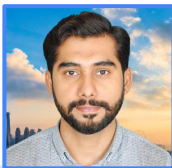# Can I trust my anomaly detection system? A case study based on eXplainable AI.

UNIVERSITÀ DI TORINO

**Elvio G. Amparore**
elviogilberto.amparore@unito.it

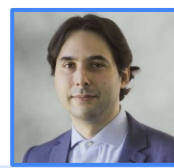**Muhammad Rashid**
muhammad.rashid@unito.it

RuleX

**Enrico Ferrari**
enrico.ferrari@rulex.ai

**Damiano Verda**
damiano.verda@rulex.ai

# Topics of the presentation

- Case study of Anomaly Detection in Industrial Quality Control System

- Use of Generative AI for Anomaly Detection

- Transparency & Trustworthiness in Anomaly Detection Systems

  ○ Review an explainable AD system architecture* that combines VAE-GAN models with the LIME and SHAP explanation methods.

  ○ Quantify the AD system efficacy using anomaly scores

  ○ Use XAI methods to determine if anomalies are indeed detected for the right reason, improving the framework of Ravi et al*.

\* Ravi, A., Yu, X., Santelices, I., Karray, F., & Fidan, B. (2021, August). **General frameworks for anomaly detection explainability: comparative study**. In *2021 IEEE International Conference on Autonomous Systems (ICAS)* (pp. 1-5). IEEE.
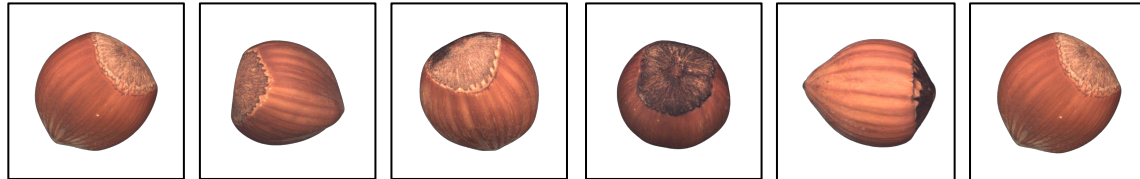
# Defining Defective & Non-defective

Consider an industrial quality control system use case.

- Non defective products are common and easy to capture and describe
- Defective products are rare and unpredictable

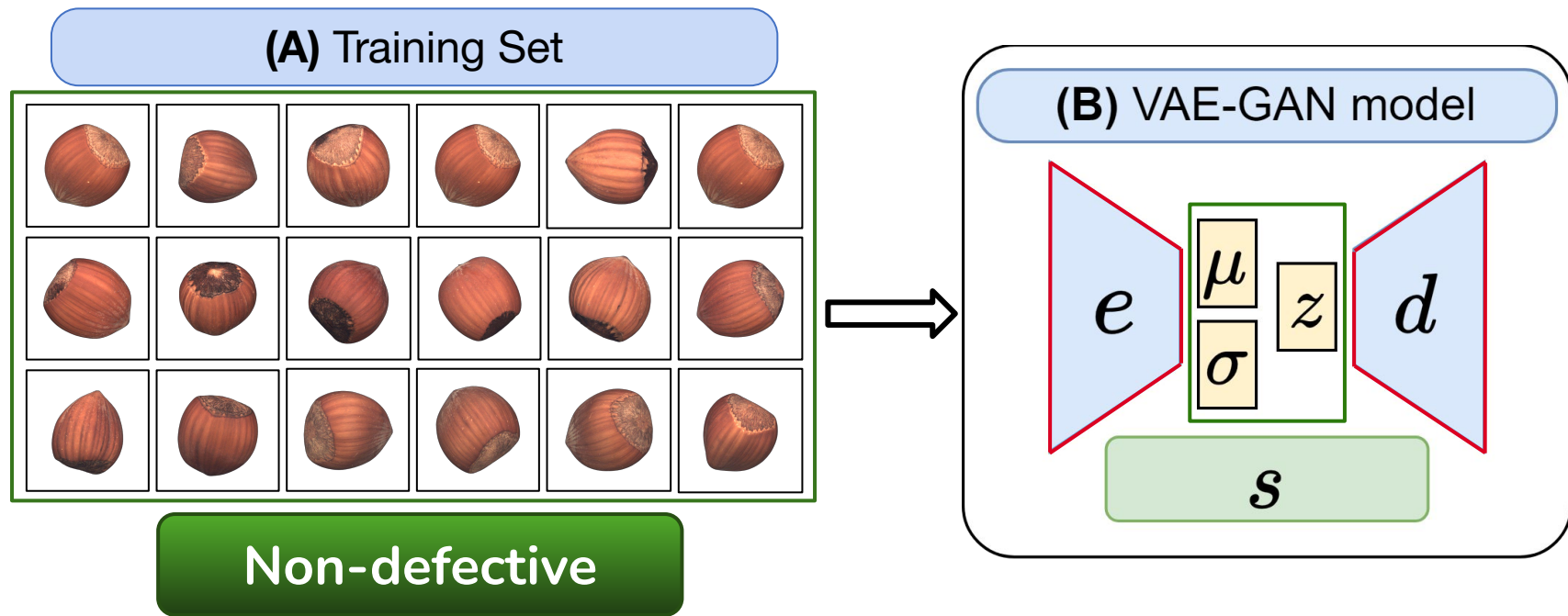→ setup for **anomaly detection**.

**Non-defective:**



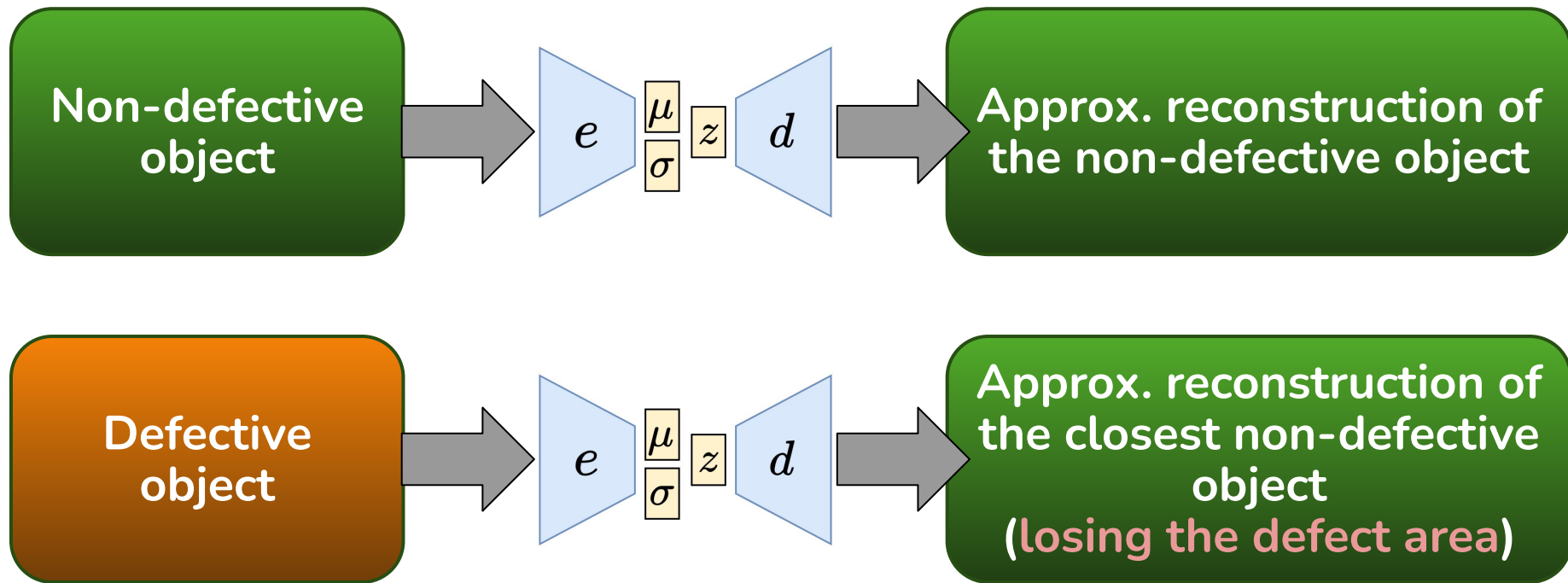**Defective:**



Bergmann, P., Batzner, K., Fauser, M., Sattlegger, D., & Steger, C. (2021). The **MVTec** anomaly detection dataset: **a comprehensive real-world dataset for unsupervised anomaly detection**. *International Journal of Computer Vision*.

## Training the **V**ariational **A**uto-**E**ncoder



(A) Training Set

**Non-defective**

(B) VAE-GAN model

$e$ $\mu$ $\sigma$ $z$ $d$

$s$

# Anomaly detection with VAE-GAN

Anomaly Detection with **V**ariational **A**uto-**E**ncoder

| Non-defective object | → | $e$ $\mu$ $\sigma$ $z$ $d$ | → | Approx. reconstruction of the non-defective object |

| Defective object | → | $e$ $\mu$ $\sigma$ $z$ $d$ | → | Approx. reconstruction of the closest non-defective object (**losing the defect area**) |

# Anomaly detection with generative AI

RuleX

## Pipeline with a **V**ariational **A**uto-**E**ncoder

**Non-defective object ξ**

**Reconstruction to the non-defective object ξ'**

**(B)** VAE-GAN model

$e$ $\mu$ $\sigma$ $z$ $d$

$s$

(D) Anomaly Map $m$

$-$

$$m = \left| gs(\xi) - gs(\xi') \right|$$
$$\alpha = \max(m)$$

max

**(E)** anomaly score $\alpha$

0.111

◆ Anomaly Map is sum of Real Anomaly & background noise (if any)

# Anomaly detection with generative AI

## Pipeline with a **V**ariational **A**uto-**E**ncoder



**Defective object ξ**

**(B)** VAE-GAN model

$e$ $\mu$ $\sigma$ $z$ $d$

$s$

**Reconstruction losing the defect area ξ'**

(D) Anomaly Map $m$

$-$

max

**(E)** anomaly score $\alpha$

0.442

$$m = \left| gs(\xi) - gs(\xi') \right|$$
$$\alpha = \max(m)$$

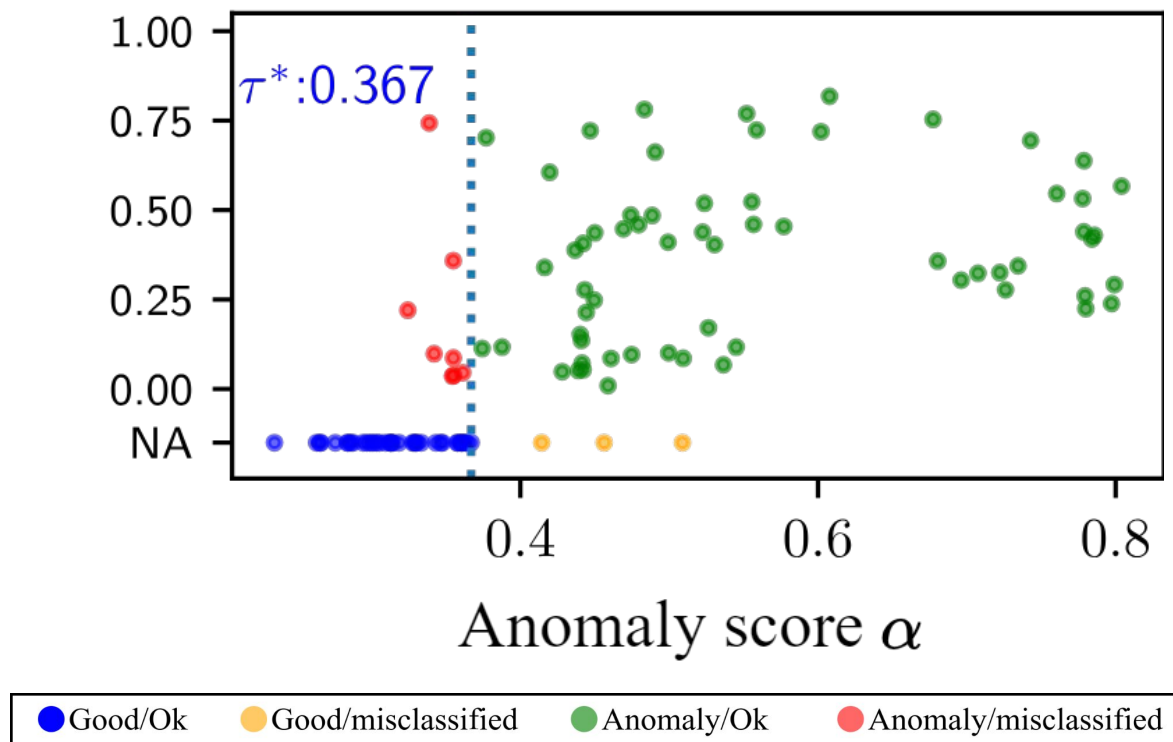◆ Anomaly Map is sum of Real Anomaly & background noise (if any)

# Anomaly detection threshold

Anomalous: if $\alpha \geq \tau*$

Finding the optimal threshold $\tau$* means solving an optimization problem.

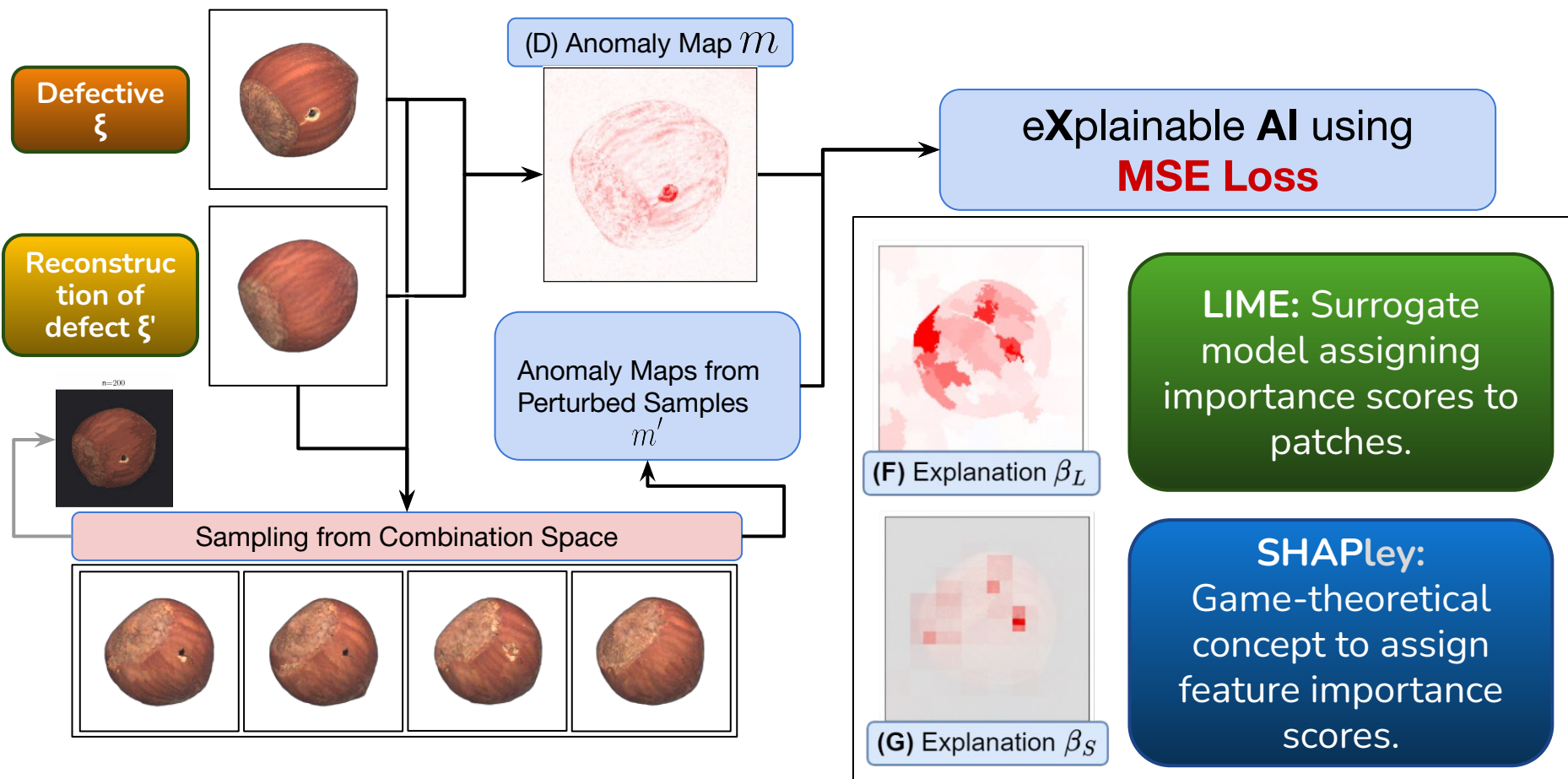$$\tau^* = \underset{\tau}{\operatorname{argmax}} \sqrt{\operatorname{TPR}(\tau) \times (1 - \operatorname{FPR}(\tau))}$$

True Positive Rate   : Anomalous as anomalous
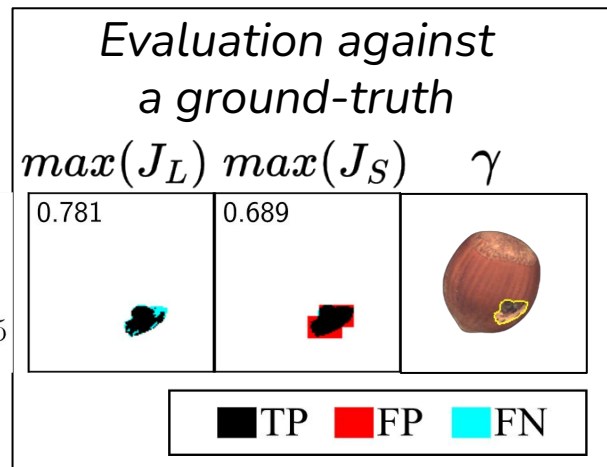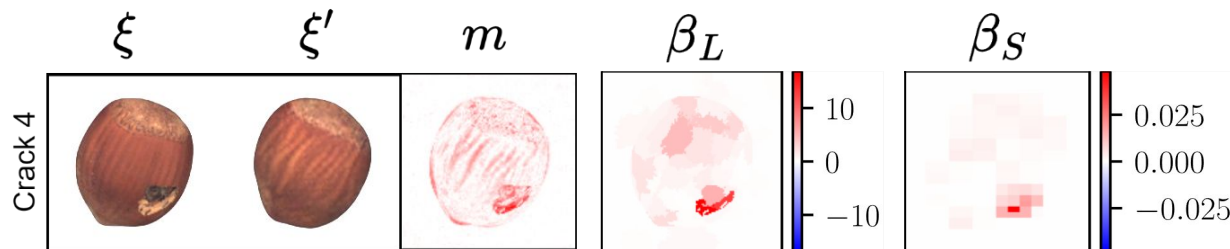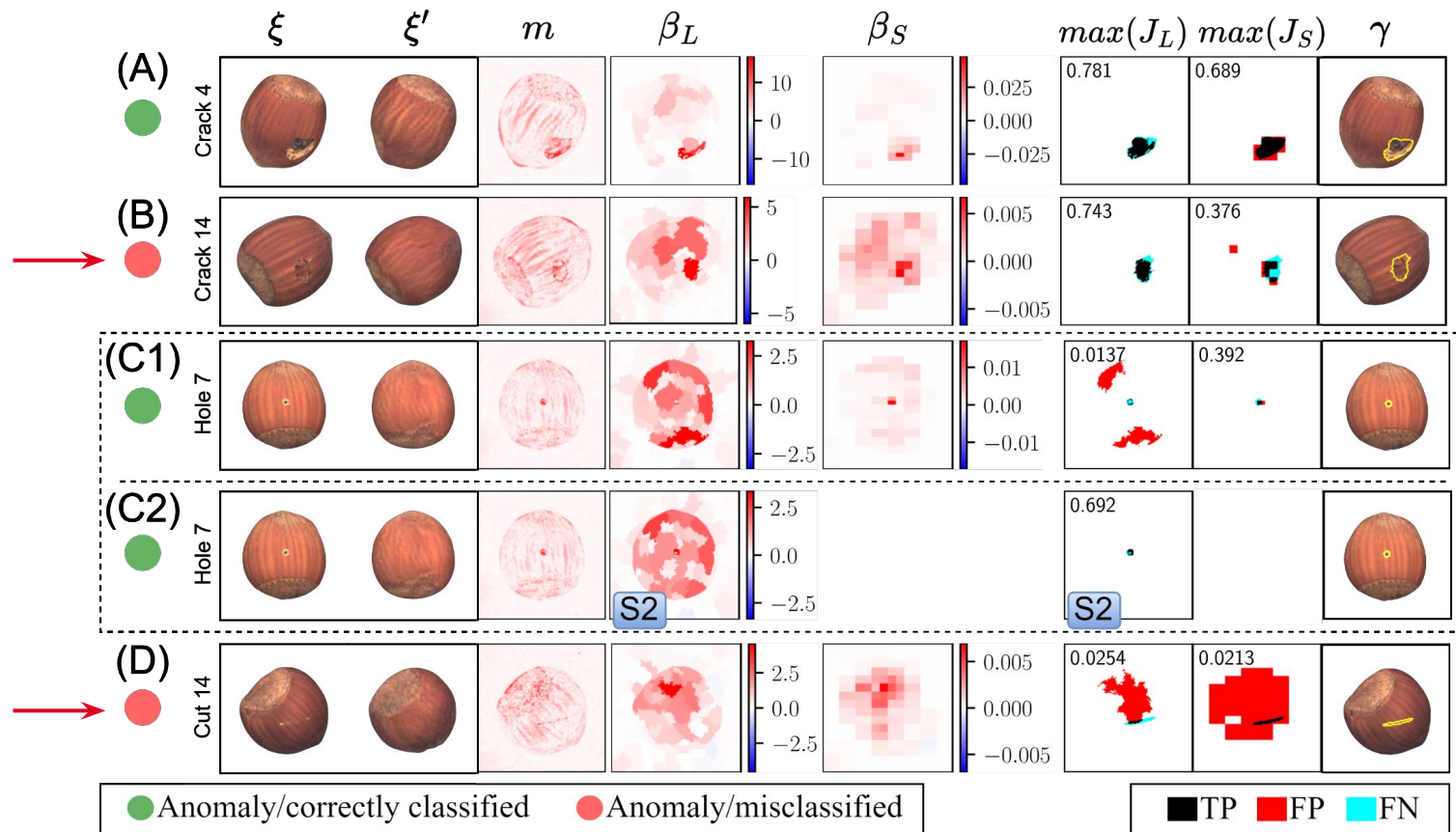False Positive Rate  : Normal as anomalous

# Anomaly detection threshold

$\tau^*$:0.367

Anomaly score $\alpha$

Good/Ok    Good/misclassified    Anomaly/Ok    Anomaly/misclassified

# Contribution of XAI in Anomaly Detection

- We have the anomaly map + a detection threshold.
  Is it enough for explaining anomaly?
  → Of course not.

- Problem:  anomaly map $\alpha$ is
    the sum of **reconstruction error** (noise) + **anomaly**  (if any)
  → Need a way to:
  - separate the **anomaly** from the **noise**;
  - and to localize the region of the **anomaly**.

- More precise information

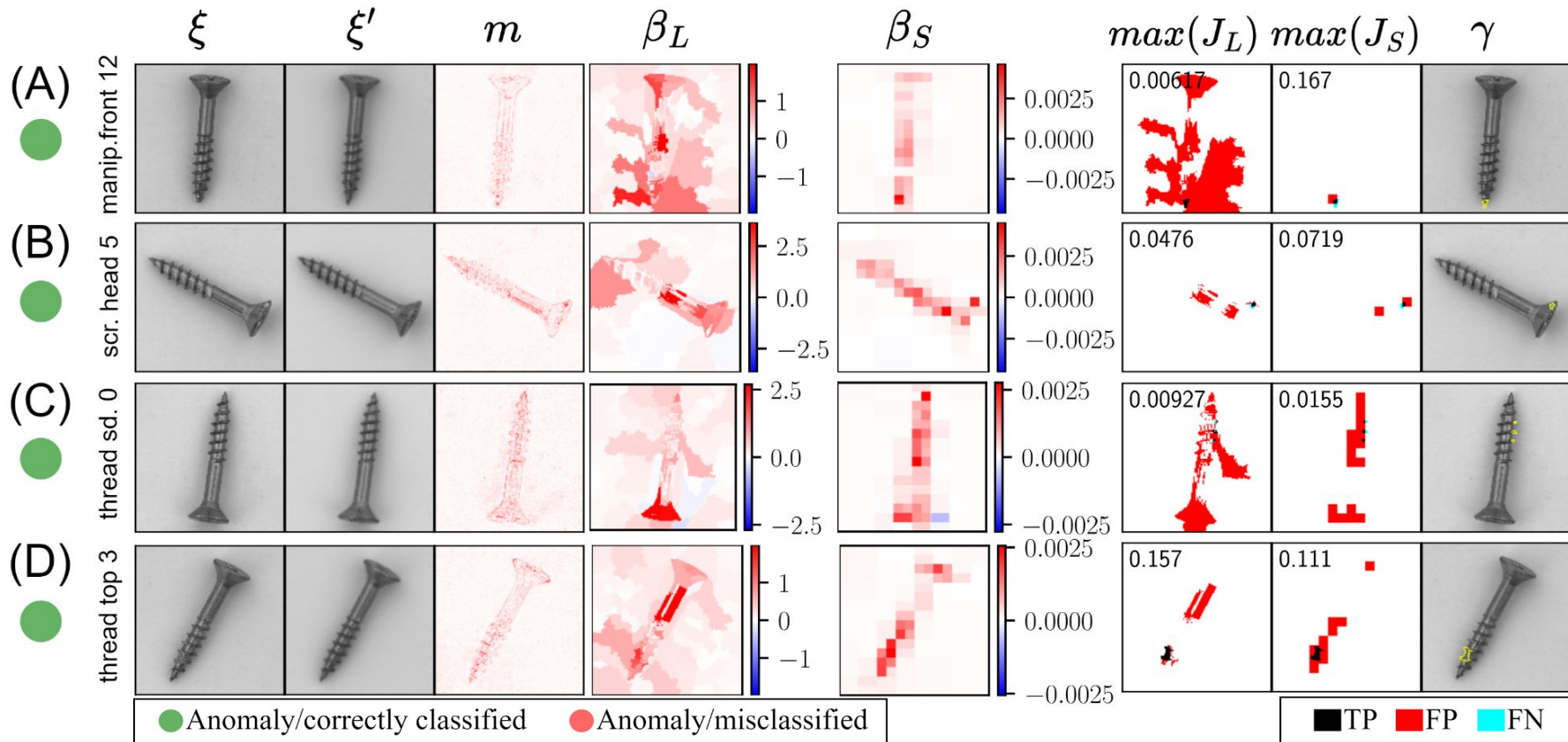- Localization of anomaly

- Is the **anomaly** a real anomaly?

Defective $\xi$

Reconstruction of defect $\xi'$

n=200

Sampling from Combination Space

(D) Anomaly Map $m$

Anomaly Maps from Perturbed Samples $m'$

eXplainable AI using MSE Loss

(F) Explanation $\beta_L$

(G) Explanation $\beta_S$

LIME: Surrogate model assigning importance scores to patches.

SHAPley: Game-theoretical concept to assign feature importance scores.

## Explaining anomaly detection with XAI

# Explaining Anomalies

# Explaining Anomalies

| | $\xi$ | $\xi'$ | $m$ | $\beta_L$ | $\beta_S$ | $max(J_L)$ | $max(J_S)$ | $\gamma$ |

(A) manip.front 12 — 0.00617 | 0.167

(B) scr. head 5 — 0.0476 | 0.0719

(C) thread sd. 0 — 0.00927 | 0.0155

(D) thread top 3 — 0.157 | 0.111

● Anomaly/correctly classified    ● Anomaly/misclassified

■ TP  ■ FP  ■ FN
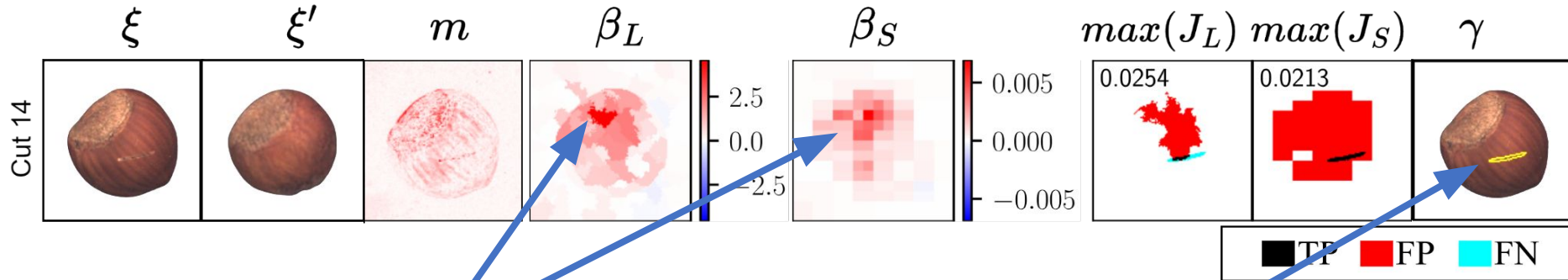
# Revealing model misbehaviours

Samples may be classified as anomalous for the wrong reason, and only XAI can reveal such behaviour.



model detect anomaly from a badly reconstructed region…

…but the anomaly is a small cut region in a different location.

# Conclusions

- XAI methods are relevant in finding the **true drivers** behind AI systems using techniques like classification and/or anomaly detection.

- Case study based on **reconstruction error maps** generated from **VAE-GAN** models.

- **Multiple XAI** techniques to separate the **reconstruction error** (noise) from the **anomaly** (if any).

- A sample may be detected as **anomalous** for the **wrong reasons**, yet this misbehaviour may not be detectable from the information provided by the anomaly detection system alone → Role of XAI!

# Can I trust my anomaly detection system?
# A case study based on e**X**plainable **AI**.

https://github.com/rashidrao-pk/anomaly_detection_trust_case_study

# Thank you! – Questions?

UNIVERSITÀ DI TORINO

**Elvio G. Amparore**
elviogilberto.amparore@**unito.it**

**Muhammad Rashid**
muhammad.rashid@**unito.it**

Rule**X**

**Enrico Ferrari**
enrico.ferrari@**rulex.ai**

**Damiano Verda**
damiano.verda@**rulex.ai**

# Supplementary Material