# RC3: Resilient Computing and Cybersecurity Center

جامعة الملك عبدالله
للعلوم والتقنية
King Abdullah University of
Science and Technology

## AI for Cybersecurity vs. Cybersecure AI: a chicken and egg problem?

https://rc3.kaust.edu.sa
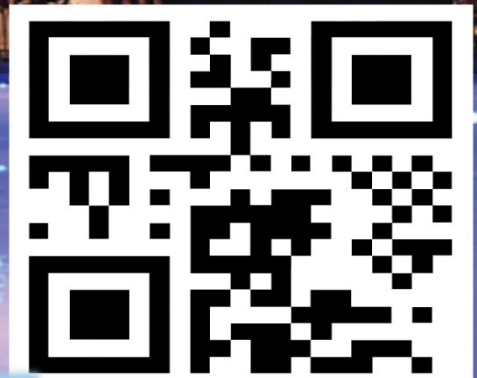
**Paulo Esteves-Veríssimo,** Professor, Director

King Abdullah University of Science and Technology, CEMSE
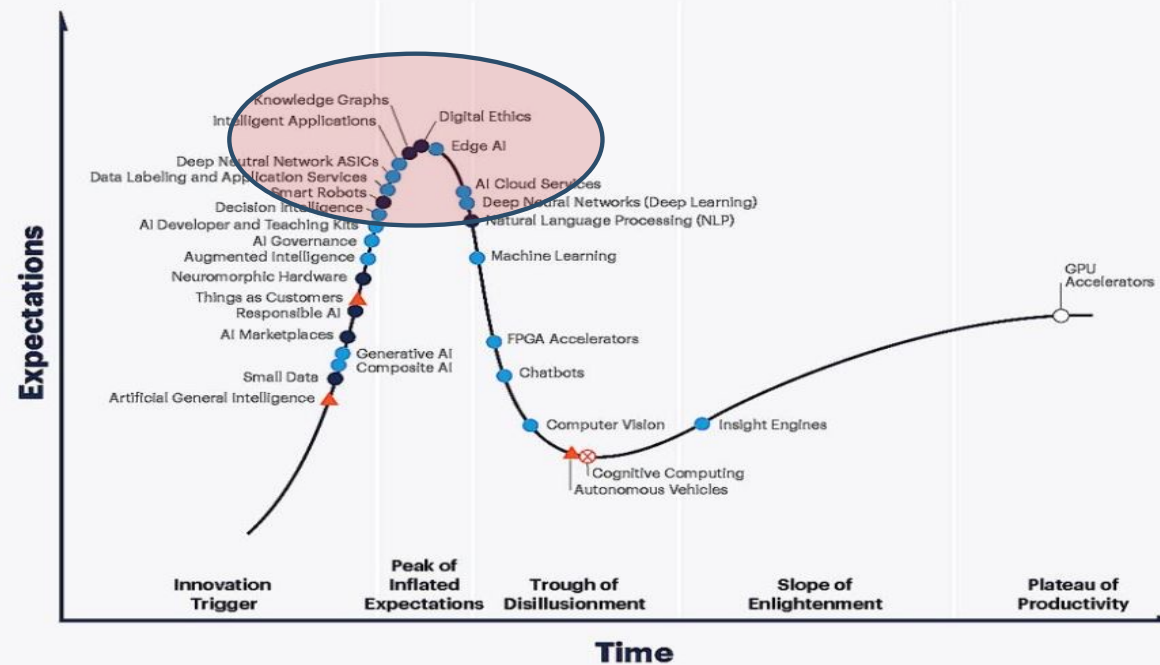Resilient Computing and Cybersecurity Center – RC3

paulo.verissimo@kaust.edu.sa

https://cemse.kaust.edu.sa/people/person/paulo-verissimo

*Int'l Workshop on Workshop on Cyber Resilient Systems, 4-7th June –Cascais, Portugal*

# Enter AI, ML, DNN, ...



Hype Cycle for Artificial Intelligence, 2020

# AI vs. Security vs. Safety



AI/ML for Cybersecurity and Safety

Cybersecure and Safe AI/ML

# Homogeneous ML/DNN-based systems cannot give strong assurance guarantees

- **Status-quo**
  - *Autonomous cars use ML-powered multi-sensor perception and complex control logic, and sometimes redundant modules to which they hand over in case of problems.*

- **Assurance**
  - *Infeasible to provide reliable figures/conclusions*
  - *Impossible to certify under current best practices*

Tesla radar did not recognize a camel, causing an accident in the UAE

# Can we leverage the best of the *security and dependability* fields?

- In particular, dependability teaches us that our techniques:
  - (i) should identify the uncertainties and weaknesses exhibited at **component level**,
  - and (ii) craft mechanisms that address them, to produce predictably correct **system-level** results.
- *Result (ii) always conditioned by how well we did (i)*

# Hybrid ML-based systems may help
## an autonomous vehicles ecosystem example

- **Redundancy –** *these components cooperate redundantly to achieve the end goal of safety*
  - *Replication*
  - *Reconfiguration, hand-over*
  - *Take-over*
  - *Diversity, for malicious faults*

- **Hybrid architecture -** *Autonomous cars having different realms running under different assumptions*
  - *Hybrid system and fault assumptions ("hierarchy of functions")*
  - *Modular*
  - *Distributed.*

- **Assurance –** *enablers of the goal*
  - *Recent Hybrid Logic of Events allows verifying architecturally hybrid systems by proof assistants.*
  - *Trusted-trustworthy hybrids anchor the global trustworthiness, through proof of the Lifting predicate*

# Hybridisation-aware distributed algorithms, models, and architectures

جامعة الملك عبدالله
للعلوم والتقنية
King Abdullah University of
Science and Technology

Hybridisation-aware algorithms: models, architect., and control

*Leveraging trusted-trustworthy components and TEE, with the right set of simple functions (failure detectors, monotonic counters, reliable timers and clocks, PRG, signatures, indelible logs, binary consensus)*

Formally defined Secure interfaces

R1

Ultimately trusted hybrids

T

*Hybridisation-aware Distributed Algorithms for RESILIENCE:* Redundancy mgt. for error detection, recovery, masking, self-healing, etc.

R2

T

R3

Just hardened q.b. Payload

T

R4

T

*[Paulo Verissimo, "Travelling through Wormholes: a new look at Distributed Systems Models", SIGACT News, vol. 37-1, Mar. 2006. ]*

# KARYON architecture:
# proof of concept of hybridisation for safety

▸ **Main Concepts:**
- Level of Service
- Data validity
- Safety kernel

KARYON Workshop, Borås, Sweden, Dec 11, 2014

# Karyon Architecture



Wireless Network
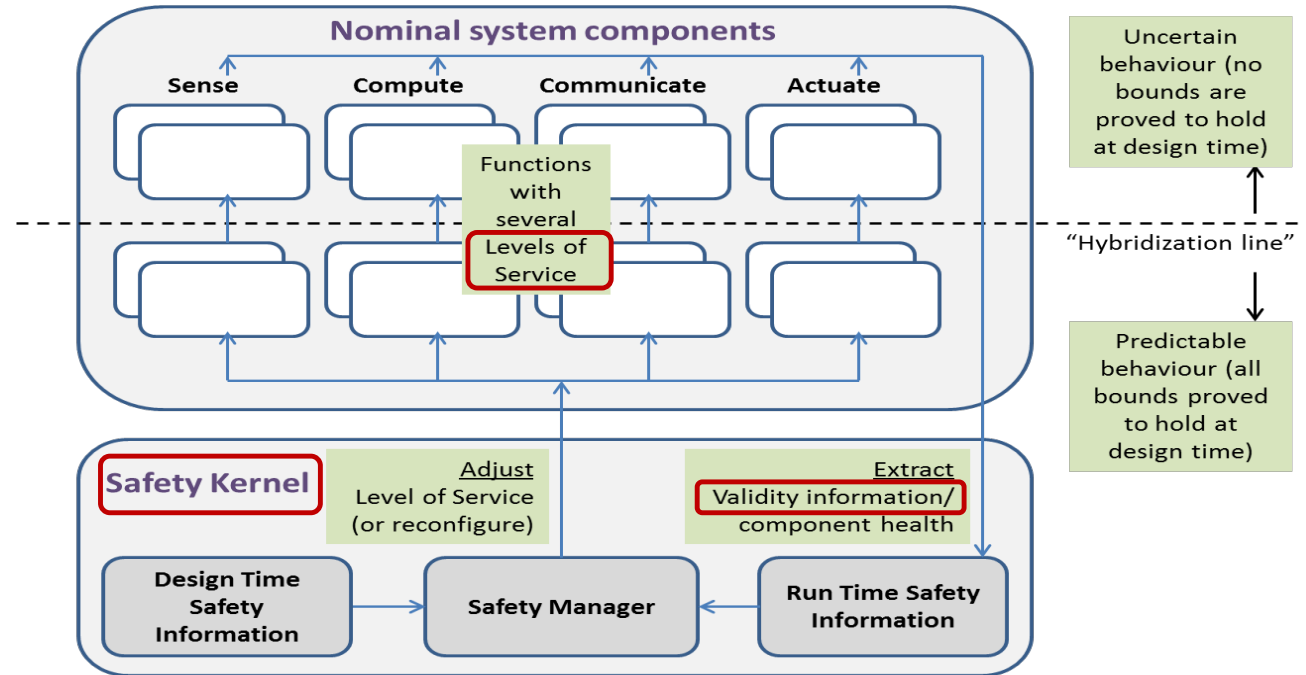
HYBRIDISATION BY LEVEL OF SERVICE
*Proven functional safety, as long as assumptions for level are satisfied*

External Smart Sensor

Perception: Information from External Sensors & Systems

State & Intentions

**Reconfigurable & Reliable Perception**

Internal Smart Sensor

Adaptive fusion

Event Selection

System Model

**Complex Control**

Functional level 3

Functional level 2

Functional level n

Bank of Control Functions

**Safety kernel**

Safety Element

Safety Element

Safety Element

guarded actuation

ABSTRACT SENSOR MODEL
*integrity attributes are evaluated against a set of different safety requirements*

SAFETY KERNEL HYBRID
*takes over, toward fail-operational or fail-safe termination*

(hidden) perception-reaction-channel

Environment

KARYON Workshop, Borås, Sweden, Dec 11, 2014