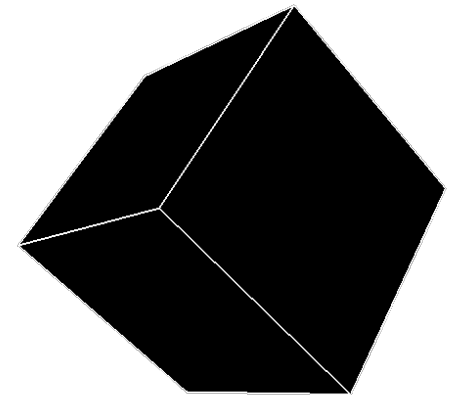


# On the Role of Knowledge Graphs in Explainable AI A Machine Learning Perspective

**Freddy Lécué**  
Inria, France  
CortAlx@Thales, Canada  
@freddylecue



**October 27<sup>th</sup>, 2019**  
**International Semantic Web Conference – Workshop on Semantic Explainability**  
**Auckland, New Zealand**

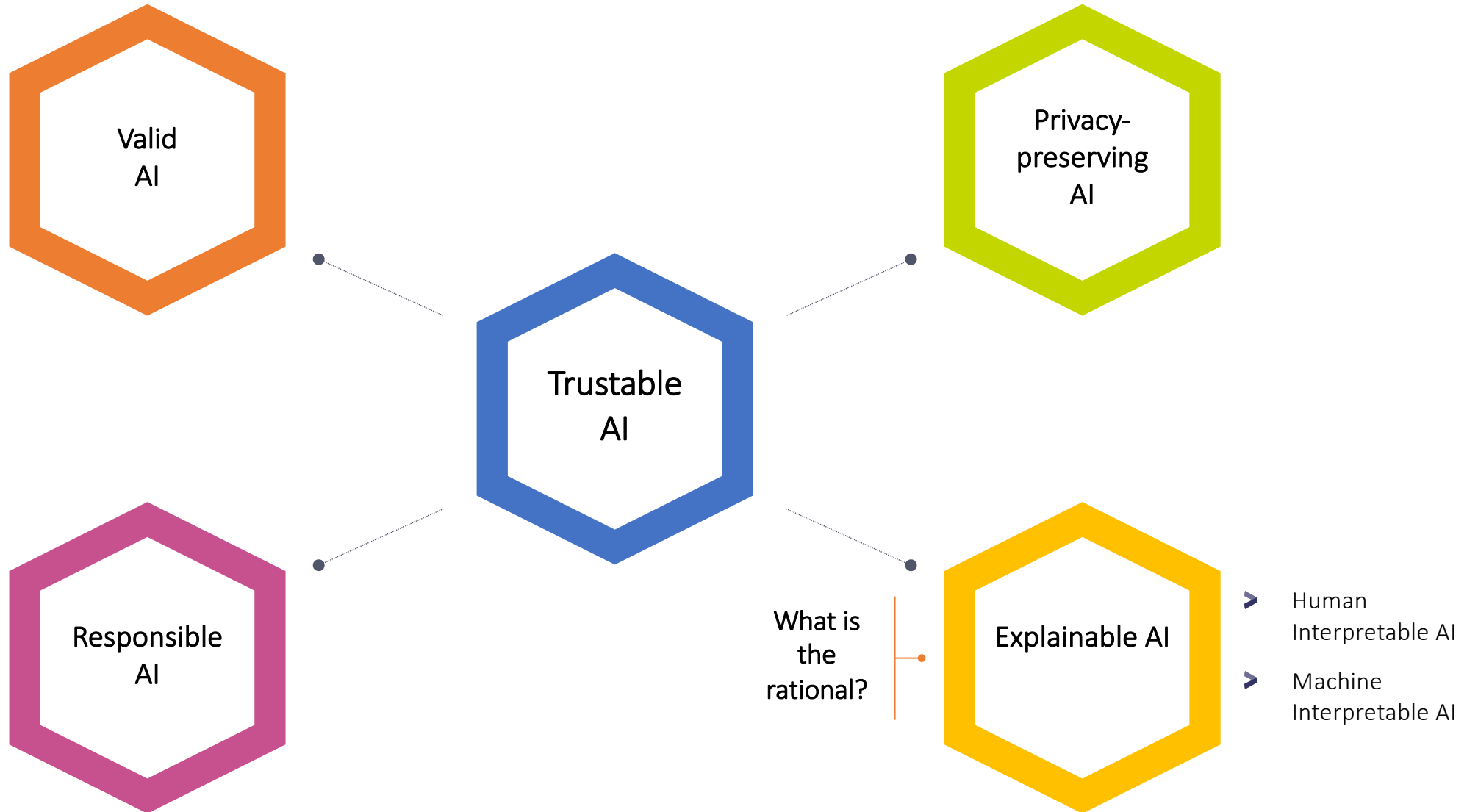


# Scope

---



# AI Adoption: Requirements



# 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 .

---

# Outline

---

# Outline

- **Explanation in Artificial Intelligence**
    - Motivation
    - Definitions
    - Evaluation (with role of the human in XAI systems)
    - The Role of Humans
    - Explanations in Different AI fields
  - **On the Role of Knowledge Graph in Explainable Machine Learning**
  - **XAI Industrial Applications using Knowledge Graphs on Machine Learning**
  - **Conclusion + Q&A**
-

# Motivation

---

# Business to Customer



Gary Chavez added a photo you might ...  
be in.  
about a minute ago · 👥



# Critical Systems

---











# Markets We Serve (Critical Systems)



Aerospace



Space



Ground Transportation



Defence



Security

**Trusted Partner** For A Safer World

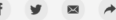
**But not Only  
Critical Systems**

---

# When a Computer Program Keeps You in Jail

By Rebecca Wexler

June 13, 2017



# COMPAS recidivism black bias

**DYLAN FUGETT**

**Prior Offense**  
1 attempted burglary

**Subsequent Offenses**  
3 drug possessions

**LOW RISK 3**

**BERNARD PARKER**

**Prior Offense**  
1 resisting arrest without violence

**Subsequent Offenses**  
None

**HIGH RISK 10**

*Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.*

# Motivation (2)

## Finance:

- Credit scoring, loan approval
- Insurance quotes



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

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

[Twitter](#) [Facebook](#) [LinkedIn](#) [Save](#)

Oliver Ralph MAY 16, 2017

[24](#)

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



# Motivation (3)

 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>

## 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.

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

Rich Caruana  
Microsoft Research  
rcaruana@microsoft.com

Yin Lou  
LinkedIn Corporation  
ylou@linkedin.com

Johannes Gehrke  
Microsoft  
johannes@microsoft.com

Paul Koch  
Microsoft Research  
paulkoch@microsoft.com

Marc Sturm  
NewYork-Presbyterian Hospital  
mas9161@nyp.org

Noémie Elhadad  
Columbia University  
noemie.elhadad@columbia.edu

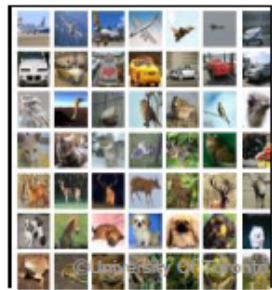
Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, Noemie Elhadad: Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission. KDD 2015: 1721-1730

# XAI in a Nutshell

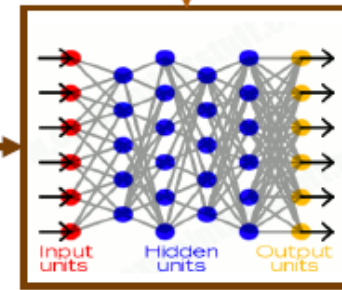
---

# XAI in a Nutshell

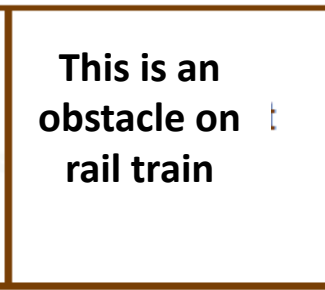
## Today



Training Data



Learned Function



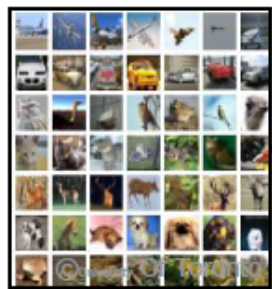
Output



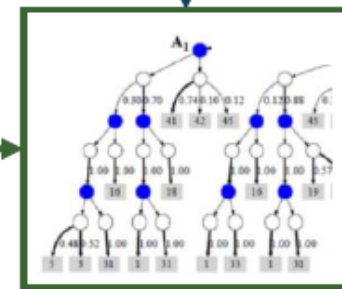
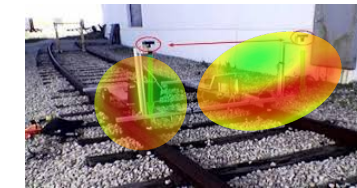
User with a Task

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

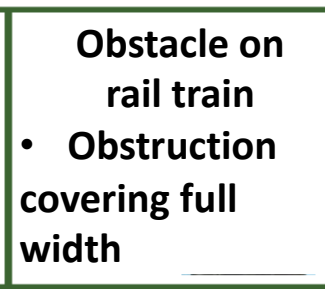
## Tomorrow



Training Data



Explainable Model



Explanation Interface

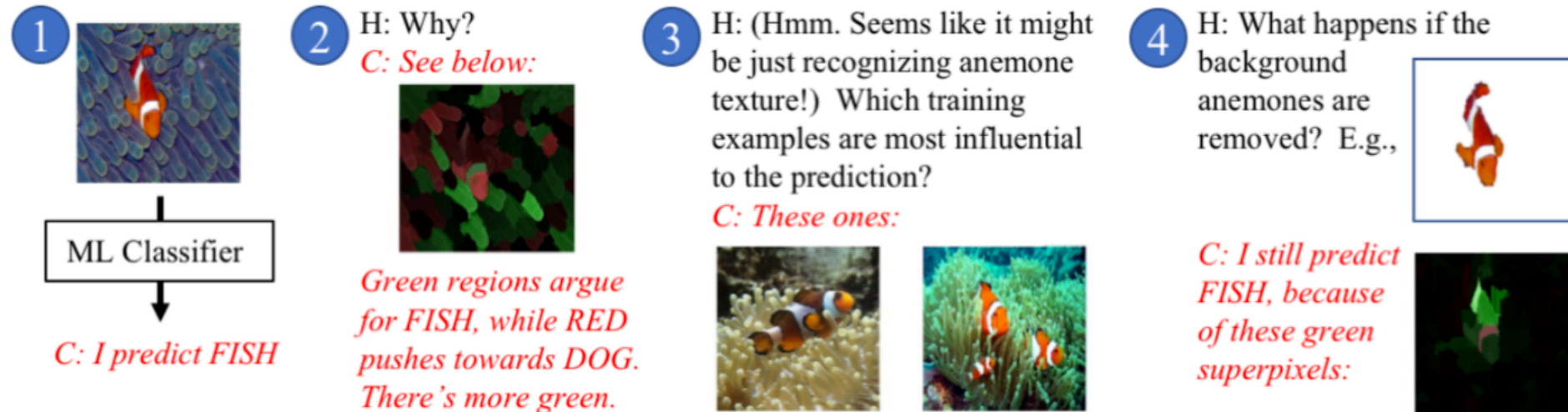


User with a Task

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



# An Example of an end-to-end XAI System

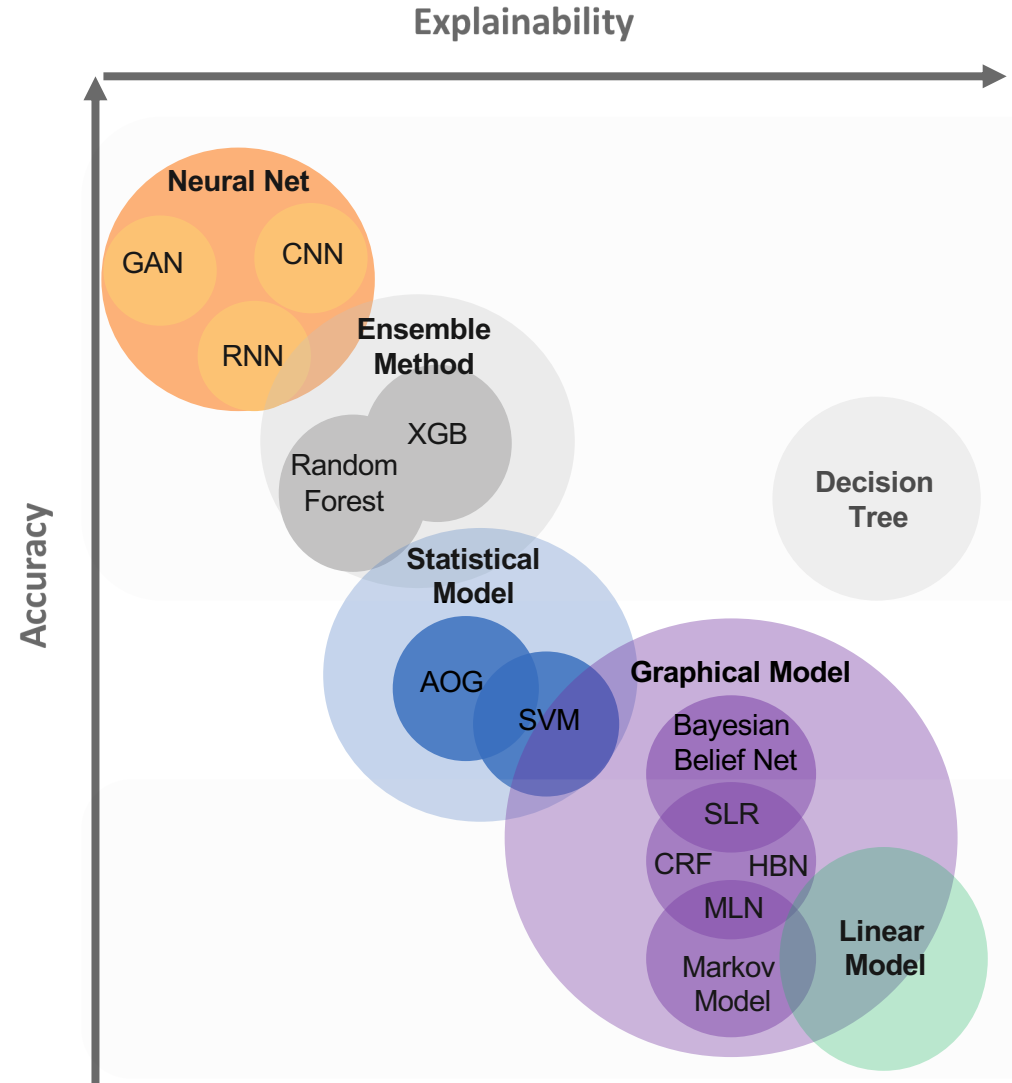


- Humans may have follow-up questions
- Explanations cannot answer all users' concerns

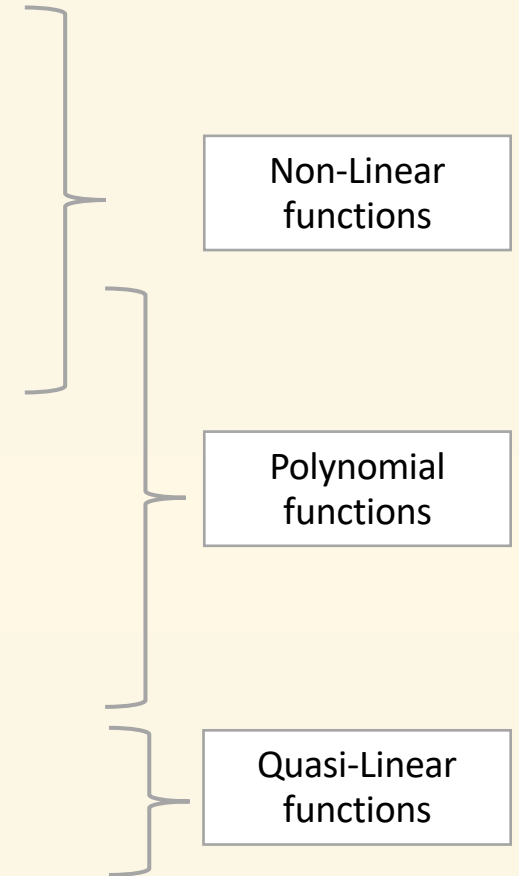
# How to Explain? Accuracy vs. Explanability

## Learning

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



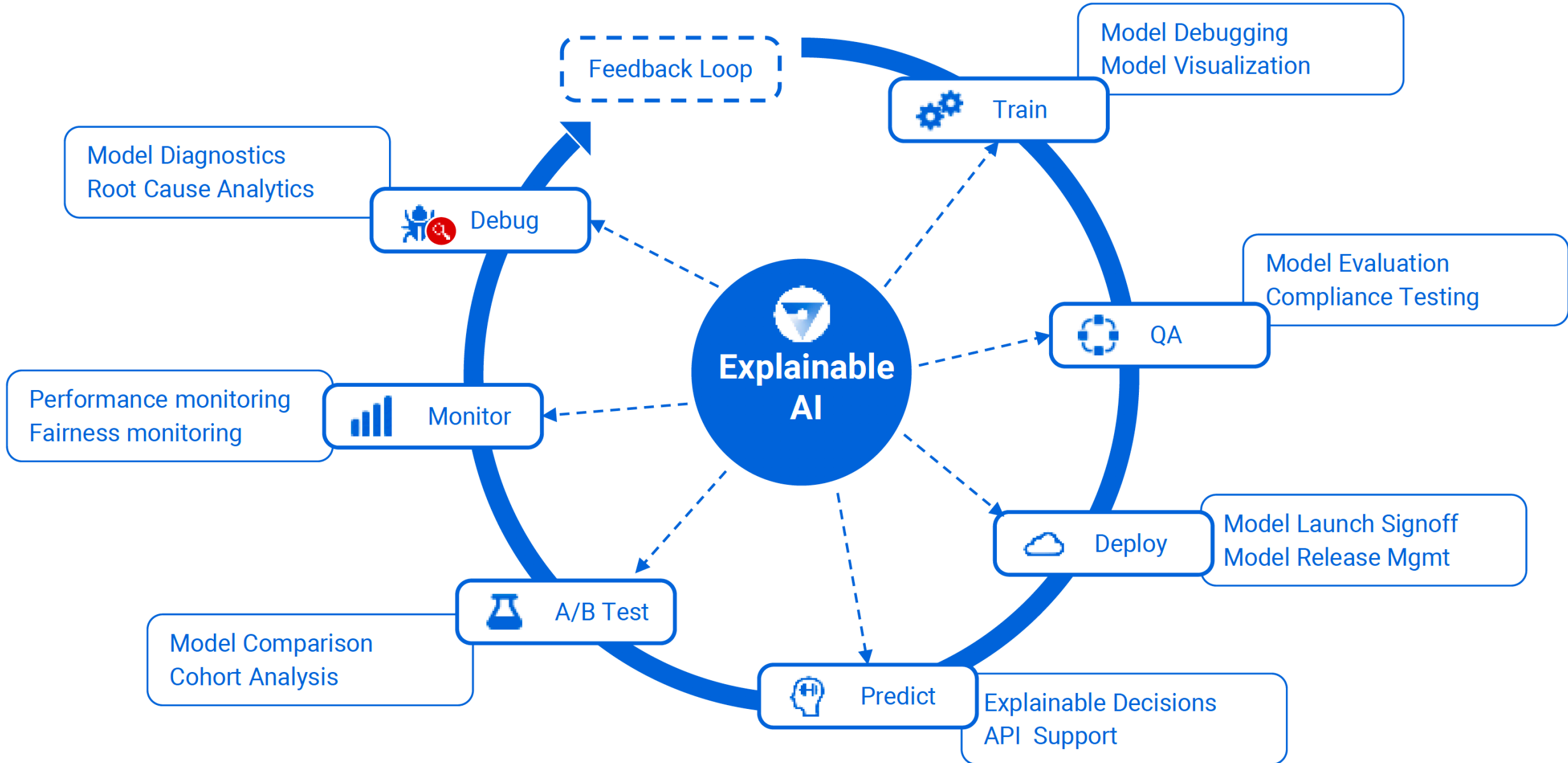
## Interpretability



**XAI Objective**  
**Supporting**  
**Industrialization of AI**  
**at Scale**

---

# Explainability by Design for AI Products



# XAI Definitions

**Explanation vs.  
Interpretability**

---

**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.*

**interpret** | ɪn'tɜːprɪt |

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

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

---

# On Role of Data In XAI

---

# Interpretable Data for Interpretable Models

Table of baby-name data  
(baby-2010.csv)

name	rank	gender	year
Jacob	1	boy	2010
Isabella	1	girl	2010
Ethan	2	boy	2010
Sophia	2	girl	2010
Michael	3	boy	2010

Field names

One row  
(4 fields)

2000 rows  
all told

Tabular

Images



Text



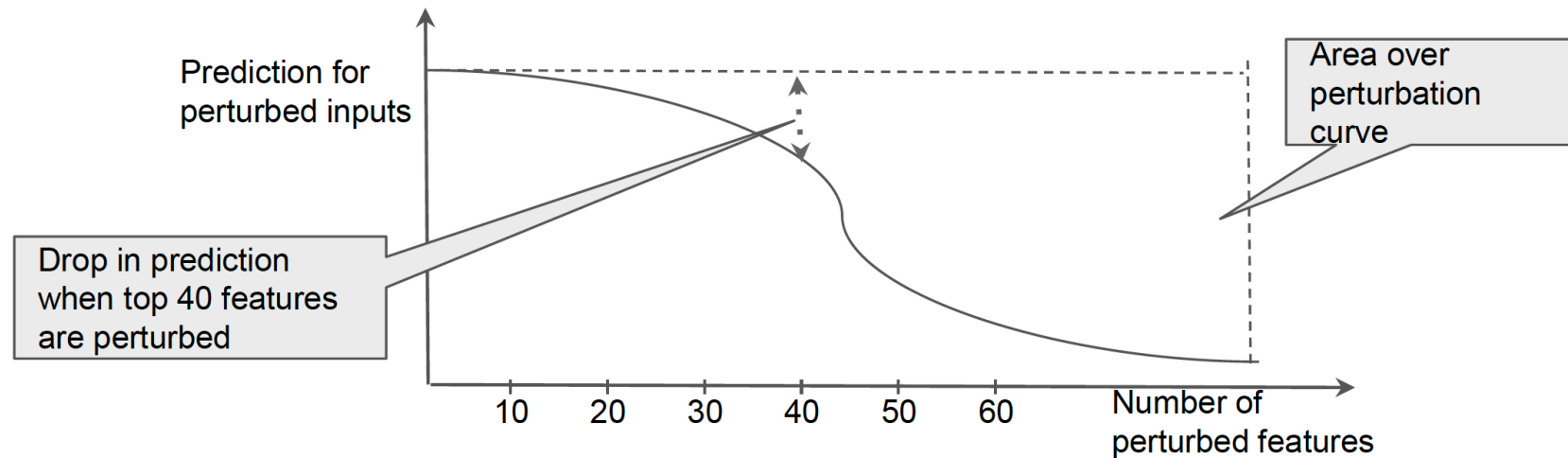
**What about the  
Evaluation?**

---

# Perturbation-based Evaluation for Feature Attribution-based Approaches

## Perturb top-k features by attribution and observe change in prediction

- Higher the change, better the method
- Perturbation may amount to replacing the feature with a random value
- Samek et al. formalize this using a metric: **Area over perturbation curve**
  - Plot the prediction for input with top-k features perturbed as a function of k
  - Take the area over this curve



# Human (Role)-based Evaluation is Essential... but too often based on size!

## **Evaluation criteria** for Explanations [Miller, 2017]

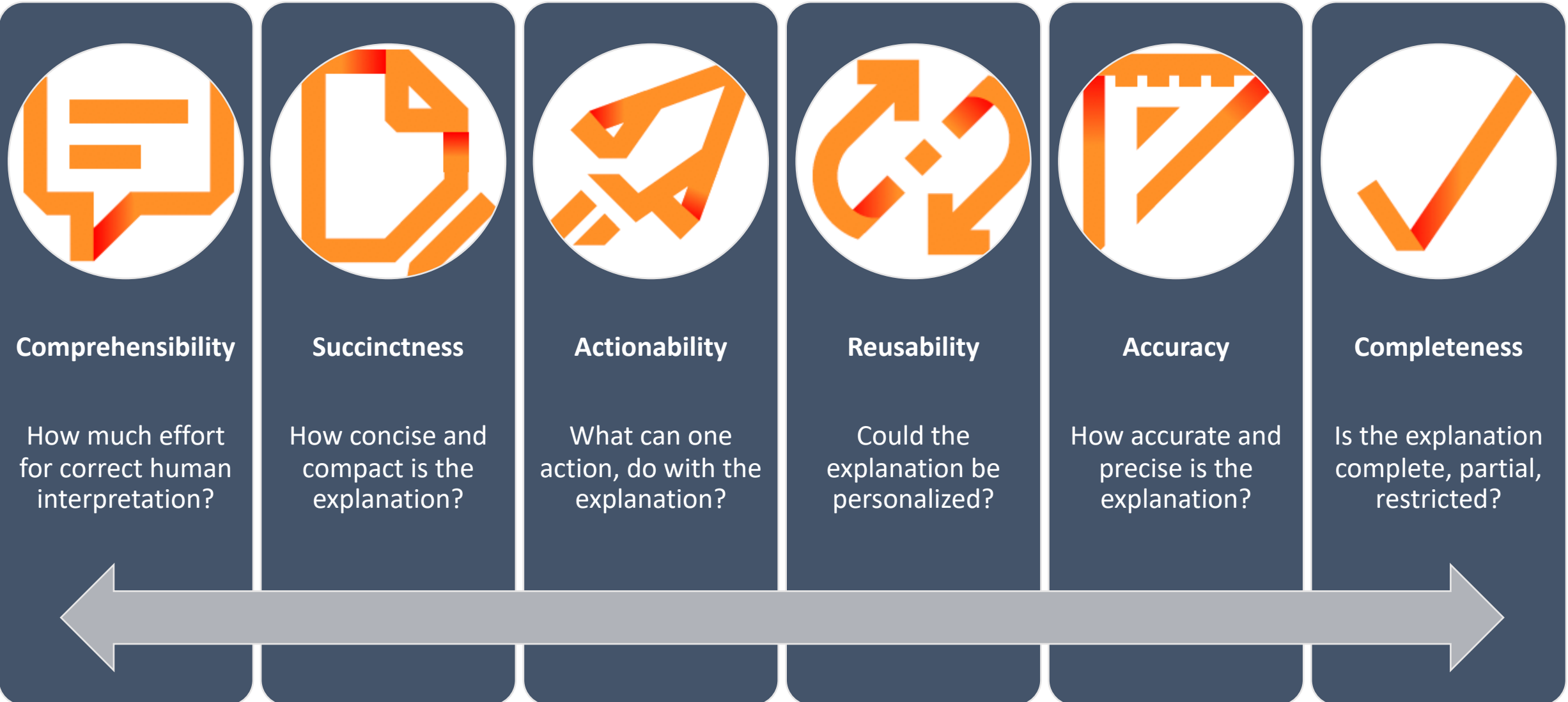
- Truth & probability
- Usefulness, relevance
- Coherence with prior belief
- Generalization

## **Cognitive chunks** = basic explanation units (for different explanation needs)

- Which basic units for explanations?
- How many?
- How to compose them?
- Uncertainty & end users?

[Doshi-Velez and Kim 2017, Poursabzi-Sangdeh 18]

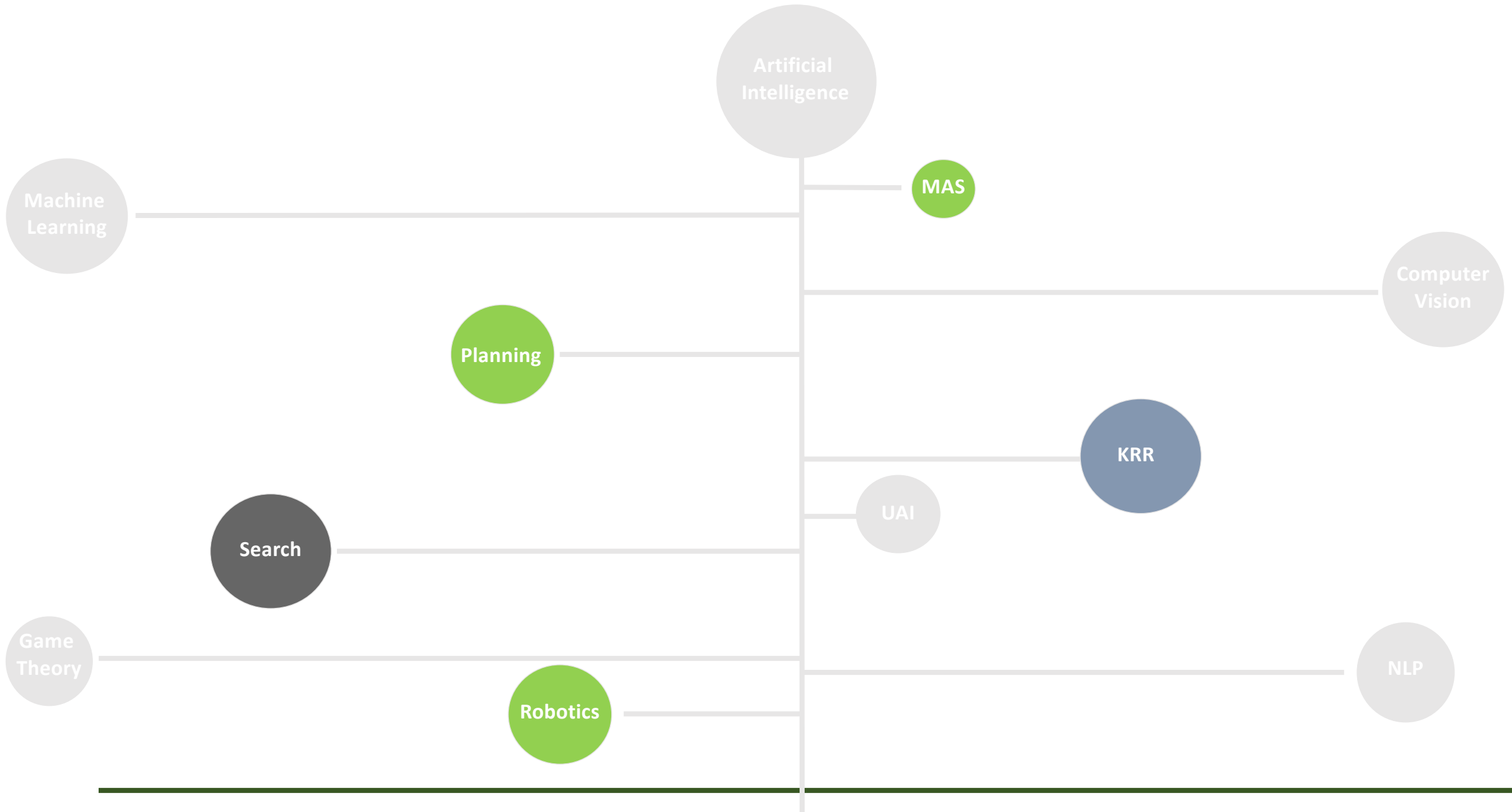
# XAI: One Objective, Many Metrics



**XAI in AI**

---

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



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

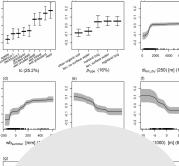
How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

Artificial Intelligence



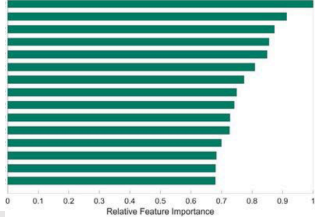
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Dependency Plot

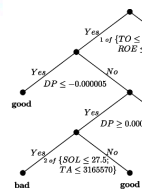


Machine Learning

Feature Importance



Surrogate Model



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

Artificial Intelligence

MAS

Which features are responsible of classification?

Computer Vision

Planning

KRR

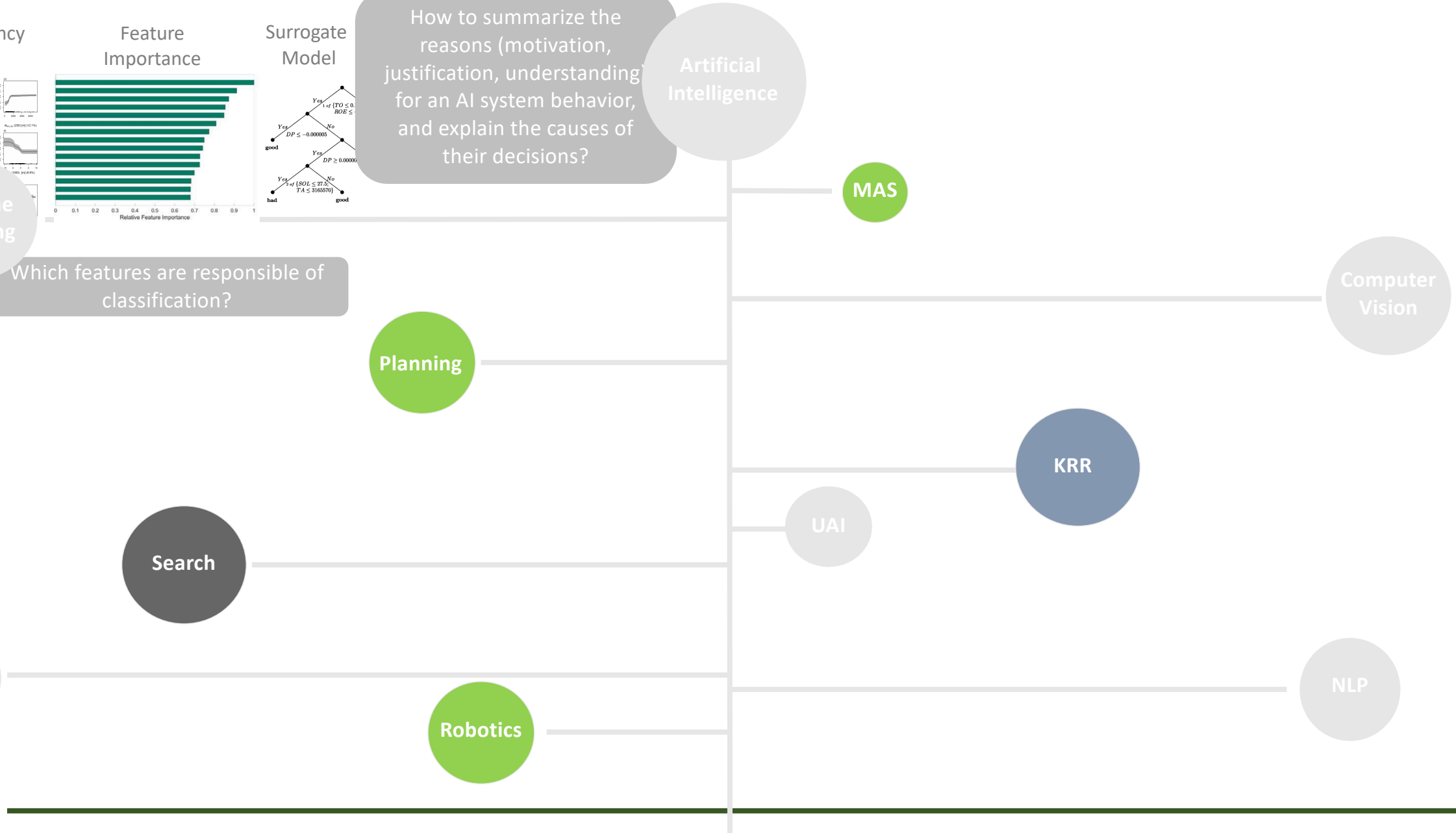
Search

UAI

Game Theory

NLP

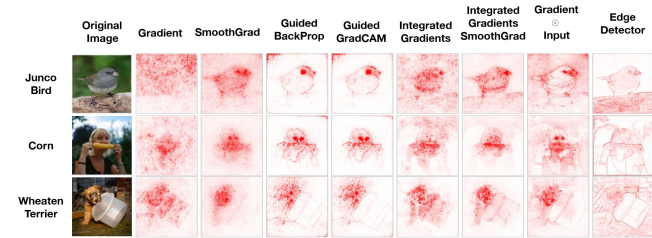
Robotics



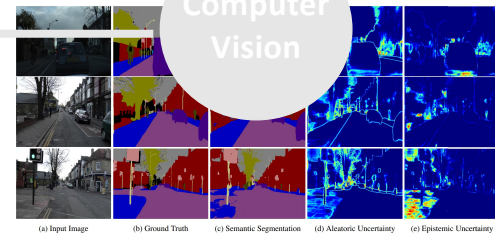


# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map



Which complex features are responsible of classification?



Uncertainty Map

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

Artificial Intelligence

MAS

Planning

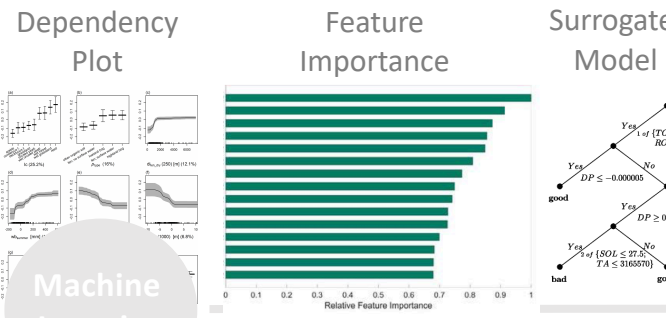
KRR

UAI

Search

NLP

Robotics



Machine Learning

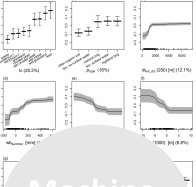
Which features are responsible of classification?

Game Theory

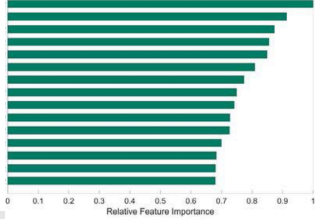
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

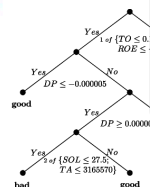
Dependency Plot



Feature Importance



Surrogate Model



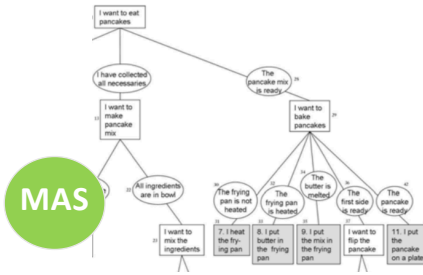
How to summarize the reasons (motivation, justification, understanding) for an AI system behavior, and explain the causes of their decisions?

Artificial Intelligence

Machine Learning

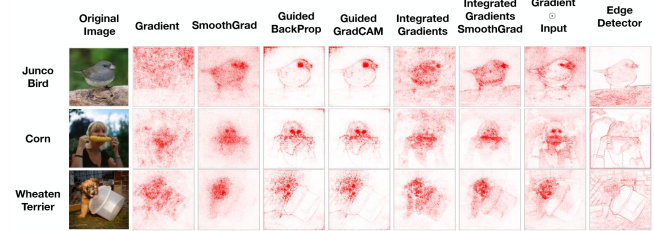
Which features are responsible of classification?

Strategy Summarization



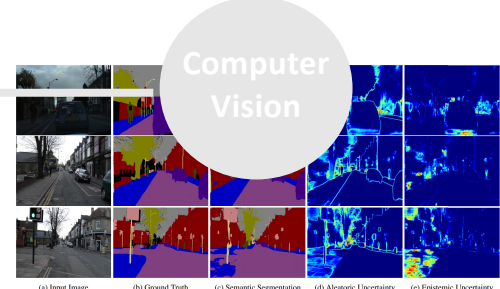
MAS

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



Which complex features are responsible of classification?

Planning



Computer Vision

(a) Input Image (b) Ground Truth (c) Semantic Segmentation (d) Aleatoric Uncertainty (e) Epistemic Uncertainty

Uncertainty Map

Search

KRR

UAI

Game Theory

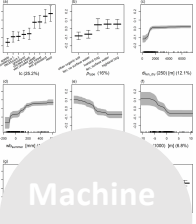
NLP

Robotics

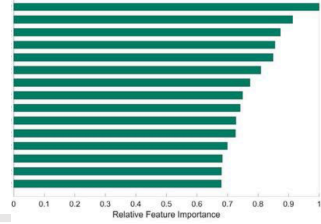
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

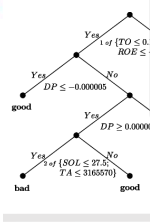
Dependency Plot



Feature Importance



Surrogate Model



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

Artificial Intelligence

Machine Learning

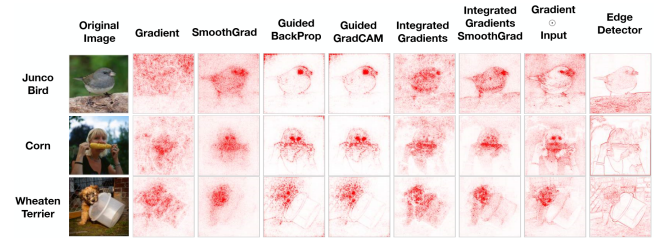
Which features are responsible of classification?

MAS

Strategy Summarization



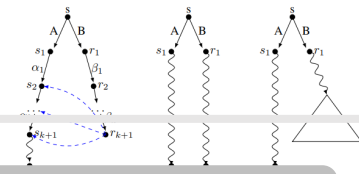
- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?



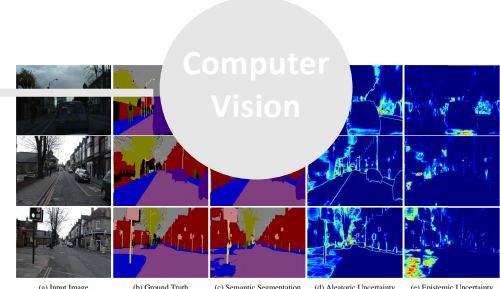
Which complex features are responsible of classification?

Planning

Plan Refinement



Which actions are responsible of a plan?



Computer Vision

Uncertainty Map

Search

KRR

Game Theory

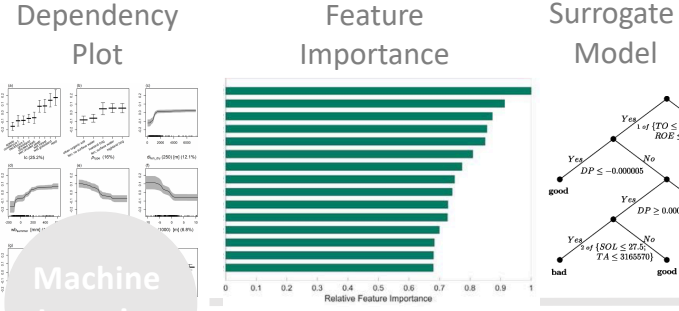
UAI

Robotics

NLP

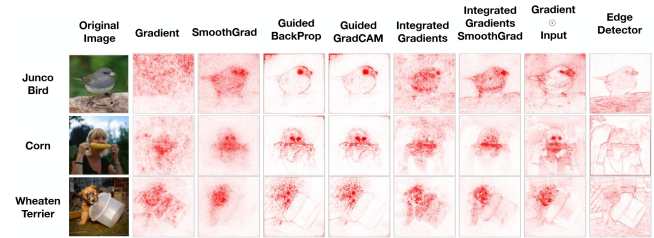
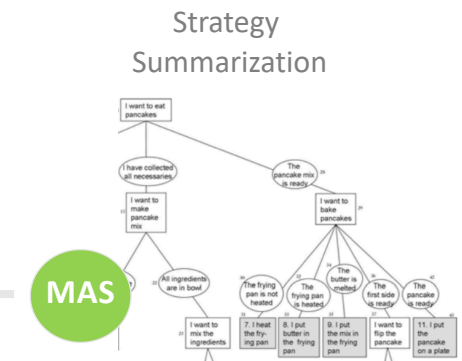
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map



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

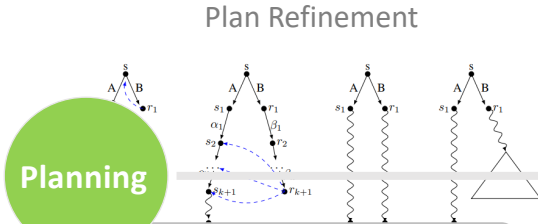
Artificial Intelligence



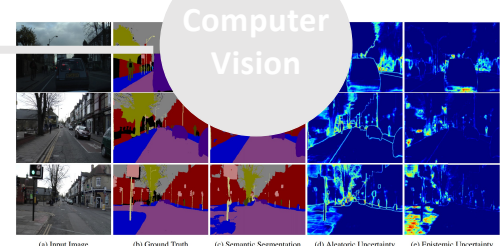
Which complex features are responsible of classification?

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

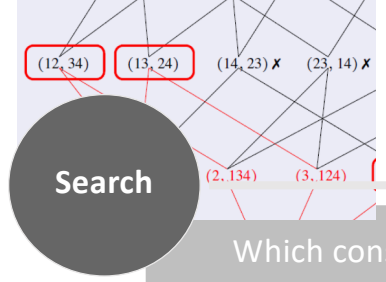
Which features are responsible of classification?



Which actions are responsible of a plan?



Uncertainty Map



Which constraints can be relaxed?

KRR

UAI

Game Theory

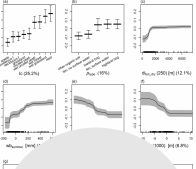
Robotics

NLP

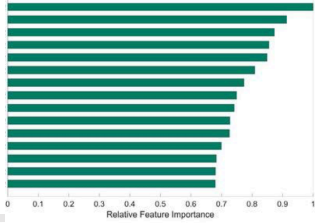
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

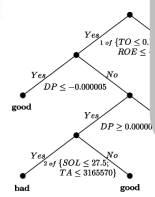
Dependency Plot



Feature Importance



Surrogate Model



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

Artificial Intelligence

Machine Learning

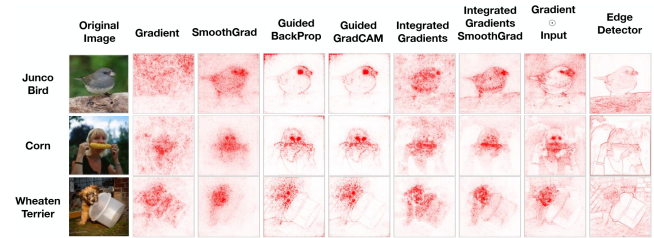
Which features are responsible of classification?

MAS

Strategy Summarization



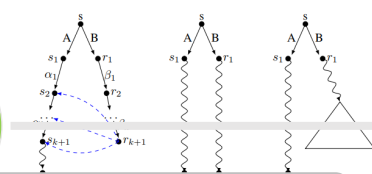
- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?



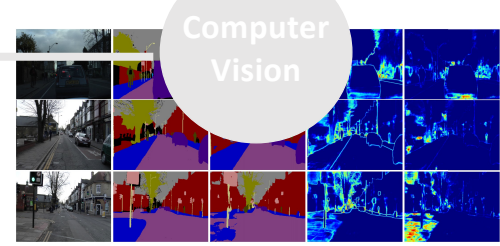
Which complex features are responsible of classification?

Planning

Plan Refinement



Which actions are responsible of a plan?

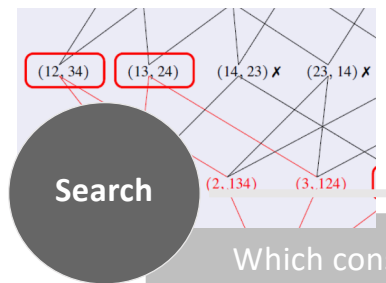


Computer Vision

Uncertainty Map

KRR

UAI



Conflicts Resolution

Search

Which constraints can be relaxed?

Game Theory

Which combination of features is optimal?

Robotics

NLP



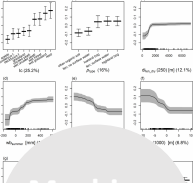
Shapely Values



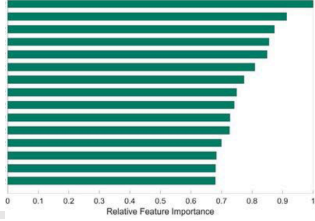
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

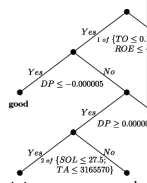
Dependency Plot



Feature Importance



Surrogate Model



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

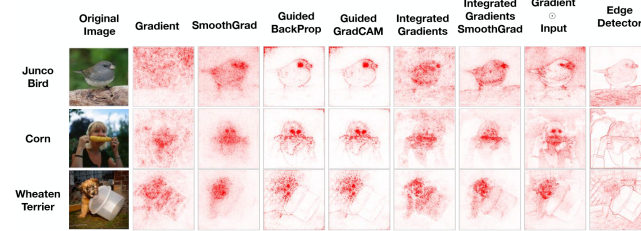
Artificial Intelligence

Machine Learning

Which features are responsible of classification?

MAS

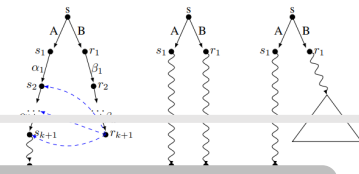
Strategy Summarization



Which complex features are responsible of classification?

Planning

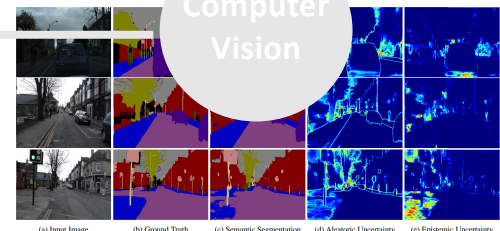
Plan Refinement



Which actions are responsible of a plan?

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

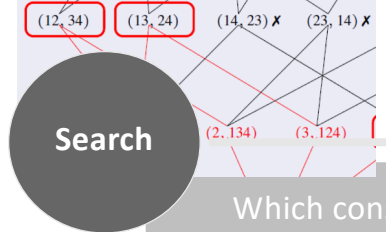
Computer Vision



Uncertainty Map

KRR

UAI



Search

Which constraints can be relaxed?

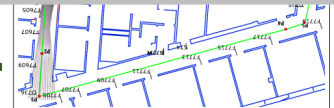
Game Theory

Which combination of features is optimal?

Robotics

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

NLP



Shapely Values

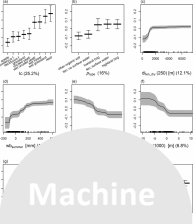
Narrative-based



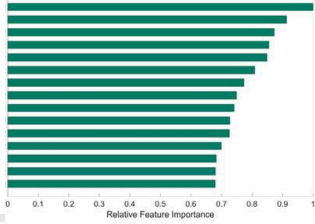
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

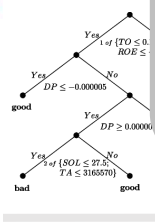
Dependency Plot



Feature Importance



Surrogate Model



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

Artificial Intelligence

Machine Learning

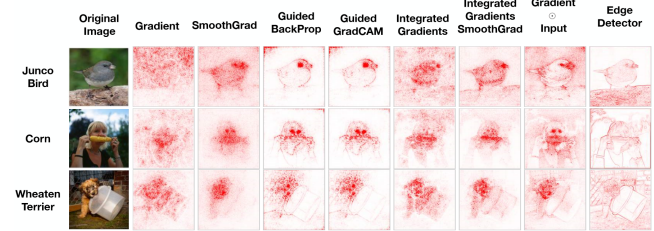
Which features are responsible of classification?

MAS

Strategy Summarization



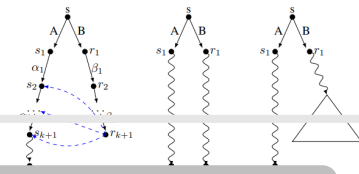
- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?



Which complex features are responsible of classification?

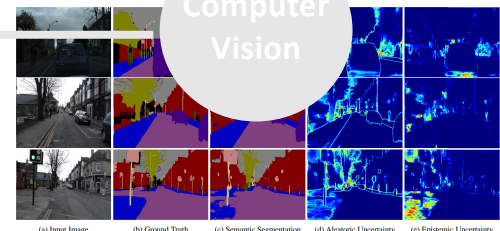
Planning

Plan Refinement



Which actions are responsible of a plan?

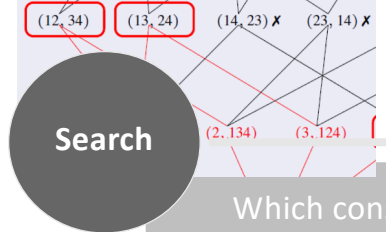
Computer Vision



Uncertainty Map

KRR

UAI



Which constraints can be relaxed?

Game Theory

Which combination of features is optimal?

Robotics

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

Machine Learning based

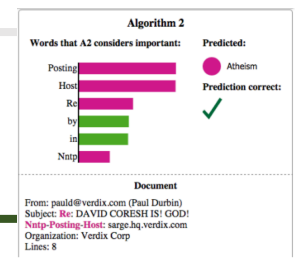
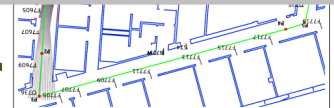
NLP

Which entity is responsible for classification?



Shapely Values

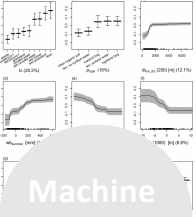
Narrative-based



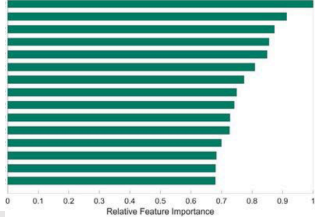
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

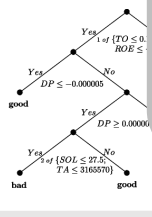
Dependency Plot



Feature Importance



Surrogate Model



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

Artificial Intelligence

Machine Learning

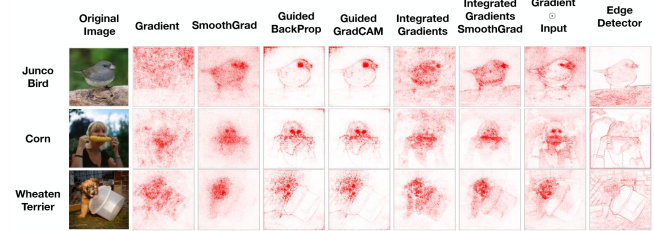
Which features are responsible of classification?

MAS

Strategy Summarization



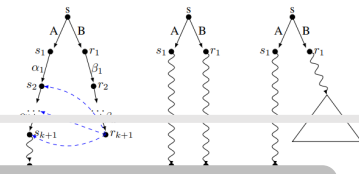
- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?



Which complex features are responsible of classification?

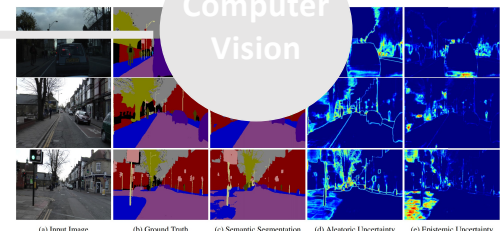
Planning

Plan Refinement



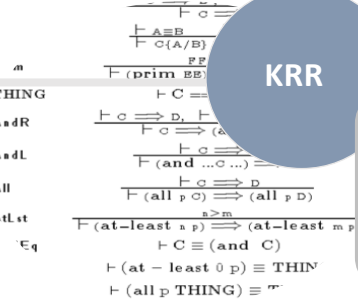
Which actions are responsible of a plan?

Computer Vision



UAI

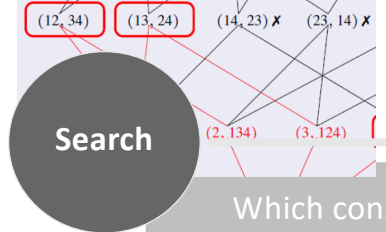
Diagnosis



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

Abduction

Uncertainty Map



Search

Which constraints can be relaxed?

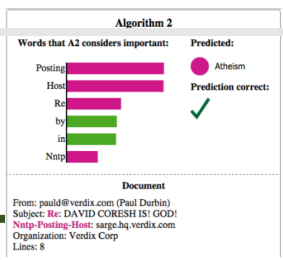
Game Theory

Which combination of features is optimal?

Robotics

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

Machine Learning based



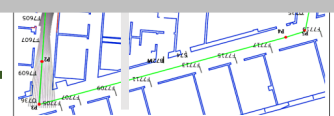
NLP

Which entity is responsible for classification?

Shapely Values



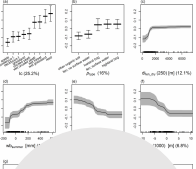
Narrative-based



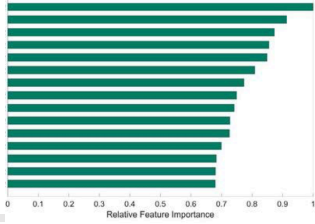
# XAI: One Objective, Many 'AI's, Many Definitions, Many Approaches

Saliency Map

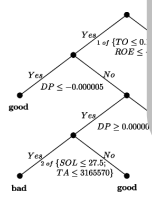
Dependency Plot



Feature Importance



Surrogate Model



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

Artificial Intelligence

Machine Learning

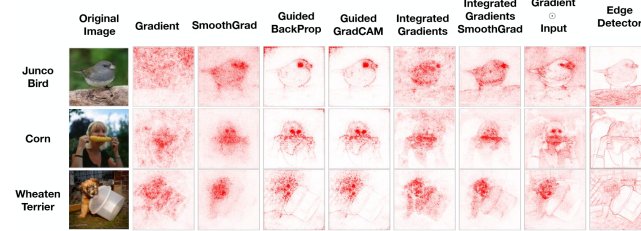
Which features are responsible of classification?

MAS

Strategy Summarization



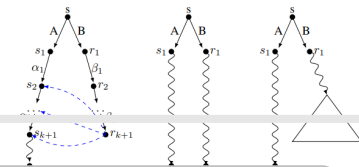
- Which agent strategy & plan ?
- Which player contributes most?
- Why such a conversational flow?



Which complex features are responsible of classification?

Planning

Plan Refinement



Which actions are responsible of a plan?

Computer Vision



Diagnosis

UAI

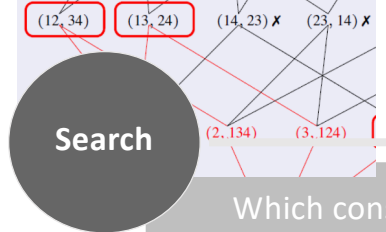
Uncertainty as an alternative to explanation

Abduction



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

Uncertainty Map



Search

Which constraints can be relaxed?

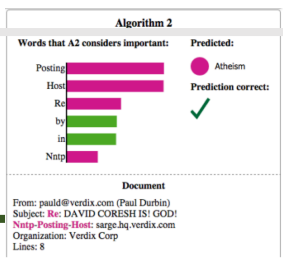
Game Theory

Which combination of features is optimal?

Robotics

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

Machine Learning based



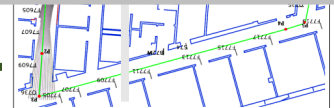
NLP

Which entity is responsible for classification?

Shapely Values



Narrative-based



# On the Role of Knowledge Graphs in Explainable AI A Machine Learning Perspective

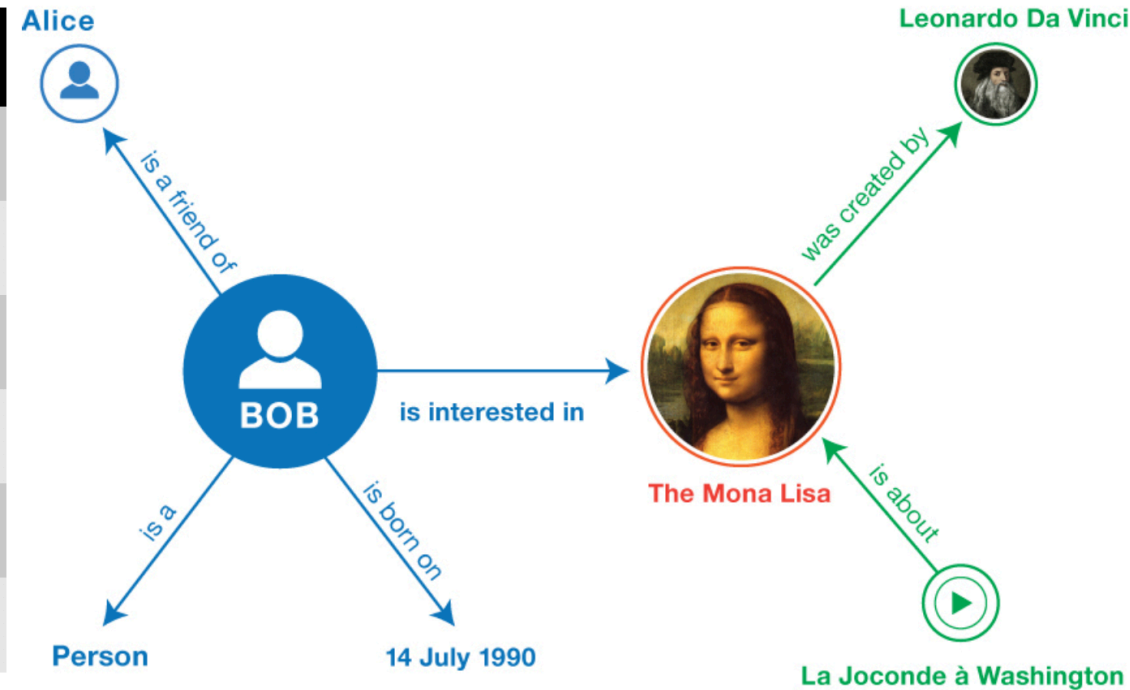
On the Role of Knowledge Graph in Explainable AI - under open review at the Semantic Web Journal -  
<http://www.semantic-web-journal.net/content/role-knowledge-graphs-explainable-ai>

---

# Knowledge Graph (1)

- Set of (*subject*, *predicate*, *object* — **SPO**) **triples** - *subject* and *object* are **entities**, and *predicate* is the **relationship** holding between them.
- Each SPO **triple** denotes a **fact**, i.e. the existence of an actual relationship between two entities.

subject	predicate	object
<i>Bob</i>	<i>is interested in</i>	<i>The Mona Lisa</i>
<i>Bob</i>	<i>is a friend of</i>	<i>Alice</i>
<i>The Mona Lisa</i>	<i>was created by</i>	<i>Leonardo Da Vinci</i>
<i>Bob</i>	<i>is a</i>	<i>Person</i>
<i>La Joconde à W.</i>	<i>is about</i>	<i>The Mona Lisa</i>
<i>Bob</i>	<i>is born on</i>	<i>14 July 1990</i>





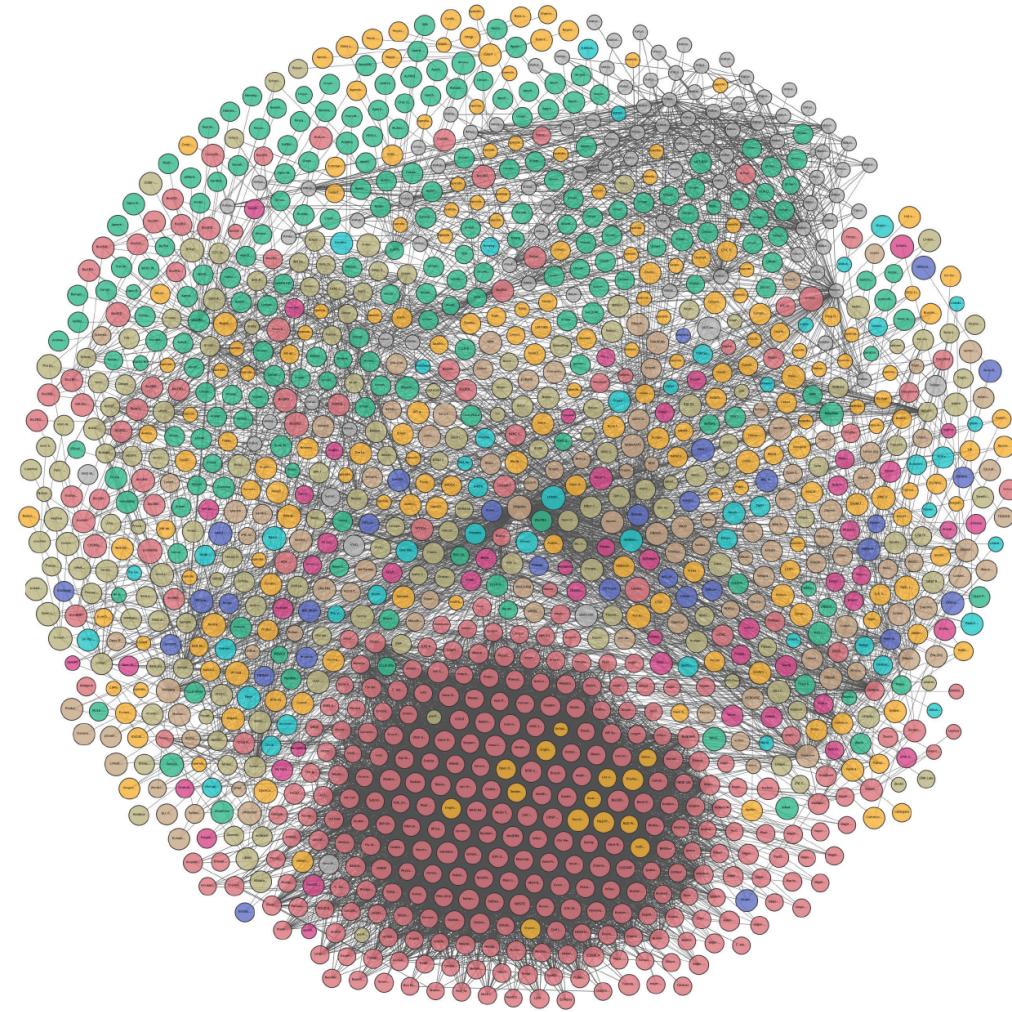
# Knowledge Graph (2)

Name	Entities	Relations	Types	Facts
Freebase	40M	35K	26.5K	637M
DBpedia (en)	4.6M	1.4K	735	580M
YAGO3	17M	77	488K	150M
Wikidata	15.6M	1.7K	23.2K	66M
NELL	2M	425	285	433K
Google KG	570M	35K	1.5K	18B
Knowledge Vault	45M	4.5K	1.1K	271M
Yahoo! KG	3.4M	800	250	1.39B

- **Manual Construction** - curated, collaborative
- **Automated Construction** - semi-structured, unstructured

Right: **Linked Open Data cloud** - over 1200 interlinked KGs encoding more than 200M facts about more than 50M entities.

Spans a variety of domains - Geography, Government, Life Sciences, Linguistics, Media, Publications, Cross-domain..





# Knowledge Graph Construction

Knowledge Graph construction methods can be classified in:

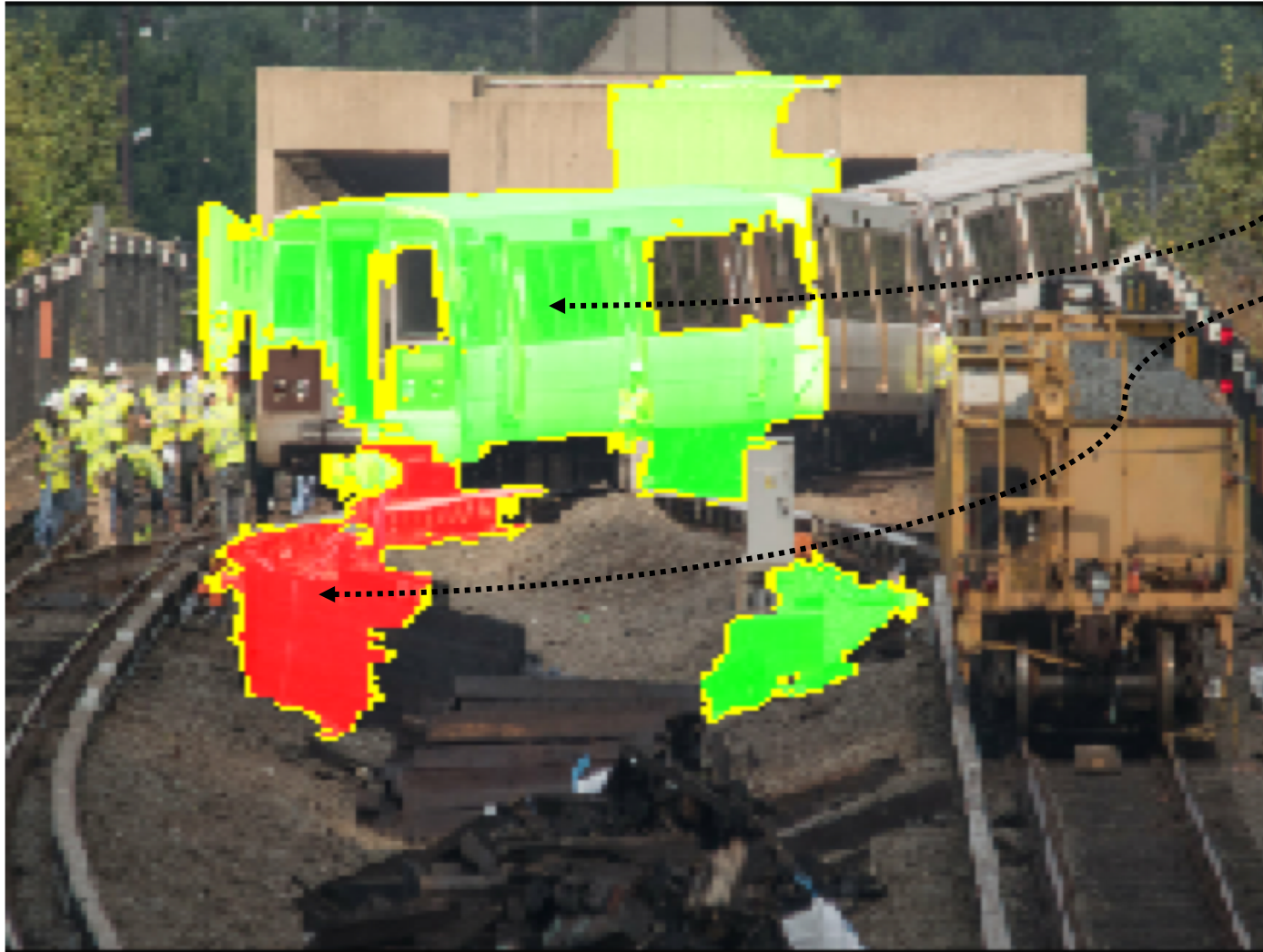
- **Manual** — curated (e.g. via experts), collaborative (e.g. via volunteers)
- **Automated** — semi-structured (e.g. from infoboxes), unstructured (e.g. from text)

Coverage is an issue:

- **Freebase** (40M entities) - 71% of persons without a birthplace, 75% without a nationality, even worse for other relation types [Dong et al. 2014]
- **DBpedia** (20M entities) - 61% of persons without a birthplace, 58% of scientists missing why they are popular [Krompaß et al. 2015]

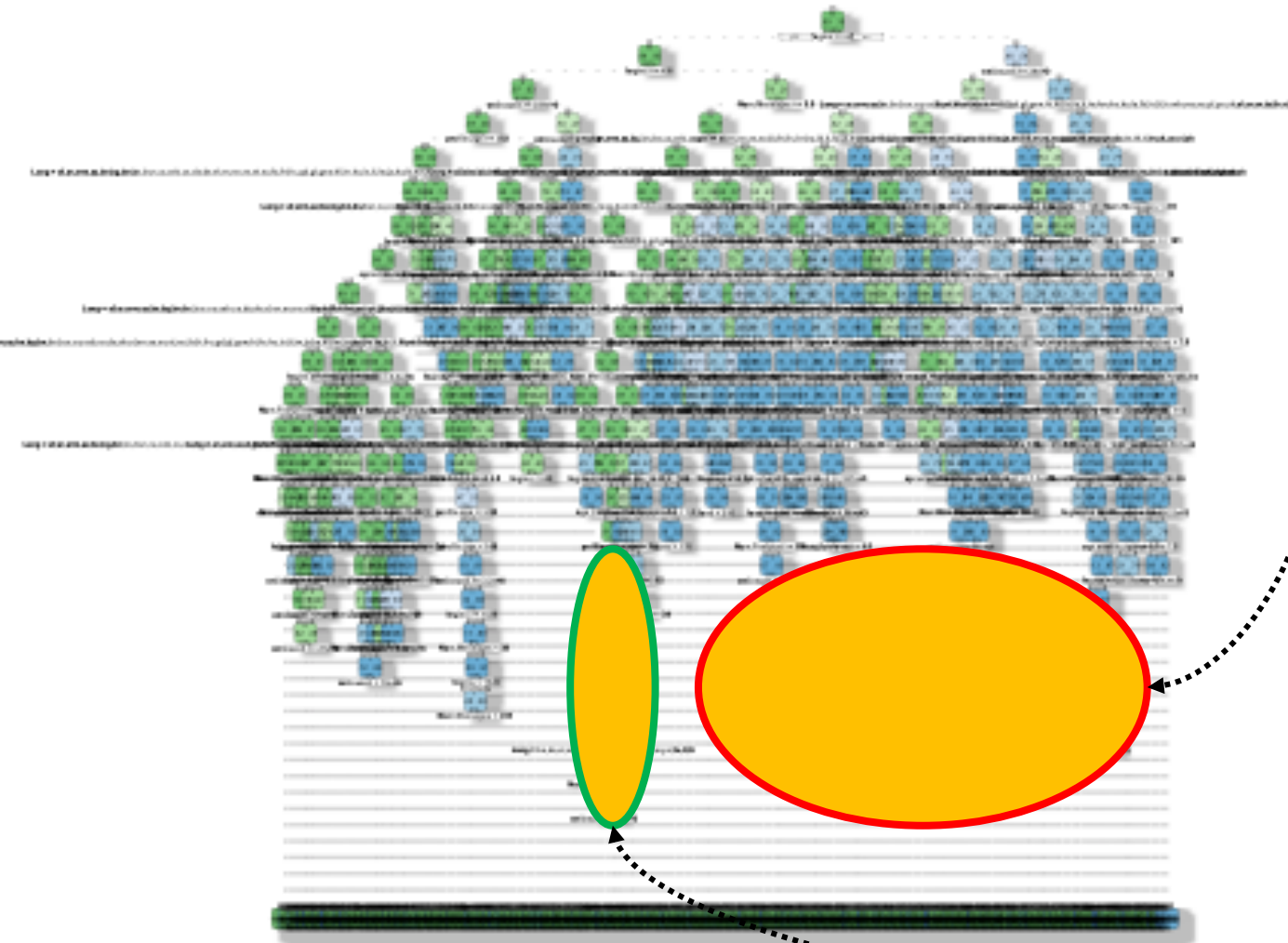
**Relational Learning** can help us overcoming these issues.

# Knowledge Graph in Machine Learning (1)



Augmenting (input) features with more semantics such as knowledge graph embeddings / entities

# Knowledge Graph in Machine Learning (2)

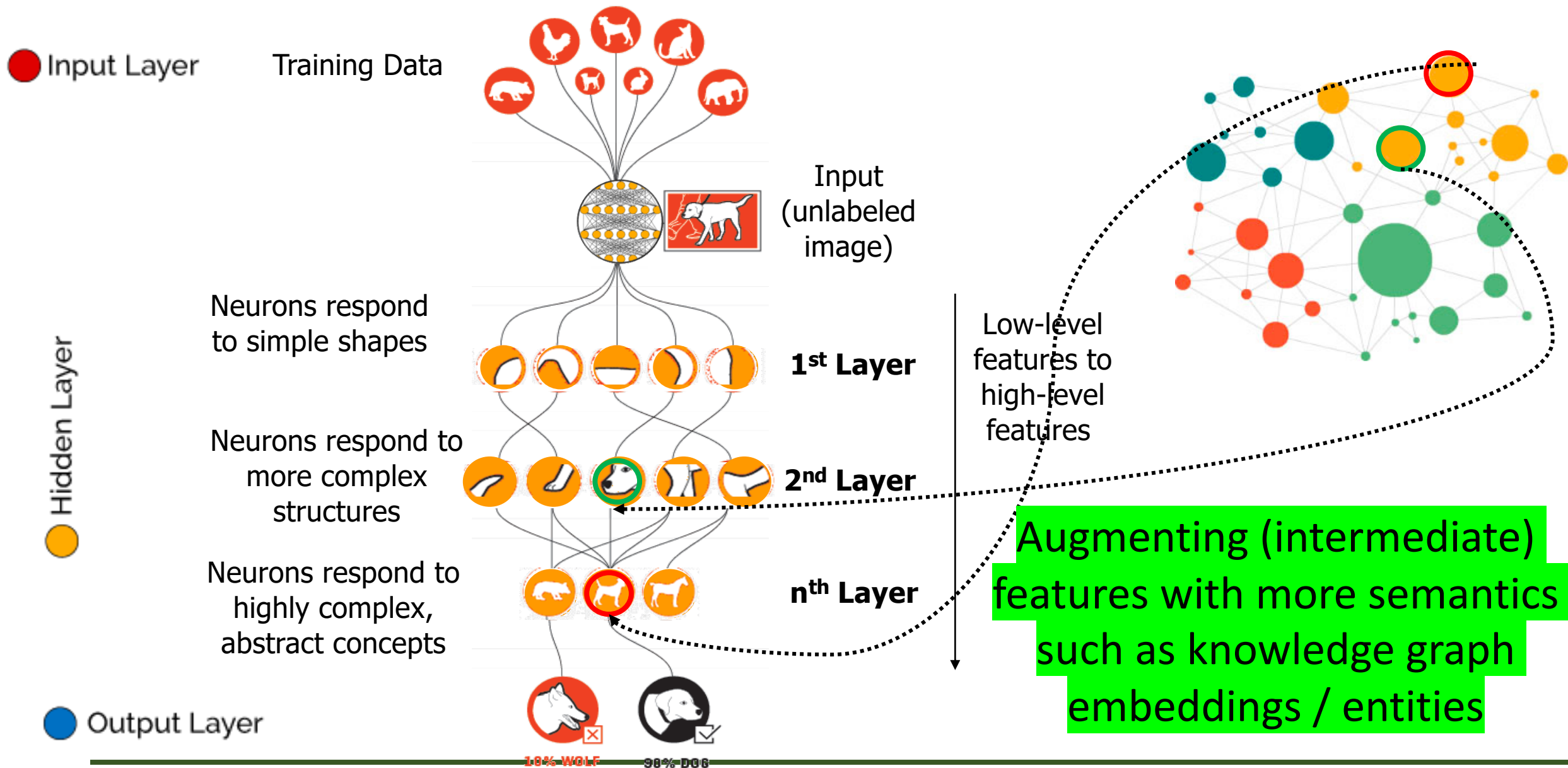


Augmenting machine learning models with more semantics such as knowledge graphs entities

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

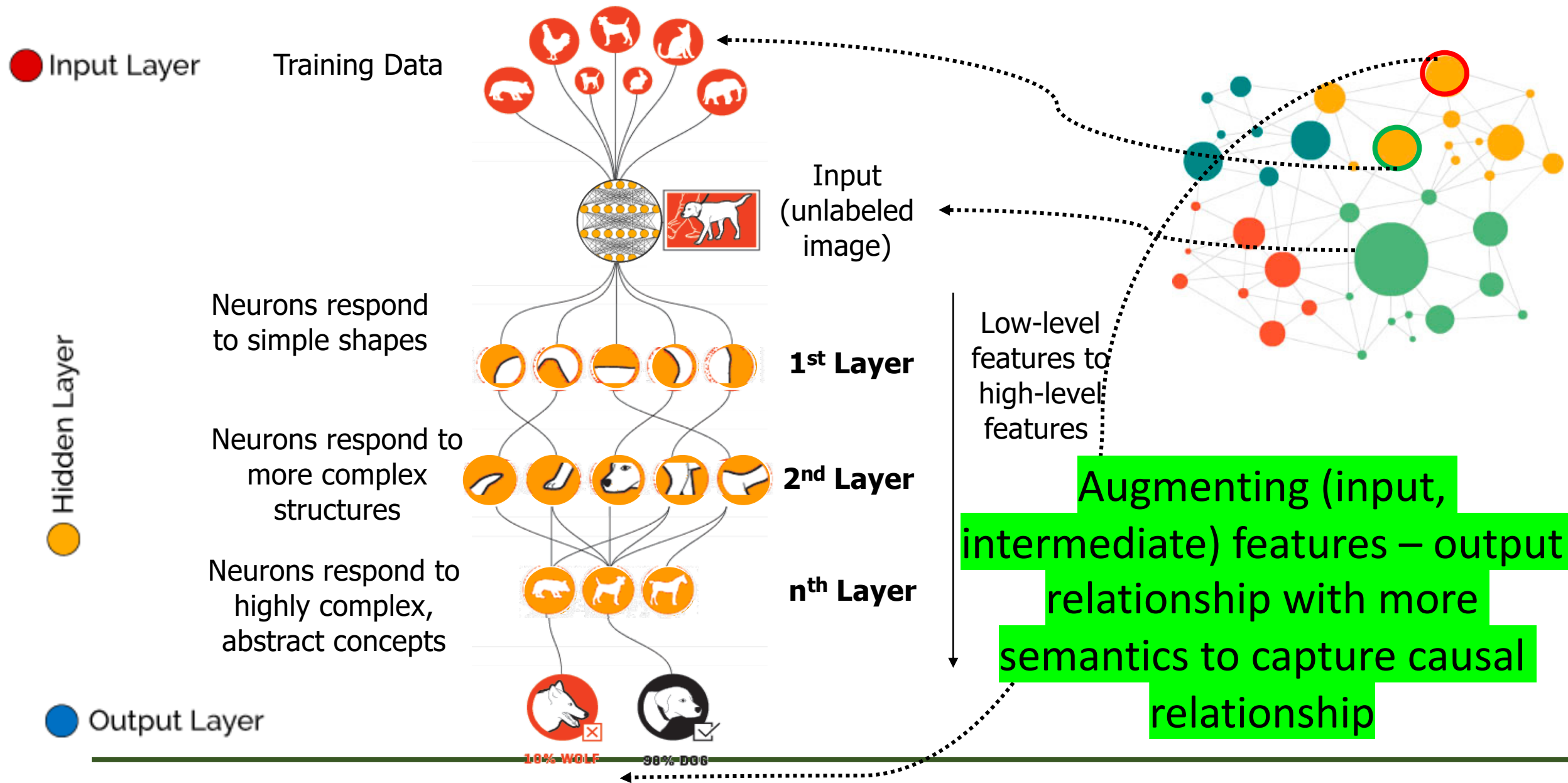
<https://stats.stackexchange.com/questions/230581/decision-tree-too-large-to-interpret>

# Knowledge Graph in Machine Learning (3)





# Knowledge Graph in Machine Learning (4)



# Knowledge Graph in Machine Learning (5)



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



Augmenting models with semantics to support personalized explanation



# Knowledge Graph in Machine Learning (6)

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

Augmenting input features and domains with semantics to support interpretable transfer learning

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-Zhejiang University Frontier Technology Research Center

#### **Freddy Lecue**

INRIA, France  
Accenture Labs, Ireland

#### **Ian Horrocks**

Department of Computer Science  
University of Oxford, UK

**On One  
Industrial  
Application in  
Thales**

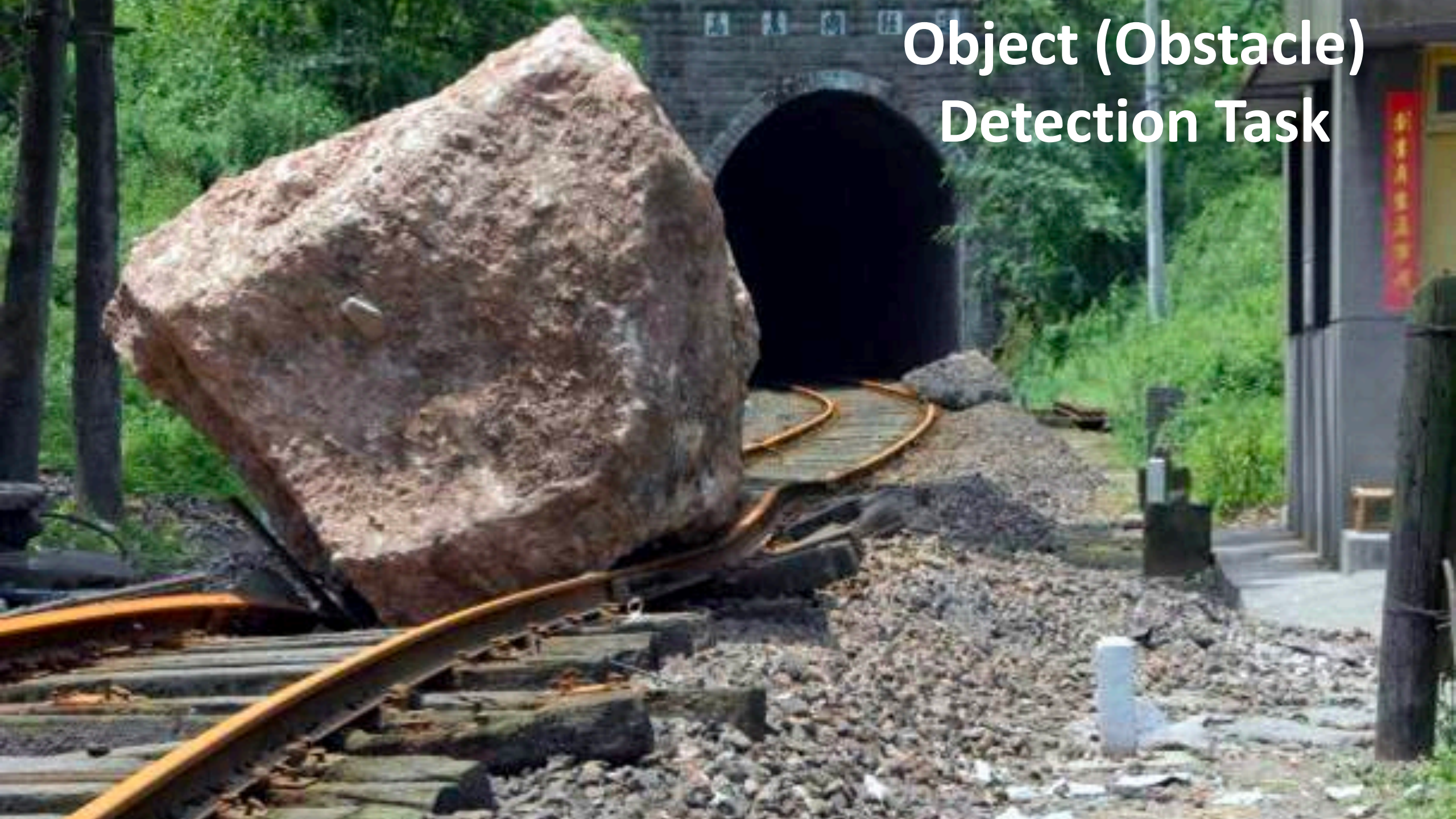
---

**State of the Art  
Machine Learning  
Applied to Critical  
Systems**

---



# Object (Obstacle) Detection Task





# Object (Obstacle) Detection Task State- of-the-art ML Result

Lumbermill - .59





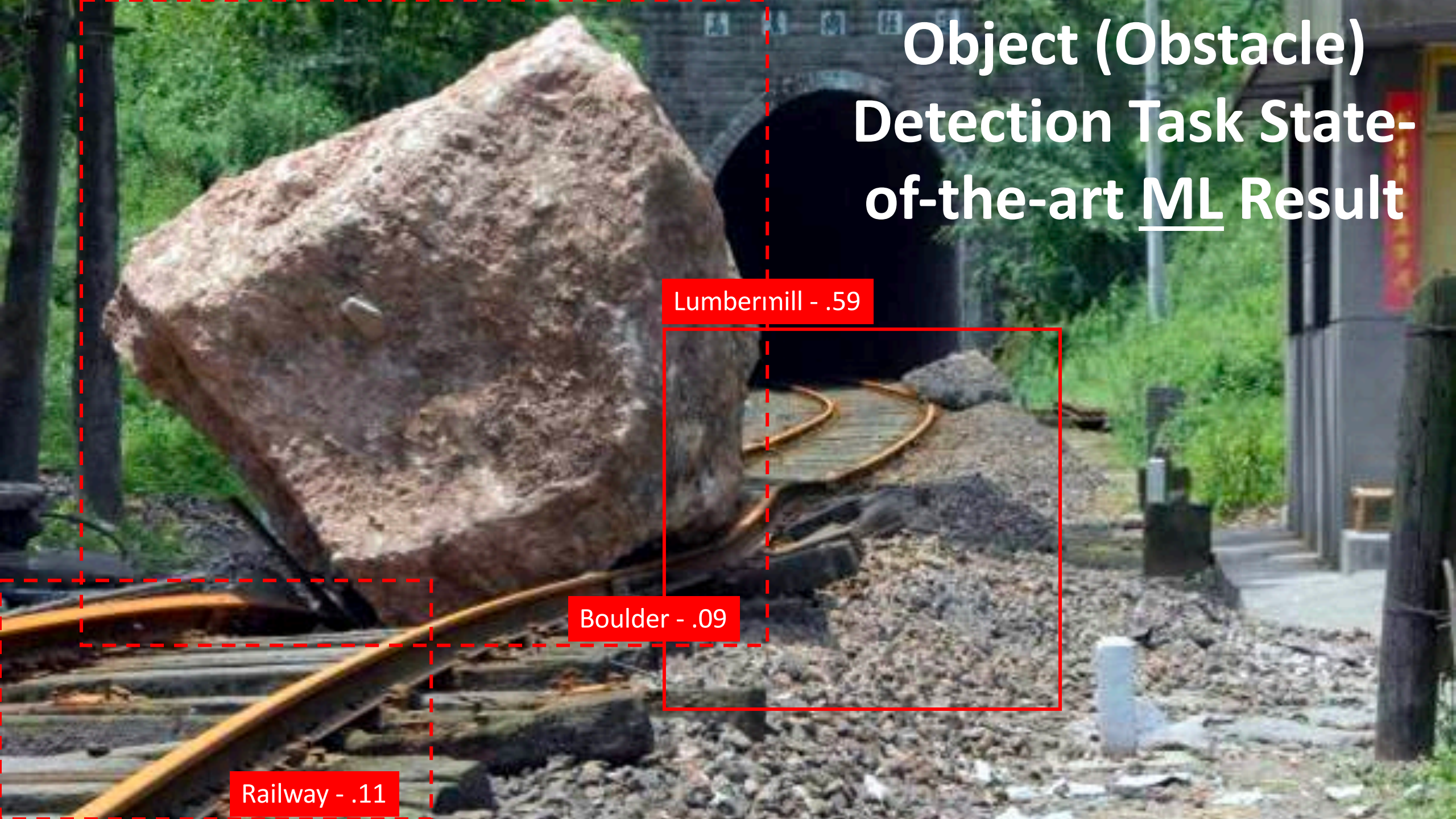
# Object (Obstacle) Detection Task State- of-the-art ML Result

Lumbermill - .59



Boulder - .09

Railway - .11





**State of the Art**

**XAI**

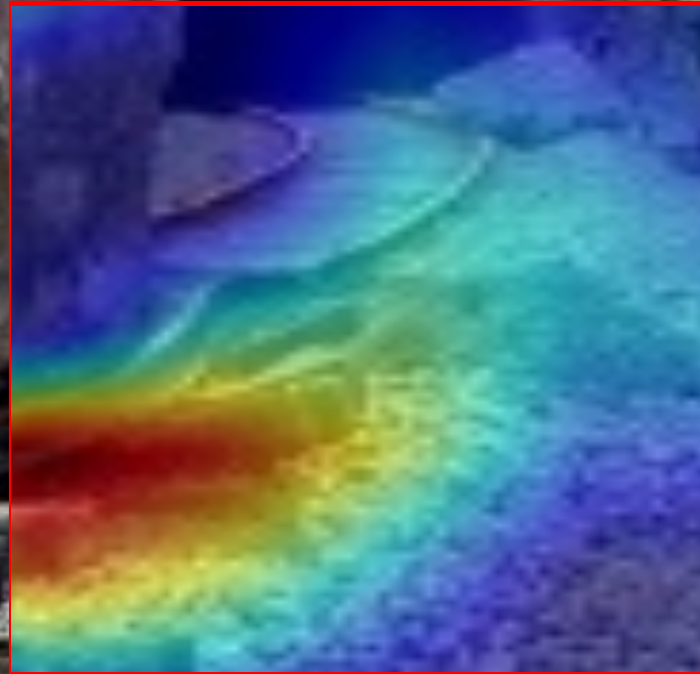
**Applied to Critical**

**Systems**

---

# Object (Obstacle) Detection Task State-of-the-art XAI Result

Lumbermill - .59



**Unfortunately, this is of  
NO use for a human  
behind the system**

---

**Let's stay back**

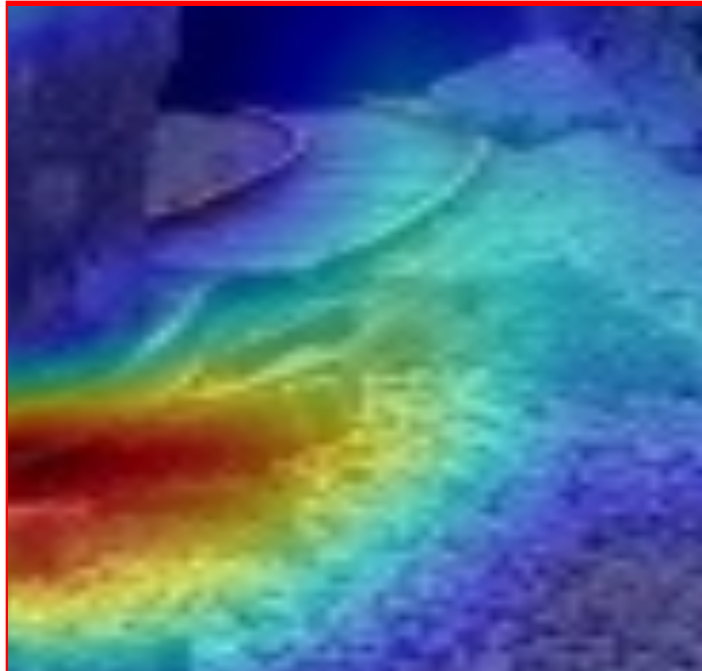
**Why this Explanation?  
(meta explanation)**

---



# After Human Reasoning...

Lumbermill - .59

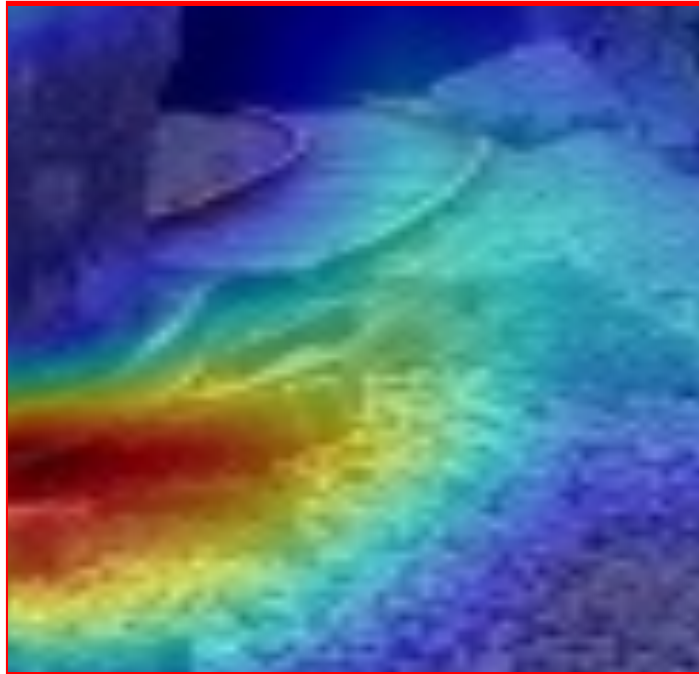


DBpedia		Browse using	Formats	Faceted Browser	Sparql Endpoint
dbo:wikiPageID	352327	(xsd:integer)			
dbo:wikiPageRevisionID	734430894	(xsd:integer)			
dct:subject	<ul style="list-style-type: none"><li>dbc:Sawmills</li><li>dbc:Saws</li><li>dbc:Ancient_Roman_technology</li><li>dbc:Timber_preparation</li><li>dbc:Timber_industry</li></ul>				
http://purl.org/linguistics/gold/hypernym	dbr:Facility				
rdf:type	<ul style="list-style-type: none"><li>owl:Thing</li><li>dbo:ArchitecturalStructure</li></ul>				
rdfs:comment	<p>A sawmill or lumber mill is a facility where logs are cut into lumber. Prior to the invention of the sawmill, boards were rived (split) and planed, or more often sawn by two men with a whipsaw, one above and another in a saw pit below. The earliest known mechanical mill is the Hierapolis sawmill, a Roman water-powered stone mill at Hierapolis, Asia Minor dating back to the 3rd century AD. Other water-powered mills followed and by the 11th century they were widespread in Spain and North Africa, the Middle East and Central Asia, and in the next few centuries, spread across Europe. The circular motion of the wheel was converted to a reciprocating motion at the saw blade. Generally, only the saw was powered, and the logs had to be loaded and moved by hand. An early improvement was the developm <sup>(en)</sup></p>				
rdfs:label	Sawmill <sup>(en)</sup>				
owl:sameAs	<ul style="list-style-type: none"><li>wikidata:Sawmill</li><li>dbpedia-cs:Sawmill</li><li>dbpedia-de:Sawmill</li><li>dbpedia-es:Sawmill</li></ul>				



# What is missing?




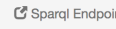
Lumbermill - .59



# Context matters

Boulder - .09

Railway - .11




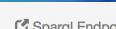
 [Browse using](#)  [Formats](#)  [Faceted Browser](#)  [Sparql Endpoint](#)

## About: Boulder

An Entity of Type : place, from Named Graph : <http://dbpedia.org>, within Data Space : [dbpedia.org](#)

In geology, a boulder is a rock fragment with size greater than 25.6 centimetres (10.1 in) in diameter. Smaller pieces are called cobbles and pebbles, depending on their "grain size". While a boulder may be small enough to move or roll manually, others are extremely massive. In common usage, a boulder is too large for a person to move. Smaller boulders are usually just called rocks or stones. The word boulder is short for boulder stone, from Middle English bulderston or Swedish bullersten. Boulder sized clasts are found in some sedimentary rocks, such as coarse conglomerate and boulder clay.

Property	Value
<a href="#">dbo:abstract</a>	<ul style="list-style-type: none"><li>In geology, a boulder is a rock fragment with size greater than 25.6 centimetres (10.1 in) in diameter. Smaller pieces are called cobbles and pebbles, depending on their "grain size". While a boulder may be small enough to move or roll manually, others are extremely massive. In common usage, a boulder is too large for a person to move. Smaller boulders are usually just called rocks or stones. The word boulder is short for boulder stone, from Middle English bulderston or Swedish bullersten. In places covered by ice sheets during Ice Ages, such as Scandinavia, northern North America, and Russia, glacial erratics are common. Erratics are boulders picked up by the ice sheet during its advance, and deposited during its retreat. They are called "erratic" because they typically are of a different rock type than the bedrock on which they are deposited. One of them is used as the pedestal of the Bronze Horseman in Saint Petersburg, Russia. Some noted rock formations involve giant boulders exposed by erosion, such as the Devil's Marbles in Australia's Northern Territory, the Horeke basalts in New Zealand, where an entire valley contains only boulders, and The Baths on the island of Virgin Gorda in the British Virgin Islands. Boulder sized clasts are found in some sedimentary rocks, such as coarse conglomerate and boulder clay. The climbing of large boulders is called bouldering. <sup>(en)</sup></li></ul>
<a href="#">dbo:thumbnail</a>	<ul style="list-style-type: none"><li><a href="#">wiki-commons:Special:FilePath/Balanced_Rock.jpg?width=300</a></li></ul>
<a href="#">dbo:wikiPageID</a>	<ul style="list-style-type: none"><li>60784 (xsd:integer)</li></ul>
<a href="#">dbo:wikiPageRevisionID</a>	<ul style="list-style-type: none"><li>743049914 (xsd:integer)</li></ul>
<a href="#">dct:subject</a>	<ul style="list-style-type: none"><li><a href="#">dbc:Rock_formation</a></li><li><a href="#">dbc:Rocks</a></li></ul>

 [Browse using](#)  [Formats](#)  [Faceted Browser](#)  [Sparql Endpoint](#)

## About: Rail transport

An Entity of Type : software, from Named Graph : <http://dbpedia.org>, within Data Space : [dbpedia.org](#)

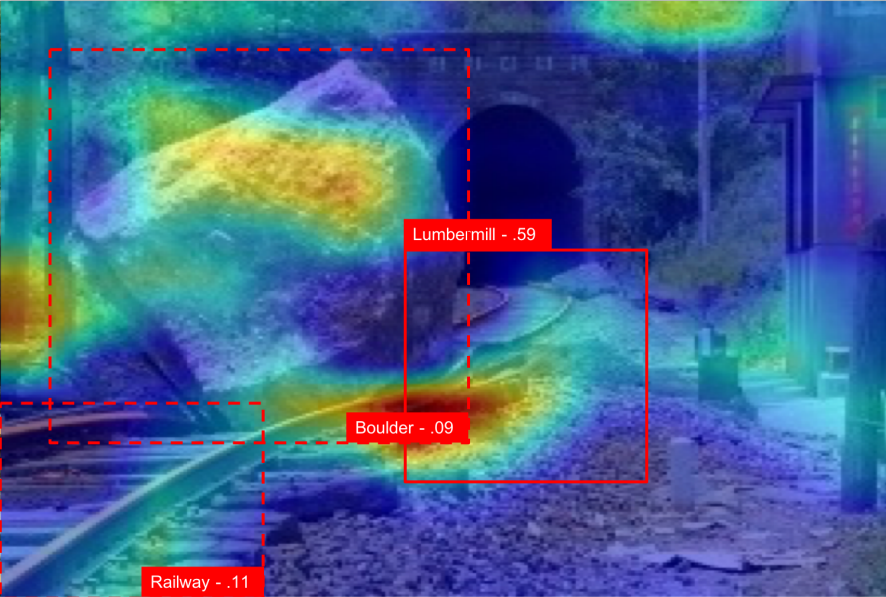
Rail transport is a means of conveyance of passengers and goods on wheeled vehicles running on rails, also known as tracks. It is also commonly referred to as train transport. In contrast to road transport, where vehicles run on a prepared flat surface, rail vehicles (rolling stock) are directionally guided by the tracks on which they run. Tracks usually consist of steel rails, installed on ties (sleepers) and ballast, on which the rolling stock, usually fitted with metal wheels, moves. Other variations are also possible, such as slab track, where the rails are fastened to a concrete foundation resting on a prepared subsurface.

Property	Value
<a href="#">dbo:abstract</a>	<ul style="list-style-type: none"><li>Rail transport is a means of conveyance of passengers and goods on wheeled vehicles running on rails, also known as tracks. It is also commonly referred to as train transport. In contrast to road transport, where vehicles run on a prepared flat surface, rail vehicles (rolling stock) are directionally guided by the tracks on which they run. Tracks usually consist of steel rails, installed on ties (sleepers) and ballast, on which the rolling stock, usually fitted with metal wheels, moves. Other variations are also possible, such as slab track, where the rails are fastened to a concrete foundation resting on a prepared subsurface. Rolling stock in a rail transport system generally encounters lower frictional resistance than road vehicles, so passenger and freight cars (carriages and wagons) can be coupled into longer trains. The operation is carried out by a railway company, providing transport between train stations or freight customer facilities. Power is provided by locomotives which either draw electric power from a railway electrification system or produce their own power, usually by diesel engines. Most tracks are accompanied by a signalling system. Railways are a safe land transport system when compared to other forms of transport. Railway transport is capable of high levels of passenger and cargo utilization and energy efficiency, but is often less flexible and more capital-intensive than road transport, when lower traffic levels are considered. The oldest, man-hauled railways date back to the 6th century BC, with Periander, one of the Seven Sages of Greece,</li></ul>

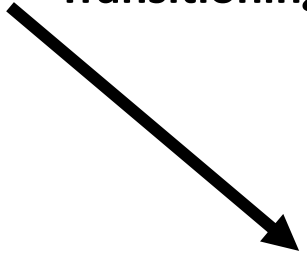
# **XAI Thales Platform**

- **Higher accuracy with no intensive fine-tuning**
  - **Human interpretable explanation**
  - **Running on the edge at inference time**
-

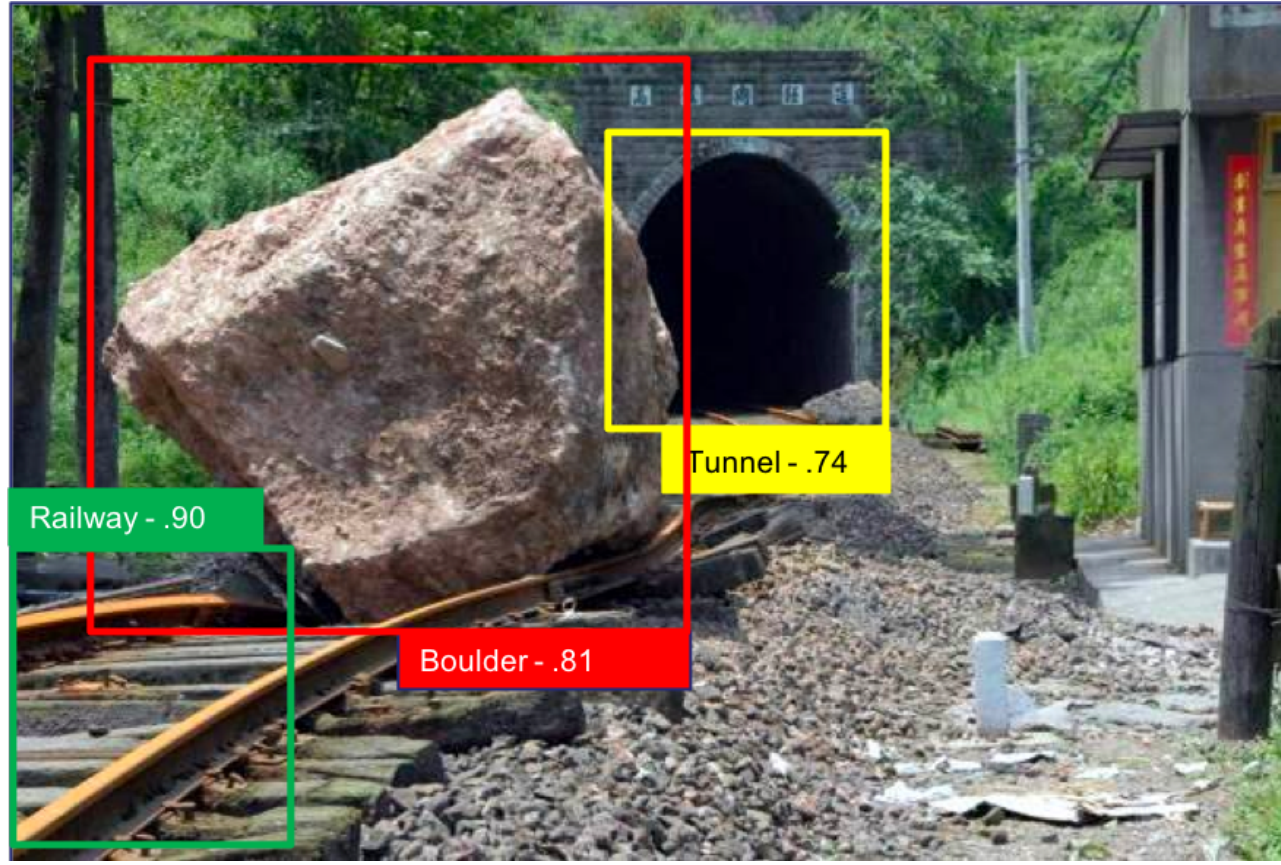




Transitioning



- **Hardware:** High performance, scalable, generic (to different FGPA family) & portable CNN dedicated programmable processor implemented on an FPGA for real-time embedded inference
- **Software:** Knowledge graph extension of object detection

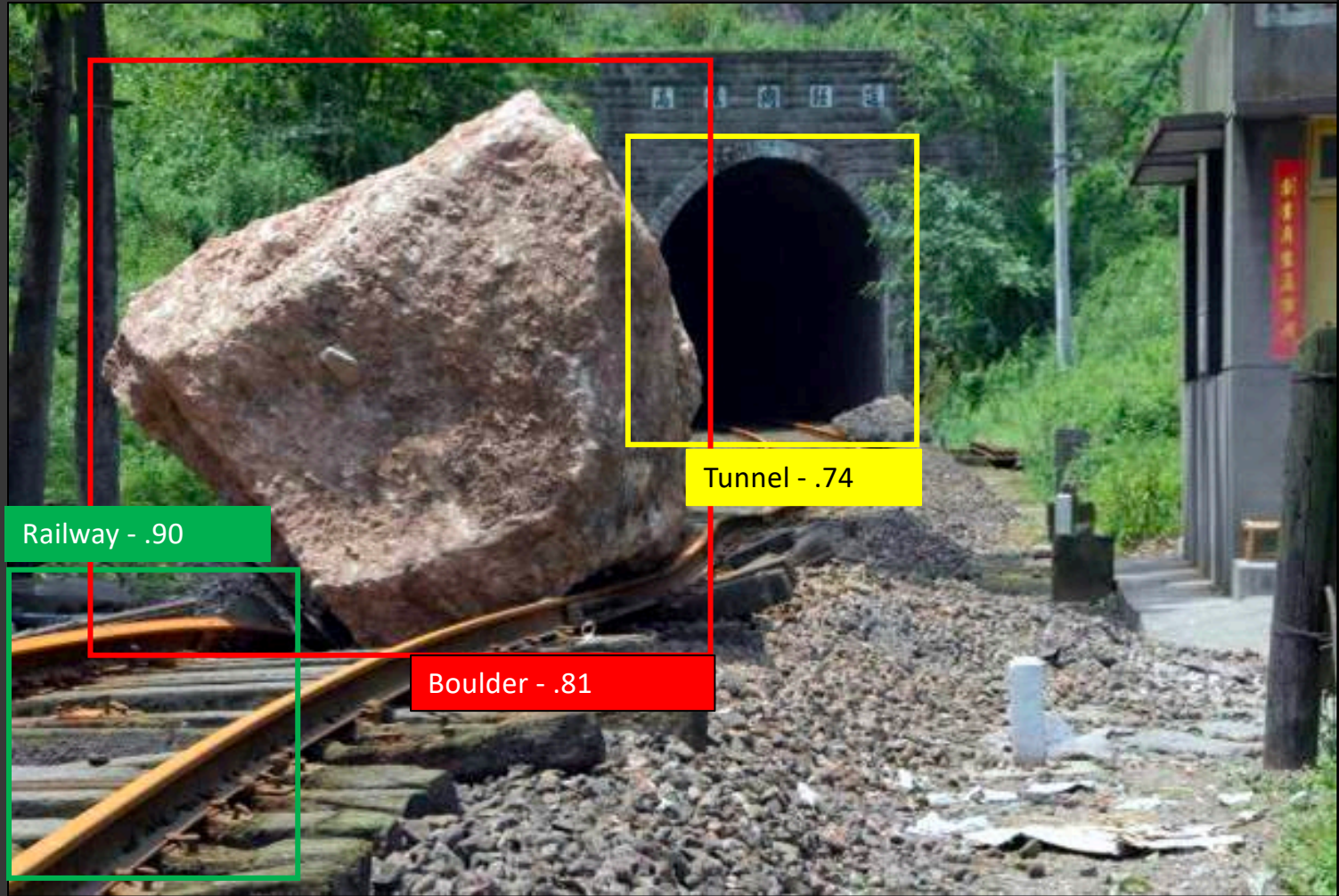
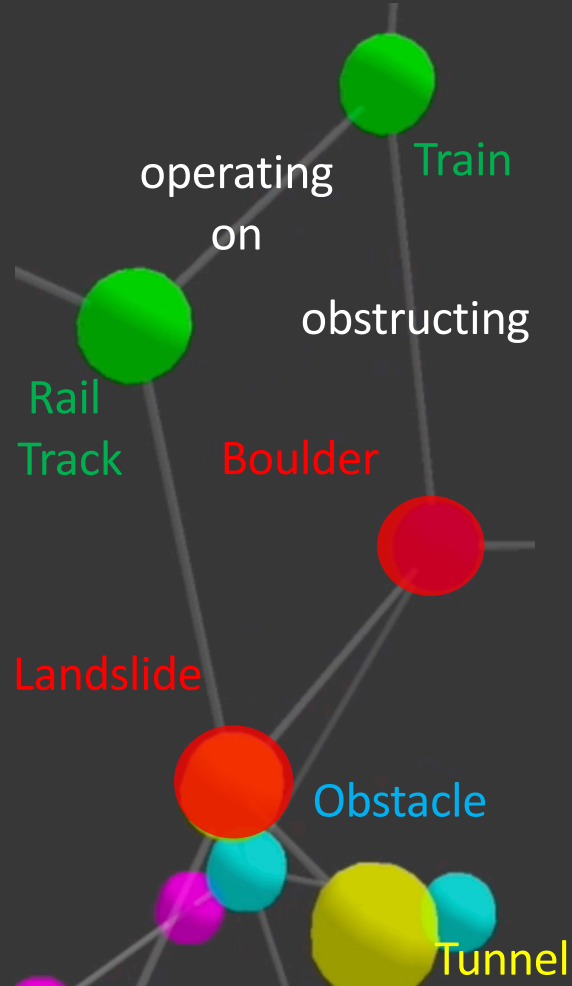
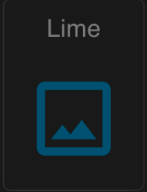


This is an **Obstacle: Boulder** obstructing the train: XG142-R on **Rail\_Track** from City: Cannes to City: Marseille at **Location: Tunnel VIX** due to **Landslide**

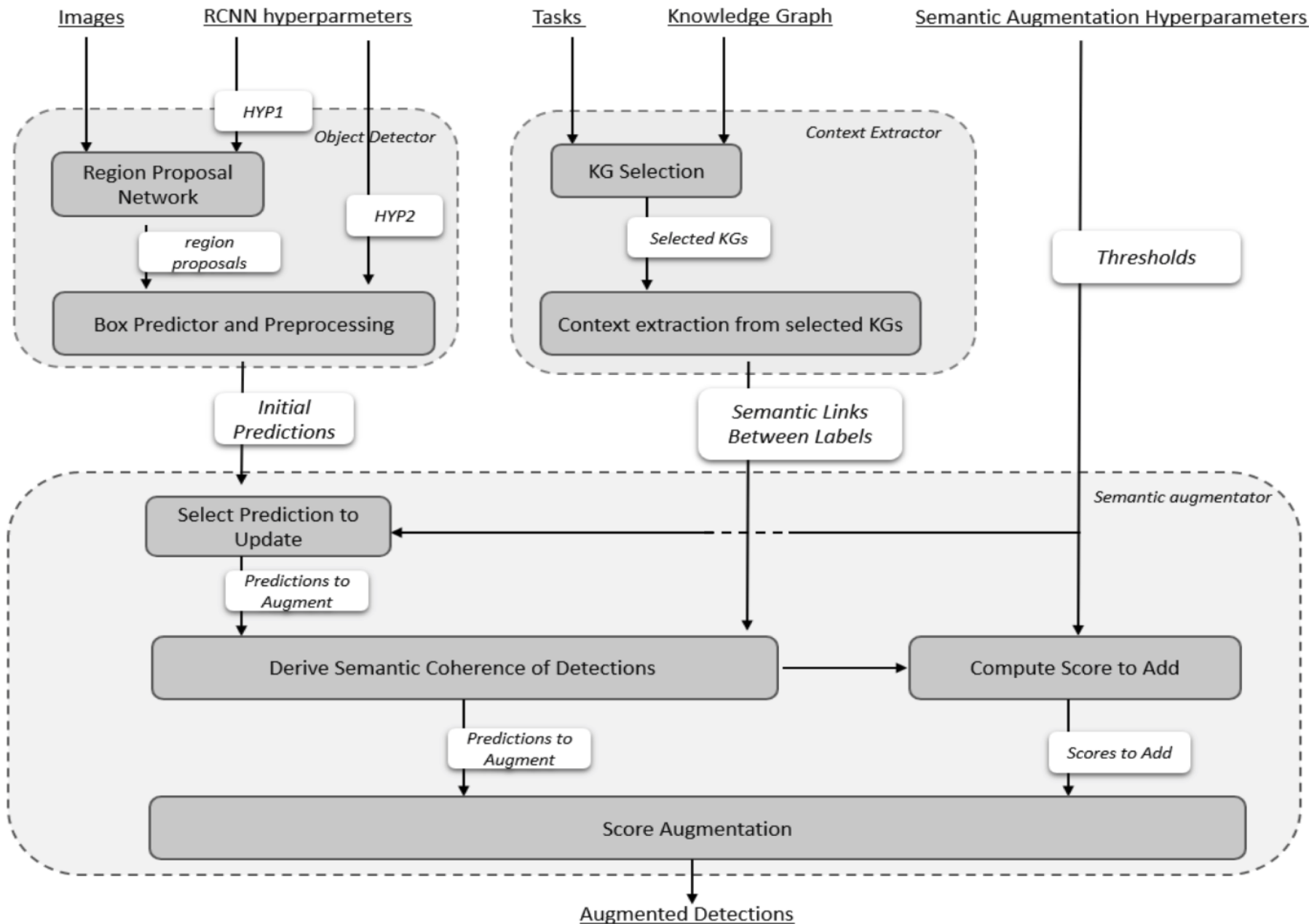


EXPLANATIONS

ResNet50 image classifier







Freddy Lécué, Jiaoyan Chen, Jeff Z. Pan, Huajun Chen: Augmenting Transfer Learning with Semantic Reasoning. IJCAI 2019: 1779-1785

Freddy Lécué, Tanguy Pommellet: Feeding Machine Learning with Knowledge Graphs for Explainable Object Detection. ISWC Satellites 2019: 277-280

Freddy Lécué, Baptiste Abeloos, Jonathan Anctil, Manuel Bergeron, Damien Dalla-Rosa, Simon Corbeil-Letourneau, Florian Martet, Tanguy Pommellet, Laura Salvan, Simon Veilleux, Maryam Ziaeeafard: Thales XAI Platform: Adaptable Explanation of Machine Learning Systems - A Knowledge Graphs Perspective. ISWC Satellites 2019: 315-316

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

**More on XAI**

---

# (Some) Tutorials, Workshops, Challenge

## Tutorial:

- AAAI 2020 Tutorial On Explainable AI: From Theory to Motivation, Applications and Limitations (#2) - <https://xaitutorial2019.github.io/> <https://xaitutorial2020.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/>
- KDD 2019 Tutorial on Explainable AI in Industry (#1) - <https://sites.google.com/view/kdd19-explainable-ai-tutorial>

## 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>
- KDD 2019 Workshop on Explainable AI for fairness, accountability, and transparency (#1) – <https://xai.kdd2019.a.intuit.com>
- 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/>
- IEEE FUZZ 2019 / Advances on eXplainable Artificial Intelligence (#2) - <https://sites.google.com/view/xai-fuzzieee2019>
- International Conference on NL Generation - Interactive Natural Language Technology for Explainable Artificial Intelligence (EU H2020 NL4XAI; #1) - <https://sites.google.com/view/nl4xai2019/>

## 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)*
-



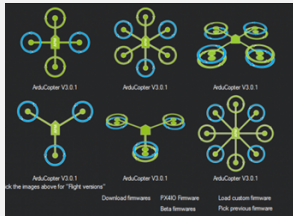
# (Some) Initiatives: XAI in USA



## Challenge Problem Areas



**Data Analytics**  
Multimedia Data

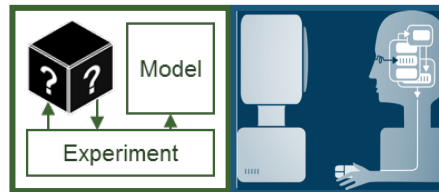
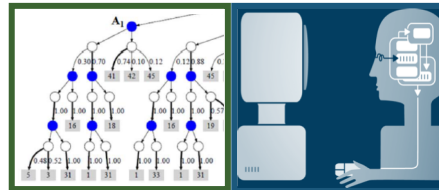
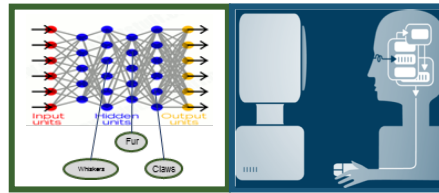


**Autonomy**  
ArduPilot &  
SITL Simulation

## 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**

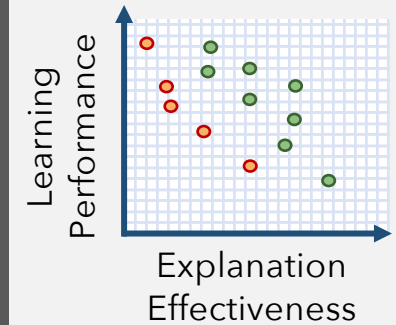
**Evaluator**

## 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

## 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

• DEEL  ab  il  (Learning) Project 2019-2024

• Research institutions

• Industrial partners



•  ers



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



# Conclusion

---

## Why do we Need XAI by the Way?

- ***To empower*** individual against undesired effects of automated decision making
  - ***To reveal*** and protect new vulnerabilities
  - ***To implement*** the “right of explanation”
  - ***To improve*** industrial standards for developing AI-powered products, increasing the trust of companies and consumers
  - ***To help*** people make better decisions
  - ***To align*** algorithms with human values
  - ***To preserve*** (and expand) human autonomy
  - **To scale and industrialize AI**
-



## Why do we Need Knowledge Graphs to Achieve XAI?

Because this is not an explanation from an intelligent system

This is even not interpretable, and then not actionable



## Conclusion

- Explainable AI is motivated by **real-world applications in AI**
  - Not a new problem – a reformulation of past research challenges in AI
  - Knowledge graphs should be foundational for XAI
  - But they are facing challenges related to their integration (data mapping)
  - **Many industrial applications already – crucial for AI adoption in critical systems**
-

# Open Research Questions for the Semantic Web / Knowledge Graph Community

- [Data] Machine learning experts do not buy the **data – knowledge mapping**
- [Explanation] There is ***no agreement*** on ***what an explanation is***
- [Explanation] There is ***not a formalism*** for ***explanations (neither model nor output)***
- [Model] ***There is very limited work in machine learning modules composability – and none from a semantics perspective***
- [Model] ***There is no work on describing and representing models***
- [Model] What are **disentangled representations** and how can its factors be quantified and detected?
- [Human-in-the-loop] There is ***no work*** that seriously addresses the problem of ***quantifying*** the grade of ***comprehensibility*** of an explanation for humans





# Job Openings

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

Thales is a global technology leader for the Defence and aerospace class technology, the combined expertise of our experts have made Thales a key player in keeping the public safe by protecting the national security interests of countries around the world.

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

This is a unique opportunity to play a key role on a world-class Technology (TRT) in Canada (Quebec and Montreal). We have applied R&T experts at five locations worldwide. We are working on intelligence technologies. Our passion is imagining and developing cutting edge AI technologies. Not only will you join a global network, but this TRT is also co-located within our new Artificial Intelligence eXpertise i.e., the new flagship program to work.

## Job Description

An AI (Artificial Intelligence) Research and Technology Scientist will be developing innovative prototypes to demonstrate artificial intelligence. To be successful in this role, one must be able to think what's new, and a strong ability to learn new technologies. You will have hands-on technical skills and be familiar with latest technologies. You will contribute as technical subject matter expert to our products and its business units. In addition to the implementation of the product, the individual will also be involved in the initial project planning, and team work is also critical for this role.

As a Research and Technology Applied AI Scientist you will be working on fast paced projects.

## Professional Skill Requirements

- Good foundation in mathematics, statistics

- 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 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
- 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

AUGUST 28TH, 2019

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

@freddylecue  
<https://tinyurl.com/freddylecue>  
Freddy.lecue.e@thalesdigital.io