



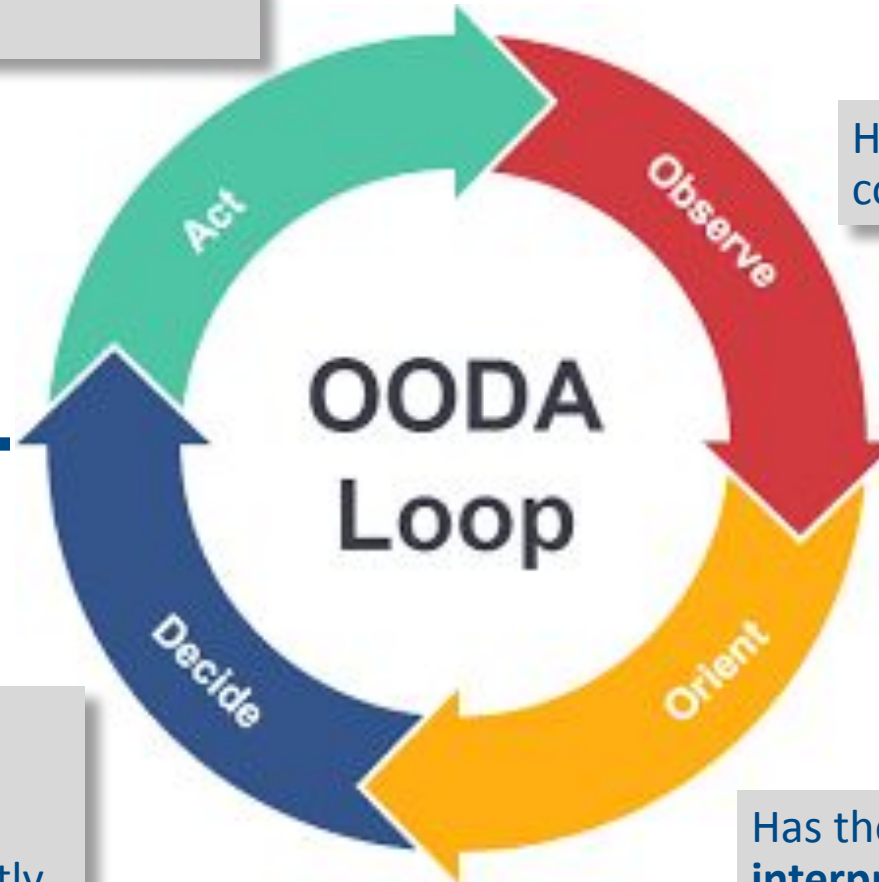
Trustworthy AI: An operator's view...

- 'Black Box' Decision Making – the challenge of assurance
- Calibrated trust – what happens when it goes wrong?
- Trustworthy AI – an enabler

Decision Making in a System of Systems

Based on a calibrated and appropriately supervised system, is the **Action** correct?

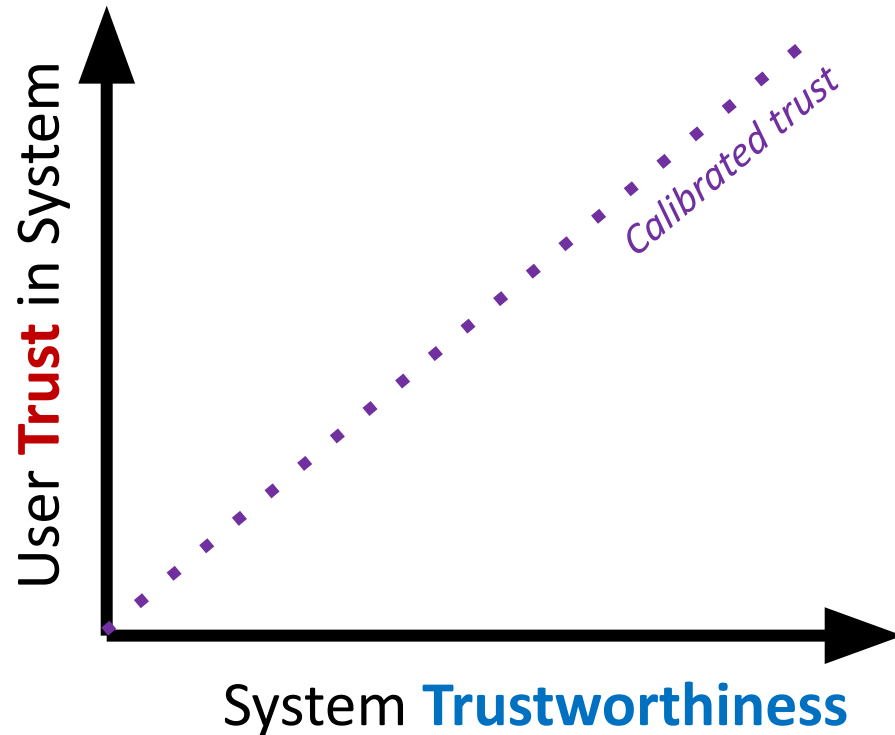
Has the system **observed** the correct things?



Has the system made the correct **interpretation** of the situation based on what it has observed?

Decision based on *trust* in Observe and Orient phases

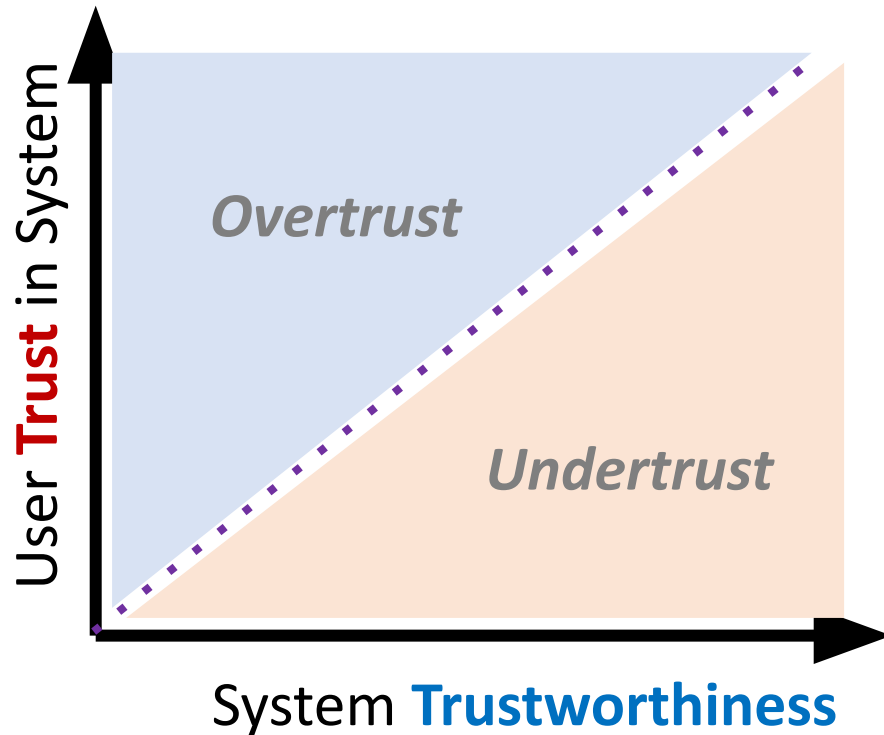
Level of trust has to be *calibrated* correctly – and *supervised* accordingly



Trust = response of a user in a situation of uncertainty or vulnerability. *Subjective*

Trustworthiness = measure of trust qualities in a system (autonomous or AI). *Objective*

User **Trust** must be commensurate with the **Trustworthiness** of the system (well calibrated)



Overtrust. Trust in the system is greater than the system can deliver:

- Over-reliance on AI/automation
- Taking inappropriate or misguided action

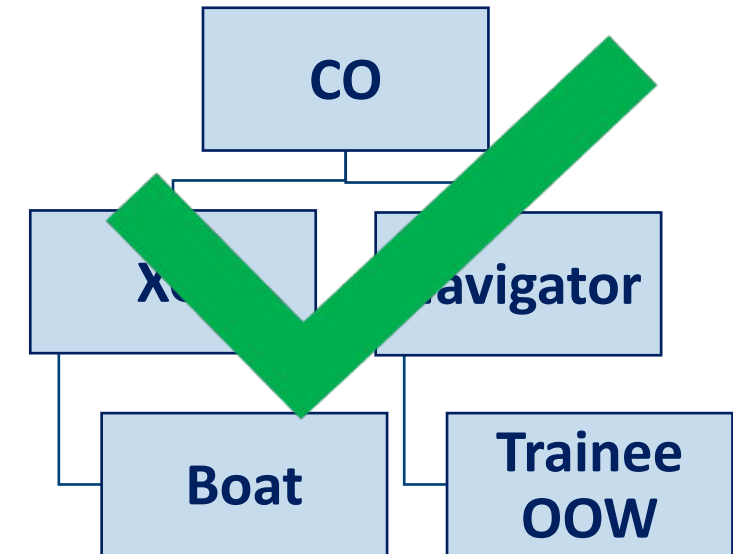
Undertrust. System performs better than supervisor allows for:

- Supervisor 'knows better'
- Taking alternative, contrary or abortive action

Human example – HMS DIAMOND, Akrotiri Bay 2018



Trust/Trustworthiness *System of Systems*



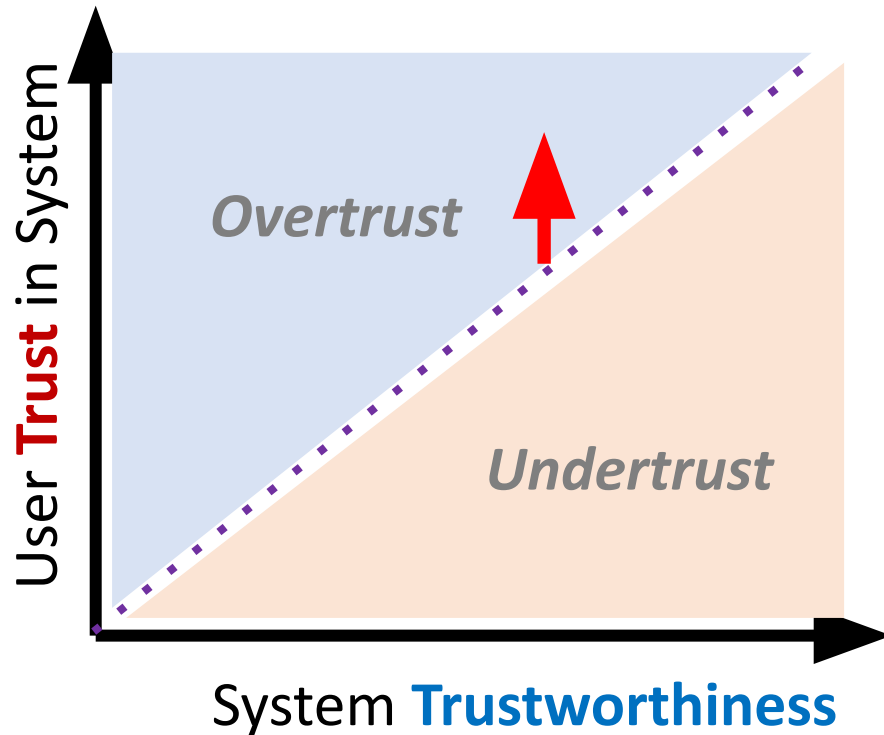
Human example – HMS DIAMOND, Akrotiri Bay 2018



Image courtesy of Google Earth 2023

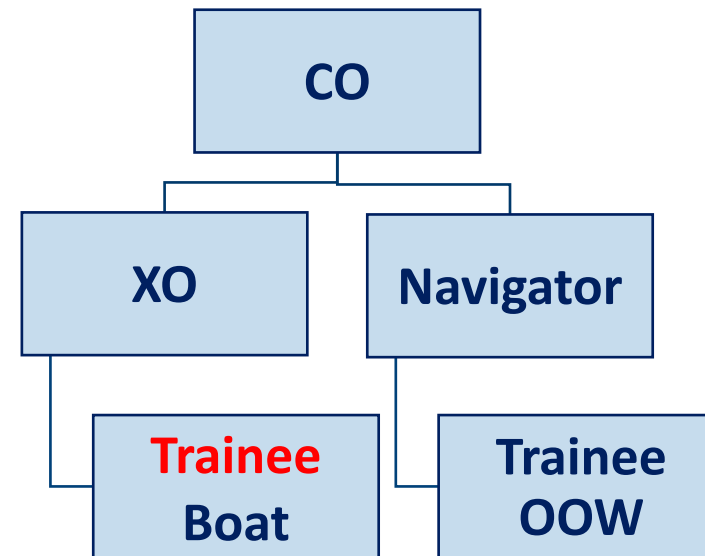
Trust/Trustworthiness *System of Systems*

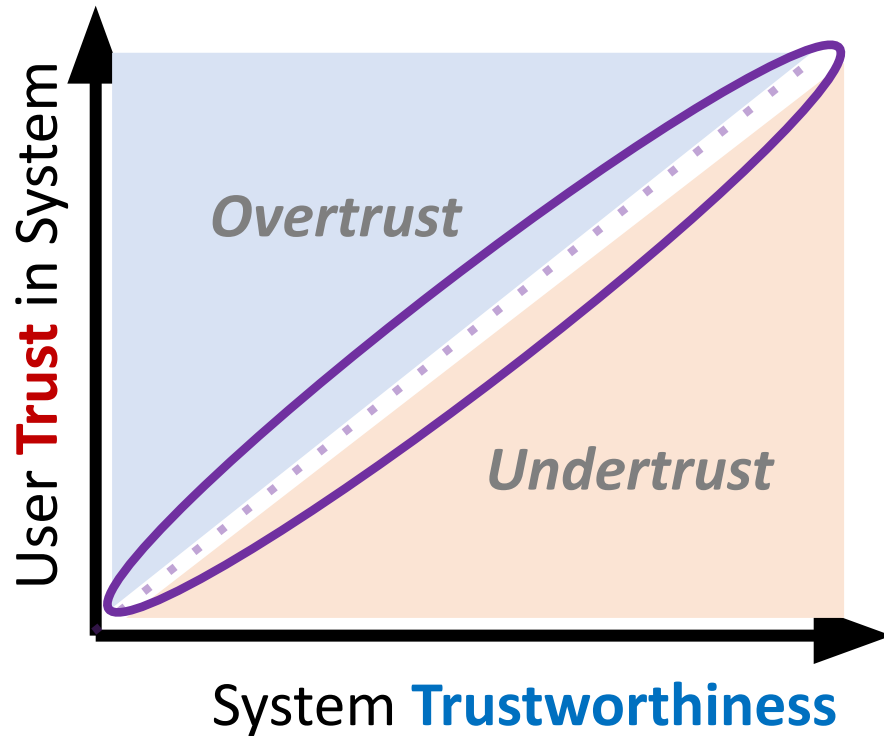




Moved from calibrated to miss-calibrated trust = **overtrust**

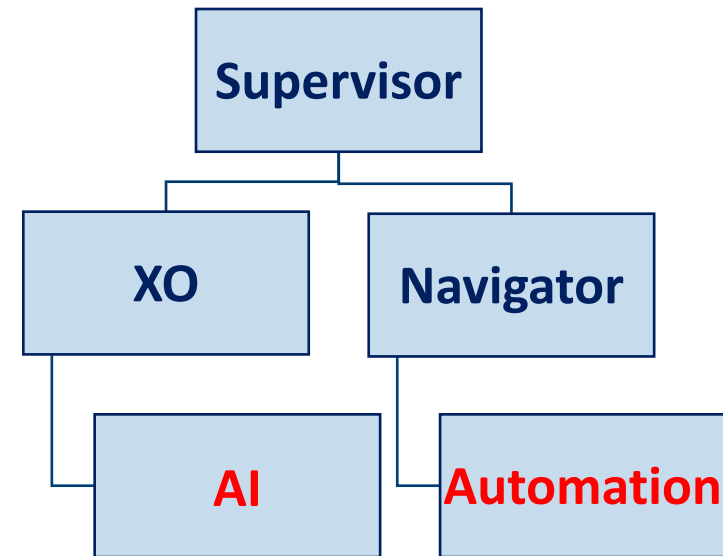
Human System: *Calibrated trust* mix of intuition, experience, qualification, external validation



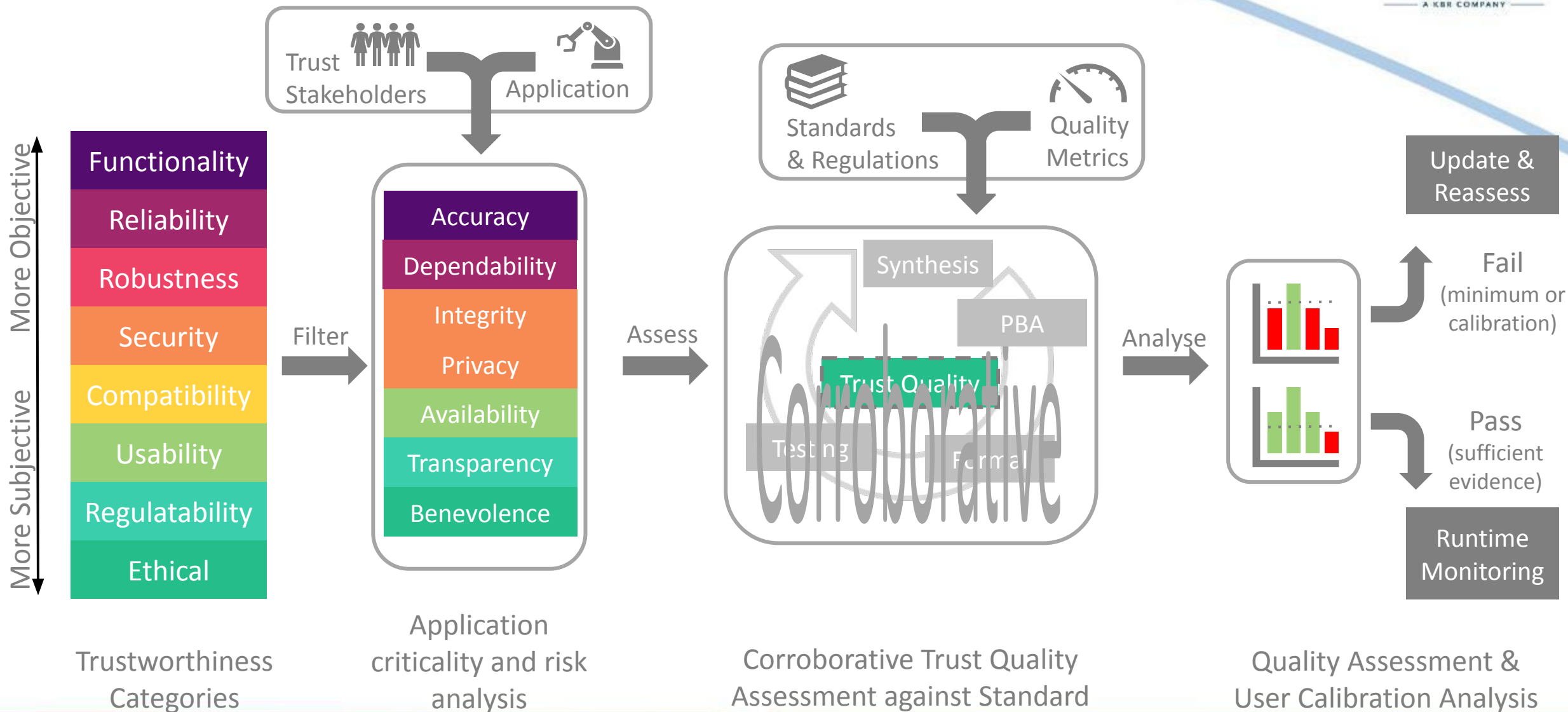


AI Applicable lessons:

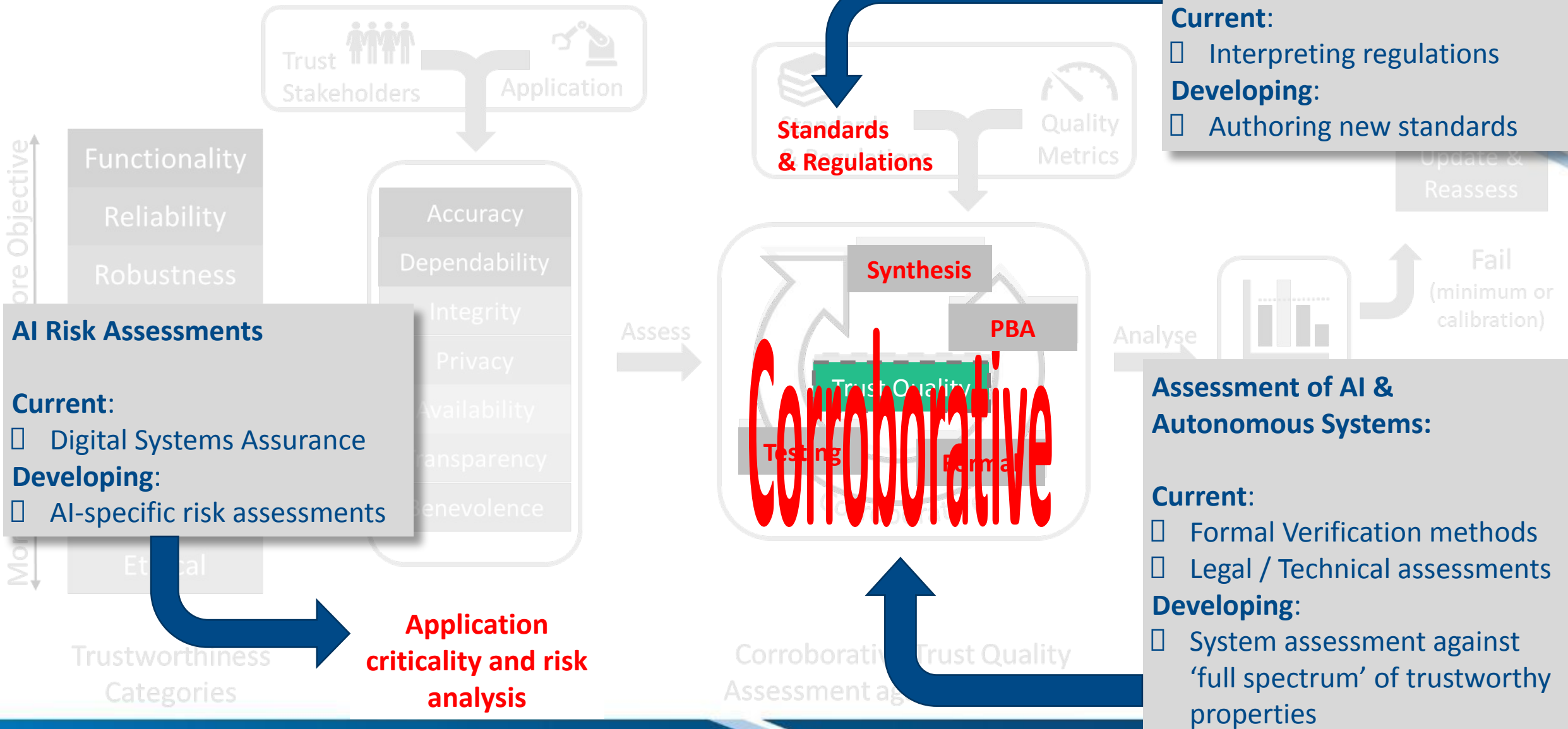
- Increase automation in the system – *calibrated trust* vital
- **But** harder to define the line = area
- Need to measure *System Trustworthiness*



Assuring AI / Autonomy: Frazer-Nash Research and Development



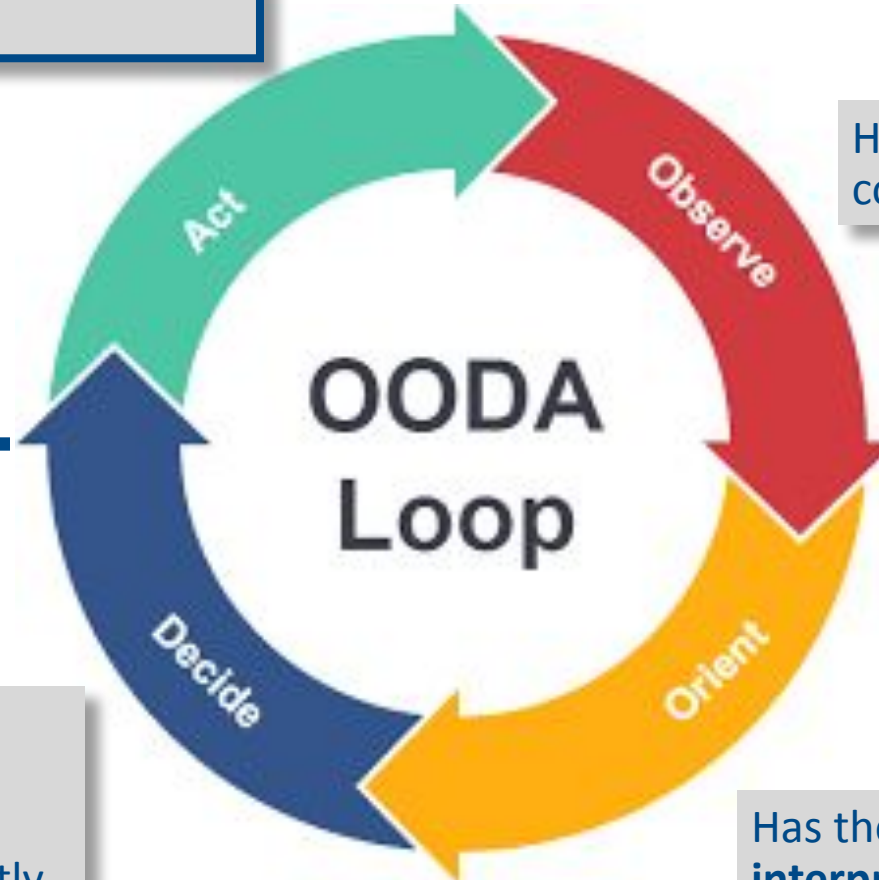
Trustworthiness as an 'AI Enabler'



Decision Making in a System of Systems

Based on a calibrated and appropriately supervised system, is the **Action** correct?

Has the system **observed** the correct things?



Has the system made the correct **interpretation** of the situation based on what it has observed?

Decision based on *trust* in Observe and Orient phases

Level of trust has to be *calibrated* correctly – and *supervised* accordingly



Thank you

Ben Keith
b.keith@fnc.co.uk