

# **Explainable AI for Fairness and Accountability**

**Leilani H. Gilpin, Ph.D.**

**Assistant Professor of Computer Science & Engineering**

**UC Santa Cruz**

**[lgilpin@ucsc.edu](mailto:lgilpin@ucsc.edu)**

**[lgilpin.com](http://lgilpin.com)**

# Talk Agenda

Brief Intro

Motivate problem: Systems are imperfect

What is explainability?

What is *actually* being explained?

How to evaluate explainability?

How to explain complex systems? (autonomous driving)

# About Me

- B.S in Computer Science, B.S. in Mathematics at UC San Diego
- M.S. in Computational Mathematics from Stanford University (2013), Ph.D. in EECS from MIT (2020).
- Industry experience
  - Xerox PARC
  - INRIA (France)
  - Sony AI
- Research: The methodologies and technologies for complex systems to explain themselves.



# Complex Systems Fail in Complex Ways



**Predictive Inequity in Object Detection**

---

**Benjamin Wilson<sup>1</sup> Judy Hoffman<sup>1</sup> Jamie Morgenstern<sup>1</sup>**

K. Eykholt et al. "Robust Physical-World Attacks on Deep Learning Visual Classification."

# Societal Need for Explanation

BUSINESS NEWS OCTOBER 9, 2018 / 11:12 PM / 2 MONTHS AGO

## Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin

8 MIN READ



### Business Impact

## An AI-Fueled Credit Formula Might Help You Get a Loan

Startup ZestFinance says it has built a machine-learning system that's smart enough to find new borrowers and keep bias out of its credit analysis.

by Nanette Byrnes February 14, 2017

# Talk Agenda

Motivate problem: Systems are imperfect

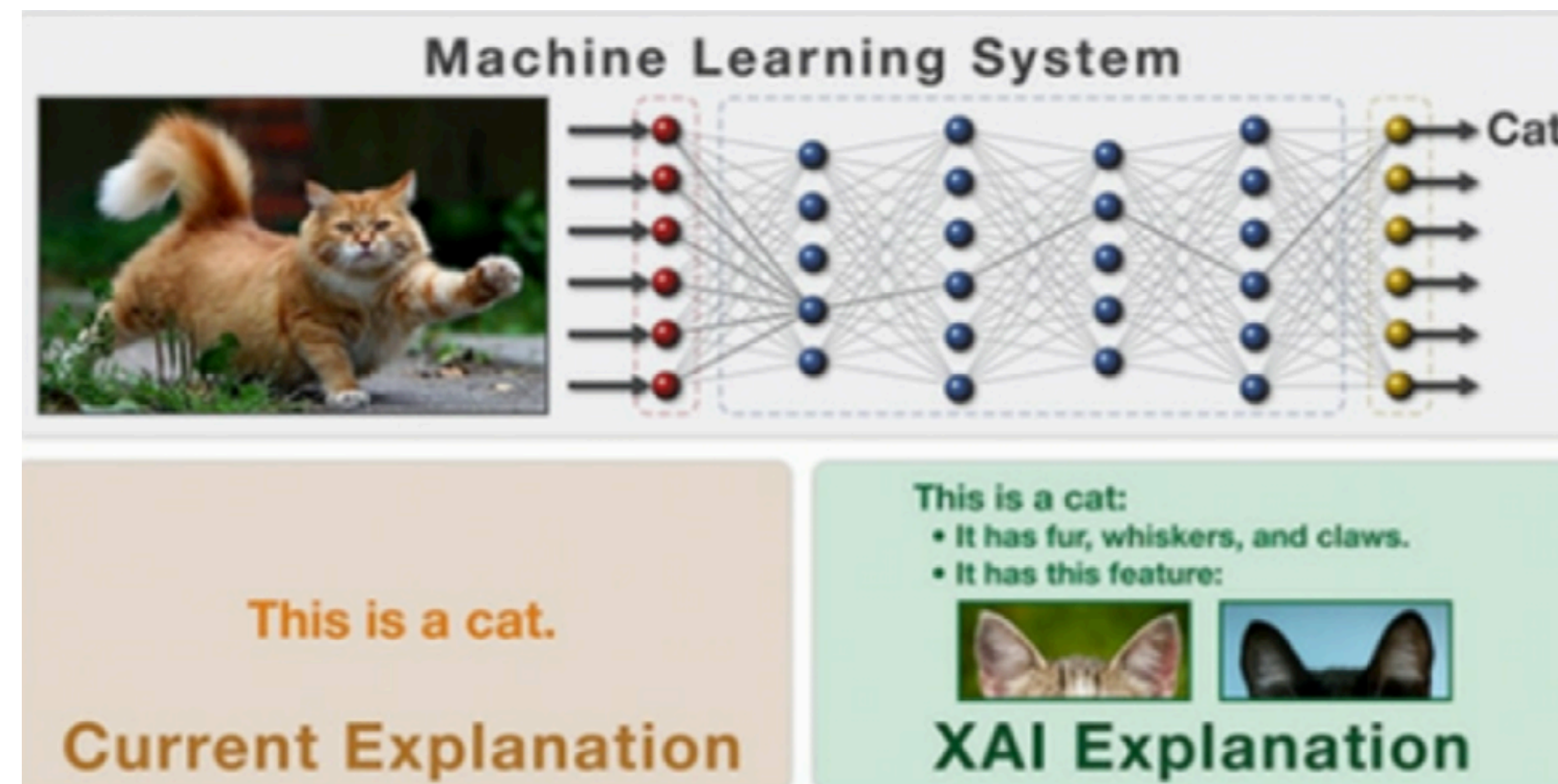
What is explainability?

What is *actually* being explained?

How to evaluate explainability?

How to explain complex systems? (autonomous driving)

# What is Explainability?



From Darpa XAI

**“Explanations...express answer to not just any questions but to questions that present the kind of intellectual difficulty...”**

Sylvain Bromberger, *On What We Know We Don't Know*



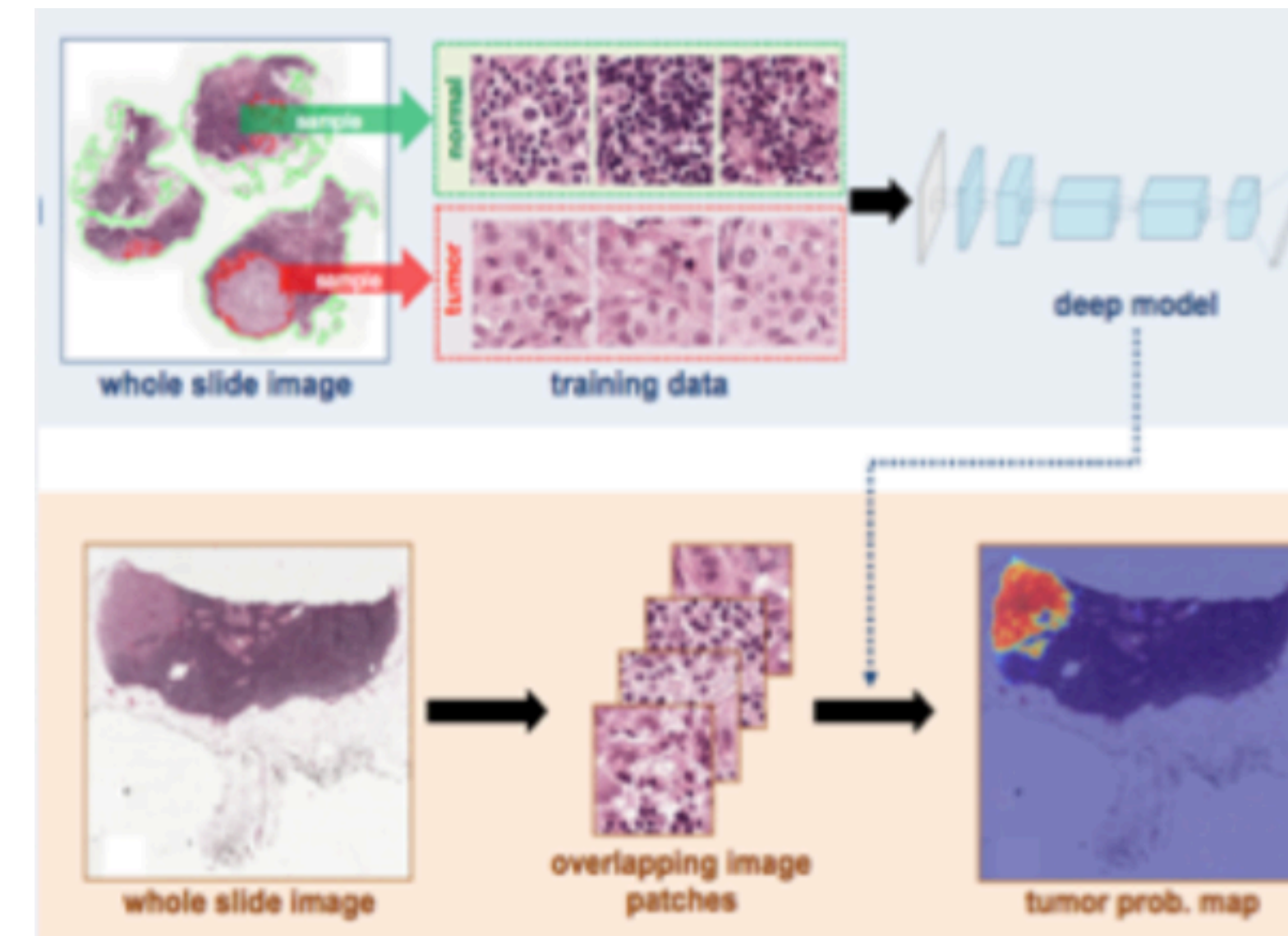
# Deep Nets are Everywhere



Self-driving Cars

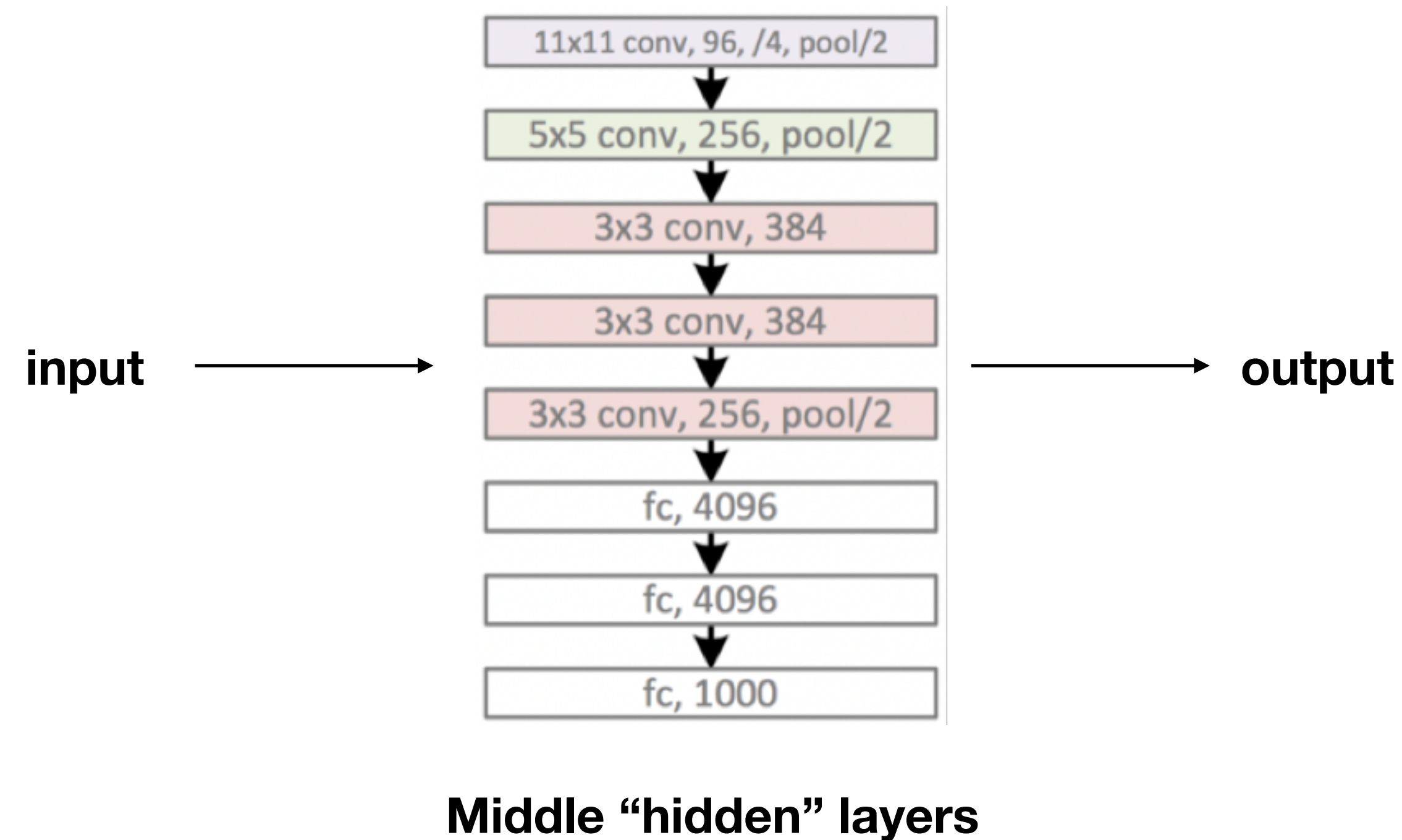


Playing Go



Making Medical Decisions

# Deep Nets are Not Understandable



Whenever correct: “whatever you did in the middle, do more.”

Whenever wrong: “whatever you did in the middle, do less.”

# Review of Research in XAI

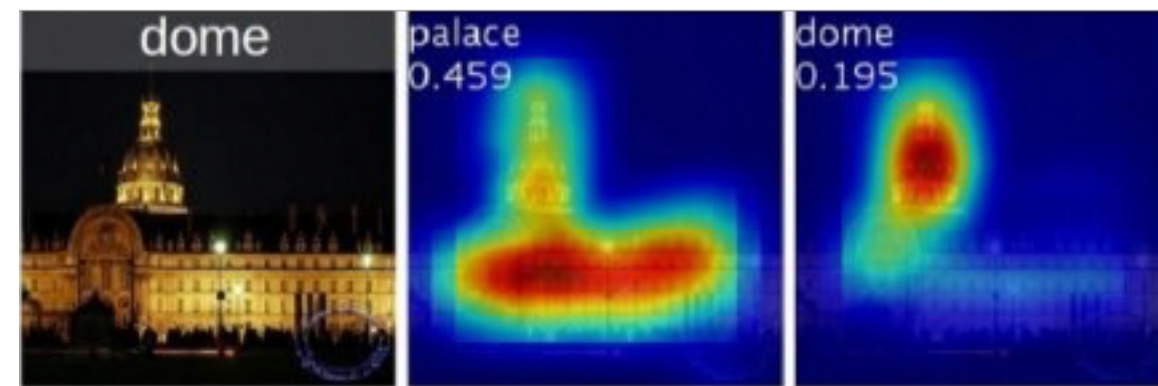
- Definitions
- Taxonomy
  - Survey: Literature review (87 papers) in computer science, artificial intelligence, and philosophy.
  - Recommendations for Evaluation
- How can explanations help (e.g. anomaly detection).
- Contributions and Future Work

# Definitions

- Explainability != Interpretability
- **Interpretability** describes the internals of a system that is *understandable* to humans.
- **Completeness** describes operation in an *accurate* way.
- An explanation needs **both**.

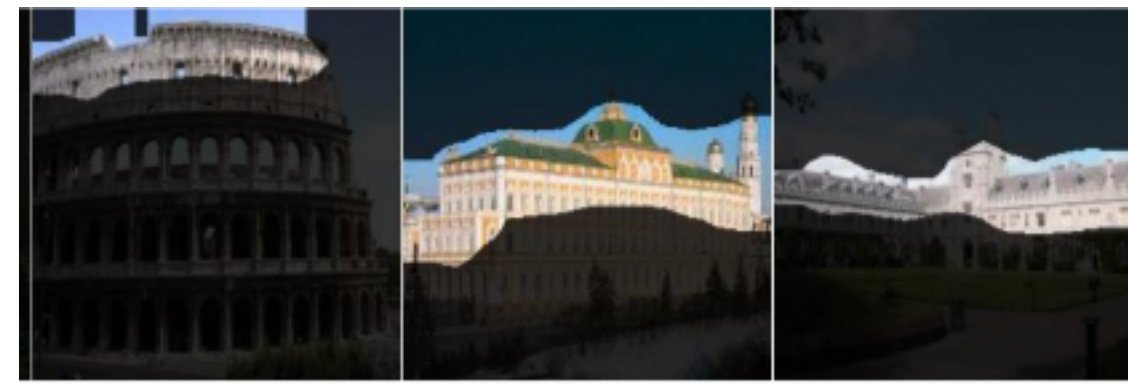
# What we Have

## Visual cues



Interpretable,  
**not complete**


## Role of individual units




Complete,  
**not interpretable**

## Attention based


*Q: Is this a healthy meal?*    Textual Justification    Visual Pointing



→ *A: No*    ...because it is a hot dog with a lot of toppings.



→ *A: Yes*    ...because it contains a variety of vegetables on the table.



Interpretable,  
**not complete**

# Why this Matters

## Interpretability

- GDPR
- Liability for decision making

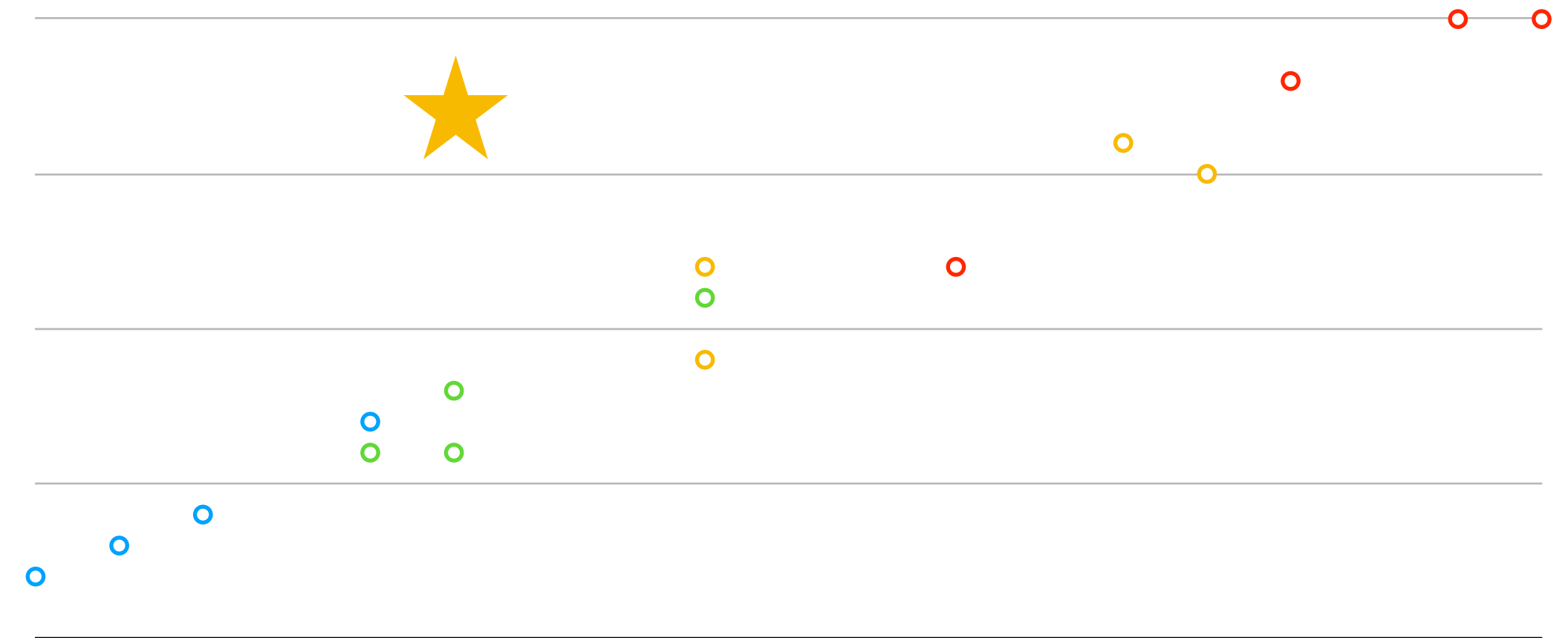


# Why this Matters

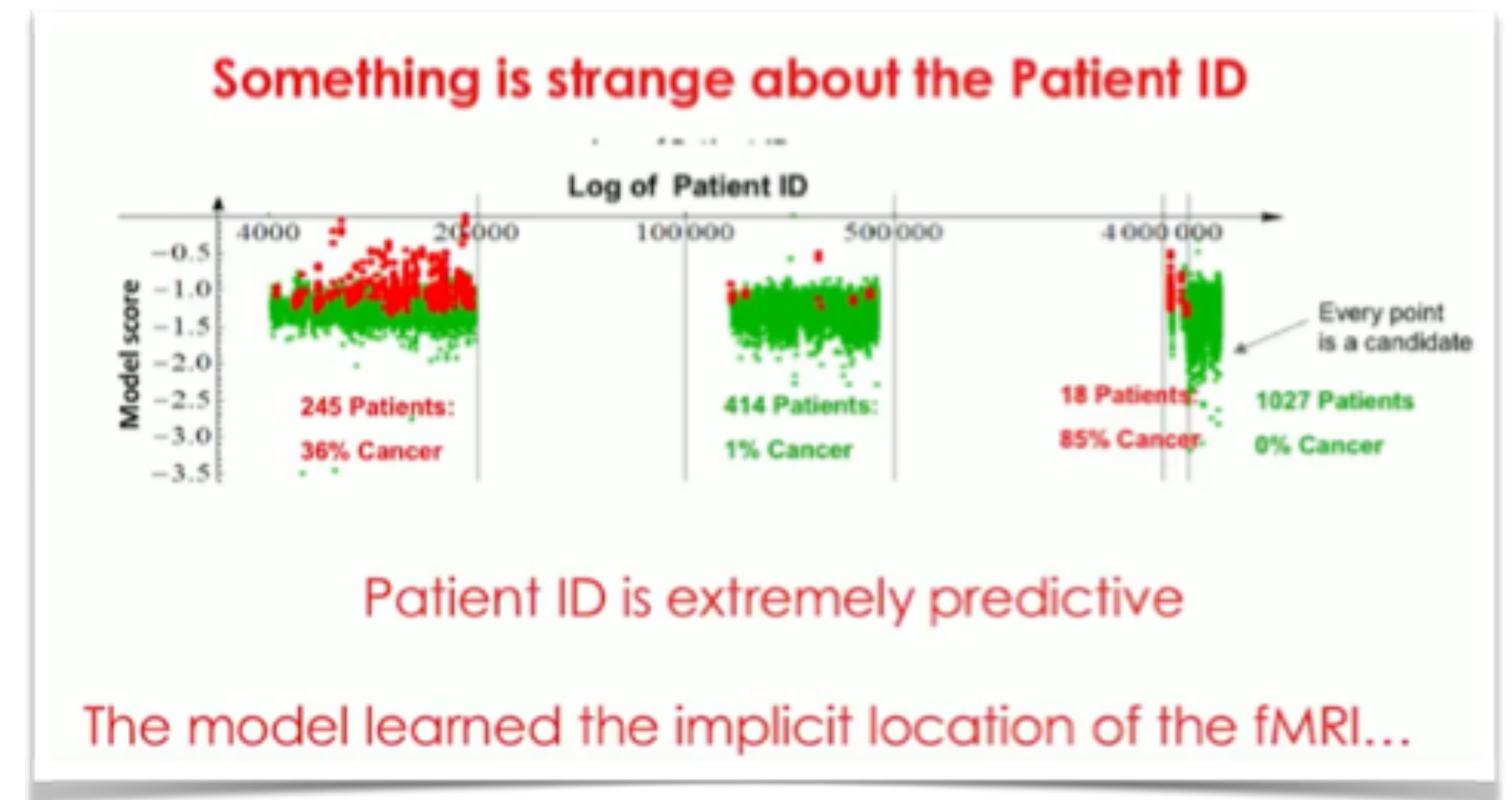
## Completeness

- Explaining the wrong thing.
- Making decisions for the wrong reasons.

Billing amount



Procedure code



From Claudia Perlich at *Women in Data Science 2018*.

# Talk Agenda

Motivate problem: Systems are imperfect

What is explainability?

What is *actually* being explained?

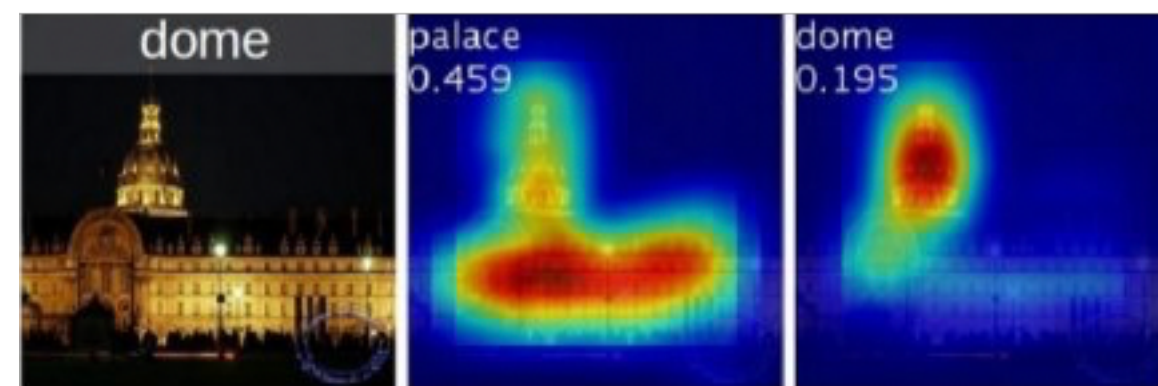
How to evaluate explainability?

How to explain complex systems? (autonomous driving)



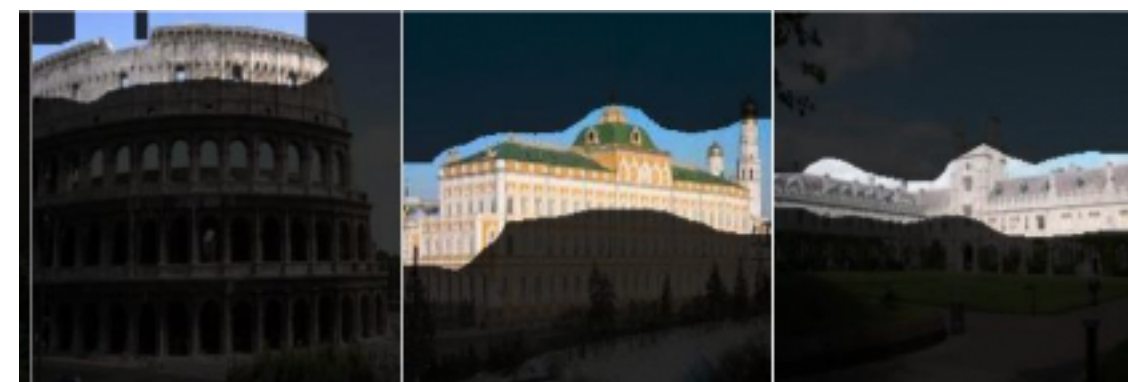
# What is Being Explained?

Visual cues



Explain  
processing

Role of individual  
units



Explain  
representation

Attention based

*Q: Is this a healthy meal?*    Textual Justification    Visual Pointing



→ *A: No*

*...because it  
is a hot dog  
with a lot of  
toppings.*



→ *A: Yes*

*...because it  
contains a  
variety of  
vegetables on  
the table.*



Explanation  
producing

# Taxonomy

	Processing	Representation	Explanation producing
Methods	Proxy Methods Decision Trees Salience Mapping Automatic-rule extraction	Role of layers Role of neurons Role of vectors	Scripted conversations Attention based Disentangled representations

# Methods that Explain Processing

## **DeepRED – Rule Extraction from Deep Neural Networks\***

Jan Ruben Zilke, Eneldo Loza Mencía, and Frederik Janssen

Technische Universität Darmstadt  
Knowledge Engineering Group

j.zilke@mail.de, {eneldo,janssen}@ke.tu-darmstadt.de

---

## **Extracting Rules from Artificial Neural Networks with Distributed Representations**

---

**Sebastian Thrun**  
University of Bonn  
Department of Computer Science III  
Römerstr. 164, D-53117 Bonn, Germany  
E-mail: thrun@carbon.informatik.uni-bonn.de

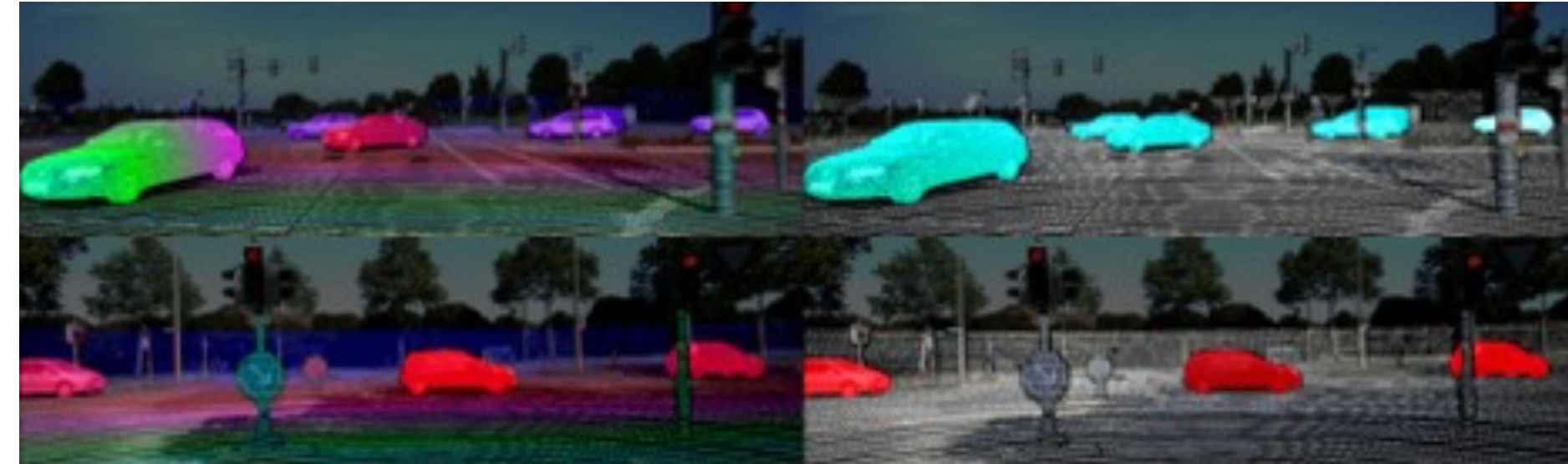
## **“Why Should I Trust You?” Explaining the Predictions of Any Classifier**

Marco Tulio Ribeiro  
University of Washington  
Seattle, WA 98105, USA  
marcotcr@cs.uw.edu

Sameer Singh  
University of Washington  
Seattle, WA 98105, USA  
sameer@cs.uw.edu

Carlos Guestrin  
University of Washington  
Seattle, WA 98105, USA  
guestrin@cs.uw.edu

# Examples of Processing Methods



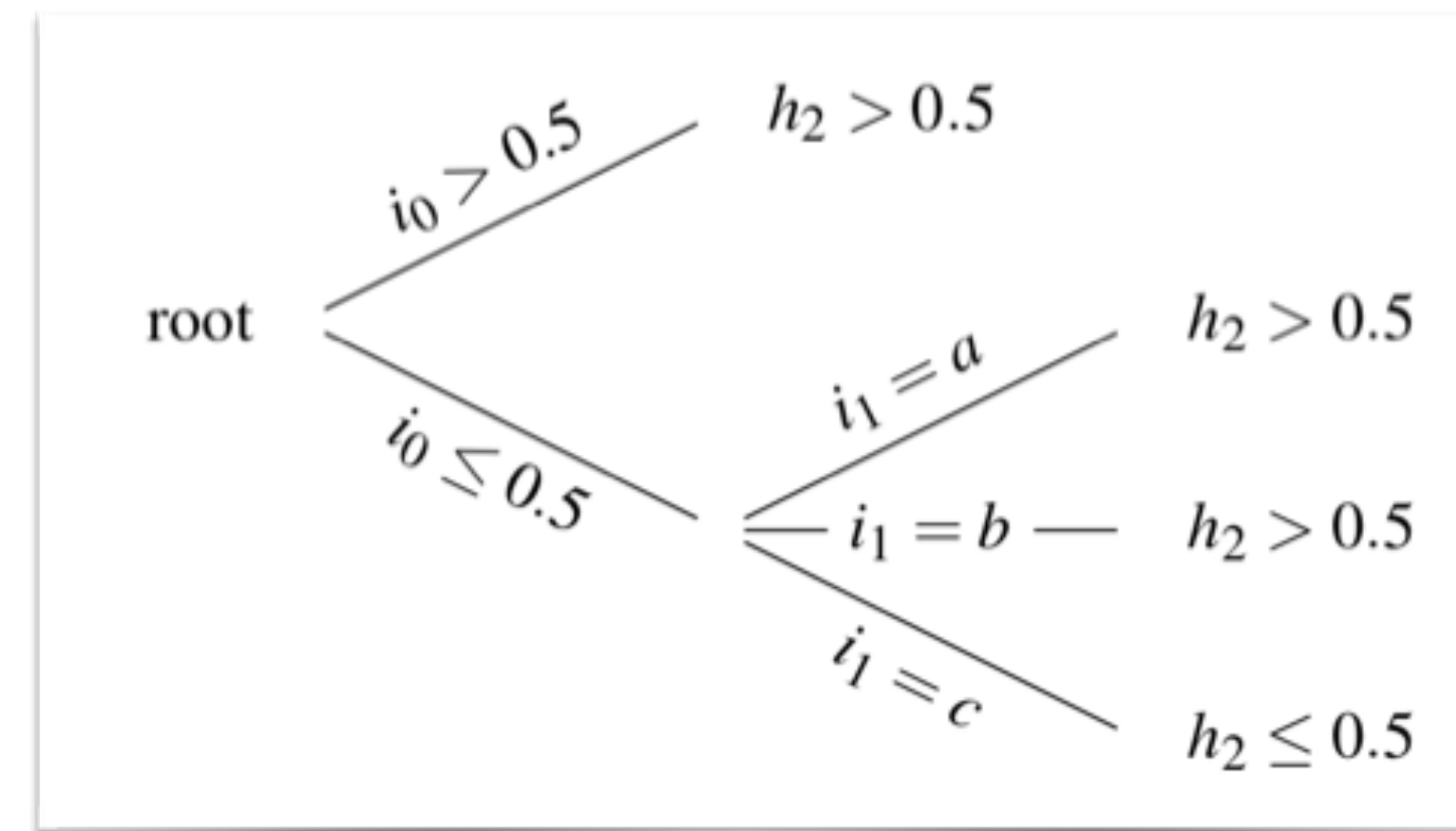
Geiger, Andreas, Philip Lenz, and Raquel Urtasun. "Are we ready for autonomous driving? The kitti vision benchmark suite." *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on. IEEE, 2012.*

## DeepRED – Rule Extraction from Deep Neural Networks\*

Jan Ruben Zilke, Eneldo Loza Mencía, and Frederik Janssen

Technische Universität Darmstadt  
Knowledge Engineering Group

j.zilke@mail.de, {eneldo,janssen}@ke.tu-darmstadt.de



Zilke, Jan Ruben et al. "DeepRED - Rule Extraction from Deep Neural Networks." *DS (2016).*

# Taxonomy

	Processing	Representation	Explanation producing
Methods	Proxy Methods Decision Trees Salience Mapping Automatic-rule extraction	Role of layers Role of neurons Role of vectors	Scripted conversations Attention based Disentangled representations

# Methods that Explain Representations

## **Network Dissection:**

### **Quantifying Interpretability of Deep Visual Representations**

David Bau\*, Bolei Zhou\*, Aditya Khosla, Aude Oliva, and Antonio Torralba  
CSAIL, MIT

{davidbau, bzhou, khosla, oliva, torralba}@csail.mit.edu

---

## **Interpretability Beyond Feature Attribution:**

### **Quantitative Testing with Concept Activation Vectors (TCAV)**

---

Been Kim Martin Wattenberg Justin Gilmer Carrie Cai James Wexler  
Fernanda Viegas Rory Sayres

## **CNN Features off-the-shelf: an Astounding Baseline for Recognition**

Ali Sharif Razavian Hossein Azizpour Josephine Sullivan Stefan Carlsson  
CVAP, KTH (Royal Institute of Technology)

Stockholm, Sweden

{razavian, azizpour, sullivan, stefanc}@csc.kth.se

# Examples of Explained Representations

**Network Dissection:  
Quantifying Interpretability of Deep Visual Representations**

David Bau\*, Bolei Zhou\*, Aditya Khosla, Aude Oliva, and Antonio Torralba  
CSAIL, MIT  
{davidbau, bzhou, khosla, oliva, torralba}@csail.mit.edu



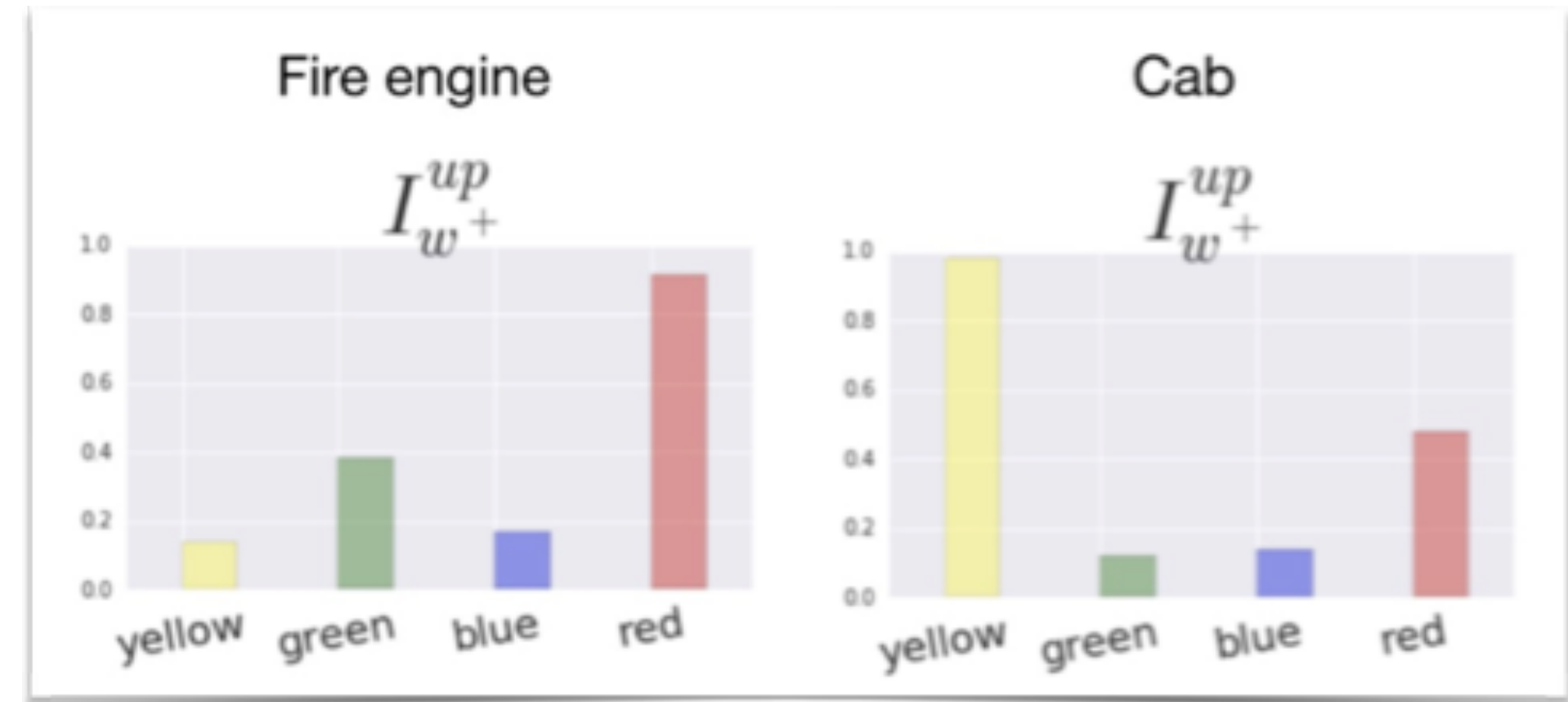
D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba, "Network dissection: Quantifying interpretability of deep visual representations," in *Computer Vision and Pattern Recognition*, 2017.

---

**Interpretability Beyond Feature Attribution:  
Quantitative Testing with Concept Activation Vectors (TCAV)**

---

Been Kim Martin Wattenberg Justin Gilmer Carrie Cai James Wexler  
Fernanda Viegas Rory Sayres



Kim, Been, et al. "Tcav: Relative concept importance testing with linear concept activation vectors." *arXiv preprint arXiv:1711.11279* (2017).

# Taxonomy

	Processing	Representation	Explanation producing
Methods	Proxy Methods Decision Trees Saliency Mapping Automatic-rule extraction	Role of layers Role of neurons Role of vectors	Scripted conversations Attention based Disentangled representations



# Methods that Produce Explanations

## Multimodal Explanations: Justifying Decisions and Pointing to the Evidence

Dong Huk Park<sup>1</sup>, Lisa Anne Hendricks<sup>1</sup>, Zeynep Akata<sup>2,3</sup>, Anna Rohrbach<sup>1,3</sup>,  
Bernt Schiele<sup>3</sup>, Trevor Darrell<sup>1</sup>, and Marcus Rohrbach<sup>4</sup>

<sup>1</sup>EECS, UC Berkeley, <sup>2</sup>University of Amsterdam, <sup>3</sup>MPI for Informatics, <sup>4</sup>Facebook AI Research

## Hierarchical Question-Image Co-Attention for Visual Question Answering

---

Jiasen Lu\*, Jianwei Yang\*, Dhruv Batra\*<sup>†</sup>, Devi Parikh\*<sup>†</sup>  
\* Virginia Tech, <sup>†</sup> Georgia Institute of Technology  
{jiasenlu, jw2yang, dbatra, parikh}@vt.edu

## InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets

---

Xi Chen<sup>†‡</sup>, Yan Duan<sup>†‡</sup>, Rein Houthoofd<sup>†‡</sup>, John Schulman<sup>†‡</sup>, Ilya Sutskever<sup>‡</sup>, Pieter Abbeel<sup>†‡</sup>  
<sup>†</sup> UC Berkeley, Department of Electrical Engineering and Computer Sciences  
<sup>‡</sup> OpenAI

# Examples that Produce Explanations

## Multimodal Explanations: Justifying Decisions and Pointing to the Evidence

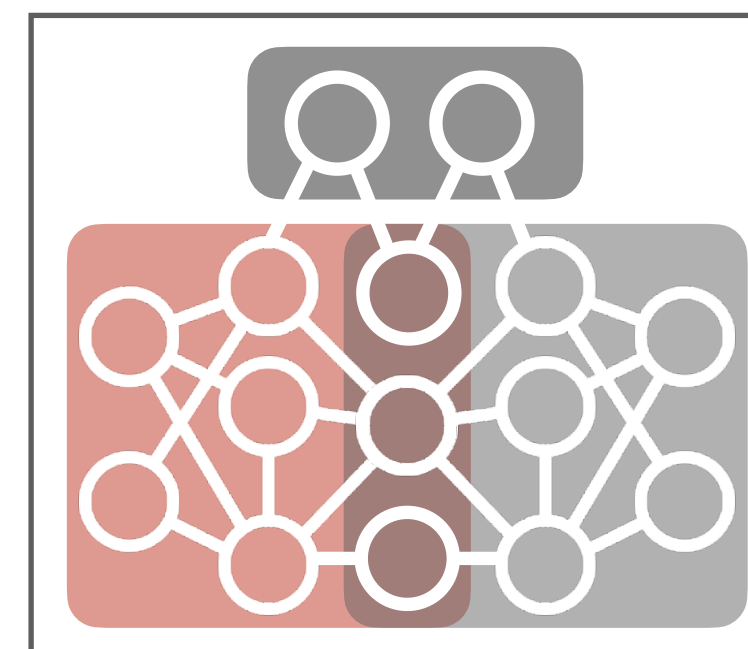
Dong Huk Park<sup>1</sup>, Lisa Anne Hendricks<sup>1</sup>, Zeynep Akata<sup>2,3</sup>, Anna Rohrbach<sup>1,3</sup>,  
Bernt Schiele<sup>3</sup>, Trevor Darrell<sup>1</sup>, and Marcus Rohrbach<sup>4</sup>

<sup>1</sup>EECS, UC Berkeley, <sup>2</sup>University of Amsterdam, <sup>3</sup>MPI for Informatics, <sup>4</sup>Facebook AI Research



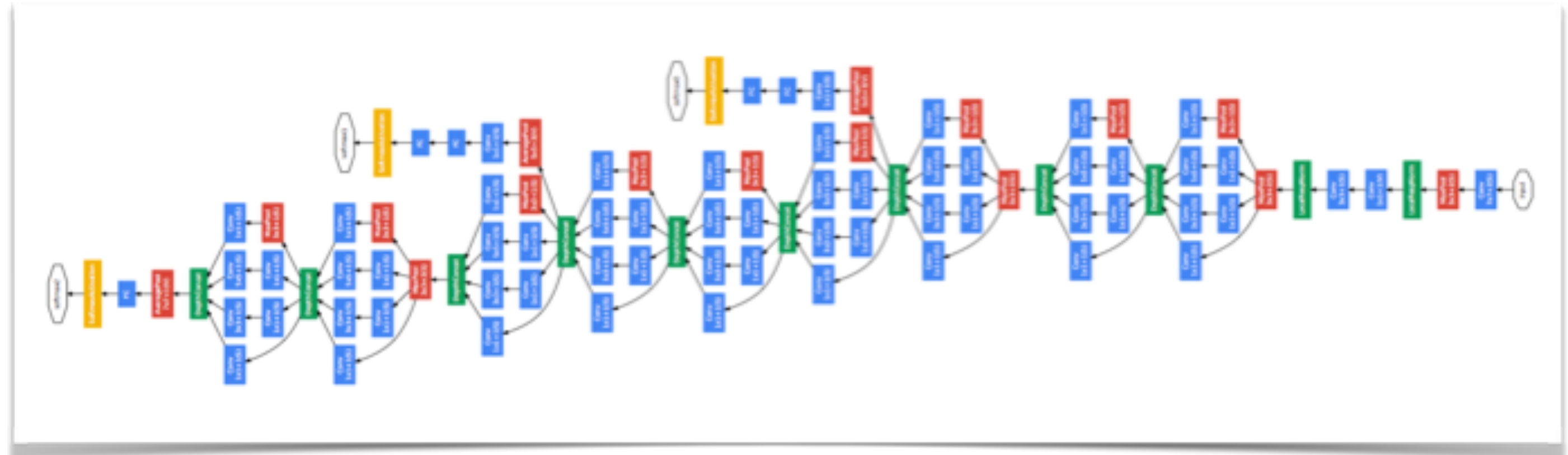
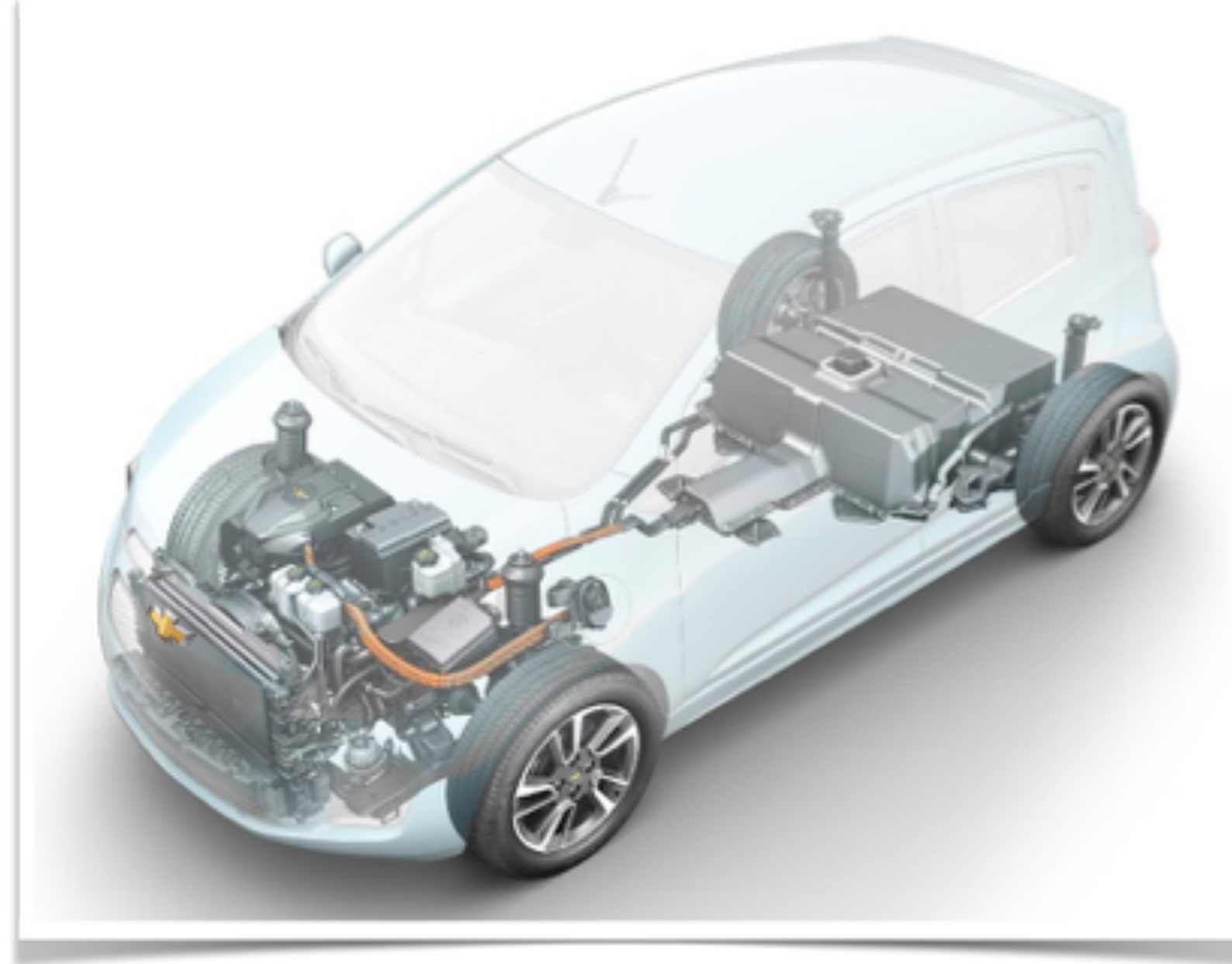
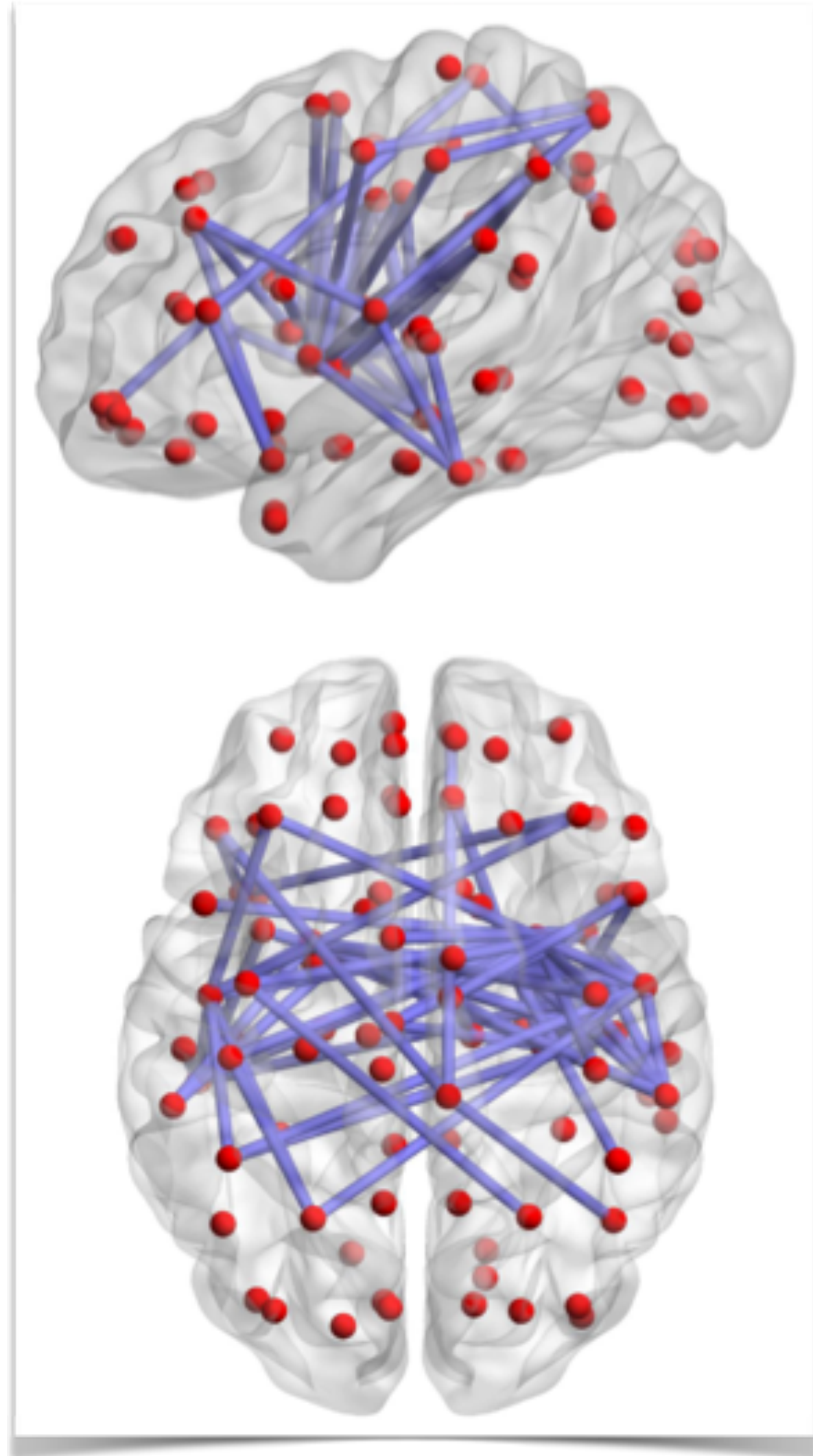
Park, Dong Huk, et al. "Multimodal Explanations: Justifying Decisions and Pointing to the Evidence." *31st IEEE Conference on Computer Vision and Pattern Recognition*. 2018.

- [1] L.H. Gilpin. Explaining possible futures for robust autonomous decision-making. Proceedings of the AAAI Fall Symposium on Anticipatory Thinking, 2019.
- [2] L.H. Gilpin, V. Penubarthi, and L. Kagal. Explaining Multimodal Errors in Autonomous Vehicles. DSAA 2021.



The best option is to veer and slow down. The vehicle is traveling too fast to suddenly stop. The vision system is inconsistent, but the lidar system has provided a reasonable and strong claim to avoid the object moving across the street.

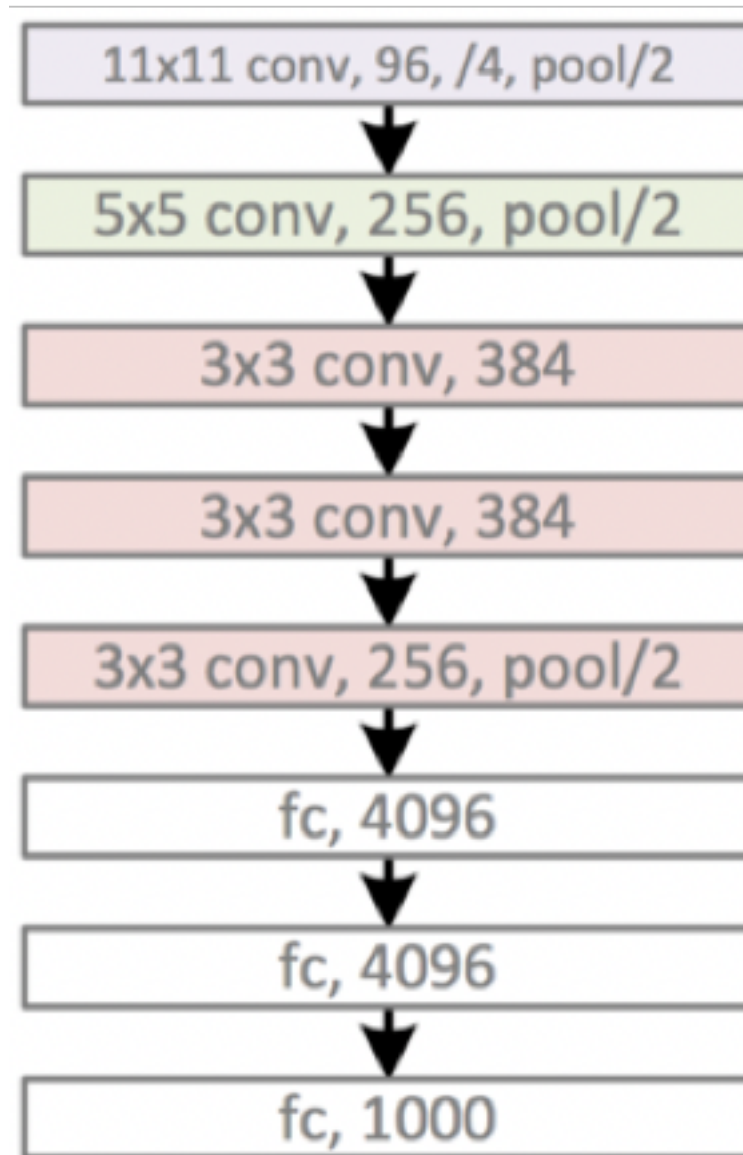
# A Problem: Insides Matter



# The More Complex (Deeper)

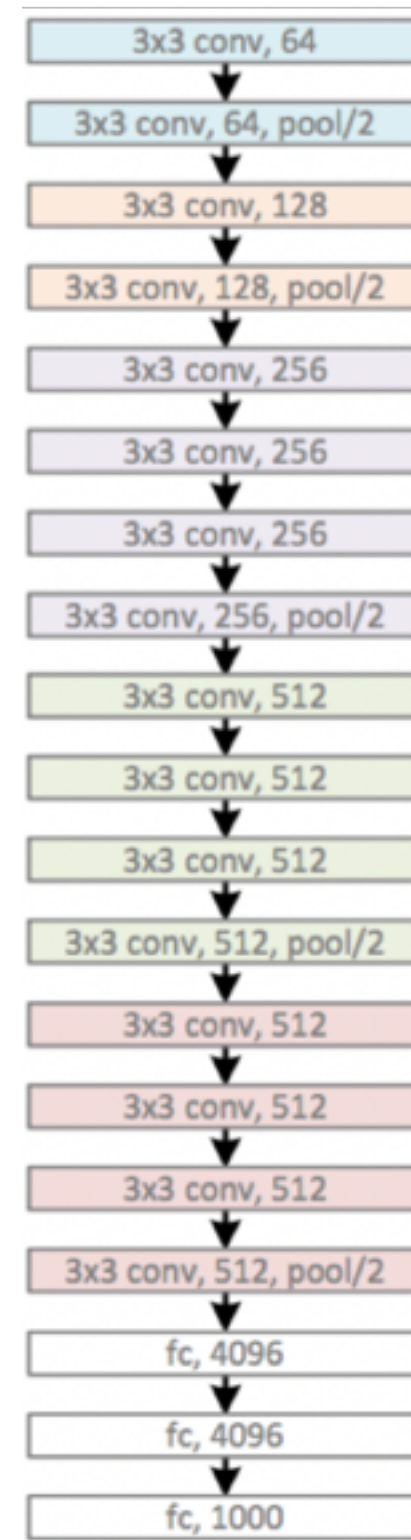
## The Deeper the Mystery

IMAGENET



AlexNet (2012)

8 layers; acc 84.7%



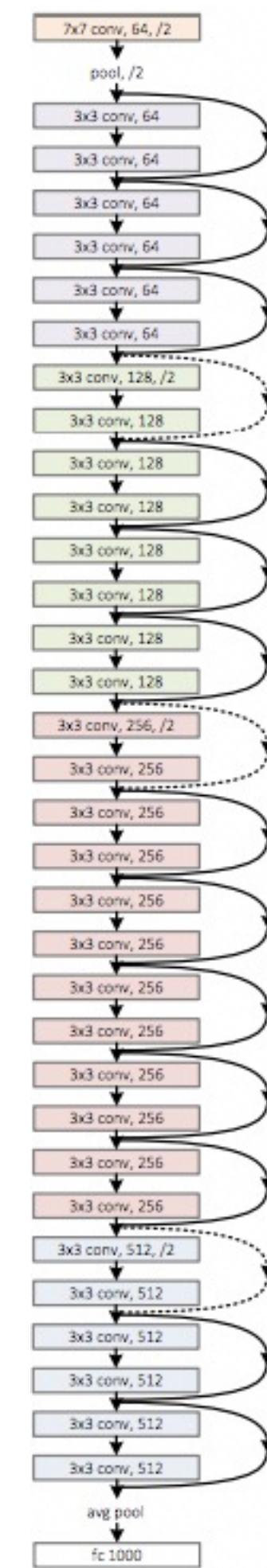
VGG (2014)

19 layers; acc 91.5%



GoogLeNet (2015)

22 layers; acc 92.2%



ResNet (2016)

152 layers; acc 95.6%

# Talk Agenda

Motivate problem: Systems are imperfect

What is explainability?

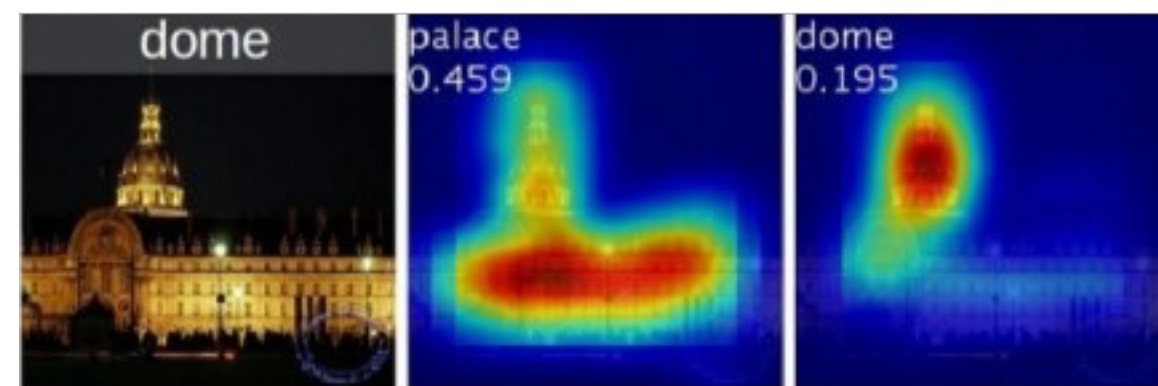
What is *actually* being explained?

How to evaluate explainability?

How to explain complex systems? (autonomous driving)

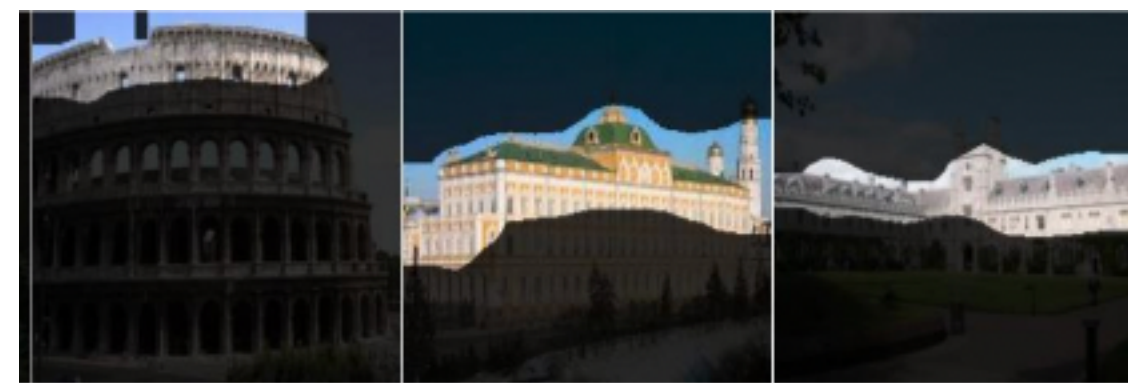
# What is Being Explained?

Visual cues



Completeness to model

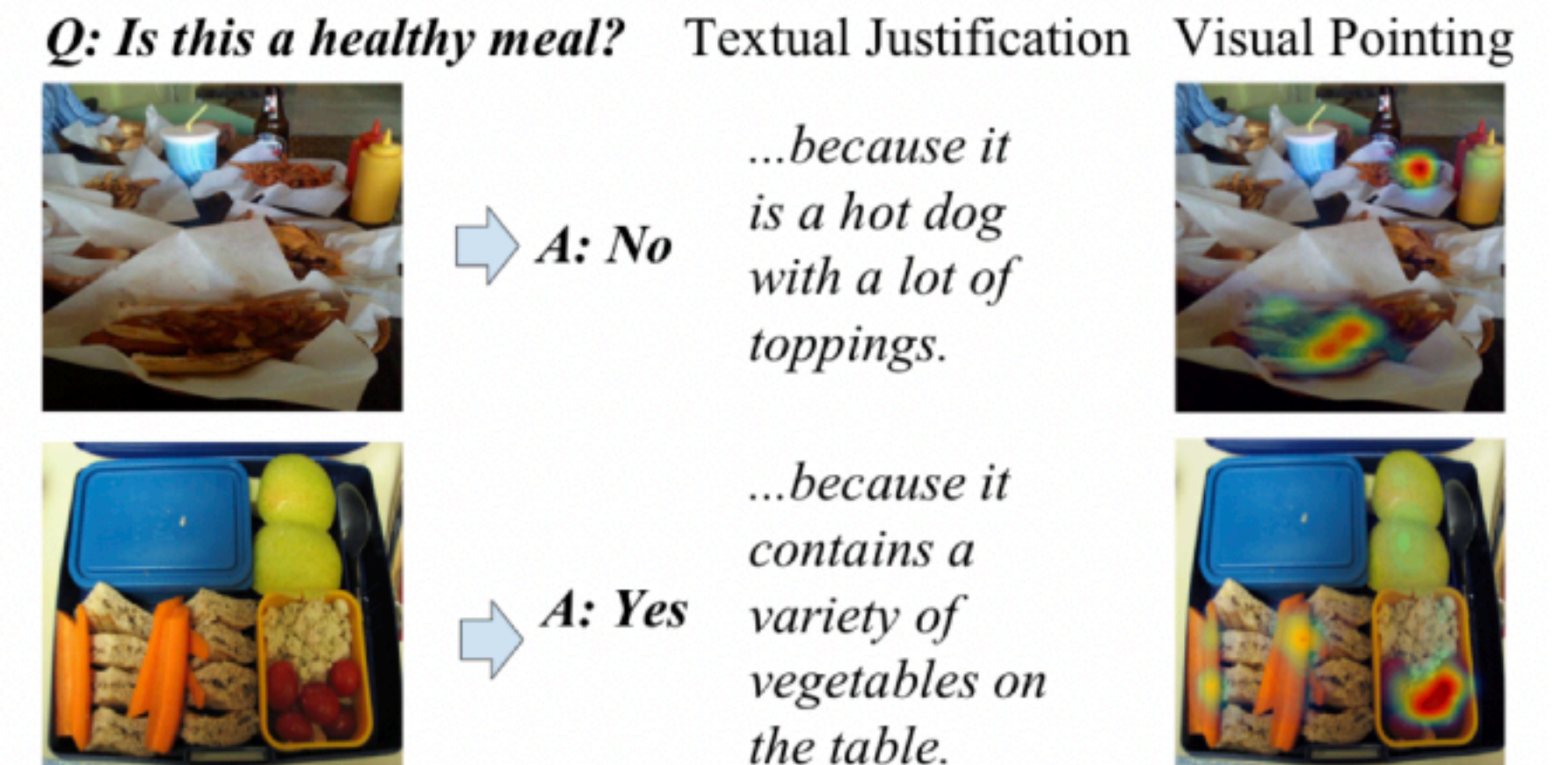
Role of individual units



Completeness on other tasks

Attention based

*Q: Is this a healthy meal?*    Textual Justification    Visual Pointing



→ *A: No*    ...because it is a hot dog with a lot of toppings.

→ *A: Yes*    ...because it contains a variety of vegetables on the table.

Human evaluation

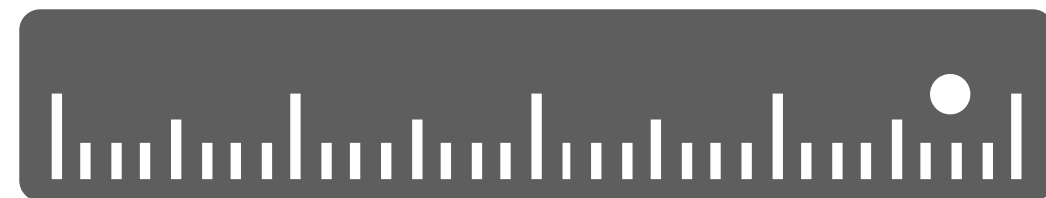
# Taxonomy

	Processing	Representation	Explanation producing
Methods	Proxy Methods Decision Trees Salience Mapping Automatic-rule extraction	Role of layers Role of neurons Role of vectors	Scripted conversations Attention based Disentangled representations

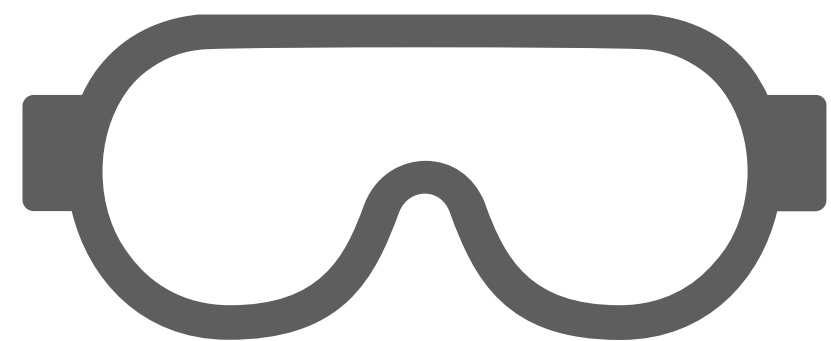
# Challenges in Explainability



- Standards and metrics for explanations
  - How to **evaluate** explanations?



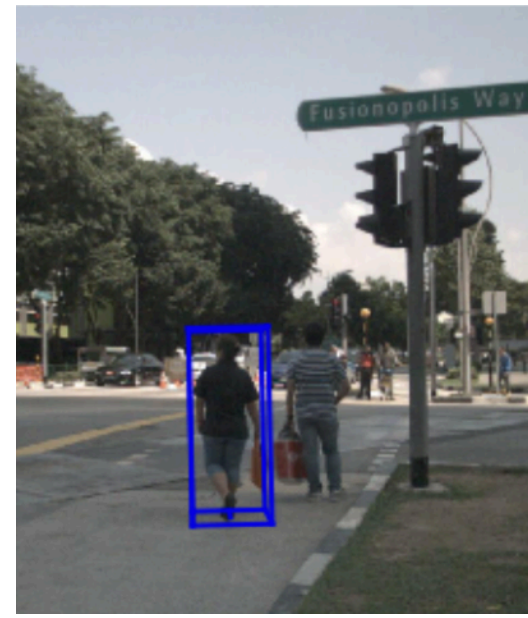
- Current metrics of evaluation are “fuzzy”
  - User based evaluations are not *always* appropriate



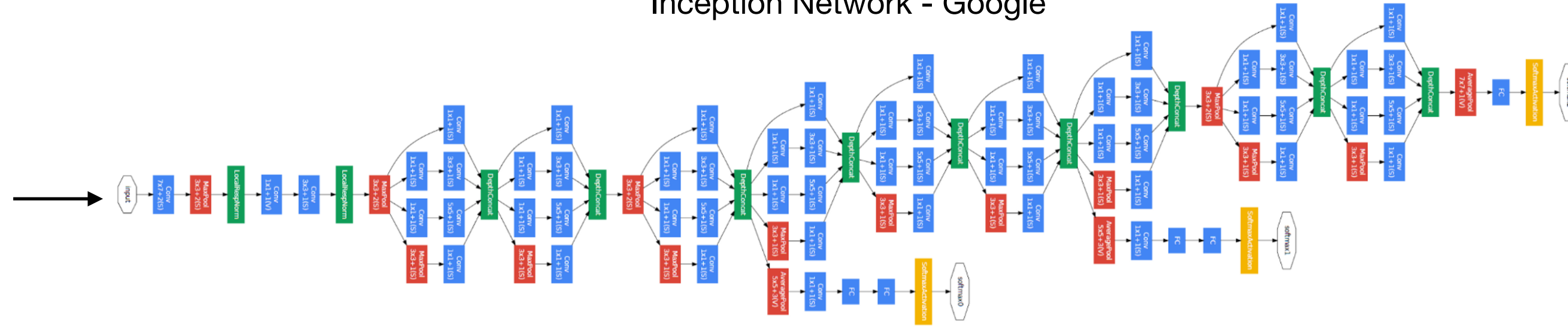
- Benchmarks for safety-critical and mission-critical tasks.



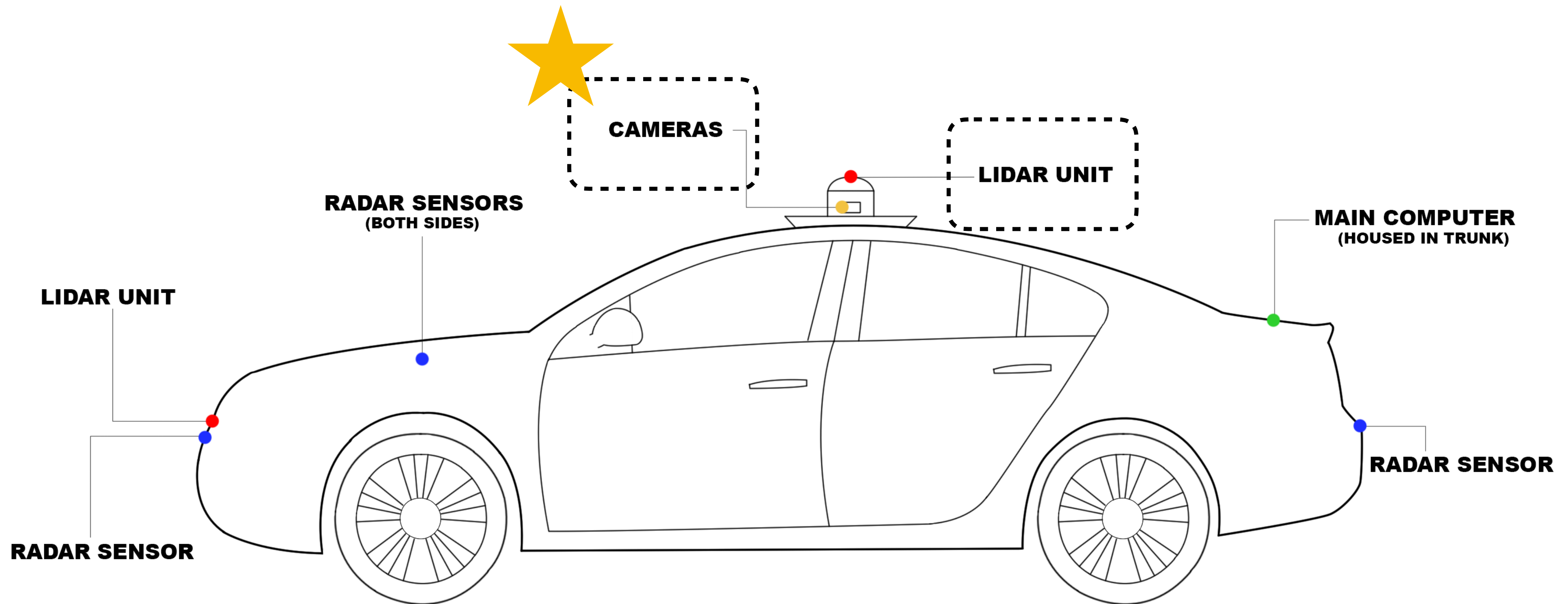
# A Neural Network Labels Camera Data



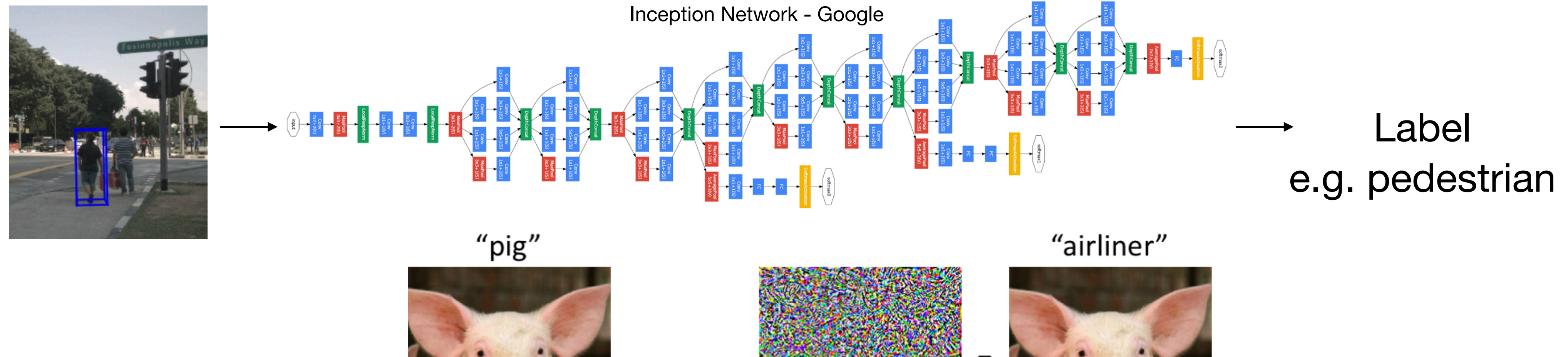
Inception Network - Google



→ Label  
e.g. pedestrian



# Problem: Neural Networks are Brittle



For self-driving, and other mission-critical, safety-critical applications, these mistakes have CONSEQUENCES.



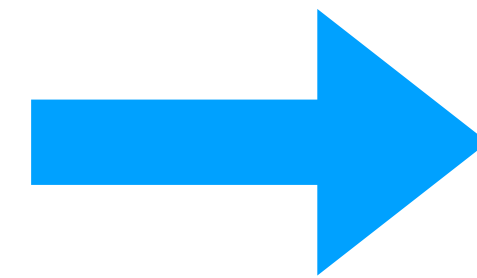
# Vision: Real World Adversarial Examples



“Realistic” Adversarial examples

# Vision: Real World Adversarial Examples

## Anticipatory Thinking Layer for Error Detection

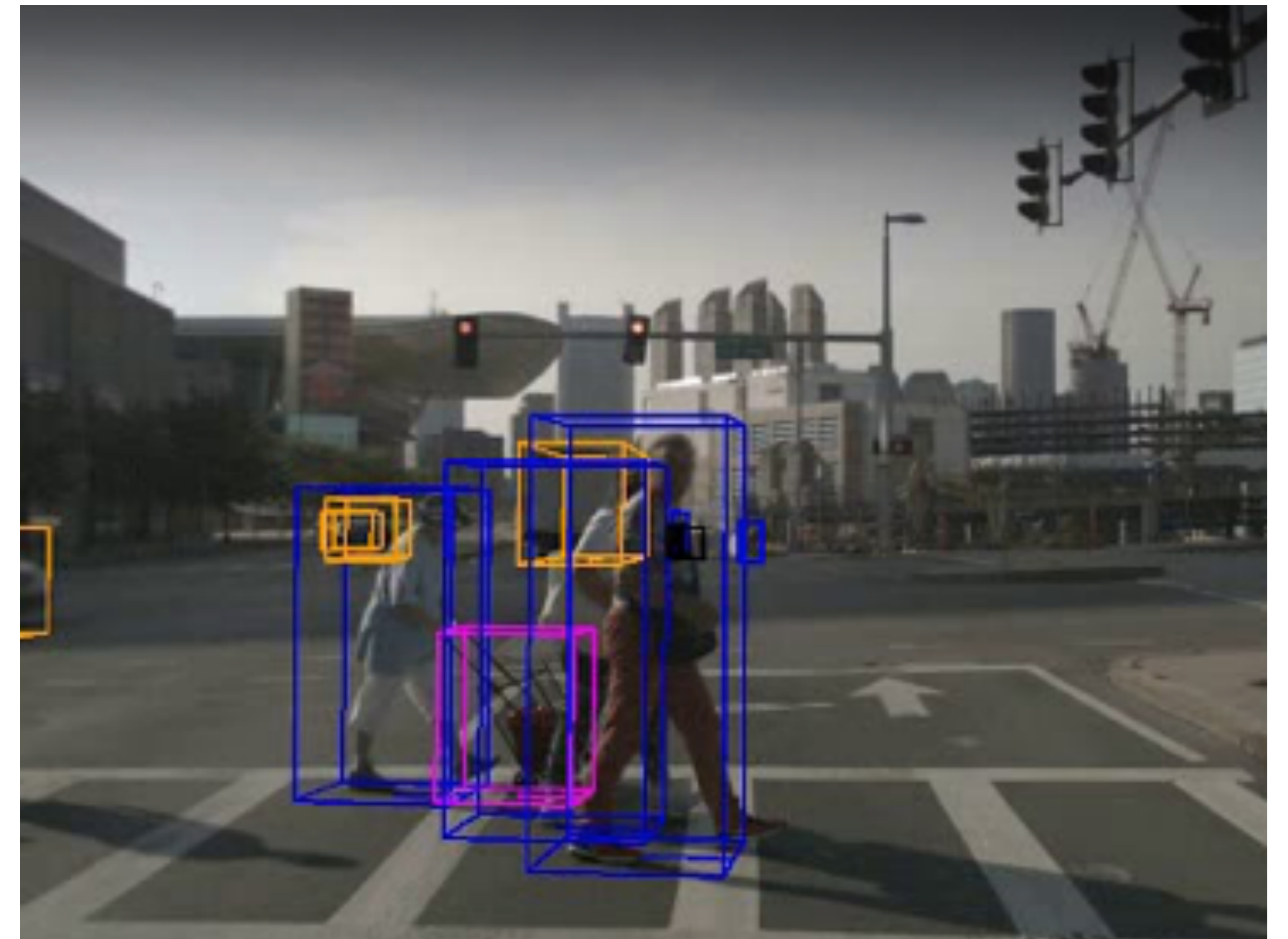


The traffic lights are on top of the truck. The lights are not illuminated. The lights are moving at the same rate as the truck, therefore this is not a “regular” traffic light for slowing down and stopping at.

“Realistic” Adversarial examples

# Lack of Data and Challenges for AVs

- Existing Challenges
  - Targeted as optimizing a mission or trajectory and not safety.
  - Data is hand-curated.
- Failure data is not available
  - Unethical to get it (cannot just drive into bad situations).
  - Want the data to be realistic (usually difficult in simulation).



Data from NuScenes

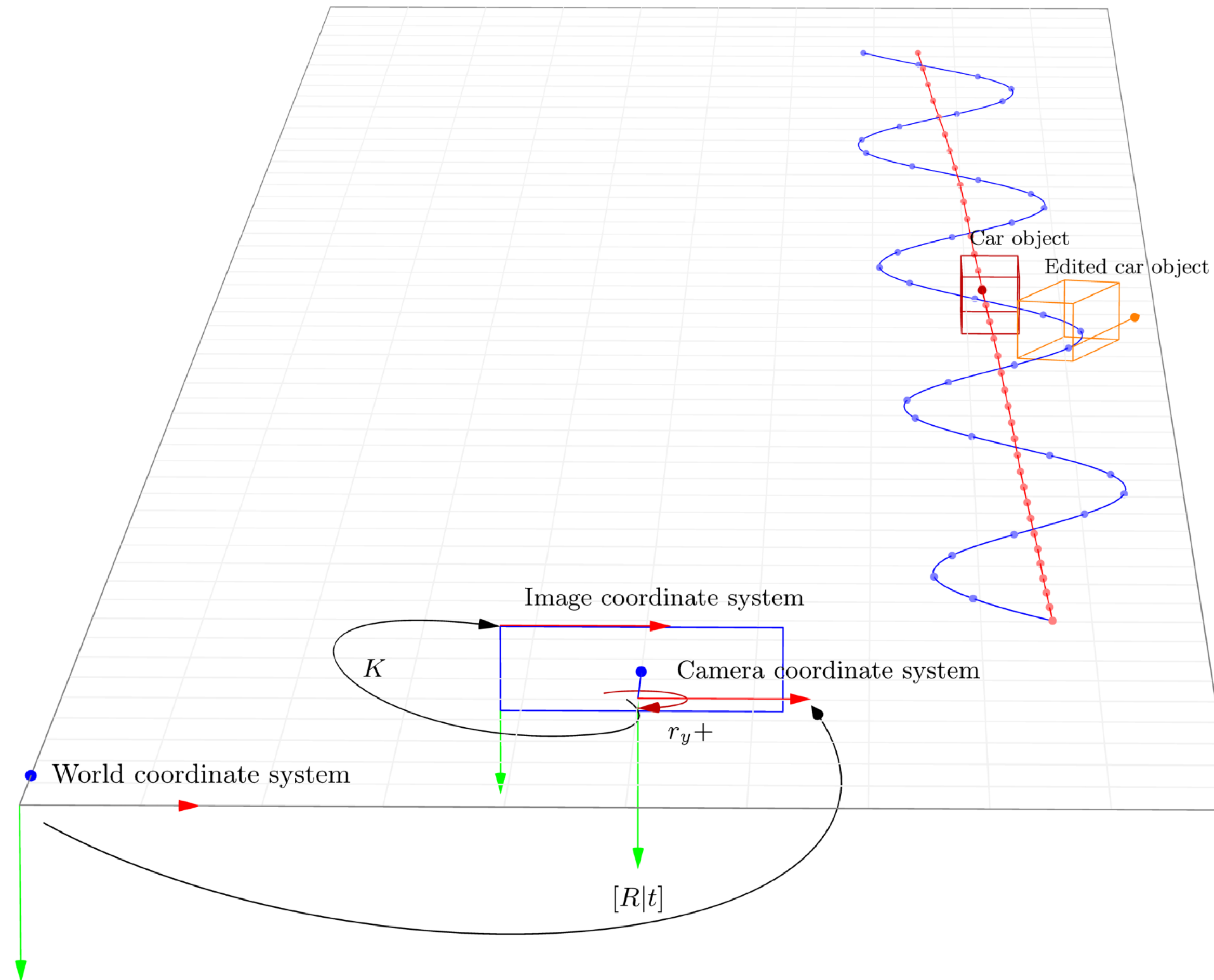
# Approach: Content Generation

## Anticipatory Thinking Layer for Error Detection

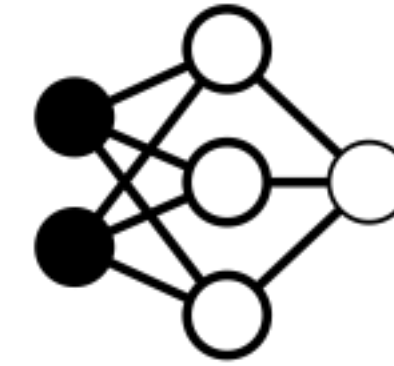


# Approach: Content Generation

## Anticipatory Thinking Layer for Error Detection



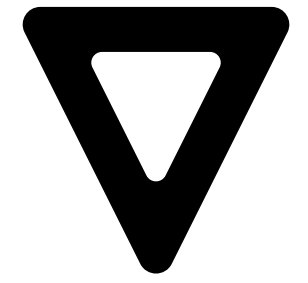
# Contributions



Opaque Systems



Autonomous Systems



Error Detection

Brief Intro

Motivate problem: Systems are imperfect

What is explainability?

What is *actually* being explained?

How to evaluate explainability?

How to critical systems? (autonomous driving)

