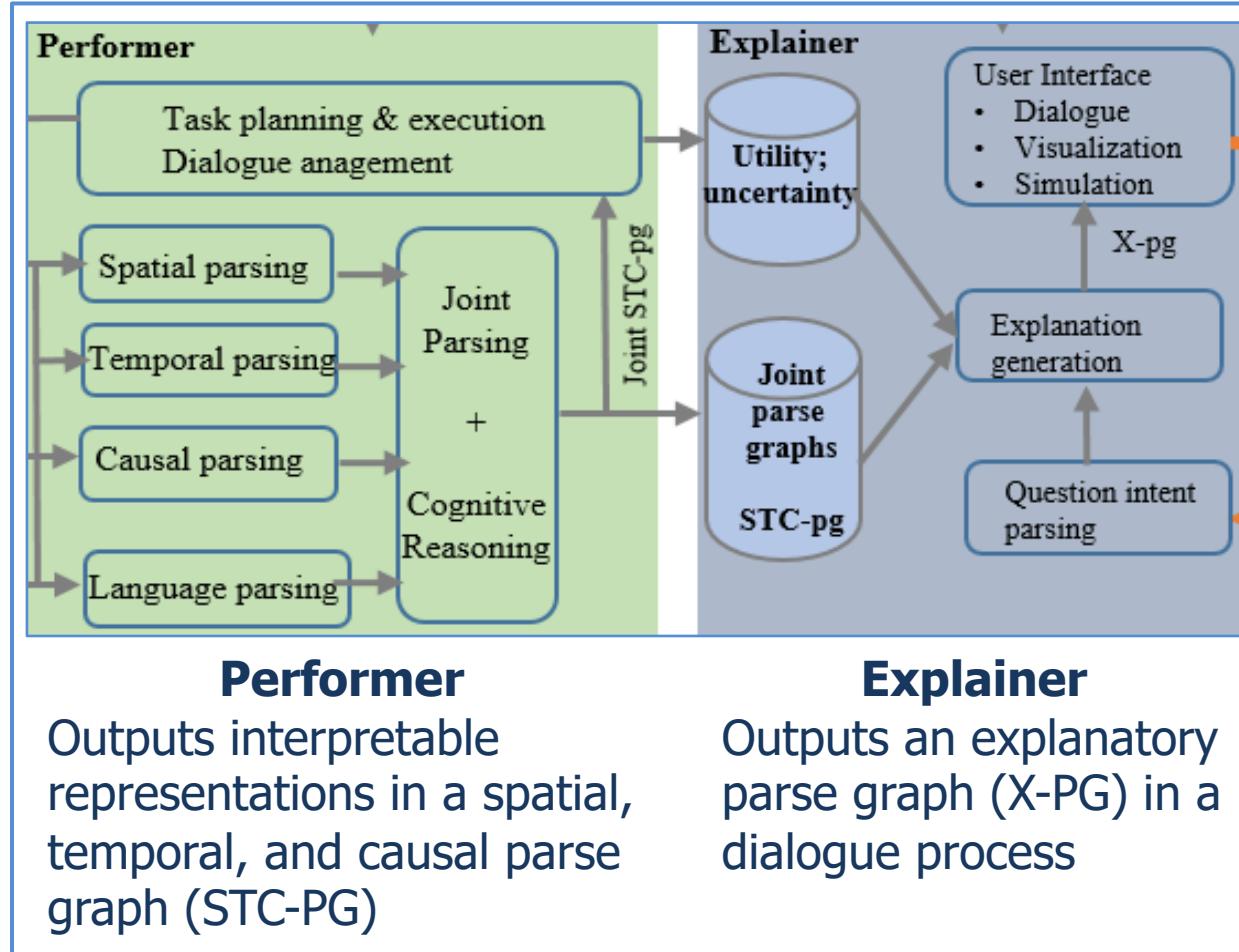
• PI: Song-Chun Zhu (UCLA)

• Ying Nian Wu (UCLA)  
• Sinisa Todorovic (OSU)

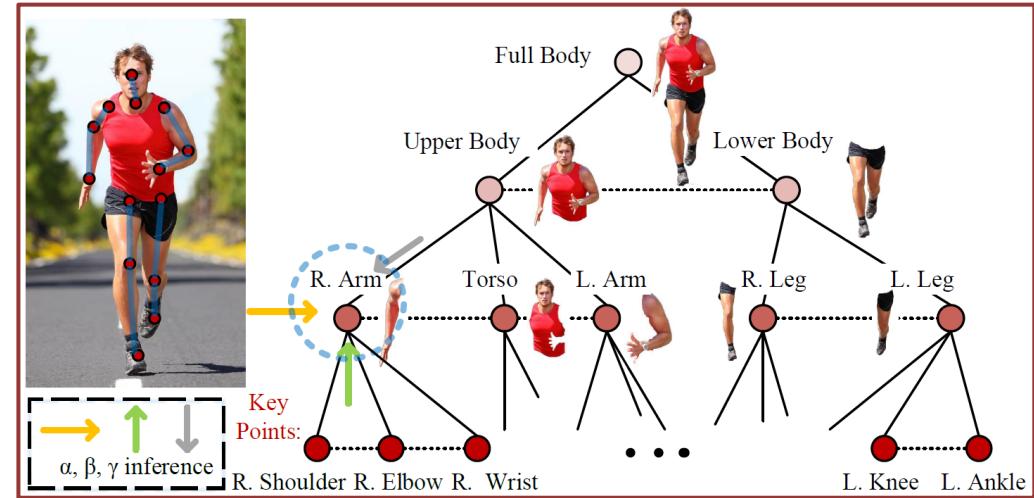
• Joyce Chai (Michigan State)

UCLA, Oregon State, Michigan State

## System Architecture



## STC Parse Graph



An attributed parse graph for a running person. Each node has 3 computing channels:

- $\alpha$ : grounding the node on DNN features;
- $\beta$ : bottom-up;
- $\gamma$ : top-down.

An explanation is represented as parse graph X-pg

## Oregon State University

### Explainable Model

#### Adaptive Programs

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

• PI: Alan Fern (OSU)

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

### Explanation Interface

#### Acceptance Testing

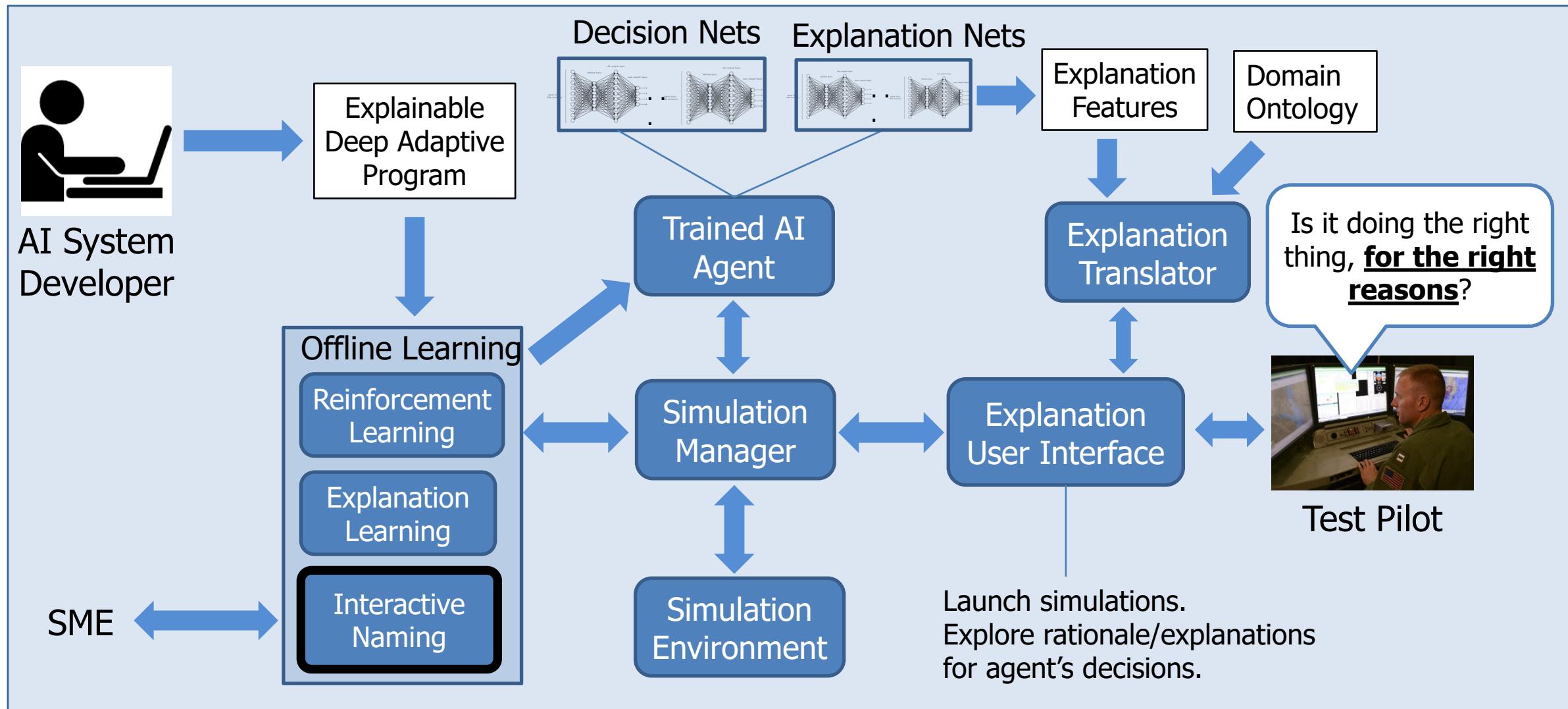
- Provides a visual and Natural Language explanation interface for acceptance testing by test pilots based on Information Foraging Theory

### Challenge Problem

#### Autonomy

- Real-time Strategy Games based on custom designed game engine designed to support explanation
- Starcraft II

## Oregon State University



## PARC, CMU, U. Edinburgh, U. Michigan, USMA, IHMC

### Explainable Model

#### Cognitive Model

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

### Explanation Interface

#### Interactive Training

- Interactive visualization of states, actions, policies, and values
- Module for test pilots to refine and train the system

### Challenge Problem

#### Autonomy

- MAVSim wrapper over ArduPilot simulation environment
- Value of Explanation framework for measuring explanation effectiveness

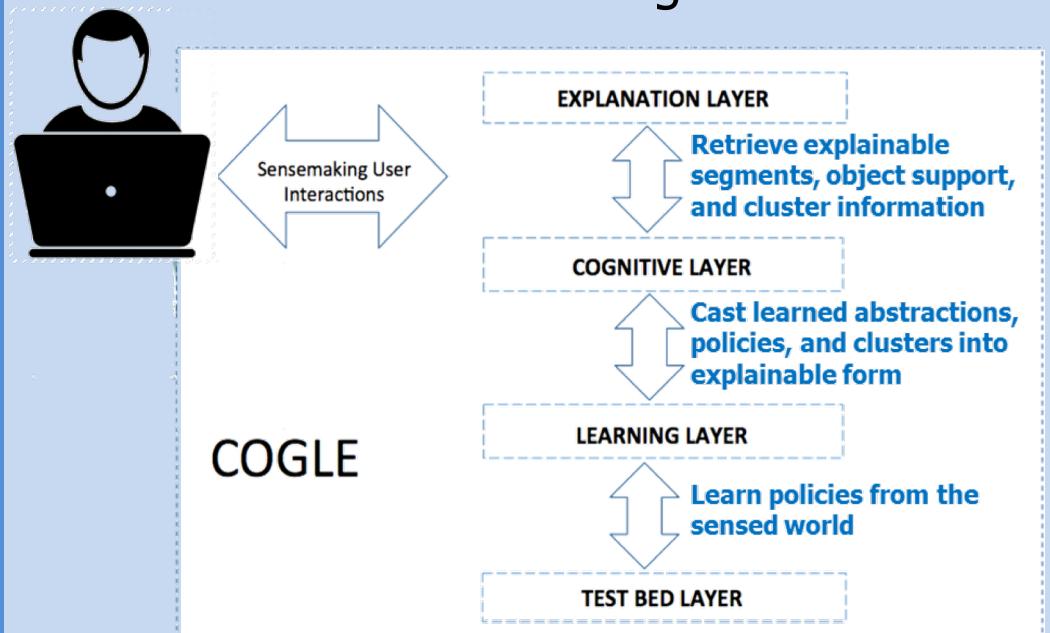
• PI: Mark Stefik (PARC)

• Honglak Lee (U. Michigan)  
• Subramanian Ramamoorthy (U. Edinburgh)

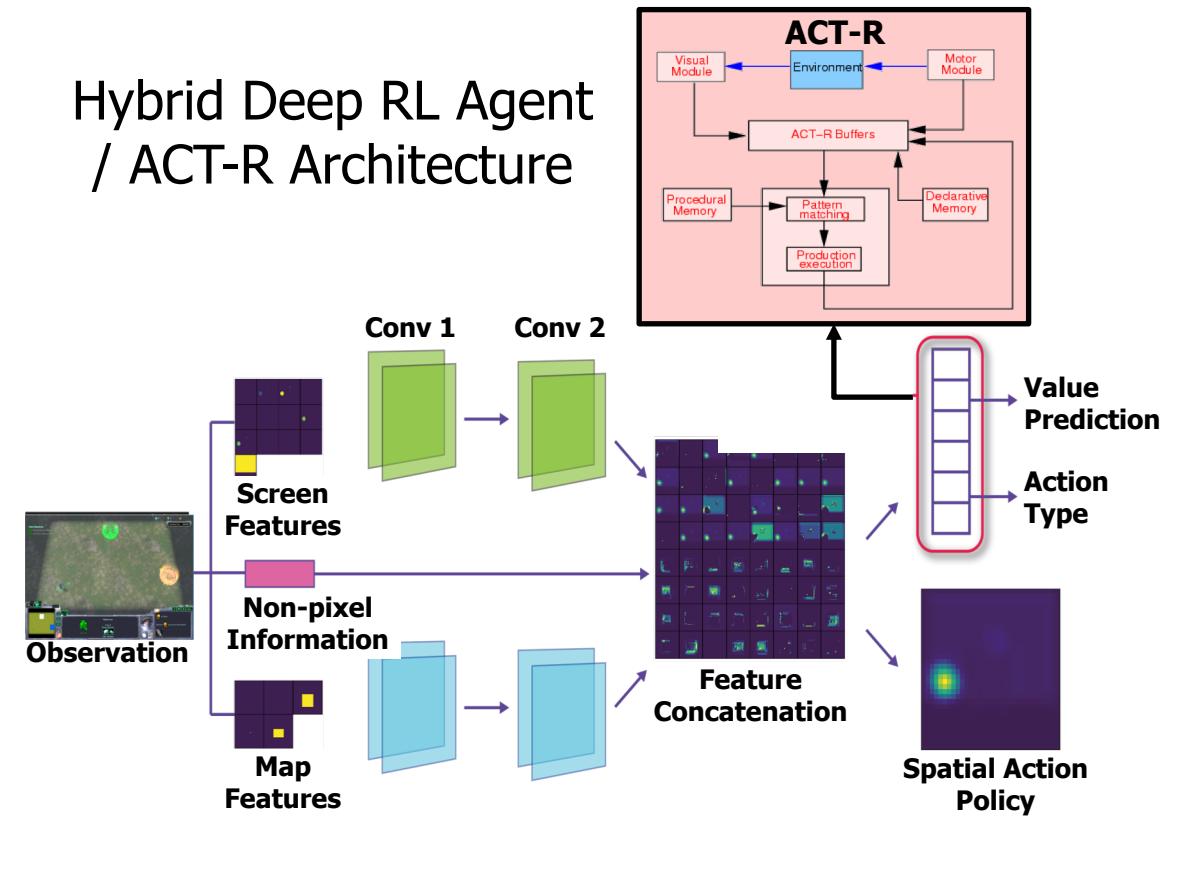
• Christian Lebiere (CMU)  
• John Anderson (CMU)  
• Robert Thomson (USMA)

• Michael Youngblood (PARC)

## PARC, CMU, U. Edinburgh, U. Michigan, USMA, IHMC

Layered Cognitive Architecture  
to Partition Explanation And  
Learning

## Hybrid Deep RL Agent / ACT-R Architecture



## Carnegie Mellon University

### Explainable Model

#### Explainable RL (XRL)

- Create a new scientific discipline for Explainable Reinforcement Learning with work on new algorithms and representations

• **PI:** Zico Kolter (CMU)

- Geoff Gordon (CMU)
- Pradeep Ravikumar (CMU)

### Explanation Interface

#### XRL Interaction

- Interactive explanations of dynamic systems
- Human-machine interaction to improve performance

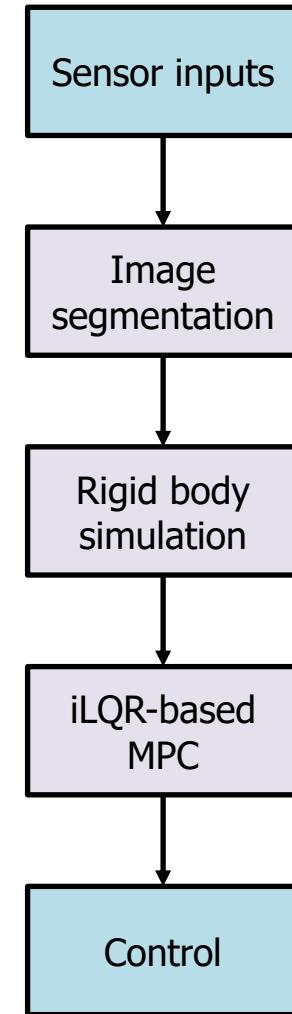
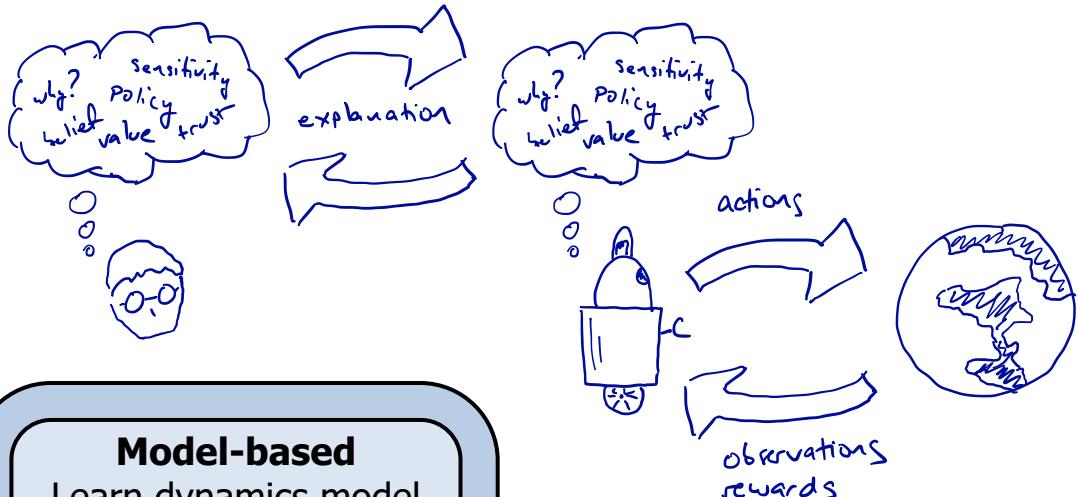
### Challenge Problem

#### Autonomy

- Open AI Gym
- Autonomy in the electrical grid
- Mobile service robots
- Self-improving educational software

## Carnegie Mellon University

Create a new discipline of explainable RL to enable dynamic human-machine interaction and adaptation for maximum team performance



**Differentiable Physics** - Applies implicit differentiation to solutions of LCP to analytically derive a backpropagation update of next state with respect to previous state, control, and model parameters

## SRI International, U. Toronto, UCSD, U. Guelph

### Explainable Model

#### Deep Learning

- Multiple deep learning techniques
- Attention-based mechanisms
  - Compositional NMNs
  - GANs

### Explanation Interface

#### Show & Tell Explanation

- DNN visualization
- Query evidence that explains DNN decisions
- Generate natural language justifications

### Challenge Problem

#### Data Analytics

- VQA
  - Visual Gnome
  - Flickr30
- MovieQA

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

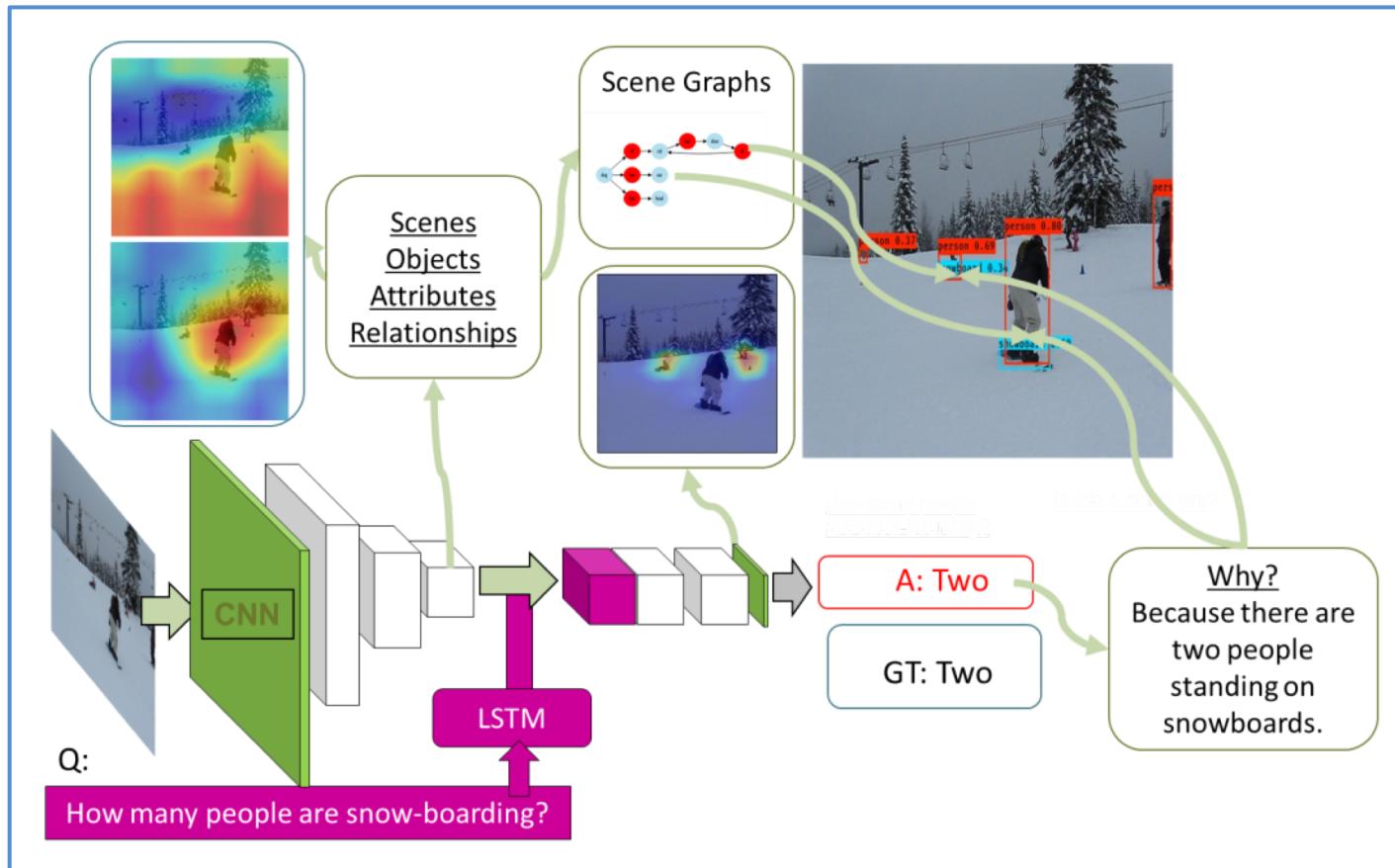
- Xiao Lin (SRI)
- Ryan Villamil (SRI)
- Dejan Jovanovic (SRI)
- Avi Ziskind (SRI)
- Michael Wessel (SRI)

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

- Jürgen Schulze (UCSD)

## SRI International, U. Toronto, UCSD, U. Guelph

Interpretable, Scene Graph-based VQA System with Active Attention



- Generate “show-and-tell” explanations with justifications of decisions accompanied by visualizations of input data used to generate inferences
- Scene and Situation Graphs, inferred from images and videos, support rich multimodal data analytics and explanations
- Scene Graphs guide attentional scanning for interpretable analytics

## Raytheon BBN, Georgia Tech, UT Austin, MIT

### Explainable Model

#### Deep Learning

- Semantic labelling of DNN neurons
- DNN audit trail construction
- Gradient-weighted Class Activation Mapping

### Explanation Interface

#### Argumentation Theory

- Comprehensive strategy based on argumentation theory
- NL generation
- DNN visualization

### Challenge Problem

#### Data Analytics

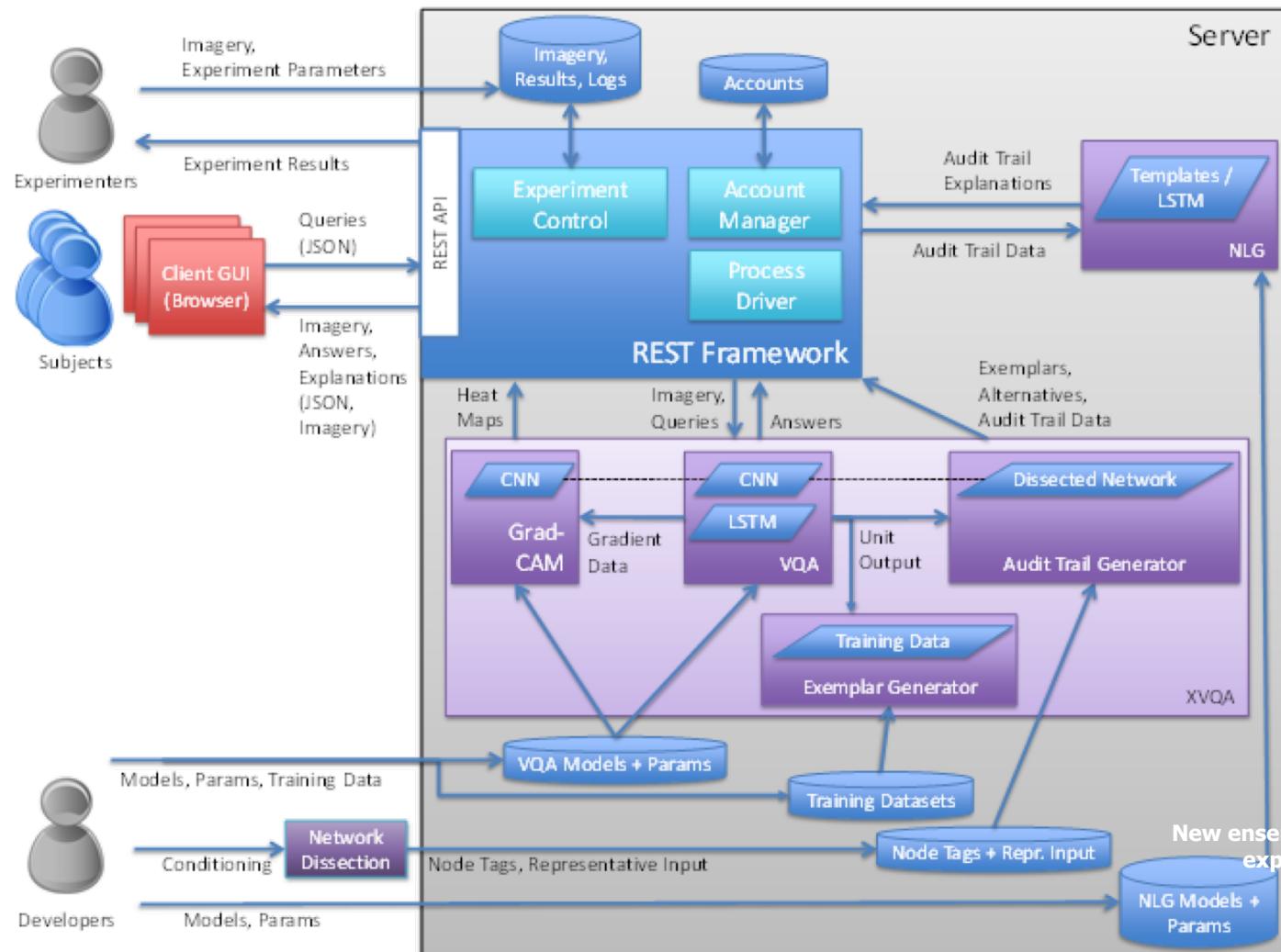
- VQA for images and video

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

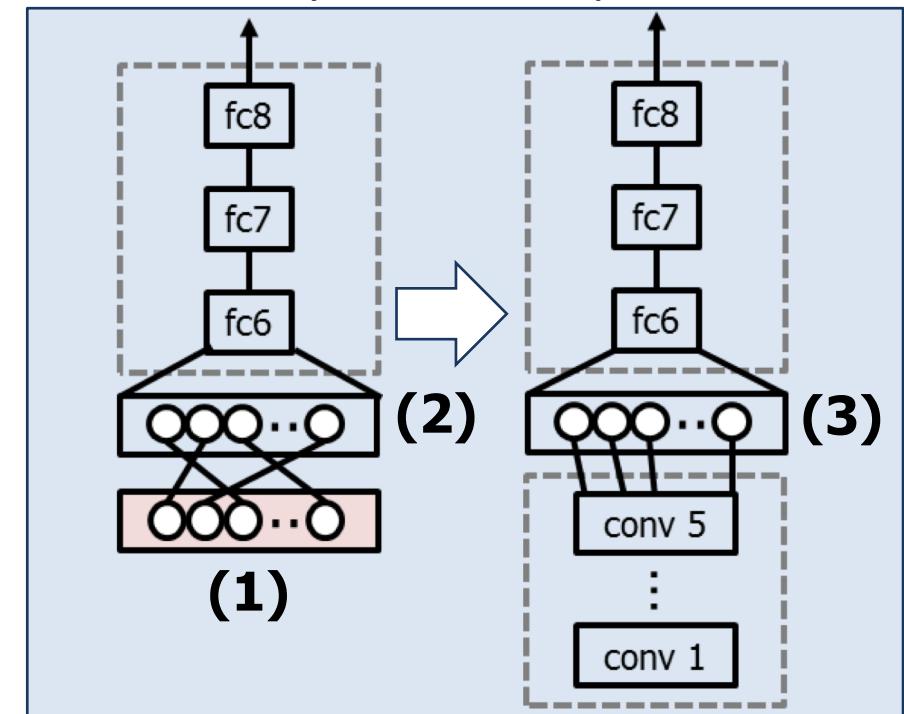
• Antonio Torralba (MIT)  
• Ray Mooney (UT Austin)

• Devi Parikh (Georgia Tech)  
• Dhruv Batra (Georgia Tech)

## Raytheon BBN, Georgia Tech, UT Austin, MIT



Improve the interpretability of units using a **new conditioning method** to retrain the network to intentionally include *concept detectors*



- 1) Pick units from standard vocabulary**
- 2) Train top part of net**
- 3) Use top to train bottom**

## UT Dallas, UCLA, Texas A&amp;M, Indian Institute of Technology

## Explainable Model

**Probabilistic Logic**

- Tractable Probabilistic Logic Models (TPLMs)
  - an important class of (non-deep learning) interpretable models

## Explanation Interface

**Decision Diagrams**

- Enables users to explore and correct the underlying model as well as add background knowledge

## Challenge Problem

**Data Analytics**

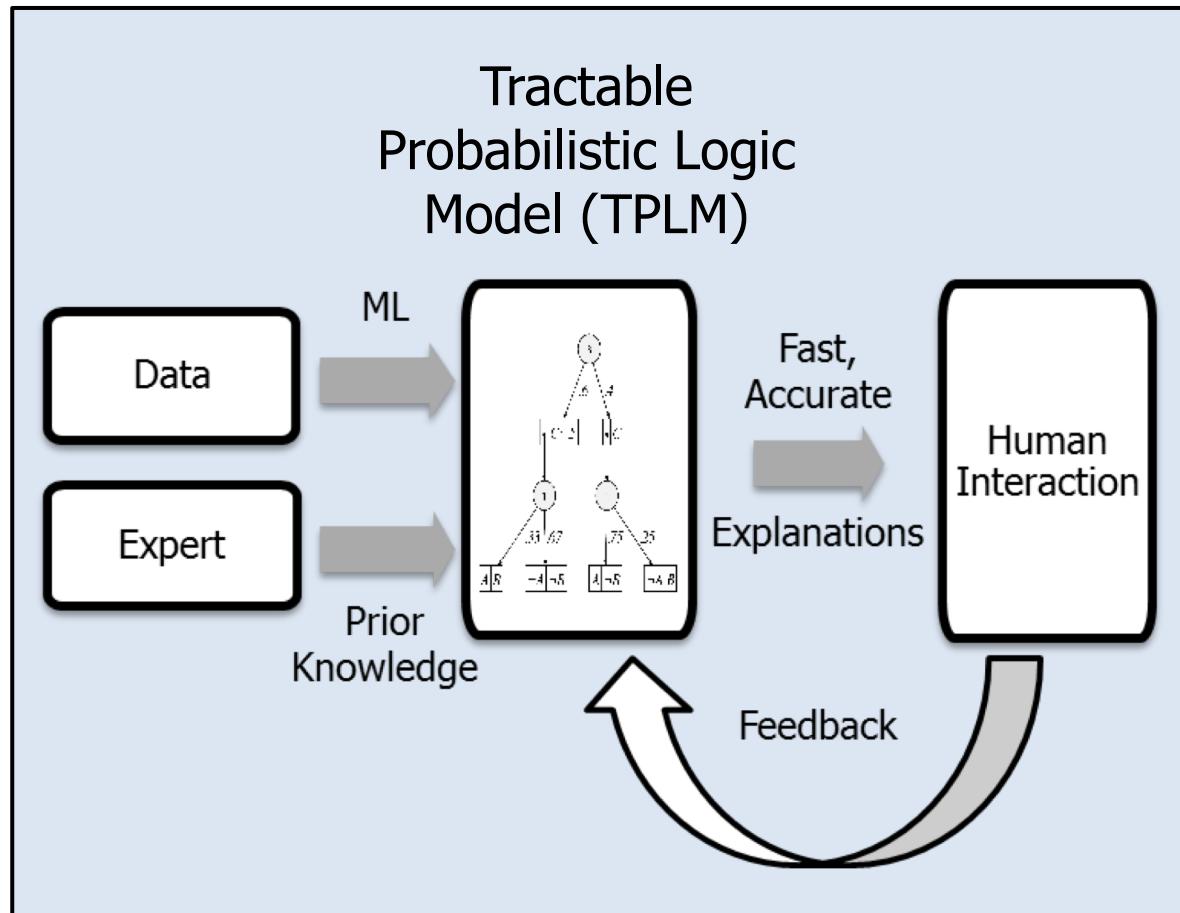
- Infer activities in multimodal data (video and text)
- Wetlab (biology) and TACoS (cooking) datasets

- **PI:** Vibhav Gogate (UT Dallas)

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

**UT Dallas, UCLA, Texas A&M, Indian Institute of Technology**

Use interpretable and tractable models based on well-founded principles from logic and probability theory



**Find all videos in which a person peels oranges**  
(explanations are captions generated by TPLMs)



Person using his hands to peel oranges. I can see the orange skin



Person using his hands. I can see the orange skin and skinless orange



Person using his hands to peel oranges. I can see the orange skin on the table and peeled oranges

## Texas A&M, Washington State

### Explainable Model

#### Mimic Learning

- Mimic learning framework combines DL models for prediction and shallow models for explanations
- Interpretable learning algorithms extract knowledge from DNNs for relevant explanations

### Explanation Interface

#### Interactive Visualization

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

### Challenge Problem

#### Data Analytics

- Multiple tasks using data from Twitter, Facebook, ImageNet, and news websites

• PI: Xia Hu (Texas A&M)

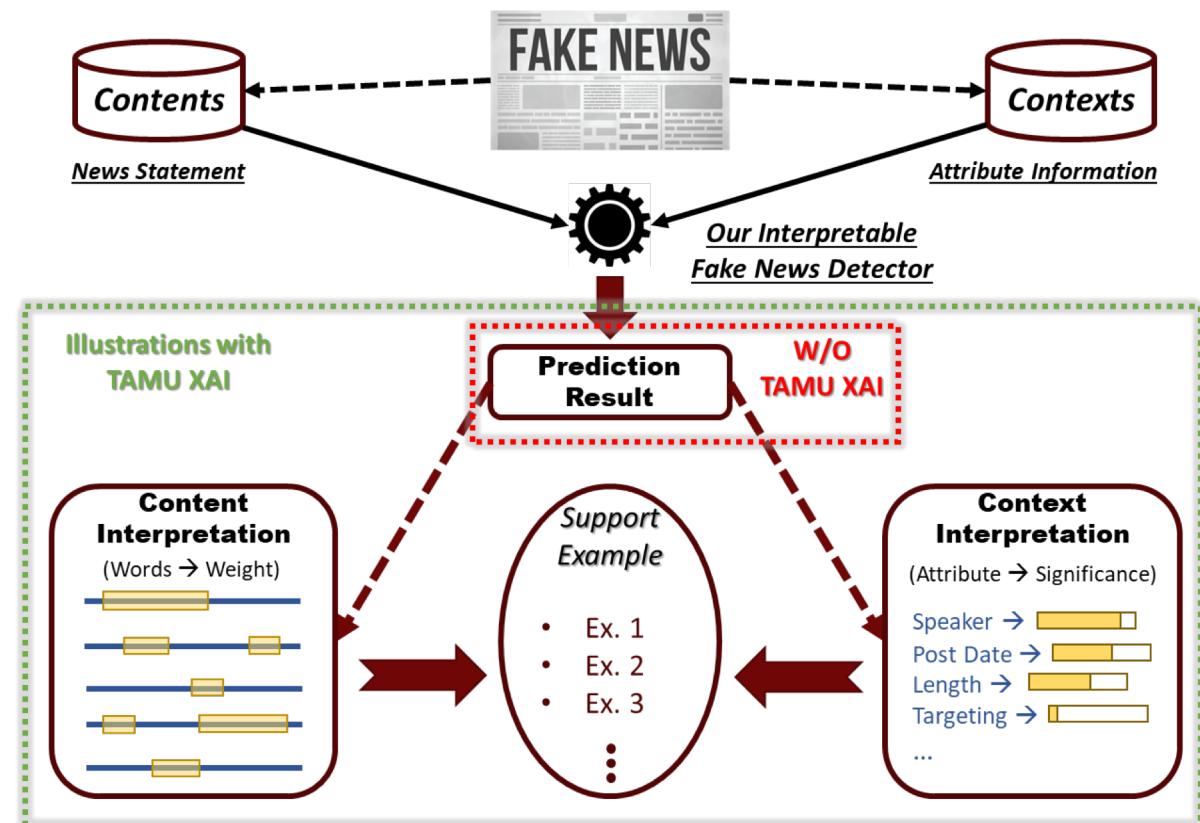
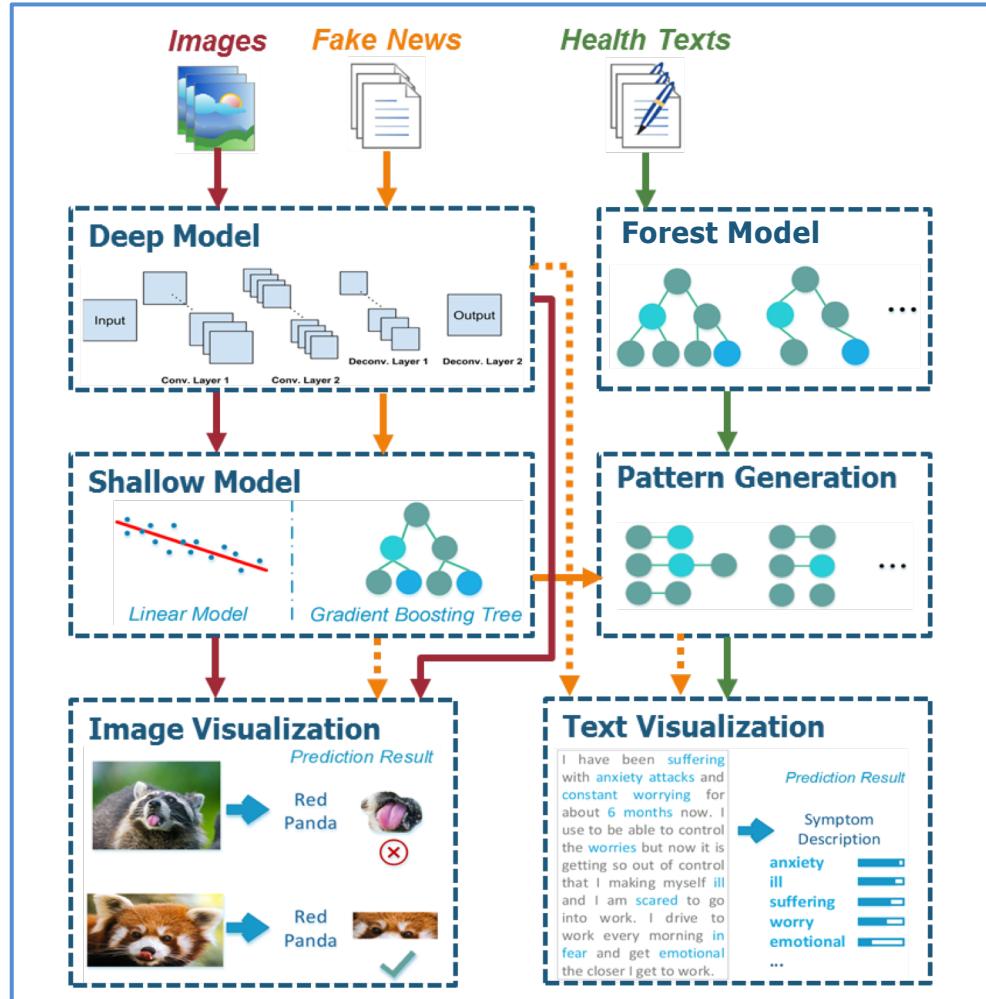
• Shuiwang Ji (Wash. State)

• Eric Ragan (Texas A&M)

# Transforming Deep Learning to Harness the Interpretability of Shallow Models

## Texas A&M, Washington State

Develop an end-to-end interpretable deep learning infrastructure with image and text datasets



## Rutgers University

### Explainable Model

#### Model Induction

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

• **PI:** Patrick Shafto (Rutgers)

• Scott Cheng-Hsin Yang (Rutgers)

### Explanation Interface

#### Bayesian Teaching

- Example-based explanation of
  - Full model
  - User-selected sub-structure
  - User submitted examples

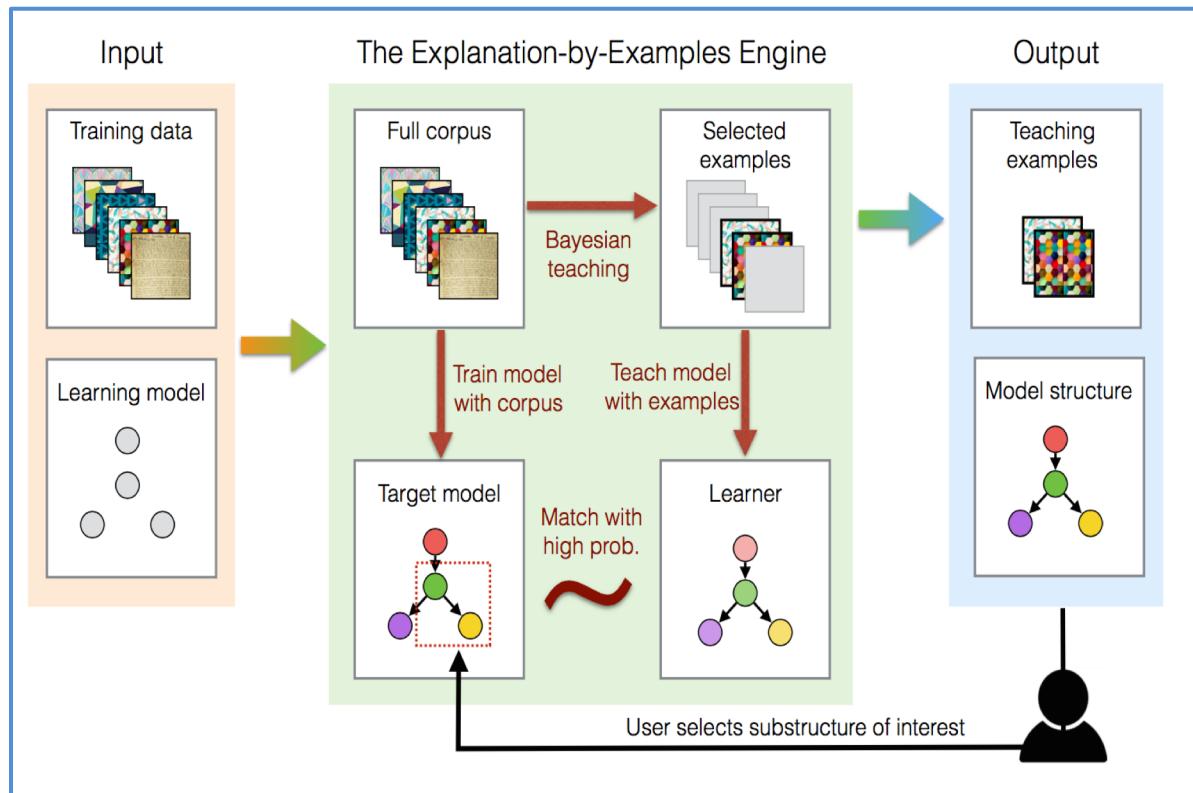
### Challenge Problem

#### Data Analytics

- Image processing
- Text corpora
- VQA
- Movie events

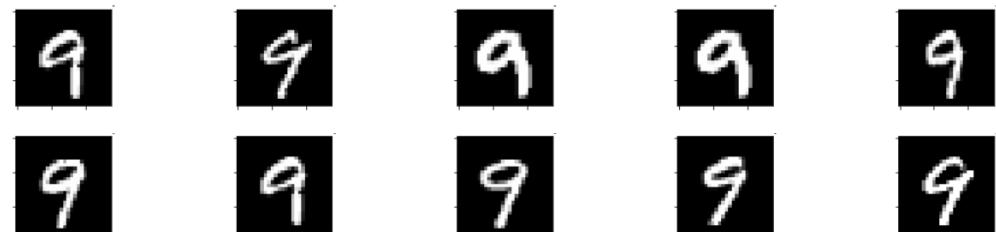
## Rutgers University

Extend Bayesian teaching to enable automatic explanation by selecting the subset of data that are most representative of the model's generative process

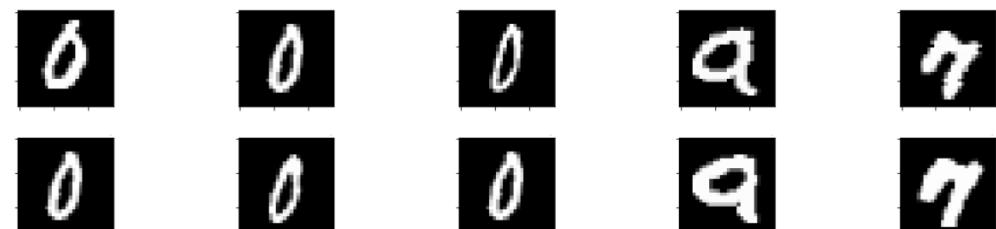


Good and bad examples for teaching a category  
(illustrates model strengths and weaknesses)

Good pairs of examples of the category 9

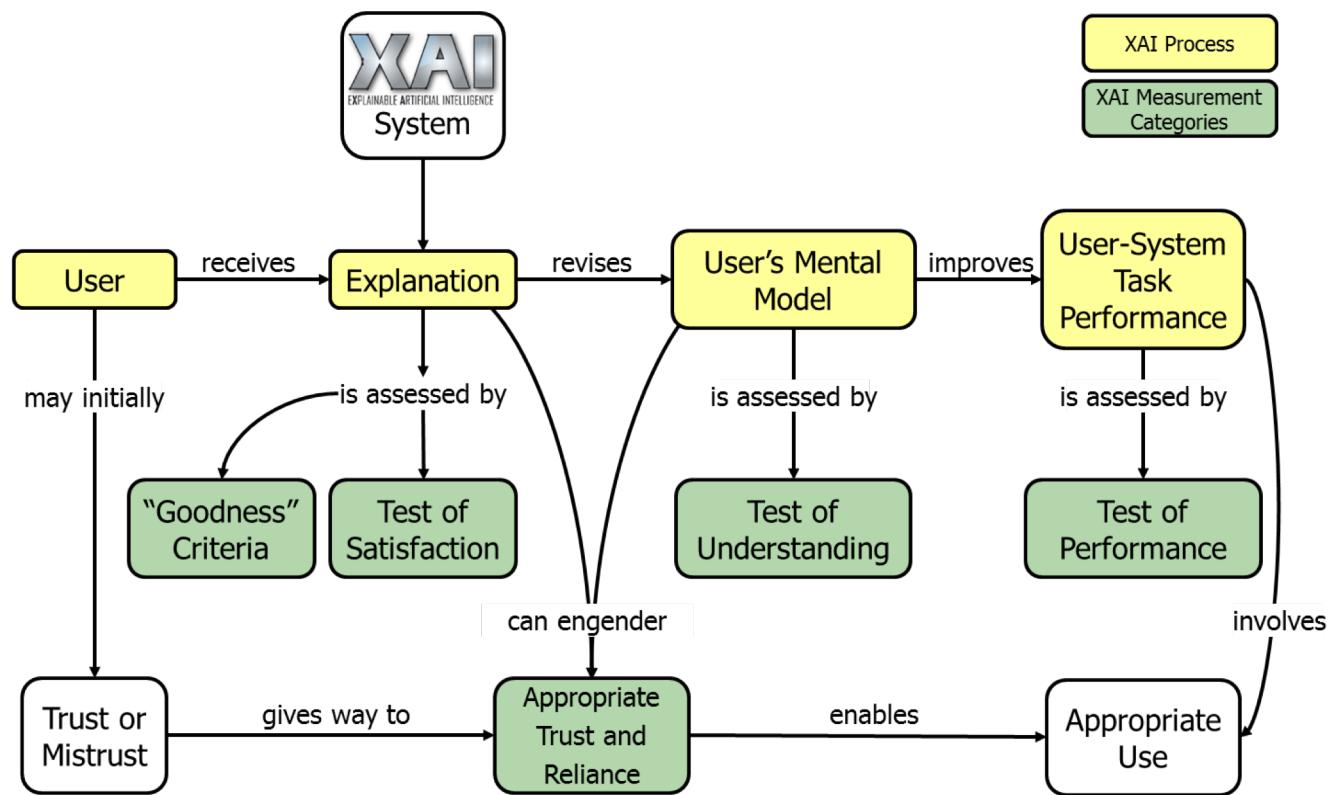


Bad pairs of examples of the category 9



CY2017		CY2018												CY2019												CY2020												CY2021											
FY2017						FY2018						FY2019						FY2020						FY2021																									
	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	JAN	FEB	MAR	APR	MAY											
<b>PHASE 1</b> <b>Technology Demonstrations</b>																		<b>PHASE 2</b> <b>Government Evaluations</b>																															
<b>Evaluator</b>	Define Evaluation Framework												Prep for Eval 1	<b>Eval 1</b>	Analyze Results	Prep for Eval 2				<b>Eval 2</b>	Analyze Results	Prep for Eval 3				<b>Eval 3</b>	Analyze Results and Accept Toolkits																						
	Develop and Demonstrate Explainable Learning Systems												<b>Eval 1</b>	Refine and Test Explainable Learning Systems				<b>Eval 2</b>	Refine and Test Explainable Learning Systems				<b>Eval 3</b>	Deliver Software Libraries and																									
<b>Explainable Learning Systems</b>	Summarize Current Psychological Theories of Explanation						Develop Computational Model of Explanation												Refine and Test Computational Model of Explanation												Deliver Computational Model Software																		
	Meetings												Kickoff												Progress Report												Tech Demos												
Phase 1 Evaluations																		Eval 1 Results																		Eval 2 Results													
Final												Final												Final												Final													

## Explanation Process & Measures



## Experimental Conditions

**Without Explanation** - The explainable learning system is used to perform a task without providing an explanation to the user

**With Explanation** - The explainable learning system is used to perform a task and generates explanations for every recommendation or decision it makes, and every action it takes

**Partial Explanation** - The explainable learning system is used to perform a task and generates only partial or ablated explanations (to assess various explanation features)

**Control** - A baseline state-of-the-art non-explainable system is used to perform a task



[www.darpa.mil](http://www.darpa.mil)