

ALMA MATER STUDIORUM  
UNIVERSITÀ DI BOLOGNA



# Risk and Exposure of XAI in Persuasion and Argumentation: The case of Manipulation

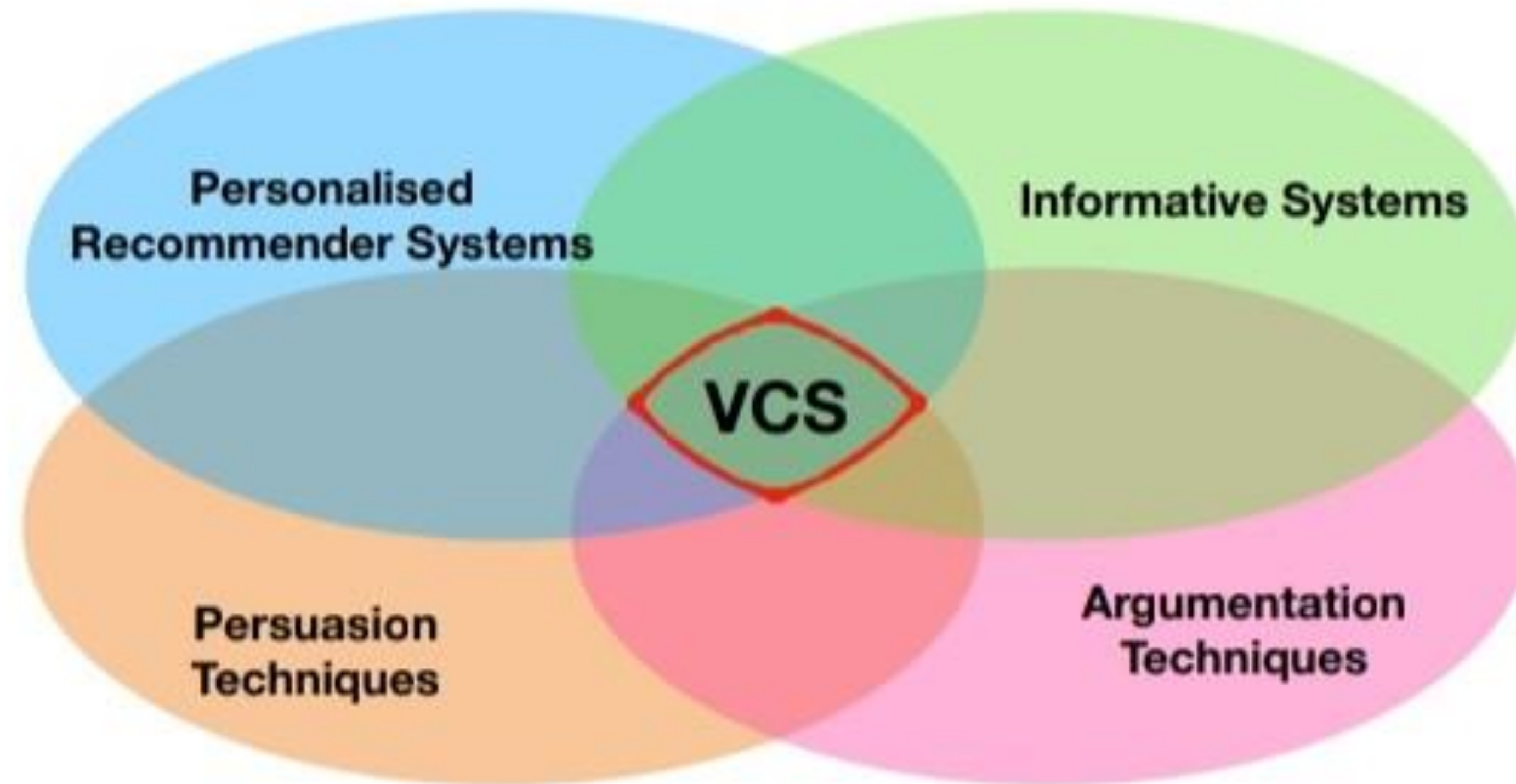
**Rachele Carli<sup>1</sup>, Amro Najjar<sup>2</sup>, Davide Calvaresi<sup>3</sup>**

<sup>1</sup>Alma Mater Research Institute for Human-Centered AI, University of Bologna & ICR Group, University of Luxembourg;

<sup>2</sup>Luxembourg Institute of Science and Technology (LIST), Luxembourg;

<sup>3</sup>University of Applied Sciences Western Switzerland, Switzerland





Transparency

Safety



Autonomy



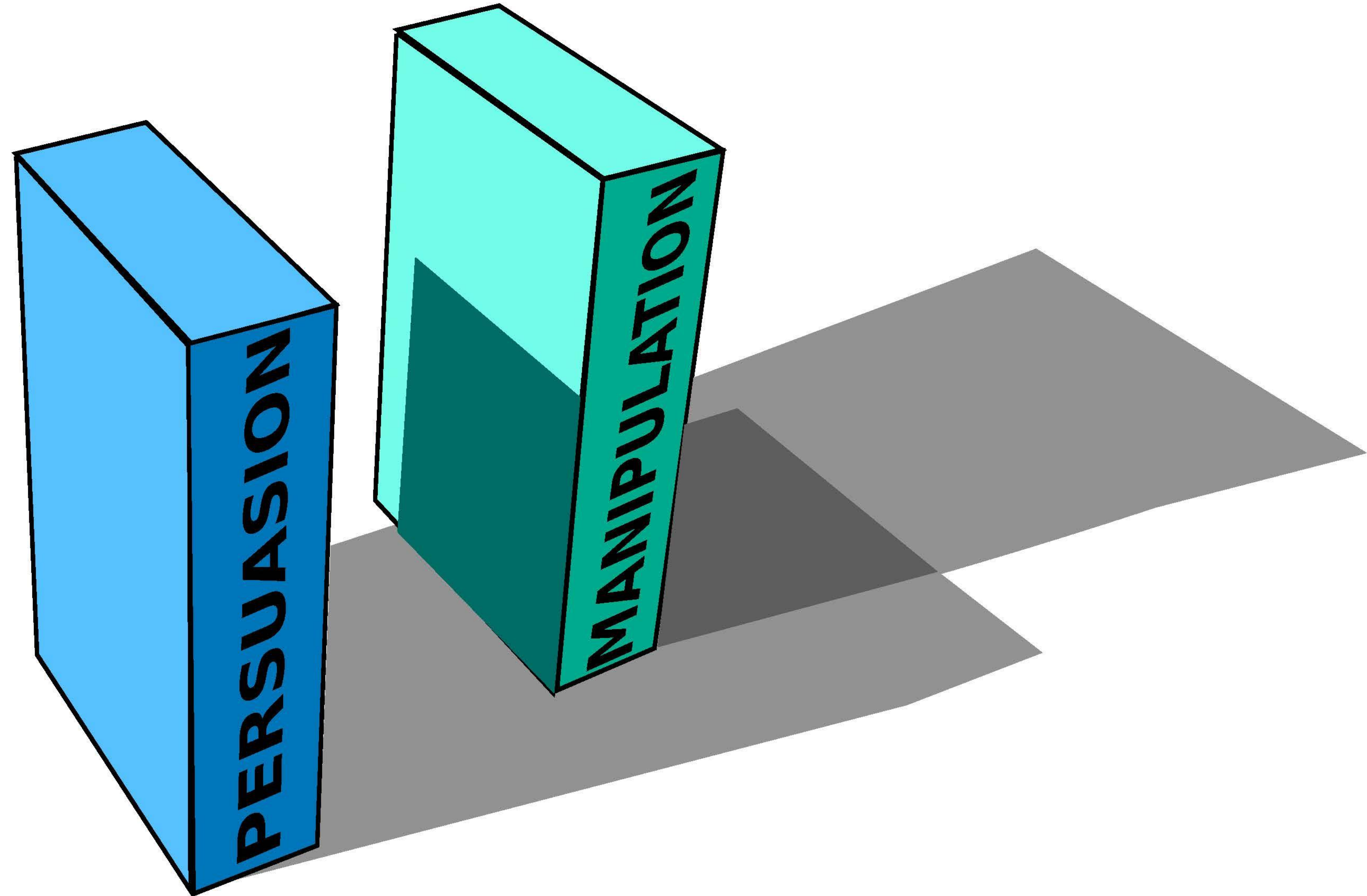
Can AI systems be **actually** transparent for non-experts?



Can the **mere** fact of giving an explanation make the system safer?



Can the explanation make the user **really aware** of system's dynamics?





**No** clear definitions



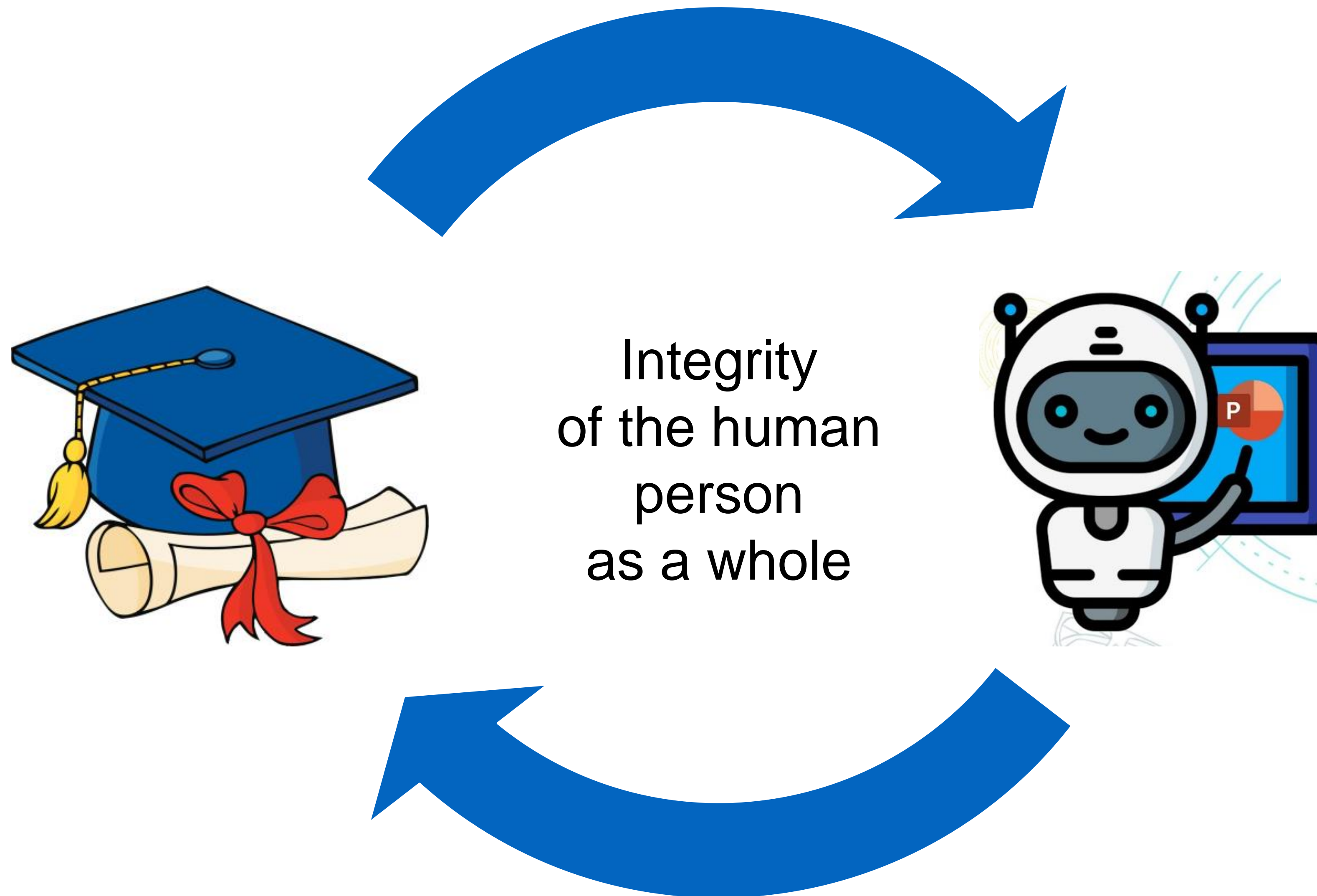
**No** clearly identifiable causal link



**No** clear boundaries



Curbing **Self-governance**



**Adaptivity**

**Granularity of explanations**

**Paternity of choice over  
autonomy**

**Order matters**

**From manipulation to  
manipulative techniques**

Specificity vs General cases:

- Context-sensitive
- Goal-sensitive
- Precise interaction-sensitive



Adaptivity

**Granularity of explanations**

Paternity of choice over  
autonomy

Order matters

From manipulation to  
manipulative techniques

For each recommendation:

- Specific expl.
- Identifiable expl.
- No hidden implications
- Easier traceability

Adaptivity

Granularity of explanations

**Paternity of choice over  
autonomy**

Order matters

From manipulation to  
manipulative techniques

- Track decision-making process
- Conscious reconstruction of decision
- Coherence with users' goal

**Adaptivity**

**Granularity of explanations**

**Paternity of choice over  
autonomy**

**Order matters**

**From manipulation to  
manipulative techniques**

- Relevance
- Highest impact on safety
- FRIA criterion

**Adaptivity**

**Granularity of explanations**

**Paternity of choice over  
autonomy**

**Order matters**

**From manipulation to  
manipulative techniques**

- Correcting the dynamics to correct the effect
- No exploitation of vulnerabilities
- Maximisation of users' utility
- Respect for users' goal/will



[rachele.carli2@unibo.it](mailto:rachele.carli2@unibo.it)

[amro.najjar@list.lu](mailto:amro.najjar@list.lu)

[davide.calvaresi@hevs.ch](mailto:davide.calvaresi@hevs.ch)