

# MATLAB EXPO

## 2021

### 인공지능 해부하기 : 설명 가능한 인공지능(eXplainable AI)

송완빈 과장



# AI models you interact with regularly

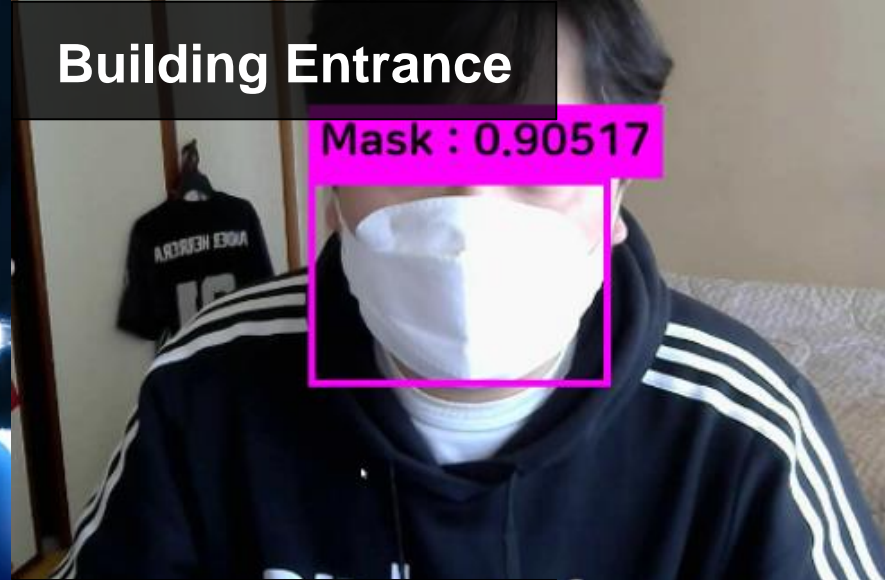
## Activity Monitoring



## Smart Devices



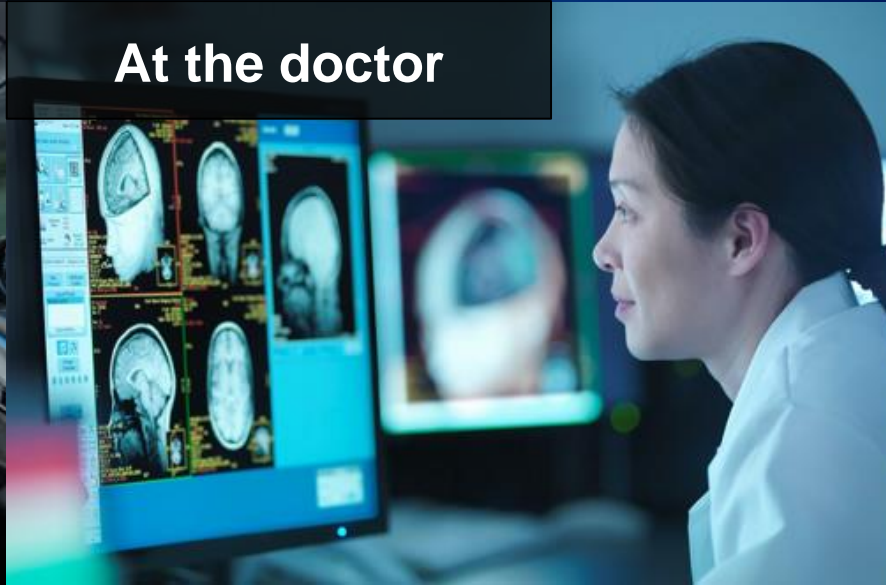
## Building Entrance



## Self-driving car



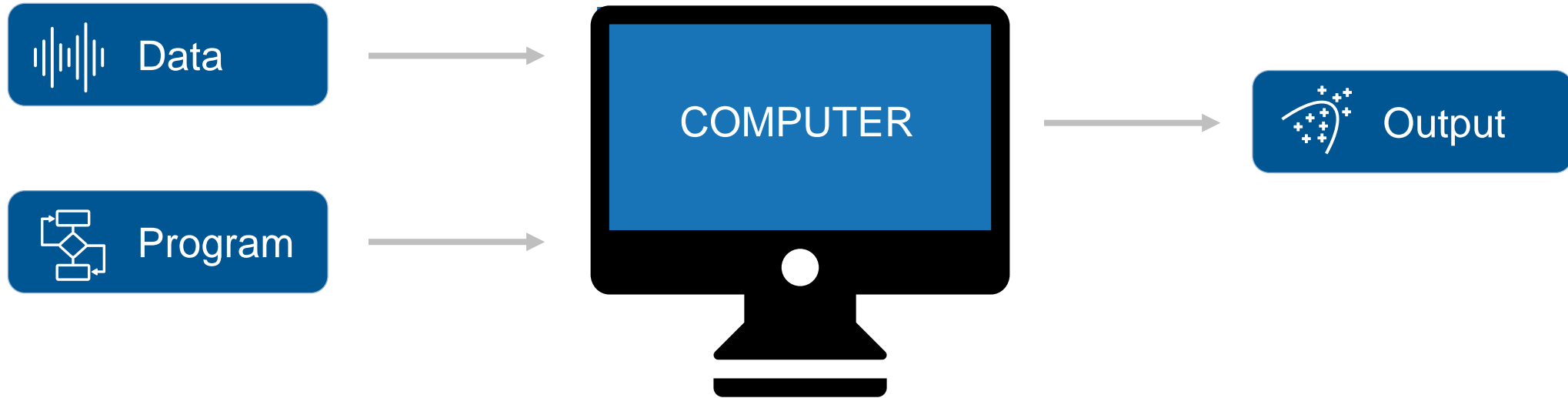
## At the doctor



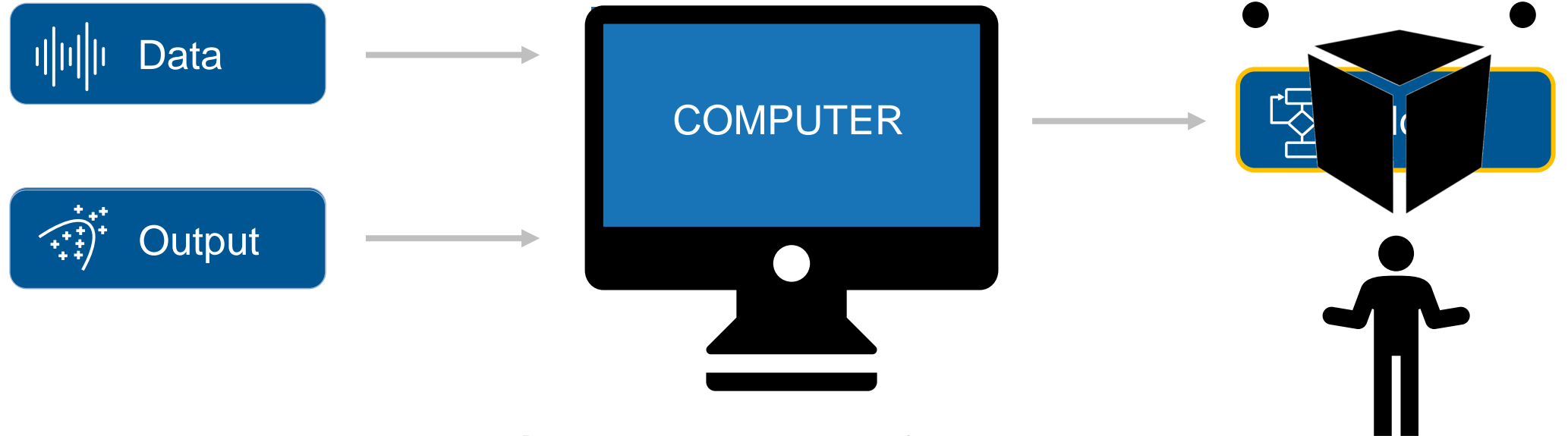
## Smart Home



# Traditional Software Development



# Machine Learning Software Development

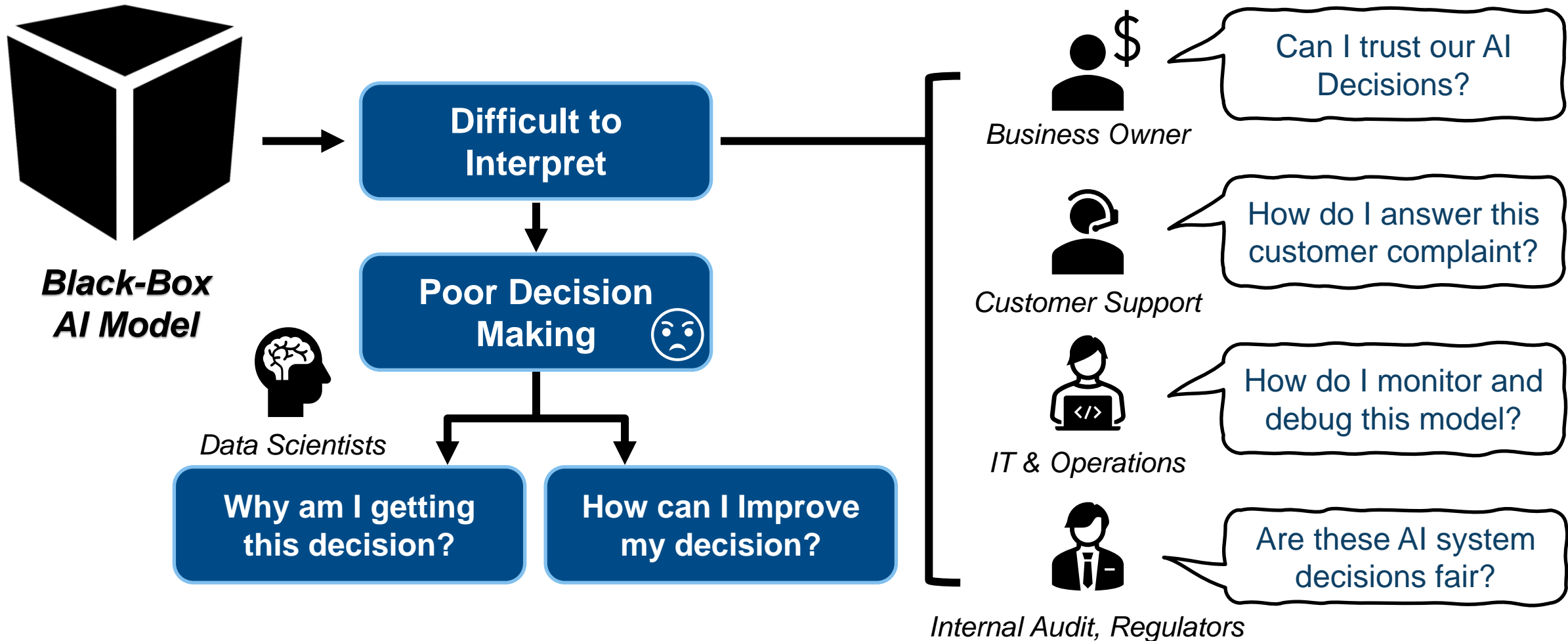


- ✓ Machine Learning models **provide the best results** for many tasks



- ✓ Desire to put ML models in production in safety critical situations.

# Problems with Black-Box AI Models



# What is “eXplainable AI”?

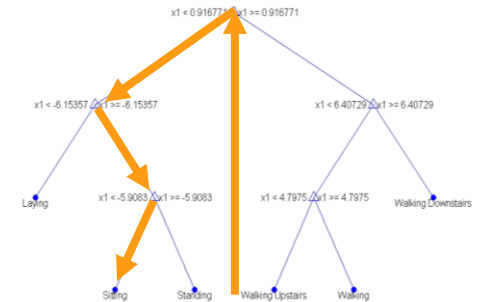
- Methods or techniques used to explain machine learning reasoning to humans



***Black-Box  
AI Model***

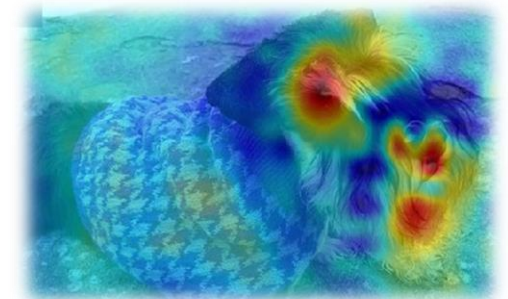
- Interpretability**

- Discern the mechanics without necessarily knowing why



- Explainability**

- Being able to quite literally explain what is happening



# The goal of eXplainable AI

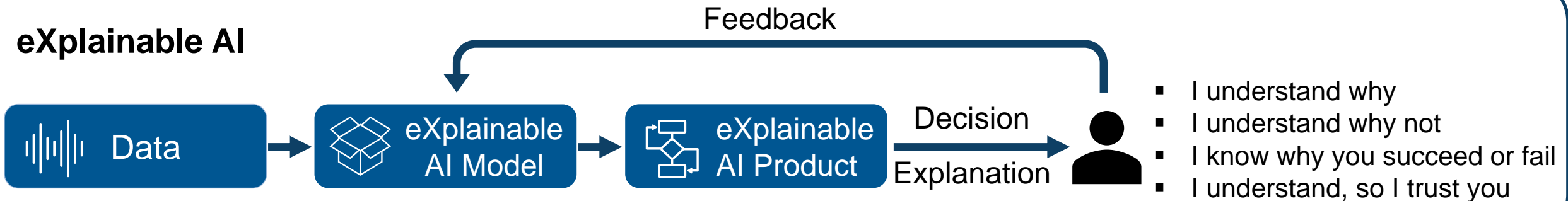
## Today

### Black Box AI

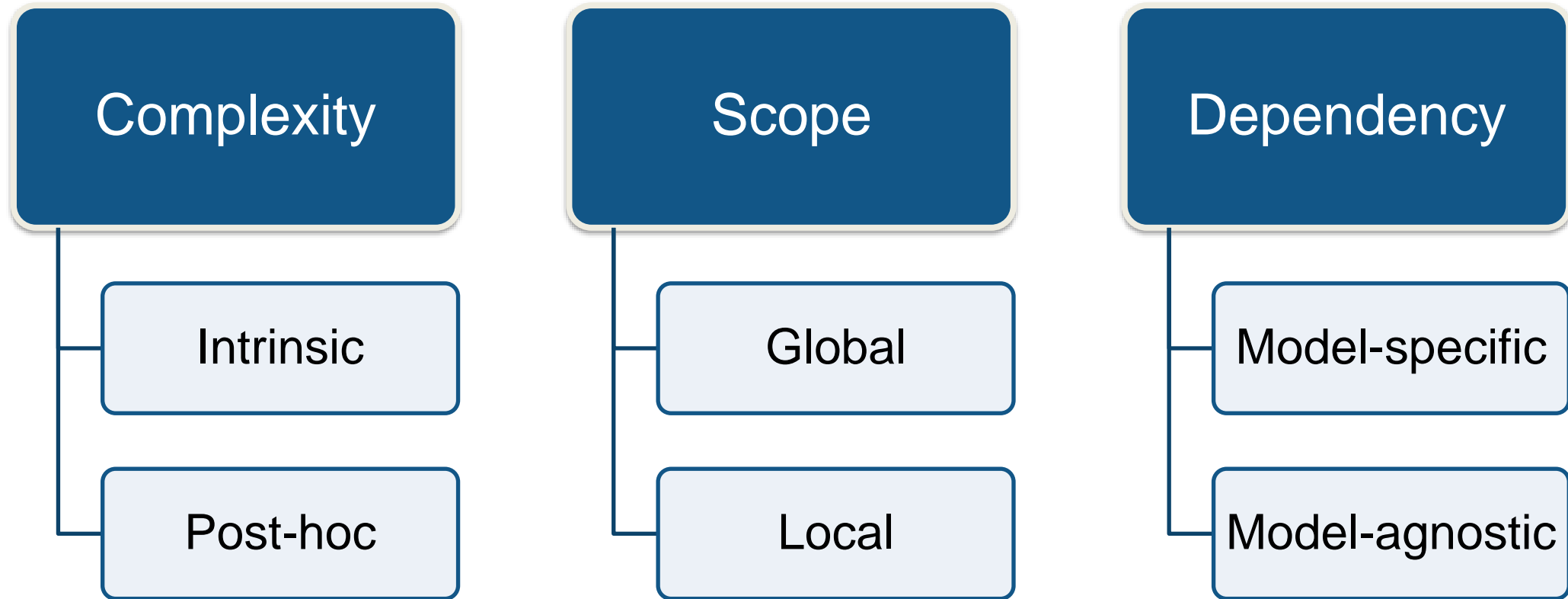


## Tomorrow

### eXplainable AI



# Categorize with respect to the viewpoint of eXplainable AI





# Viewpoint of eXplainable AI

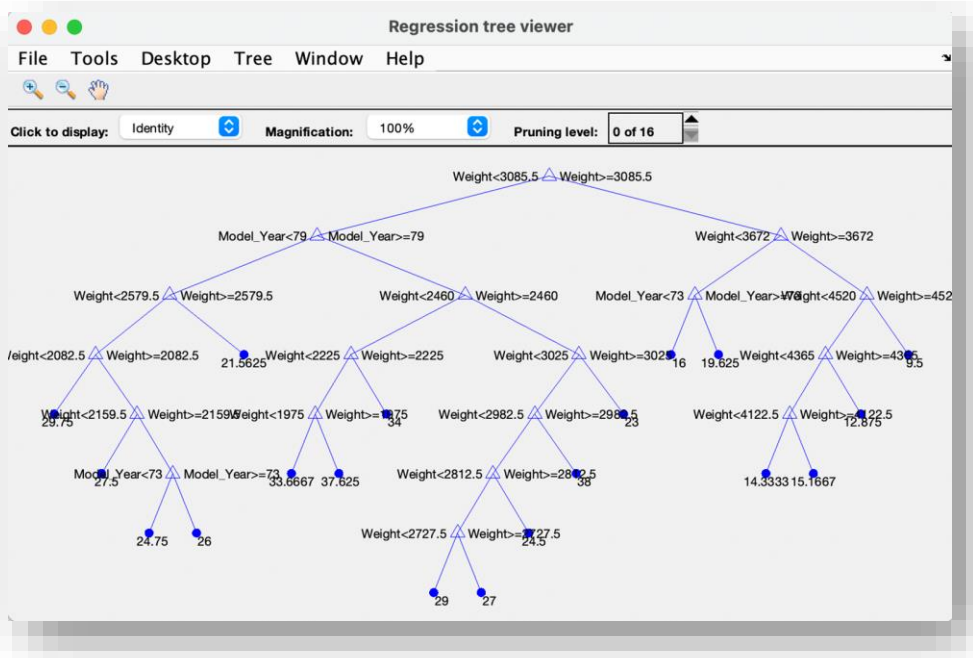
Complexity

Intrinsic

Post-hoc

- Intrinsic / Transparency

- Model itself already have transparency and interpretable



Linear regression model:  
 $Model\_Year \sim 1 + Weight + MPG$

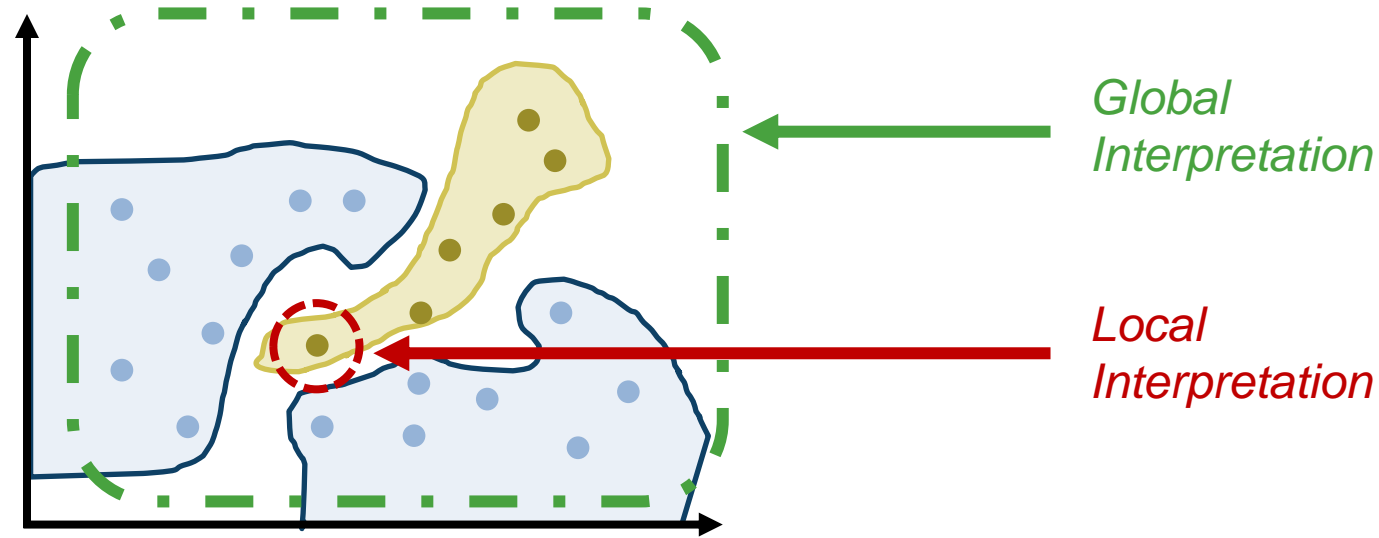
Estimated Coefficients:

	Estimate	SE	tStat	pValue
<b>(Intercept)</b>	49.758	4.1453	12.003	1.8646e-20
<b>Weight</b>	0.0032738	0.00080383	4.0728	9.9032e-05
<b>MPG</b>	0.70279	0.080191	8.7639	9.8456e-14

- Post-hoc

- Normally use complex model which has large predictive power and interpret later

# Viewpoint of eXplainable AI



Scope

Global

Local

- Global
  - Provides an overview of the most influential variables in the model, based on the data input and the predicted variable.
- Local
  - Explains conditional interaction between input/output with respect to single prediction result

# Viewpoint of eXplainable AI

Dependency

Model-specific

Model-agnostic

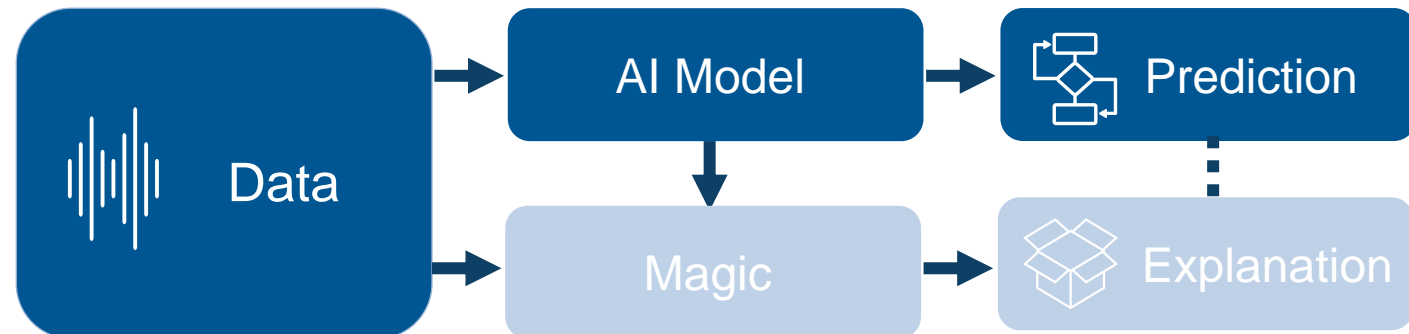
- Model-specific

- Deals with inner working of model to interpret its result

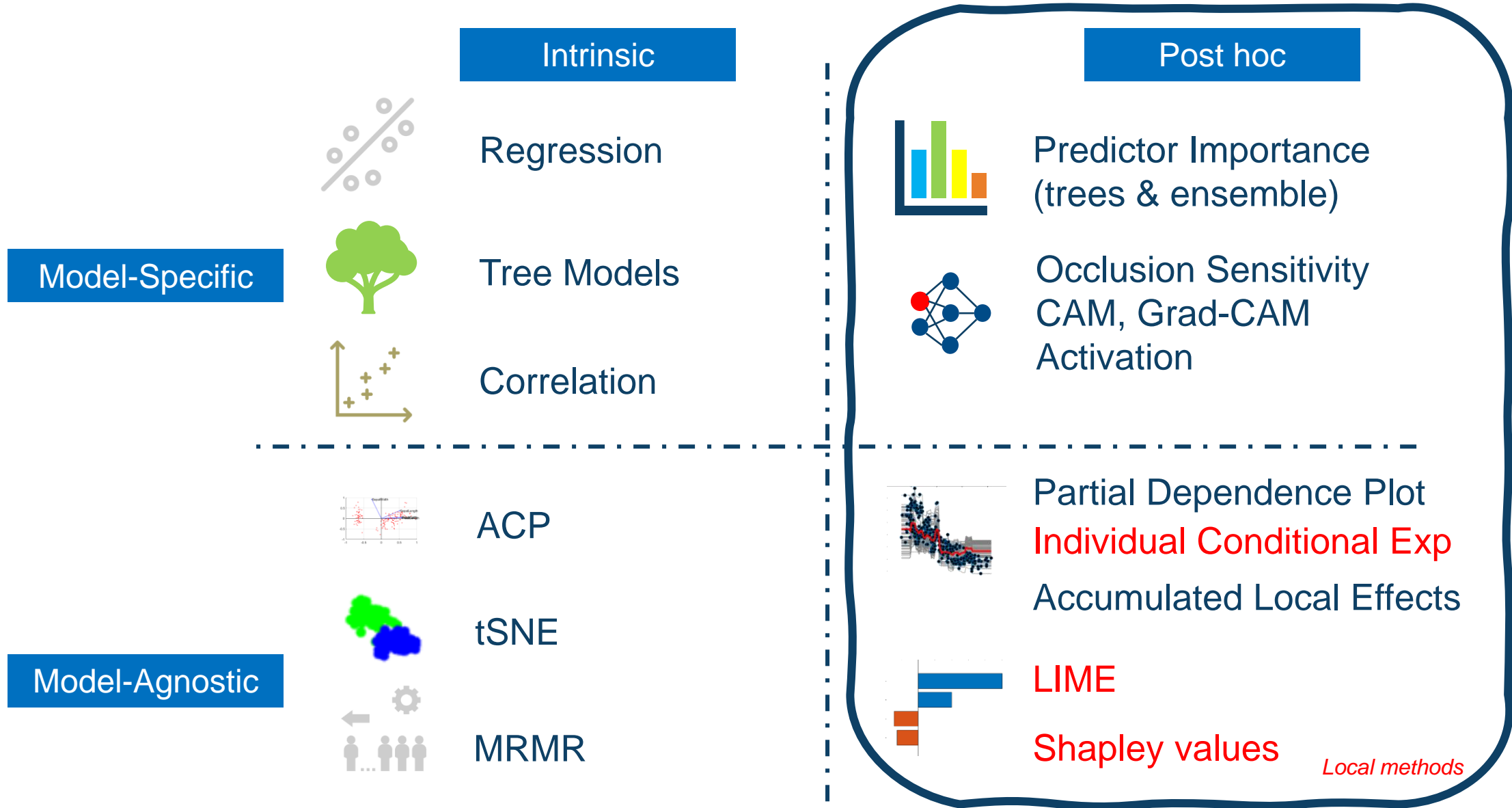


- Model-agnostic

- Deals with analyzing the feature and its relationships with its output



# Different Interpretability methods



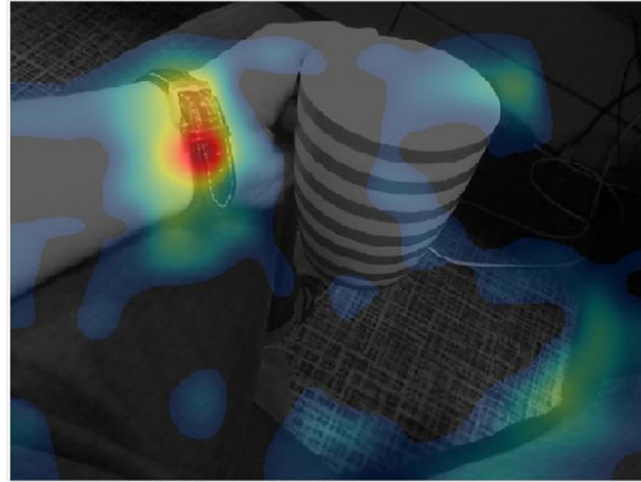


# Deep Learning interpretability methods

## Grad-CAM & Occlusion sensitivity

Post hoc

Model-Specific



Truth:	Coffee mug
AI:	Buckle (15%) ❌

AI classifies incorrectly as "buckle" due to the watch

```
scoreMap = gradCAM(net,X,label)
```



```
scoreMap = occlusionSensitivity(net,X,label)
```

# Deep Learning interpretability methods

## Grad-CAM & Occlusion sensitivity

✓ **Initial investigation:** why are my salad pictures misclassified as pizza?



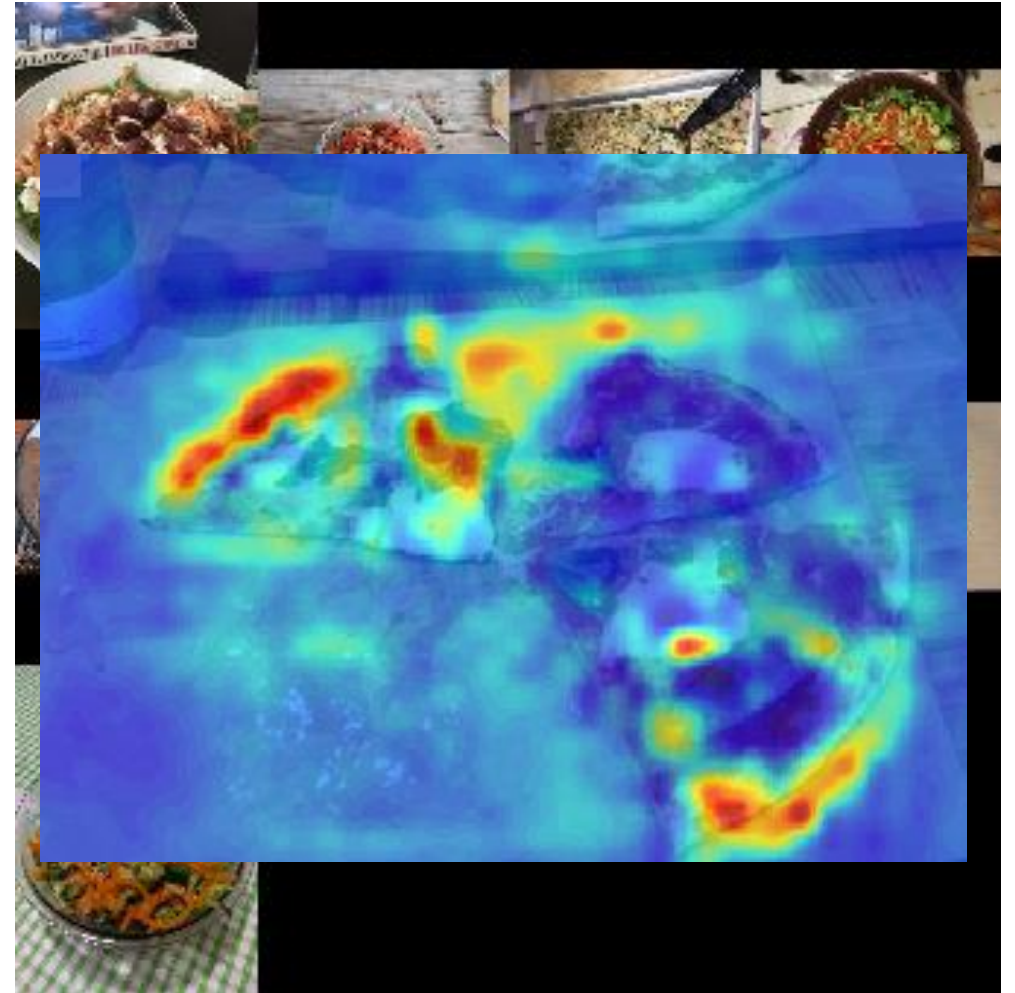
✓ **Hypothesis:** the network is focused on the curving edges of pizzas



✓ **Test:** does this work for other data  
pizza images?



✓ **Fix:** add more pizza slice images and  
salads on plates with curved edges



# Musashi Seimitsu Industry Uses Deep Learning for Visual Inspection of Automotive Parts

## Challenge

Reduce the workload and cost for manually operated visual inspection of 1.3 million automotive parts per month, by implementing an anomaly detection system using deep learning.

## Solution

Musashi Seimitsu built a camera connection setup, preprocess images, create a custom annotation tool, and improve the model accuracy. They generated code for the trained model using GPU Coder™, implemented it on NVIDIA® Jetson.

## Benefits of using MATLAB and Simulink

- Enable a seamless development workflow from image capture to implementation on embedded GPU
- Estimate and visualize the defect area using **Class Activation Mapping**
- Create custom user interfaces with App Designer for improving labeling efficiency
- Leverage consulting services to maximize the benefits of using MATLAB

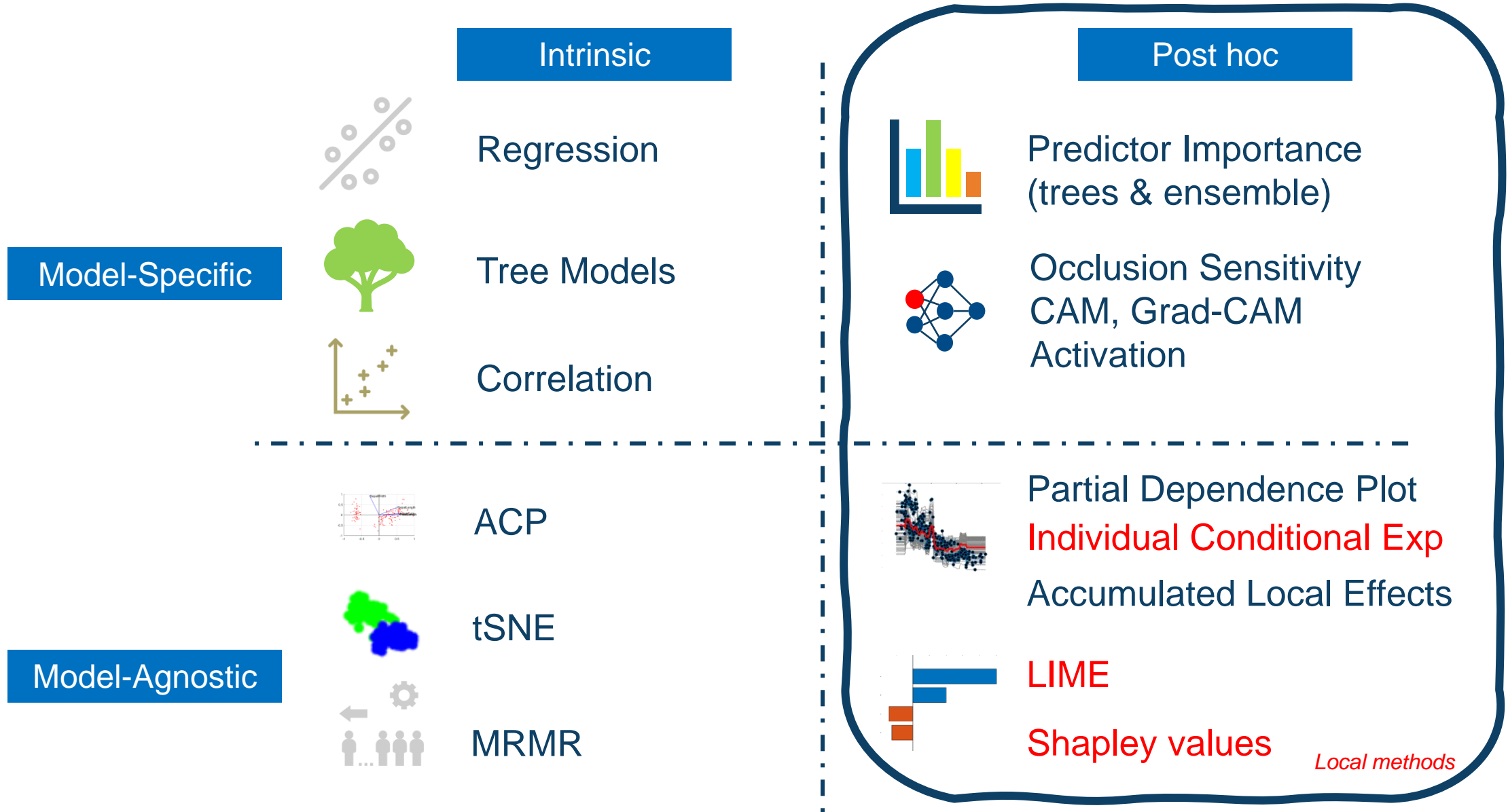


*Using camera connection, preprocessing, and various pretrained models in MATLAB enabled us to work on the entire workflow. Through discussions with consultants, our team gained many tips for solving problems, growing the skills of our engineers.*





# Different Interpretability methods

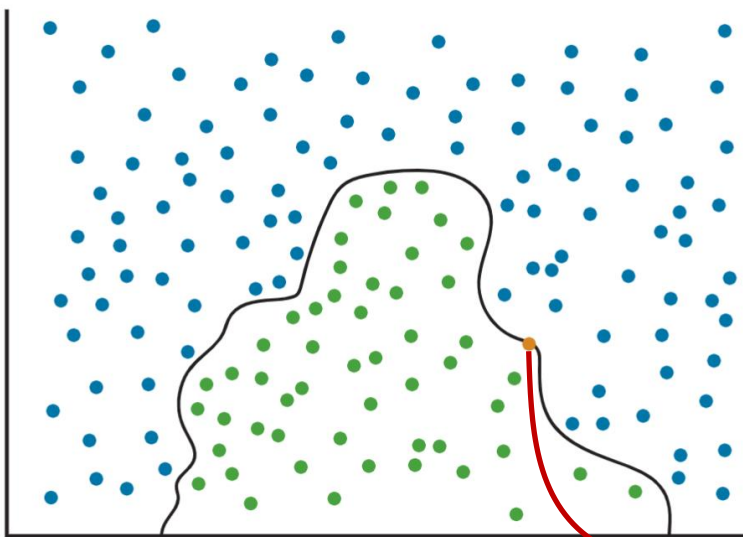


# LIME(Local Interpretable Model-agnostic Explanation)

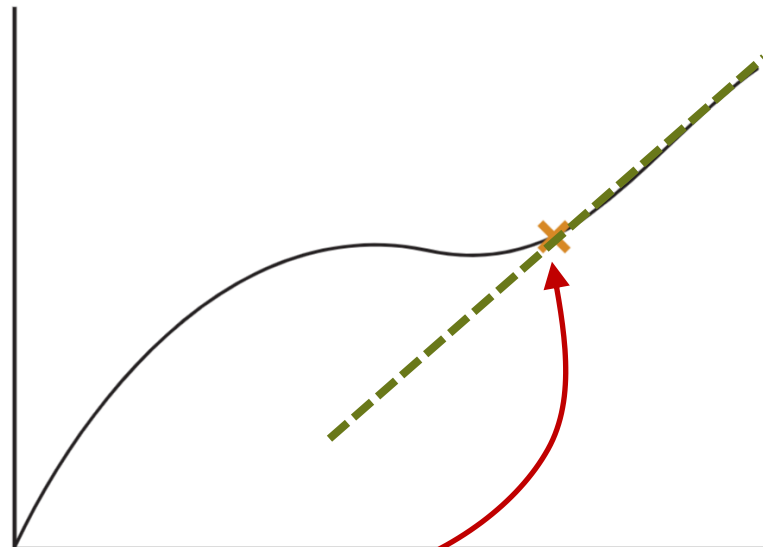
Post hoc

Model-Agnostic

- Fit a simple interpretable model for 1 query point

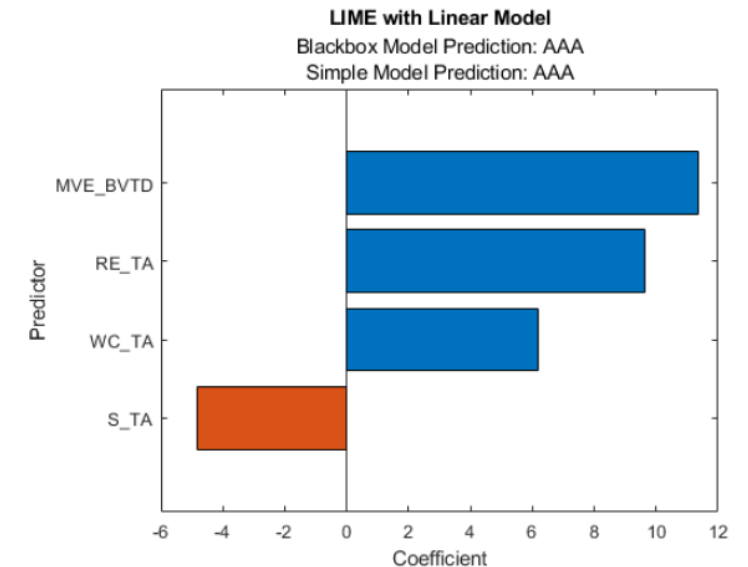


- ✓ Complex model and Point of Interest



- ✓ Approximate with Simple Model

```
results = lime(blackbox)
```



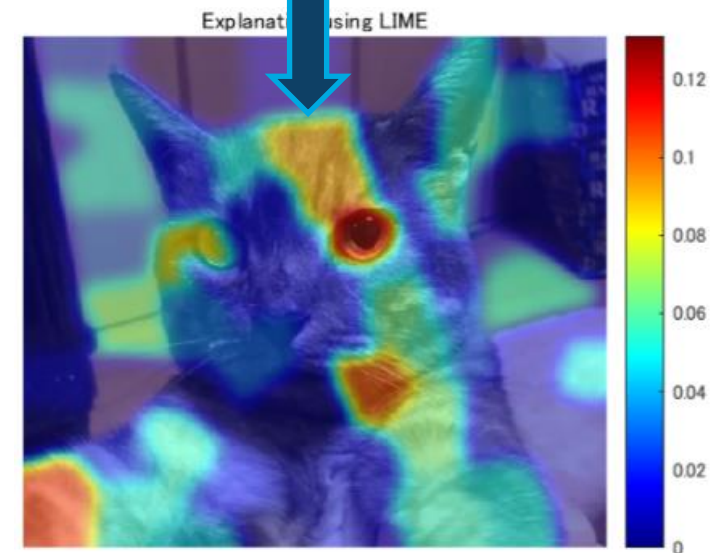
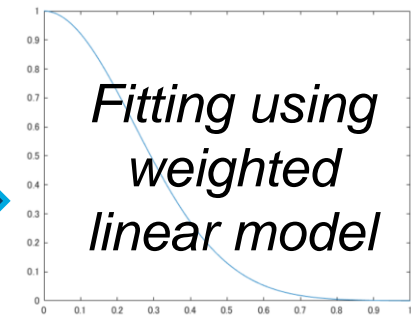
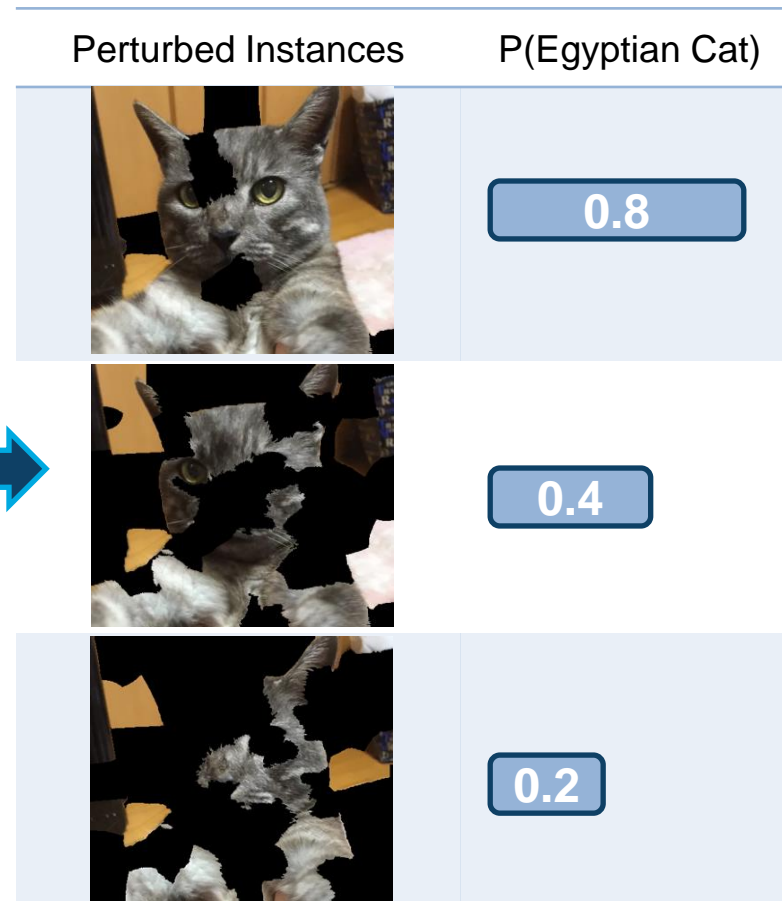
- ✓ 'Explain' = drivers within area of interest

# LIME(Local Interpretable Model-agnostic Explanation)

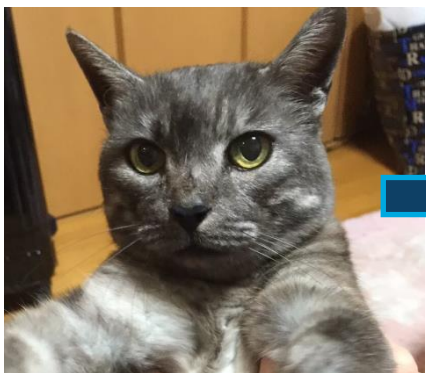
Post hoc

Model-Agnostic

- Fit a simple interpretable model for perturbed instances



```
scoreMap = imageLIME(net,X,label)
```



Original Image




Interpretable Components

# By whitening your AI model, you can expect

**Explainable AI**

**MATLAB**



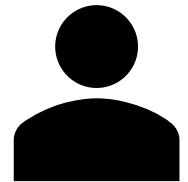
*PDP, ICE  
LIME, Shap  
Occlusion Sensitivity, CAM,  
Grad-CAM, Predictor  
importance, etc.*

Use MATLAB functions to **explain** your model



Data Scientist

I can understand my models & debug it easily



Manager

I trust & understand the data scientist's models



Regulator

I can validate model fairness and trustworthy

# Where can I use Knowledge learned from XAI?

- Decision Critical Application

*Fraud detection*



*Self driving*



*Heath care*



- AI Software = **Code** + **Data**  
(Model/Algorithm)

Scope Project

Data  
Preparation

AI Modeling

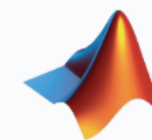
Simulation &  
Test

Deployment

*Iterative Improvement*

***Big Data to Good Data***

# Start whitening your AI model with Golden References Today



MATLAB Deep Learning

[mathworks.github.io](https://mathworks.github.io)

<https://www.mathworks.com/solutions/deep-le...>

**Grad-CAM**

**LIME**

**Test Image**

**Grad-CAM: Road**

**LSTM Activations**

**Max pooling activations**

**Final conv activations**

**Softmax activations**

**Golden Reference**

**Information**

**ESU**

**VL**

**Image LIME (golden retriever - top 4 features)**

**Understand Network Predictions Using Occlusion**

Use occlusion sensitivity to understand why a deep neural network makes a classification decision. Occlusion sensitivity...

**Visualize In-Classification Maximal Activation**

Use a data set to visualize which parts of an image activate the convolutional layers of a neural network to understand how it makes a classification decision.

**Investigate Saliency Classification**

Use locally interpretable model-agnostic explanations (LIME) to investigate the role of individual convolutional layers in a deep neural network.

**Deep Dream Image with GoogLeNet**

Generate images using the deepDreamImage function with a pretrained convolutional neural network GoogLeNet.

**Understand Network Predictions Using LIME**

Use locally interpretable model-agnostic explanations (LIME) to understand why a deep neural network makes a classification decision.

[Open Live Script](#)

Understanding Network Predictions for Image Classification (UNPIC)

Image Data | Accuracy | Predict | Prediction Explorer | **Features** | t-SNE

Visualize which parts of an image are most important for classification. Visualize for a chosen image file, or a random image from the test data.

Choose image file:

Random image from class: french fries

True class: french fries

Grad-CAM Settings

Target class: french fries

Feature map: inception\_5b-output

# MATLAB EXPO 2021

감사합니다

