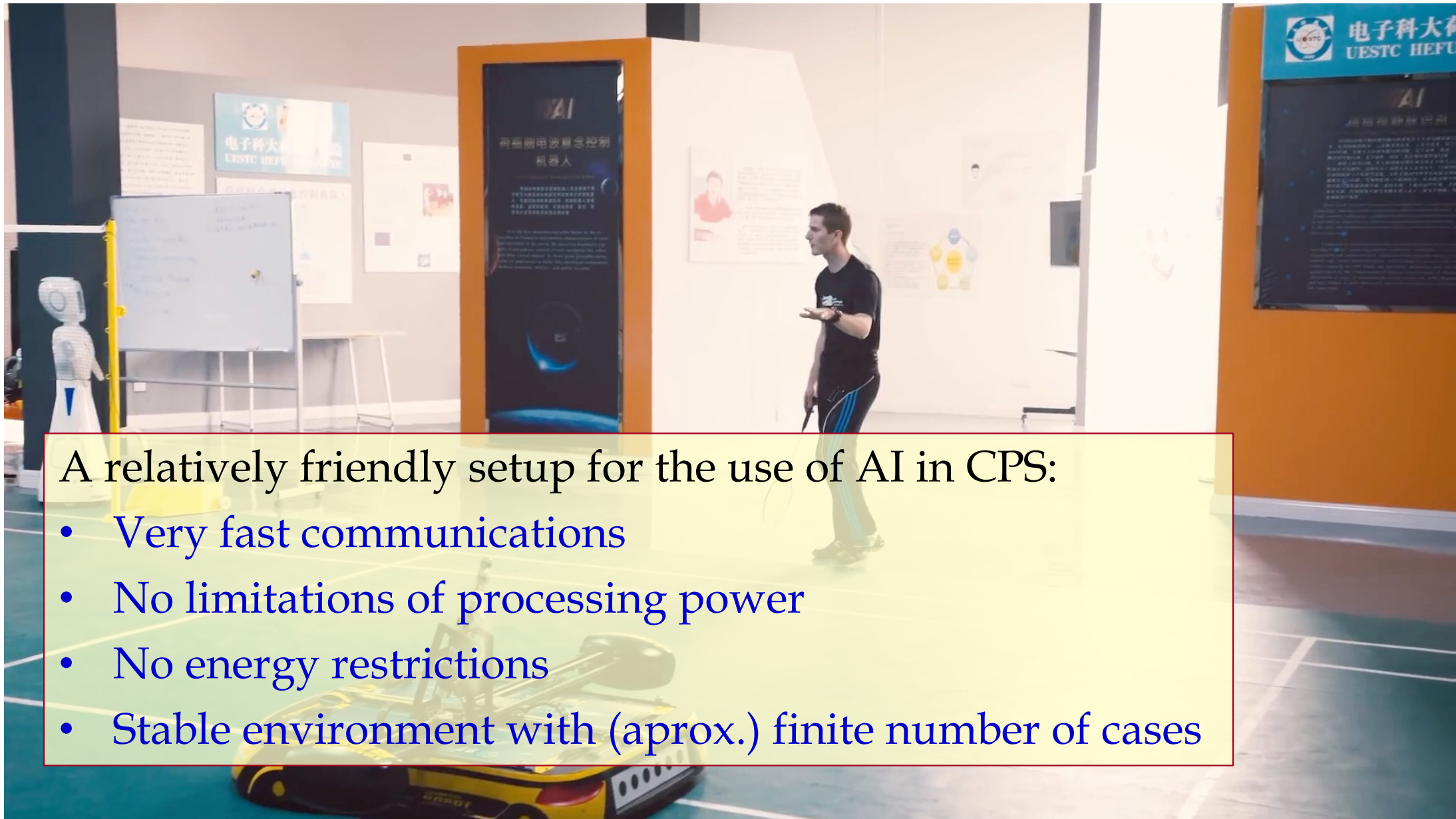


***Safety and security
in AI enabled critical applications:
who is going to solve the problems?***

**Henrique Madeira
University of Coimbra, Portugal**

80th Meeting of the IFIP 10.4 Working Group on
Dependable Computing and Fault Tolerance
Virtual - 25 June 2021 – 27 June 2021









A relatively friendly setup for the use of AI in CPS:

- Very fast communications
- No limitations of processing power
- No energy restrictions
- Stable environment with (aprox.) finite number of cases



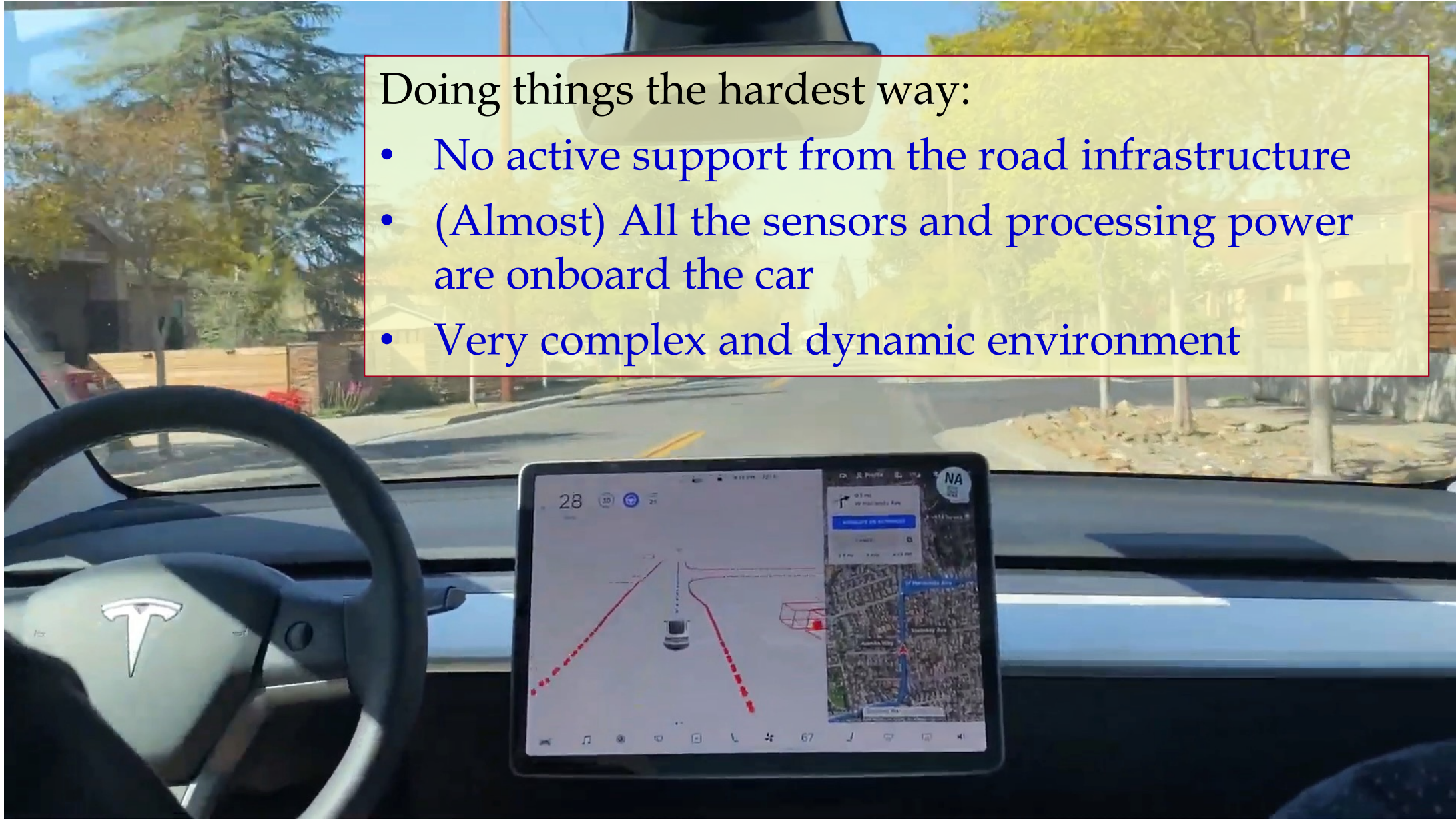
SAE AUTOMATION LEVELS

Full Automation

					
0	1	2	3	4	5
No Automation	Driver Assistance	Partial Automation	Conditional Automation	High Automation	Full Automation
Zero autonomy; the driver performs all driving tasks.	Vehicle is controlled by the driver, but some driving assist features may be included in the vehicle design.	Vehicle has combined automated functions, like acceleration and steering, but the driver must remain engaged with the driving task and monitor the environment at all times.	Driver is a necessity, but is not required to monitor the environment. The driver must be ready to take control of the vehicle at all times with notice.	The vehicle is capable of performing all driving functions under certain conditions. The driver may have the option to control the vehicle.	The vehicle is capable of performing all driving functions under all conditions. The driver may have the option to control the vehicle.

Doing things the hardest way:

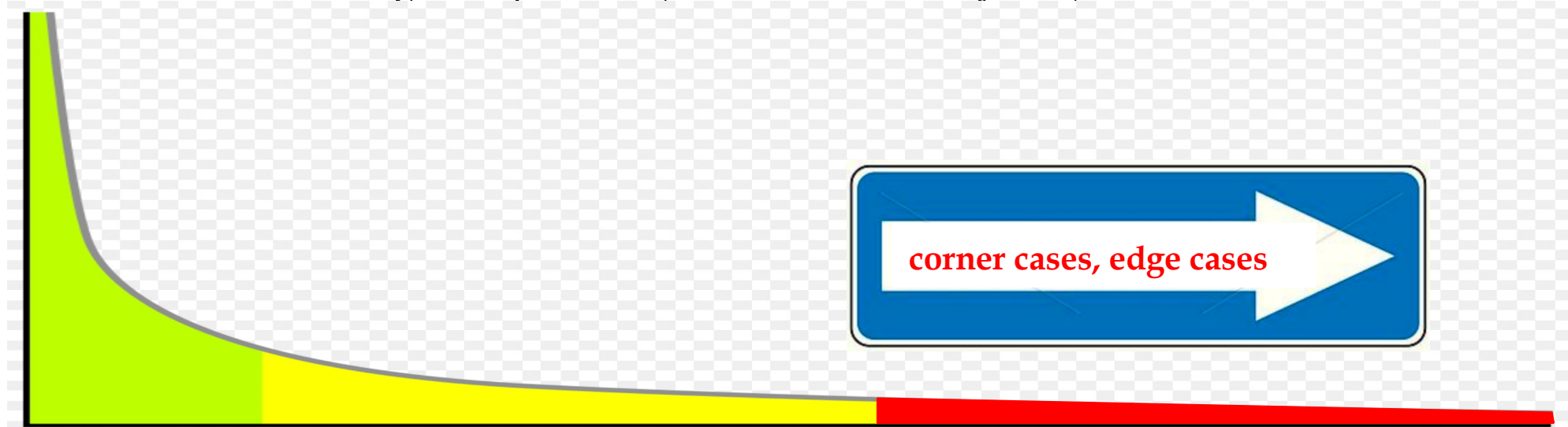
- No active support from the road infrastructure
- (Almost) All the sensors and processing power are onboard the car
- Very complex and dynamic environment



Difficulties in AI enabled critical applications

AI/ML used in safety-critical functions:

- ◆ Lack of clear functional specifications
- ◆ Non-deterministic and probabilistic outputs
- ◆ Limitations of the training data
- ◆ Non-explainable ML (i.e., black box)
- ◆ Exhaustive testing is impossible (as usual in ordinary SW) but in addition to that ML



Additional (classic) difficulties

- **Software faults**

- ◆ Defect densities remain nearly the same (i.e., high) for decades
- ◆ Many CPS have now millions of lines of code. A modern car, for example, has > 100 millions lines of code.

- **Hardware faults**

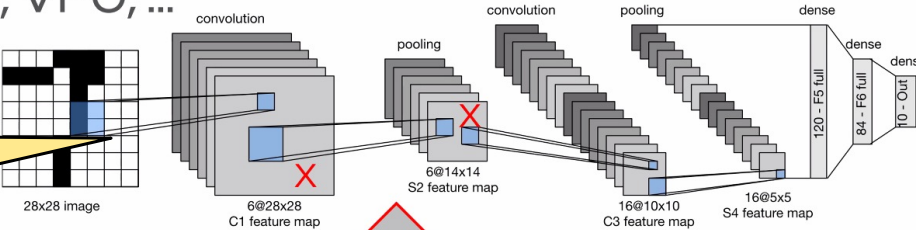
- ◆ “Silicon defects are getting worse”, Michael Paulitsch (Intel). The International Technology Roadmap for Semiconductors (ITRS) says the same.
- ◆ AI in safely critical applications needs massive hardware → increases the rate of hardware faults

Additional (classic) difficulties

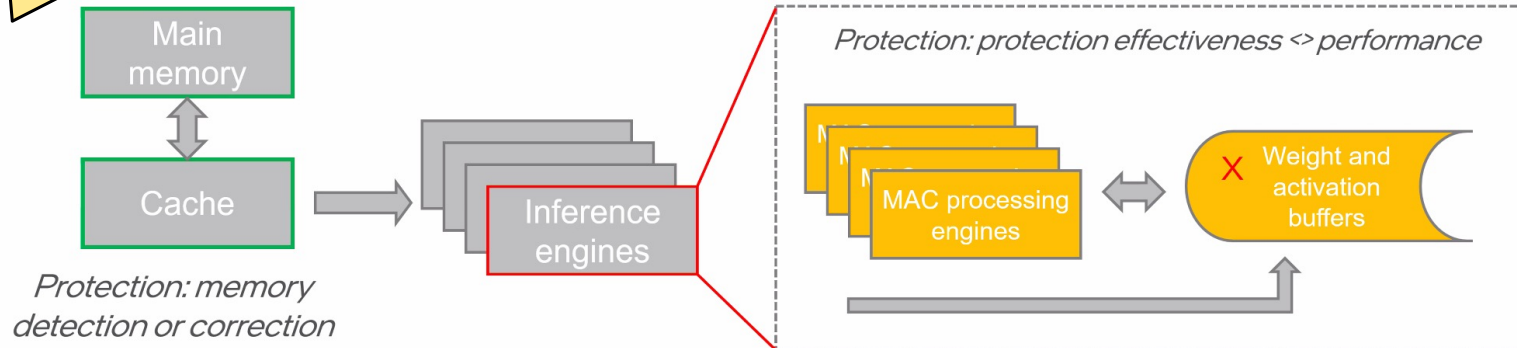
Common Elements in Intel CNN AI/ML accelerator HW

CPU extensions, GPGPU, VPU, ...

From Michael Paulitsch's slides



459



chips in
duced by

ernational
es the rate

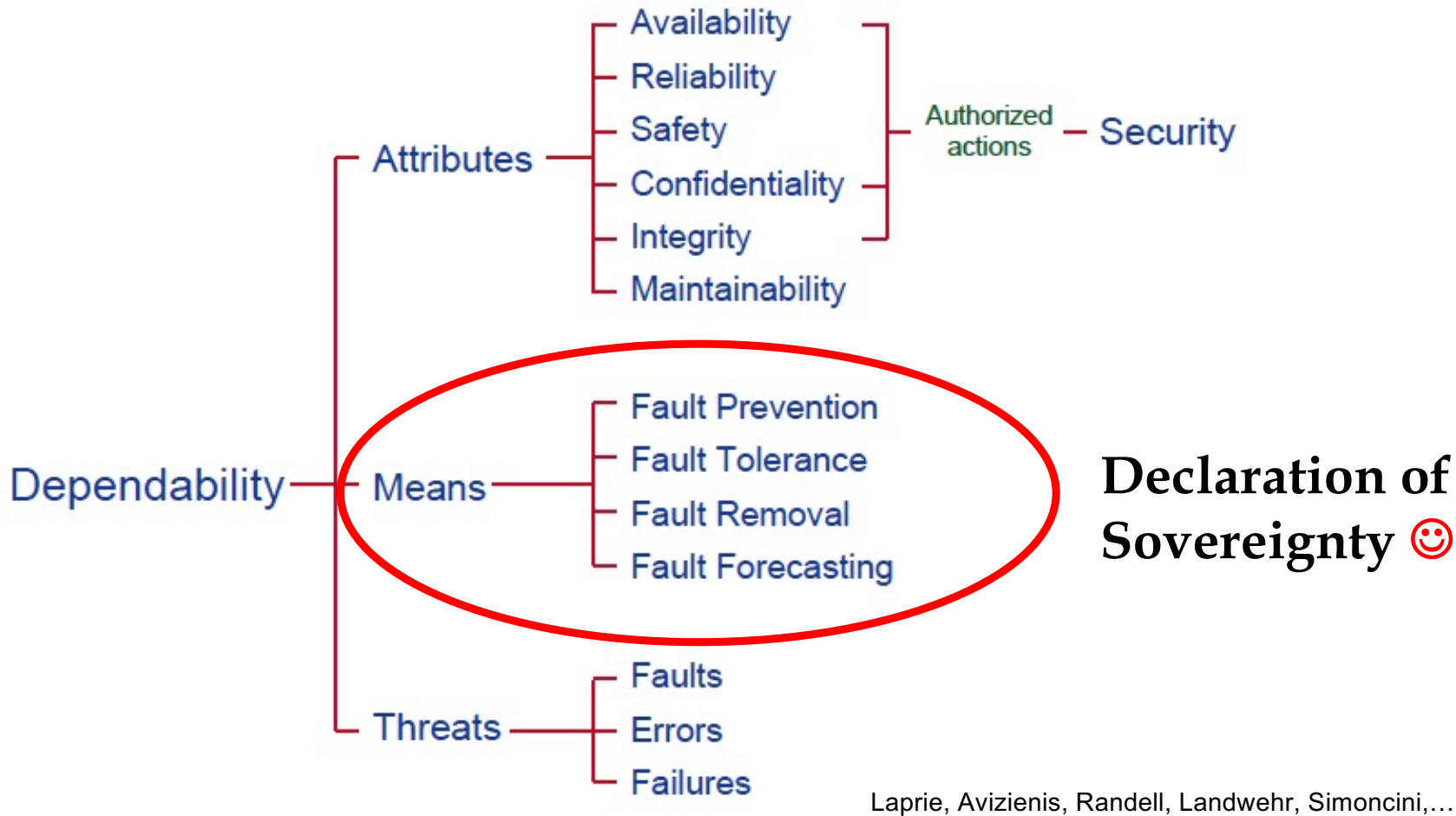
Full self-driving: is it a classic problem?

- **Building dependable systems using components that are not perfect looks like a classic problem:**
 - ◆ The output accuracy of AI components (in the absence of faults) is probabilistic (specially for black-box AI)
 - ◆ All components are subject to software faults
 - ◆ HW faults must be considered
 - ◆ AI used in safety-critical applications is an interdisciplinary problem, no matter the application area (automotive, medical devices, industry 4.0, avionics, etc.).

Who is going to solve the problems?

(What is expected from our research community?)

Cartography of (our) Dependability World



Cartography of (our) Dependability World



Who is going to solve the problems?

Problems:

- ♦ **How to assure safety and security** in AI enabled safety-critical applications?
- ♦ **How to demonstrate** that one can trust on AI enabled safety-critical applications?

Can we do that for self-driving cars?

- Millions of vehicles
- Billions of driving hours
- Huge pressure to cut cost
- Very high criticality

Who is going to solve the problems?

Artificial intelligence



The solution is more and better AI

Problems:

- ◆ **How to assure safety and security** in AI enabled safety-critical applications?
- ◆ **How to demonstrate** that one can trust on AI enabled safety-critical applications?

More

- Robust AI models
- Non-symbolic AI
 - Larger training data sets and bigger and more complex neural networks
 - Interpolation vs extrapolation
- Explainable AI
- Ensembles
- ...

Who is going to solve the problems?

Problems:

- ◆ How to assure safety and security in AI enabled safety-critical applications?
- ◆ How to demonstrate that one can trust on AI enabled safety-critical applications?

The solution is more and better AI

The solution is in the process

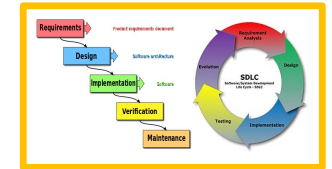
The solution is in models and tools

- The goal is to be able to measure software reliability
- Software reliability growth models... and AI

Artificial intelligence



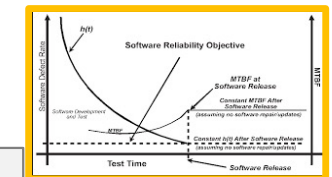
Software Engineering



Standards



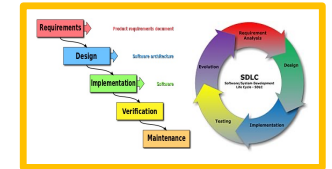
Software Reliability



Who is going to solve the problems?

Artificial intelligence

Software Engineering



The solution is more and better AI

The solution is in the process

Standards

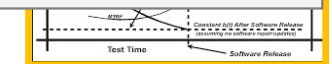


Problems:

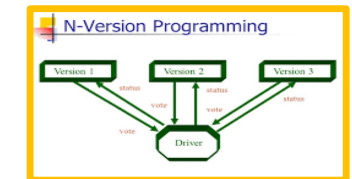
- ♦ How to assure safety and security in AI enabled safety-critical applications?
- ♦ How to demonstrate that one can trust on AI enabled safety-critical applications?

- Build dependable systems with unreliable components
- Can we use old recipes?

The solution is in the architecture



Dependability



Architecture as (part of) the solution

Black-Box Monitoring

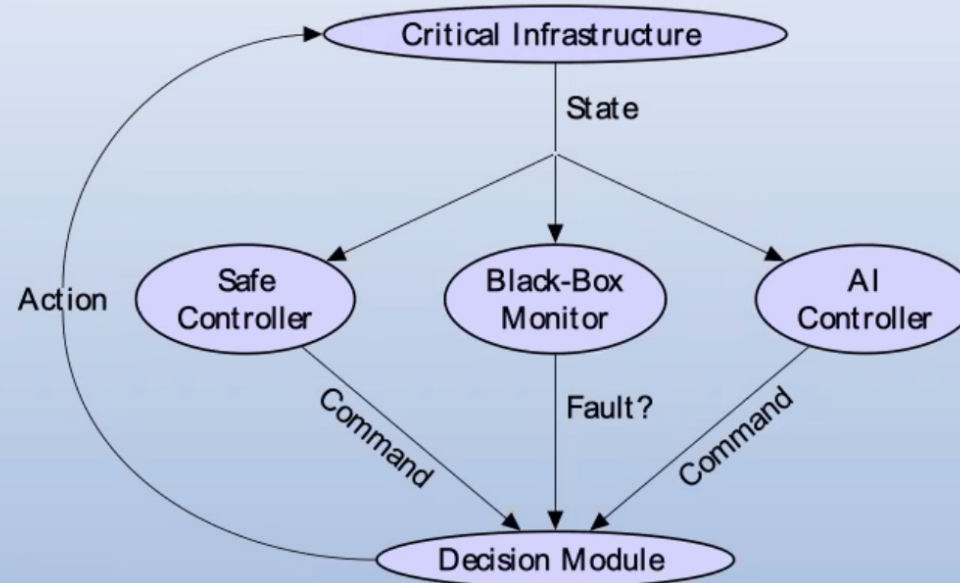
Black-Box monitoring is a standard approach to create dependable systems

The systems work roughly as follows:

State is collected and passed to a trusted controller, an AI controller and a monitor.

Each controller proposes an action

The decision module uses the output of the monitor to determine which action should be performed

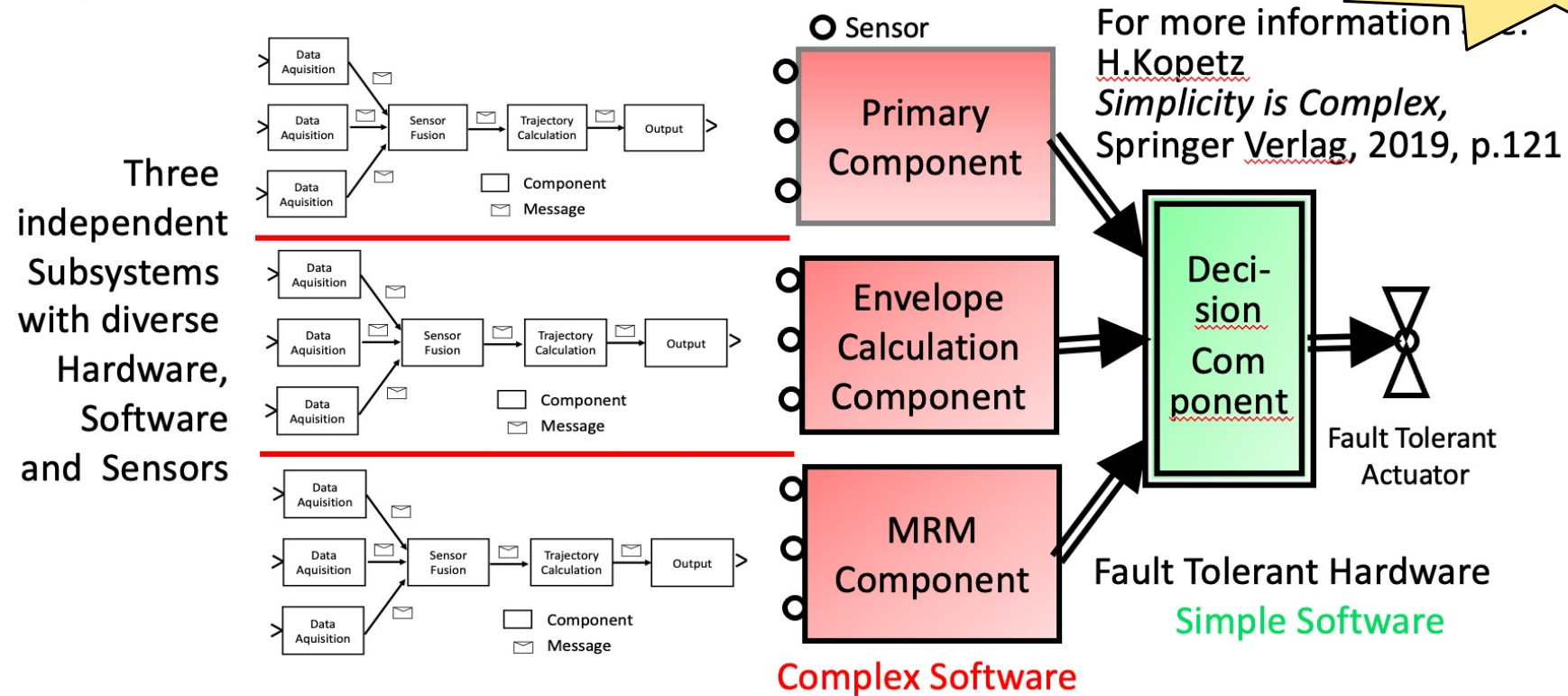


From Yair Amir's slides

Architecture as (part of) the solution

Sketch of an ADS for Unsupervised Autonomous Control

From Herman Kopetz's slides



Market and industry are not always rational

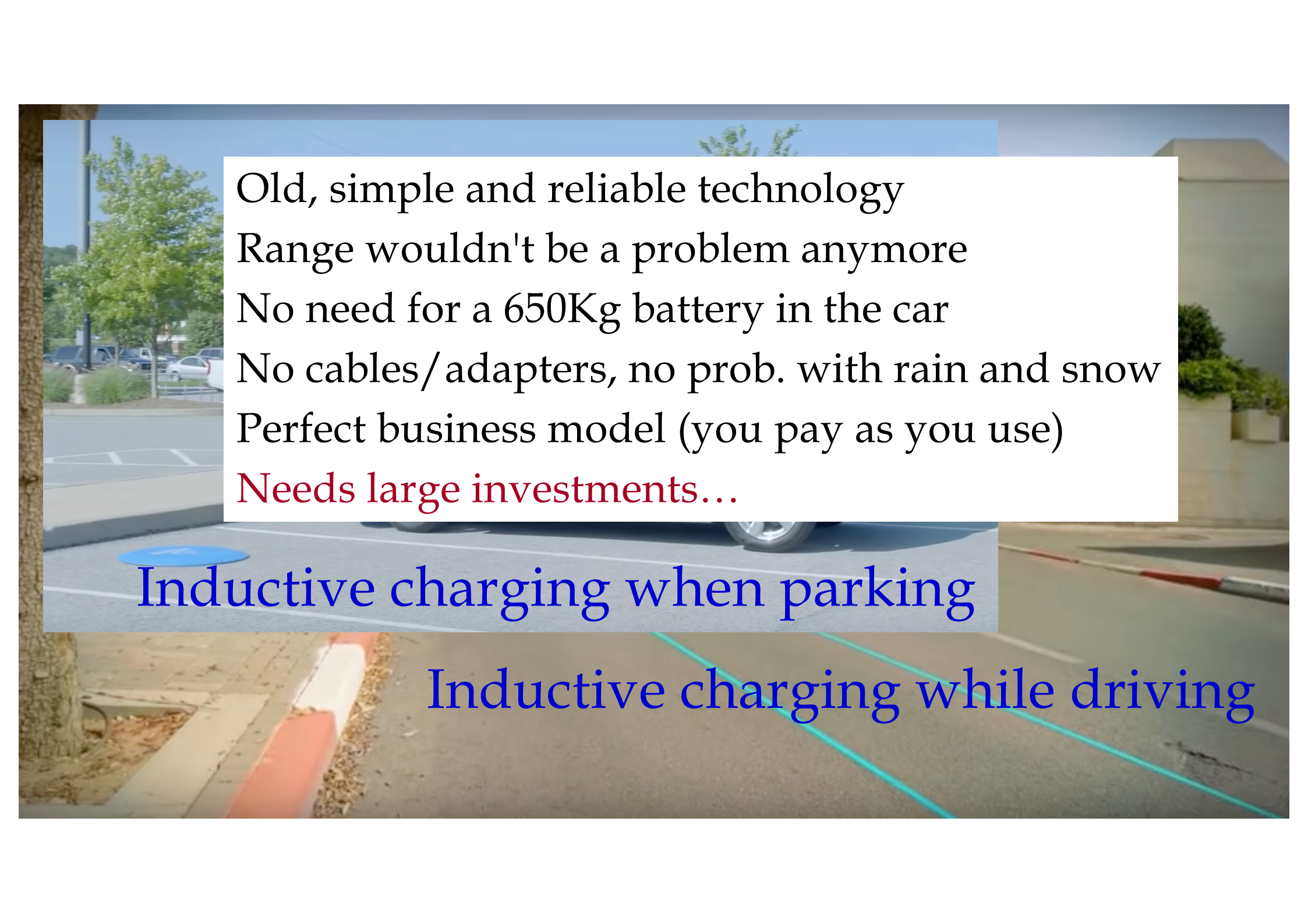




Inductive charging when parking



Inductive charging while driving



Old, simple and reliable technology
Range wouldn't be a problem anymore
No need for a 650Kg battery in the car
No cables/adapters, no prob. with rain and snow
Perfect business model (you pay as you use)
Needs large investments...

Inductive charging when parking

Inductive charging while driving

Market and industry are not always rational

The quest for long-range and fast-charging in EVs:

- Largely unregulated and uncooperative.
- Support from an inductive charging infrastructure would make the task a lot easier (and more environmentally friendly).

Similarly...

The quest for full self-driving:

- Largely uncooperative (and not so regulated in some geographies).
- Support from the road infrastructure would make the task a lot easier.
- ODD should also define where and when full self-driving can be used.



The

-
-

Simi



Market and industry are not always rational



BUBBLES: Defining the BUilding Basic BLocks for a U-Space SEparation Management Service

The diagram illustrates the progression of U-space services through four maturity levels, U1 to U4, which are represented by circles of increasing size and color intensity (from light blue to dark blue). A red arrow labeled "Communications capability" points from U1 to U4, indicating that this capability is a key factor in the progression. Below each maturity level is a list of services, with some items circled in red in the original image.

U1	U2	U3	U4
U-Space foundation services	U-Space initial Services	U-Space enhanced Services	U-Space full services
<ul style="list-style-type: none">• e-registration• e-identification• geofencing	<ul style="list-style-type: none">• flight planning• flight approval• tracking• airspace dynamic information• procedural interface with ATC	<ul style="list-style-type: none">• capacity management• assistance for conflict detection	<ul style="list-style-type: none">• integrated interfaces with manned aviation• additional new services

SORA - Specific Operations Risk Assessment
MEDUSA - U-Space Safety Assessment

Full self-driving: is it a classic problem? *(a sort of conclusion)*

- **Building dependable systems using components that are not perfect looks like a classic problem:**
 - ◆ The output accuracy of AI components (in the absence of faults) is probabilistic (specially for black-box AI)
 - ◆ All components are subject to software faults
 - ◆ HW faults must be considered
 - ◆ AI used in safety-critical applications is an interdisciplinary problem, no matter the application area (automotive, medical devices, industry 4.0, avionics, etc.).

Who is going to solve the problems?

(What is expected from our research community?)

Extra slides

Cartography of our World

(Portuguese view circa 1500)



Cartography of (our) Dependability World

