# SAFEXPLAIN

## Making certifiable AI a reality for critical systems:
## Core Demo

As fully autonomous systems become more and more crucial for advanced functions, AI takes on an increasingly more dominant role in critical systems. However, current AI software is developed as a black-box that is hard to verify, lacks clear requirements, and clashes with the rigorous, certifiable and explainable processes required by safety-critical systems.

The SAFEXPLAIN project is working to close this gap through its next-generation open software platform designed to show how AI can meet the stringent requirements of functional safety in autonomous technologies. The SAFEXPLAIN Core Demo is a direct, small-scale concretization of our end-to-end approach to making certifiable AI a reality for critical systems by ensuring that systems are safe by construction and follow the latest relevant standards.

This small scale, fully functional and configurable teaser shows how SAFEXPLAIN technology can accommodate scenarios with critical functionalities in three selected 'toy' examples from the automotive, rail and space domains. SAFEXPLAIN offers a pivotal concept in platform-level support. Its middleware works to ensure compliance with safety patterns by design while preserving modularity.

The SAFEXPLAIN technology deployed in this demo is beneficial for developers and decision-makers focused on transport and mobility in Critical Autonomous AI-Based Systems (CAIS) who need evidence of what type of evaluation and support will be available in the near future regarding the certifiability of their CAIS products.

The core demo provides a generic skeleton that accomodates simple, functionally relevant examples. It focuses on 'Safety Pattern 2' where an AI/ML constitutent partially affects the decision-process. Figure 1 visualizes this relationship.

## FuSa Compliant Safety Pattern 2

**Execution platform (HW, OS, libraries...)**

**AI-based sub-system**

**AI/ML Constituent**

DL Framework

DL Model

**Traditional sub-system**

Traditional SW item

---

**Human-machine interaction only potentially affects a safety function**

- Traditional FuSa risk factors (systematic/random)
- HP MPSoCs platform integration risk factors
- AI performance insufficiency
- AI FuSa risk factors (LOW/MED integrity levels)



Figure 1: Core demo skeleton

## Plug-in 'functional' code

- *ML nodes*
- *Supervisor*
- *Decision function*
- *Control logic*
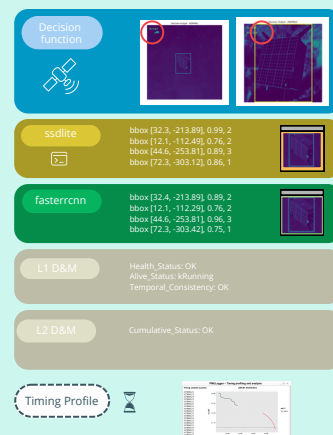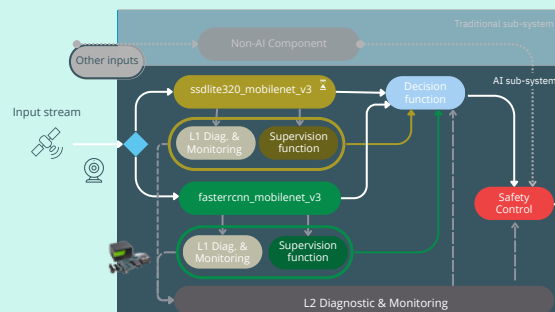
**Join the final event**

# How it works: Core demo instantiated in 3 cases

## SPACE DOMAIN

Goal: *Identify the target, estimate its pose, and monitor the agent position, to signal potential drifts, sensor faults, etc*

- AIKO open models + Input images
- L1 D&M Temporal consistency
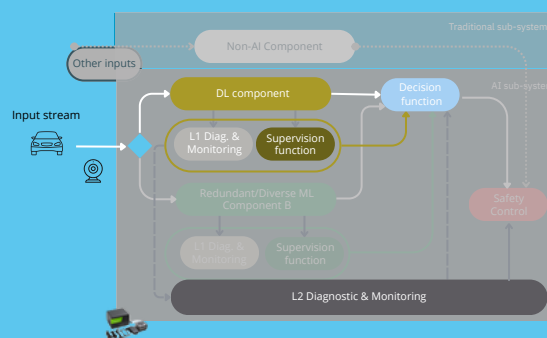- Supervision VAE
- Decision Ensemble



Non-AI Component · Traditional sub-system
Other inputs
Input stream
ssdlite320_mobilenet_v3
L1 Diag. & Monitoring · Supervision function
fasterrcnn_mobilenet_v3
L1 Diag. & Monitoring · Supervision function
AI sub-system
Decision function
Safety Control
Actuator
L2 Diagnostic & Monitoring

Decision function
ssdlite — bbox [32.3, -213.89], 0.99, 2 / bbox [12.1, -112.49], 0.76, 2 / bbox [44.6, -253.81], 0.89, 3 / bbox [72.3, -303.12], 0.86, 1
fasterrcnn — bbox [32.4, -213.89], 0.89, 2 / bbox [12.1, -112.29], 0.76, 2 / bbox [44.6, -253.81], 0.96, 3 / bbox [72.3, -303.42], 0.75, 1
L1 D&M — Health_Status: OK / Alive_Status: kRunning / Temporal_Consistency: OK
L2 D&M — Cumulative_Status: OK
Timing Profile

Target <---> Laptop over the network

## AUTOMOTIVE DOMAIN

Goal**:** *Validate the system's capacity to detect pedestrians, issue warnings, and perform emergency braking*

- YOLO v11 pretrained model + frames from CARLA scenario
- L1 D&M Temporal consistency
- Supervision Function from EXPLib (anomaly detectors)
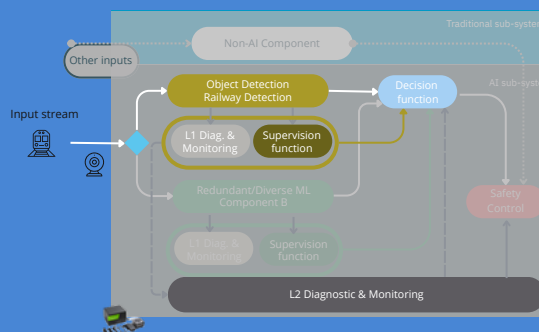- Decision Function (mainly visualization oriented)



Non-AI Component · Traditional sub-system
Other inputs
Input stream
DL component
L1 Diag. & Monitoring · Supervision function
Redundant/Diverse ML Component B
L1 Diag. & Monitoring · Supervision function
AI sub-system
Decision function
Safety Control
Actuator
L2 Diagnostic & Monitoring

Decision function
Supervision function — Anomaly_detected: YES / Anomaly_Type: 21 / Anomaly_Score: 0.785
Object Detection — bbox [32.3, -213.89], 0.99, 2 / bbox [12.1, -112.49], 0.76, 2 / bbox [44.6, -253.81], 0.89, 3 / bbox [72.3, -303.12], 0.86, 1
L1 D&M — Health_Status: OK / Alive_Status: kRunning / Temporal_Consistency: OK
L2 D&M — Cumulative_Status: OK

Target <---> Laptop over the network

## RAIL DOMAIN

Goal: *Validate the system's capacity to detect cars, issue warnings, and perform service braking*

- YoloV8 pretrained model + frames from Unreal Engine
- L1 D&M Temporal consistency Bounding boxes for detected obstacles
- Data-quality and temporal-consistency scores



Non-AI Component · Traditional sub-system
Other inputs
Input stream
Object Detection Railway Detection
L1 Diag. & Monitoring · Supervision function
Redundant/Diverse ML Component B
L1 Diag. & Monitoring · Supervision function
AI sub-system
Decision function
Safety Control
Actuator
L2 Diagnostic & Monitoring

Object Detection Railway Detection — bbox [32.3, -213.89], 0.99, 2 / bbox [12.1, -112.49], 0.76, 2 bbox [44.6, -253.81], 0.89, 3 / bbox [72.3, -303.12], 0.86, 1
L1 D&M — Health_Status: OK / Alive_Status: kRunning / Temporal_Consistency: OK
L2 D&M — Cumulative_Status: OK

Target <---> Laptop over the network