

# Trustworthy AI for defense : engineering challenges

Patricia Besson  
Thales/cortAixLabs France

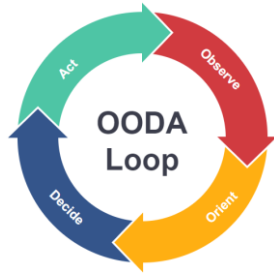
[www.thalesgroup.com](http://www.thalesgroup.com)



# Artificial Intelligence for defense : new capabilities in a constrained environment

## > Artificial intelligence ⇔ artificial capabilities

- ▶ **Acceleration** and **automatisation** of the **OODA loop**  
(Observe, Orient, Decide, Act)



- ▶ **Autonomous** or **semi-autonomous** systems



## > AI in defense needs to operate in an embedded AND safety critical world. Systems should work :

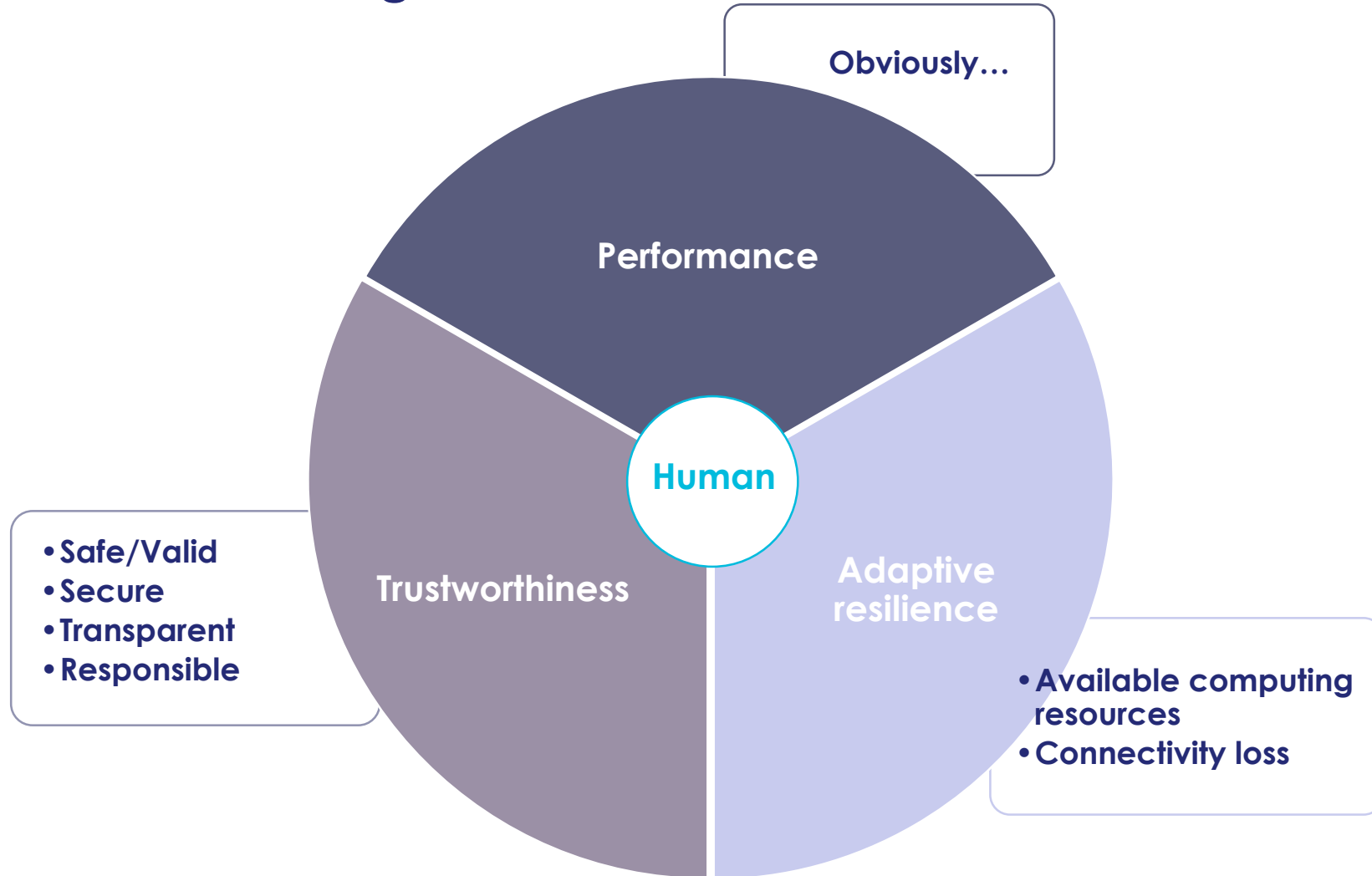
All the time, in the time –  
usual product lifespan 10-20 years

Using **little/specific** data, potentially with **controlled access** –  
Sovereign & classified data requiring specific access rights

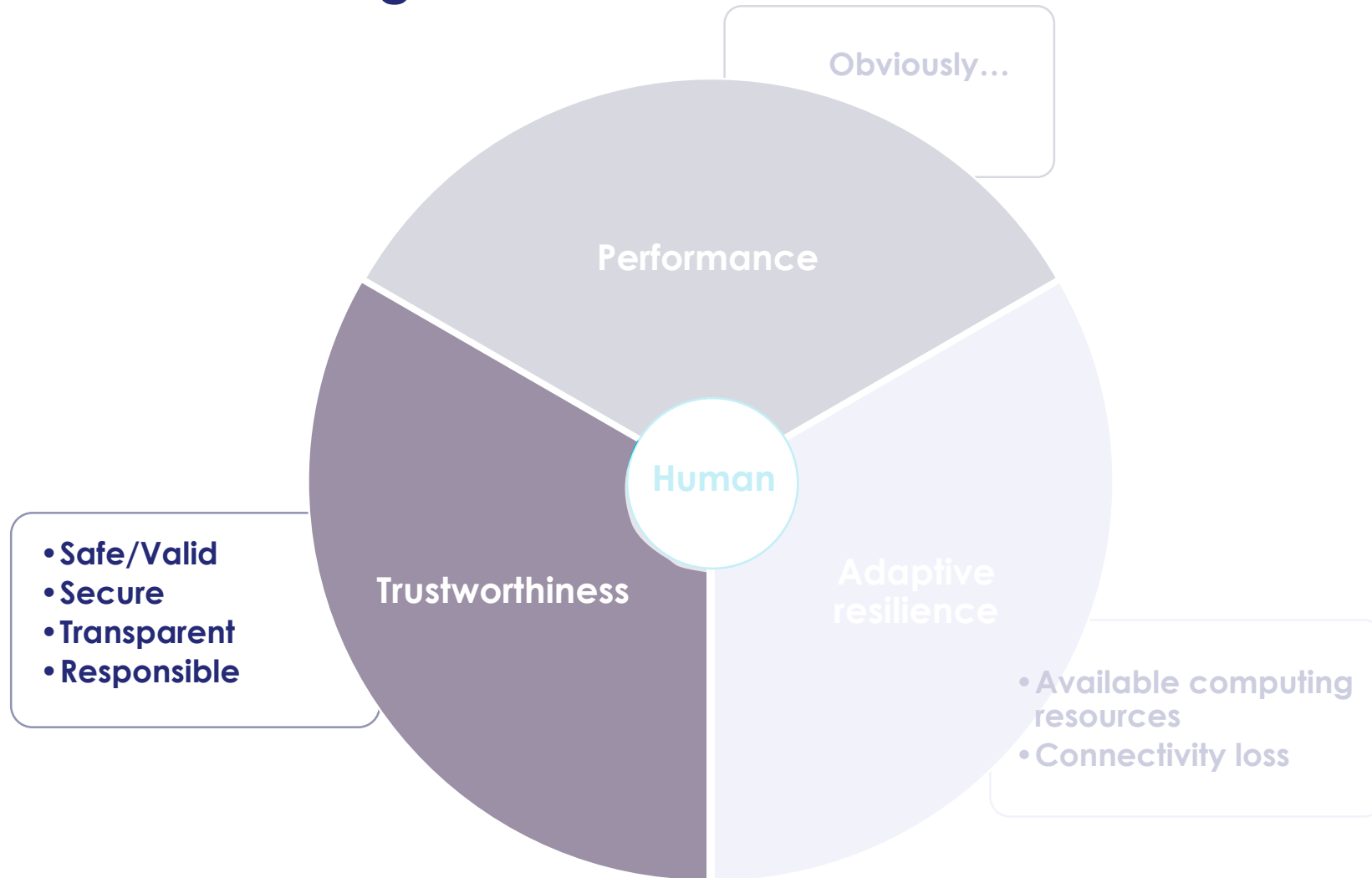
In **changing conditions** –  
including **extreme/rare environment**: °C, radiations, vibrations...

Constrained in **size, weight, power and cost** –  
Field & Edge deployment

# AI in Products : Technological Levers for Critical Defense and Security Systems



# AI in Products : Technological Levers for Critical Defense and Security Systems



# Trustworthy AI

**VALID**

**Doing all and only  
what it is meant to do**

**SECURE**

**Resilient and robust vs  
adversarial conditions**

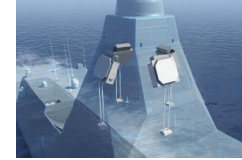
**TRANSPARENT**

**Explainable,  
understable, providing  
context justifications**

**RESPONSIBLE**

**Compliant with  
regulation, legal,  
ethical frameworks**

# Trustworthy AI : it all starts with the system



SYSTEM

IVVQ viewpoint

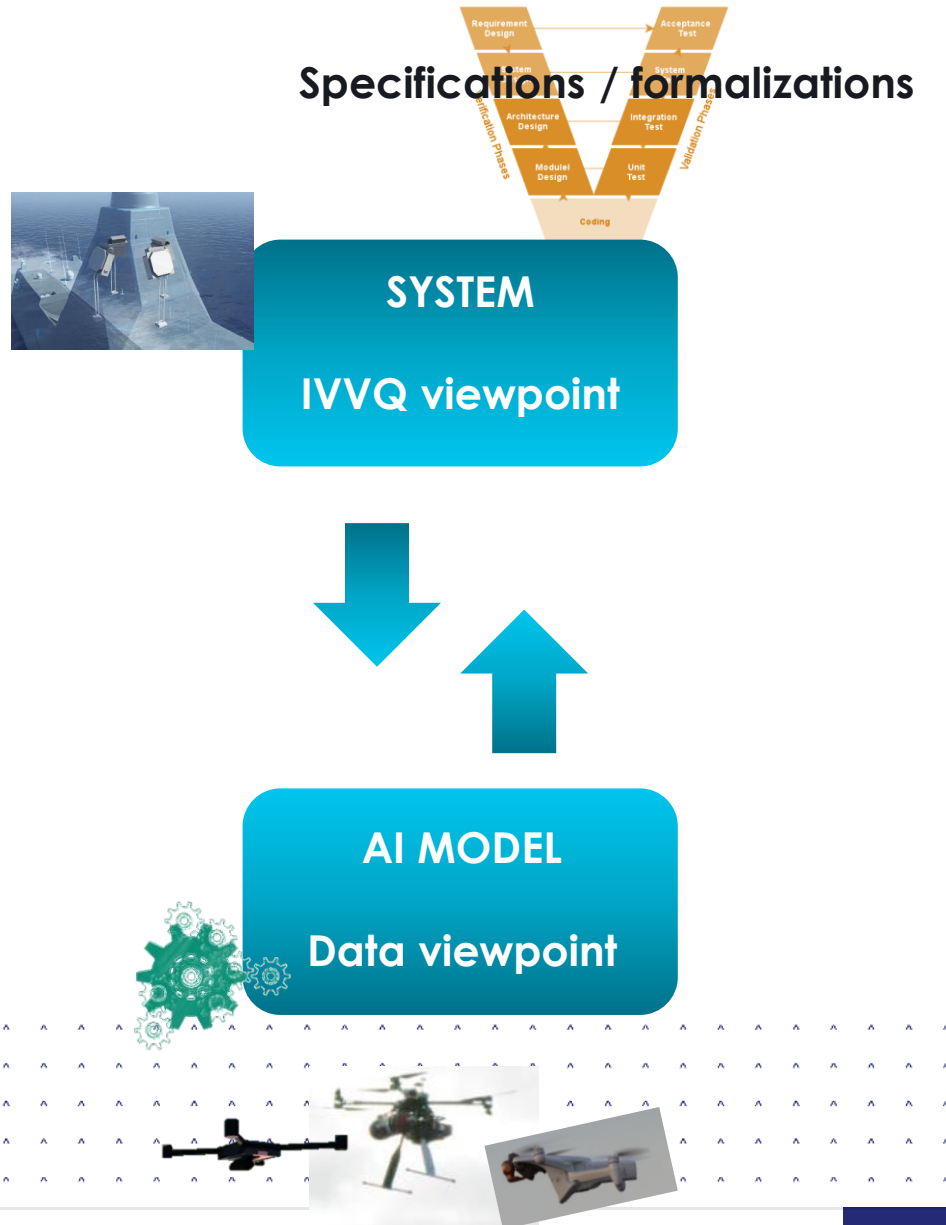
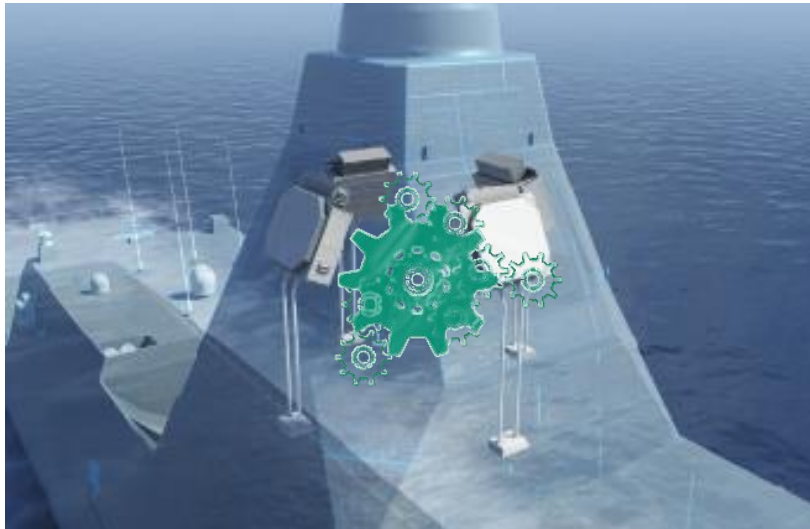


AI MODEL


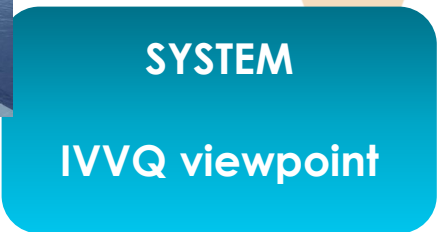
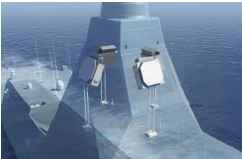
Data viewpoint



# Trustworthy AI : it all starts with the system









# Evaluating and testing AI models

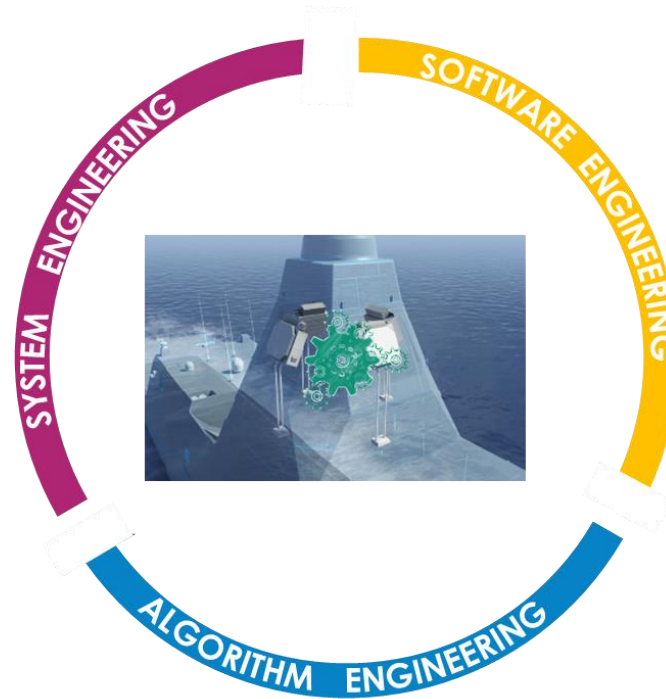
> **Valid model** ⇔ **doing all and only what it is meant to do**

> **Approaches to evaluating and testing AI models :**

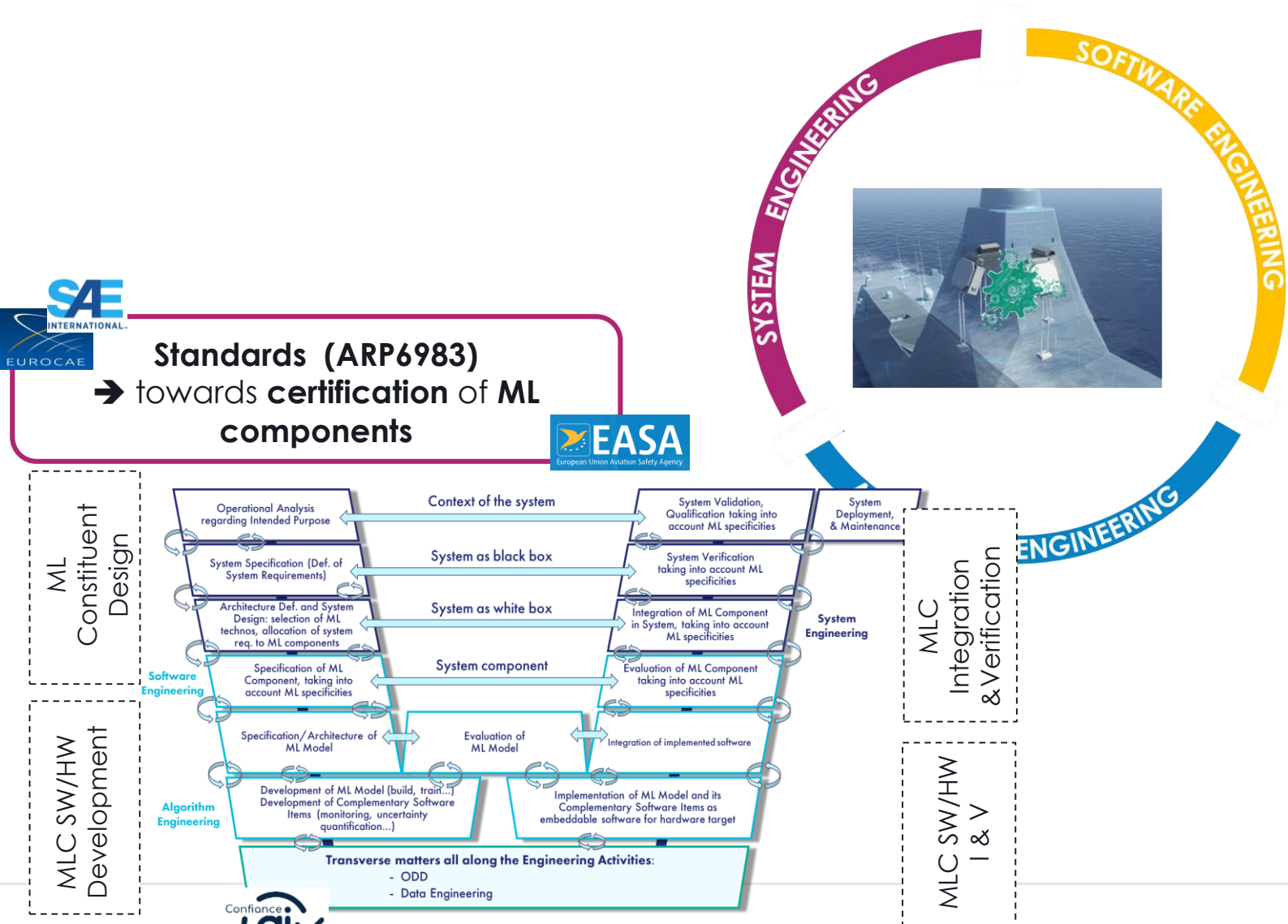
- Against a **pre-defined evaluation data base** → Constituted how and by whom? Representativity?
- **By usage** : in-situ user feedback → only possible for some systems (decision-aiding systems)
- Typical **system engineering approach** : testing and documenting explicit properties and behaviours of the model → mandatory for certification

**No universally accepted testing standards to date → inconsistencies in evaluation methods and results**

# A triple-engineering approach



# Reconciling the system and the model viewpoints



# Reconciling the system and the model viewpoints

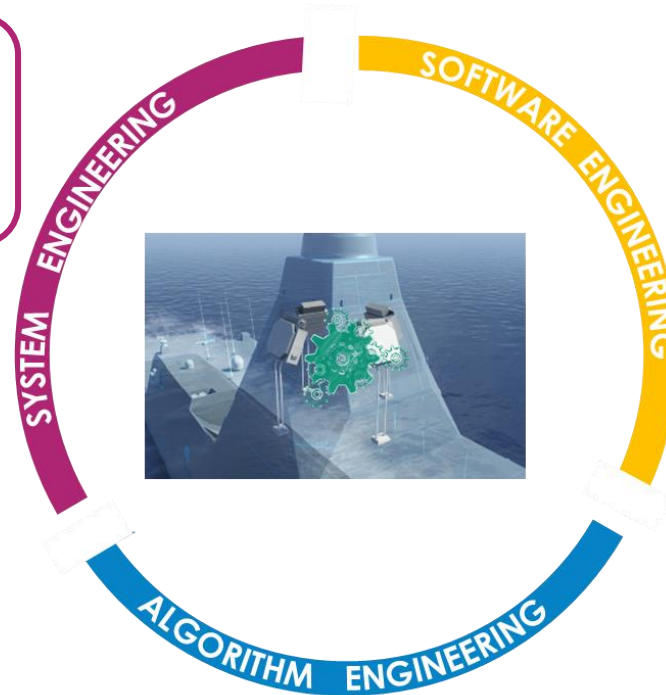
Reconciling **Data** and **Functional Intent**

- Operational Design Domain (**ODD**)
- Rigorous methodologies and tools (problem formalization, symbolic AI,...)



**Standards (ARP6983)**

- towards **certification** of **ML components**



Operational Design Domain (**ODD**)

Voluntary **restriction of the Operational Domain** (specific operating conditions ) within which an **AI constituent** within a **given system** is **intended to function**

# Anchor the approach in Software Engineering

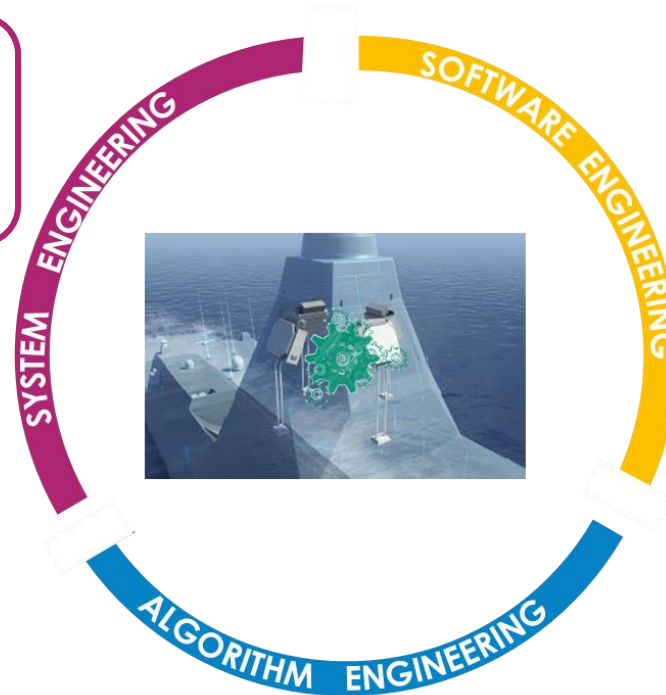
Reconciling **Data** and **Functional Intent**

- ➔ Operational Design Domain (**ODD**)
- ➔ Rigorous methodologies and tools (problem formalization, symbolic AI,...)



**Standards (ARP6983)**

- ➔ towards **certification** of **ML components**



- Data & Knowledge **design guidelines**
- Continuous Integration
- ➔ **ModelOps**

**Formal verification/  
proof of code**

# Tackle the Algorithm Engineering deadlocks

## Reconciling **Data** and **Functional Intent**

- ➔ Operational Design Domain (**ODD**)
- ➔ Rigorous methodologies and tools (problem formalization, symbolic AI,...)

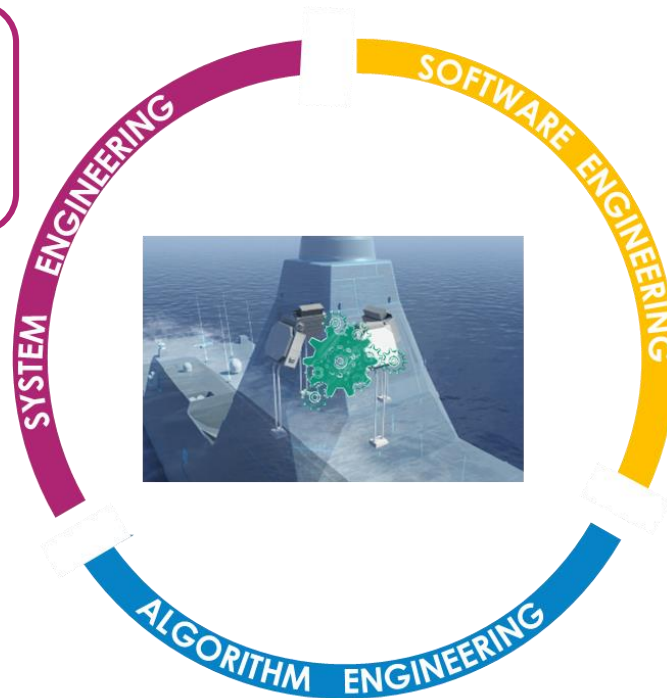


## Standards (ARP6983)

- ➔ towards **certification** of **ML** components

## Monitoring

in **operational environment**  
(abnormality detection, XAI)



- Data & Knowledge **design guidelines**
- Continuous Integration
- ➔ **ModelOps**

**Formal verification/  
proof of code**

## Metric Definition

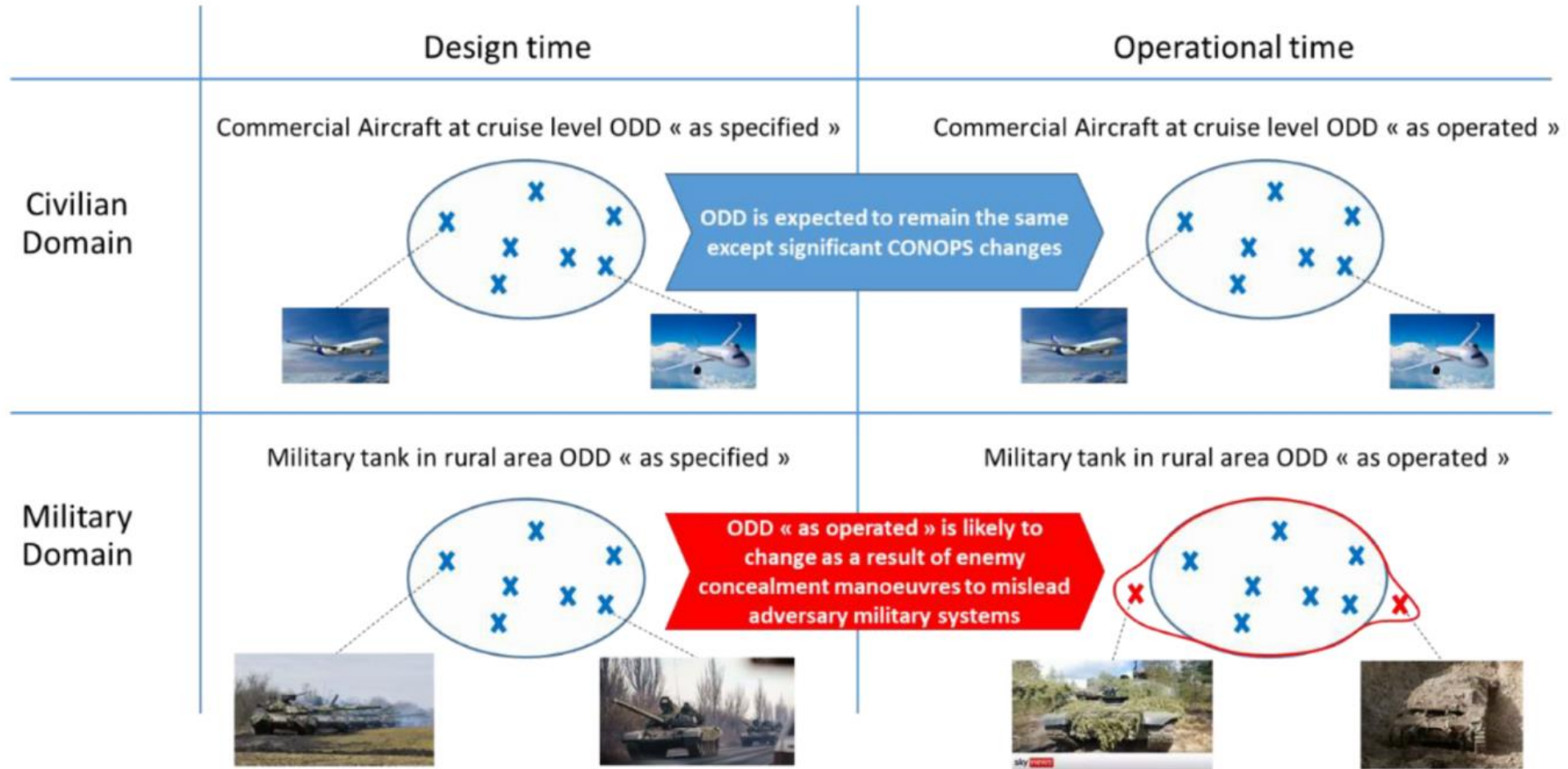
to **assess** and **monitor** algorithm **fit-for-purpose**

## Evaluation & tests :

- **Formal** methods
- **Axiomatic** proofs (robust by-design models)
- **Experimental** approach (Acceptable for IVVQ as long as **analyses are complete, traceable and quantitative**)
- Explainable AI (**XAI**)
- **Cybersecurity** tests



# Adaptive ODD: Incremental Learning & Qualification



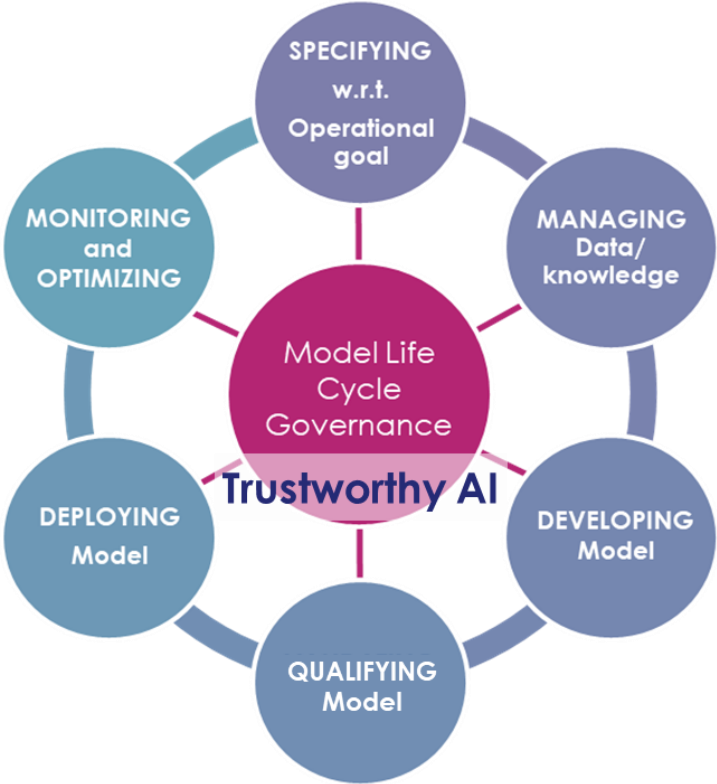
Differences between civilian and military ODDs [From EDA TAID white Paper]

# End-to-End Trustworthy AI Engineering Lifecycle

*Critical needs in defense*

At run time

Runtime Assurance  
(adaptive ODD)



At design time

Development assurance

# Conclusion : Trustworthy AI is at hand

**Interdisciplinary approach is a key to success:** system engineering, safety, security, software...

- IA can't be trustworthy per se
- The goal is to quantify and manage the risk

## Some resources:

- ❖ Take a look at the Confiance.AI project outputs

<https://www.confiance.ai/>



- ❖ We will release soon an **open source soft** for benchmarking neural networks based on **confidence scores** on the ThalesGroup GitHub

<https://github.com/ThalesGroup>

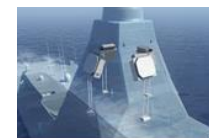


- ❖ **EUROCAE/SAE ARP6983** to be released in early **2026**



- ❖ Take a look at the freshly published **EDA white paper** on **Trustworthy AI for Defense**

<https://eda.europa.eu/docs/default-source/brochures/taid-white-paper-final-09052025.pdf>



**SYSTEM**

**IVVQ viewpoint**



**AI MODEL**

**Data viewpoint**





---

**Patricia Besson**

Thales / cortAlx Labs France

 [patricia.besson@thalesgroup.com](mailto:patricia.besson@thalesgroup.com)