# EXPLAINABILITY FIRST!

## COUSTEAUING THE DEPTHS OF NEURAL NETWORKS

ES4CPS@Dagstuhl – Jan 7, 2019
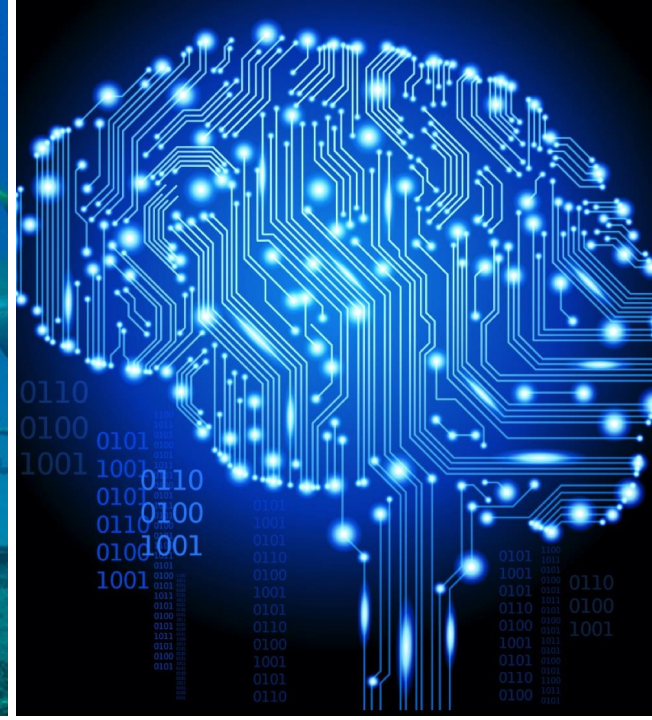
@mrksbrg

mrksbrg.com

Markus Borg

**RISE Research Institutes of Sweden AB**
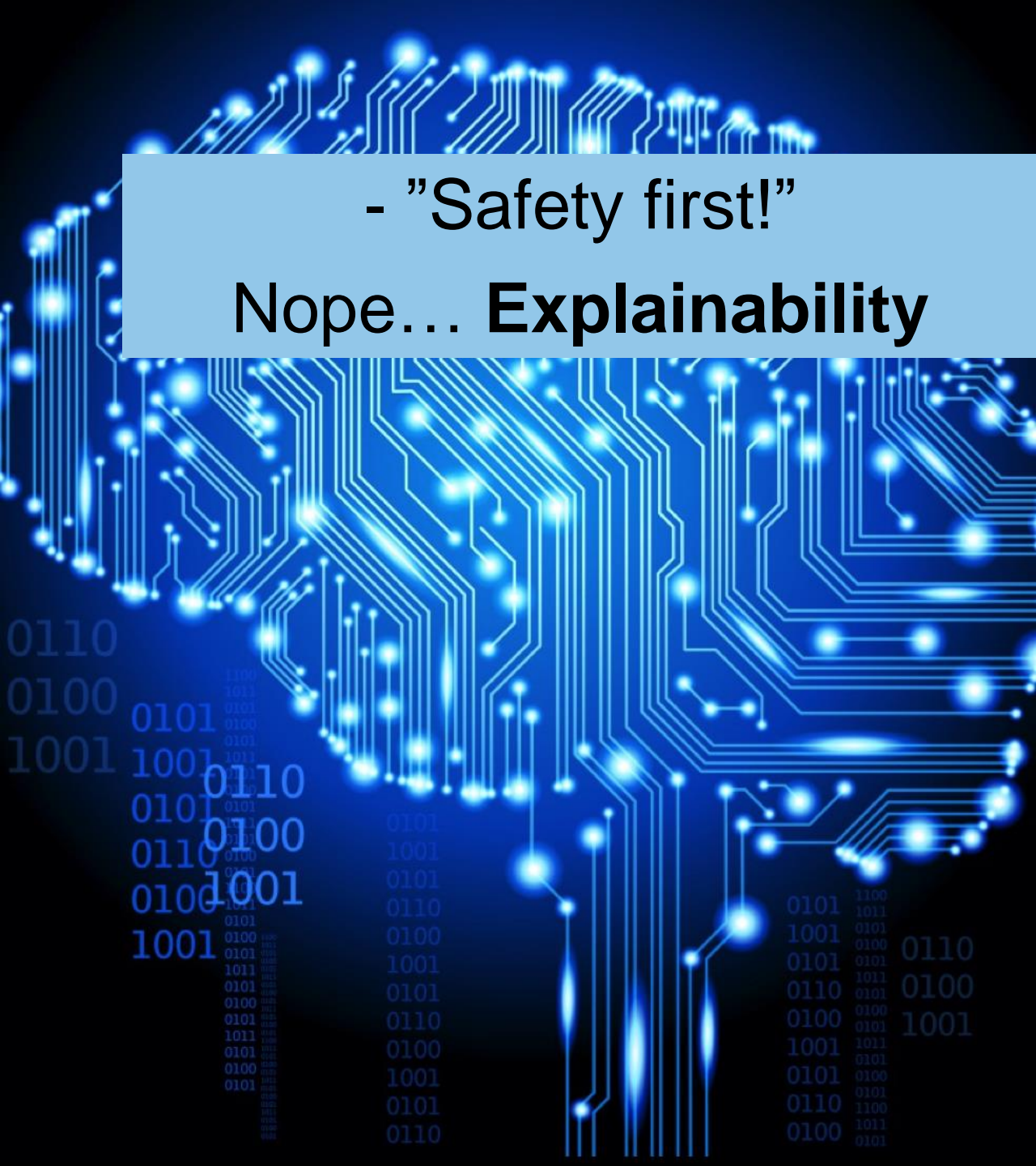
- "Safety first!"
Nope… **Explainability**

- "Aller voir!"

# Who is Markus Borg?

**Development engineer, ABB, Malmö, Sweden** 2007-2010

- Editor and compiler development
- Safety-critical systems

**PhD student, Lund University, Sweden** 2010-2015

- Machine learning for software engineering
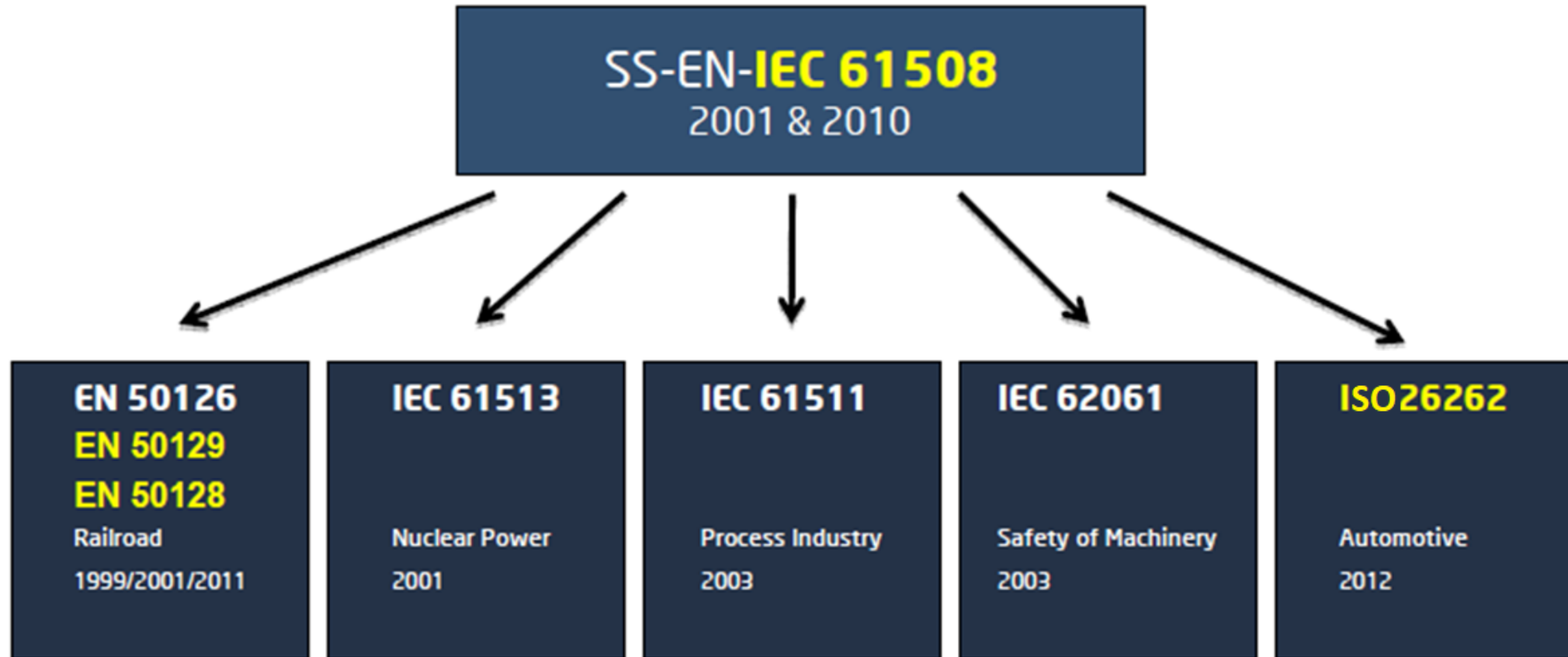- Bug reports and traceability

**Senior researcher, RISE AB, Lund, Sweden** 2015-

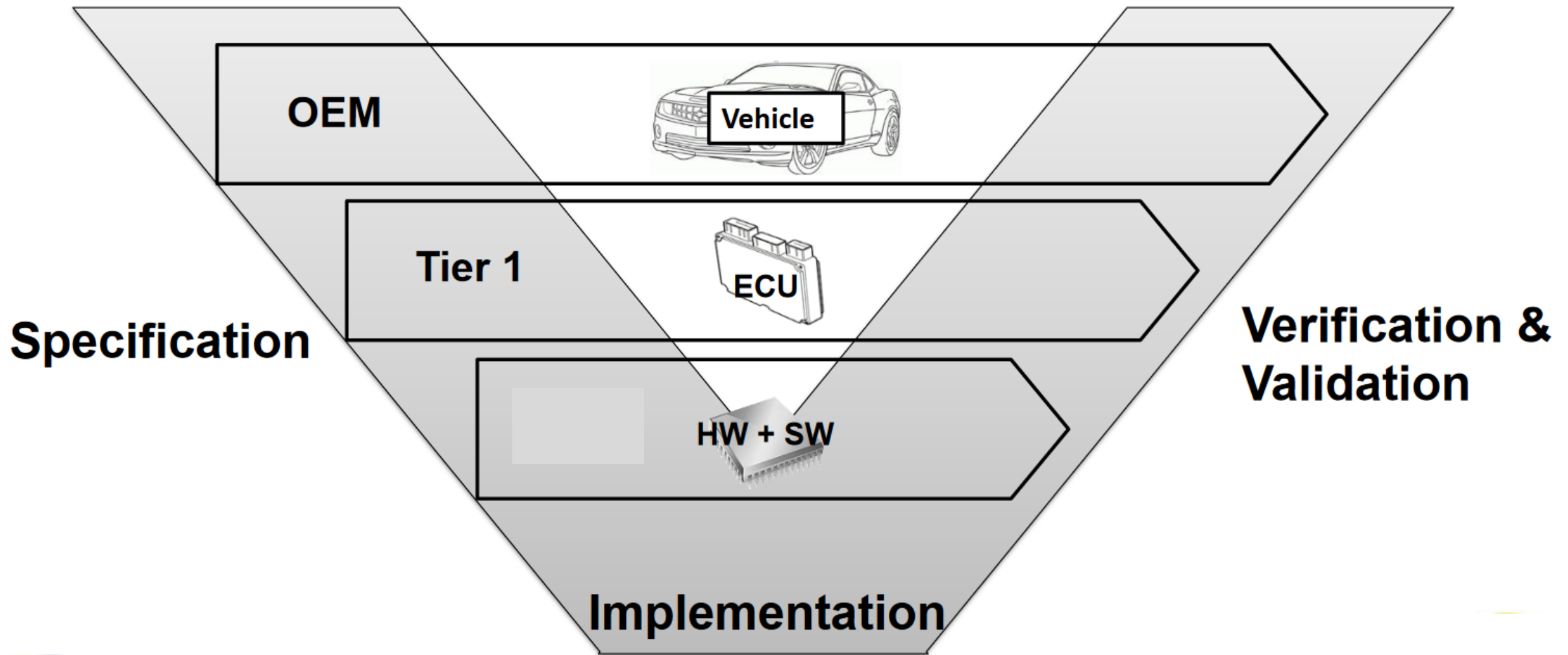- Software engineering for machine learning
- Software testing and V&V

# Background

# Functional Safety Standards

Specification

OEM

Vehicle

Tier 1

ECU

Verification &
Validation

HW + SW

Implementation

RI.
SE

# Achieving Safety in Software Systems

1. Develop understanding of situations that lead to safety-related failures
   - Hazard = system state that could lead to an accident

2. Design software so that such failures do not occur
   - Fault tree analysis

The system shall never hurt anyone
- even if the system does not conform to its specification

RI.
SE

# Safety certification => Put evidence on the table!

- Safety requirement: "Stop for crossing pedestrians"
- How do you argue in the safety case?

# Safety evidence – In a nutshell

- System specifications
  - and why we believe it is valid
- Comprehensive V&V process descriptions
  - and its results
  - coverage testing for all critical code
- Software process descriptions
  - hazard register and safety requirements
  - code reviews
  - traceability from safety requirements to code and tests
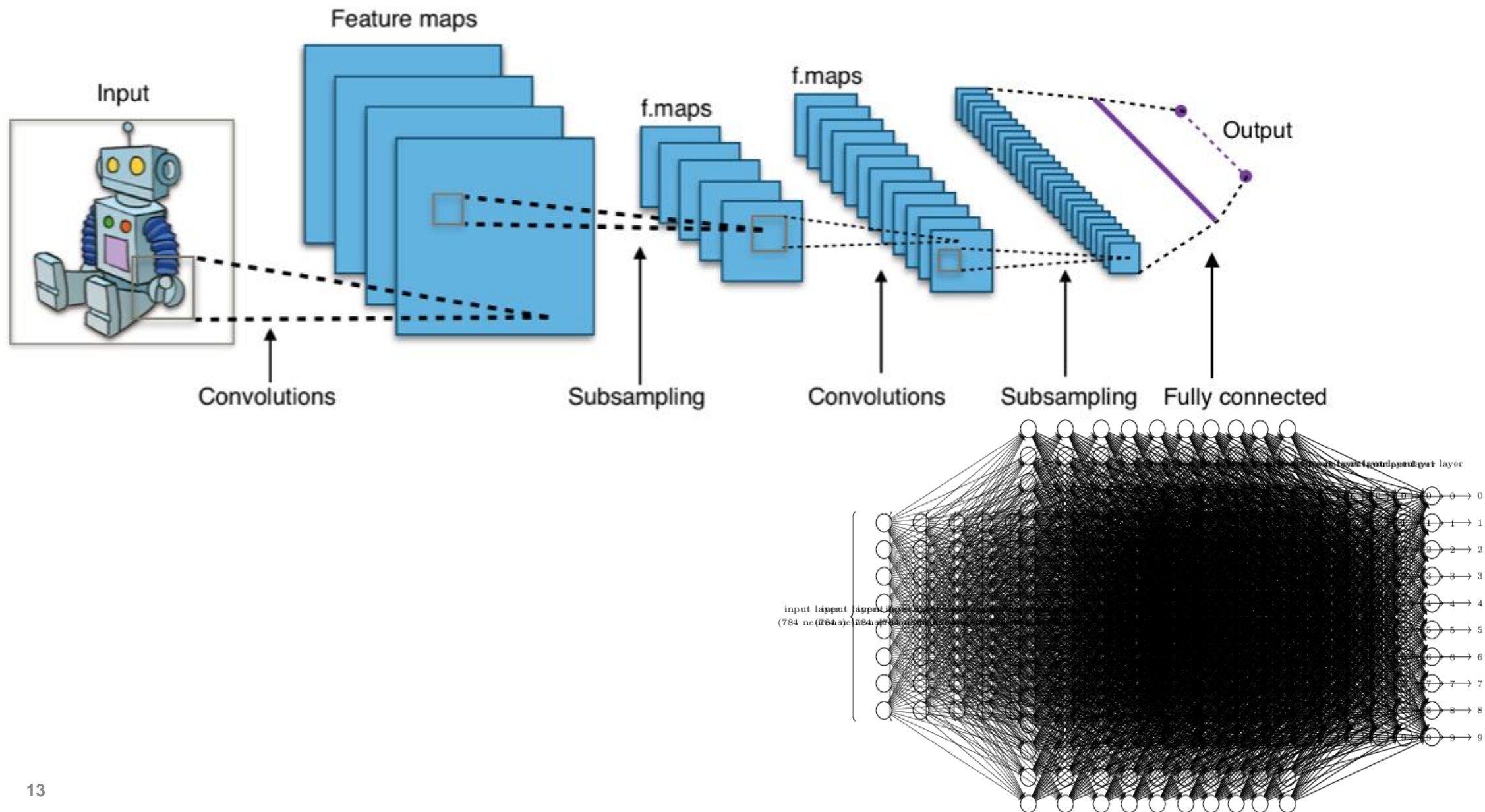  - ...

RI.
SE

# Application context

Safe-Req-A1:
In autonomous highway mode A, the vehicle shall keep a minimum safe distance of 50 m to preceding traffic

Realize vehicular perception using deep learning

# Autonomous Driving thanks to Convolutional Neural Networks
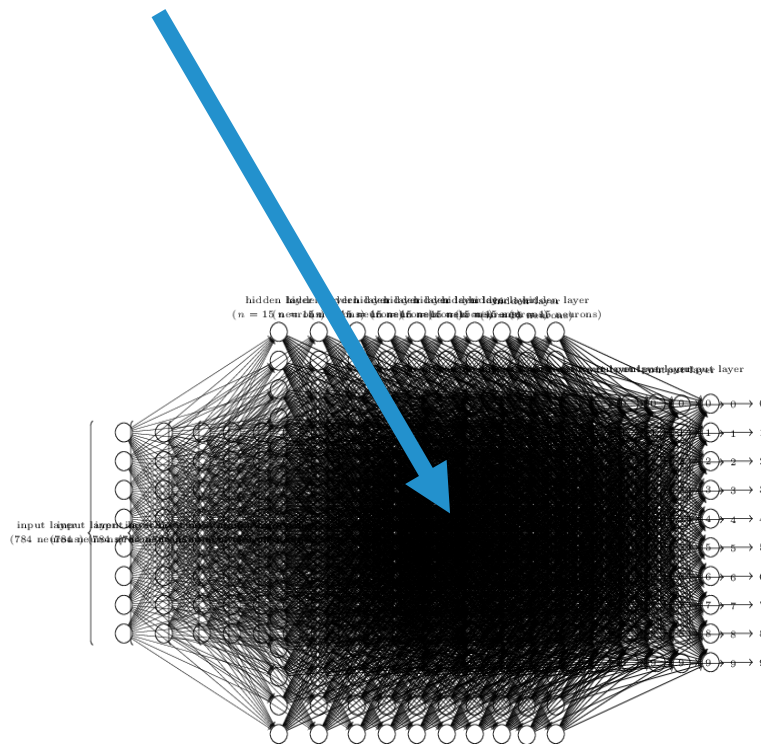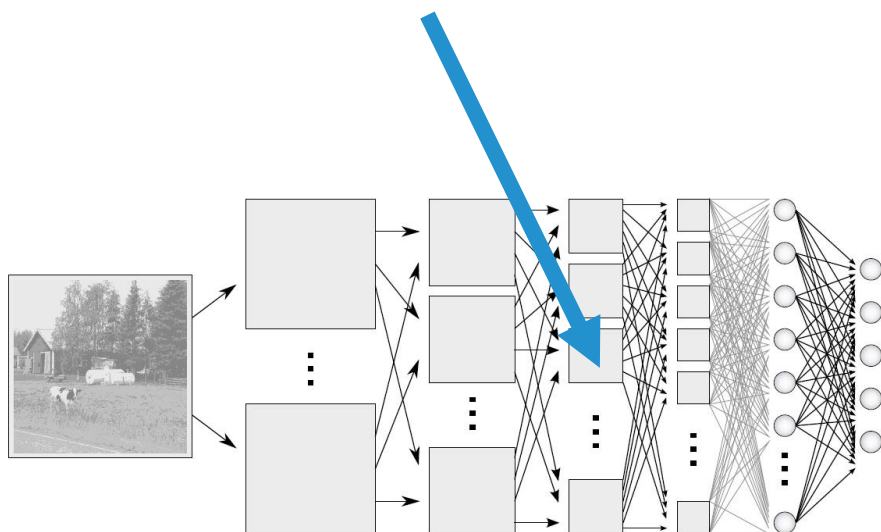
# Trace from Safe-Req-A1 to... what?

"Aller voir!"

# Trace from Safe-Req-A1 to… what?

**3) in training examples used to train and test the deep learning model**

**1) inside a human-interpretable model of a deep neural network**

**2) parameter values in a trained deep learning model**

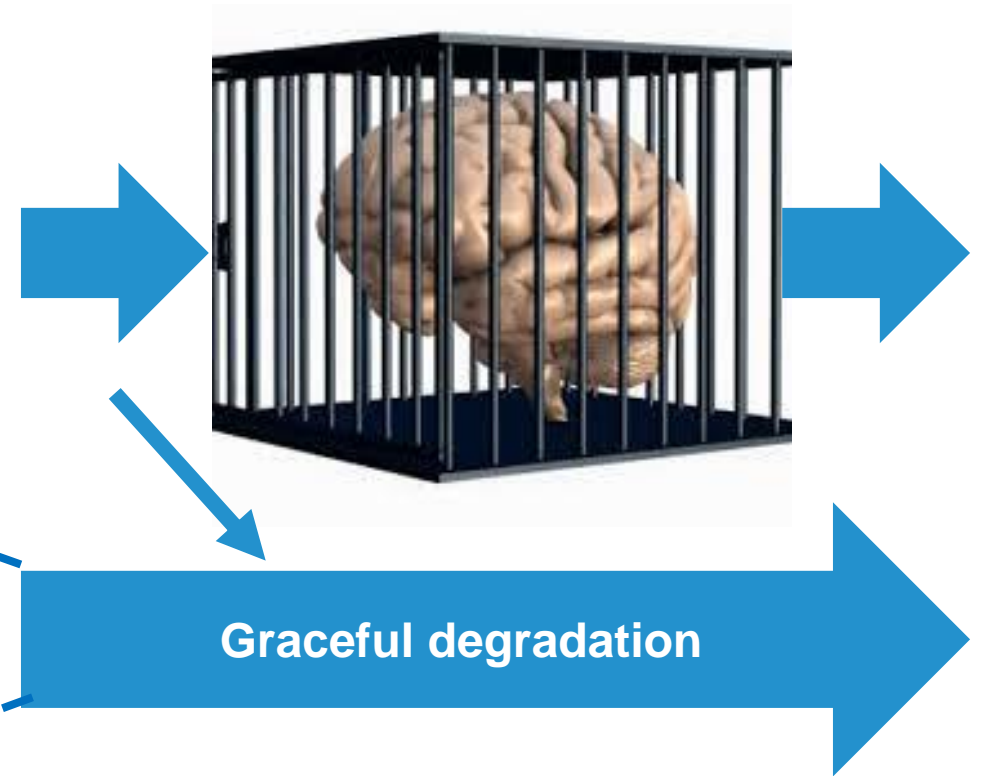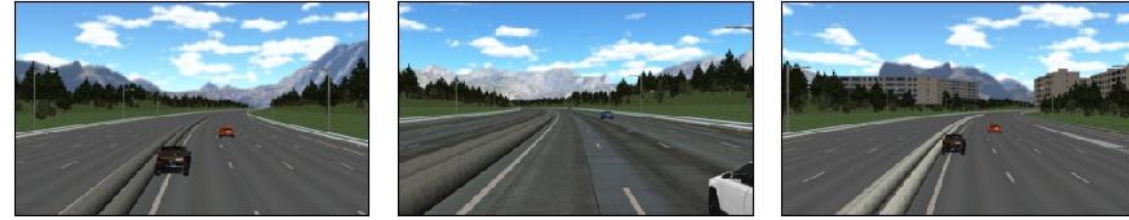# Open challenge

LET´S PUT OUR HEADS TOGETHER. TO KEEP AHEAD.

RI.
SE

# System feature - Autonomous highway driving
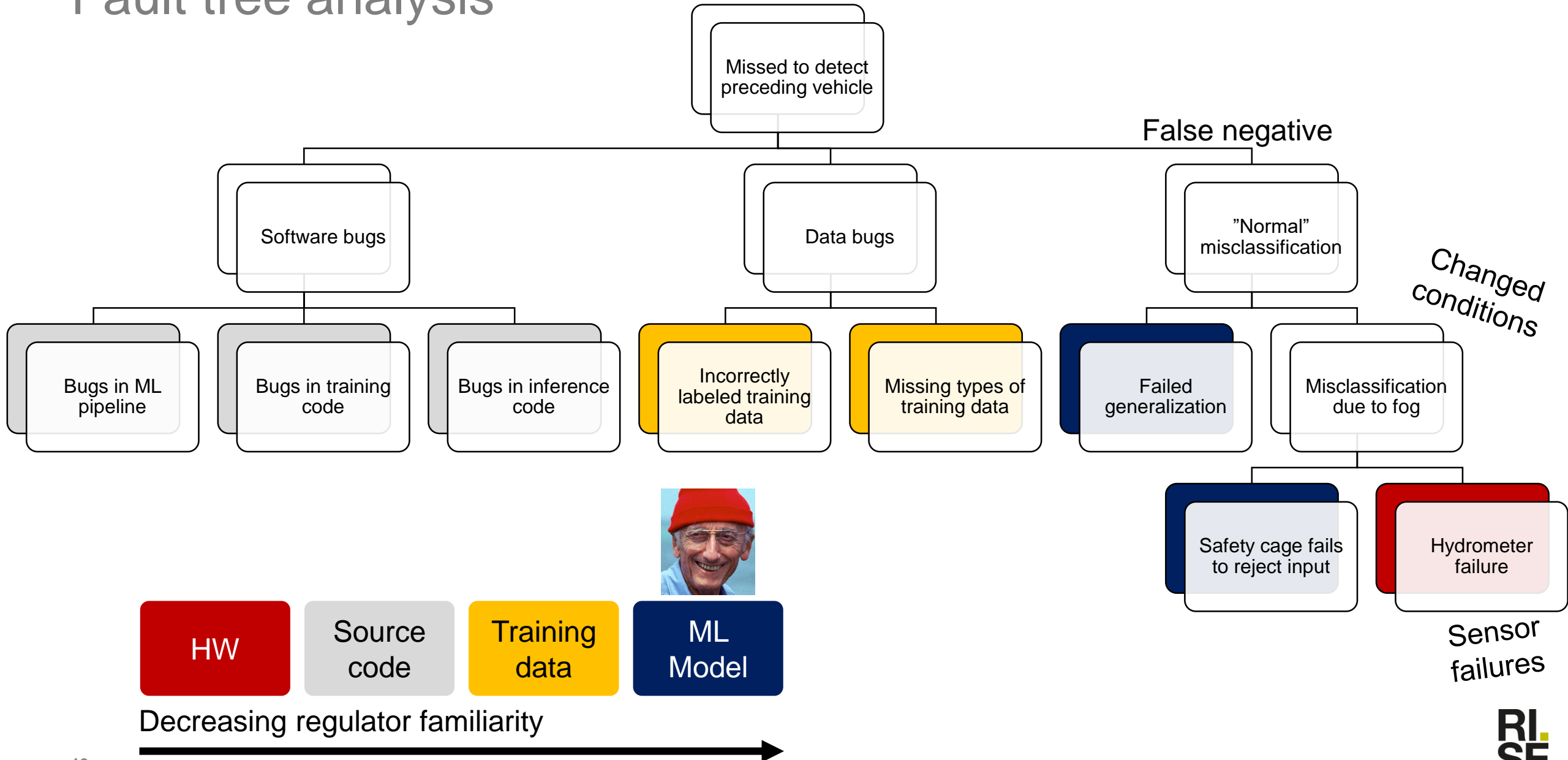
- FR1: … shall have an autonomous mode … in normal conditions…
- FR2: If the conditions change … shall request manual mode …
- FR3: If the driver does not comply … perform graceful degradation

RI.
SE

# Safety cage architecture

- Add reject option for deep network
  - Novelty detection

- Graceful degradation
  - turn on hazard lights
  - slow down
  - attempt to pull over

**Graceful degradation**

# Fault tree analysis

# Explainability additions

POST-DEPL ARNING...

- System specifications
  - CNN architecture, safety cage architecture
  - description of training data
- V&V process descriptions
  - training-validation-test split
  - neuron coverage
  - approach to simulation
- Software process extensions
  - new ML hazards advarsarial example mitigation strategy
  - traceability from all safety requirements to **data** and **code** and **tests**
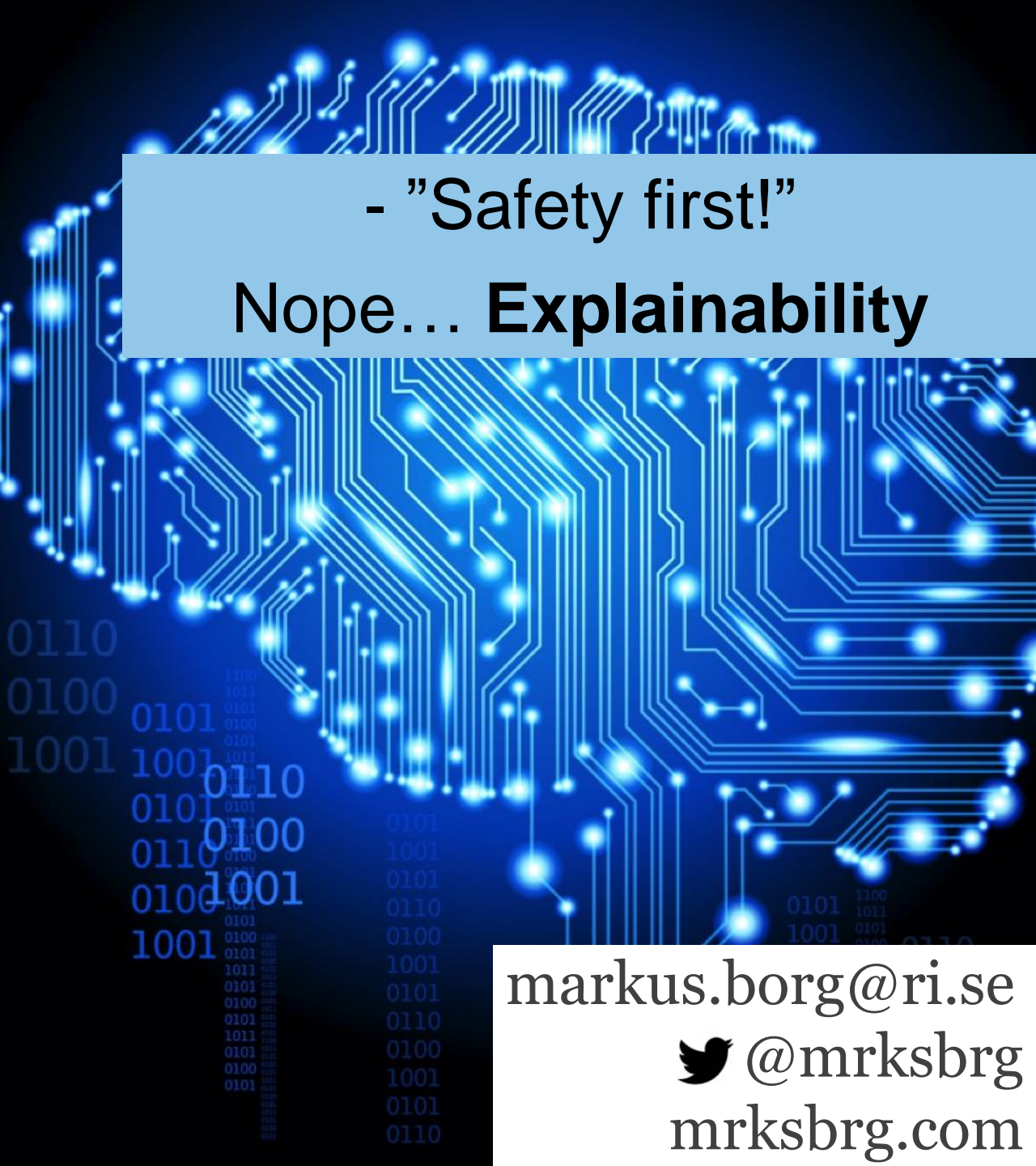  - staff ML training

Safe-Req-A1

- "Safety first!"
Nope... **Explainability**

- "Aller voir"

markus.borg@ri.se
@mrksbrg
mrksbrg.com