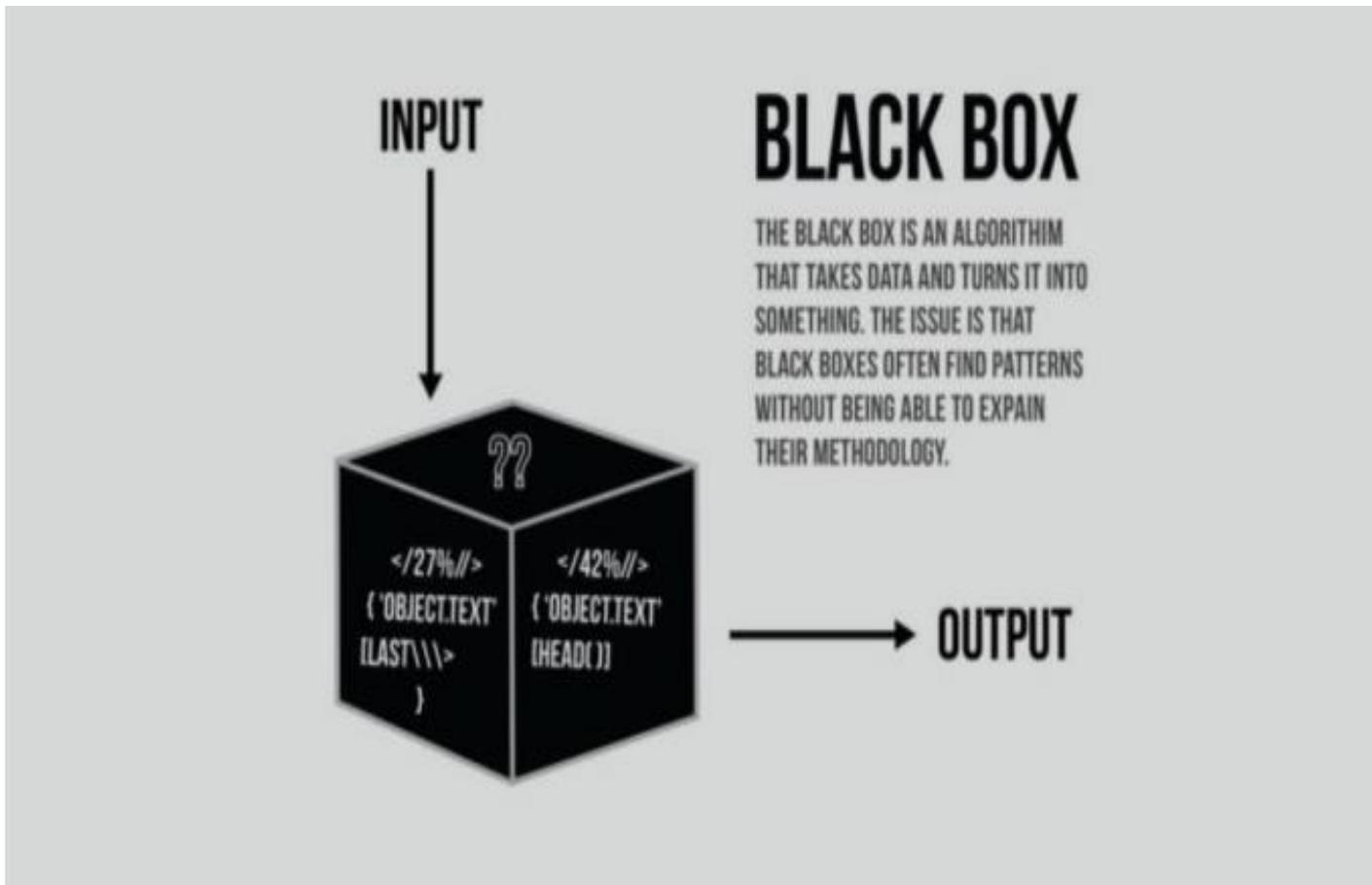


# Outline

- Interpretability = Transparency ≠ Explainability?
- DARPA XAI Grant program
- Results of the DARPA XAI Grant program
- Key takeaways and discussions



“ If you can't explain it simply, you don't understand it well enough.

— Albert Einstein

**Table 2**

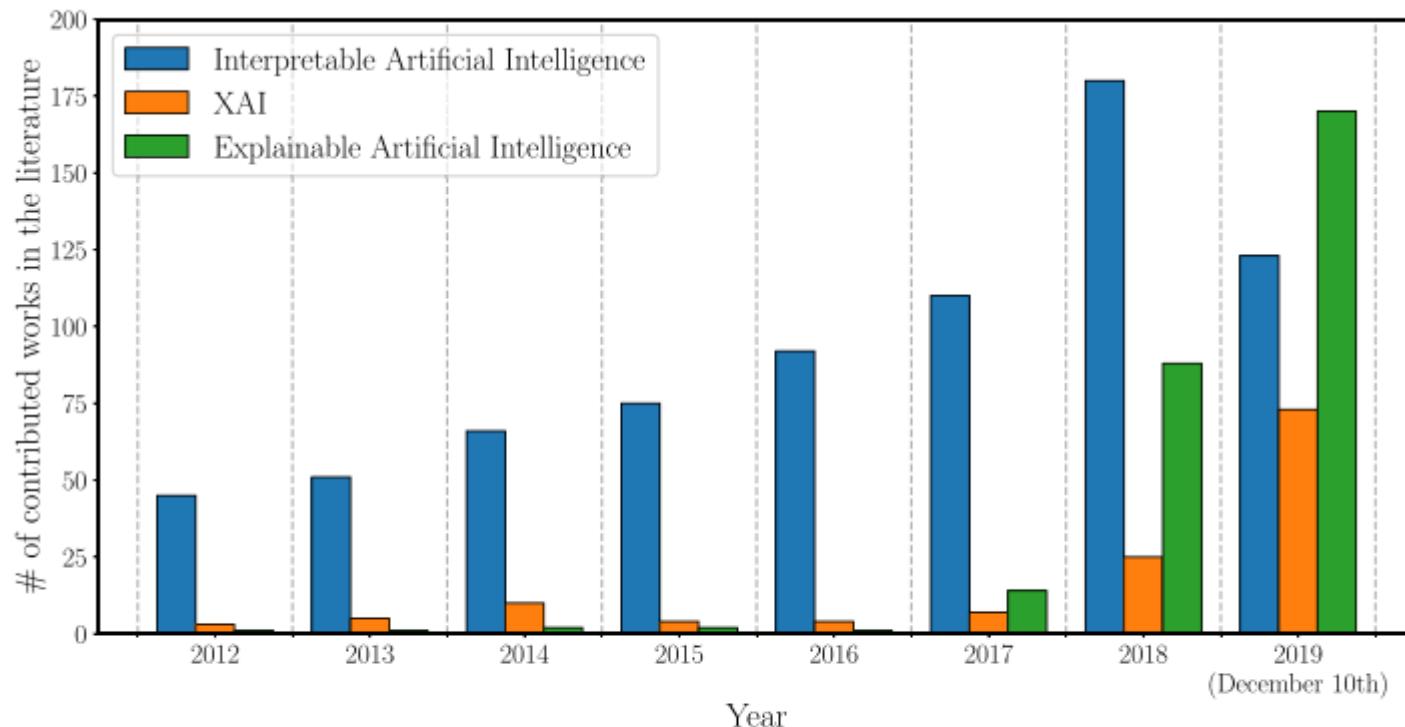
Overall picture of the classification of ML models attending to their level of explainability.

| Model                        | Transparent ML Models   |   |  | Post-hoc analysis  |
|------------------------------|---|---|--|--|
|                              | Simulability  | Decomposability   | Algorithmic Transparency   |  |
| Linear/Logistic Regression   | Predictors are human readable and interactions among them are kept to a minimum   | Variables are still readable, but the number of interactions and predictors involved in them have grown to force decomposition  | Variables and interactions are too complex to be analyzed without mathematical tools   | Not needed   |
| Decision Trees               | A human can simulate and obtain the prediction of a decision tree on his/her own, without requiring any mathematical background                                 | The model comprises rules that do not alter data whatsoever, and preserves their readability  | Human-readable rules that explain the knowledge learned from data and allows for a direct understanding of the prediction process  | Not needed   |
| K-Nearest Neighbors          | The complexity of the model (number of variables, their understandability and the similarity measure under use) matches human naive capabilities for simulation | The amount of variables is too high and/or the similarity measure is too complex to be able to simulate the model completely, but the similarity measure and the set of variables can be decomposed and analyzed separately | The similarity measure cannot be decomposed and/or the number of variables is so high that the user has to rely on mathematical and statistical tools to analyze the model | Not needed   |
| Rule Based Learners          | Variables included in rules are readable, and the size of the rule set is manageable by a human user without external help                                      | The size of the rule set becomes too large to be analyzed without decomposing it into small rule chunks   | Rules have become so complicated (and the rule set size has grown so much) that mathematical tools are needed for inspecting the model behaviour                           | Not needed   |
| General Additive Models      | Variables and the interaction among them as per the smooth functions involved in the model must be constrained within human capabilities for understanding      | Interactions become too complex to be simulated, so decomposition techniques are required for analyzing the model   | Due to their complexity, variables and interactions cannot be analyzed without the application of mathematical and statistical tools                                       | Not needed   |
| Bayesian Models              | Statistical relationships modeled among variables and the variables themselves should be directly understandable by the target audience                         | Statistical relationships involve so many variables that they must be decomposed in marginals so as to ease their analysis  | Statistical relationships cannot be interpreted even if already decomposed, and predictors are so complex that model can be only analyzed with mathematical tools          | Not needed   |
| Tree Ensembles               | x   | x   | x  | Needed: Usually <i>Model simplification or Feature relevance</i> techniques                |
| Support Vector Machines      | x   | x   | x  | Needed: Usually <i>Model simplification or Local explanations</i> techniques               |
| Multi-layer Neural Network   | x   | x   | x  | Needed: Usually <i>Model simplification, Feature relevance or Visualization</i> techniques |
| Convolutional Neural Network | x   | x   | x  | Needed: Usually <i>Feature relevance or Visualization</i> techniques                       |
| Recurrent Neural Network     | x   | x   | x  | Needed: Usually <i>Feature relevance</i> techniques  |

# Why need explanation?

- “**Right to an explanation**”, suggested e.g. in the GDPR art. 13, 14, and 22, regarding automated decision making? ([Goodman & Flaxman, 2017](#))
- Social issues in explainability ([Felten, 2017](#)):
  - **Confidentiality** (e.g. algorithm is trade secret or security risk)
  - **Complexity** (e.g. well understood but highly complex to give)
  - **Unreasonableness** (e.g. the use of rationally justified information to make decision that are not reasonable, discriminative or unfair)
  - **Injustice** (e.g. inconsistent with legal or moral code)

## Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI ([Arrieta, et.al., 2020](#))



**Fig. 1.** Evolution of the number of total publications whose title, abstract and/or keywords refer to the field of XAI during the last years. Data retrieved from Scopus® (December 10th, 2019) by using the search terms indicated in the legend when querying this database. It is interesting to note the latent need for interpretable AI models over time (which conforms to intuition, as interpretability is a requirement in many scenarios), yet it has not been until 2017 when the interest in techniques to explain AI models has permeated throughout the research community.

# Interpretability = Explainability



Artificial Intelligence  
Volume 267, February 2019, Pages 1-38



## Explanation in artificial intelligence: Insights from the social sciences

Tim Miller

### 2.1.5. Interpretability and justification

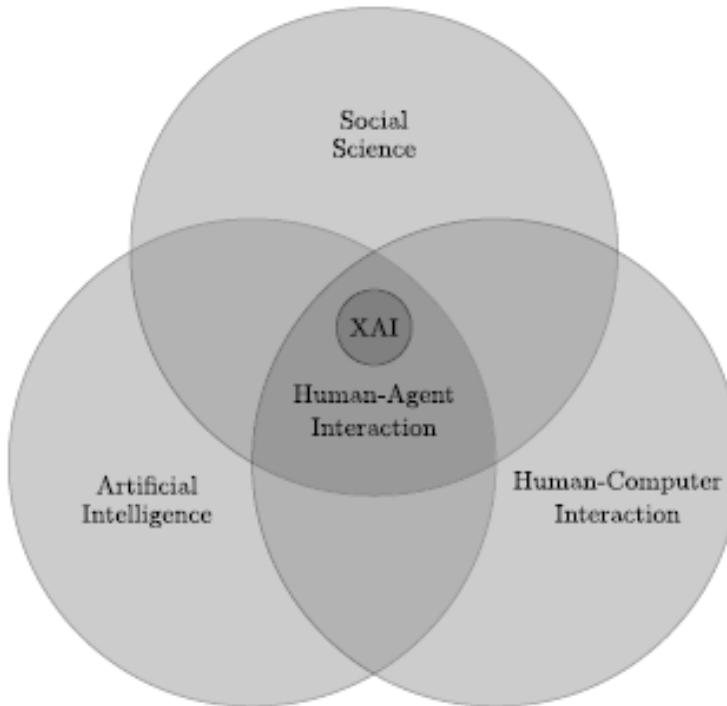
Here, we briefly address the distinction between *interpretability*, *explainability*, *justification*, and *explanation*, as used in this article; and as they seem to be used in artificial intelligence.

Lipton [103] provides a taxonomy of the desiderata and methods for interpretable AI. This paper adopts Lipton's assertion that explanation is post-hoc interpretability. I use Biran and Cotton [9]'s definition of interpretability of a model as: the degree to which an observer can understand the cause of a decision. Explanation is thus one mode in which an observer may obtain understanding, but clearly, there are additional modes that one can adopt, such as making decisions that are inherently easier to understand or via introspection. I equate 'interpretability' with 'explainability'.

A justification explains why a decision is good, but does not necessarily aim to give an explanation of the actual decision-making process [9].

# What makes a good explanation?

T. Miller / Artificial Intelligence 267 (2019) 1–38



**Fig. 1.** Scope of explainable artificial intelligence.

“.....models of how humans explain decisions and behaviour to each other are a good way to start analysing the problem ....”

# Interpretability $\neq$ Explainability

## Explaining Explanations: An Overview of Interpretability of Machine Learning

Leilani H. Gilpin, David Bau, Ben Z. Yuan, Ayesha Bajwa, Michael Specter and Lalana Kagal  
Computer Science and Artificial Intelligence Laboratory  
Massachusetts Institute of Technology  
Cambridge, MA 02139  
{lgilpin, davidbau, bzy, abajwa, specter, lkagal}@ mit.edu

We take the stance that **interpretability** alone is insufficient. In order for humans to trust black-box methods, we need **explainability** – models that are able to summarize the reasons for neural network behavior, gain the trust of users, or produce insights about the causes of their decisions. While interpretability is a substantial first step, these mechanisms need to *also* be complete, with the capacity to defend their actions, provide relevant responses to questions, and be audited. Although interpretability and explainability have been used interchangeably, we argue there are important reasons to distinguish between them. Explainable models are interpretable by default, but the reverse is not always true.

# DARPA's Explainable Artificial Intelligence Program

*David Gunning, David W. Aha*

■ Dramatic success in machine learning has led to a new wave of AI applications (for example, transportation, security, medicine, finance, defense) that offer tremendous benefits but cannot explain their decisions and actions to human users. DARPA's explainable artificial intelligence (XAI) program endeavors to create AI systems whose learned models and decisions can be understood and appropriately trusted by end users. Realizing this goal requires methods for learning more explainable models, designing effective explanation interfaces, and understanding the psychologic requirements for effective explanations. The XAI developer teams are addressing the first two challenges by creating ML techniques and developing principles, strategies, and human-computer interaction techniques for generating effective explanations. Another XAI team is addressing the third challenge by summarizing, extending, and applying psychologic theories of explanation to help the XAI evaluator define a suitable evaluation framework, which the developer teams will use to test their systems. The XAI teams completed the first of this 4-year program in May 2018. In a series of ongoing evaluations, the developer teams are assessing how well their XAM systems' explanations improve user understanding, user trust, and user task performance.

Advances in machine learning (ML) techniques promise to produce AI systems that perceive, learn, decide, and act on their own. However, they will be unable to explain their decisions and actions to human users. This lack is especially important for the Department of Defense, whose challenges require developing more intelligent, autonomous, and symbiotic systems. Explainable AI will be essential if users are to understand, appropriately trust, and effectively manage these artificially intelligent partners. To address this, DARPA launched its explainable artificial intelligence (XAI) program in May 2017. DARPA defines explainable AI as AI systems that can explain their rationale to a human user, characterize their strengths and weaknesses, and convey an understanding of how they will behave in the future. Naming this program explainable AI (rather than interpretable, comprehensible, or transparent AI, for example) reflects DARPA's objective to create more human-understandable AI systems through the use of effective explanations. It also reflects the XAI team's interest in the human psychology of explanation, which draws on the vast body of research and expertise in the social sciences.

# Whos'who in XAI?

<https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8466590>

## Funding agency and academia



Funding agency, interest on  
intelligent and autonomous  
product

A banner for the ACM Conference on Fairness, Accountability, and Transparency (ACM FAccT). The background is light blue with abstract geometric shapes. The title "ACM Conference on Fairness, Accountability, and Transparency (ACM FAccT)" is displayed prominently in large, dark blue serif font. A subtitle below it reads: "A computer science conference with a cross-disciplinary focus that brings together researchers and practitioners interested in fairness, accountability, and transparency in socio-technical systems." There are small icons of a person, a gear, and a leaf at the bottom of the banner.

ACM Conference on  
Fairness, Accountability,  
and Transparency (ACM  
FAccT)

A computer science conference with a cross-disciplinary focus that brings together researchers and practitioners interested in fairness, accountability, and transparency in socio-technical systems.

ACM FAccT is the new acronym for the ACM Conference on Fairness, Accountability, and Transparency!

ACM FAccT will be held online March 3-10, 2021. More information on this year's conference is posted on the [2021 ACM FAccT webpage](#).

"FAT" academia (fairness accountability, and transparency): promoting and enabling explainability and **fairness** in algorithmic decision-making system with **social and commercial impact**

Company/business community,  
focusing on the demand by  
**industry regulation**



"Our mission is to build Explainable AI products and solutions that help to optimize human cognitive performance. A cornerstone of that mission is never to operate as a 'black box,'" said Welsh. "Explainable AI means that the system can justify its reasoning. Kyndi's product exists because Deep Learning is a 'black box' and cannot be used in regulated industries where organizations are required to explain the reasons for any advice on any decision."



# Outline

- Interpretability = Transparency ≠ Explainability?
- DARPA XAI Grant program
- Results of DARPA XAI Grant program
- Key takeaways and discussions



DEFENSE ADVANCED  
RESEARCH PROJECTS AGENCY

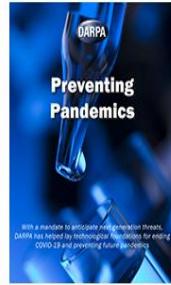
ABOUT US / OUR RESEARCH /

› Defense Advanced Research Projects Agency › Advancing National Security Through Fundamental Research

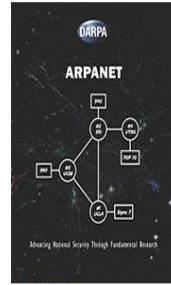
## Advancing National Security Through Fundamental Research

For sixty years, DARPA has held to a singular and enduring mission: to make pivotal investments in breakthrough technologies for national security.

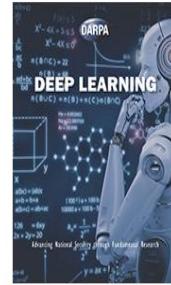
The genesis of that mission and of DARPA itself dates to the launch of Sputnik in 1957, and a commitment by the United States that, from that time forward, it would be the initiator and not the victim of strategic technological surprises. Working with innovators inside and outside of government, DARPA has repeatedly delivered on that mission, transforming revolutionary concepts and even seeming impossibilities into practical capabilities. The ultimate results have included not only game-changing military capabilities such as precision weapons and stealth technology, but also such icons of modern civilian society such as the Internet, automated voice recognition and language translation, and Global Positioning System receivers small enough to embed in myriad consumer devices.



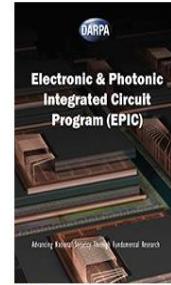
Preventing Pandemics



ARPANET



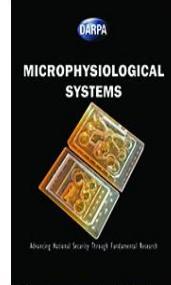
Deep Learning



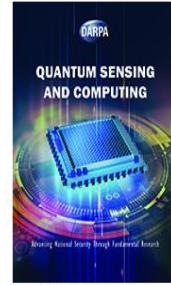
EPIC



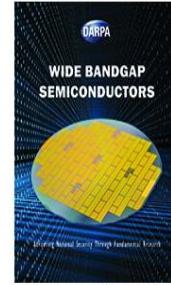
MAKE-IT



Microphysiological Systems



Quantum Sensing



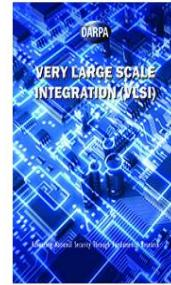
Wide Bandgap Semiconductor



Stealth



Sigma



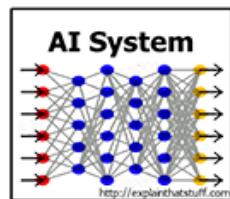
Very Large Scale Integrated Circuits



> Defense Advanced Research Projects Agency > Our Research > Explainable Artificial Intelligence

# Explainable Artificial Intelligence (XAI)

Dr. Matt Turek



- 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?

## RESOURCES

[DARPA-BAA-16-53](#)

[DARPA-BAA-16-53: Proposers Day Slides](#)

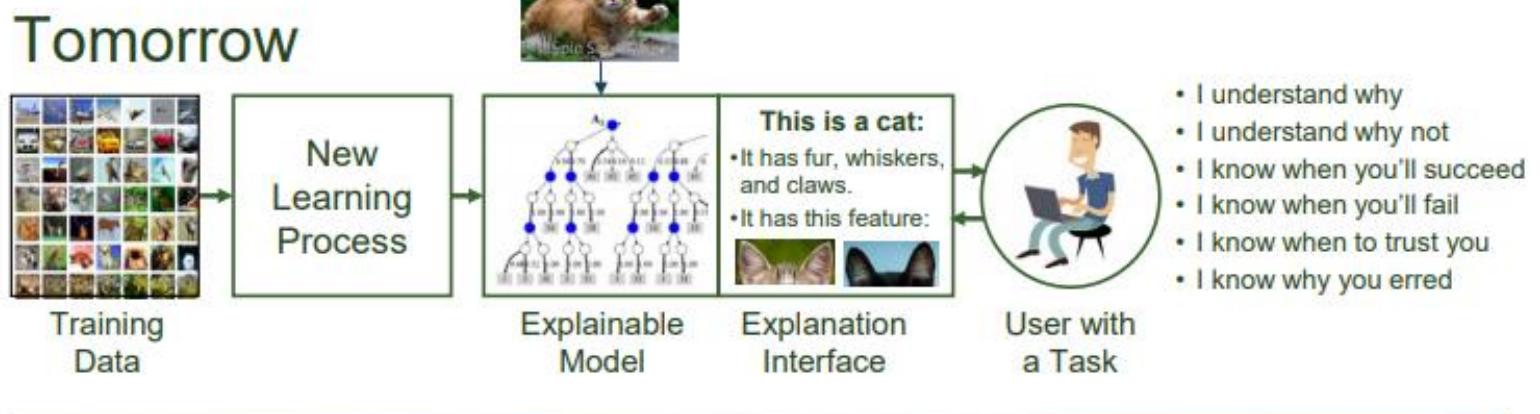
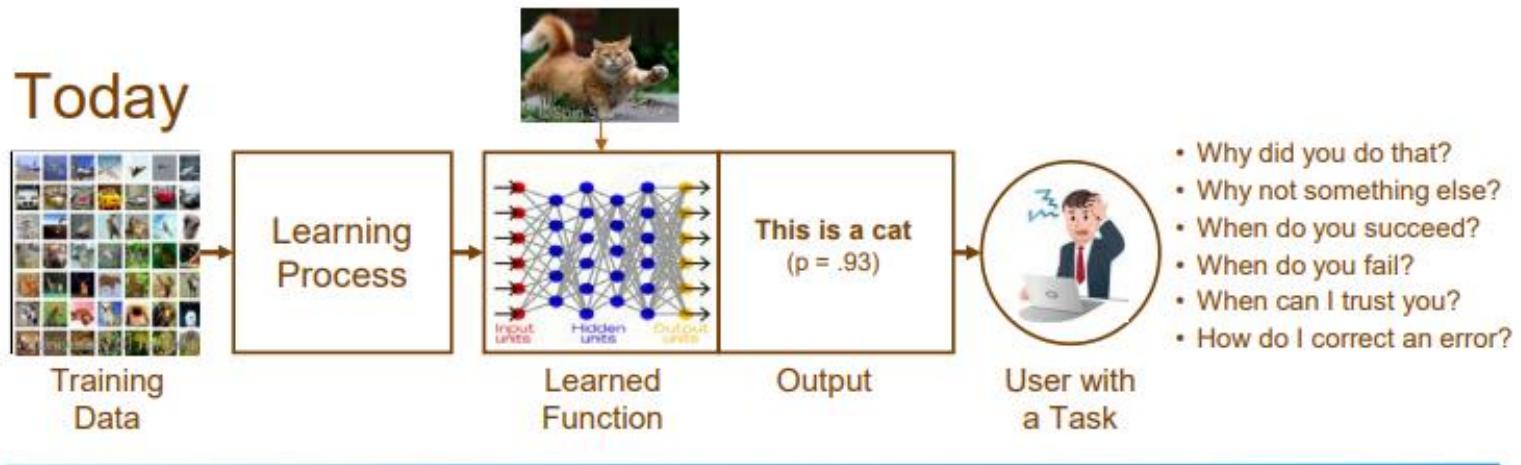
[XAI Program Portfolio](#)

Figure 1. The Need for Explainable AI

## PART I: OVERVIEW INFORMATION

- **Federal Agency Name:** Defense Advanced Research Projects Agency (DARPA), Information Innovation Office (I2O)
- **Funding Opportunity Title:** Explainable Artificial Intelligence (XAI)
- **Announcement Type:** Initial Announcement
- **Funding Opportunity Number:** DARPA-BAA-16-53
- **Catalog of Federal Domestic Assistance Numbers (CFDA):** 12.910 Research and Technology Development
- **Dates**
  - Posting Date: August 10, 2016
  - Proposers Day: August 11, 2016
  - Abstract Due Date: September 1, 2016, 12:00 noon (ET)
  - Proposal Due Date: November 1, 2016, 12:00 noon (ET)
  - BAA Closing Date: November 1, 2016, 12:00 noon (ET)
- **Anticipated Individual Awards:** Multiple awards are anticipated.
- **Types of Instruments that May be Awarded:** Procurement Contracts, Cooperative Agreements, or Other Transactions (OTs). No grants will be awarded under this solicitation.
- **Agency Contacts**
  - **Technical POC:** David Gunning, Program Manager, DARPA/I2O
  - **BAA Email:** [XAI@darpa.mil](mailto:XAI@darpa.mil)
  - **BAA Mailing Address:**  
DARPA/I2O  
ATTN: DARPA-BAA-16-53  
675 North Randolph Street  
Arlington, VA 22203-2114
  - **I2O Solicitation Website:** <http://www.darpa.mil/work-with-us/opportunities>

What are they trying to do.....



# Program Scope

Deep explanation refers to modified or hybrid DL techniques that learn more explainable features or representations or that include explanation generation facilities. Several design choices might produce

<https://www.darpa.mil/attachments/DARPA-BAA-16-53.pdf>

Interpretable models are ML techniques that learn more structured, interpretable, or causal models. Early examples included Bayesian rule lists (Letham et al. 2015), Bayesian program learning, learning models of causal relationships, and use of stochastic grammars to learn more interpretable structure.

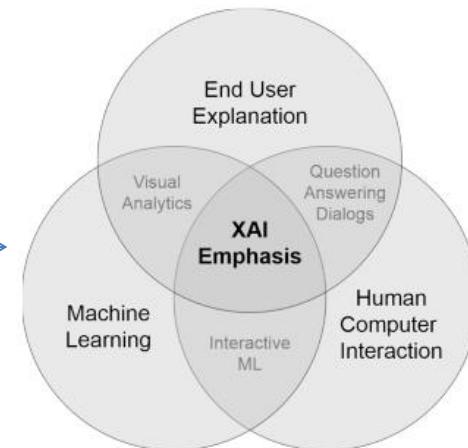
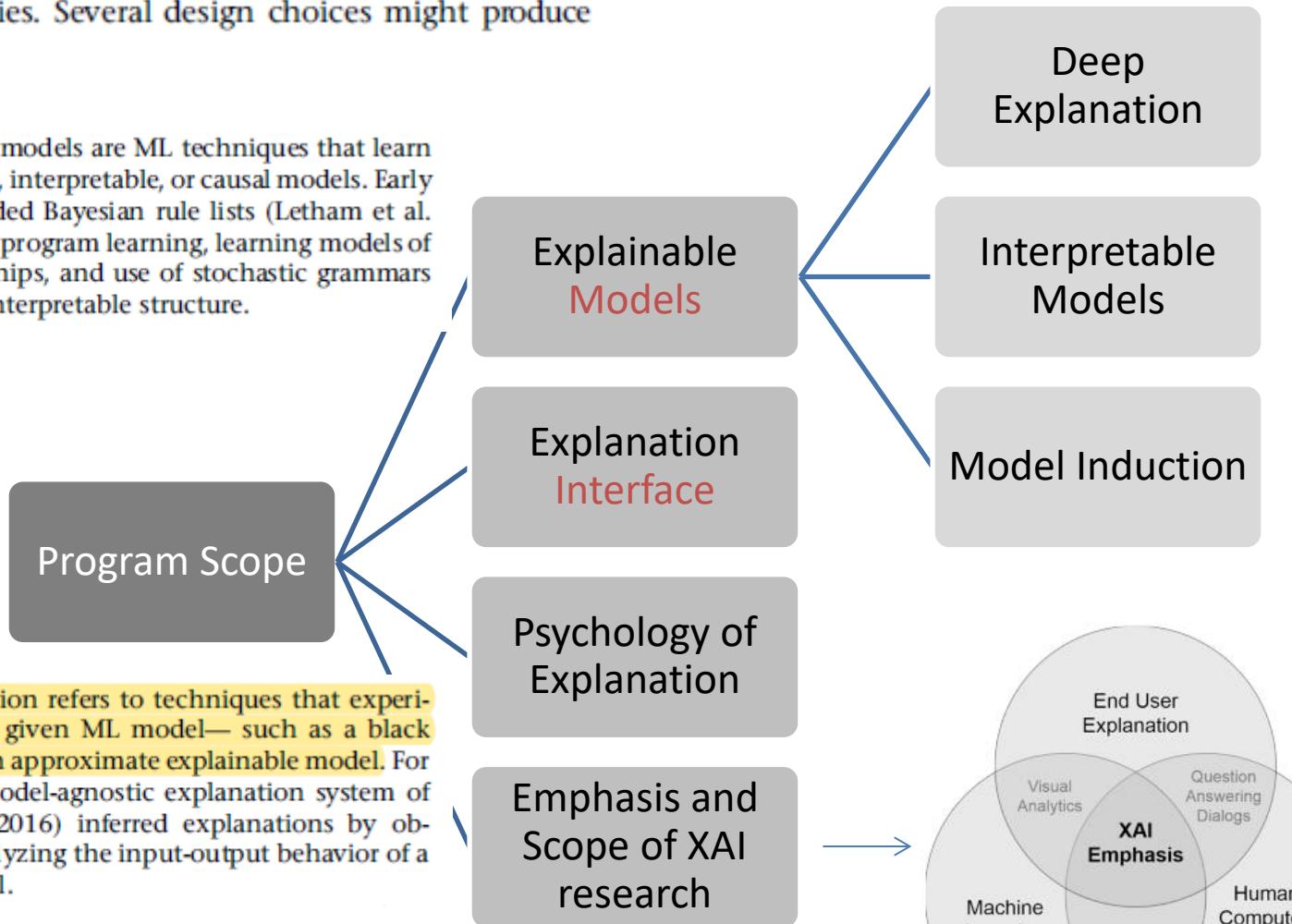


Figure 2: XAI Emphasis

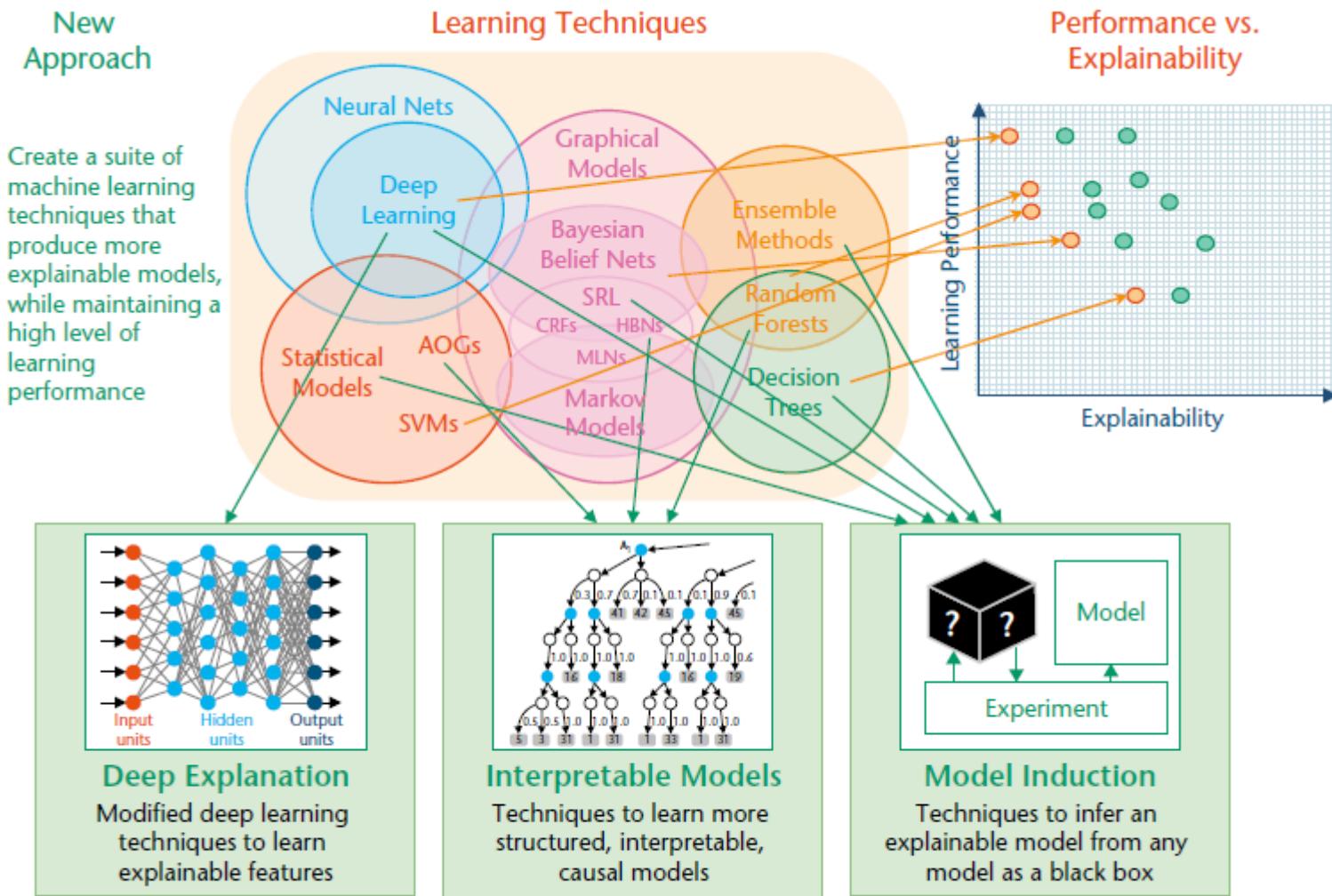
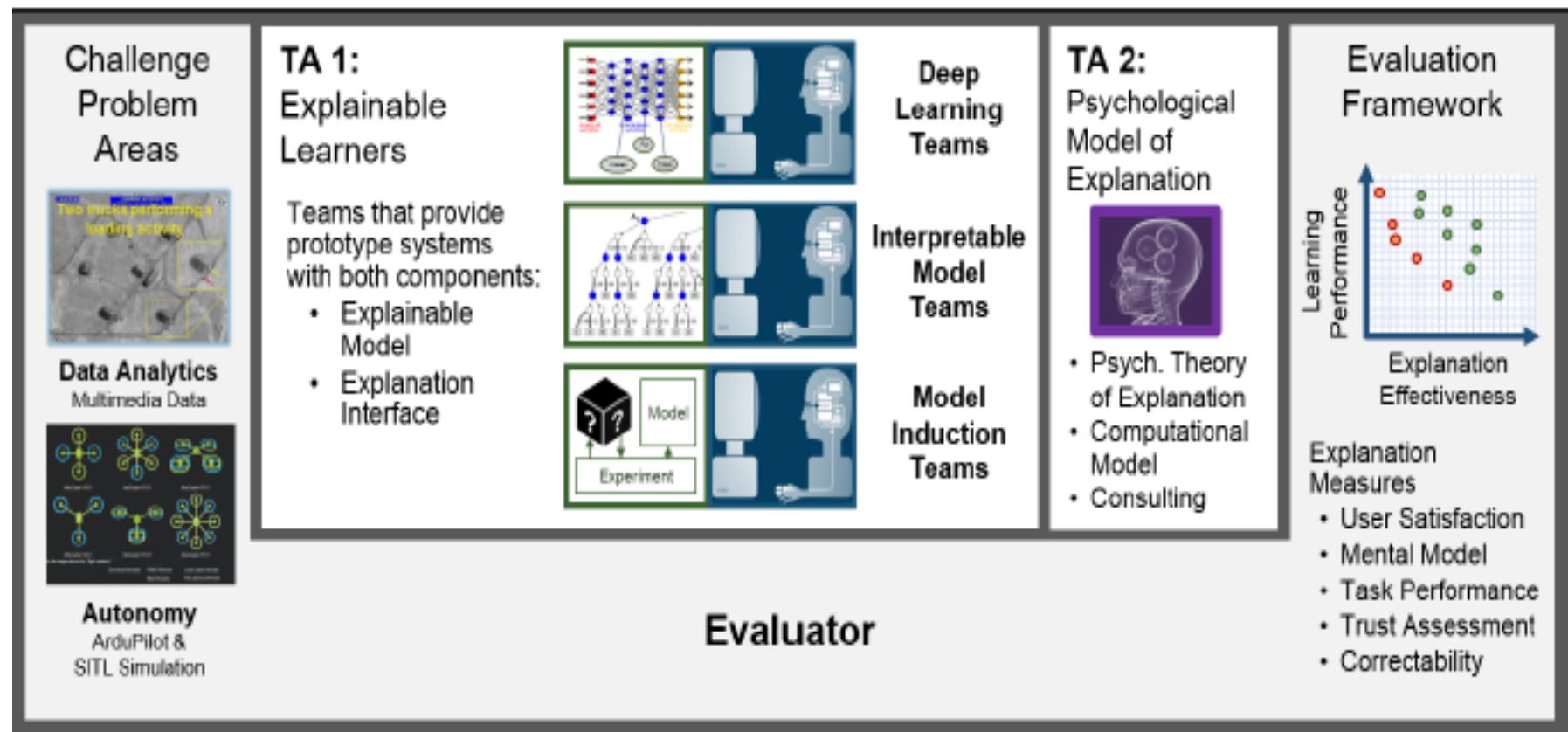


Figure 2. Strategies for Developing Explainable Models.

## Program Structure

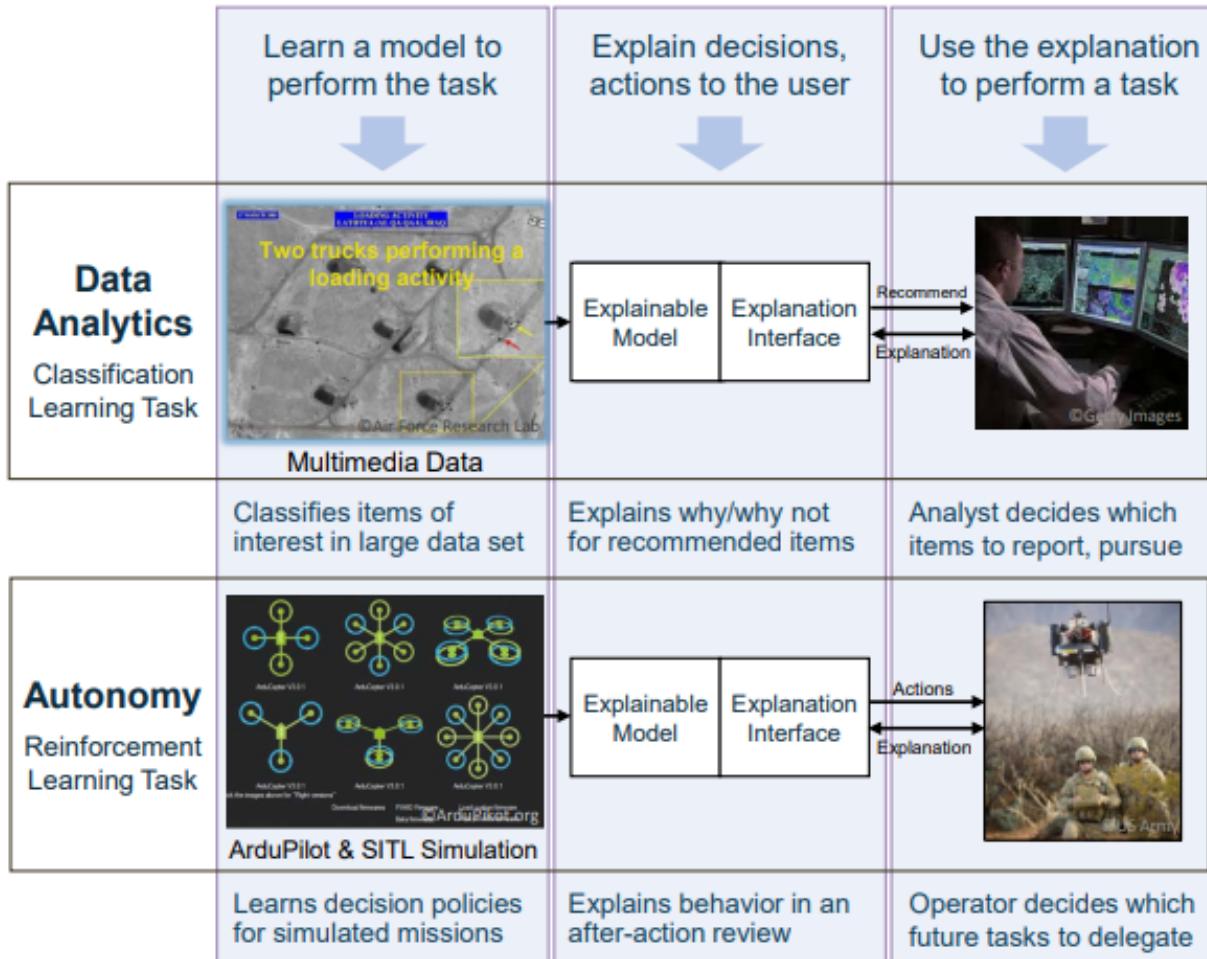


Data analytics: classification of events of interest in heterogeneous multimedia data

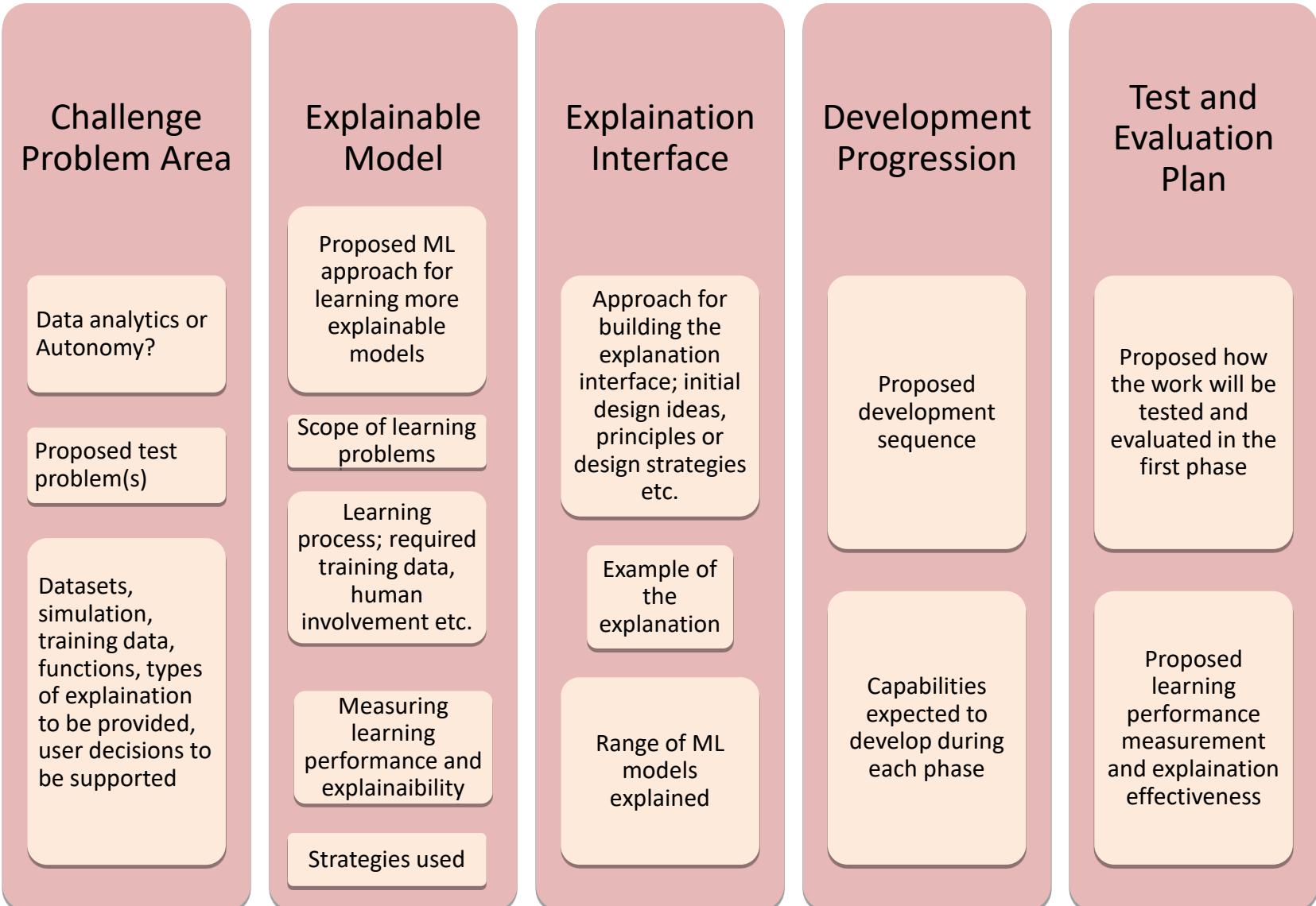
Autonomy: decision policies for autonomous systems



## Explainable AI – Challenge Problem Areas



## TA 1: Explainable Learners



## TA 2: Psychological Model of Explanation

### Theories of Explanation

Summary of current psy theories of explanation and how to further develop and refine it

How the theories will inform the development of TA 1

### Computational Model

How to develop and implement a computational model of explanation from theories

Model element from the theory

Identify predictions to be tested

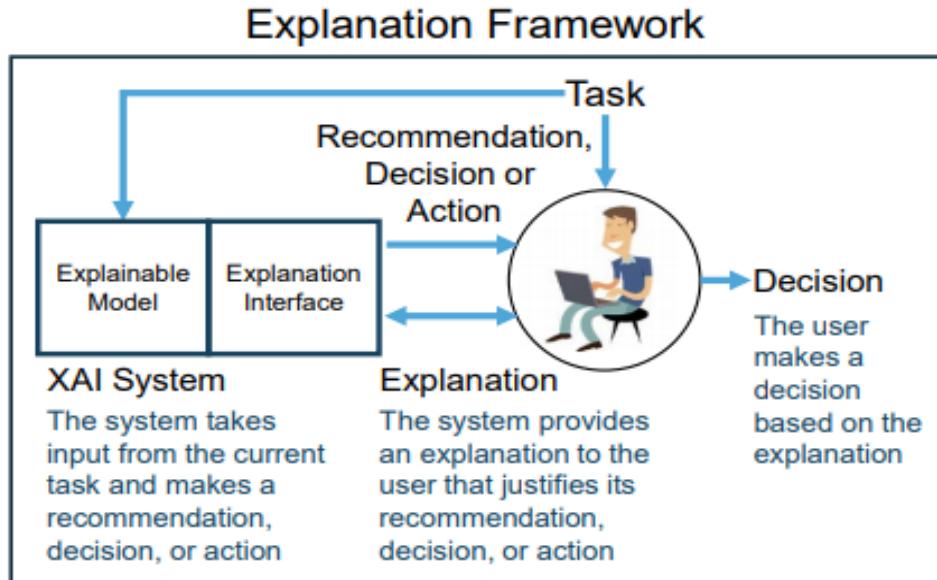
Explain how to test and refine the model

### Model Validation

Validating the computational model of the theory and how to test the model against TA 1 evaluation results in Phase 2



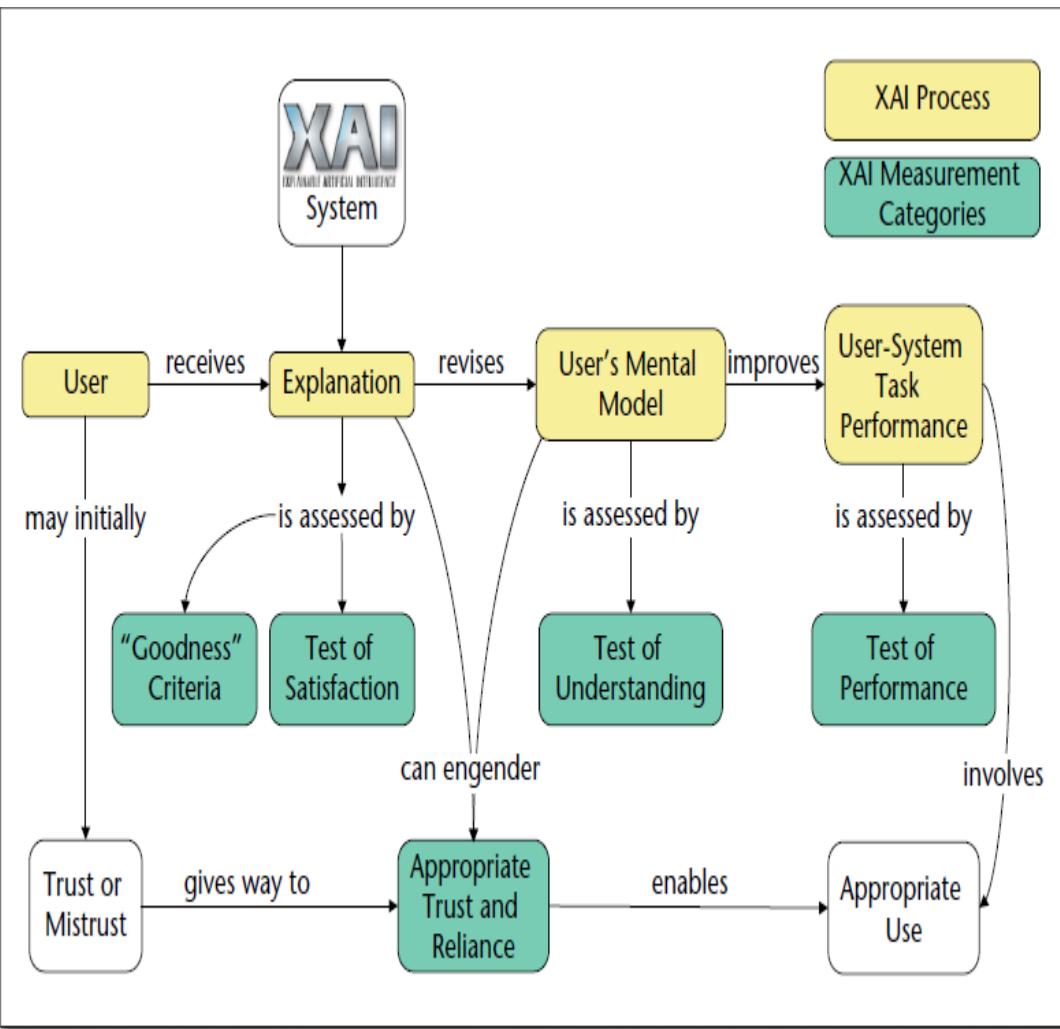
## Explainable AI – Measuring Evaluation Effectiveness



| Measure of Explanation Effectiveness |  |
|--------------------------------------|--|
| <b>User Satisfaction</b>             | <ul style="list-style-type: none"><li>• Clarity of the explanation (user rating)</li><li>• Utility of the explanation (user rating)</li></ul>  |
| <b>Mental Model</b>                  | <ul style="list-style-type: none"><li>• Understanding individual decisions</li><li>• Understanding the overall model</li><li>• Strength/weakness assessment</li><li>• 'What will it do' prediction</li><li>• 'How do I intervene' prediction</li></ul> |
| <b>Task Performance</b>              | <ul style="list-style-type: none"><li>• Does the explanation improve the user's decision, task performance?</li><li>• Artificial decision tasks introduced to diagnose the user's understanding</li></ul>  |
| <b>Trust Assessment</b>              | <ul style="list-style-type: none"><li>• Appropriate future use and trust</li></ul>   |
| <b>Correctability</b>                | <ul style="list-style-type: none"><li>• Identifying errors</li><li>• Correcting errors</li><li>• Continuous training</li></ul>   |

| Measure  | Description   |
|--|---|
| <b>ML Model performance</b>                              |   |
| Various measures (on a per-challenge problem area basis) | Accuracy/performance of the ML model in its given domain (to understand whether performance improved or degraded relative to state-of-the-art nonexplainable baselines) |
| <b>Explanation Effectiveness</b>                         |   |
| Explanation goodness                                     | Features of explanations assessed against criteria for explanation goodness   |
| Explanation satisfaction                                 | User's subjective rating of explanation completeness, usefulness, accuracy, and satisfaction  |
| Mental model understanding                               | User's understanding of the system and the ability to predict the system's decisions/behavior in new situations   |
| User task performance                                    | Success of the user performing the tasks for which the system is designed to support  |
| Appropriate Trust and Reliance                           | User's ability to know when to, and when not to, trust the system's recommendations and decisions   |

*Table 1. Measurement Categories.*



## 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

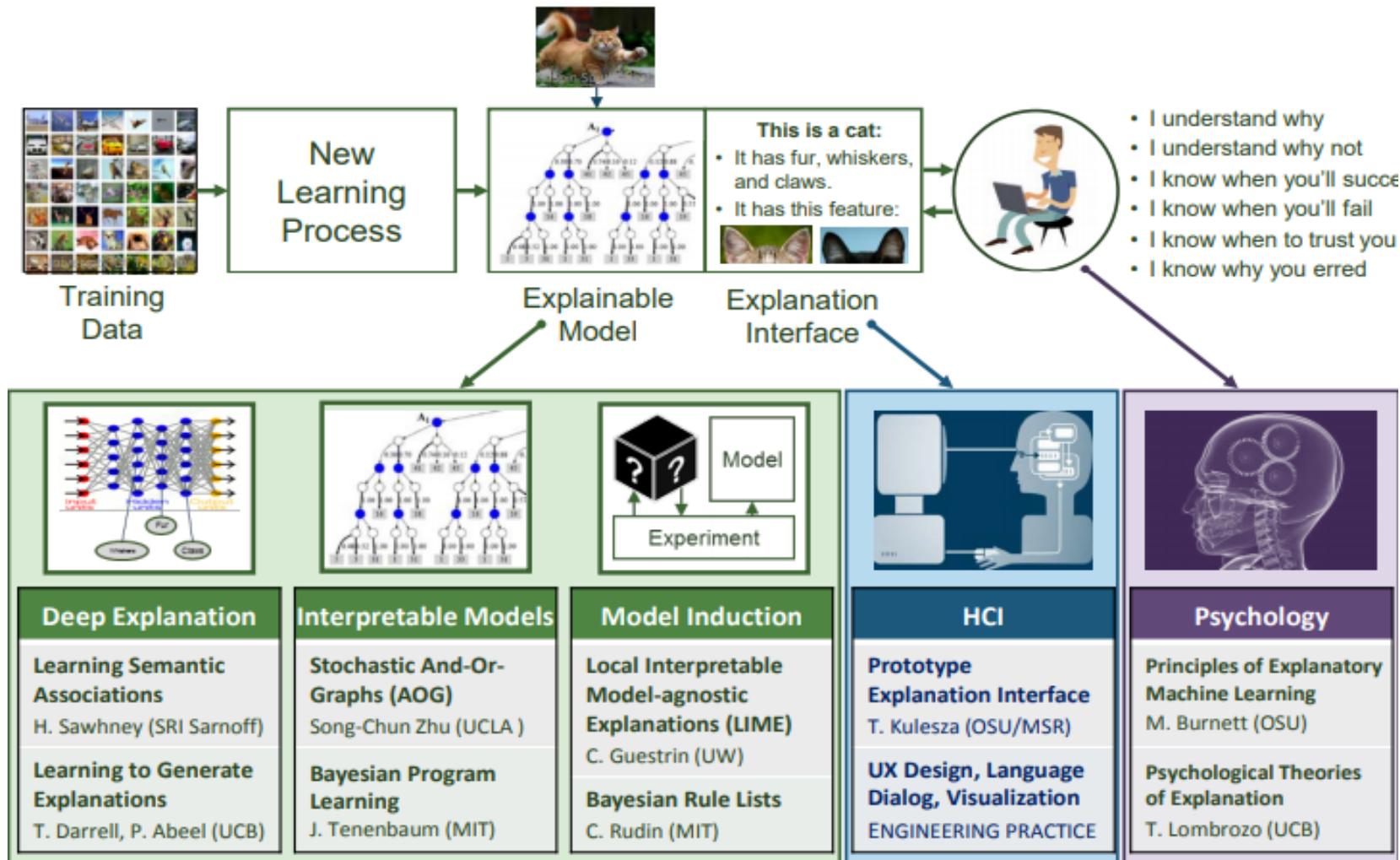
**Partial Explanation** - The explainable learning system is used to perform a task and generates only partial or ablated explanations (to assess various explanation features)

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

Figure 8. Initial Model of the Explanation Process and Explanation Effectiveness Measurement Categories.



## Explainable AI – Why Do You Think It Will Be Successful?



## Program Schedule

| 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 | 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><br><b>Technology Demonstrations</b>  |   |  |     |     |     |     |     |     |     |     |     |  |  |                 |  |     |     |     |     |     |     |     |     |   | <b>PHASE 2</b><br><b>Government Evaluations</b>   |  |        |  |                    |     |     |     |     |     |     |                                    |     |        |   |     |  |   |     |     |     |     |     |        |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Evaluator (NRL)   |   |  |     |     |     |     |     |     |     |     |     |  |  |                 |  |     |     |     |     |     |     |     |     |   | Explainable Learning Systems (11 Developer Teams) |  |        |  |                    |     |     |     |     |     |     |                                    |     |        |   |     |  |   |     |     |     |     |     |        |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Define Evaluation Framework   |   |  |     |     |     |     |     |     |     |     |     | Prepare for Eval 1                       | Eval 1                                       | Analyze Results | Prepare for Eval 2   |     |     |     |     |     |     |     |     |   |   |  | Eval 2 | Analyze Results  | Prepare for Eval 3 |     |     |     |     |     |     |                                    |     |        |   |     | Eval 3   | Analyze Results; Accept Software Libraries/Toolkits |     |     |     |     |     |        |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Develop and Demonstrate Explainable Learning Systems                                      |   |  |     |     |     |     |     |     |     |     |     | Eval 1                                   | Refine and Test Explainable Learning Systems |                 |  |     |     |     |     |     |     |     |     |   | Eval 2  | Refine and Test Explainable Learning Systems |        |  |                    |     |     |     |     |     |     |                                    |     | Eval 3 | Deliver Software Libraries and Toolkits |     |  |   |     |     |     |     |     |        |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Summarize Current Psychological Theories of Explanation                                   |   |  |     |     |     |     |     |     |     |     |     | Develop Theoretical Model of Explanation |  |                 |  |     |     |     |     |     |     |     |     | Refine and Test Model of Explanation to Support System Development and Evaluation |   |  |        |  |                    |     |     |     |     |     |     | Deliver Final Model of Explanation |     |        |   |     |  |   |     |     |     |     |     | Final  |     |     |     |     |     |     |     |     |     |     |     |     |     |
|  Kickoff |  Progress Report |  Tech Demos |     |     |     |     |     |     |     |     |     |  |  |                 |  Eval 1 Results |     |     |     |     |     |     |     |     |   |   |  |        |  Eval 2 Results |                    |     |     |     |     |     |     |                                    |     |        |   |     |  Eval 3 Results |   |     |     |     |     |     |        |     |     |     |     |     |     |     |     |     |     |     |     |     |

Figure 6. XAI Program Schedule.

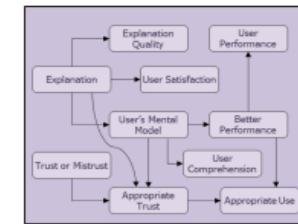
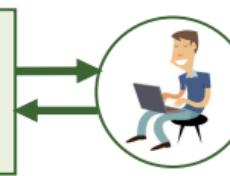
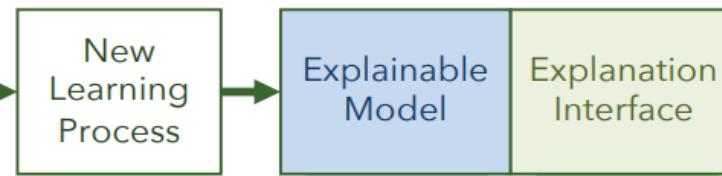
# Outline

- Interpretability = Transparency ≠ Explainability?
- DARPA XAI Grant program
- [\*\*Results: Grantees and reports\*\*](#)
- Key takeaways and discussions

## The Grantees!



## XAI Developers and Technical Approaches



IHMC  
Psychological Models  
of Explanation

| 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          |

UC Berkeley



# 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)

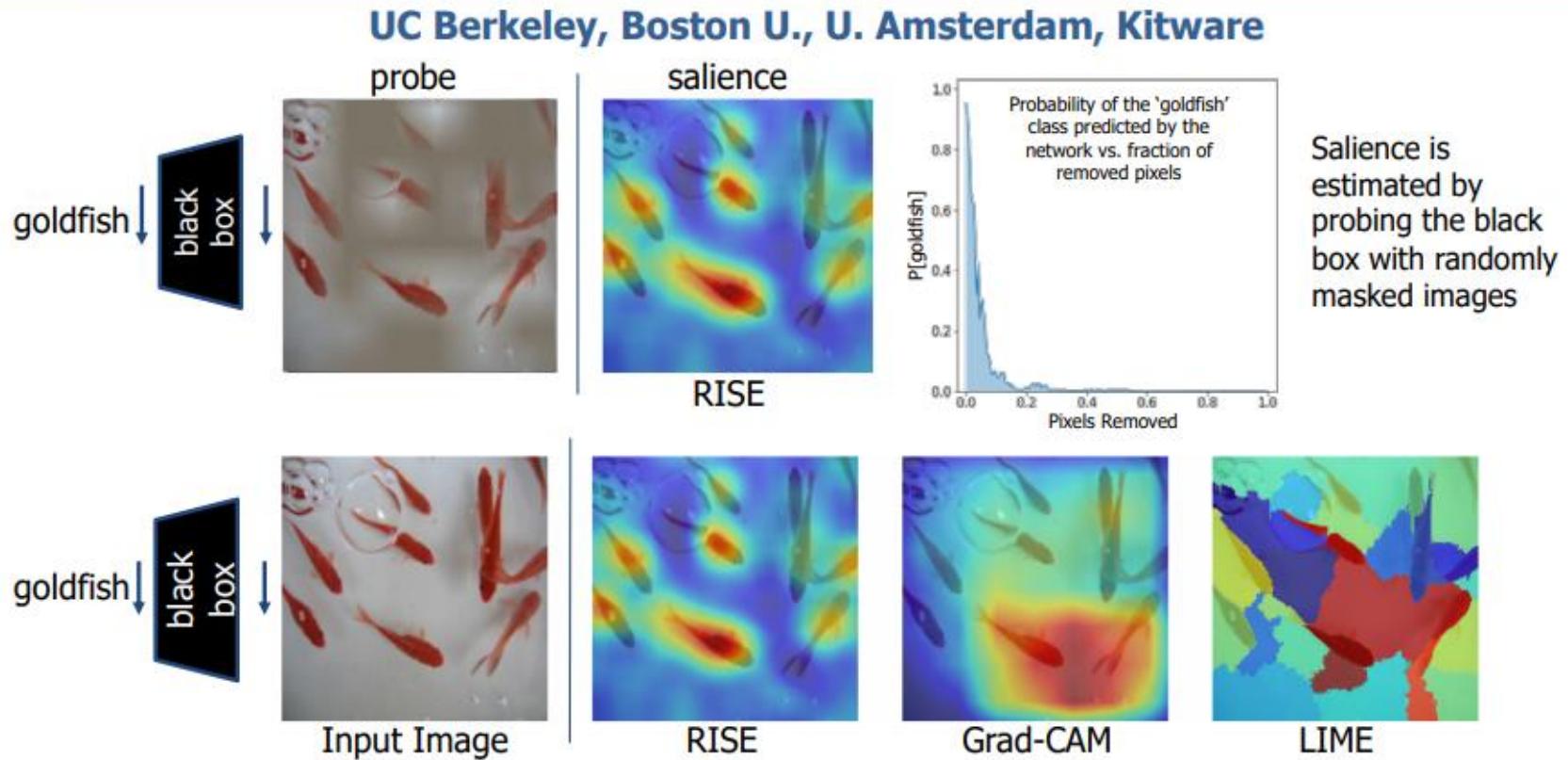
### Data Analytics

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

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

- |                               |                             |                           |
|-------------------------------|-----------------------------|---------------------------|
| • Pieter Abbeel (UC Berkeley) | • Dan Klein (UC Berkeley)   | • Anthony Hoogs (Kitware) |
| • Tom Griffiths (UC Berkeley) | • John Canny (UC Berkeley)  |                           |
| • Kate Saenko (Boston U.)     | • Anca Dragan (UC Berkeley) |                           |
| • Zeynep Akata (U. Amsterdam) |                             |                           |

## UC Berkeley



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

Approved for public release: distribution unlimited.

## UC Berkeley

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?

Accuracy of Users Judgement

No

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

## Charles River Analytics



# CAMEL: Causal Models to Explain Learning



### 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

#### Explanation Interface

##### Narrative Generation

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

#### Challenge Problem

##### Autonomy

- Atari
- Starcraft

##### Data Analytics

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

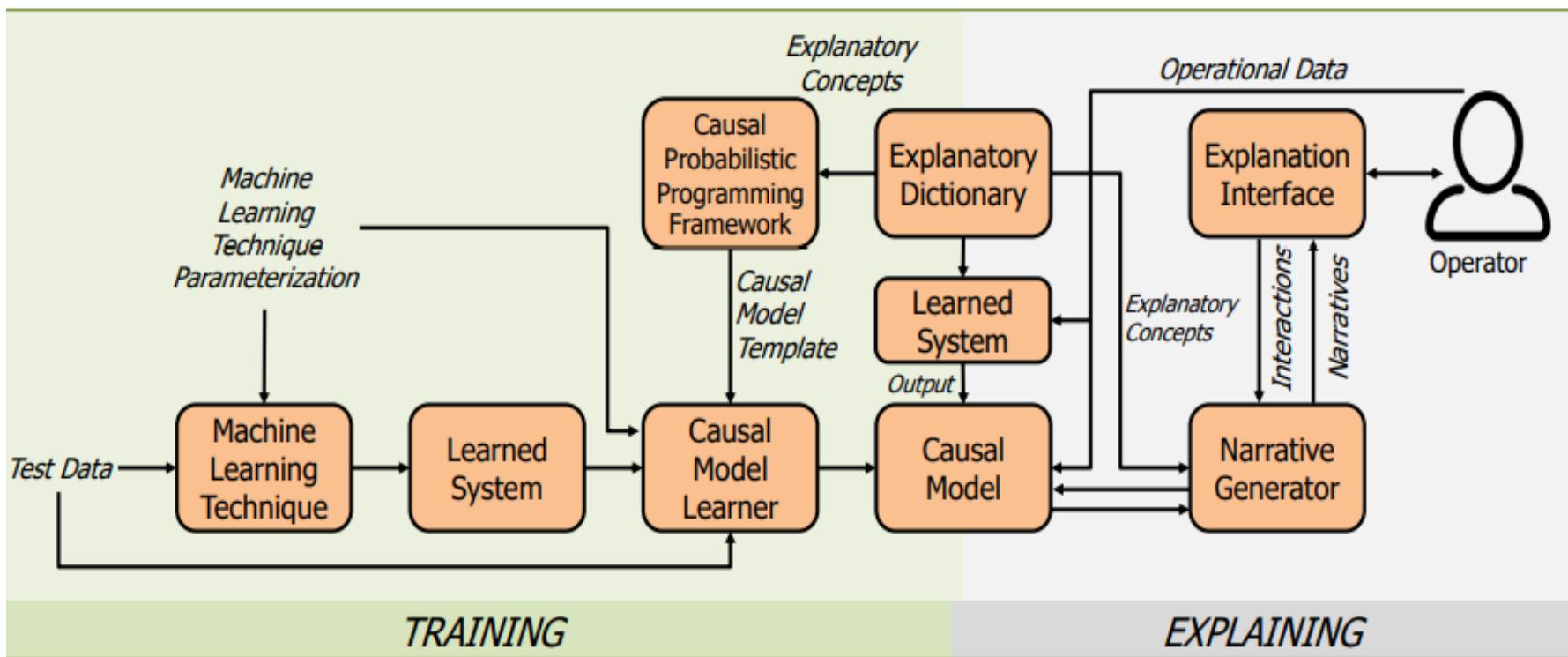
- PI: James Tittle (CRA)

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

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

## Charles River Analytics

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



## Learning and Communicating Explainable Representations



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

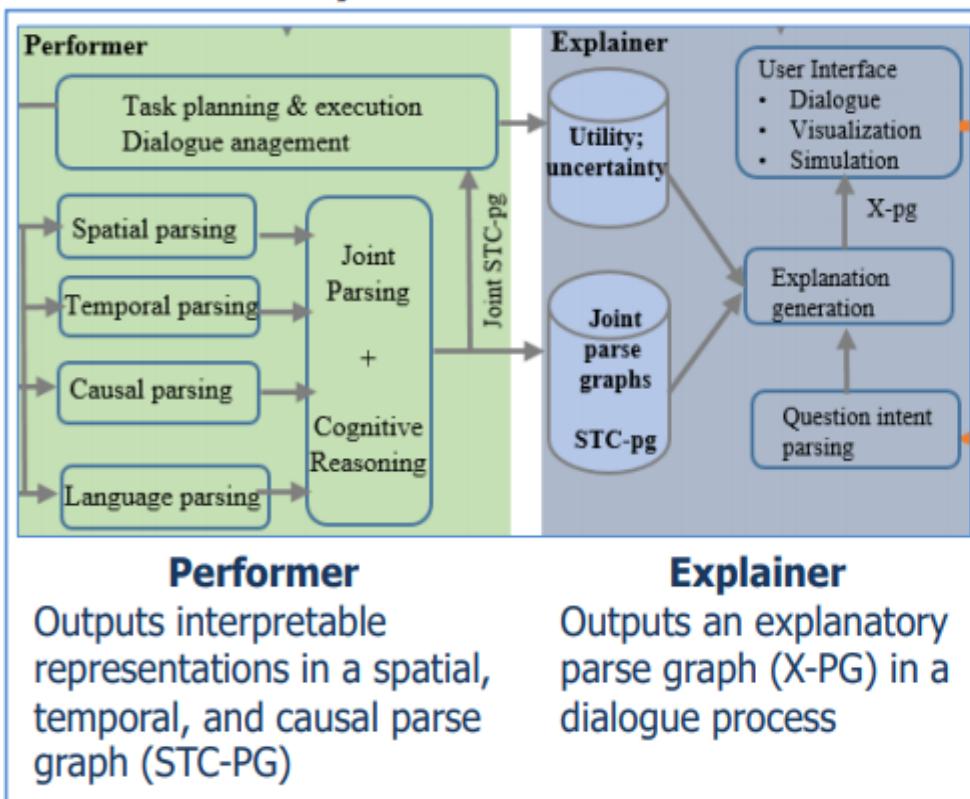
- PI: Song-Chun Zhu (UCLA)

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

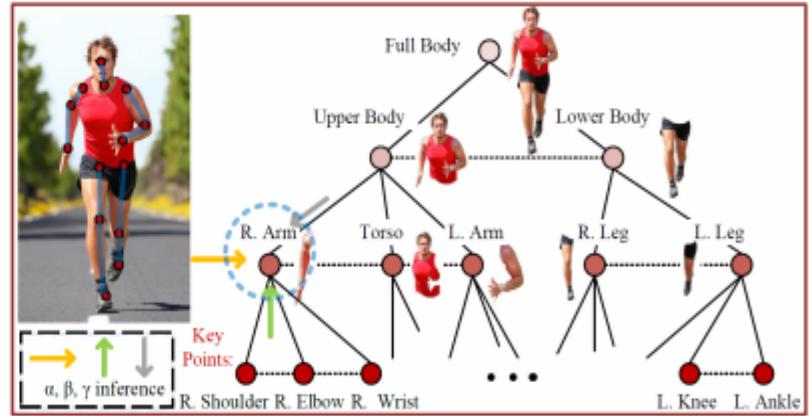
- Joyce Chai (Michigan State)

UCLA

## System Architecture



## STC Parse Graph



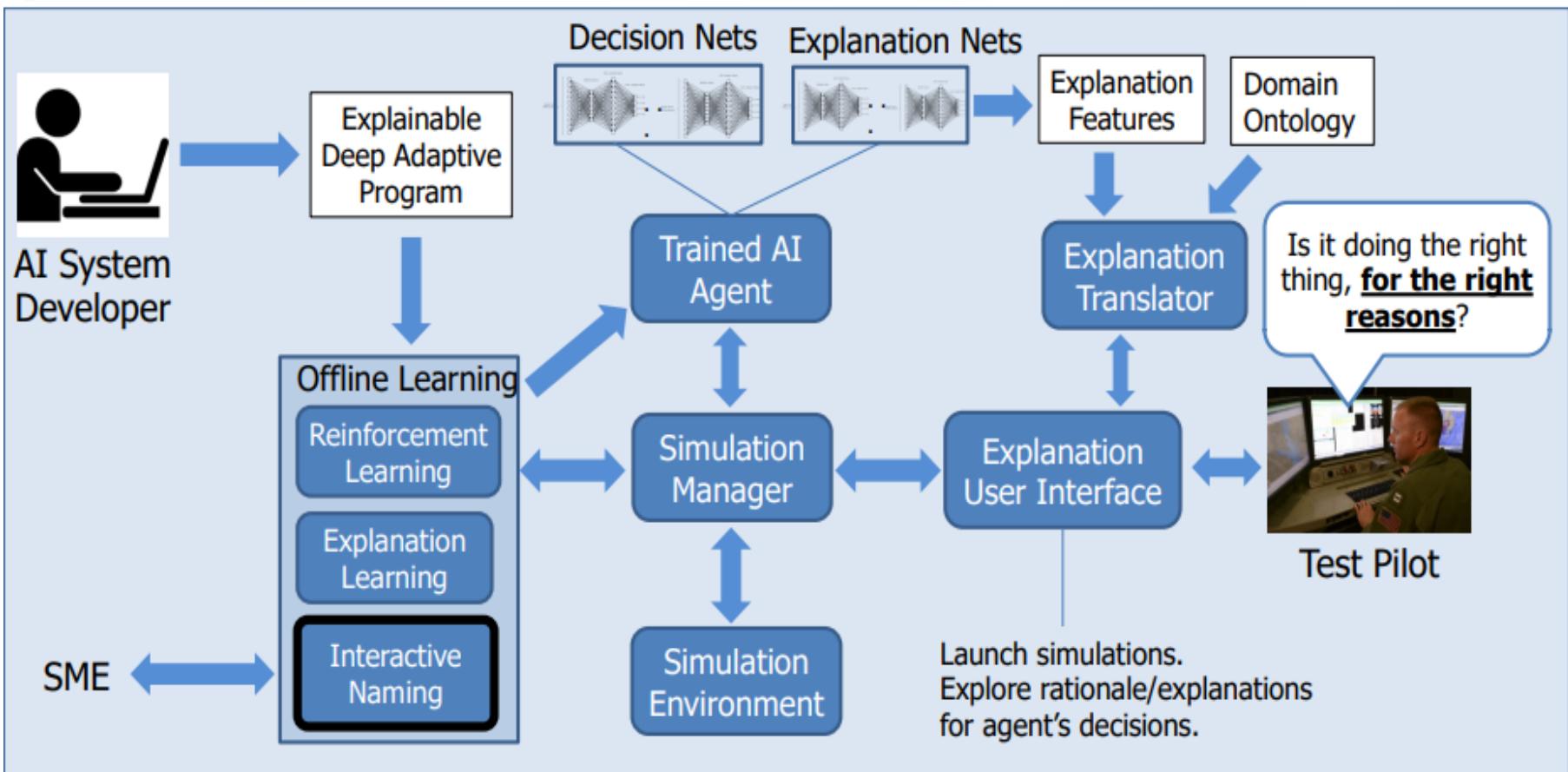
An attributed parse graph for a running person. Each node has 3 computing channels:

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

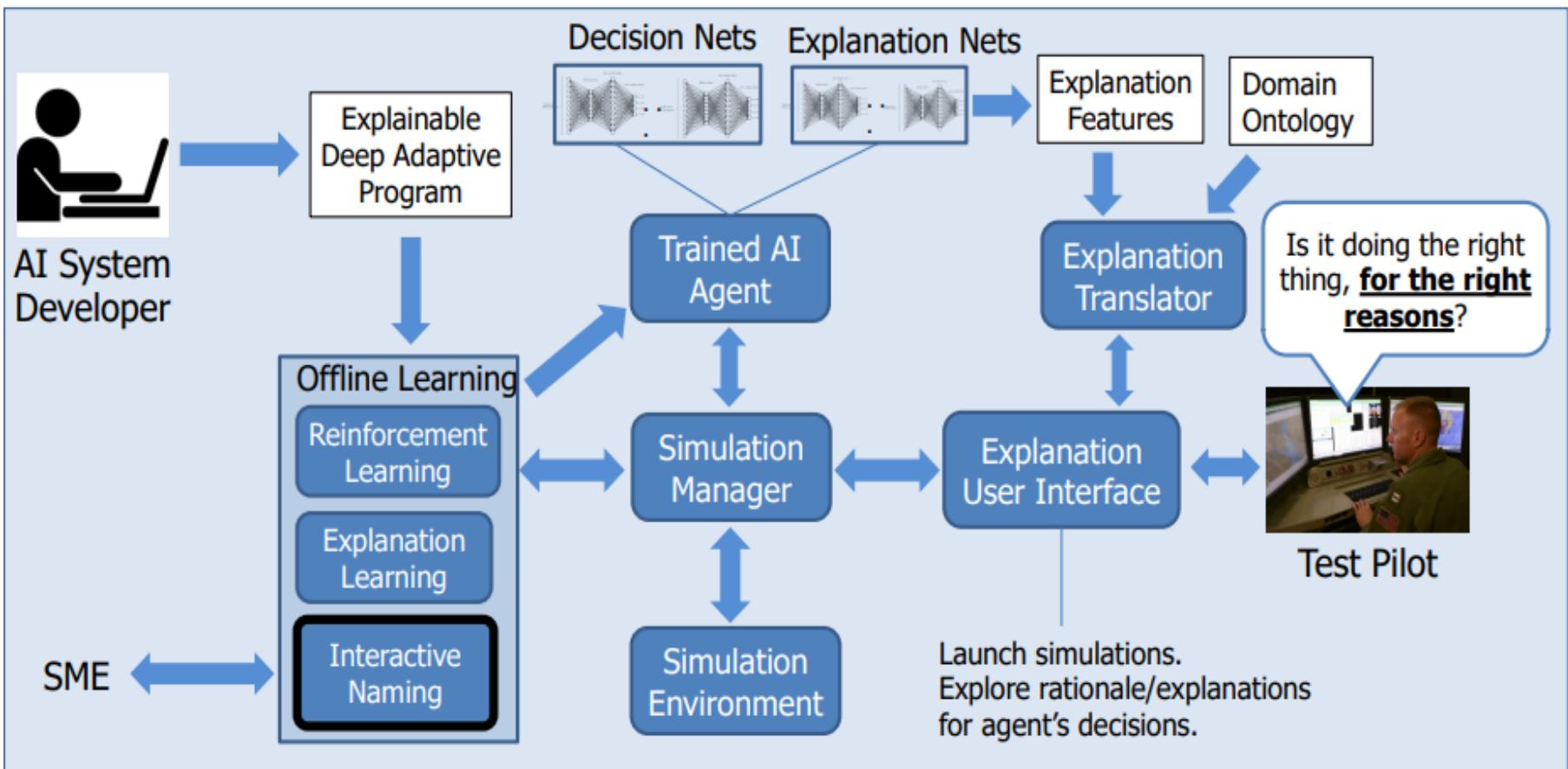
An explanation is represented as parse graph X-pg

Oregon State

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



## Oregon State



PARC



## COGLE: Common Ground Learning and Explanation



**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 Stefk (PARC)

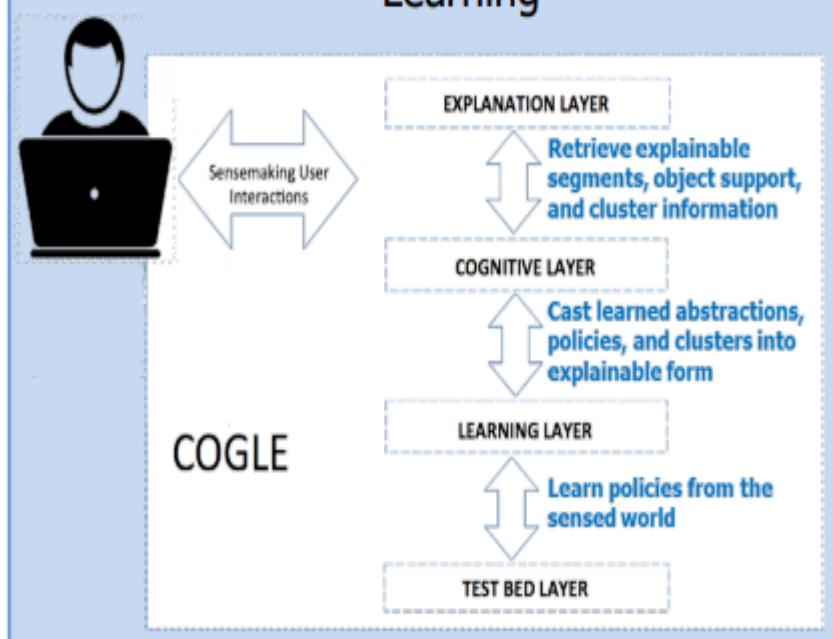
- Honglak Lee (U. Michigan)
- Subramanian Ramamoorthy (U. Edinburgh)

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

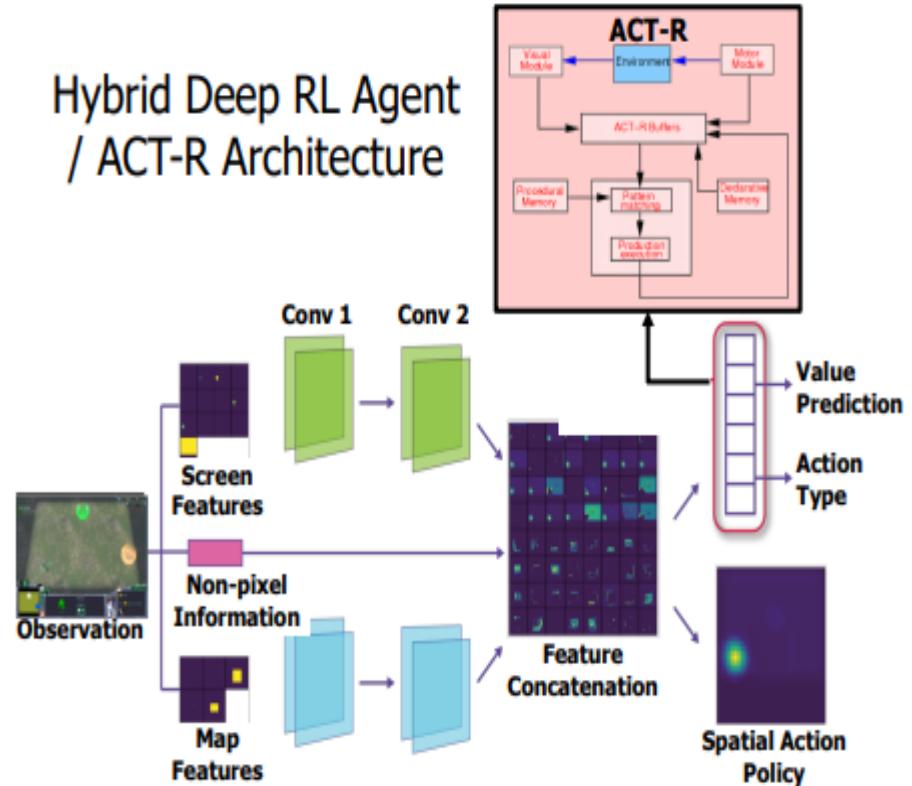
- Michael Youngblood (PARC)

## PARC

### Layered Cognitive Architecture to Partition Explanation And Learning



### Hybrid Deep RL Agent / ACT-R Architecture



CMU



# XRL: Explainable Reinforcement Learning



Carnegie Mellon University

## Explainable Model

### Explainable RL (XRL)

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

## Explanation Interface

### XRL Interaction

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

## Challenge Problem

### Autonomy

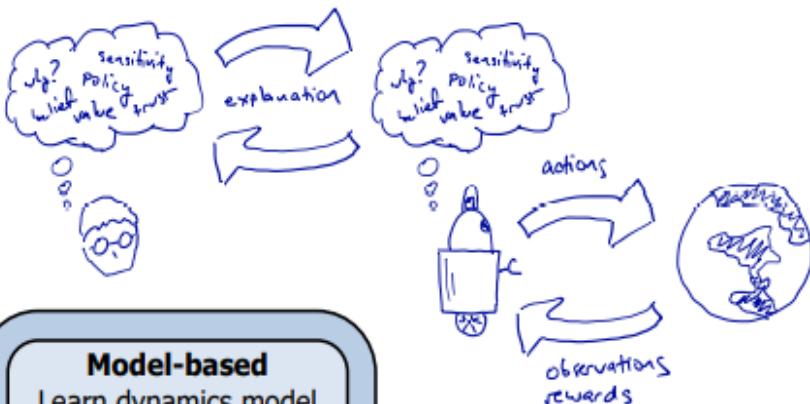
- Open AI Gym
- Autonomy in the electrical grid
- Mobile service robots
- Self-improving educational software

• PI: Zico Kolter (CMU)

• Geoff Gordon (CMU)  
• Pradeep Ravikumar (CMU)

## CMU

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

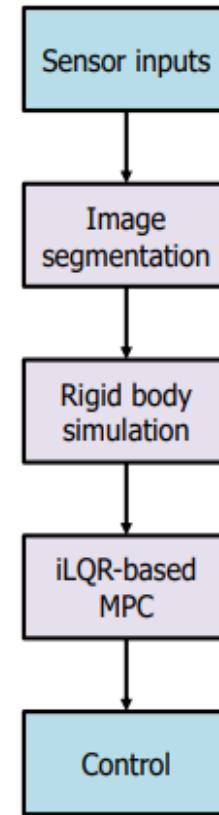


**Model-based**  
Learn dynamics model of environment, plan actions in model, and execute in real system

**Model-free**  
Directly learn value and/or policy for the environment

→ Improve model learning/ creation for RL agents to capture benefits of model-based approach

→ For any type of RL approach, provide an explanation of why an agent acted in a certain way



**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



## DARE: Deep Attention-based Representations for Explanation



### 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)

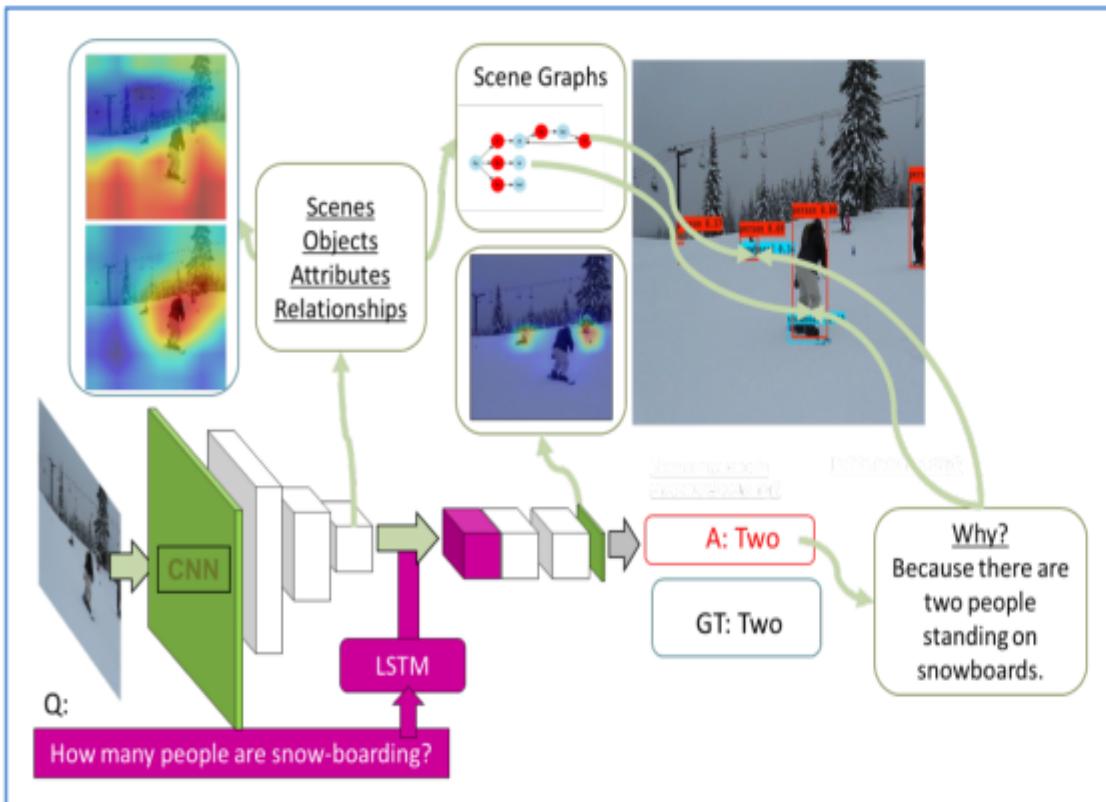
- Xiao Lin (SRI)
- Ryan Villamil (SRI)
- Dejan Jovanovic (SRI)
- Avi Ziskind (SRI)
- Michael Wessel (SRI)

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

• Jürgen Schulze (UCSD)

## SRI International

### 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



## EQUAS: Explainable QUestion Answering System



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

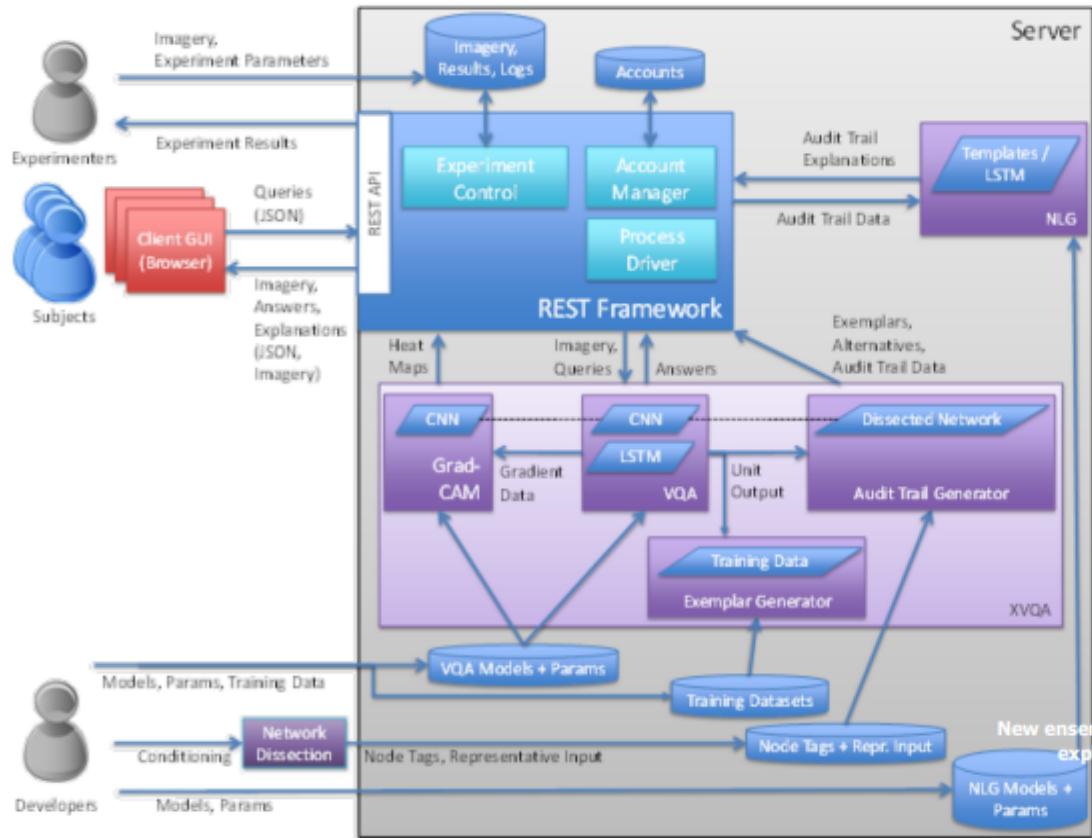
- VQA for images and video

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

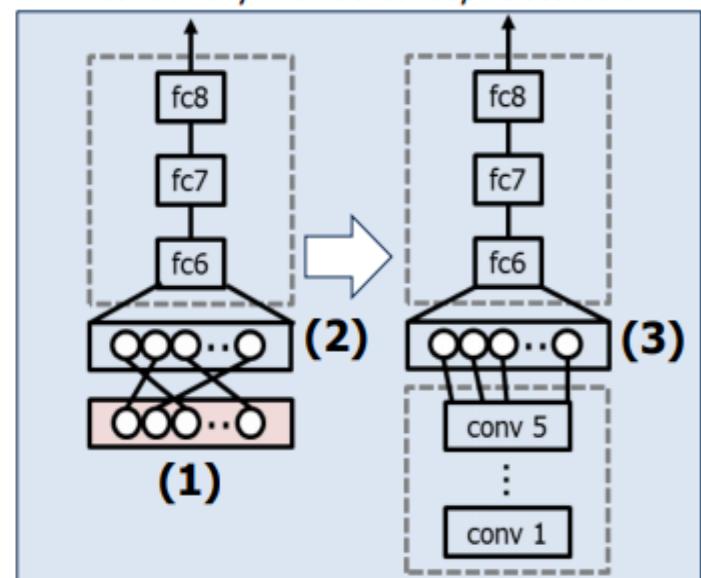
• Antonio Torralba (MIT)  
• Ray Mooney (UT Austin)

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

## Raytheon BBN



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



# Tractable Probabilistic Logic Models



UT Dallas, UCLA, Texas A&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

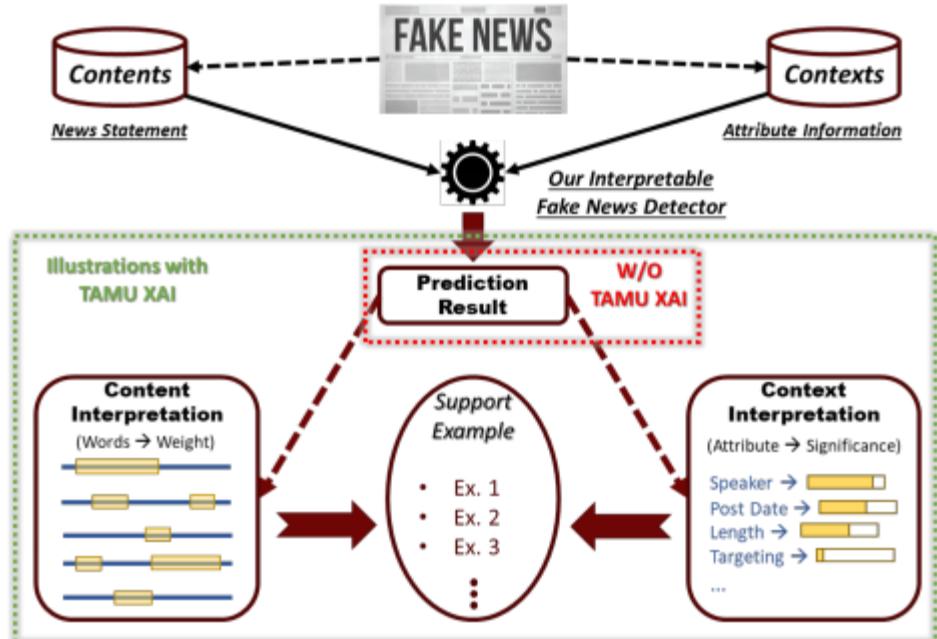
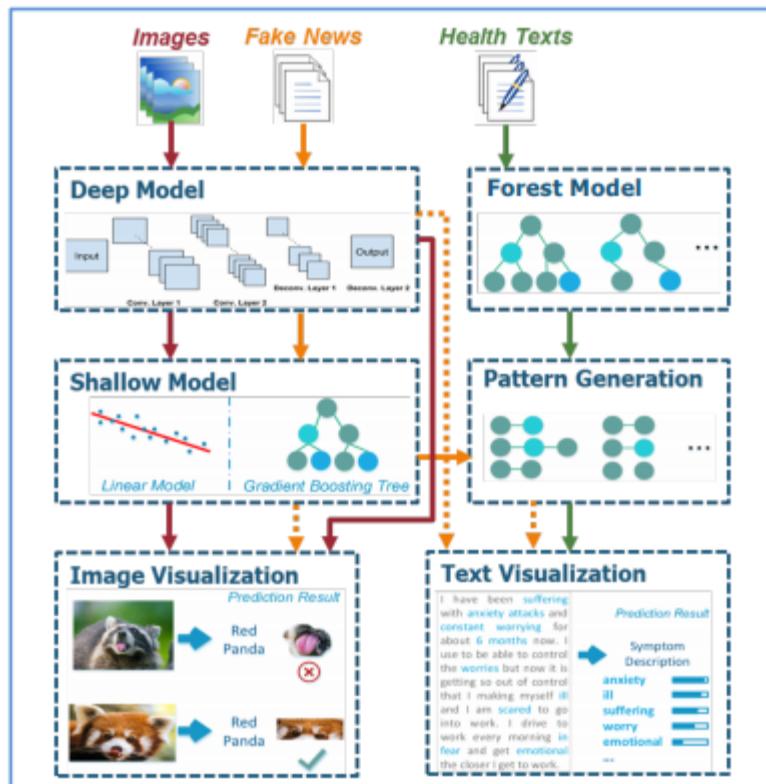
- Infer activities in multimodal data (video and text)
- Wetlab (biology) and TACoS (cooking) datasets

- 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

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



## Rutgers



## Model Explanation by Optimal Selection of Teaching Examples



### Rutgers University

#### 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
  - Full model
  - User-selected sub-structure
  - User submitted examples

#### Challenge Problem

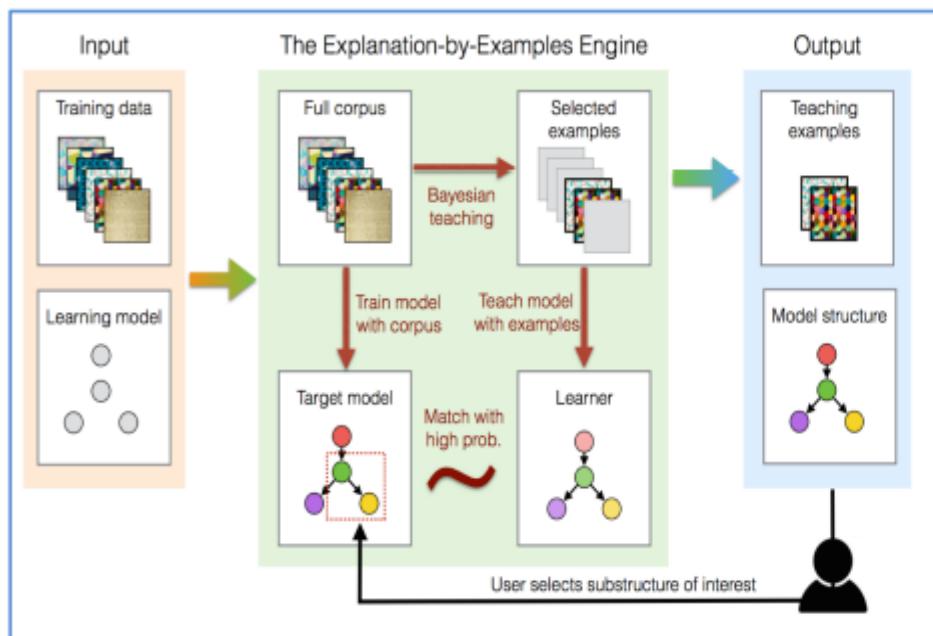
##### **Data Analytics**

- Image processing
- Text corpora
- VQA
- Movie events

- **PI:** Patrick Shafto (Rutgers)
- Scott Cheng-Hsin Yang (Rutgers)

## Rutgers

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

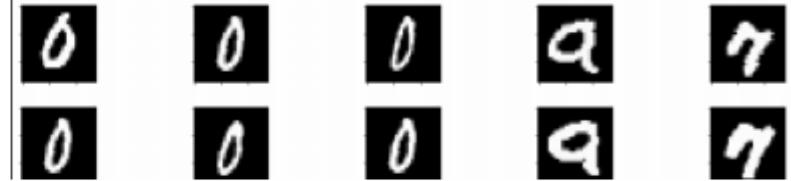


Good and bad examples for teaching a category  
(illustrates model strengths and weaknesses)

Good pairs of examples of the category 9



Bad pairs of examples of the category 9



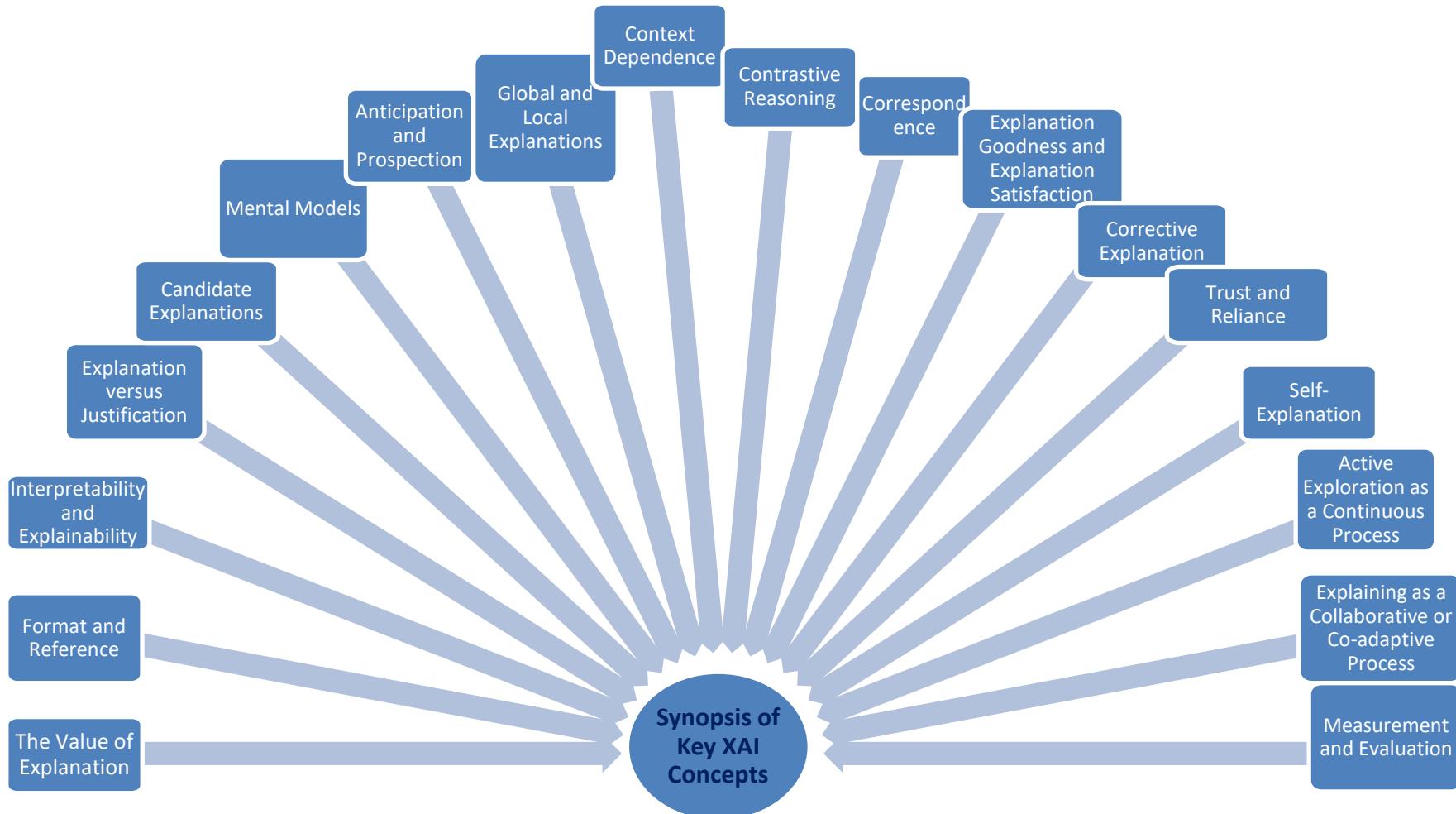
## Report from TA 2: Psychological Explanation Theories

### What Makes a Good Explanation (Literature Review)

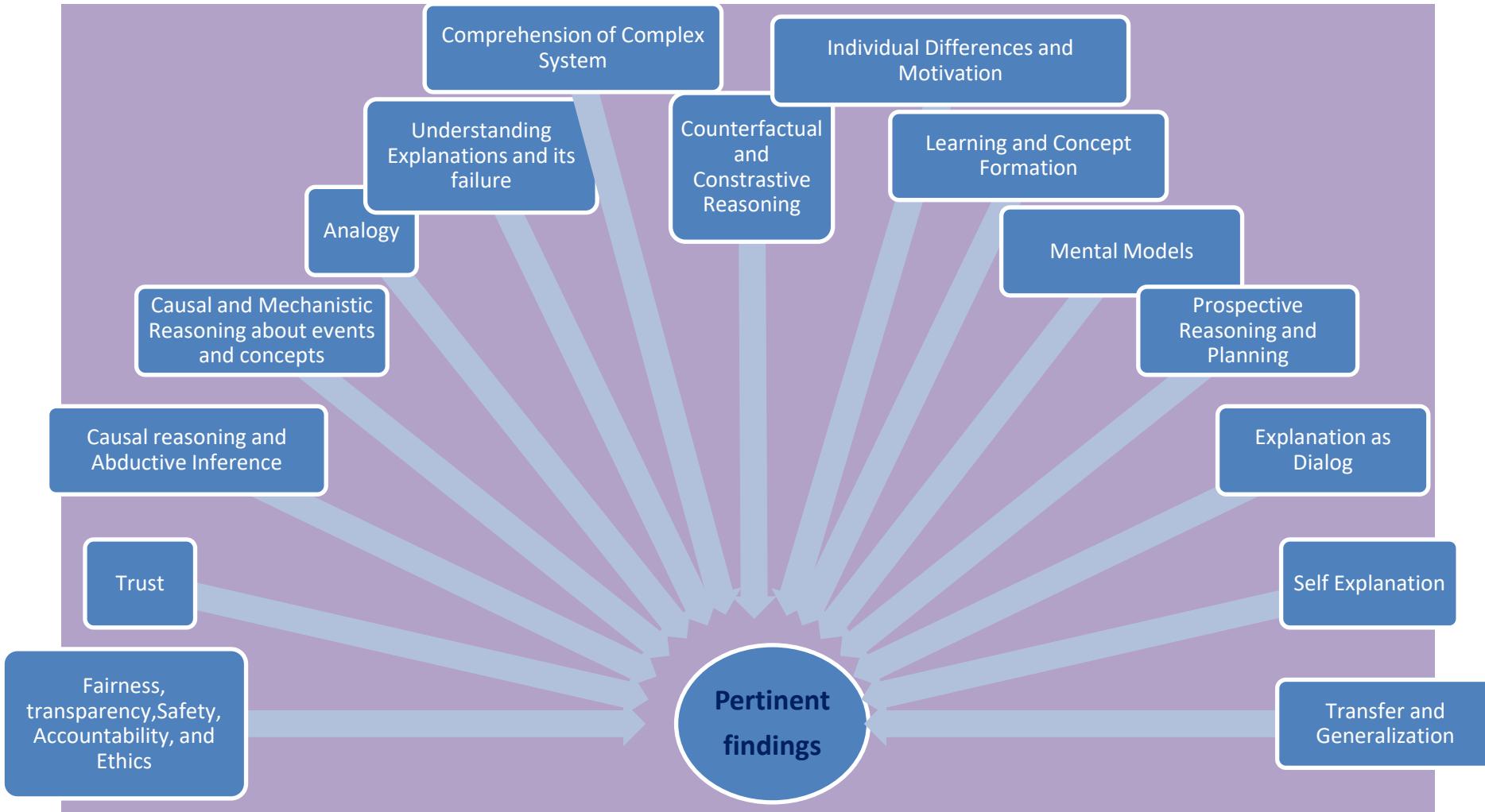
#### Recommendation

Insert detail of empirical/experimental methods as in psychological research:  
detail participants, instructions, procedures, tasks, DV, IV and Control

## Report from TA 2: Psychological Explanation Theories



## Report from TA 2: Psychological Explanation Theories



# Outline

- Interpretability = Transparency ≠ Explainability?
- DARPA XAI Grant program
- Results: Grantees and reports
- Key takeaways and discussions

- DARPA XAI "challenge" program has been working in a greater extent in explaining explainable AI compare to other initiatives e.g. EU AI guidelines
- Emphasis more on trusted machine rather than fairness → some regulation issues need more attention?
- Much left to do in psychological explanation areas in the interaction with non end-user (e.g. tackling the issue of fairness and inequality)

# discussions