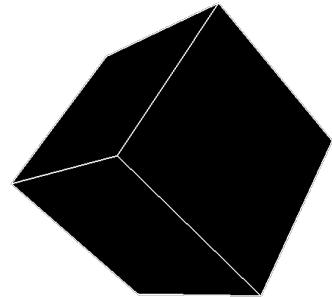


## Explainable AI – The Story So Far

August 29th, 2019 - @Inha University

Freddy Lecue  
Chief AI Scientist, CortAIx, Thales, Montreal – Canada  
Inria, Sophia Antipolis - France

@freddylecue  
<https://tinyurl.com/freddylecue>



# Motivation

# Business to Customer

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.



**Gary Chavez** added a photo you might ...  
be in.

about a minute ago ·



# Critical Systems





# Markets We Serve (Critical Systems)

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.



Aerospace



Space



Ground Transportation



Defence



Security

Trusted Partner For A Safer World

# But not Only Critical Systems

# Motivation (1)

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

## Criminal Justice

- People wrongly denied parole
- Recidivism prediction
- Unfair Police dispatch



GET UPDATES / DONATE



**STATEMENT OF CONCERN ABOUT PREDICTIVE  
POLICING BY ACLU AND 16 CIVIL RIGHTS PRIVACY,  
RACIAL JUSTICE, AND TECHNOLOGY  
ORGANIZATIONS**



[aclu.org/other/statement-concern-about-predictive-policing-aclu-and-16-civil-rights-privacy-racial-justice](http://aclu.org/other/statement-concern-about-predictive-policing-aclu-and-16-civil-rights-privacy-racial-justice)

Opinion

The New York Times

OP-ED CONTRIBUTOR

## When a Computer Program Keeps You in Jail

By Rebecca Wexler

June 13, 2017



[nytimes.com/2017/06/13/opinion/how-computers-are-harming-criminal-justice.html](http://nytimes.com/2017/06/13/opinion/how-computers-are-harming-criminal-justice.html)

## How We Analyzed the COMPAS Recidivism Algorithm

by Jeff Larson, Surya Mattu, Lauren Kirchner and Julia Angwin

May 23, 2016

[propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm](http://propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm)

# Motivation (2)

## Finance:

- Credit scoring, loan approval
- Insurance quotes

The Big Read Artificial intelligence [+ Add to myFT](#)

## Insurance: Robots learn the business of covering risk

Artificial intelligence could revolutionise the industry but may also allow clients to calculate if they need protection

Save

Oliver Ralph MAY 16, 2017

24

<https://www.ft.com/content/e07cee0c-3949-11e7-821a-6027b8a20f23>

FICO®  
COMMUNITY

Explainable Machine Learning Challenge

[community.fico.com/s/explainable-machine-learning-challenge](http://community.fico.com/s/explainable-machine-learning-challenge)



## Healthcare

- Applying ML methods in medical care is problematic.
- AI as 3<sup>rd</sup>-party actor in physician-patient relationship
- Responsibility, confidentiality?
- Learning must be done with available data.

Cannot randomize cares given to patients!

- Must validate models before use.

[Email](#) [Tweet](#)

Researchers say use of artificial intelligence in medicine raises ethical questions

In a perspective piece, Stanford researchers discuss the ethical implications of using machine-learning tools in making health care decisions for patients.

Patricia Hannon ,<https://med.stanford.edu/news/all-news/2018/03/researchers-say-use-of-ai-in-medicine-raises-ethical-questions.html>

## Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission

Rich Caruana  
Microsoft Research  
[rcaruana@microsoft.com](mailto:rcaruana@microsoft.com)

Paul Koch  
Microsoft Research  
[paulkoch@microsoft.com](mailto:paulkoch@microsoft.com)

Yin Lou  
LinkedIn Corporation  
[y lou@linkedin.com](mailto:y lou@linkedin.com)

Marc Sturm  
NewYork-Presbyterian Hospital  
[mas9161@nyp.org](mailto:mas9161@nyp.org)

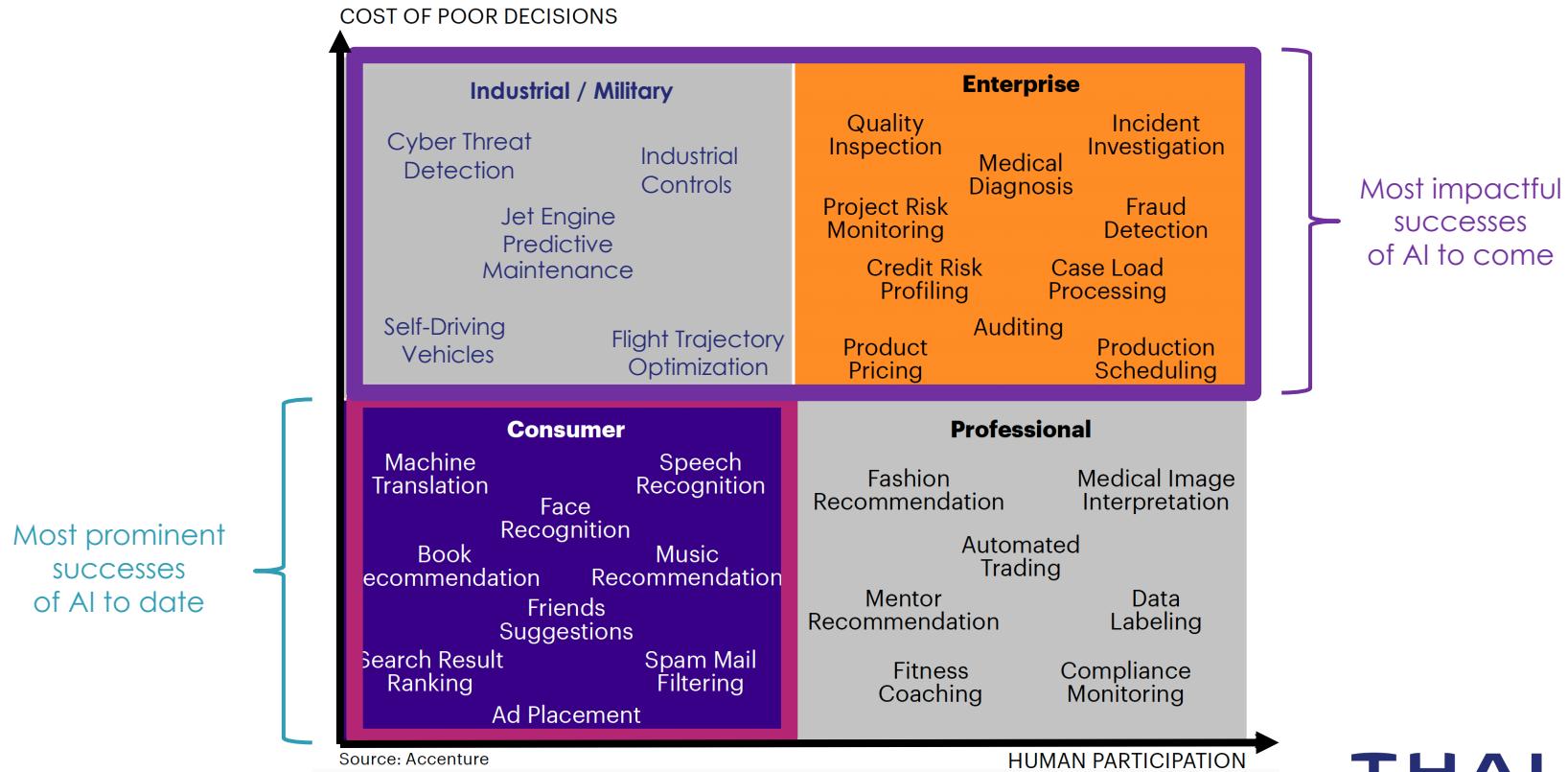
Johannes Gehrke  
Microsoft  
[johannes@microsoft.com](mailto:johannes@microsoft.com)

Noémie Elhadad  
Columbia University  
[noemie.elhadad@columbia.edu](mailto:noemie.elhadad@columbia.edu)

Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, Noémie Elhadad: Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission. KDD 2015: 1721-1730

# Trustable AI and eXplainable AI: a Reality Need

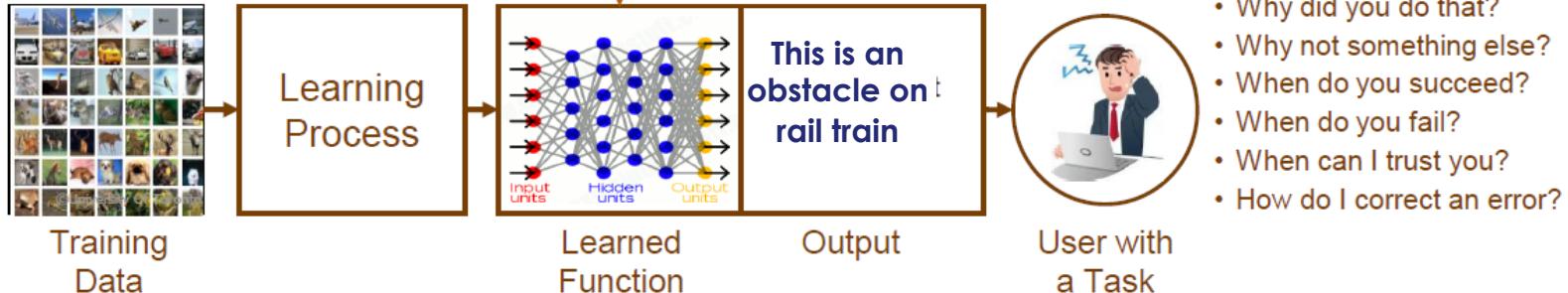
## The need for explainable AI rises with the potential cost of poor decisions



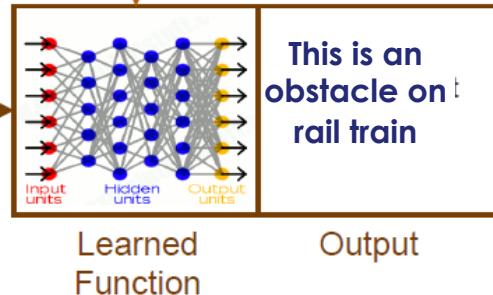
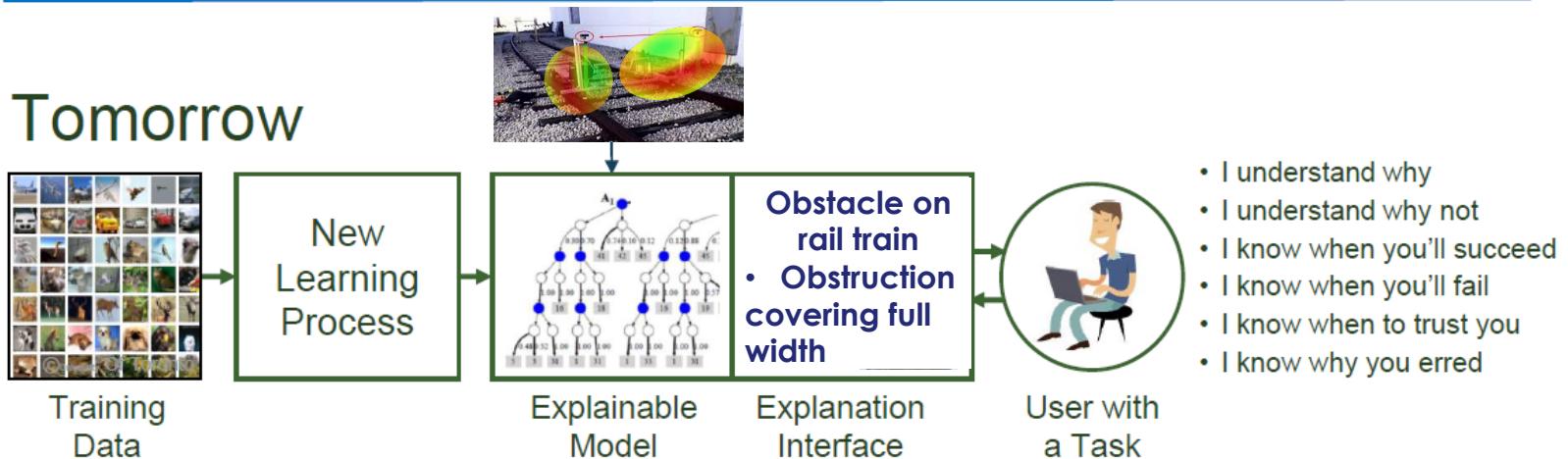
# XAI in a Nutshell

# XAI in a Nutshell

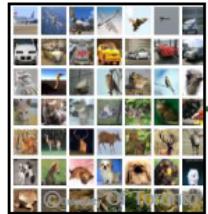
## Today



## Tomorrow



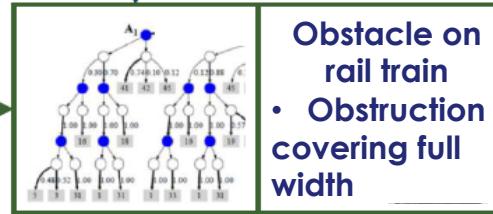
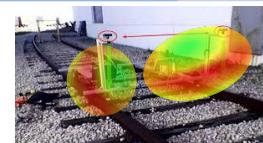
User with  
a Task



Training  
Data



New  
Learning  
Process



Explainable  
Model

Obstacle on  
rail train  
• Obstruction  
covering full  
width



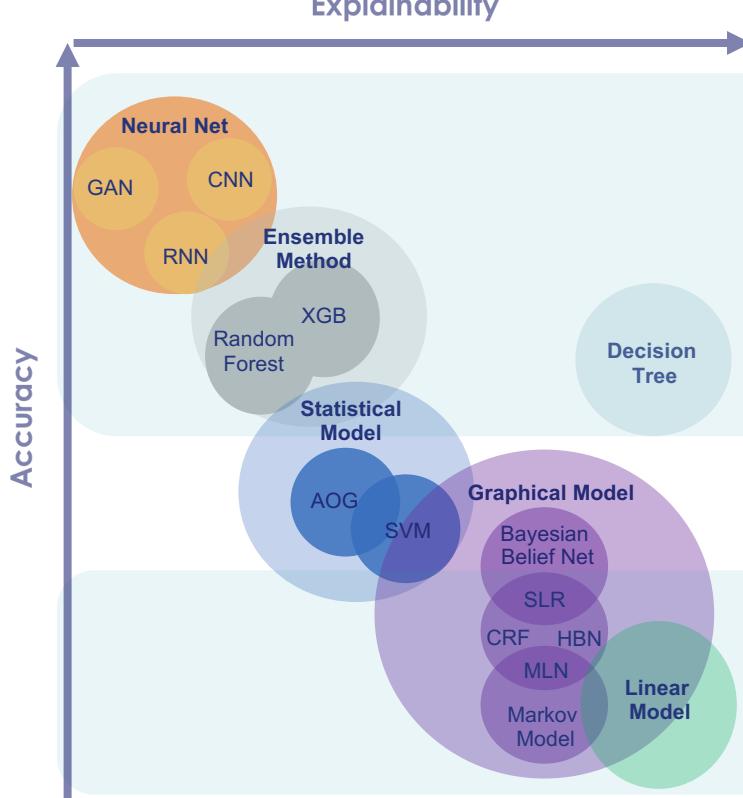
User with  
a Task

# How to Explain? Accuracy vs. Explanability

## Learning

- Challenges:
  - Supervised
  - Unsupervised learning
- Approach:
  - Representation Learning
  - Stochastic selection
- Output:
  - **Correlation**
  - **No causation**

## Explainability



## Interpretability

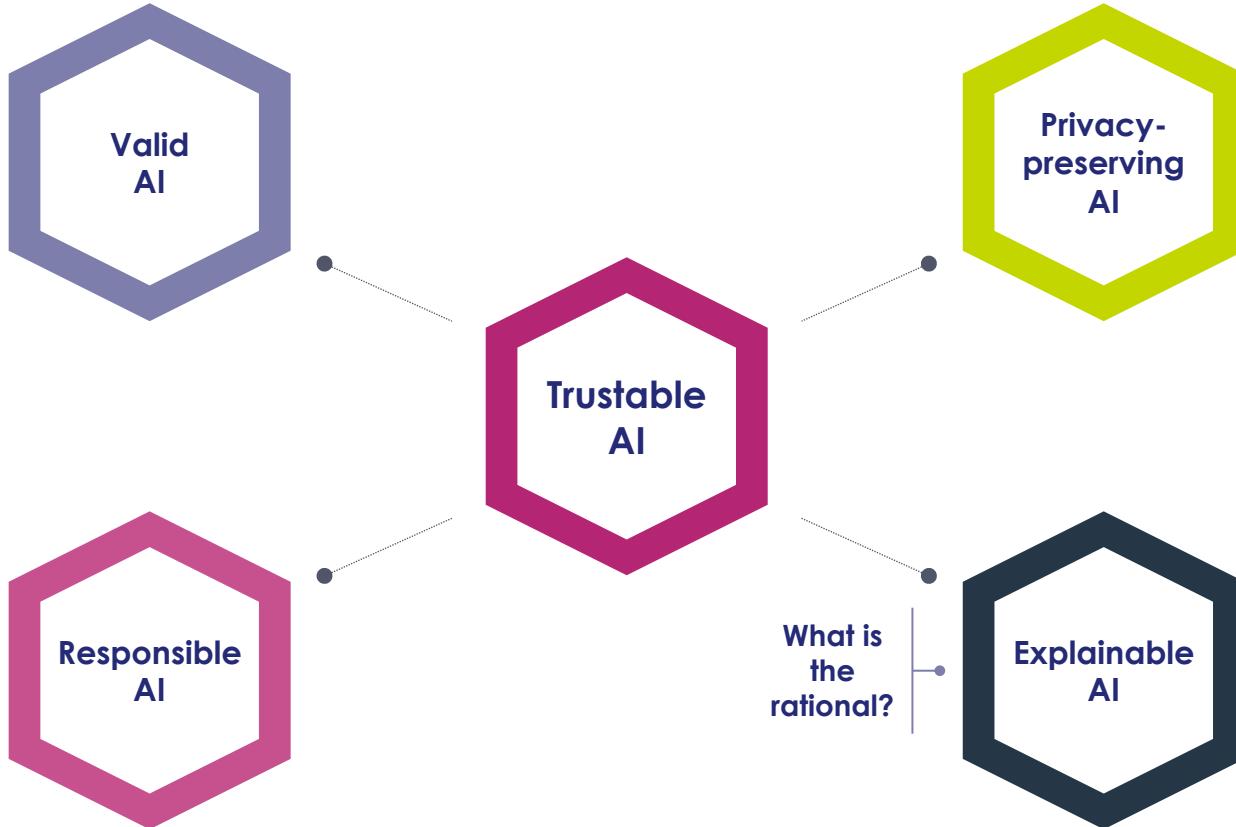
Non-Linear functions

Polynomial functions

Quasi-Linear functions

# Trustable AI

# AI Adoption: Requirements



- Human Interpretable AI
- Machine Interpretable AI

# Definitions

# Explanation in AI

**Explanation in AI aims to create a suite of techniques that produce more explainable models, while maintaining a high level of searching, learning, planning, reasoning performance: optimization, accuracy, precision; and enable human users to understand, appropriately trust, and effectively manage the emerging generation of AI systems.**

## explanation | ɛksplə'neɪʃ(ə)n |

noun

a statement or account that makes something clear: *the birth rate is central to any explanation of population trends.*

Models, Outputs of the Intelligent System

## interpret | ɪn'tə:pɪt |

verb (**interprets, interpreting, interpreted**) [with object]

1 explain the meaning of (information or actions): *the evidence is difficult to interpret.*

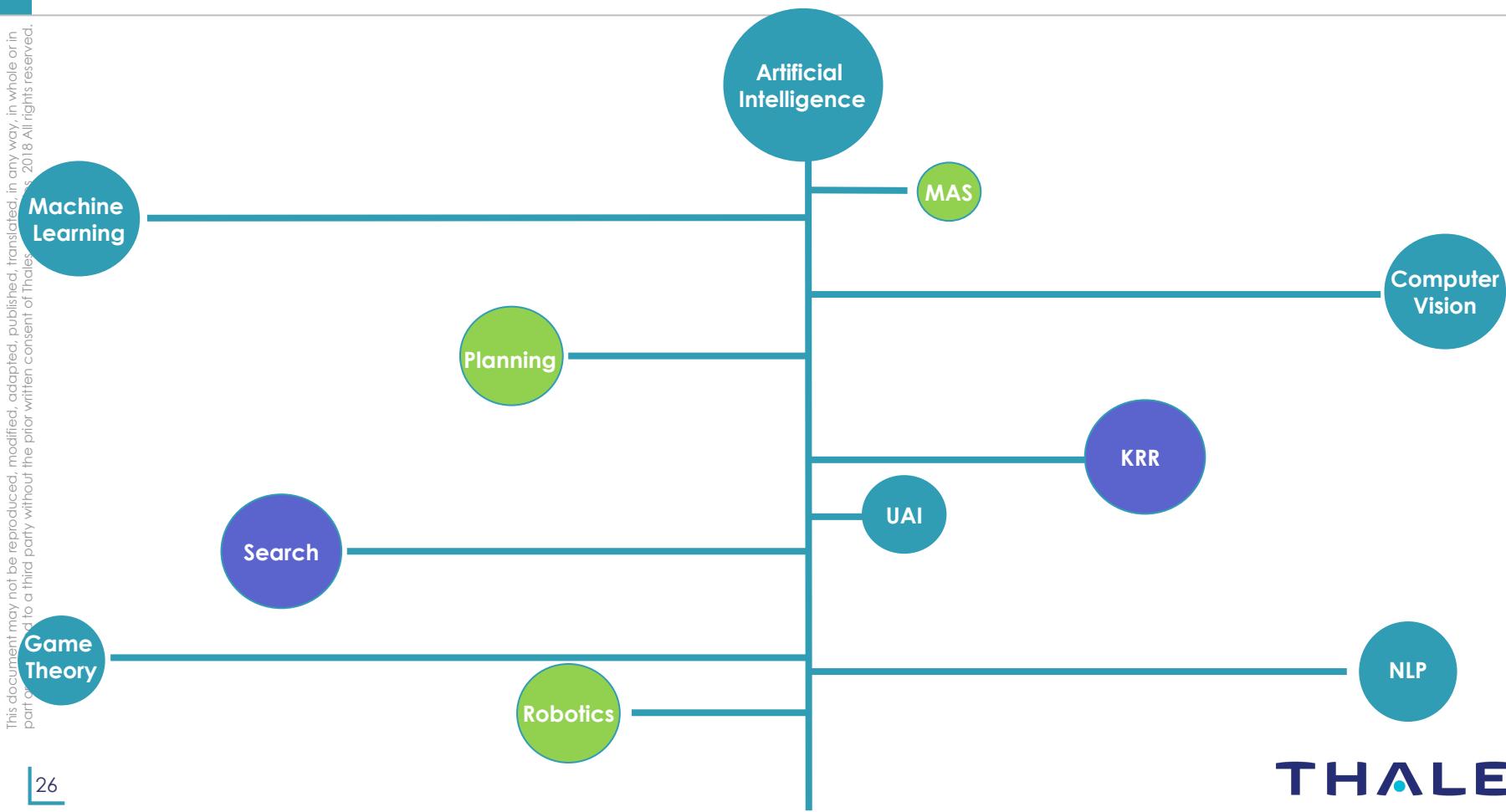
Models, Outputs of the Intelligent System

# XAI in AI

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part, or given to a third party without the prior written consent of Thales. © Thales Group 2018. All rights reserved.

26



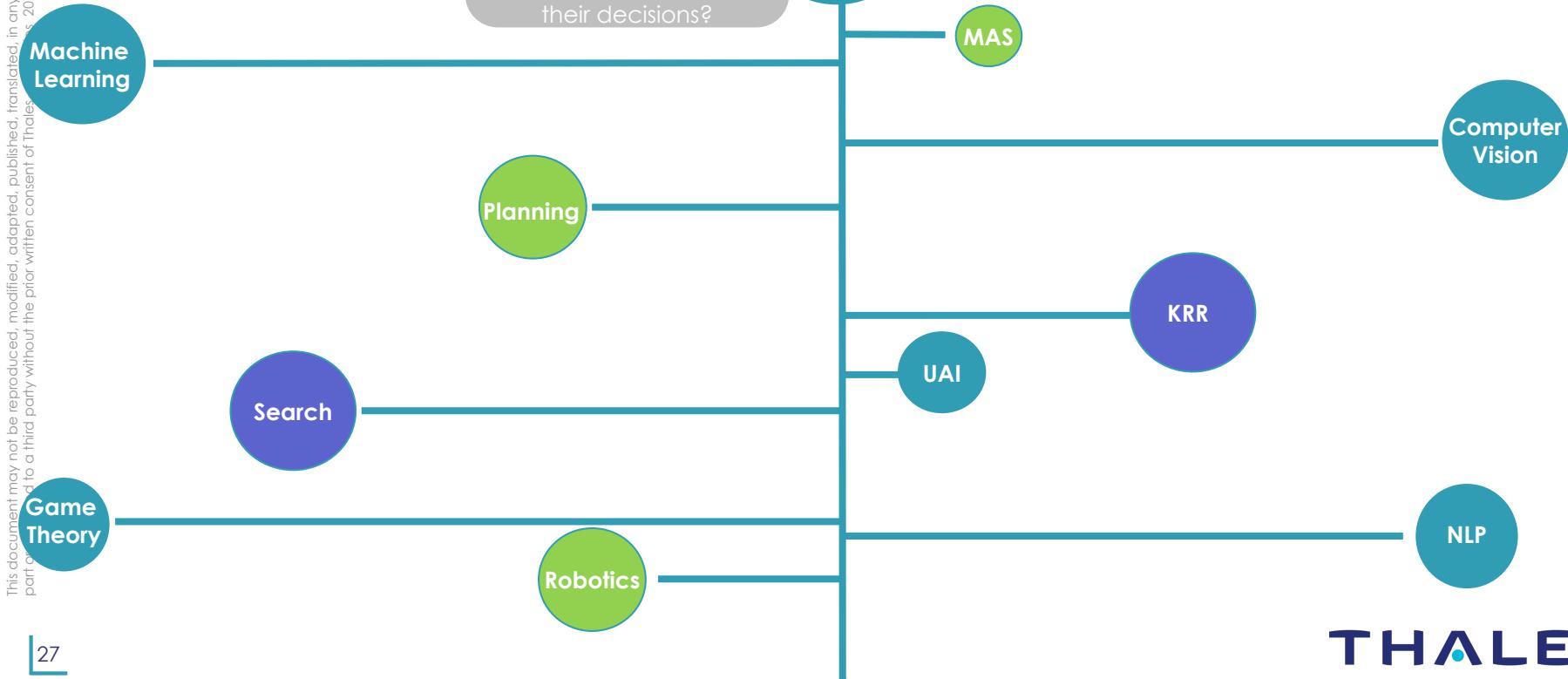
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part, or given to a third party without the prior written consent of Thales. © Thales Group 2018. All rights reserved.

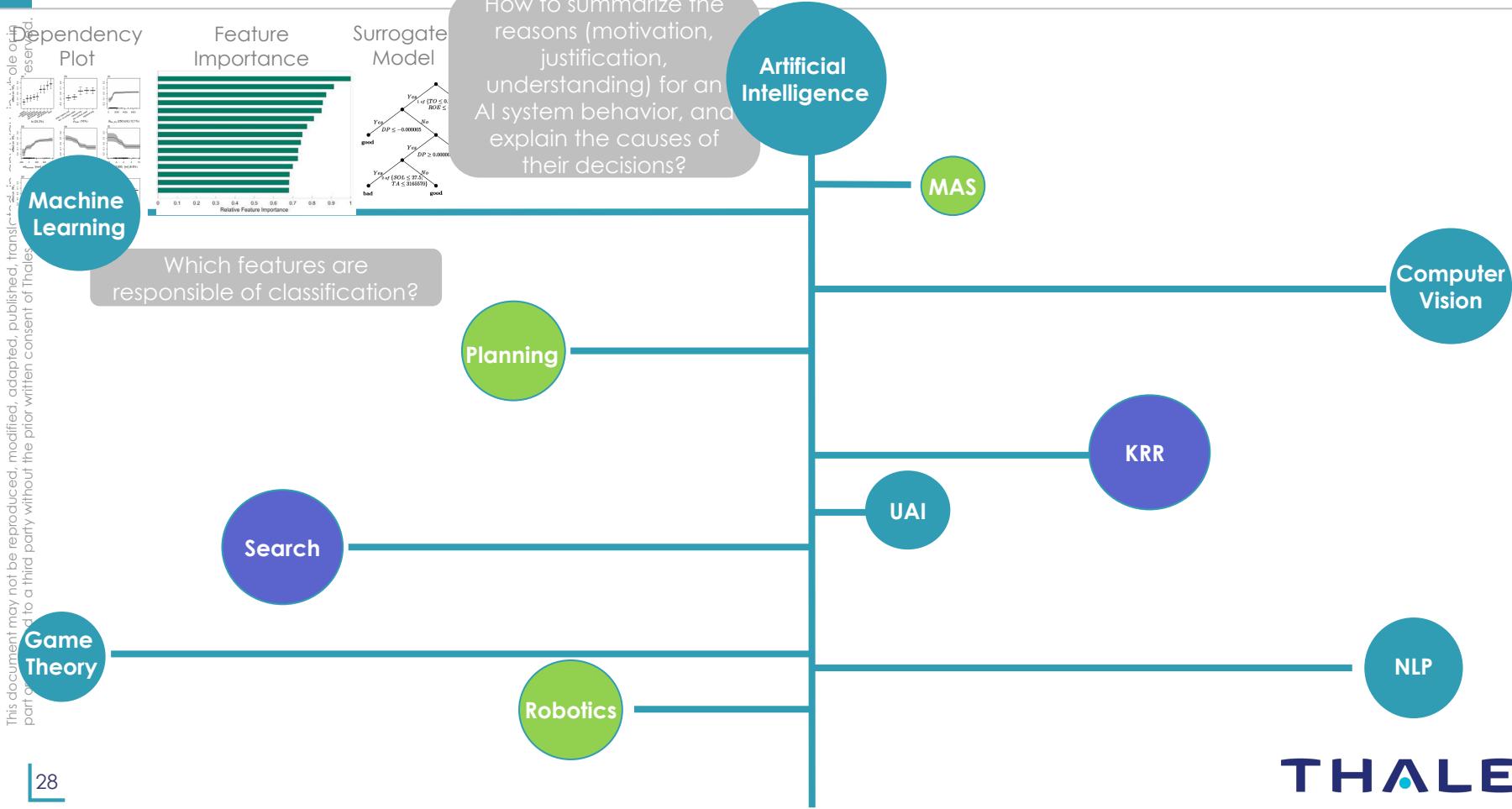
27

27

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

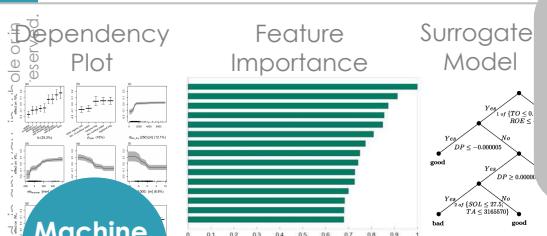


# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches



# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map



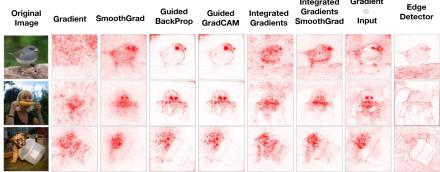
Machine Learning

Which features are responsible of classification?

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

Artificial Intelligence

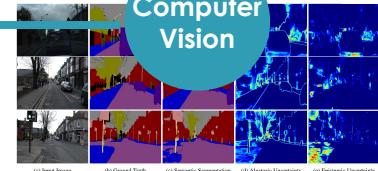
MAS



Which complex features are responsible of classification?

Planning

Computer Vision



Uncertainty Map

Search

KRR

UAI

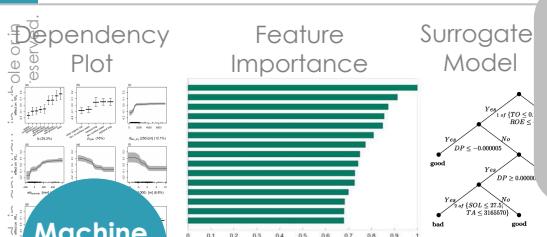
Game Theory

Robotics

NLP

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map



Machine Learning

Which features are responsible of classification?

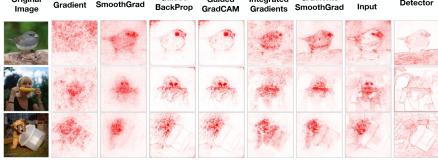
How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

Artificial Intelligence

Strategy  
Summarization

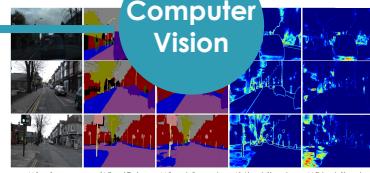
MAS

- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?



Which complex features are responsible of classification?

Computer Vision



Uncertainty Map

Planning

KRR

UAI

Search

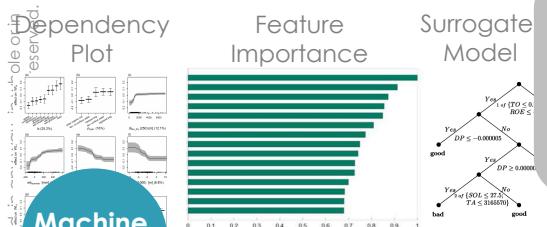
Game Theory

Robotics

NLP

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map



Machine Learning

Which features are responsible of classification?

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

Artificial Intelligence

Strategy  
Summarization

MAS

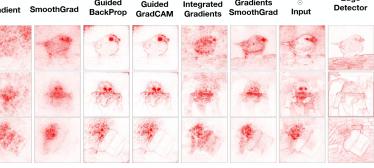
Plan Refinement

Planning

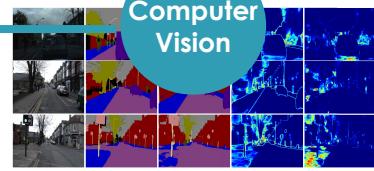
Which actions are responsible of a plan?

- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?

Computer Vision



Which complex features are responsible of classification?



Uncertainty Map

This document may not be reproduced, modified, adapted, published, transmitted or loaned to a third party without the prior written consent of Thales

Search

UAI

KRR

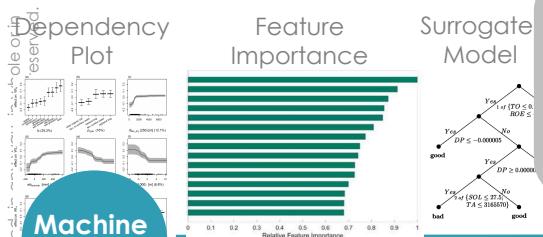
NLP

Game Theory

Robotics

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map



Machine Learning

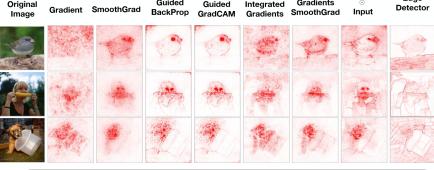
Which features are responsible of classification?

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

Artificial Intelligence

Strategy  
Summarization

MAS



Which complex features are responsible of classification?

- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?

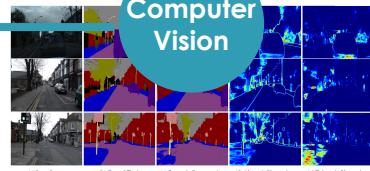
Plan Refinement

Planning

Which actions are responsible of a plan?

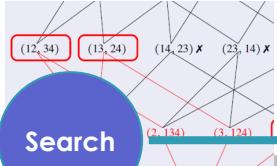
KRR

Computer Vision



Uncertainty Map

Conflicts Resolution



Search

Game Theory

Which constraints can be relaxed?

UAI

Robotics

NLP

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

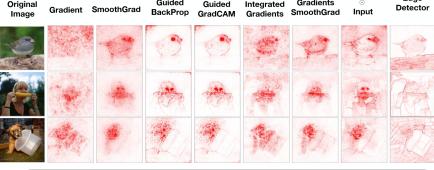


Machine Learning

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

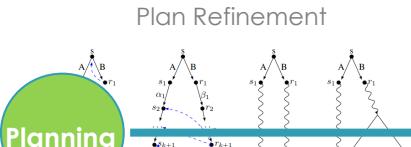
Artificial Intelligence

Strategy  
Summarization



Which complex features are responsible of classification?

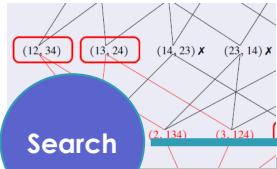
Which features are responsible of classification?



Planning

Which actions are responsible of a plan?

Conflicts Resolution



Search

Which constraints can be relaxed?

Game Theory

Which combination of features is optimal?



Shapley Values

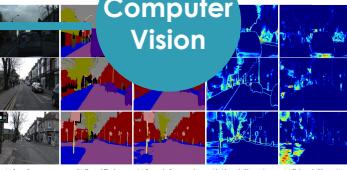
MAS

- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?

KRR

UAI

NLP



Uncertainty Map

THALES

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

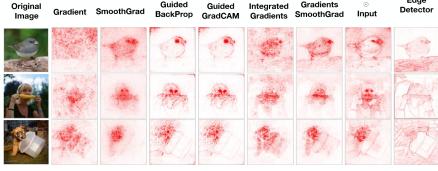


Machine Learning

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

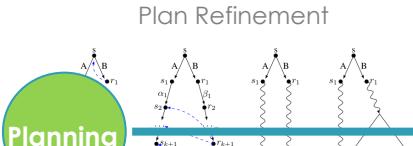
Artificial Intelligence

Strategy  
Summarization



Which complex features are responsible of classification?

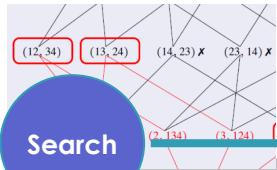
Which features are responsible of classification?



Planning

Which actions are responsible of a plan?

Conflicts Resolution



Search

Which constraints can be relaxed?

Game Theory

Which combination of features is optimal?



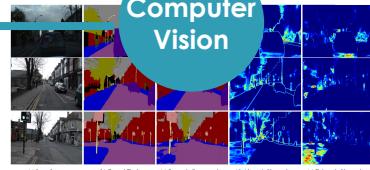
Shapley Values

MAS

- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?

KRR

Computer Vision



Uncertainty Map

UAI

Which decisions, combination of multimodal decisions lead to an action?

Robotics

Narrative-based

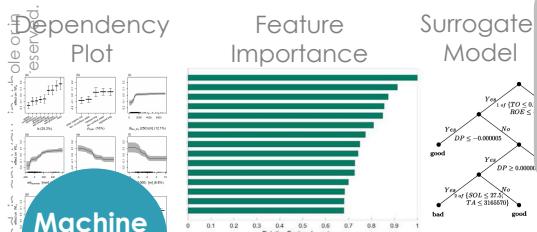


NLP

THALES

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

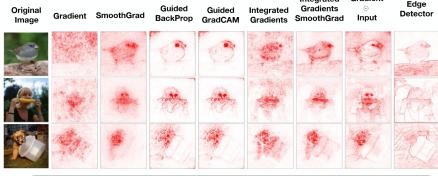


Machine Learning

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

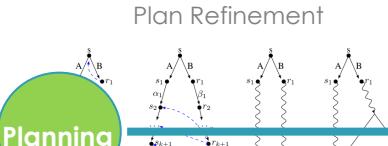
Artificial Intelligence

Strategy  
Summarization



Which complex features are responsible of classification?

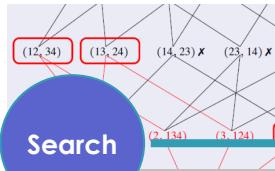
Which features are responsible of classification?



Planning

Which actions are responsible of a plan?

Conflicts Resolution



Search

Which constraints can be relaxed?

Game Theory

This document may not be reproduced, modified, adapted, published, translated or distributed to a third party without the prior written consent of Thales

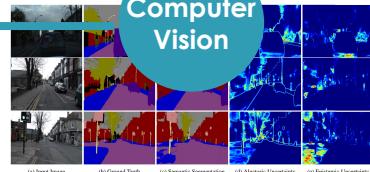
Which combination of features is optimal?



MAS

- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?

KRR



Uncertainty Map

Robotics

Narrative-based

Which decisions, combination of multimodal decisions lead to an action?



Machine Learning based

Algorithm 2	
	Predicted:
Posting	Atheism
Hiring	✓
Religious	
Secular	
Non	

Document  
From: paul@verisai.com (Paul Durbin)  
Subject: Re: DAVID CORBIN IS GOD!  
To: Net-Posting-Host: sepp2@verisai.com  
Organization: Verisai Corp  
Lines: 8

NLP

Which entity is responsible for classification?

THALES

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

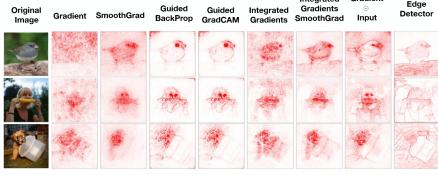


Machine Learning

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

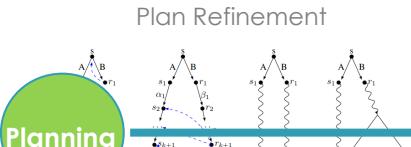
Artificial Intelligence

Strategy  
Summarization



Which complex features are responsible of classification?

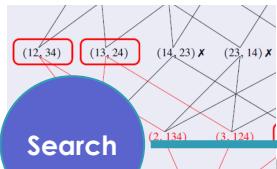
Which features are responsible of classification?



Planning

Which actions are responsible of a plan?

Conflicts Resolution



Search

Which constraints can be relaxed?

Game Theory

Which combination of features is optimal?



Robotics

Which decisions, combination of multimodal decisions lead to an action?

Narrative-based



UAI

Diagnosis

KRR

Abduction

Uncertainty Map

- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?



Computer Vision

- Which axiom is responsible of inference (e.g., classification)?
- Abduction/Diagnostic: Find the **right** root causes (abduction)?

Machine Learning based

Algorithm 2	
Posting	Predicted: Athlete
Hiring	Prediction correct: ✓
Residence	
Employment	
Non	

Document

From: paul@verisai.com (Paul Durbin)  
Subject: Re: DAVID CORBIN IS GOD!  
To: Paul@Verisai.com  
Organization: Verisai Corp

NLP

Which entity is responsible for classification?

THALES

# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

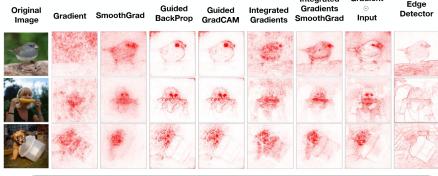


Machine Learning

How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

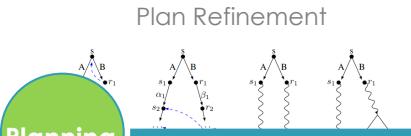
Artificial Intelligence

Strategy  
Summarization



Which complex features are responsible of classification?

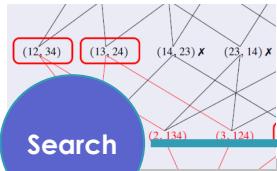
Which features are responsible of classification?



Planning

Which actions are responsible of a plan?

Conflicts Resolution



Search

Which constraints can be relaxed?

Game Theory

Which combination of features is optimal?



Robotics

Which decisions, combination of multimodal decisions lead to an action?

Narrative-based



UAI

Uncertainty as an alternative to explanation

Machine Learning based

Algorithm 2	
Posting	Predicted: Athlete
Hiring	Prediction correct: ✓
Rel	
Int	
Non	

Document

From: paul@verisai.com (Paul Durbin)  
Subject: Re: DAVID CORBIN IS GOD!  
To: Ntp-Posting-Host: sepi2.verity.com  
Organization: Verisai Corp  
Lines: 8

NLP

Which entity is responsible for classification?

THALES

Diagnosis



Abduction

Uncertainty Map

KRR

- Which axiom is responsible of inference (e.g., classification)?
- Abduction/Diagnostic: Find the **right** root causes (abduction)?

# Deep Dive

# Overview of explanation in different AI fields (1)

## Machine Learning (except Artificial Neural Network)

### Interpretable Models:

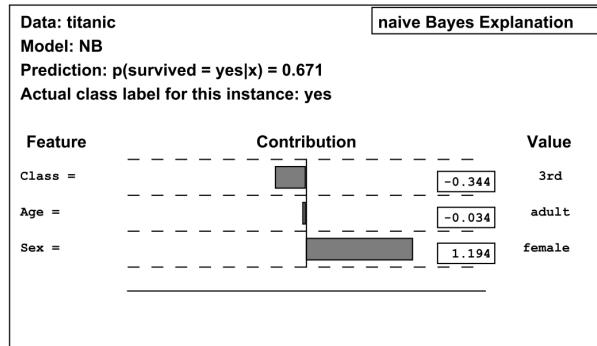
- Linear regression,
- Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs

# Overview of explanation in different AI fields (1)

## Machine Learning (except Artificial Neural Network)

### Interpretable Models:

- Linear regression,
- Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs



### Naive Bayes model

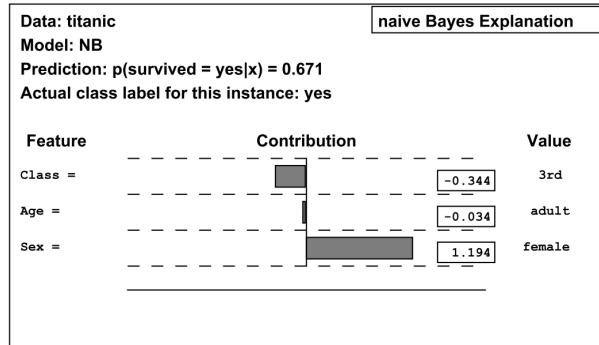
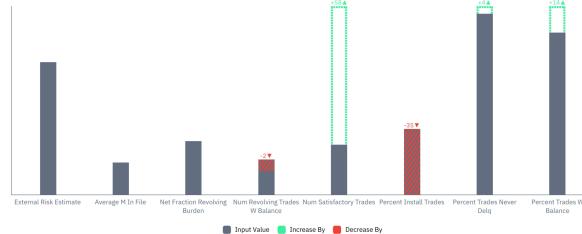
Igor Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. Artificial Intelligence in Medicine, 23:89–109, 2001.

# Overview of explanation in different AI fields (1)

## Machine Learning (except Artificial Neural Network)

### Interpretable Models:

- Linear regression,
- Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs



### Naive Bayes model

Igor Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. Artificial Intelligence in Medicine, 23:89–109, 2001.

### Counterfactual What-if

Brent D. Mittelstadt, Chris Russell, Sandra Wachter:  
Explaining Explanations in AI. FAT 2019: 279–288

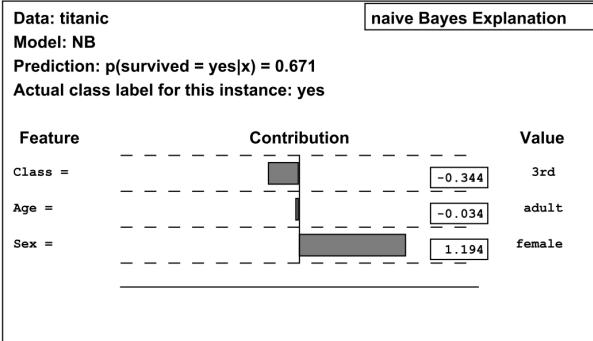
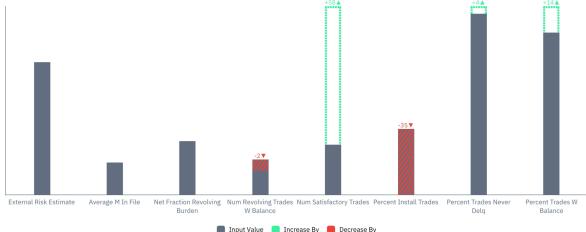
Rory McGrath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. CoRR abs/1811.05245 (2018)

# Overview of explanation in different AI fields (1)

## Machine Learning (except Artificial Neural Network)

### Interpretable Models:

- Linear regression,
- Logistic regression,
- Decision Tree,
- GLMs,
- GAMs
- KNNs



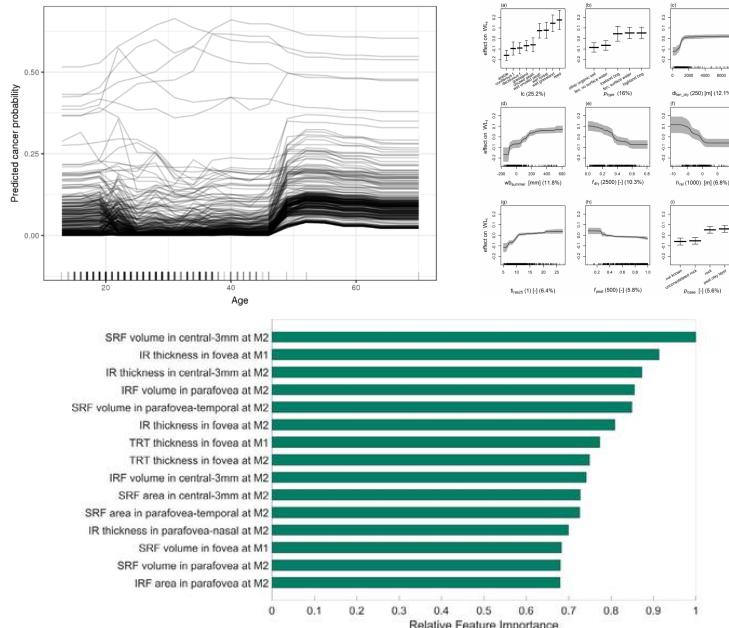
### Naive Bayes model

Igor Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. Artificial Intelligence in Medicine, 23:89–109, 2001.

### Counterfactual What-if

Brent D. Mittelstadt, Chris Russell, Sandra Wachter: Explaining Explanations in AI. FAT 2019: 279–288

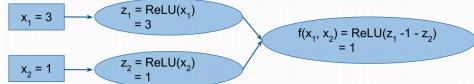
Rory McGrath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. CoRR abs/1811.05245 (2018)



Feature Importance  
Partial Dependence Plot  
Individual Conditional Expectation  
Sensitivity Analysis

# Overview of explanation in different AI fields (2)

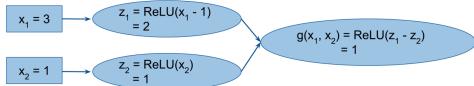
## Machine Learning (only Artificial Neural Network)



Network  $f(x_1, x_2)$

Attributions at  $x_1 = 3, x_2 = 1$

<b>Integrated gradients</b>	$x_1 = 1.5, x_2 = -0.5$
DeepLift	$x_1 = 1.5, x_2 = -0.5$
LRP	$x_1 = 1.5, x_2 = -0.5$



Network  $g(x_1, x_2)$

Attributions at  $x_1 = 3, x_2 = 1$

<b>Integrated gradients</b>	$x_1 = 1.5, x_2 = -0.5$
DeepLift	$x_1 = 2, x_2 = -1$
LRP	$x_1 = 2, x_2 = -1$

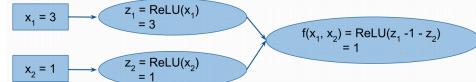
## Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319–3328, 2017.

Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation Differences. ICML 2017: 3145-3153

# Overview of explanation in different AI fields (2)

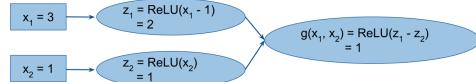
## | Machine Learning (only Artificial Neural Network)



Network  $f(x_1, x_2)$

Attributions at  $x_1 = 3, x_2 = 1$

**Integrated gradients**  $x_1 = 1.5, x_2 = -0.5$   
DeepLift  $x_1 = 1.5, x_2 = -0.5$   
LRP  $x_1 = 1.5, x_2 = -0.5$



Network  $g(x_1, x_2)$

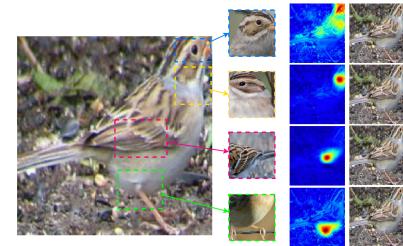
Attributions at  $x_1 = 3, x_2 = 1$

**Integrated gradients**  $x_1 = 1.5, x_2 = -0.5$   
DeepLift  $x_1 = 2, x_2 = -1$   
LRP  $x_1 = 2, x_2 = -1$

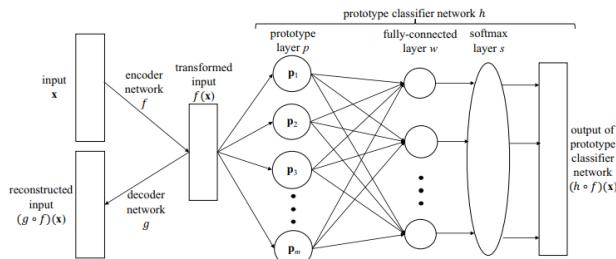
### Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319–3328, 2017.

Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation Differences. ICML 2017: 3145-3153



Chaofan Chen, Oscar Li, Alina Barnett, Jonathan Su, Cynthia Rudin: This looks like that: deep learning for interpretable image recognition. CoRR abs/1806.10574 (2018)

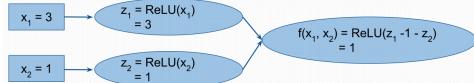


### Auto-encoder / Prototype

Oscar Li, Hao Liu, Chaofan Chen, Cynthia Rudin: Deep Learning for Case-Based Reasoning Through Prototypes: A Neural Network That Explains Its Predictions. AAAI 2018: 3530-3537

# Overview of explanation in different AI fields (2)

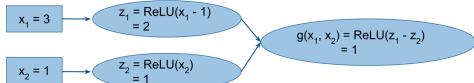
## Machine Learning (only Artificial Neural Network)



Network  $f(x_1, x_2)$

Attributions at  $x_1 = 3, x_2 = 1$

**Integrated gradients**  $x_1 = 1.5, x_2 = -0.5$   
DeepLift  $x_1 = 1.5, x_2 = -0.5$   
LRP  $x_1 = 1.5, x_2 = -0.5$



Network  $g(x_1, x_2)$

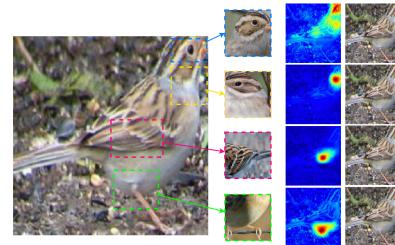
Attributions at  $x_1 = 3, x_2 = 1$

**Integrated gradients**  $x_1 = 1.5, x_2 = -0.5$   
DeepLift  $x_1 = 2, x_2 = -1$   
LRP  $x_1 = 2, x_2 = -1$

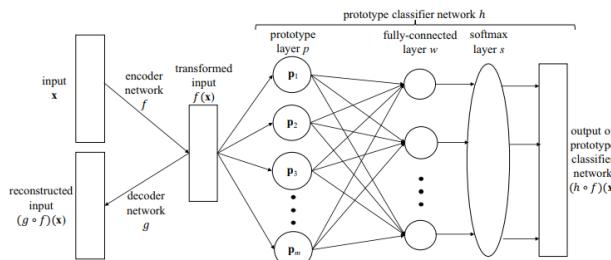
### Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319–3328, 2017.

Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation Differences. ICML 2017: 3145-3153

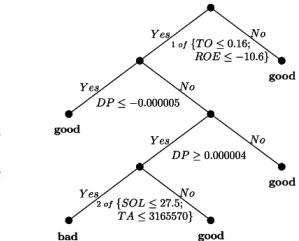
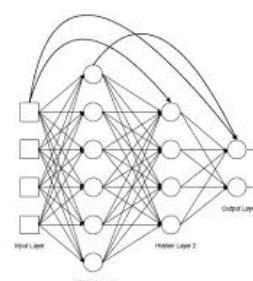


Chaofan Chen, Oscar Li, Alina Barnett, Jonathan Su, Cynthia Rudin: This looks like that: deep learning for interpretable image recognition. CoRR abs/1806.10574 (2018)



### Auto-encoder / Prototype

Oscar Li, Hao Liu, Chaofan Chen, Cynthia Rudin: Deep Learning for Case-Based Reasoning Through Prototypes: A Neural Network That Explains Its Predictions. AAAI 2018: 3530-3537



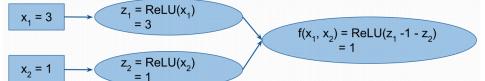
### Surrogate Model

Mark Craven, Jude W. Shavlik: Extracting Tree-Structured Representations of Trained Networks. NIPS 1995: 24-30

# Overview of explanation in different AI fields (2)

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

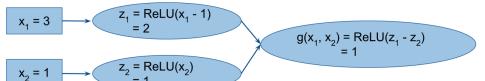
## Machine Learning (only Artificial Neural Network)



Network  $f(x_1, x_2)$

Attributions at  $x_1 = 3, x_2 = 1$

**Integrated gradients**  $x_1 = 1.5, x_2 = -0.5$   
DeepLift  $x_1 = 1.5, x_2 = -0.5$   
LRP  $x_1 = 1.5, x_2 = -0.5$



Network  $g(x_1, x_2)$

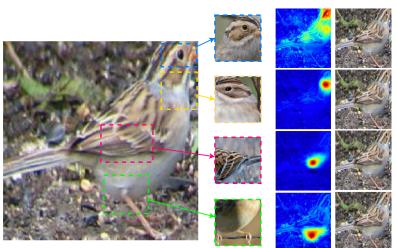
Attributions at  $x_1 = 3, x_2 = 1$

**Integrated gradients**  $x_1 = 1.5, x_2 = -0.5$   
DeepLift  $x_1 = 2, x_2 = -1$   
LRP  $x_1 = 2, x_2 = -1$

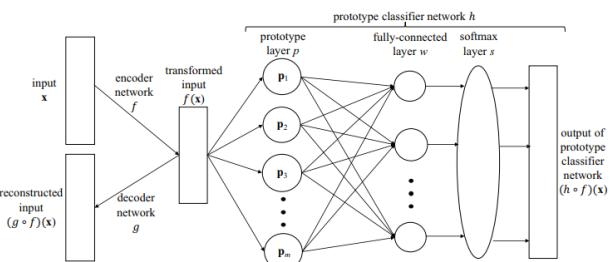
## Attribution for Deep Network (Integrated gradient-based)

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In ICML, pp. 3319–3328, 2017.

Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation Differences. ICML 2017: 3145-3153

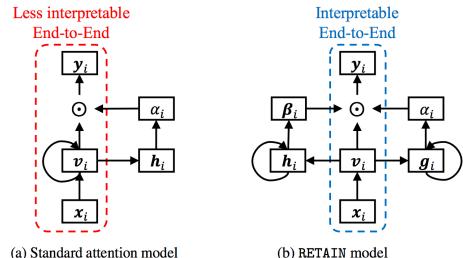


Chaofan Chen, Oscar Li, Alina Barnett, Jonathan Su, Cynthia Rudin: This looks like that: deep learning for interpretable image recognition. CoRR abs/1806.10574 (2018)



## Auto-encoder / Prototype

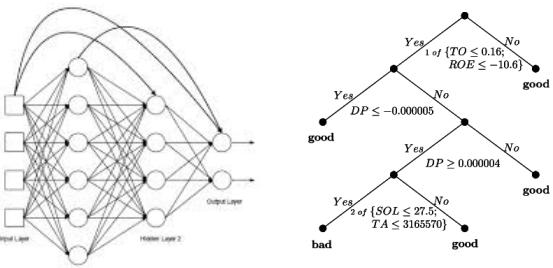
Oscar Li, Hao Liu, Chaofan Chen, Cynthia Rudin: Deep Learning for Case-Based Reasoning Through Prototypes: A Neural Network That Explains Its Predictions. AAAI 2018: 3530-3537



## Attention Mechanism

Edward Choi, Mohammad Taha Bahadori, Jimeng Sun, Joshua Kulas, Andy Schuetz, Walter F. Stewart: RETAIN: An Interpretable Predictive Model for Healthcare using Reverse Time Attention Mechanism. NIPS 2016: 3504-3512

D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. International Conference on Learning Representations, 2015



## Surrogate Model

Mark Craven, Jude W. Shavlik: Extracting Tree-Structured Representations of Trained Networks. NIPS 1995: 24-30

# Overview of explanation in different AI fields (3)

## I Computer Vision

### Train

res5c unit 924



res5c unit 2001



inception\_5b unit 626



inception\_5b unit 415



### Airplane

res5c unit 1243



res5c unit 1379



inception\_4e unit 92



## Interpretable Units

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017: 3319-3327

# Overview of explanation in different AI fields (3)

## Computer Vision

Train

res5c unit 924



res5c unit 2001



inception\_5b unit 626



inception\_5b unit 415



Airplane

res5c unit 1243



res5c unit 1379

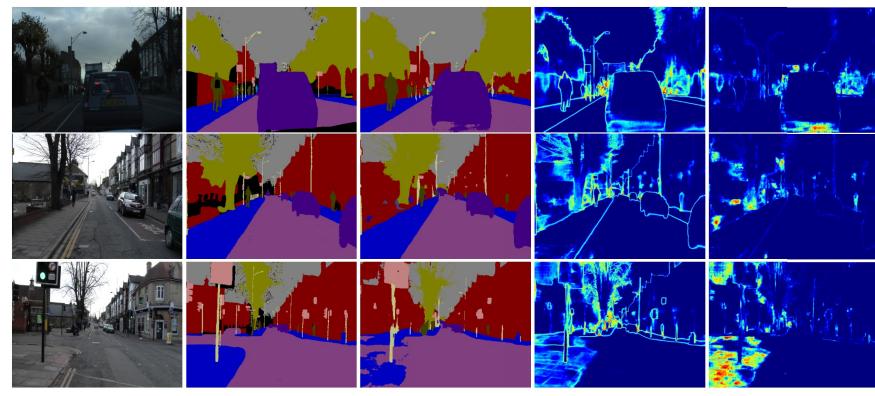


inception\_4e unit 92



## Interpretable Units

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017: 3319-3327



## Uncertainty Map

Alex Kendall, Yarin Gal: What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? NIPS 2017: 5580-5590

# Overview of explanation in different AI fields (3)

## Computer Vision

### Train

res5c unit 924



res5c unit 2001



inception\_5b unit 626

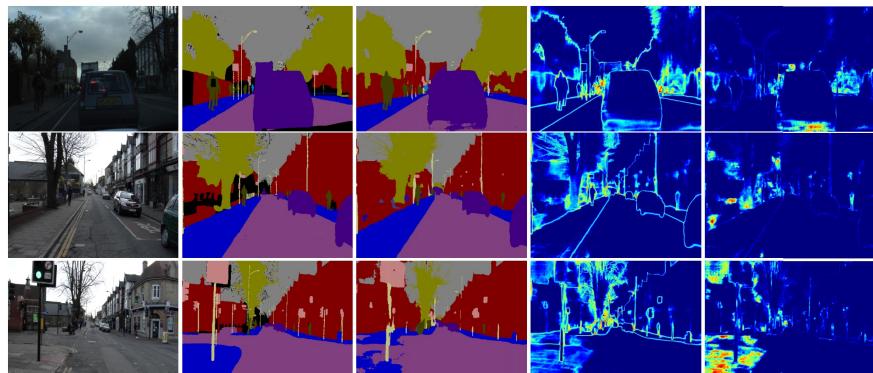


inception\_5b unit 415



## Interpretable Units

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017: 3319-3327



(a) Input Image

(b) Ground Truth

(c) Semantic Segmentation

(d) Aleatoric Uncertainty

(e) Epistemic Uncertainty

## Uncertainty Map

Alex Kendall, Yarin Gal: What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? NIPS 2017: 5580-5590

### Airplane

res5c unit 1243



res5c unit 1379



inception\_4b unit 92



### Western Grebe



Description: This is a large bird with a white neck and a black back in the water.

Class Definition: The *Western Grebe* is a waterbird with a yellow pointy beak, white neck and belly, and black back.

Explanation: This is a *Western Grebe* because this bird has a long white neck, pointy yellow beak and red eye.

### Laysan Albatross



Description: This is a large flying bird with black wings and a white belly.

Class Definition: The *Laysan Albatross* is a large seabird with a hooked yellow beak, black back and white belly.

Visual Explanation: This is a *Laysan Albatross* because this bird has a large wingspan, hooked yellow beak, and white belly.

### Laysan Albatross



Description: This is a large bird with a white neck and a black back in the water.

Class Definition: The *Laysan Albatross* is a large seabird with a hooked yellow beak, black back and white belly.

Visual Explanation: This is a *Laysan Albatross* because this bird has a hooked yellow beak white neck and black back.

## Visual Explanation

Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, Trevor Darrell: Generating Visual Explanations. ECCV (4) 2016: 3-19

# Overview of explanation in different AI fields (3)

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

## Computer Vision

### Train

res5c unit 924



res5c unit 2001



inception\_5b unit 626

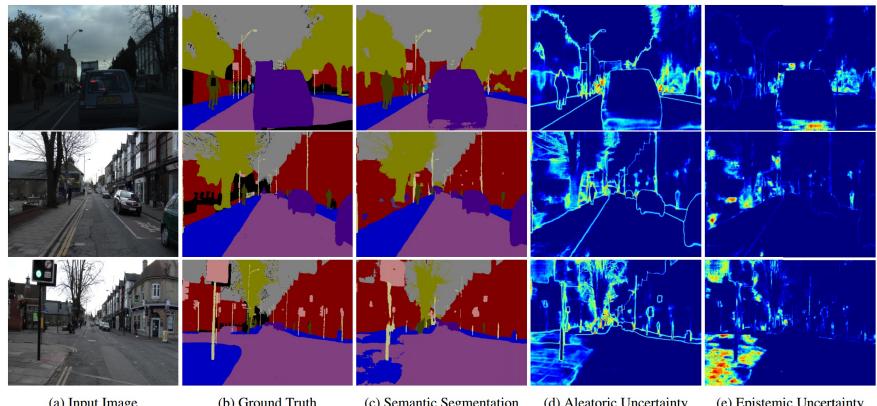


inception\_5b unit 415



### Interpretable Units

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba: Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017: 3319-3327



### Uncertainty Map

Alex Kendall, Yarin Gal: What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? NIPS 2017: 5580-5590

### Airplane

res5c unit 1243



res5c unit 1379



inception\_4b unit 92



### Western Grebe



Description: This is a large bird with a white neck and a black back in the water.

Class Definition: The *Western Grebe* is a waterbird with a yellow pointy beak, white neck and belly, and black back.

Explanation: This is a *Western Grebe* because this bird has a long white neck, pointy yellow beak and red eye.

### Laysan Albatross



Description: This is a large flying bird with black wings and a white belly.

Class Definition: The *Laysan Albatross* is a large seabird with a hooked yellow beak, black back and white belly.

Visual Explanation: This is a *Laysan Albatross* because this bird has a large wingspan, hooked yellow beak, and white belly.

### Laysan Albatross



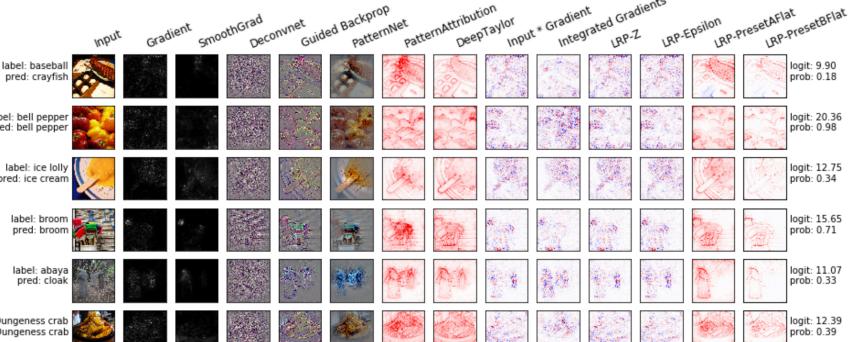
Description: This is a large bird with a white neck and a black back in the water.

Class Definition: The *Laysan Albatross* is a large seabird with a hooked yellow beak, black back and white belly.

Visual Explanation: This is a *Laysan Albatross* because this bird has a hooked yellow beak white neck and black back.

## Visual Explanation

Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, Trevor Darrell: Generating Visual Explanations. ECCV (4) 2016: 3-19

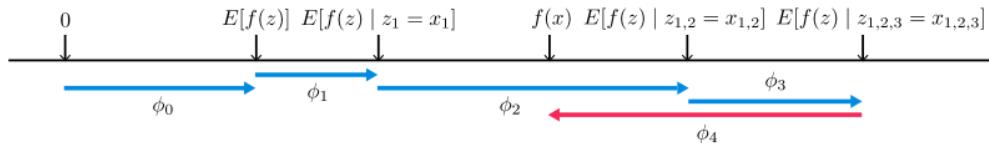
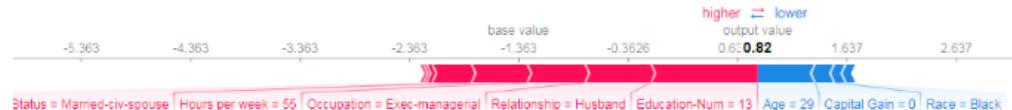


## Saliency Map

Julius Adebayo, Justin Gilmer, Michael Muelly, Ian J. Goodfellow, Moritz Hardt, Been Kim: Sanity Checks for Saliency Maps. NeurIPS 2018: 925-936

# Overview of explanation in different AI fields (4)

## I Game Theory



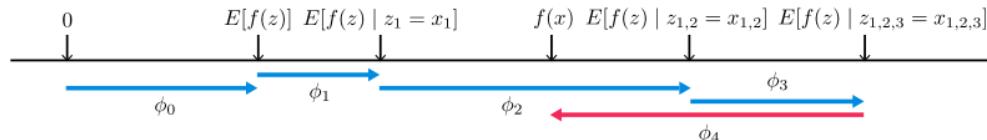
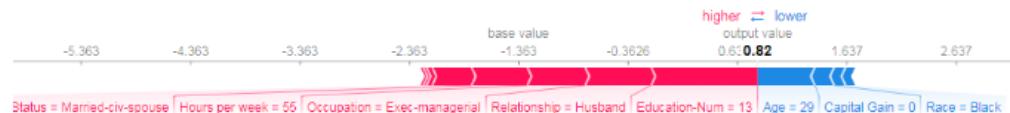
## Shapley Additive Explanation

Scott M. Lundberg, Su-In Lee: A Unified Approach to Interpreting Model Predictions.  
NIPS 2017: 4768-4777

# Overview of explanation in different AI fields (4)

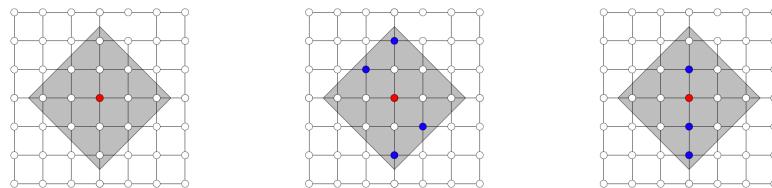
This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

## I Game Theory



## Shapley Additive Explanation

Scott M. Lundberg, Su-In Lee: A Unified Approach to Interpreting Model Predictions.  
NIPS 2017: 4768-4777



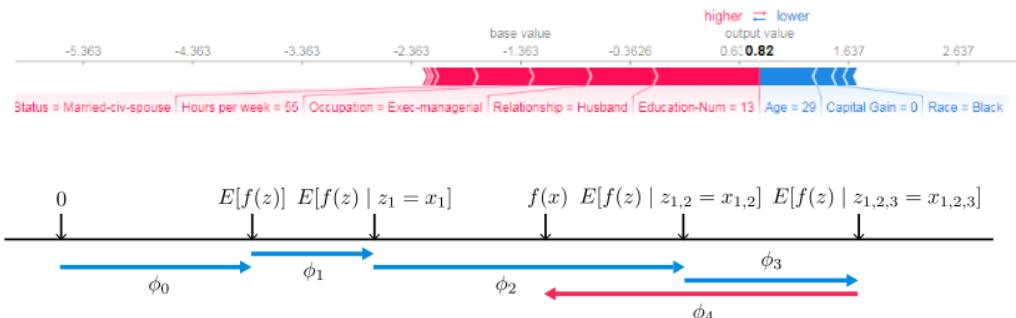
## L-Shapley and C-Shapley (with graph structure)

Jianbo Chen, Le Song, Martin J. Wainwright, Michael I. Jordan: L-Shapley and C-Shapley: Efficient Model Interpretation for Structured Data. ICLR 2019

# Overview of explanation in different AI fields (4)

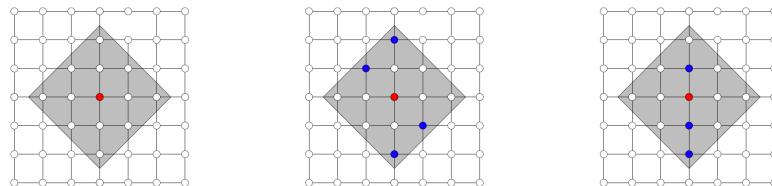
This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

## Game Theory



## Shapley Additive Explanation

Scott M. Lundberg, Su-In Lee: A Unified Approach to Interpreting Model Predictions.  
NIPS 2017: 4768-4777



## L-Shapley and C-Shapley (with graph structure)

Jianbo Chen, Le Song, Martin J. Wainwright, Michael I. Jordan: L-Shapley and C-Shapley: Efficient Model Interpretation for Structured Data. ICLR 2019

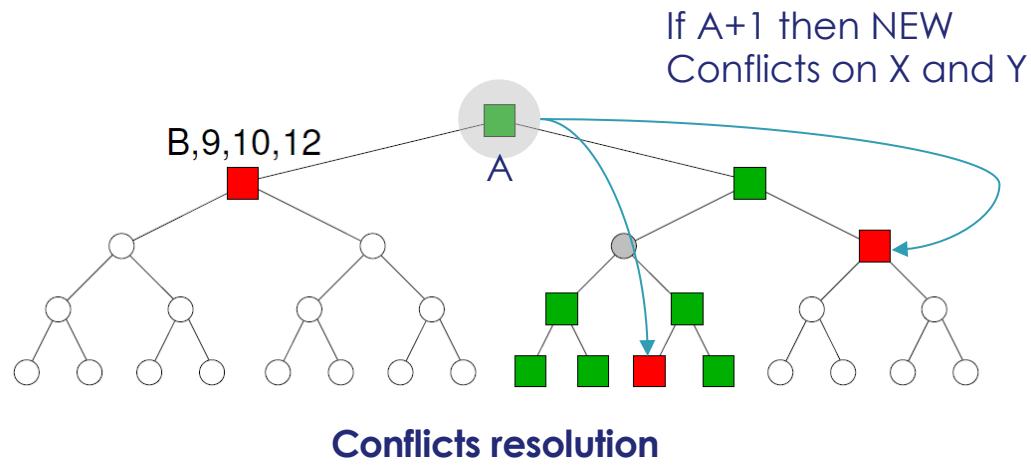
~ instancewise  
feature importance  
(causal influence)

Erik Štrumbelj and Igor Kononenko. An efficient explanation of individual classifications using game theory. Journal of Machine Learning Research, 11:1–18, 2010.

Anupam Datta, Shayak Sen, and Yair Zick. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In Security and Privacy (SP), 2016 IEEE Symposium on, pp. 598–617. IEEE, 2016.

# Overview of explanation in different AI fields (5)

## | Search and Constraint Satisfaction



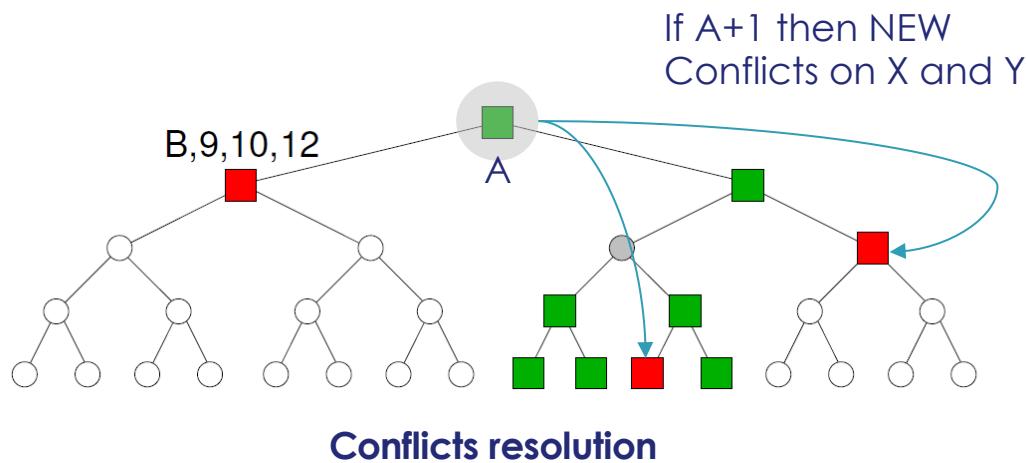
Barry O'Sullivan, Alexandre Papadopoulos, Boi Faltings, Pearl Pu: Representative Explanations for Over-Constrained Problems. AAAI 2007: 323-328

## Robustness Computation

Hebrard, E., Hnich, B., & Walsh, T. (2004, July). Robust solutions for constraint satisfaction and optimization. In ECAI (Vol. 16, p. 186).

# Overview of explanation in different AI fields (5)

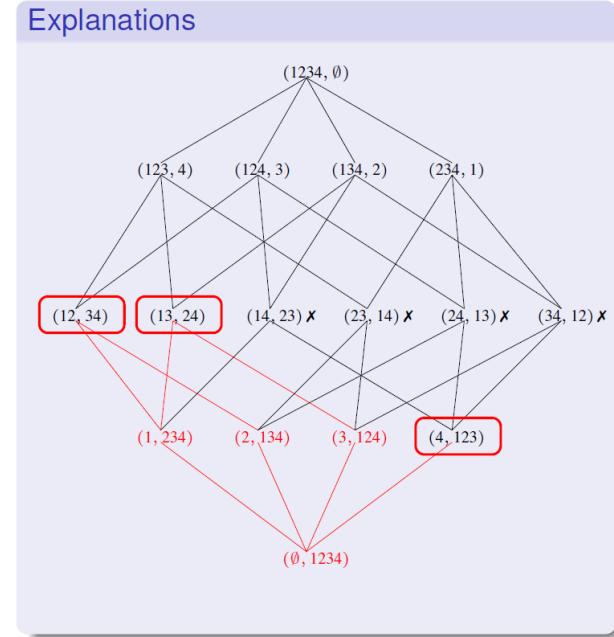
## Search and Constraint Satisfaction



Barry O'Sullivan, Alexandre Papadopoulos, Boi Faltings, Pearl Pu: Representative Explanations for Over-Constrained Problems. AAAI 2007: 323-328

## Robustness Computation

Hebrard, E., Hnich, B., & Walsh, T. (2004, July). Robust solutions for constraint satisfaction and optimization. In ECAI (Vol. 16, p. 186).



## Constraints relaxation

Ulrich Junker: QUICKPLAIN: Preferred Explanations and Relaxations for Over-Constrained Problems. AAAI 2004: 167-172

# Overview of explanation in different AI fields (6)

## Knowledge Representation and Reasoning

Ref	$\vdash C \implies C$	
Trans	$\vdash c \implies d, \vdash d \implies e \quad \vdash c \implies e$	
Eq	$\vdash A=B \quad \vdash c \implies d \quad \vdash c(A/B) \implies d(A/B)$	
Prim	$\vdash (\text{prim } EE) \implies (\text{prim } FP)$	$FF \subset EE$
THING	$\vdash C \equiv \text{THING}$	
AndR	$\vdash c \implies d, \vdash c \implies (\text{and } EE) \quad \vdash c \implies (\text{and } d \text{ } EE)$	
AndL	$\vdash c \implies e \quad \vdash (\text{and } ...c ...) \implies e$	
All	$\vdash c \implies d \quad \vdash (\text{all } p \text{ } c) \implies (\text{all } p \text{ } d)$	
AtLst	$\vdash (\text{at-least } n \text{ } p) \implies (\text{at-least } m \text{ } p)$	$n > m$
AndEq	$\vdash C \equiv (\text{and } C)$	
AtLst0	$\vdash (\text{at - least } 0 \text{ } p) \equiv \text{THING}$	
All-thing	$\vdash (\text{all } p \text{ } \text{THING}) \equiv \text{THING}$	
All-and	$\vdash (\text{and } (\text{all } p \text{ } C) (\text{all } p \text{ } D) \dots ) \equiv (\text{and } (\text{all } p \text{ } (\text{and } C \text{ } D)) \dots )$	

1.  $(\text{at-least } 3 \text{ grape}) \implies (\text{at-least } 2 \text{ grape})$  AtLst
2.  $(\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim } \text{GOOD WINE}) \implies (\text{at-least } 2 \text{ grape})$  AndL,1
3.  $(\text{prim } \text{GOOD WINE}) \implies (\text{prim } \text{WINE})$  Prim
4.  $(\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim } \text{GOOD WINE}) \implies (\text{prim } \text{WINE})$  AndL,3
5.  $A \equiv (\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim } \text{GOOD WINE})$  Told
6.  $A \implies (\text{prim } \text{WINE})$  Eq,4,5
7.  $(\text{prim } \text{WINE}) \equiv (\text{and } (\text{prim } \text{WINE}))$  AndEq
8.  $A \implies (\text{and } (\text{prim } \text{WINE}))$  Eq,7,6
9.  $A \implies (\text{at-least } 2 \text{ grape})$  Eq,5,2
10.  $A \implies (\text{and } (\text{at-least } 2 \text{ grape})) (\text{prim } \text{WINE})$  AndR,9,8

$A \equiv (\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim } \text{GOOD WINE})$

## Explaining Reasoning (through Justification) e.g., Subsumption

Deborah L. McGuinness, Alexander Borgida: Explaining Subsumption in Description Logics. IJCAI (1) 1995: 816-821

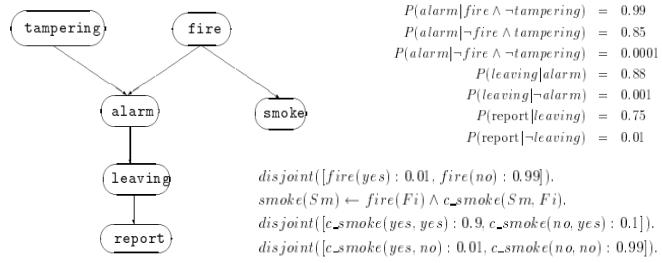
# Overview of explanation in different AI fields (6)

## Knowledge Representation and Reasoning

Ref	$\vdash C \Rightarrow C$
Trans	$\vdash c \Rightarrow d, \vdash d \Rightarrow e \quad \vdash c \Rightarrow e$
Eq	$\vdash A=B \quad \vdash c \Rightarrow d \quad \vdash c(A/B) \Rightarrow d(A/B)$
Prim	$\vdash (\text{prim } EE) \Rightarrow (\text{prim } FF)$
THING	$\vdash C \Rightarrow \text{THING}$
AndR	$\vdash c \Rightarrow d, \vdash c \Rightarrow (\text{and } EE) \quad \vdash c \Rightarrow (\text{and } d \text{ } EE)$
AndL	$\vdash c \Rightarrow e \quad \vdash (\text{and } ... c ...) \Rightarrow e$
All	$\vdash c \Rightarrow d \quad \vdash (\text{all } p \text{ } c) \Rightarrow (\text{all } p \text{ } d)$
AtLst	$\vdash_{i>m} (\text{at-least } n \text{ } p) \Rightarrow (\text{at-least } m \text{ } p)$
AndEq	$\vdash C \equiv (\text{and } C)$
AtL0	$\vdash (\text{at - least } 0 \text{ } p) \equiv \text{THING}$
All-thing	$\vdash (\text{all } p \text{ } \text{THING}) \equiv \text{THING}$
All-and	$\vdash (\text{and } (\text{all } p \text{ } C) (\text{all } p \text{ } D) ...) \equiv (\text{and } (\text{all } p \text{ } (\text{and } C \text{ } D)) ...)$

1.  $(\text{at-least } 3 \text{ grape}) \Rightarrow (\text{at-least } 2 \text{ grape})$  AtLst
2.  $(\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim GOOD WINE}) \Rightarrow (\text{at-least } 2 \text{ grape})$  AndL,1
3.  $(\text{prim GOOD WINE}) \Rightarrow (\text{prim WINE})$  Prim
4.  $(\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim GOOD WINE}) \Rightarrow (\text{prim WINE})$  AndL,3
5.  $A \equiv (\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim GOOD WINE})$  Told
6.  $A \Rightarrow (\text{prim WINE})$  Eq,4,5
7.  $(\text{prim WINE}) \equiv (\text{and } (\text{prim WINE}))$  AndEq
8.  $A \Rightarrow (\text{and } (\text{prim WINE}))$  Eq,7,6
9.  $A \Rightarrow (\text{at-least } 2 \text{ grape})$  Eq,5,2
10.  $A \Rightarrow (\text{and } (\text{at-least } 2 \text{ grape})) (\text{prim WINE})$  AndR,9,8

$A \equiv (\text{and } (\text{at-least } 3 \text{ grape})) (\text{prim GOOD WINE})$



## Abduction Reasoning (in Bayesian Network)

David Poole: Probabilistic Horn Abduction and Bayesian Networks. Artif. Intell. 64(1): 81-129 (1993)

## Explaining Reasoning (through Justification) e.g., Subsumption

Deborah L. McGuinness, Alexander Borgida: Explaining Subsumption in Description Logics. IJCAI (1) 1995: 816-821

# Overview of explanation in different AI fields (6)

## Knowledge Representation and Reasoning

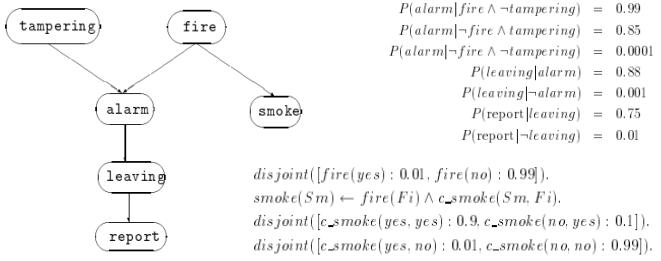
Ref	$\vdash C \Rightarrow C$
Trans	$\vdash c \Rightarrow d, \vdash d \Rightarrow e \quad \vdash c \Rightarrow e$
Eq	$\vdash A=B \quad \vdash c(A/B) \quad \vdash c \Rightarrow d \quad \vdash c(A/B) \Rightarrow d(A/B)$
Prim	$\vdash (\text{prim } EE) \Rightarrow (\text{prim } FP)$
THING	$\vdash c \Rightarrow \text{THING}$
AndR	$\vdash c \Rightarrow d, \vdash c \Rightarrow (\text{and } EE) \quad \vdash c \Rightarrow (\text{and } d \text{ } EE)$
AndL	$\vdash c \Rightarrow e \quad \vdash (\text{and } \dots c \dots) \Rightarrow e$
All	$\vdash c \Rightarrow d \quad \vdash (\text{all } p \text{ } c) \Rightarrow (\text{all } p \text{ } d)$
AtLst	$\vdash_{i>m} (\text{at-least } n \text{ } p) \Rightarrow (\text{at-least } m \text{ } p)$
AndEq	$\vdash C \equiv (\text{and } C)$
AtLs0	$\vdash (\text{at least } 0 \text{ } p) \equiv \text{THING}$
All-thing	$\vdash (\text{all } p \text{ } \text{THING}) \equiv \text{THING}$
All-and	$\vdash (\text{and } (\text{all } p \text{ } C) (\text{all } p \text{ } D) \dots) \equiv (\text{and } (\text{all } p \text{ } (\text{and } C \text{ } D)) \dots)$

1.  $(\text{at-least } 3 \text{ grape}) \Rightarrow (\text{at-least } 2 \text{ grape})$  AtLst
2.  $(\text{and } (\text{at-least } 3 \text{ grape}) (\text{prim GOOD WINE})) \Rightarrow (\text{at-least } 2 \text{ grape})$  AndL,1
3.  $(\text{prim GOOD WINE}) \Rightarrow (\text{prim WINE})$  Prim
4.  $(\text{and } (\text{at-least } 3 \text{ grape}) (\text{prim GOOD WINE})) \Rightarrow (\text{prim WINE})$  AndL,3
5.  $A \equiv (\text{and } (\text{at-least } 3 \text{ grape}) (\text{prim GOOD WINE}))$  Told
6.  $A \Rightarrow (\text{prim WINE})$  Eq,4,5
7.  $(\text{prim WINE}) \equiv (\text{and } (\text{prim WINE}))$  AndEq
8.  $A \Rightarrow (\text{and } (\text{prim WINE}))$  Eq,7,6
9.  $A \Rightarrow (\text{at-least } 2 \text{ grape})$  Eq,5,2
10.  $A \Rightarrow (\text{and } (\text{at-least } 2 \text{ grape}) (\text{prim WINE}))$  AndR,9,8

$A \equiv (\text{and } (\text{at-least } 3 \text{ grape}) (\text{prim GOOD WINE}))$

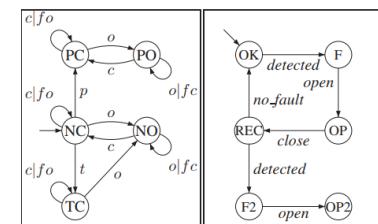
## Explaining Reasoning (through Justification) e.g., Subsumption

Deborah L. McGuinness, Alexander Borgida: Explaining Subsumption in Description Logics. IJCAI (1) 1995: 816-821



## Abduction Reasoning (in Bayesian Network)

David Poole: Probabilistic Horn Abduction and Bayesian Networks. Artif. Intell. 64(1): 81-129 (1993)



## Diagnosis Inference

Alban Grastien, Patrik Haslum, Sylvie Thiébaut: Conflict-Based Diagnosis of Discrete Event Systems: Theory and Practice. KR 2012

# Overview of explanation in different AI fields (7)

## Multi-agent Systems

MAS INFRASTRUCTURE	INDIVIDUAL AGENT INFRASTRUCTURE
<b>MAS INTEROPERATION</b> Translation Services   Interoperation Services	<b>INTEROPERATION</b> Interoperation Modules
<b>CAPABILITY TO AGENT MAPPING</b> Middle Agents	<b>CAPABILITY TO AGENT MAPPING</b> Middle Agents Components
<b>NAME TO LOCATION MAPPING</b> ANS	<b>NAME TO LOCATION MAPPING</b> ANS Component
<b>SECURITY</b> Certificate Authority   Cryptographic Services	<b>SECURITY</b> Security Module   private/public Keys
<b>PERFORMANCE SERVICES</b> MAS Monitoring   Reputation Services	<b>PERFORMANCE SERVICES</b> Performance Services Modules
<b>MULTIAGENT MANAGEMENT SERVICES</b> Logging, Activity Visualization, Launching	<b>MANAGEMENT SERVICES</b> Logging and Visualization Components
<b>ACL INFRASTRUCTURE</b> Public Ontology   Protocols Servers	<b>ACL INFRASTRUCTURE</b> ACL Parser   Private Ontology   Protocol Engine
<b>COMMUNICATION INFRASTRUCTURE</b> Discovery   Message Transfer	<b>COMMUNICATION MODULES</b> Discovery Component   Message Transfer Module
<b>OPERATING ENVIRONMENT</b> Machines, OS, Network   Multicast   Transport Layer: TCP/IP, Wireless, Infrared, SSL	

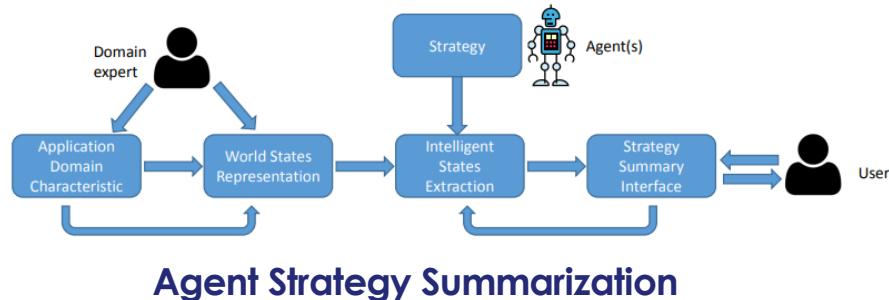
## Explanation of Agent Conflicts & Harmful Interactions

Katia P. Sycara, Massimo Paolucci, Martin Van Velsen, Joseph A. Giampaapa: The RETSINA MAS Infrastructure. Autonomous Agents and Multi-Agent Systems 7(1-2): 29-48 (2003)

# Overview of explanation in different AI fields (7)

## Multi-agent Systems

MAS INFRASTRUCTURE		INDIVIDUAL AGENT INFRASTRUCTURE	
MAS INTEROPERATION Translation Services	Interoperation Services	INTEROPERATION Interoperation Modules	
CAPABILITY TO AGENT MAPPING Middle Agents		CAPABILITY TO AGENT MAPPING Middle Agents Components	
NAME TO LOCATION MAPPING ANS		NAME TO LOCATION MAPPING ANS Component	
SECURITY Certificate Authority	Cryptographic Services	SECURITY Security Module	private/public Keys
PERFORMANCE SERVICES MAS Monitoring	Reputation Services	PERFORMANCE SERVICES Performance Services Modules	
MULTIAGENT MANAGEMENT SERVICES Logging, Activity Visualization, Launching		MANAGEMENT SERVICES Logging and Visualization Components	
ACL INFRASTRUCTURE Public Ontology	Protocols Servers	ACL INFRASTRUCTURE ACL Parser Private Ontology Protocol Engine	
COMMUNICATION INFRASTRUCTURE Discovery	Message Transfer	COMMUNICATION MODULES Discovery Component Message Transfer Module	
OPERATING ENVIRONMENT Machines, OS, Network		Multicast Transport Layer: TCP/IP, Wireless, Infrared, SSL	



## Agent Strategy Summarization

Ofra Amir, Finale Doshi-Velez, David Sarne: Agent Strategy Summarization. AAMAS 2018: 1203-1207

## Explanation of Agent Conflicts & Harmful Interactions

Katia P. Sycara, Massimo Paolucci, Martin Van Velsen, Joseph A. Giampaapa: The RETSINA MAS Infrastructure. Autonomous Agents and Multi-Agent Systems 7(1-2): 29-48 (2003)

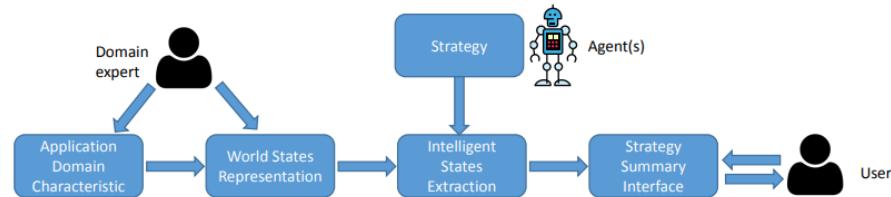
# Overview of explanation in different AI fields (7)

## Multi-agent Systems

MAS INFRASTRUCTURE		INDIVIDUAL AGENT INFRASTRUCTURE	
MAS INTEROPERATION	Translation Services Interoperation Services	INTEROPERATION	Interoperation Modules
CAPABILITY TO AGENT MAPPING	Middle Agents	CAPABILITY TO AGENT MAPPING	Middle Agents Components
NAME TO LOCATION MAPPING	ANS	NAME TO LOCATION MAPPING	ANS Component
SECURITY	Certificate Authority Cryptographic Services	SECURITY	Security Module private/public Keys
PERFORMANCE SERVICES	MAS Monitoring Reputation Services	PERFORMANCE SERVICES	Performance Services Modules
MULTIAGENT MANAGEMENT SERVICES	Logging, Activity Visualization, Launching	MANAGEMENT SERVICES	Logging and Visualization Components
ACL INFRASTRUCTURE	Public Ontology Protocols Servers	ACL INFRASTRUCTURE	ACL Parser Private Ontology Protocol Engine
COMMUNICATION INFRASTRUCTURE	Discovery Message Transfer	COMMUNICATION MODULES	Discovery Component Message Transfer Module
OPERATING ENVIRONMENT		Operating Environment	
Machines, OS, Network		Multicast Transport Layer: TCP/IP, Wireless, Infrared, SSL	

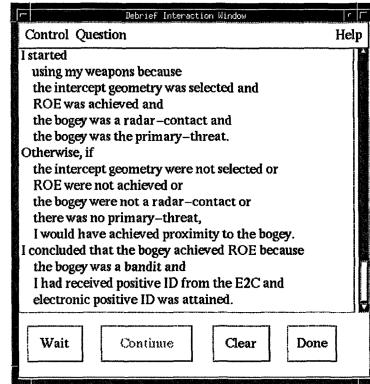
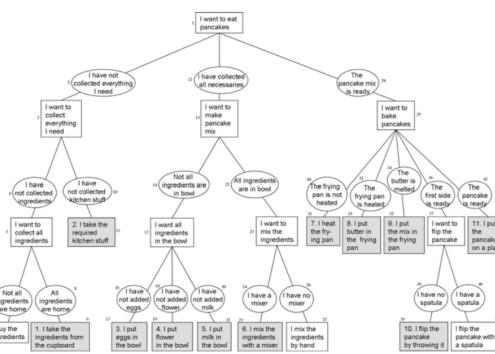
## Explanation of Agent Conflicts & Harmful Interactions

Katia P. Sycara, Massimo Paolucci, Martin Van Velsen, Joseph A. Giampaapa: The RETSINA MAS Infrastructure. Autonomous Agents and Multi-Agent Systems 7(1-2): 29-48 (2003)



## Agent Strategy Summarization

Ofra Amir, Finale Doshi-Velez, David Sarne: Agent Strategy Summarization. AAMAS 2018: 1203-1207



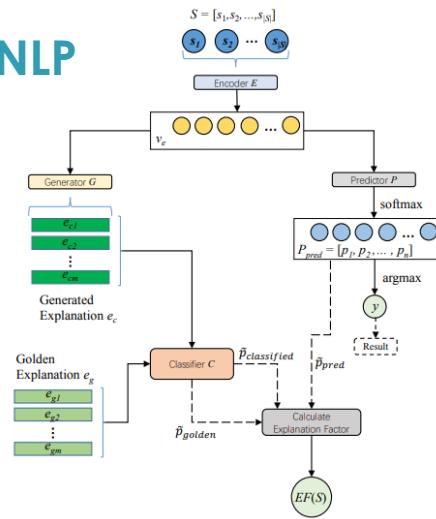
## Explainable Agents

Joost Broekens, Maaike Harbers, Koen V. Hindriks, Karel van den Bosch, Catholijn M. Jonker, John-Jules Ch. Meyer: Do You Get It? User-Evaluated Explainable BDI Agents. MATES 2010: 28-39

W. Lewis Johnson: Agents that Learn to Explain Themselves. AAAI 1994: 1257-1263

# Overview of explanation in different AI fields (8)

## NLP



Fine-grained explanations are in the form of:

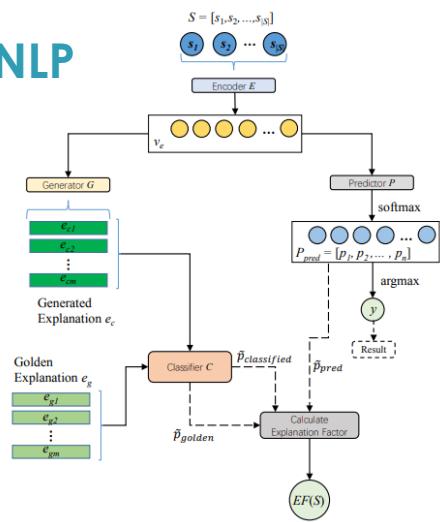
- texts in a real-world dataset;
- Numerical scores

## Explainable NLP

Hui Liu, Qingyu Yin, William Yang Wang: Towards Explainable NLP: A Generative Explanation Framework for Text Classification. CoRR abs/1811.00196 (2018)

# Overview of explanation in different AI fields (8)

## NLP



## Explainable NLP

Hui Liu, Qingyu Yin, William Yang Wang: Towards Explainable NLP: A Generative Explanation Framework for Text Classification. CoRR abs/1811.00196 (2018)

- Fine-grained explanations are in the form of:
- texts in a real-world dataset;
  - Numerical scores

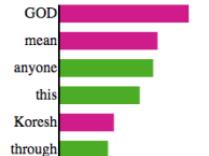
Example #3 of 6

True Class: Atheism

Instructions Previous Next

### Algorithm 1

Words that A1 considers important:



Predicted:

Atheism

Prediction correct:



### Document

From: pauld@verdix.com (Paul Durbin)  
Subject: Re: DAVID CORESH IS! **GOD!**  
Nntp-Posting-Host: sarge.hq.verdix.com  
Organization: Verdix Corp  
Lines: 8

### Algorithm 2

Words that A2 considers important:



Predicted:

Atheism

Prediction correct:



### Document

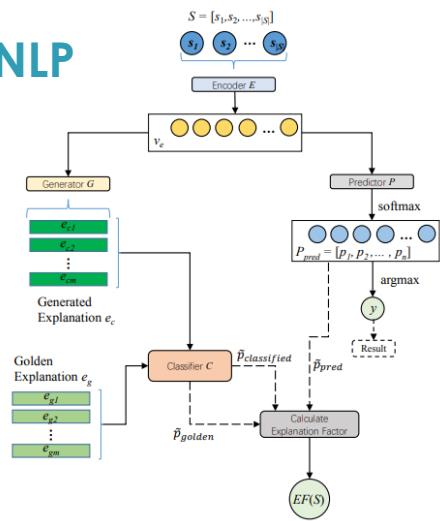
From: pauld@verdix.com (Paul Durbin)  
Subject: **Re: DAVID CORESH IS! GOD!**  
**Nntp-Posting-Host:** sarge.hq.verdix.com  
Organization: Verdix Corp  
Lines: 8

## LIME for NLP

Marco Túlio Ribeiro, Sameer Singh, Carlos Guestrin: "Why Should I Trust You?": Explaining the Predictions of Any Classifier. KDD 2016: 1135-1144

# Overview of explanation in different AI fields (8)

## NLP



### Explainable NLP

Hui Liu, Qingyu Yin, William Yang Wang: Towards Explainable NLP: A Generative Explanation Framework for Text Classification. CoRR abs/1811.00196 (2018)

Hendrik Strobelt, Sebastian Gehrmann, Hanspeter Pfister, Alexander M. Rush: LSTMVis: A Tool for Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks. IEEE Trans. Vis. Comput. Graph. 24(1): 667-676 (2018)

Hendrik Strobelt, Sebastian Gehrmann, Michael Behrisch, Adam Perer, Hanspeter Pfister, Alexander M. Rush: Seq2seq-Vis: A Visual Debugging Tool for Sequence-to-Sequence Models. IEEE Trans. Vis. Comput. Graph. 25(1): 353-363 (2019)

- Fine-grained explanations are in the form of:
- texts in a real-world dataset;
  - Numerical scores

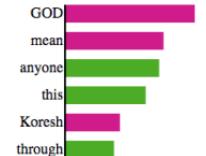
Example #3 of 6

True Class: Atheism

Instructions Previous Next

#### Algorithm 1

Words that A1 considers important:



Predicted:

Atheism

Prediction correct:

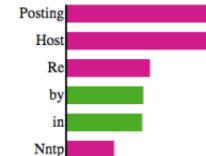


#### Document

From: pauld@verdix.com (Paul Durbin)  
Subject: Re: DAVID CORESH IS! GOD!  
Nntp-Posting-Host: sarge.hq.verdix.com  
Organization: Verdix Corp  
Lines: 8

#### Algorithm 2

Words that A2 considers important:



Predicted:

Atheism

Prediction correct:



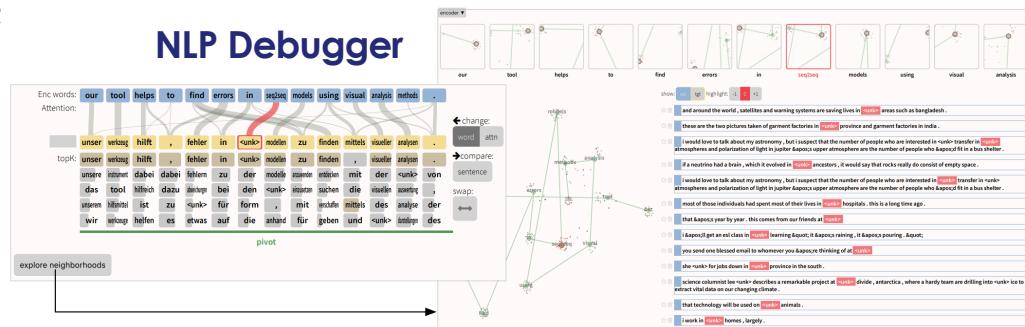
#### Document

From: pauld@verdix.com (Paul Durbin)  
Subject: Re: DAVID CORESH IS! GOD!  
Nntp-Posting-Host: sarge.hq.verdix.com  
Organization: Verdix Corp  
Lines: 8

## LIME for NLP

Marco Túlio Ribeiro, Sameer Singh, Carlos Guestrin: "Why Should I Trust You?": Explaining the Predictions of Any Classifier. KDD 2016: 1135-1144

## NLP Debugger

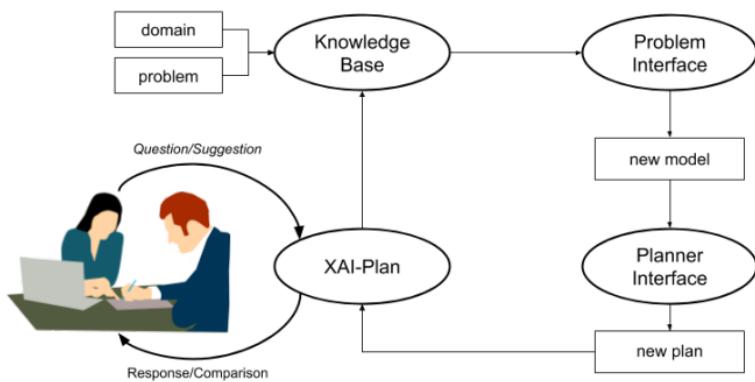


# Overview of explanation in different AI fields (9)

## Planning and Scheduling

Explanation Type	R1	R2	R3	R4
Plan Patch Explanation / VAL	✗	✓	✗	✓
Model Patch Explanation	✓	✗	✓	✓
Minimally Complete Explanation	✓	✓	✗	?
Minimally Monotonic Explanation	✓	✓	✓	?
(Approximate) Minimally Complete Explanation	✗	✓	✗	✓

Rita Borgo, Michael Cashmore, Daniele Magazzeni: Towards Providing Explanations for AI Planner Decisions. CoRR abs/1810.06338 (2018)



### XAI Plan

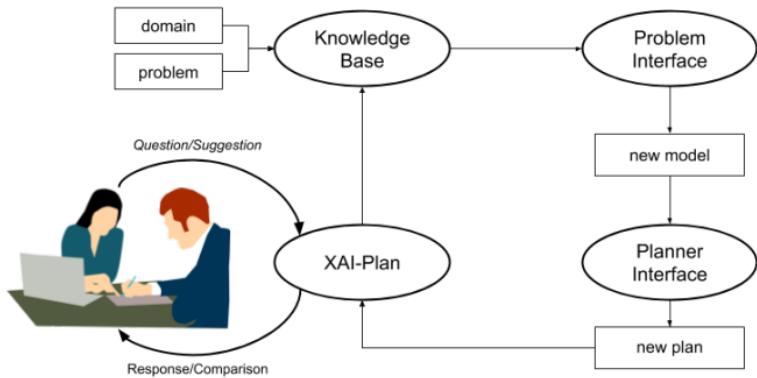
Rita Borgo, Michael Cashmore, Daniele Magazzeni: Towards Providing Explanations for AI Planner Decisions. CoRR abs/1810.06338 (2018)

# Overview of explanation in different AI fields (9)

## Planning and Scheduling

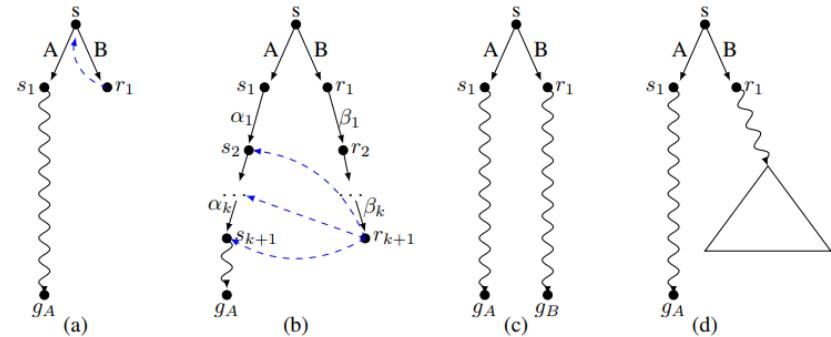
Explanation Type	R1	R2	R3	R4
Plan Patch Explanation / VAL	X	✓	X	✓
Model Patch Explanation	✓	X	✓	✓
Minimally Complete Explanation	✓	✓	X	?
Minimally Monotonic Explanation	✓	✓	✓	?
(Approximate) Minimally Complete Explanation	X	✓	X	✓

Rita Borgo, Michael Cashmore, Daniele Magazzeni: Towards Providing Explanations for AI Planner Decisions. CoRR abs/1810.06338 (2018)



## XAI Plan

Rita Borgo, Michael Cashmore, Daniele Magazzeni: Towards Providing Explanations for AI Planner Decisions. CoRR abs/1810.06338 (2018)



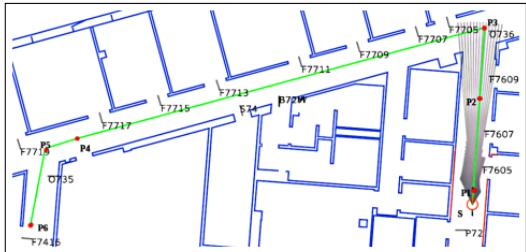
## Human-in-the-loop Planning

Maria Fox, Derek Long, Daniele Magazzeni: Explainable Planning. CoRR abs/1709.10256 (2017)

## (Manual) Plan Comparison

# Overview of explanation in different AI fields (10)

## Robotics



		Abstraction, A			
Specificity, S		Level 1	Level 2	Level 3	Level 4
	General Picture	Start and finish point of the complete route	Total distance and time taken for the complete route	Total distance and time taken for the complete route	Starting and ending landmark of complete route
	Summary	Start and finish point for subroute on each floor of each building	Total distance and time taken for subroute on each floor of each building	Total distance and angles for subroute on each floor of each building	Starting and ending landmark for subroute on each floor of each building
	Detailed Narrative	Start and finish points of complete route plus time taken for each edge of route	Angle turned at each point plus total distance and time taken for each edge of route	Turn direction at each point plus total distance for each edge of route	All landmarks encountered on the route

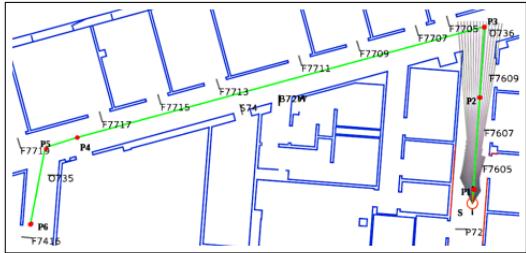
## Narration of Autonomous Robot Experience

Stephanie Rosenthal, Sai P Selvaraj, and Manuela Veloso. Verbalization: Narration of autonomous robot experience. In IJCAI, pages 862–868. AAAI Press, 2016.

Daniel J Brooks et al. 2010. Towards State Summarization for Autonomous Robots.. In AAAI Fall Symposium: Dialog with Robots, Vol. 61. 62.

# Overview of explanation in different AI fields (10)

## Robotics



Abstraction, A				
	Level 1	Level 2	Level 3	Level 4
General Picture	Start and finish point of the complete route	Total distance and time taken for the complete route	Total distance and time taken for the complete route	Starting and ending landmark of complete route
Summary	Start and finish point for subroute on each floor of each building	Total distance and time taken for subroute on each floor of each building	Total distance and angles for subroute on each floor of each building	Starting and ending landmark for subroute on each floor of each building
Detailed Narrative	Start and finish points of complete route plus time taken for each edge of route	Angle turned at each point plus the total distance and time taken for each edge of route	Turn direction at each point plus total distance for each edge of route	All landmarks encountered on the route

## Narration of Autonomous Robot Experience

Stephanie Rosenthal, Sai P Selvaraj, and Manuela Veloso. Verbalization: Narration of autonomous robot experience. In IJCAI, pages 862–868. AAAI Press, 2016.

Daniel J Brooks et al. 2010. Towards State Summarization for Autonomous Robots.. In AAAI Fall Symposium: Dialog with Robots, Vol. 61. 62.

**Robot:** I have decided to turn left.

**Human:** Why did you do that?

**Robot:** I believe that the correct action is to turn left  
BECAUSE:

I'm being asked to go forward

AND This area in front of me was 20 cm higher than me  
\*highlights area\*

AND the area to the left has maximum protrusions of less  
than 5 cm \*highlights area\*

AND I'm tilted to the right by more than 5 degrees.

Here is a display of the path through the tree that lead to  
this decision. \*displays tree\*

**Human:** How confident are you in this decision?

**Robot:** The distribution of actions that reached this leaf node is shown in this histogram. \*displays histogram\*  
This action is predicted to be correct 67% of the time.

**Human:** Where did the threshold for the area in front come from?

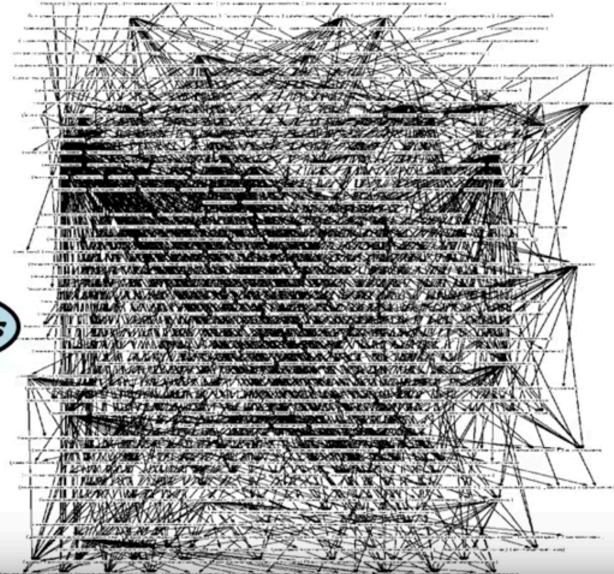
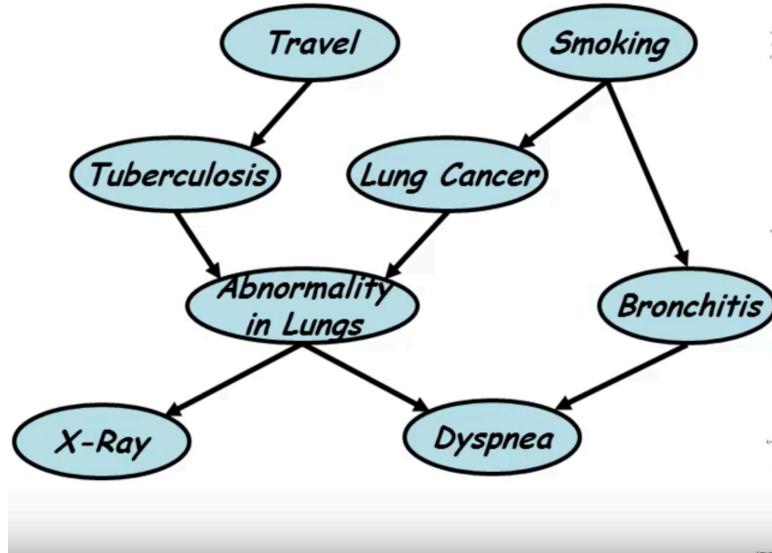
**Robot:** Here is the histogram of all training examples that  
reached this leaf. 80% of examples where this area was  
above 20 cm predicted the appropriate action to be “drive  
forward”.

## From Decision Tree to human-friendly information

Raymond Ka-Man Sheh: "Why Did You Do That?" Explainable Intelligent Robots. AAAI Workshops 2017

# Overview of explanation in different AI fields (11)

## | Reasoning under uncertainty



Probabilistic Graphical Models

Daphne Koller, Nir Friedman: Probabilistic Graphical Models - Principles and Techniques. MIT Press 2009, ISBN 978-0-262-01319-2, pp. I-XXXV, 1-1231

# Evaluation

# XAI: One Objective, Many Metrics



## Comprehensibility

How much effort for correct human interpretation?



## Succinctness

How concise and compact is the explanation?



## Actionability

What can one action, do with the explanation?



## Reusability

Could the explanation be personalized?



## Accuracy

How accurate and precise is the explanation?



## Completeness

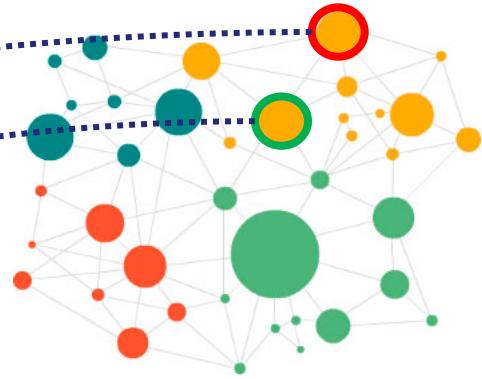
Is the explanation complete, partial, restricted?



# On the role of Knowledge Graphs in Explainable Machine Learning

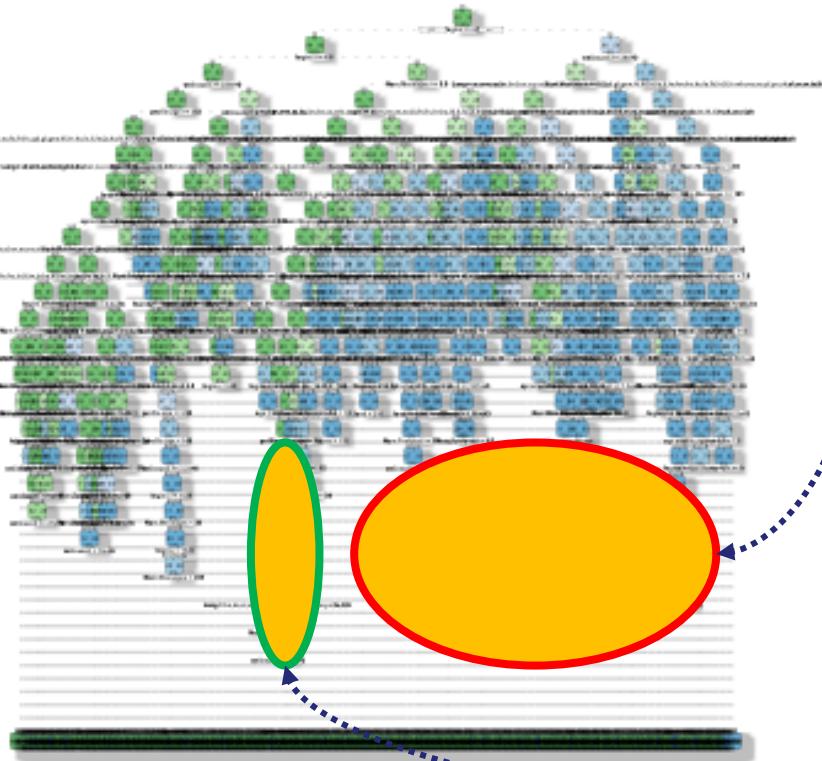
# Knowledge Graph Embeddings in Machine Learning

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.



<https://stats.stackexchange.com/questions/230581/decision-tree-too-large-to-interpret>

# Knowledge Graph for Decision Trees



Rattle 2016-Aug-18 16:15:42 sklisarov

<https://stats.stackexchange.com/questions/230581/decision-tree-too-large-to-interpret>

# Knowledge Graph for Deep Neural Network (1)

● Input Layer

Training Data

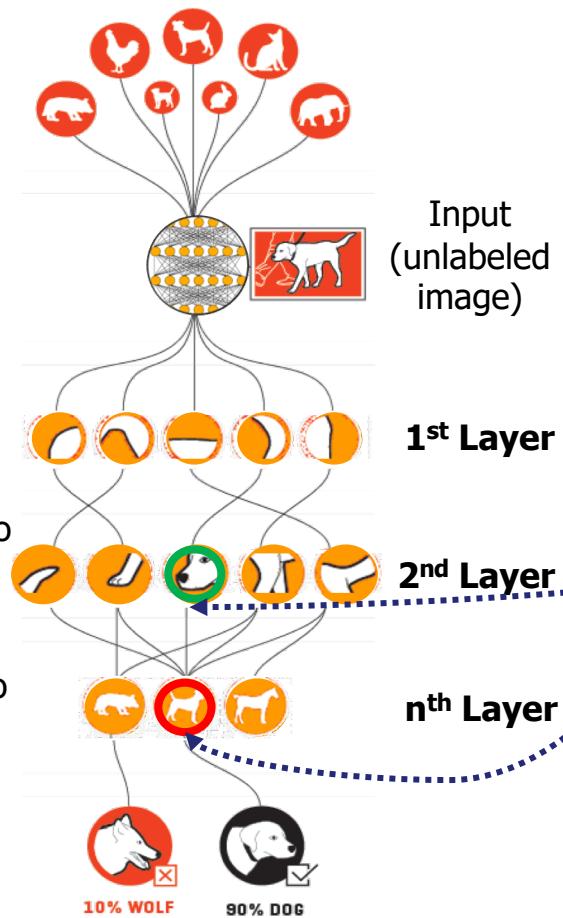
Neurons  
respond to  
simple shapes

Neurons respond to  
more complex  
structures

Neurons respond to  
highly complex,  
abstract concepts

77

● Output Layer



Input  
(unlabeled  
image)

1<sup>st</sup> Layer

2<sup>nd</sup> Layer

n<sup>th</sup> Layer

Low-level  
features to  
high-level  
features



# Knowledge Graph for Deep Neural Network (2)

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

Input Layer

Training Data

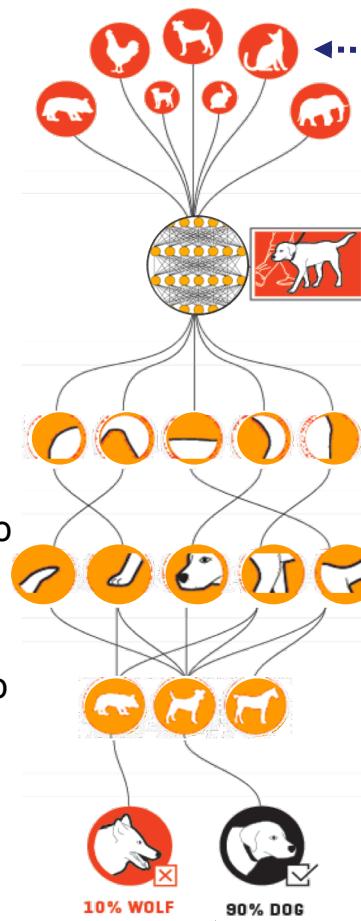
Neurons  
respond to  
simple shapes

Neurons respond to  
more complex  
structures

Neurons respond to  
highly complex,  
abstract concepts

Hidden Layer

Output Layer



Input  
(unlabeled  
image)

1<sup>st</sup> Layer

2<sup>nd</sup> Layer

n<sup>th</sup> Layer

Low-level  
features to  
high-level  
features

What is the causal  
relationship  
between the input /  
hidden / output  
layers

# Knowledge Graph for Personalized XAI



Description 1: This is an orange train accident

Description 2: This is an train accident between two speed merchant trains of characteristics X43-B and Y33-C in a dry environment

Description 3: This is a public transportation accident



## “How to explain transfer learning with appropriate knowledge representation?

Proceedings of the Sixteenth International Conference on Principles of Knowledge Representation and Reasoning (KR 2018)

### Knowledge-Based Transfer Learning Explanation

**Jiaoyan Chen**  
Department of Computer Science  
University of Oxford, UK

**Jeff Z. Pan**  
Department of Computer Science  
University of Aberdeen, UK

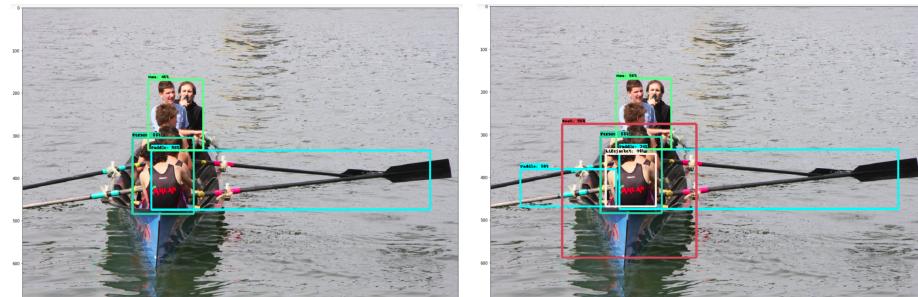
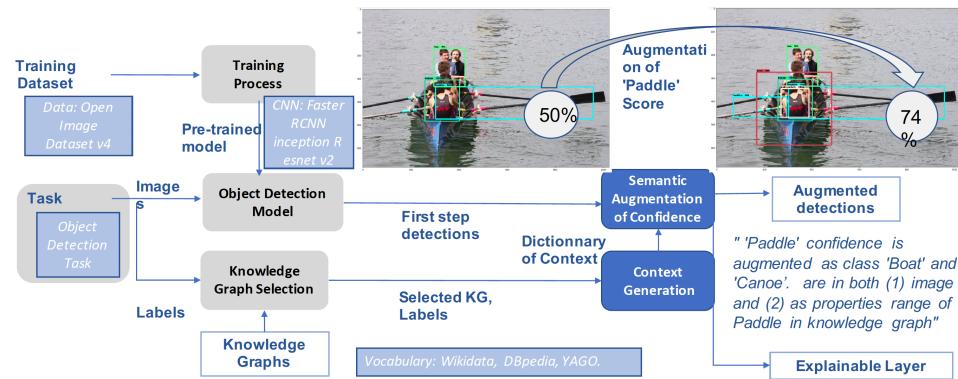
**Huajun Chen**  
College of Computer Science, Zhejiang University, China  
Alibaba-Zhejian University Frontier Technology Research Center

**Freddy Lecue**  
INRIA, France  
Accenture Labs, Ireland

**Ian Horrocks**  
Department of Computer Science  
University of Oxford, UK

# Applications

# Explainable Boosted Object Detection – Industry Agnostic



**Fig. 2.** Left image: results from baseline Faster RCNN: Paddle: 50% confidence, Person: 66%, Man: 46%. Right image: results from the semantic augmentation: Paddle: 74% confidence, Person: 66%, Man: 56%, Boat: 58% with explanation: Person, Paddle, Water as part of the context in the image and knowledge graph of concept Boat. (color print).

**Challenge:** Object detection is usually performed from a large portfolio of Artificial Neural Networks (ANNs) architectures trained on large amount of labelled data. Explaining object detections is rather difficult due to the high complexity of the most accurate ANNs.

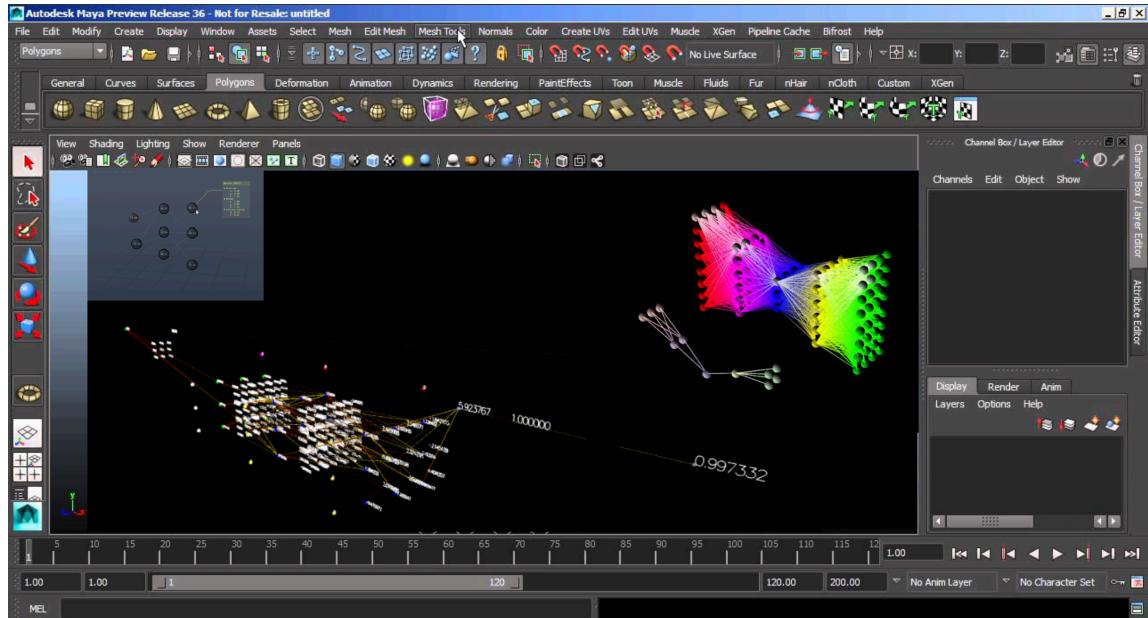
**AI Technology:** Integration of AI related technologies i.e., Machine Learning (Deep Learning / CNNs), and knowledge graphs / linked open data.

**XAI Technology:** Knowledge graphs and Artificial Neural Networks

**THALES**

# Debugging Artificial Neural Networks – Industry Agnostic

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.



Zetane.com

**Challenge:** Designing Artificial Neural Network architectures requires lots of experimentation (i.e., training phases) and parameters tuning (optimization strategy, learning rate, number of layers...) to reach optimal and robust machine learning models.

**AI Technology:** Artificial Neural Network

**XAI Technology:** Artificial Neural Network, 3D Modeling and Simulation Platform For AI

# Explaining Visual Question Answering – Industry Agnostic

## Tabular QA

Rank	Nation	Gold	Silver	Bronze	Total
1	India	102	58	37	197
2	Nepal	32	10	24	65
3	Sri Lanka	16	42	62	120
4	Pakistan	10	36	30	76
5	Bangladesh	2	10	35	47
6	Bhutan	1	6	7	14
7	Maldives	0	0	4	4

Q: How many medals did India win?  
A: 197

Neural Programmer (2017) model  
33.5% accuracy on WikiTableQuestions

## Visual QA



Q: How symmetrical are the white bricks on either side of the building?  
A: very

Kazemi and Elqursh (2017) model.  
61.1% on VQA 1.0 dataset  
(state of the art = 66.7%)

## Reading Comprehension

Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager

Q: Name of the quarterback who was 38 in Super Bowl XXXIII?  
A: John Elway

Yu et al (2018) model.  
84.6 F-1 score on SQuAD (state of the art)

Q: How symmetrical are the white bricks on either side of the building?  
A: very

Q: How **asymmetrical** are the white bricks on either side of the building?  
A: very

Q: How **big** are the white bricks on either side of the building?  
A: very

Q: How **fast** are the **bricks speaking** on either side of the building?  
A: very

**Challenge:** What is the robustness of Visual Question Answering models? What is the impact of semantics?

**AI Technology:** Artificial Neural Networks.

**XAI Technology:** Integrated Gradients



What is the **man** doing? → What is the **tweet** doing?  
How many **children** are there? → How many **tweet** are there?

VQA model's response remains the same 75.6% of the time on questions that it originally answered correctly

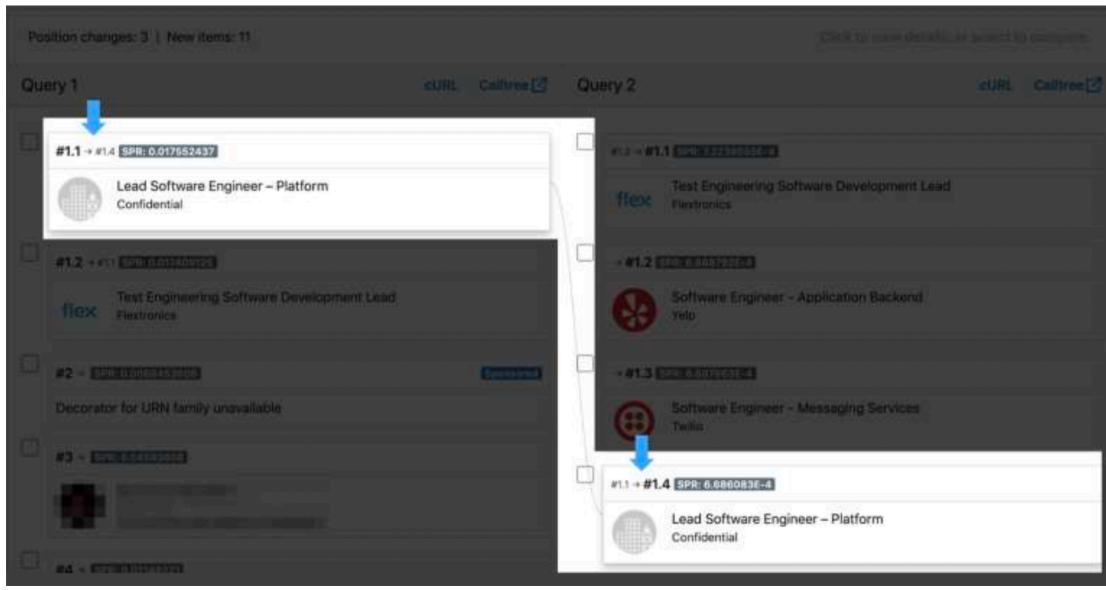
# Relevance Debugging and Explaining – Industry Agnostic



**Challenge:** A Machine Learning system can fail in many different points e.g., data features selection, construction, inconsistencies. How to debug bad performance in machine learning models and prediction?

**AI Technology:** Artificial Neural Networks.

**XAI Technology:** Model / Prediction comparison

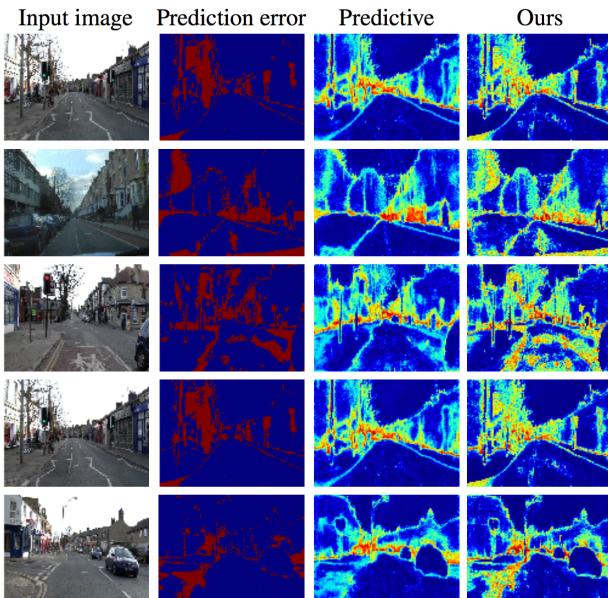


# Obstacle Identification Certification (Trust) - Transportation

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.



THALES



**Challenge:** Public transportation is getting more and more self-driving vehicles. Even if trains are getting more and more autonomous, the human stays in the loop for critical decision, for instance in case of obstacles. In case of obstacles trains are required to provide recommendation of action i.e., go on or go back to station. In such a case the human is required to validate the recommendation through an explanation exposed by the train or machine.

**AI Technology:** Integration of AI related technologies i.e., Machine Learning (Deep Learning / CNNs), and semantic segmentation.

**XAI Technology:** Deep learning and Epistemic uncertainty

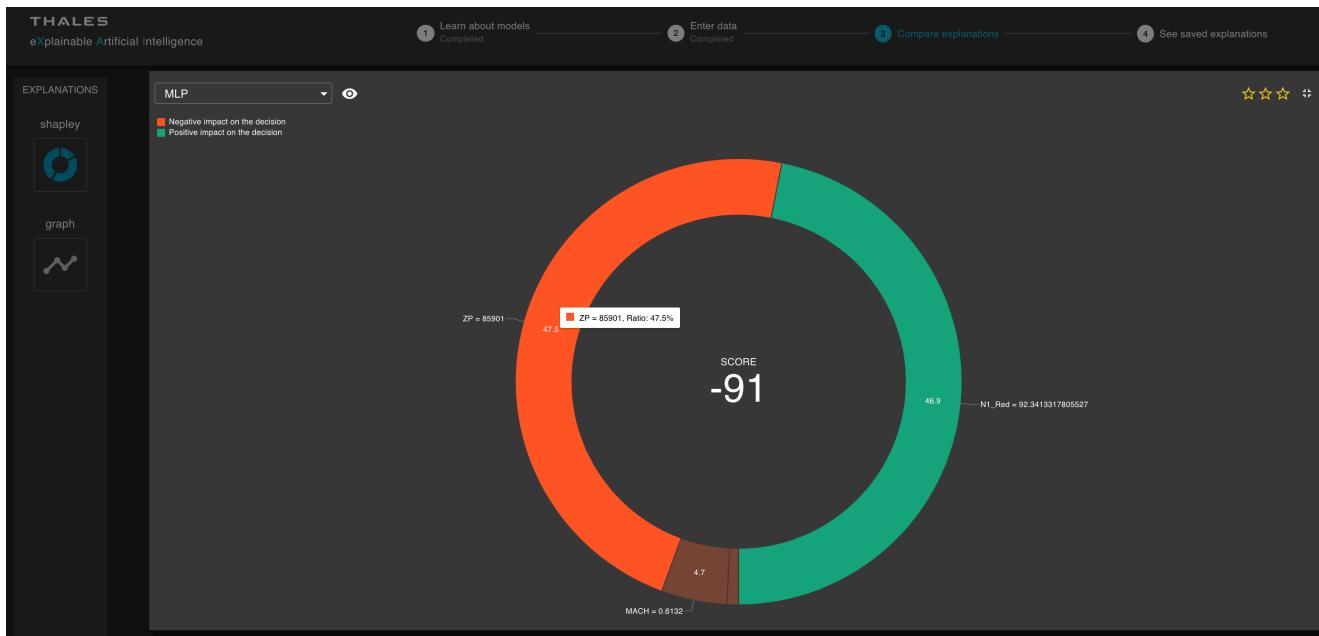


THALES

# Explaining Flight Performance- Transportation

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

THALES



**Challenge:** Predicting and explaining aircraft engine performance

**AI Technology:** Artificial Neural Networks

**XAI Technology:** Shapely Values

THALES

# Explainable On-Time Performance - Transportation

## KLM / Transavia Flight Delay Prediction

PLANE INFO		ARRIVAL			TURNAROUND			DEPARTURE				
Status / Aircraft	Flight	ETA	Status	Delay Code	Gate	Slot	Progress	Milestones	Flight	ETA	Status	Delay Code
✓ <a href="#">urtwet</a> ✓	4567	18:30	Scheduled	-	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	-
⚠ <a href="#">idsfew</a> ✓	4567	18:30	Delayed	ABC, DEF, GHI	345345	1	<div style="width: 0%; background-color: red;"></div>		5678	19:00	Delayed	ABC, DEF, GHI
✓ <a href="#">passidb</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
🚫 <a href="#">kshdbs</a> ✓	4567	-	Cancelled	ABC, DEF, GHI	-	-	<div style="width: 0%; background-color: grey;"></div>		5678	-	Cancelled	ABC, DEF, GHI
⚠ <a href="#">woowdfs</a> ✓	4567	18:35	Delayed	ABC, DEF, GHI	345345	1	<div style="width: 25%; background-color: yellow;"></div>		5678	19:00	Delayed	ABC, DEF, GHI
⚠ <a href="#">pdigts</a> ✓	4567	18:30	Delayed	ABC, DEF, GHI	345345	1	<div style="width: 0%; background-color: red;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI
✓ <a href="#">aedbsc</a> ✓	4567	18:30	Scheduled	ABC, DEF, GHI	345345	1	<div style="width: 50%; background-color: green;"></div>		5678	19:00	Scheduled	ABC, DEF, GHI

Jiaoyan Chen, Freddy Lécué, Jeff Z. Pan, Ian Horrocks, Huajun Chen: Knowledge-Based Transfer Learning Explanation. KR 2018: 349-358

Nicholas McCarthy, Mohammad Karzand, Freddy Lecue: Amsterdam to Dublin Eventually Delayed? LSTM and Transfer Learning for Predicting Delays of Low Cost Airlines: AAAI 2019

**Challenge:** Globally 323,454 flights are delayed every year. Airline-caused delays totaled 20.2 million minutes last year, generating huge cost for the company. Existing in-house technique reaches 53% accuracy for **predicting flight delay**, does not provide any time estimation (in **minutes** as opposed to True/False) and is unable to capture the underlying reasons (explanation).

**AI Technology:** Integration of AI related technologies i.e., Machine Learning (Deep Learning / Recurrent neural Network), Reasoning (through semantics-augmented case-based reasoning) and Natural Language Processing for building a robust model which can (1) predict flight delays in minutes, (2) explain delays by comparing with historical cases.

**XAI Technology:** Knowledge graph embedded Sequence Learning using LSTMs

THALES

INNOVATION ARCHITECTURE:  
**ACCENTURE**  
LABS

THALES

# Model Explanation for Sales Prediction - Sales

① What are top driver features for a certain company to have high/low probability to upsell/churn?

① Feature Contributor

② Which top driver features can be perturbed if we want to increase/decrease probability for a certain company?

② Feature Influencer



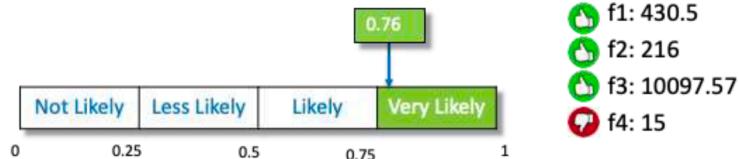
**Challenge:** How to predict and explain upsell / churn for a company?

**AI Technology:** Artificial Neural Networks.

**XAI Technology:** Features importance (contribution, influence), LIME.

Company: CompanyX

Upsell LCP (LinkedIn Career Page)



## Top Feature Contributor

- f1: 430.5
- f2: 216
- f3: 10097.57
- f4: 15

## Top Feature Influencer (Positive)

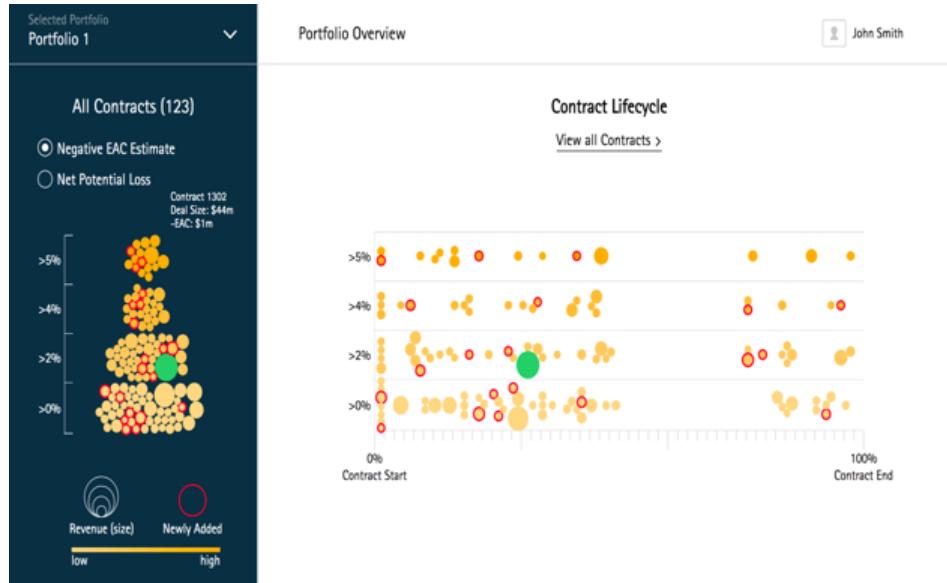
- f5: 0 → 5.4, 0.03
- f6: 168 → 0, 0.03
- f7: 0 → 0.24, 0.02

## Top Feature Influencer (Negative)

- f1: 430.5 → 148.7, -0.20
- f2: 216 → 0, -0.17
- f8: 423 → 146.0, -0.07

# Explainable Risk Management - Finance

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.



**Challenge:** Accenture is managing every year more than 80,000 opportunities and 35,000 contracts with an expected revenue of \$34.1 billion. Revenue expectation does not meet estimation due to the complexity and risks of critical contracts. This is, in part, due to the (1) large volume of projects to assess and control, and (2) the existing non-systematic assessment process.

**AI Technology:** Integration of AI technologies i.e., Machine Learning, Reasoning, Natural Language Processing for building a robust model which can (1) predict revenue loss, (2) recommend corrective actions, and (3) explain why such actions might have a positive impact.

**XAI Technology:** Knowledge graph embedded Random Forrest

Jiewen Wu, Freddy Lécué, Christophe Guéret, Jer Hayes, Sara van de Moosdijk, Gemma Gallagher, Peter McCanney, Eugene Eichelberger: Personalizing Actions in Context for Risk Management Using Semantic Web Technologies. International Semantic Web Conference (2) 2017: 367-383

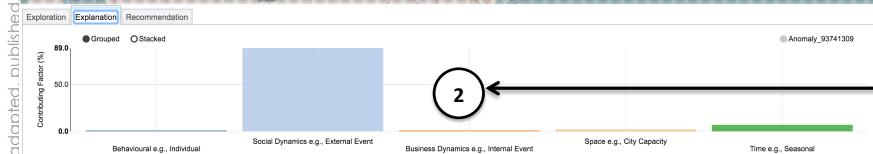
# Explainable Anomaly Detection – Finance (Compliance)

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part.

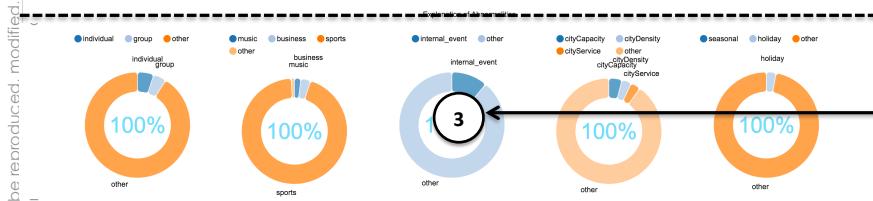
92



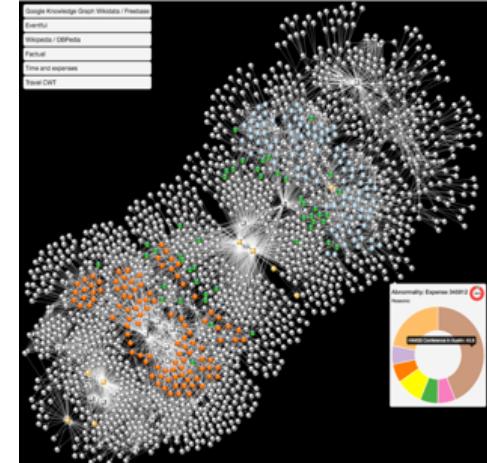
Data analysis for spatial interpretation of abnormalities: abnormal expenses



Semantic explanation (structured in classes: fraud, events, seasonal) of abnormalities



Detailed semantic explanation (structured in sub classes e.g. categories for events)



Freddy Lécué, Jiewen Wu: Explaining and predicting abnormal expenses at large scale using knowledge graph based reasoning. J. Web Sem. 44: 89-103 (2017)

**Challenge:** Predicting and explaining abnormally employee expenses (as high accommodation price in 1000+ cities).

**AI Technology:** Various techniques have been matured over the last two decades to achieve excellent results. However most methods address the problem from a statistic and pure data-centric angle, which in turn limit any interpretation. We elaborated a web application running live with real data from (i) travel and expenses from Accenture, (ii) external data from third party such as Google Knowledge Graph, DBpedia (relational DataBase version of Wikipedia) and social events from Eventful, for explaining abnormalities.

**THALES**

# Counterfactual Explanations for Credit Decisions (1) - Finance

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

**Local, post-hoc, contrastive explanations of black-box classifiers**

**Required minimum change in input vector to flip the decision of the classifier.**

**Interactive Contrastive Explanations** THALES

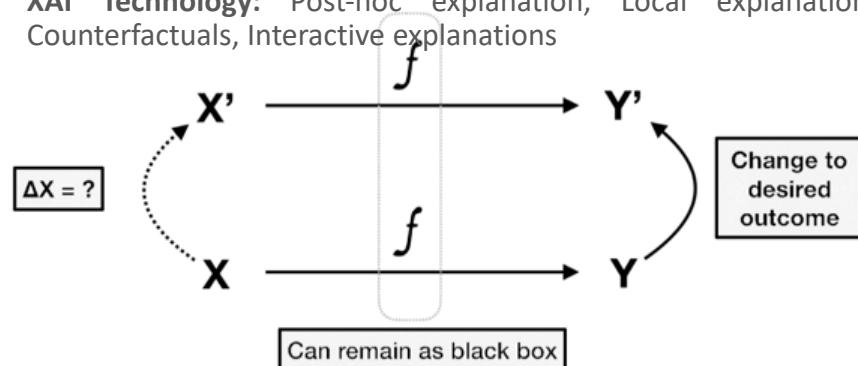


Rory McGrath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, Freddy Lécué: Interpretable Credit Application Predictions With Counterfactual Explanations. FEAP-AI4fin workshop, NeurIPS, 2018.

**Challenge:** We predict loan applications with off-the-shelf, interchangeable black-box estimators, and we explain their predictions with counterfactual explanations. In counterfactual explanations the model itself remains a black box; it is only through changing inputs and outputs that an explanation is obtained.

**AI Technology:** Supervised learning, binary classification.

**XAI Technology:** Post-hoc explanation, Local explanation, Counterfactuals, Interactive explanations



# Counterfactual Explanations for Credit Decisions (2) - Finance



## Sorry, your loan application has been rejected.

Our analysis:

The following features **were too high**:

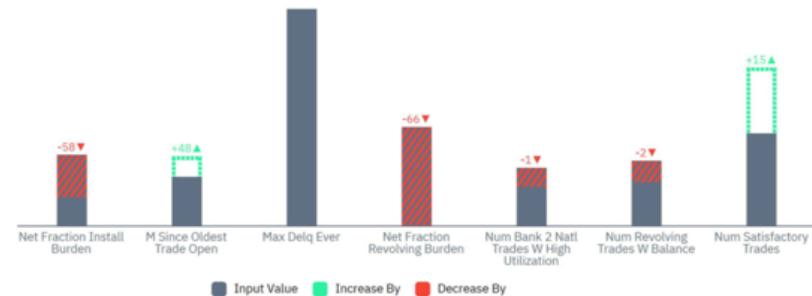
PercentInstallTrad... NetFractionRevolv... NetFractionInstall...  
NumRevolvingTra... NumBank2NatlTra... PercentTradesWB...

The following features **were too low**:

MSinceOldestTrad... AverageMInFile NumTotalTrades

The following features **require changes**:

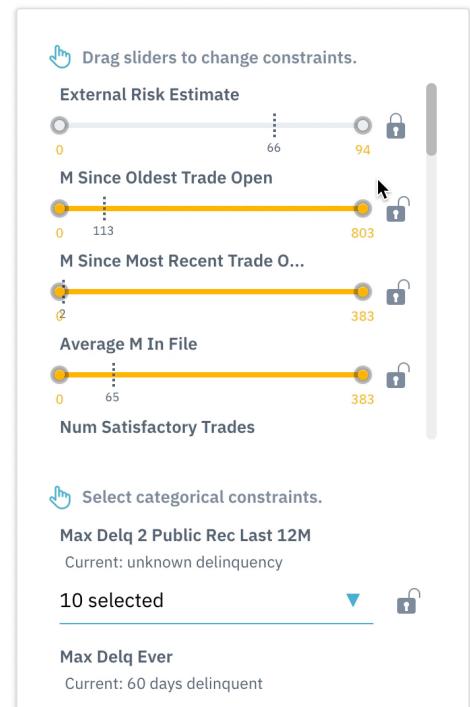
MaxDelq2PublicR... MaxDelqEver



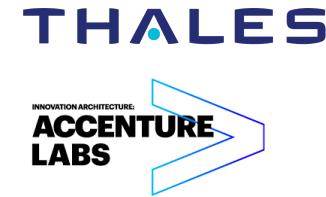
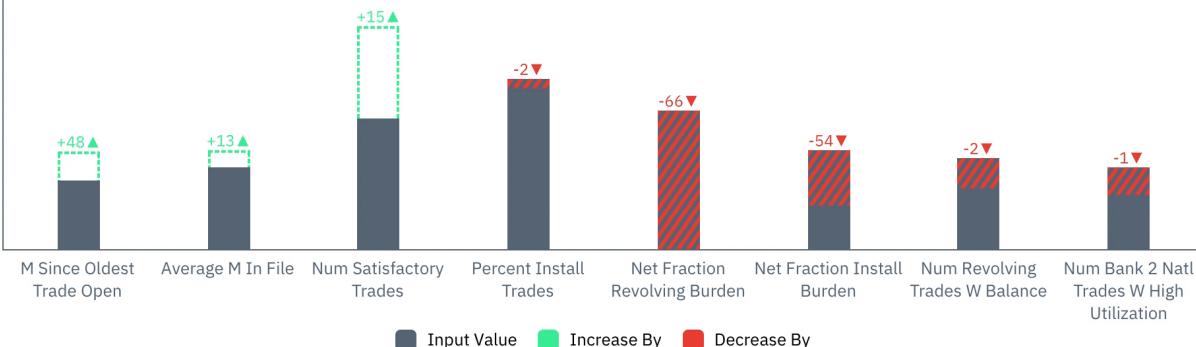
Counterfactuals suggest where to increase (green, dashed) or decrease (red, striped) each feature.

# Counterfactual Explanations for Credit Decisions (3) - Finance

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed.



## RECOMMENDED CHANGES



# Explaining Talent Search Results – Human Resources



**Challenge:** How to rationalize a talent search for a recruiter when looking for candidates for a given role. Features are dynamic and costly to compute. Recruiters are interested in discriminating between two candidates to make a selection.

**AI Technology:** Generalized Linear Mixed Models, Artificial Neural Networks, XGBoost

**XAI Technology:** Generalized Linear Mixed Models (inherently explainable), Integrated Gradient, Features Importance in XGBoost

Feature	Description	Difference (1 vs 2)	Contribution
Feature.....	Description.....	-2.0476928	-2.144455602
Feature.....	Description.....	-2.3223877	1.903594618
Feature.....	Description.....	0.11666667	0.2114946752
Feature.....	Description.....	-2.1442587	0.2060414469
Feature.....	Description.....	-14	0.1215354111
Feature.....	Description.....	1	0.1000282466
Feature.....	Description.....	-92	-0.085286277
Feature.....	Description.....	0.9333333	0.0568533262
Feature.....	Description.....	-1	-0.051796317
Feature.....	Description.....	-1	-0.050895940

# Explanation of Medical Condition Relapse – Health

THALES



**Challenge:** Explaining medical condition relapse in the context of oncology.

**AI Technology:** Relational learning

**XAI Technology:** Knowledge graphs and Artificial Neural Networks

Knowledge  
graph parts  
explaining  
medical  
condition relapse

THALES

# Breast Cancer Survival Rate Prediction - Health



modified, adapted, published, translated, in any way, in whole or in part. All rights reserved.

Age at diagnosis	<input type="button" value="–"/> <input type="text" value="69"/> <input type="button" value="+"/> Age must be between 25 and 85
Post Menopausal?	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Unknown
ER status	<input type="radio"/> Positive <input type="radio"/> Negative
HER2 status	<input type="radio"/> Positive <input type="radio"/> Negative <input type="radio"/> Unknown
Ki-67 status	<input type="radio"/> Positive <input type="radio"/> Negative <input type="radio"/> Unknown Positive means more than 10%
Tumour size (mm)	<input type="button" value="–"/> <input type="text" value="7"/> <input type="button" value="+"/>
Tumour grade	<input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3
Detected by	<input type="radio"/> Screening <input type="radio"/> Symptoms <input type="radio"/> Unknown
Positive nodes	<input type="button" value="–"/> <input type="text" value="2"/> <input type="button" value="+"/>
Micrometastases	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Unknown Enabled when positive nodes is zero

## Results

Table Curves Chart Texts Icons  
New recording

These results are for women who have already had surgery. This table shows the percentage of women who survive at least    years after surgery, based on the information you have provided.

Treatment	Additional Benefit	Overall Survival %
Surgery only	-	72%
+ Hormone therapy	0%	72%

If death from breast cancer were excluded, 82% would survive at least 10 years.

Show ranges?  Yes  No

David Spiegelhalter, Making Algorithms trustworthy, NeurIPS 2018 Keynote  
[predict.nhs.uk/tool](http://predict.nhs.uk/tool)

**Challenge:** Predict is an online tool that helps patients and clinicians see how different treatments for early invasive breast cancer might improve survival rates after surgery.

**AI Technology:** competing risk analysis

**XAI Technology:** Interactive explanations, Multiple representations.

# More on XAI

# (Some) Tutorials, Workshops, Challenge

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

## Tutorial:

- | AAAI 2019 Tutorial on On Explainable AI: From Theory to Motivation, Applications and Limitations (#1) - <https://xaitutorial2019.github.io/>
- | ICIP 2018 / EMBC 2019 Interpretable Deep Learning: Towards Understanding & Explaining Deep Neural Networks (#2) - <http://interpretable-ml.org/icip2018tutorial/> - <http://interpretable-ml.org/embc2019tutorial/>
- | ICCV 2019 Tutorial on Interpretable Machine Learning for Computer Vision (#2) - <https://interpretablevision.github.io/>

## Workshop:

- | ISWC 2019 Workshop on Semantic Explainability (#1) - <http://www.semantic-explainability.com/>
- | IJCAI 2019 Workshop on Explainable Artificial Intelligence (#3) - <https://sites.google.com/view/xai2019/home> 55 paper submitted in 2019
- | IJCAI 2019 Workshop on Optimisation and Explanation in AI (#1) - <https://www.doc.ic.ac.uk/~kc2813/OXAI/>
- | SIGIR 2019 Workshop on Explainable Recommendation and Search (#2) <https://ears2019.github.io/>
- | ICAPS 2019 Workshop on Explainable Planning (#2)- [https://kcl-planning.github.io/XAIP-Workshops/ICAPS\\_2019](https://kcl-planning.github.io/XAIP-Workshops/ICAPS_2019) 23 papers submitted in 2019  
<https://openreview.net/group?id=icaps-conference.org/ICAPS/2019/Workshop/XAIP>
- | ICCV 2019 Workshop on Interpreting and Explaining Visual Artificial Intelligence Models (#1) - <http://xai.unist.ac.kr/workshop/2019/>
- | NeurIPS 2019 Workshop on Challenges and Opportunities for AI in Financial Services: the Impact of Fairness, Explainability, Accuracy, and Privacy - <https://sites.google.com/view/feap-ai4fin-2018/>
- | CD-MAKE 2019 – Workshop on Explainable AI (#2) - <https://cd-make.net/special-sessions/make-explainable-ai/>
- | AAAI 2019 / CVPR 2019 Workshop on Network Interpretability for Deep Learning (#1 and #2) - <http://networkinterpretability.org/> - <https://explainai.net/>

## Challenge:

- | 2018: FICO Explainable Machine Learning Challenge (#1) - <https://community.fico.com/s/explainable-machine-learning-challenge>

# (Some) Software Resources

- | DeepExplain: perturbation and gradient-based attribution methods for Deep Neural Networks interpretability. [github.com/marcoancona/DeepExplain](https://github.com/marcoancona/DeepExplain)
- | iNNvestigate: A toolbox to iNNvestigate neural networks' predictions. [github.com/albermax/innvestigate](https://github.com/albermax/innvestigate)
- | SHAP: SHapley Additive exPlanations. [github.com/slundberg/shap](https://github.com/slundberg/shap)
- | Microsoft Explainable Boosting Machines. <https://github.com/Microsoft/interpret>
- | GANDissect: Pytorch-based tools for visualizing and understanding the neurons of a GAN. <https://github.com/CSAILVision/GANDissect>
- | ELI5: A library for debugging/inspecting machine learning classifiers and explaining their predictions. [github.com/TeamHG-Memex/eli5](https://github.com/TeamHG-Memex/eli5)
- | Skater: Python Library for Model Interpretation/Explanations. [github.com/datascienceinc/Skater](https://github.com/datascienceinc/Skater)
- | Yellowbrick: Visual analysis and diagnostic tools to facilitate machine learning model selection. [github.com/DistrictDataLabs/yellowbrick](https://github.com/DistrictDataLabs/yellowbrick)
- | Lucid: A collection of infrastructure and tools for research in neural network interpretability. [github.com/tensorflow/lucid](https://github.com/tensorflow/lucid)
- | LIME: Agnostic Model Explainer. <https://github.com/marcotcr/lime>
- | Sklearn\_explain: model individual score explanation for an already trained scikit-learn model. [https://github.com/antoinecarme/sklearn\\_explain](https://github.com/antoinecarme/sklearn_explain)
- | Heatmapping: Prediction decomposition in terms of contributions of individual input variables
- | Deep Learning Investigator: Investigation of Saliency, Deconvnet, GuidedBackprop and more. <https://github.com/albermax/innvestigate>
- | Google PAIR What-if: Model comparison, counterfactual, individual similarity. <https://pair-code.github.io/what-if-tool/>
- | Google tf-explain: <https://tf-explain.readthedocs.io/en/latest/>
- | IBM AI Fairness: Set of fairness metrics for datasets and ML models, explanations for these metrics. <https://github.com/IBM/aif360>
- | Blackbox auditing: Auditing Black-box Models for Indirect Influence. <https://github.com/algofairness/BlackBoxAuditing>
- | Model describer: Basic statistical metrics for explanation (visualisation for error, sensitivity). <https://github.com/DataScienceSquad/model-describer>
- | AXA Interpretability and Robustness: <https://axa-rev-research.github.io/> (more on research resources – not much about tools)

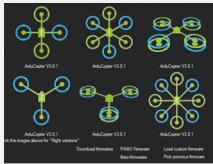
# (Some) Initiatives: XAI in USA



## Challenge Problem Areas



### Data Analytics Multimedia Data

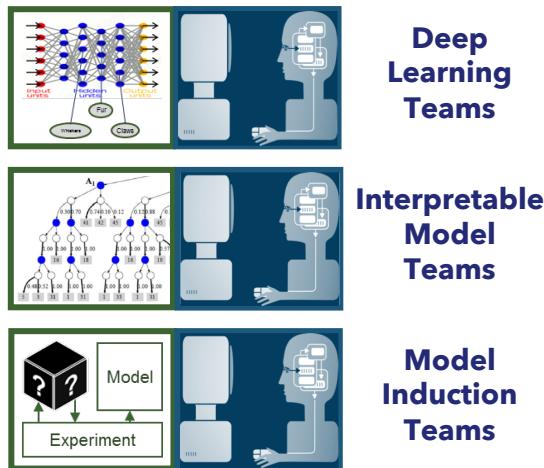


### Autonomy ArduPilot & SITL Simulation

## TA 1: Explainable Learners

Teams that provide prototype systems with both components:

- Explainable Model
- Explanation Interface



### Deep Learning Teams

### Interpretable Model Teams

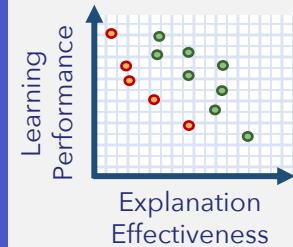
### Model Induction Teams

## TA 2: Psychological Model of Explanation



- Psych. Theory of Explanation
- Computational Model
- Consulting

## Evaluation Framework



### Explanation Measures

- User Satisfaction
- Mental Model
- Task Performance
- Trust Assessment
- Correctability

## Evaluator

### TA1: Explainable Learners

- Explainable learning systems that include both an explainable model and an explanation interface

### TA2: Psychological Model of Explanation

- Psychological theories of explanation and develop a computational model of explanation from those theories

# (Some) Initiatives: XAI in Canada

This document may not be reproduced, modified, adapted, published, translated, in any way, in whole or in part or disclosed to a third party without the prior written consent of Thales - © Thales 2018 All rights reserved.

## | DEEL (Dependable Explainable Learning) Project 2019-2024

### > Research institutions



### > Industrial partners



### > Academic partners

- Science and technology to develop new methods towards Trustable and Explainable AI



#### System Robustness

- To biased data
- Of algorithm
- To change
- To attacks

#### Certificability

- Structural warranties
- Risk auto evaluation
- External audit

#### Explicability & Interpretability

#### Privacy by design

- Differential privacy
- Homomorphic coding
- Collaborative learning
- To attacks

# (Some) Initiatives: XAI in EU



- | Explainable AI is motivated by real-world applications in AI
- | Not a new problem – a reformulation of past research challenges in AI
- | Multi-disciplinary: multiple AI fields, HCI, social sciences (multiple definitions)
- | In AI (in general): many interesting / complementary approaches
- | Many industrial applications already – crucial for AI adoption in critical systems

# Future Challenges

- | Creating awareness! Success stories!
  - | Foster multi-disciplinary collaborations in XAI research.
  - | Help shaping industry standards, legislation.
  - | More work on transparent design.
  - | Investigate symbolic and sub-symbolic reasoning.
- 
- | **Evaluation:**
    - > We need benchmark - Shall we start a task force?
    - > We need an XAI challenge - Anyone interested?
    - > Rigorous, agreed upon, human-based evaluation protocols

Wherever safety and Security are Critical, Thales can build smarter solutions. Everywhere.

# Job Openings

Technology leader for the Defence & Security, the combined expertise can have made Thales a key player in keeping the public protecting the national security interests of count

Established in 1972, Thales Canada has over 1,800 employees in Toronto and Vancouver working in Defence, Avionics and Space.

This is a unique opportunity to play a key role on the Research and Technology (RTT) in Canada (Quebec and Montreal). We are looking for applied R&T experts at five locations worldwide. We are looking for individuals who are interested in developing cutting edge AI technologies. Our passion is imagining and creating new technologies. Not only will you join our network, but this RTT is also co-located within Cognitec (Intelligence eXpertise) i.e., the new flagship program to work.

## Job Description

An AI (Artificial Intelligence) Research and Technology (RTT) scientist will be responsible for developing innovative prototypes to demonstrate the potential of artificial intelligence. To be successful in this role, one must be able to work independently, be creative, and have a strong ability to learn new technical skills and be familiar with latest developments in the field. The individual will contribute as technical subject matter experts to various projects and its business units. In addition to the implementation of AI solutions, the individual will also be involved in the initial project planning, design, and team work is also critical for this role.

As a Research and Technology Applied AI Scientist, you will be involved in the development of AI solutions for various projects, including the design, implementation, and testing of AI models.

## Professional Skill Requirements

- Good foundation in mathematics, statistics, and computer science

- Strong knowledge of Machine Learning foundations
- Strong development skills with Machine Learning frameworks e.g., Scikit-learn, TensorFlow, PyTorch, Theano
- Knowledge of mainstream Deep Learning architectures (MLP, CNN, RNN, etc).
- Strong Python programming skills
- Working knowledge of Linux OS
- Eagerness to contribute in a team-oriented environment
- Demonstrated leadership abilities in school, civil or business organisations
- Ability to work creatively and analytically in a problem-solving environment
- Proven verbal and written communication skills in English (talks, presentations, publications, etc.)

## Basic Qualifications

- Master's degree in computer science, engineering or mathematics fields
- Prior experience in artificial intelligence, machine learning, natural language processing, or advanced analytics

## Preferred Qualifications

- Minimum 3 years of analytic experience Python with interest in artificial intelligence with working with structured and unstructured data (SQL, Cassandra, MongoDB, Hive, etc.)
- A track record of outstanding AI software development with Github (or similar) evidence
- Demonstrated abilities in designing large scale AI systems
- Demonstrated interest in Explainable AI and/or relational learning (circled in red)
- Work experience with programming languages such as C, C++, Java, scripting languages (Perl/Python/Ruby) or similar
- Hands-on experience with data visualization, analytics tools/languages
- Demonstrated teamwork and collaboration in professional settings
- Ability to establish credibility with clients and other team members

MAY 2ND, 2019

Freddy Lecue  
Chief AI Scientist, CortAIx, Thales, Montreal – Canada

@freddylecue

<https://tinyurl.com/freddylecue>

Freddy.lecue.e@thaledigital.io