# Towards using ML in critical applications

#### András Pataricza Budapest University of Technology and Economics Dept. Measurement and information Systems With contributions of István Majzik, Zoltán Micskei, András Vörös and András Földvári

The research reported in this presentation was supported by the BME-Artificial Intelligence FIKP grant of EMMI (BME FIKP-MI)





Budapest University of Technology and Economics Department of Measurement and Information Systems

## Focal problem

- Al solutions are highly effective
- BUT: can we trust them?
- Black box?
- What are the impacts
  - Integrity of the decision



- Representativeness of the teaching set
- Coverage
- **IN CRITICAL APPLICATIONS**





## Challenges



Source: Andrej Karpathy: Building the Software 2.0 Stack

Dataset:

- Teaching
- Testing
- Debugging:

Interpretability/Explainability

- Using a 90% accurate model we understand, or
- 99% accurate model we don't.
- Faults/attacks, faulttolerance/defenses
- Safety
- Security





## **Hippocratic Oath**

I swear by <u>Apollo</u> the Healer, by <u>Asclepius</u>, by <u>Hygieia</u>, by <u>Panacea</u>, and by all the gods and goddesses, making them my witnesses, that I will carry out, according to my ability and judgment, this oath and this indenture.

- To hold my teacher in this art equal to my own parents;
- I will use treatment to help the sick according to my ability and judgment, but never with a view to injury and wrong-doing







# Explainable Artificial Intelligence (XAI)

#### DARPA-BAA-16-53

Mr. David Gunning : https://www.darpa.mil/program/explainable-artificial-intelligence





#### The Need for Explainable AI





#### XAI Concept



RG

T





MÚEGYETEM





#### Requirements

**Cyber-Physical Systems** 





## Cyber-Physical Systems definition

 "Cyber-Physical Systems or "smart" systems are co-engineered interacting networks of physical and computational components. These systems will provide the foundation of our critical infrastructure, form the basis of emerging and future smart services, and improve our quality of life in many areas."



National Institute of Standards and Technology U.S. Department of Commerce





Let's reach an unlimited intelligence by the synergy of

- intelligence in the cyber space and
- ES interfacing them to the physical world



# THE NEW ERA: V INTERNET OF THINGS AKA CYBER-PHYSICAL SYSTEMS



#### Self-\* properties – dynamic challenges



## **Critical applications**

- What we specified?
  - Safety function requirements: Function which is intended to achieve or maintain a safe state
  - Safety integrity requirements: Probability of a safetyrelated system satisfactorily performing the required safety functions (i.e., without failure)
- Safety Integrity Level and component fault rates
  - SIL 4: 10<sup>-8</sup> ...10<sup>-9</sup> faults per hour
  - Typical electronic components: 10<sup>-5</sup>...10<sup>-</sup>??3ults/hour
  - Typical software: 1..10 faults per 1000 line of code

OAI?



## Goals

- Requirements in critical systems: Safety, dependability
- Architecture design (patterns) in critical systems
- Focus: Design of system architecture to ...
  - Maintain safety
  - Even in the case if AI fails
- Fault-tolerant computing has 40+ years of experience building dependable systems out of unreliable components







## Objectives of architecture design



- Stopping (switch-off)
  is a safe state
- In case of a detected error the system has to be

stopped

Error detection is required

- Stopping (switch-off)
  is not a safe state
- Service is needed even in case of a detected error
  - full service
  - degraded (but safe) service
- Fault tolerance is required





#### Typical architectures for fail-stop operation





Budapest University of Technology and Economics Department of Measurement and Information Systems

#### Single channel architecture with built-in self-test

- Single processing flow with error detection
- Online self-checking



#### Two-channels architecture with comparison



- Shared input
- Comparison of outputs
- Stopping in case of deviation





#### **N-version**





м Ú Е С У Е Т Е М 1782

### **Recovery blocks**

- Passive redundancy: Activation only in case of faults
  - The primary variant is executed first
  - Acceptance checking on the output of the variants
    EXPLANATION
  - In case of a detected error another variant (TRADIONAL) is executed







#### Wrappers as by-products of testing?





#### Example assumption: data quality



#### Easy to check







#### Analysis techniques overview





4



#### **Run-time verification**

- Test Generator
  ORunning System
- System Under Test

System
 Components

Test Oracle



OMG Standard

**Extensions** 

- Publish-subscribe
- 15 QoS properties



RG

Т

#### Summary



Budapest University of Technology and Economics Department of Measurement and Information Systems



#### Summary

- Proper wrappers can eliminate ML induced failures
  - Fallback to traditional

Checking explanations is easier

Fault -tolerant computing has libraries



