# Using Explanations to Identify Problems and Limitations in AI Models used for Intelligent Maintenance

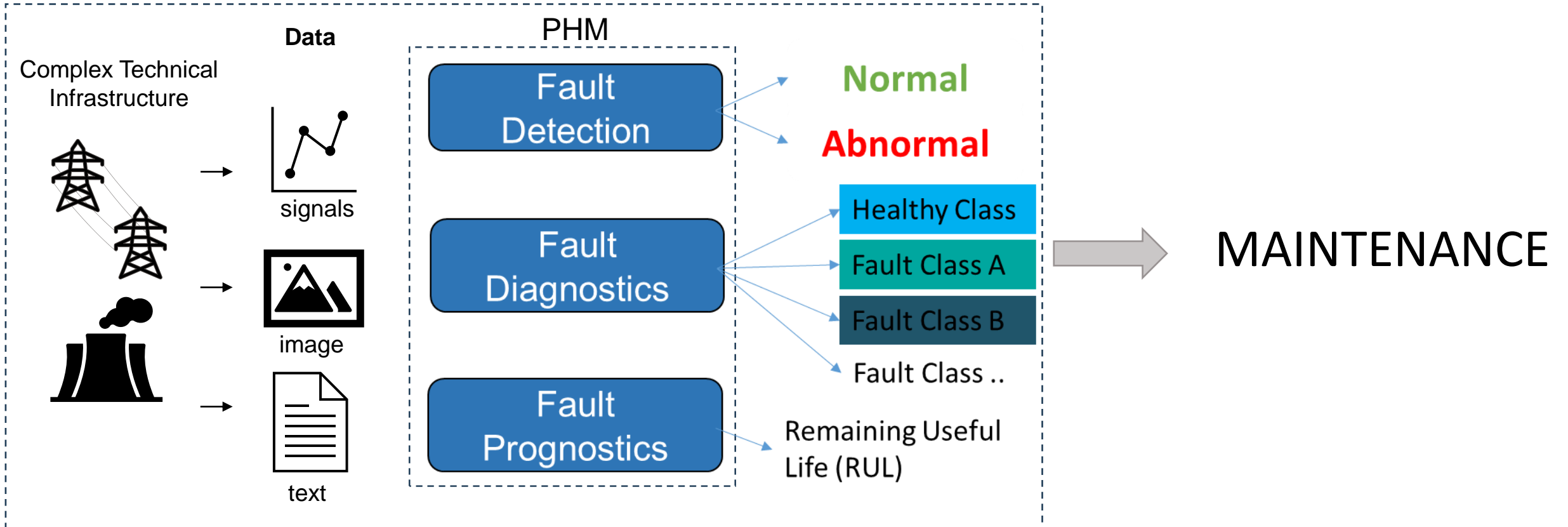Giovanni FLOREALE[1], Piero BARALDI[1], Enrico ZIO[2,1], Olga FINK[3],

*1 Department of Energy, Politecnico di Milano, Milan, Italy*

*2 MINES Paris-PSL University, Centre de Recherche sur les Risques et les Crises (CRC), Sophia Antipolis, France*
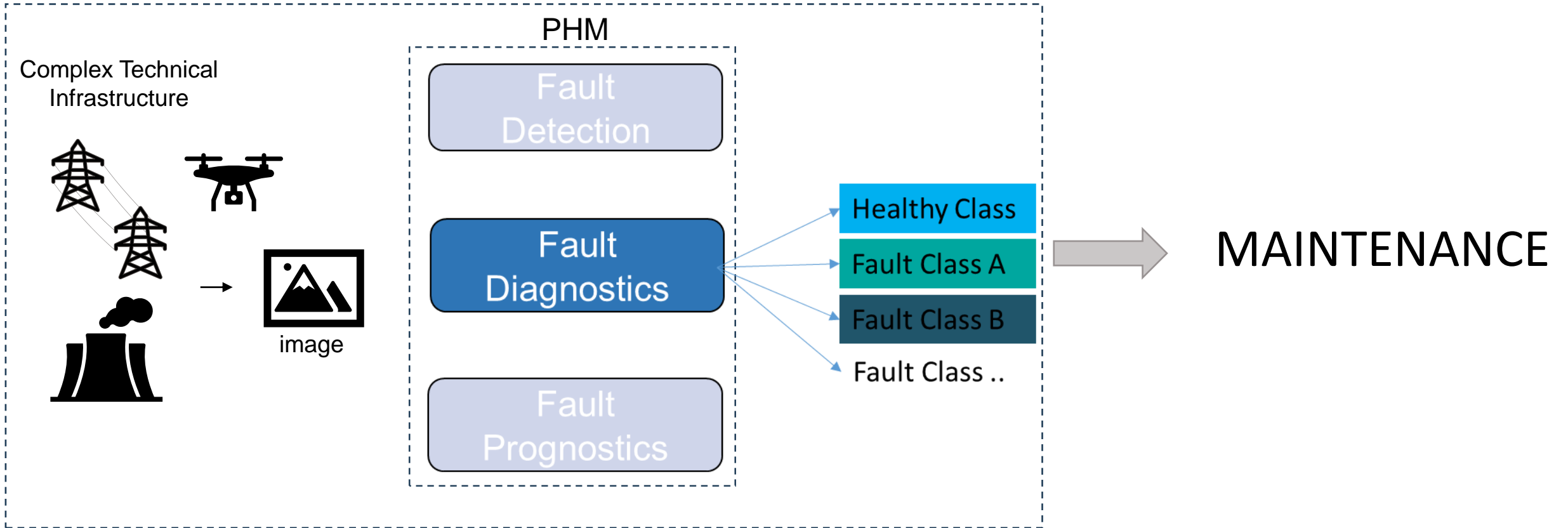
*3 IMOS, EPFL, Lausanne, Switzerland*
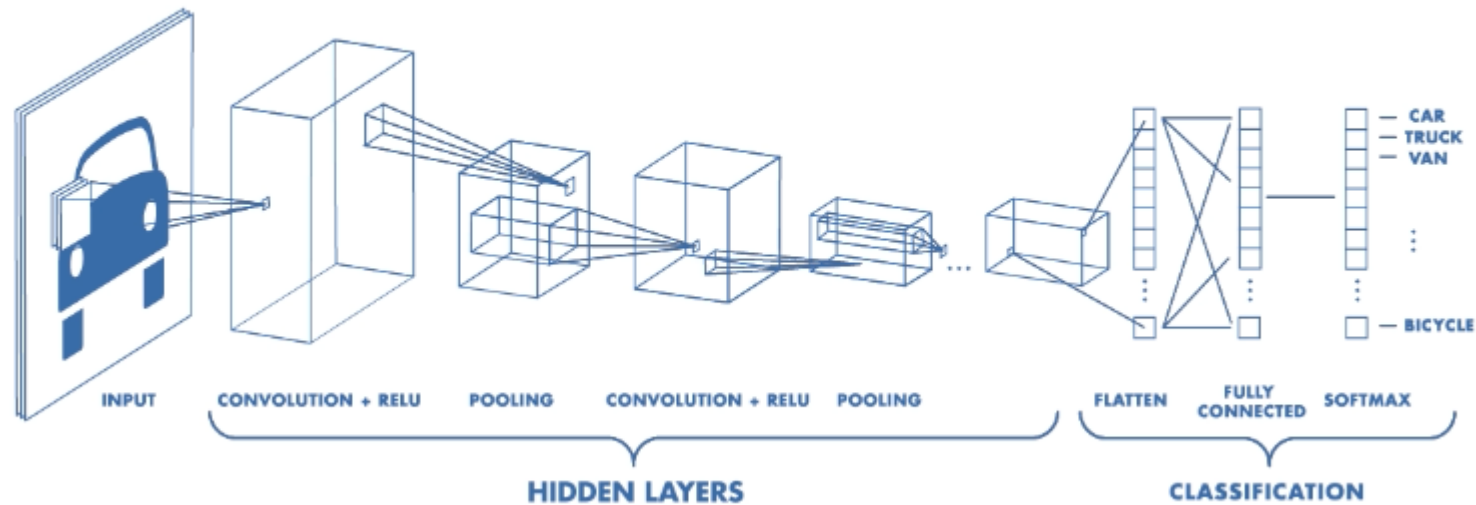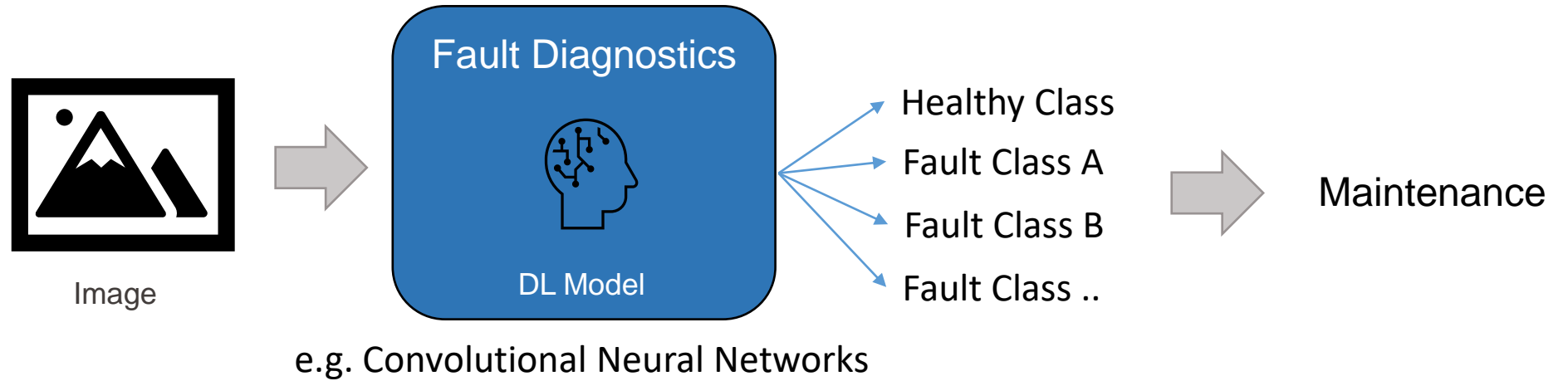
Lausanne, September 12th, 2023

# Context: What is the Problem?

# Context: What is the Problem?

EPFL  École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# How is Fault Diagnostics from Images done?



Image

Fault Diagnostics

DL Model

e.g. Convolutional Neural Networks

Healthy Class

Fault Class A

Fault Class B

Fault Class ..

Maintenance



INPUT  CONVOLUTION + RELU  POOLING  CONVOLUTION + RELU  POOLING  FLATTEN  FULLY CONNECTED  SOFTMAX

CAR
TRUCK
VAN

BICYCLE

HIDDEN LAYERS

CLASSIFICATION

Giovanni Floreale

EPFL  École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# What are the Technical and Scientific Challenges?



Image

Fault Diagnostics

DL Model

**Challenges:**
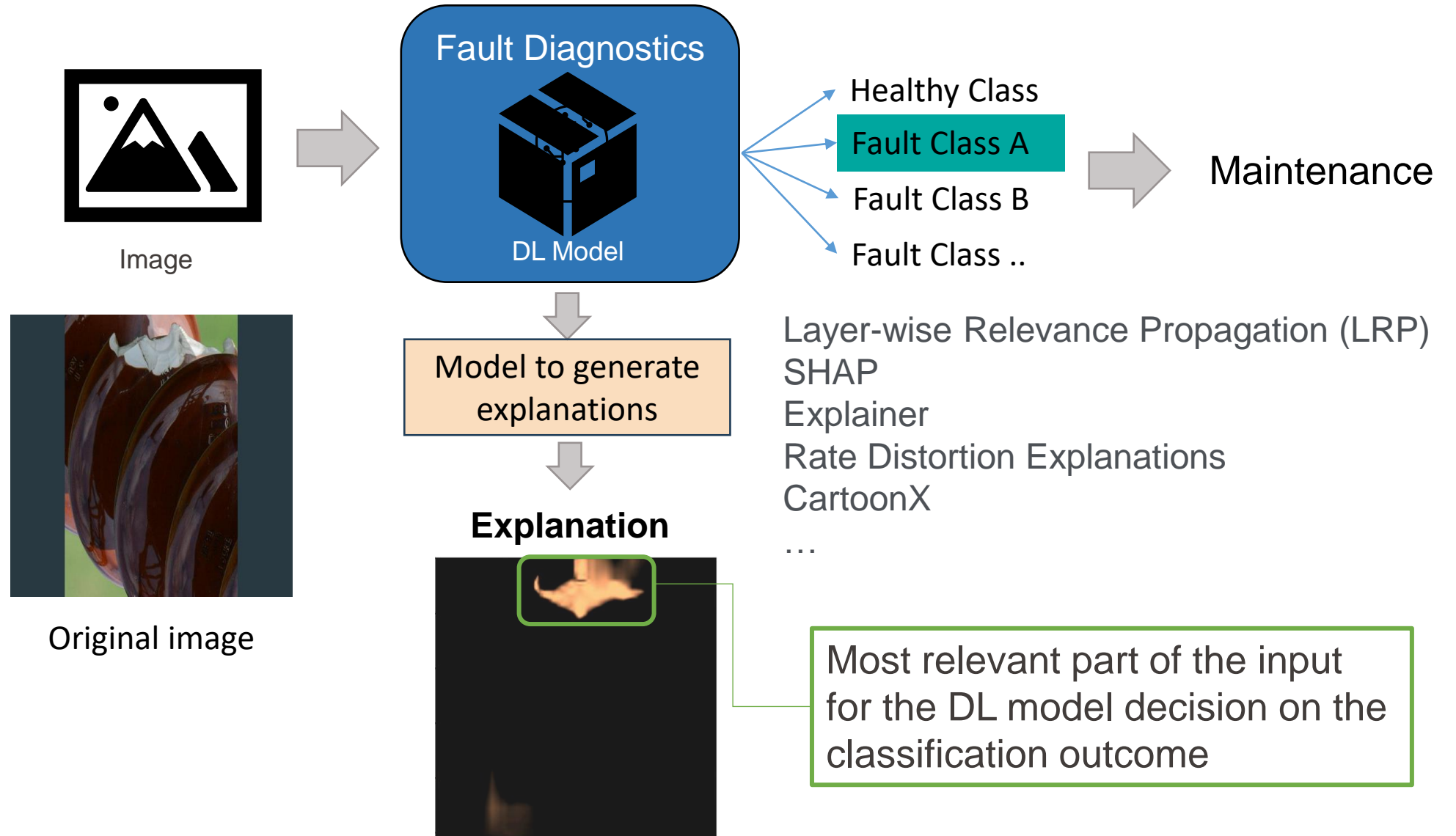- Black Box
- Performance

Classification Accuracy

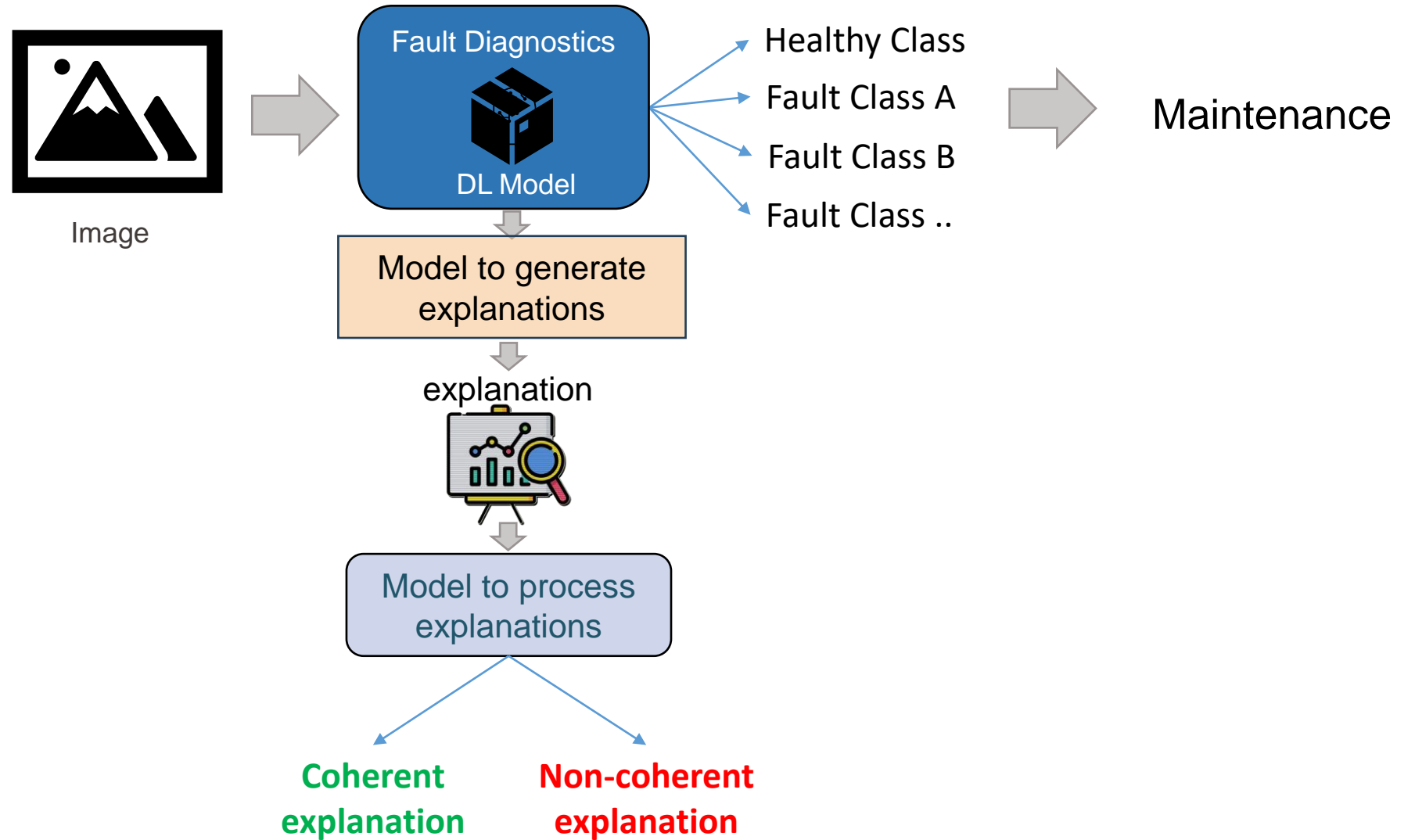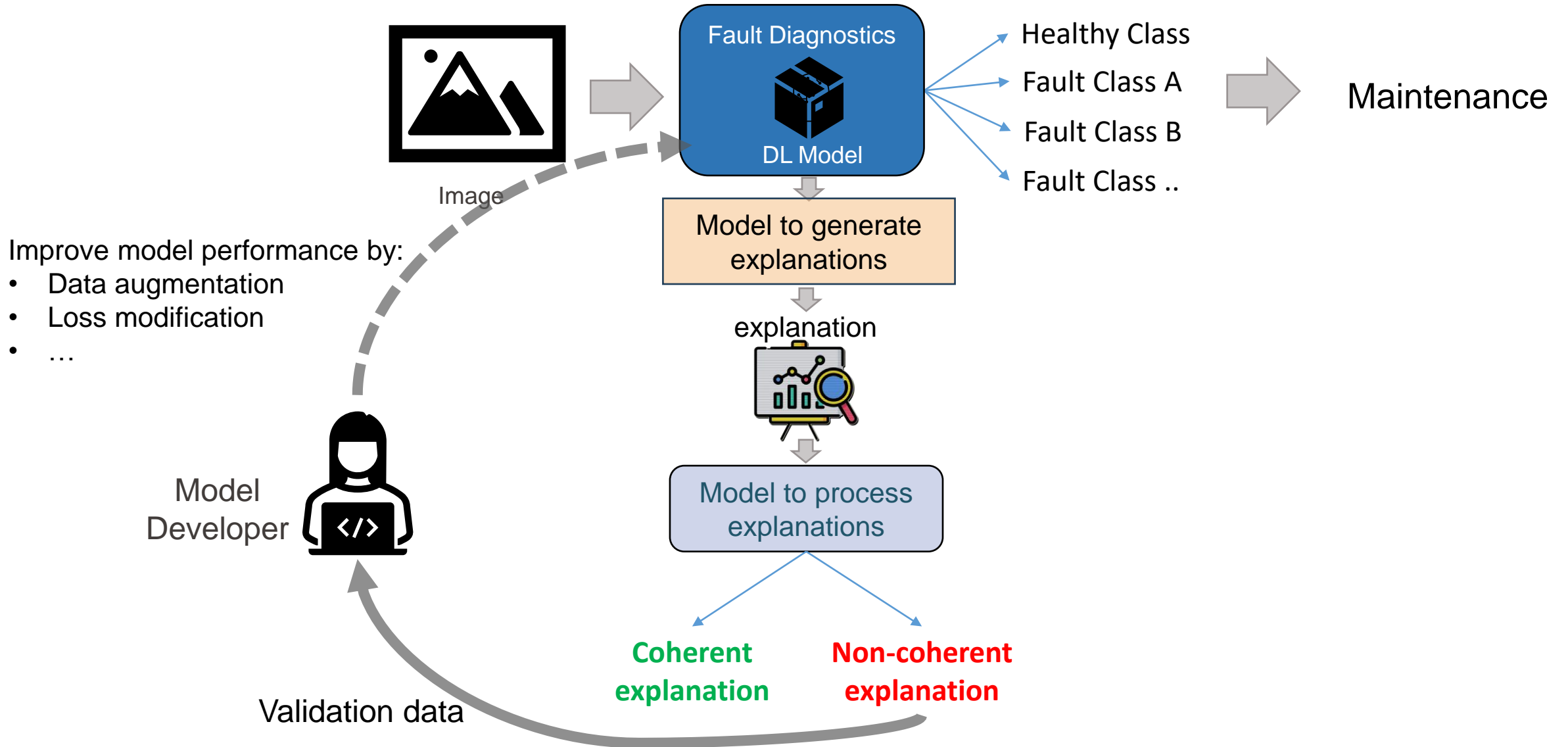Training set

Test set

In field data

0

Giovanni Floreale

EPFL École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# Opportunities: What Can Be Done?



Image

Original image

Fault Diagnostics

DL Model

Healthy Class
Fault Class A
Fault Class B
Fault Class ..

Maintenance

Model to generate explanations

Layer-wise Relevance Propagation (LRP)
SHAP
Explainer
Rate Distortion Explanations
CartoonX
…

**Explanation**

Most relevant part of the input for the DL model decision on the classification outcome

Giovanni Floreale

EPFL  École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# Objective: What are we Trying to Do?



Image

Fault Diagnostics
DL Model

Healthy Class
Fault Class A
Fault Class B
Fault Class ..

Maintenance

Model to generate explanations

explanation

Model to process explanations

**Coherent explanation**

**Non-coherent explanation**

# Relevance: Why is it Useful? (1)



Image

Fault Diagnostics
DL Model

Healthy Class
Fault Class A
Fault Class B
Fault Class ..

Maintenance

Model to generate explanations

explanation

Model to process explanations

**Coherent explanation**

**Non-coherent explanation**

Improve model performance by:
• Data augmentation
• Loss modification
• …

Model Developer

Validation data

# Relevance: Why is it Useful? (2)



Image

Fault Diagnostics
DL Model

Healthy Class
Fault Class A
Fault Class B
Fault Class ..

Maintenance

Improve model performance by:
- Data augmentation
- Loss modification
- ...

Model to generate explanations

explanation

Model to process explanations

Coherent explanation

Non-coherent explanation

Model Developer

Validation data

Informed decision making

Maintenance Operator

In-field data

Giovanni Floreale

POLITECNICO MILANO 1863

# Developed Method

Giovanni Floreale

EPFL École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# Developed Method



Image → Fault Diagnostics DL Model $\phi$ → Healthy Class / Fault Class A / Fault Class B / Fault Class .. → Intelligent Maintenance

Model to generate explanations → explanation $x_i$

Embedding space definition: loss function

🟢 Correct classifications:
minimize the distance from the centre

coherent

Embedding Space

$$\min_{\mathcal{W}} \quad \frac{1}{n+m} \sum_{i=1}^{n} \|\phi(x_i; \mathcal{W}) - c\|^2$$

Giovanni Floreale

# Developed Method



Embedding space definition: loss function

- 🔴 Classifications errors:
  maximize the distance from the centre
  (hyperparameter $\eta$)

$$\min_{\mathcal{W}} \quad \frac{1}{n+m}\sum_{i=1}^{n}\|\phi(x_i;\mathcal{W})-c\|^2 + \frac{\eta}{n+m}\sum_{j=1}^{m}\left(\|\phi(\tilde{x}_j;\mathcal{W})-c\|^2\right)^{\tilde{y}_j}.$$
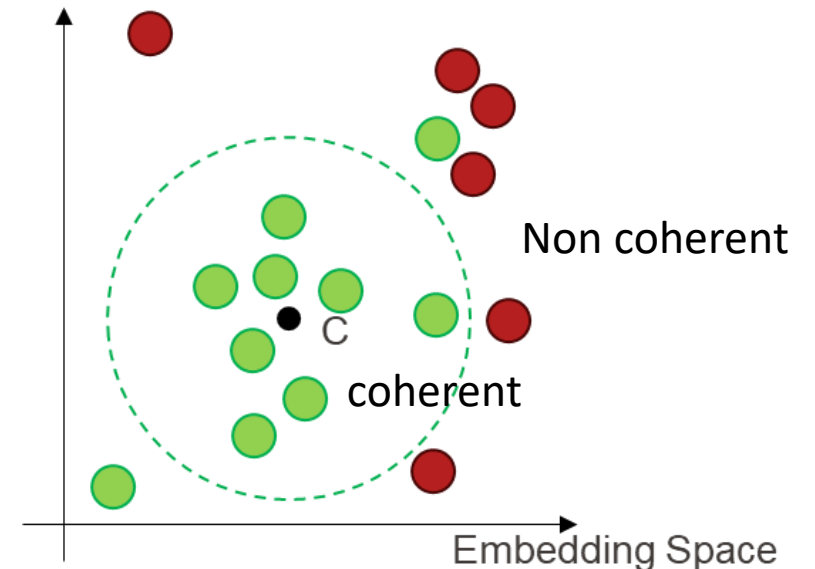
# Developed Method



Embedding space definition: loss function

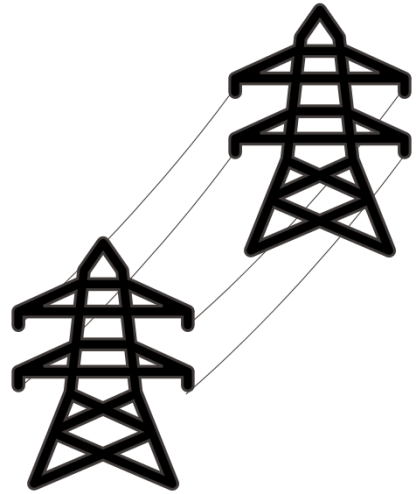Regularization term to avoid overfitting (hyperparameter $\lambda$)

$$\min_{\mathcal{W}} \quad \frac{1}{n+m}\sum_{i=1}^{n}\|\phi(x_i;\mathcal{W})-c\|^2 + \frac{\eta}{n+m}\sum_{j=1}^{m}\left(\|\phi(\tilde{x}_j;\mathcal{W})-c\|^2\right)^{\tilde{y}_j} + \frac{\lambda}{2}\sum_{\ell=1}^{L}\|W^\ell\|_F^2.$$
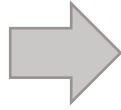
# Case Study

Critical Components:



Power Grid                    Insulators                    Drones                    Shells' images
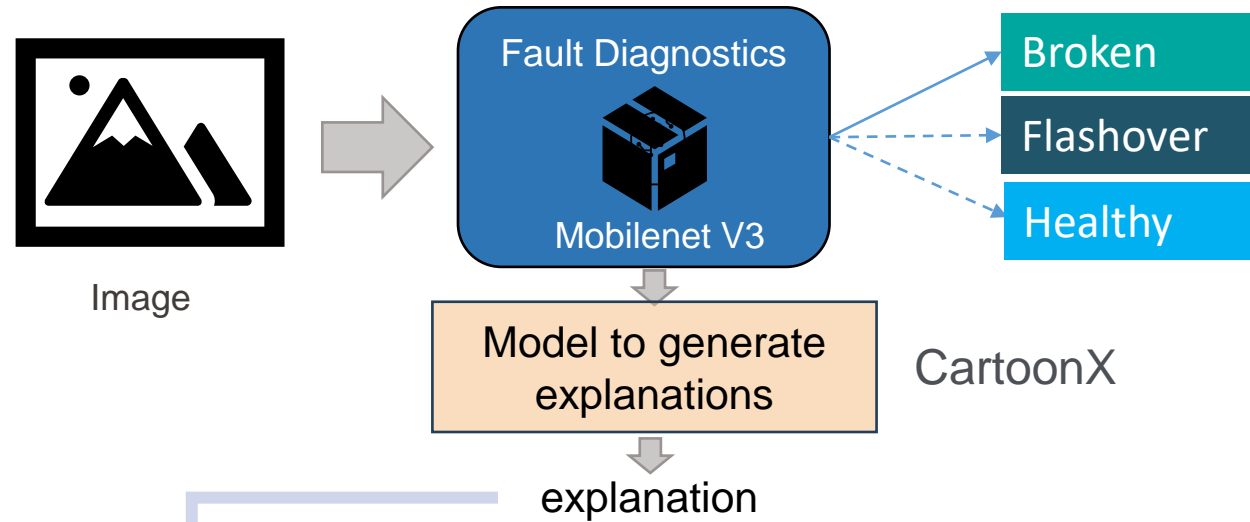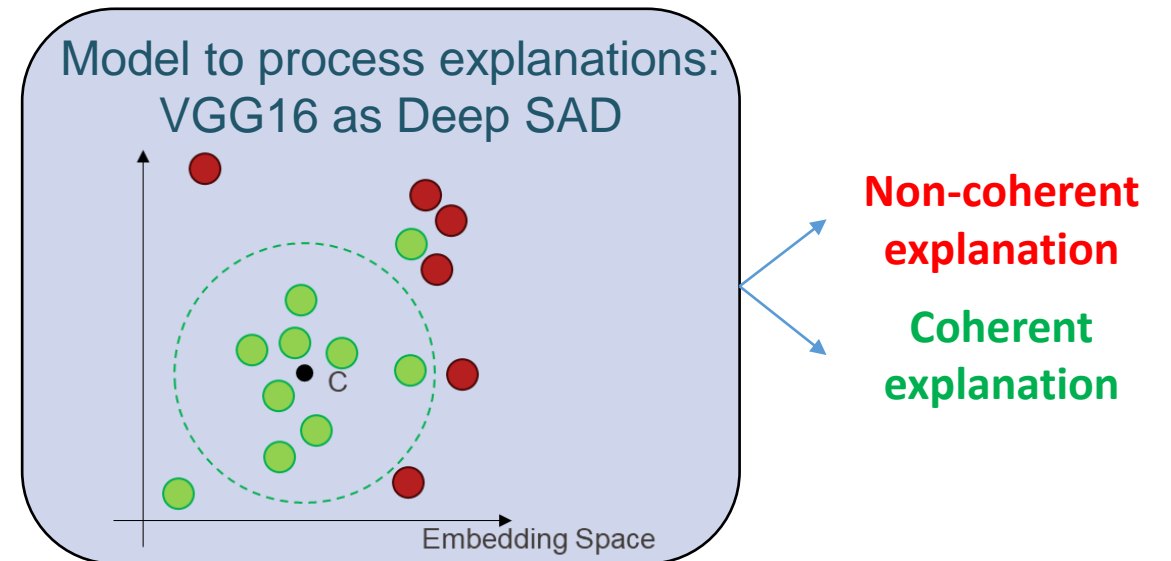
e.g. swiss power grid:
6700 km long → 12 000 pylons

# Case Study

# Relevance: Why is it Useful? (1)



Giovanni Floreale    EPFL École polytechnique fédérale de Lausanne    POLITECNICO MILANO 1863

# Results: Proposed Approach on Validation Set



Model developer is guided for model improvement

e.g. data augmentation

short-cut!

Non coherent

Image

Fault Diagnostics

DL Model

Broken    Correct!
Flashover
Healthy

Model to generate explanations

explanation

Deep SAD to process explanations

Embedding Space

Validation set

# Results: Proposed Approach on Validation Set

# Relevance: Why is it Useful? (2)



Giovanni Floreale

EPFL École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# Results: DL model



In-field image

Fault Diagnostics

DL Model

Broken

Flashover

Healthy

Maintenance

Model to generate explanations

explanation

| Accuracy of DL model | 89.20% |
|---|---|

# Results: Proposed Approach In-field Application



| | |
|---|---|
| Accuracy of DL model | 89.20% |
| Fraction of explanations that is revised by experts | 19.52% |
| Accuracy of the proposed approach | 94.87% |

+5.67%

In-field image

Fault Diagnostics
DL Model

Broken
Flashover
Healthy

Maintenance

Model to generate explanations

explanation

Deep SAD to process explanations

Non coherent

Embedding Space

Broken
Flashover
Healthy

Operator is asked to correct the DL model classification

# Conclusion



**Challenges:**
- Black Box
- In-field performance drop

Giovanni Floreale

EPFL  École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# Conclusion



Image

Fault Diagnostics

DL Model

Broken

Flashover

Healthy

Maintenance

CartoonX to generate explanations

Deep SAD to process explanations

Coherent explanation

Non-coherent explanation

Model Developer

✓ short-cut identification

Validation data

EPFL École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863

# Conclusion



Giovanni Floreale — EPFL École polytechnique fédérale de Lausanne — POLITECNICO MILANO 1863

# Conclusion

**Future Works:**

Explanation-based performance improvement



Image

Fault Diagnostics

DL Model

Broken

Flashover

Healthy

Maintenance

CartoonX to generate explanations

Deep SAD to process explanations

**Coherent explanation**     **Non-coherent explanation**

Model Developer

short-cut identification

Validation data

Reduced expert effort

Improved performance

| Accuracy of DL model | 89.20% | +5.67% |
| Fraction of explanations that is revised by experts | 19.52% | |
| Accuracy of the proposed approach | 94.87% | |

Operator asked to reclassify non-coherent explanations

Broken

Flashover

Healthy

In-field data

Giovanni Floreale

EPFL  École polytechnique fédérale de Lausanne

POLITECNICO MILANO 1863