

# ON THE **CHALLENGES** RAISED BY ROBOTS POWERED BY **ARTIFICIAL INTELLIGENCE**

Fulvio Mastrogiovanni

Professor of Robotics and Artificial Intelligence, University of Genoa  
Founder and Chief Science Officer, Teseo srl  
Coordinator of the Scientific Committee, Digital Innovation Hub Liguria

$10^{-9}$

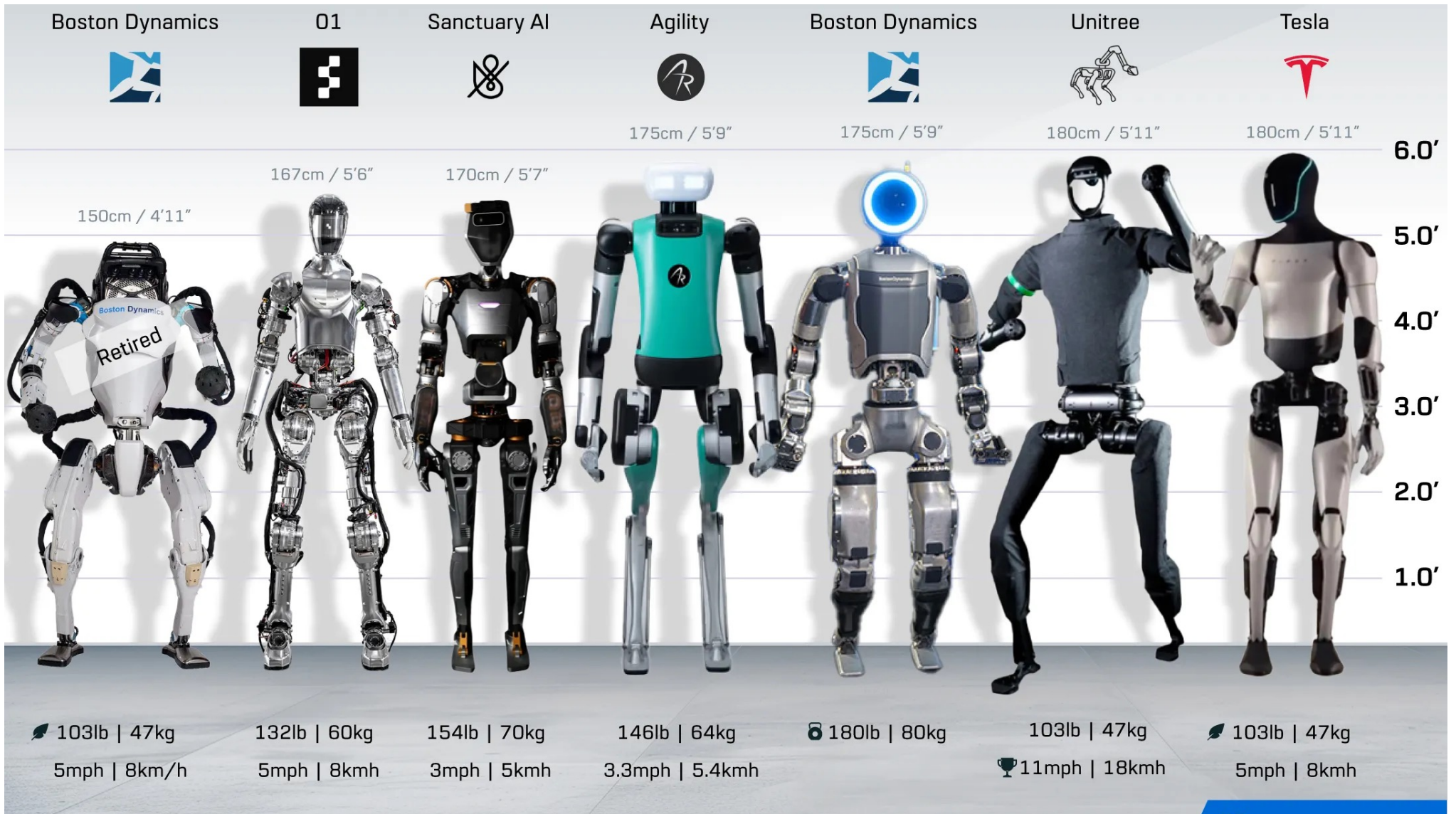


ROBOTS OPERATING

./ IN AUTONOMY

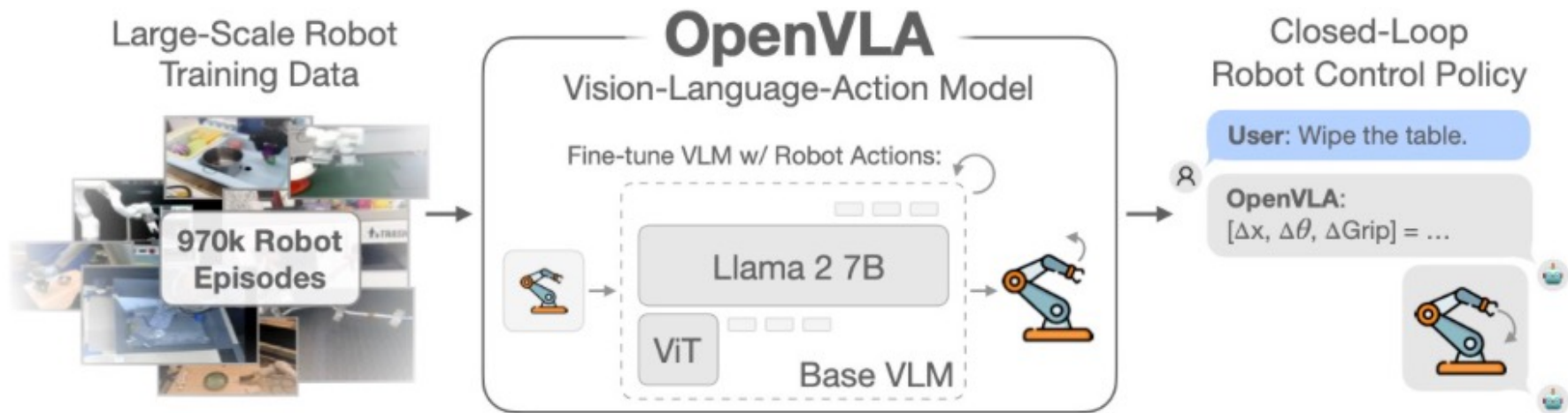
./ JOINTLY WITH HUMANS



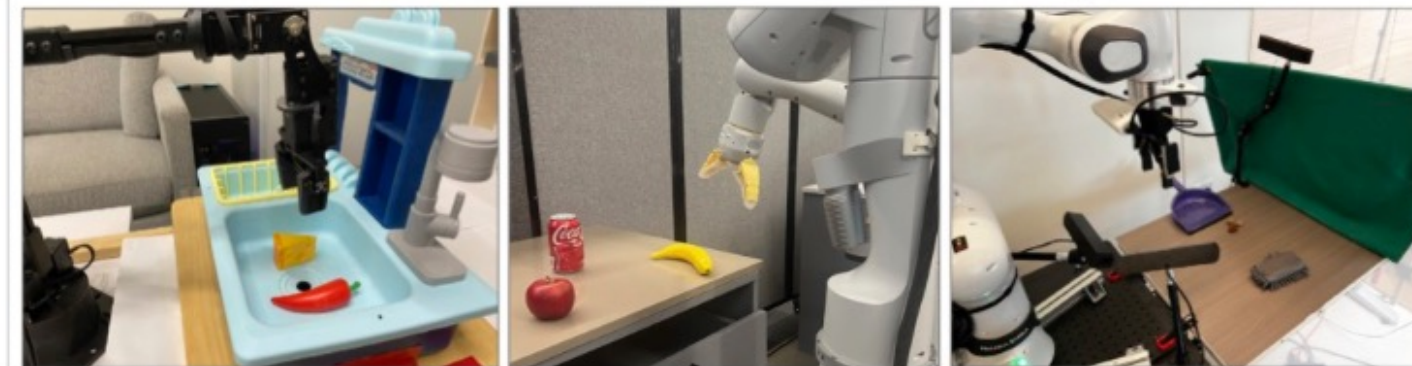


ROBOT INTELLIGENCE = BODY + AI

ROBOT INTELLIGENCE = BODY + VLA






**Multi-Robot Control & Efficient Fine-Tuning**



**Fully**

**Open-Source**

-  Data
-  Weights
-  Code



Put Eggplant into Pot



Put Yellow Corn on Pink Plate



0.16

ENGINEERING



$10^{-9}$

0.16



LEARNING

ROBOT INTELLIGENCE = F(BODY, IA, ...)

ROBOT INTELLIGENCE = F(BODY, IA, ...)

A. NEWELL (1990)





# COGNITIVE ARCHITECTURES FOR ROBOTS

## WHAT A ROBOT COGNITIVE ARCHITECTURE SHOULD EXHIBIT/

1/ REACTIVE EXECUTION.

2/ KNOWLEDGE REPRESENTATION AND INFERENCE.

3/ REASONING.

4/ HIERARCHICAL AND RECURSIVE REASONING.

5/ LEARNING FROM EXPERIENCE.

6/ MEMORY MODELS: SEMANTIC MEMORY, EPISODIC MEMORY, MENTAL IMAGERY.

7/ INTERACTION MODELS.

8/ SOCIAL MODELS OF INTERACTION.

...

## KEY QUESTIONS/

**Q1/** ARE CURRENT PRINCIPLES IN COMPUTATION **COMPATIBLE** WITH COGNITION IN ROBOTS?

**Q2/** WHAT ARE THE (MAYBE TRIVIAL BUT HIDDEN) **ASSUMPTIONS** IN INTEGRATING AI TECHNIQUES IN COGNITIVE ARCHITECTURES FOR ROBOTS?

**Q3/** IS ROBOT INTELLIGENCE A PROPERTY **ASSOCIATED WITH** THE ALGORITHMS, OR IS IT CONNECTED TO THE COGNITIVE ARCHITECTURE **AS A WHOLE**?

**Q4/** IS THERE A **MORE FUNDAMENTAL DESCRIPTION LEVEL** FOR ROBOT BEHAVIOURS?

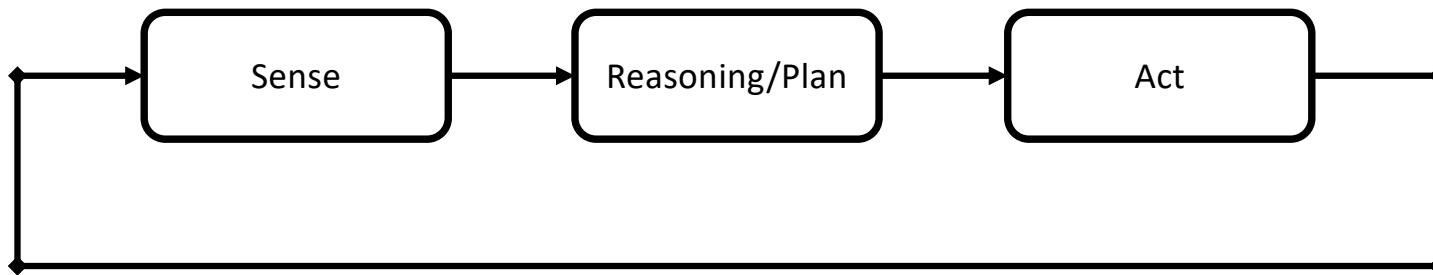
## THREE+ ARCHITECTURAL FAMILIES/

1/ SENSE-PLAN-ACT.

2/ REACTIVE.

2+/ BEHAVIOUR-BASED.

3/ HYBRID REACTIVE-DELIBERATIVE (COGNITIVE?).

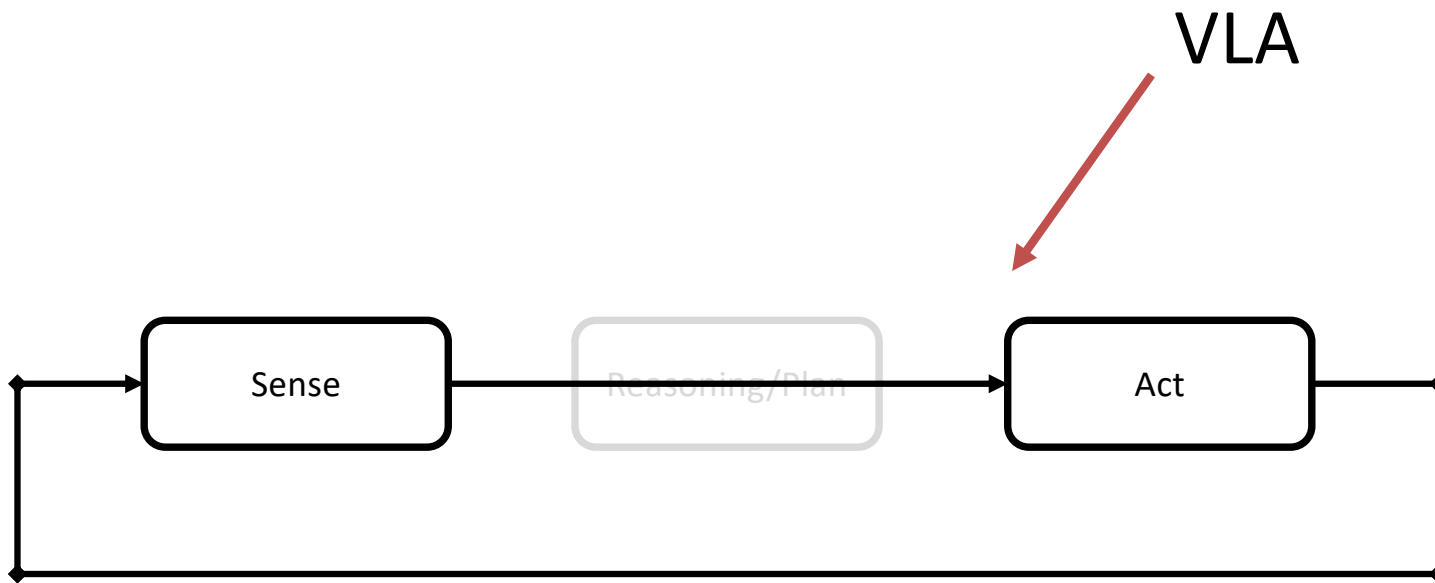


SENSE-PLAN-ACT

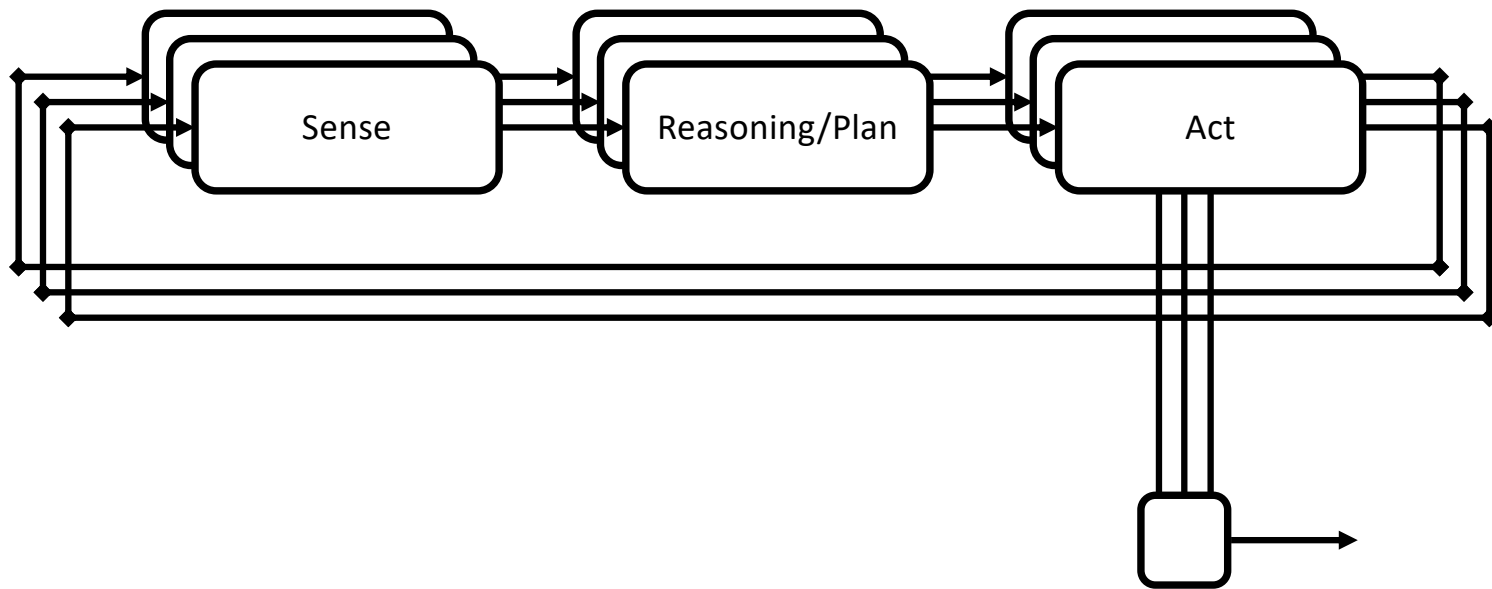




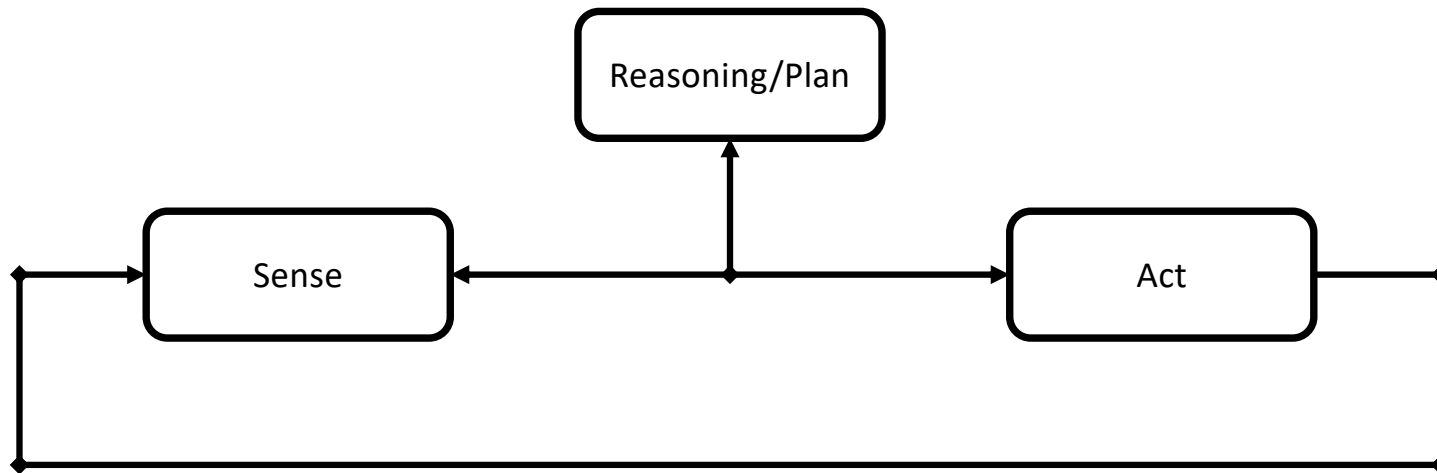
REACTIVE



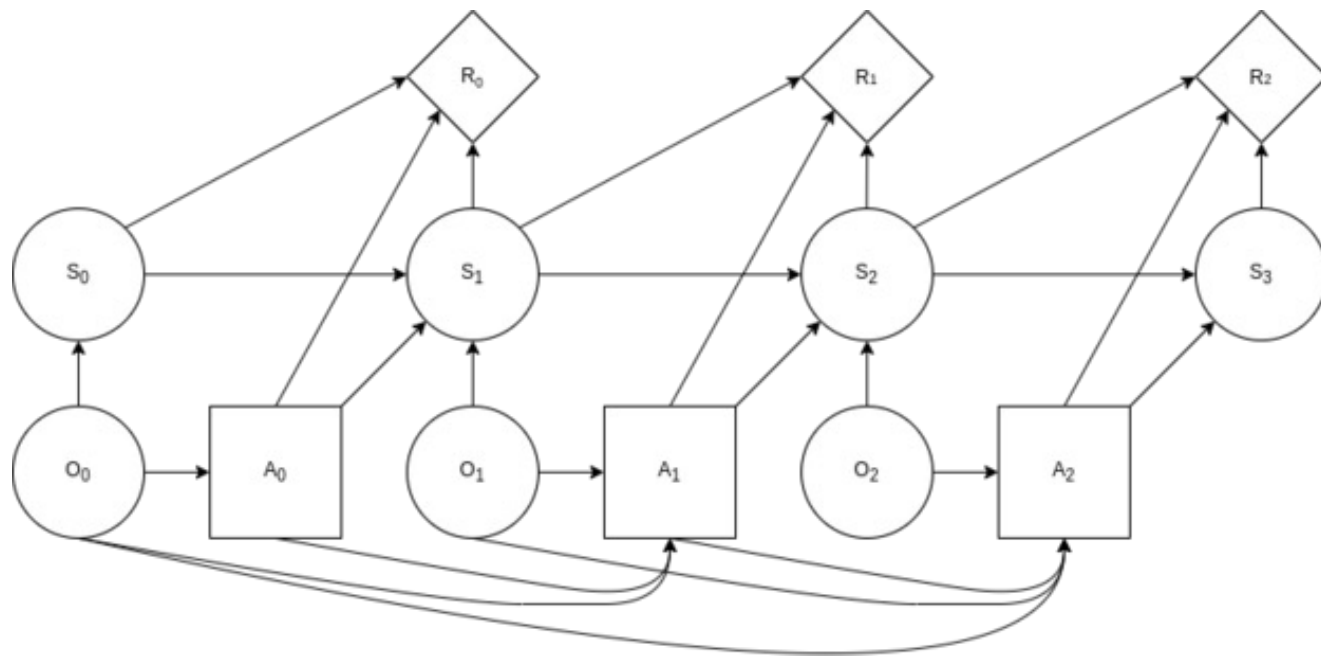
REACTIVE



BEHAVIOUR-BASED

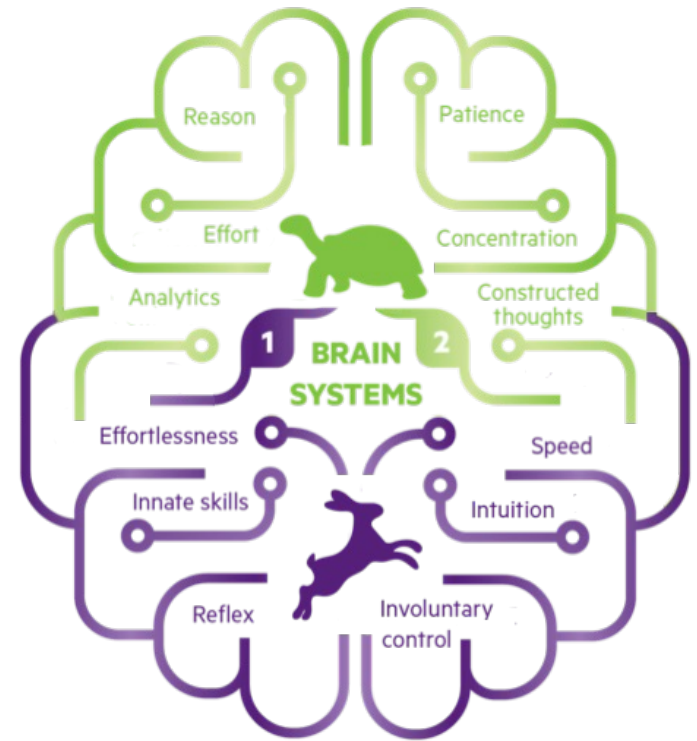
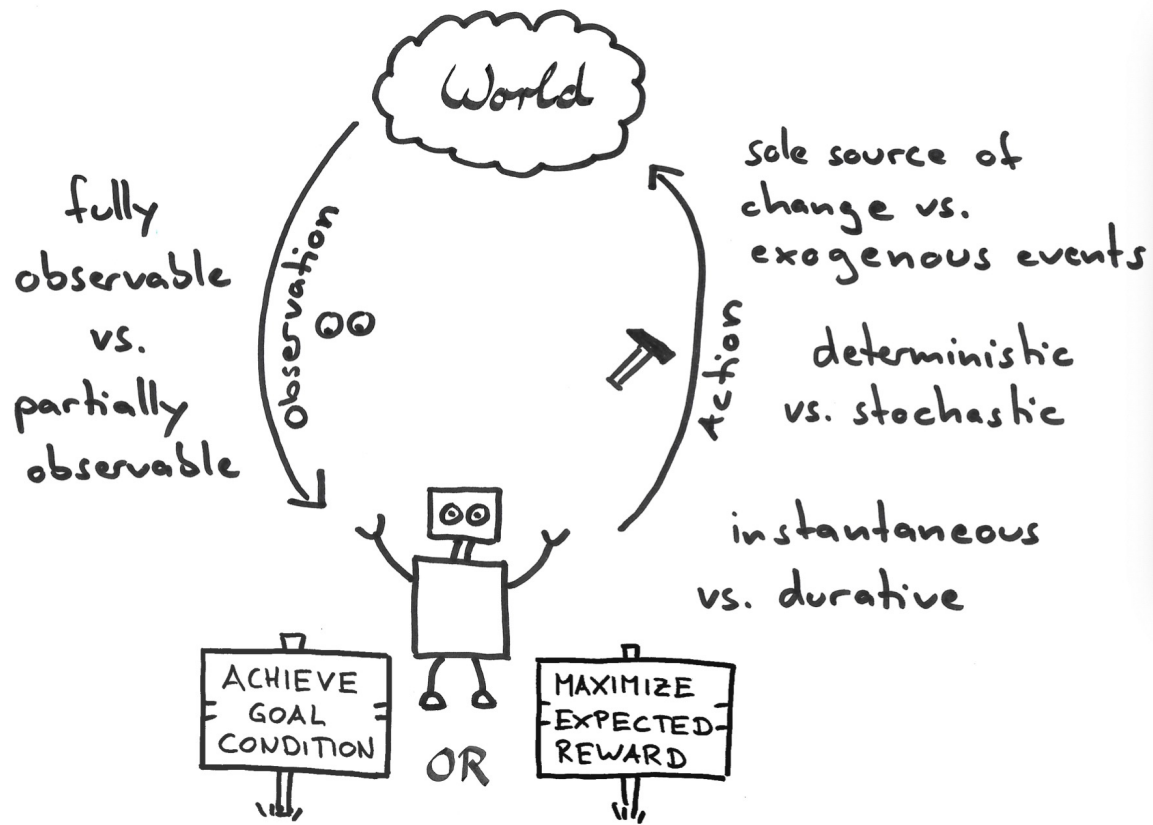


HYBRID REACTIVE-DELIBERATIVE

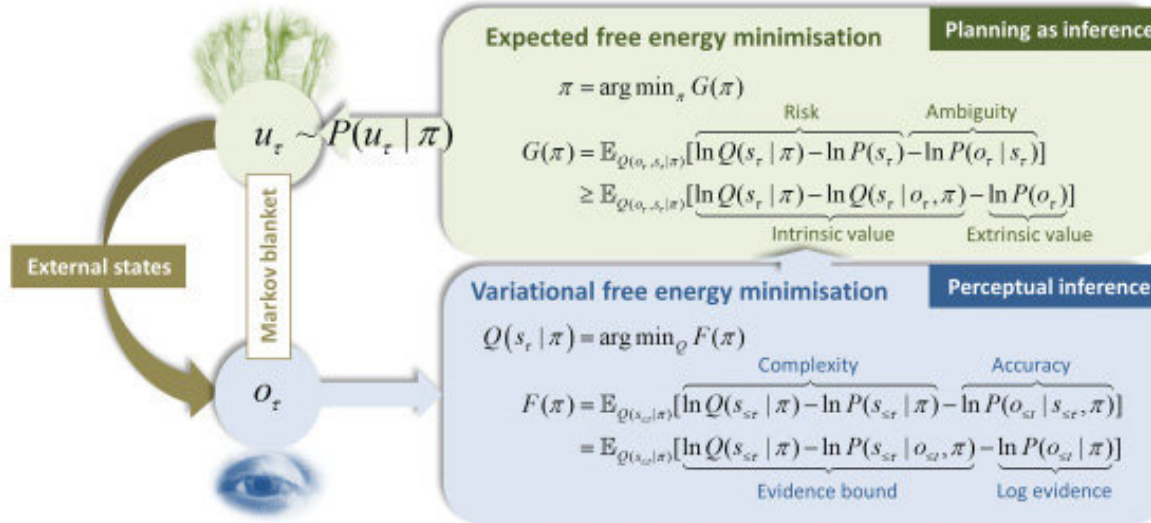




$$b'(s') = \eta O(o | s', a) \sum_{s \in \mathcal{S}} T(s' | s, a) b(s)$$



## Active inference



**No priors**

$$\mathbb{E}_{\phi} [D[Q(s_\tau | o_\tau, \pi) || Q(s_\tau | \pi)]]$$

$$=$$

$$D_{KL}[Q(s_\tau, o_\tau | \pi) || Q(s_\tau | \pi)Q(o_\tau | \pi)]$$

Bayesian surprise  
Optimal Bayesian design  
Intrinsic motivation  
Infomax principle



**No ambiguity**

$$D_{KL}[Q(s_\tau | \pi) || P(s_\tau)]$$

$$=$$

$$D_{KL}[Q(o_\tau | \pi) || P(o_\tau)]$$

Risk-sensitive policies  
KL control  
Occam's principle



**No intrinsic value**

$$\mathbb{E}_{\phi} [\ln P(s_\tau)]$$

$$=$$

$$\mathbb{E}_{\phi} [\ln P(o_\tau)]$$

Bayesian decision theory  
Expected utility theory



**No ambiguity or priors**

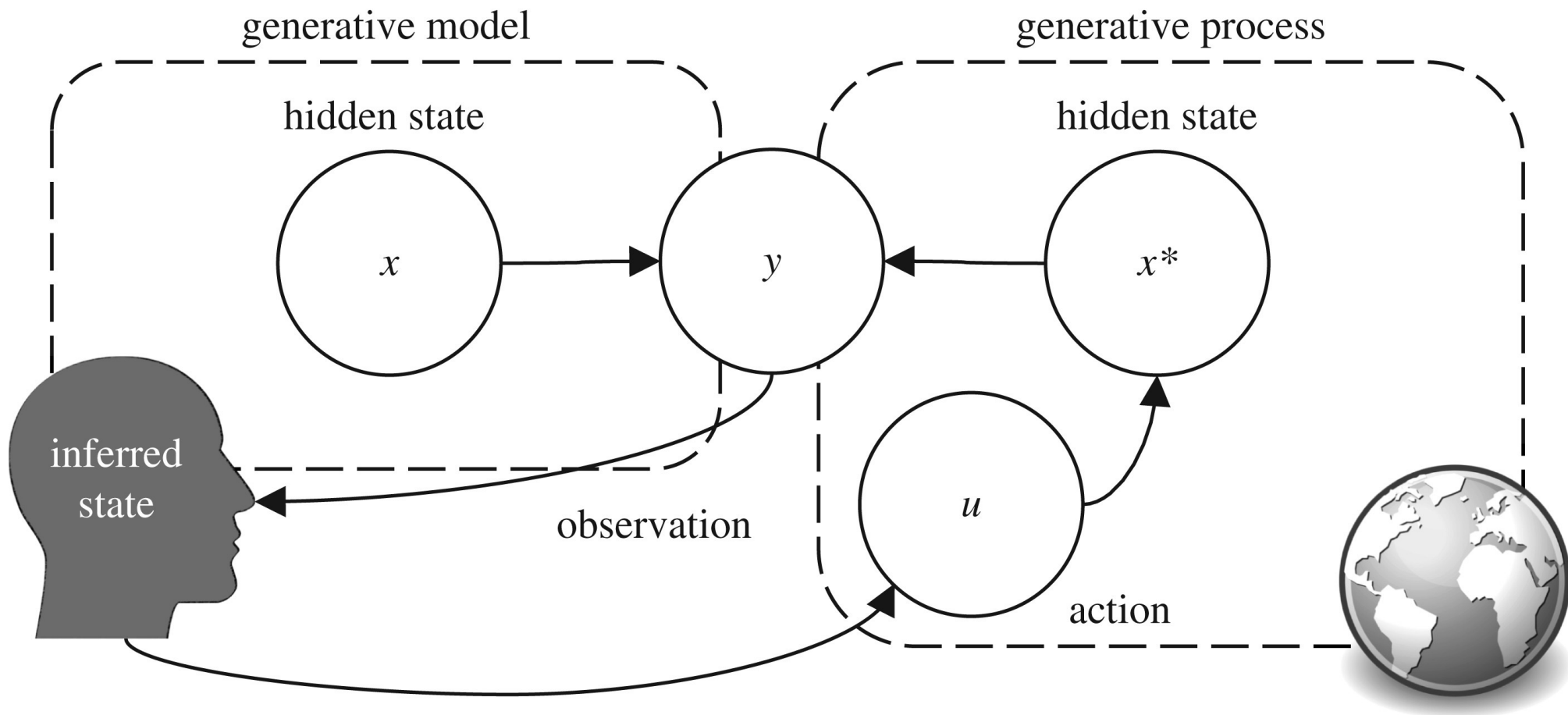
$$\mathbb{E}_{\phi} [-\ln Q(s_\tau | \pi)]$$

$$=$$

$$H[Q(s_\tau | \pi)]$$

Maximum entropy principle





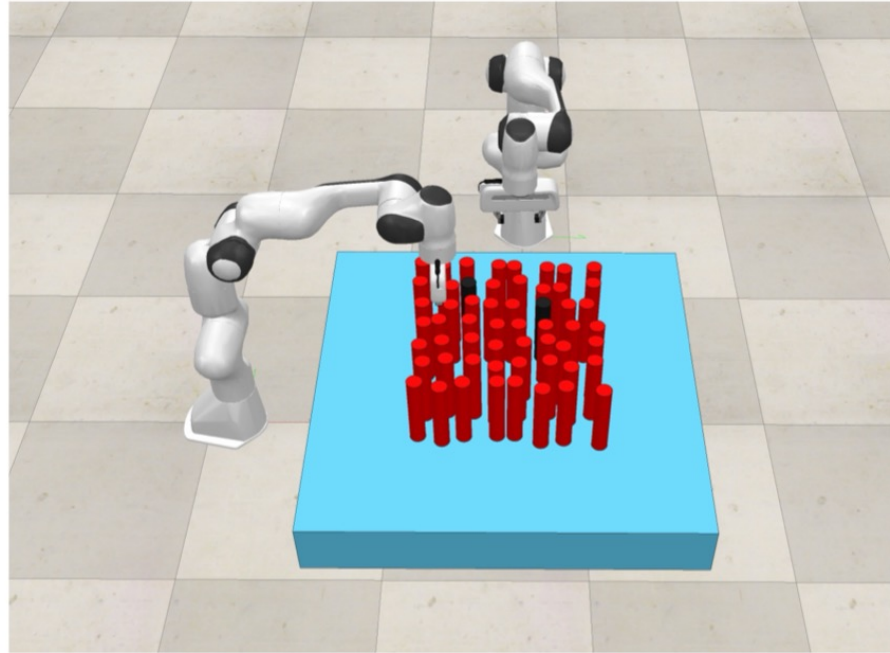
## KEY DIRECTIONS/

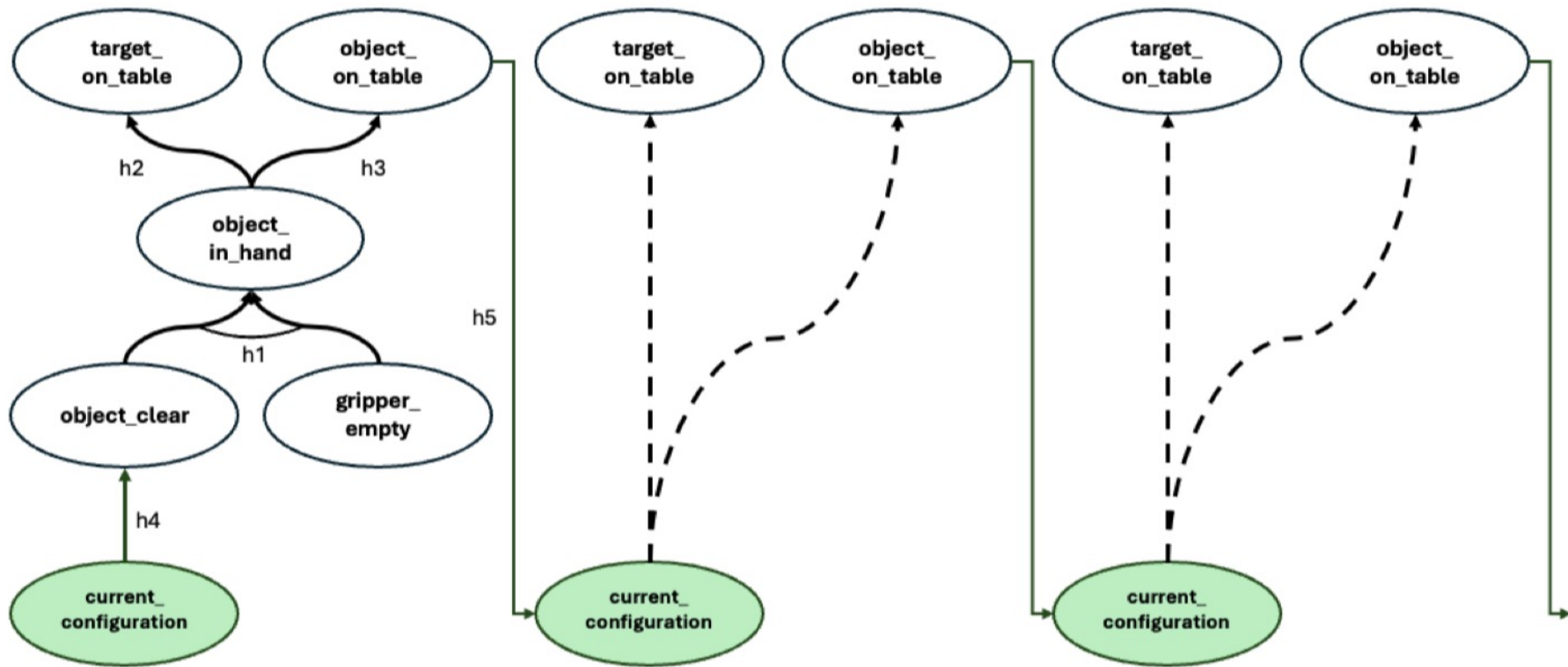
**D1/** ROBOT COGNITIVE ARCHITECTURES BASED ON BRAIN-LIKE PRINCIPLES, SUCH AS MODULARITY, STRUCTURE, HIERARCHY AND RECURSIVENESS.

**D2/** USE OF NEURO-SYMBOLIC APPROACHES TO AI, WHICH INTEGRATE DATA-DRIVEN TECHNIQUES WITHIN MODEL-BASED ARCHITECTURES.

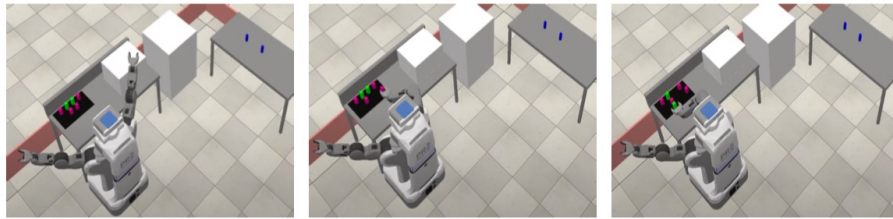
**D3/** DISTRIBUTION OF THE “INTELLIGENCE ARCHITECTURE” TO BOTH THE ROBOT “BRAIN” AND THE ROBOT “BODY”.

**D4/** ALIGNMENT BETWEEN THE ENCODED BEHAVIOURS AND HUMAN VALUES.

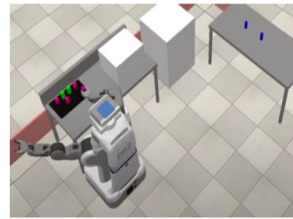




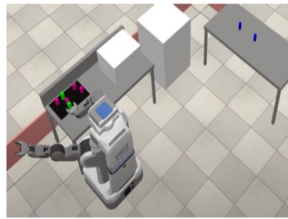




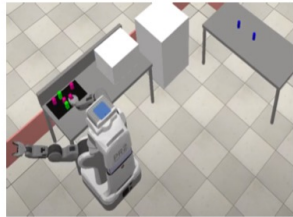
(a)



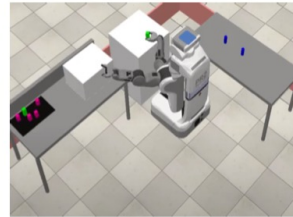
(b)



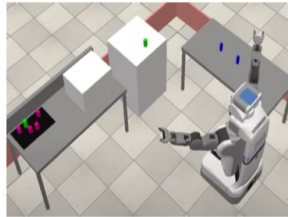
(c)



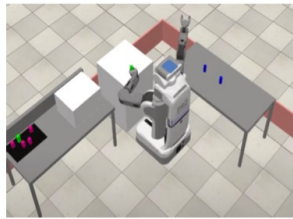
(d)



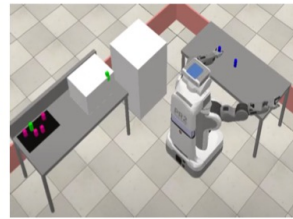
(e)



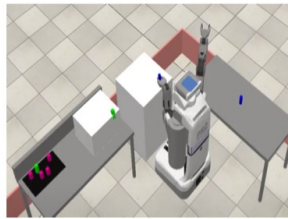
(f)



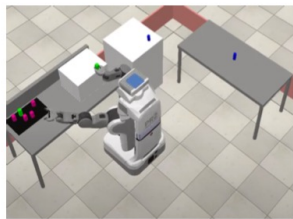
(g)



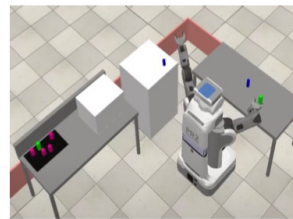
(h)



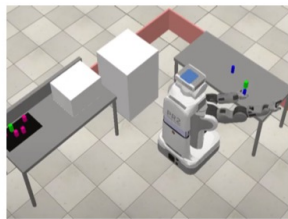
(i)



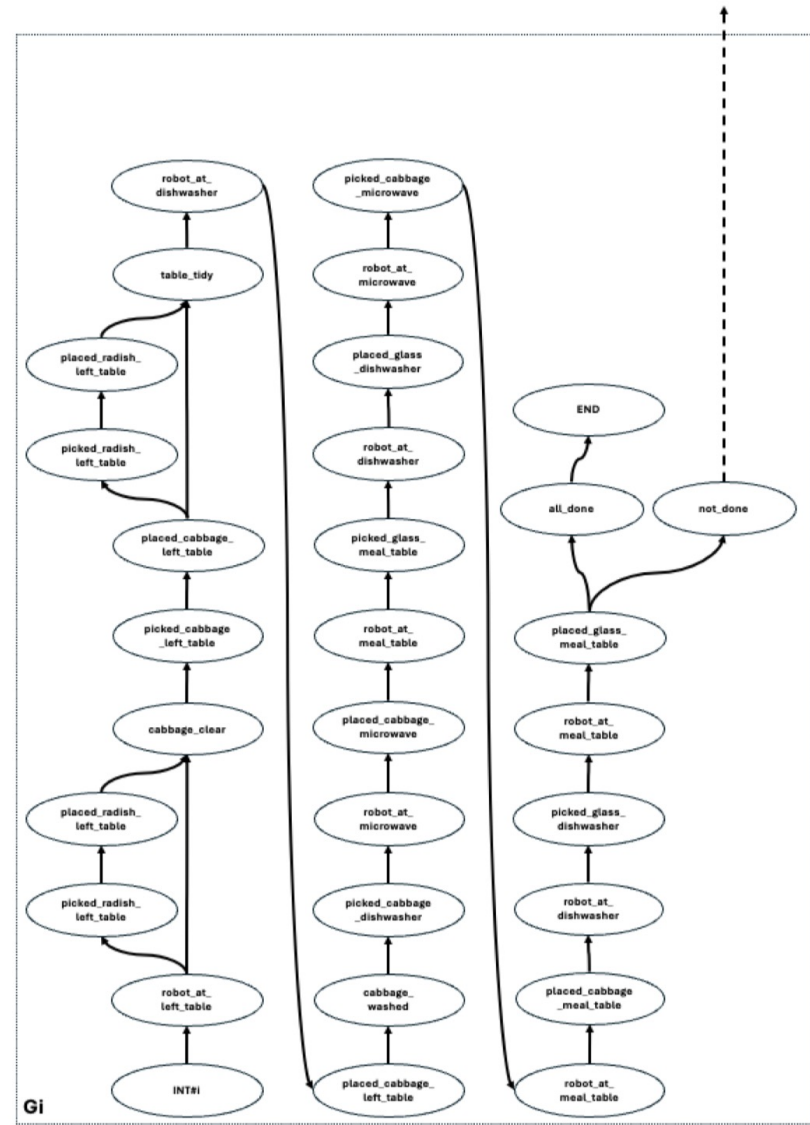
(j)



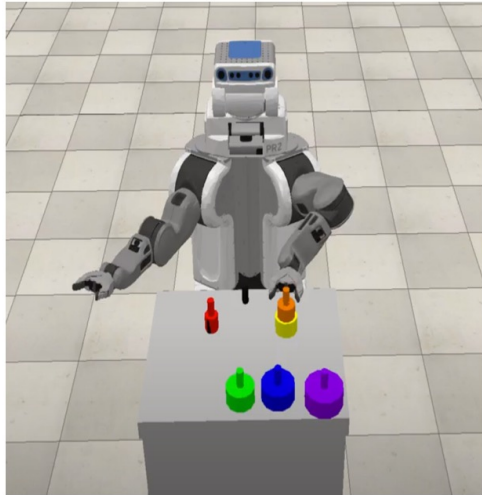
(k)



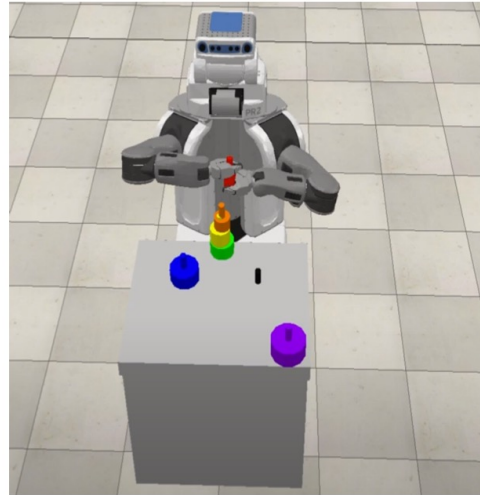
(l)



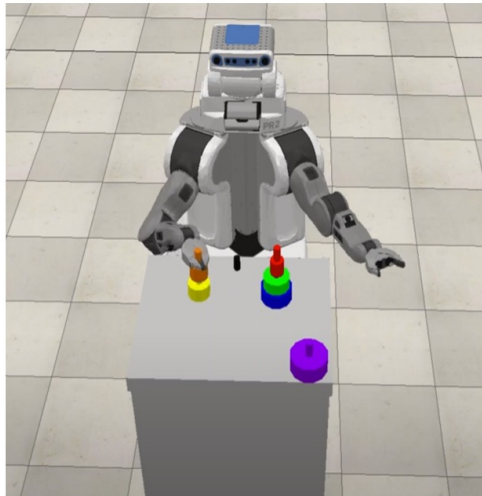




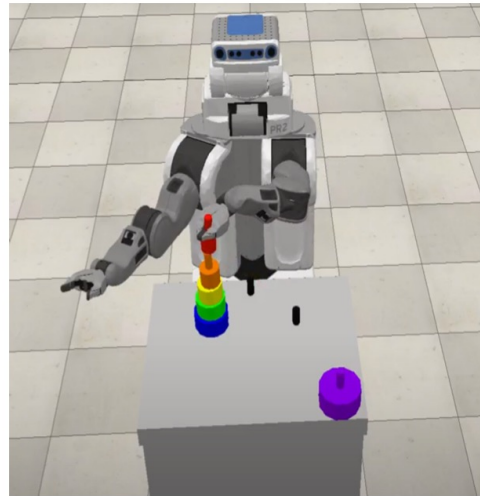
(a)



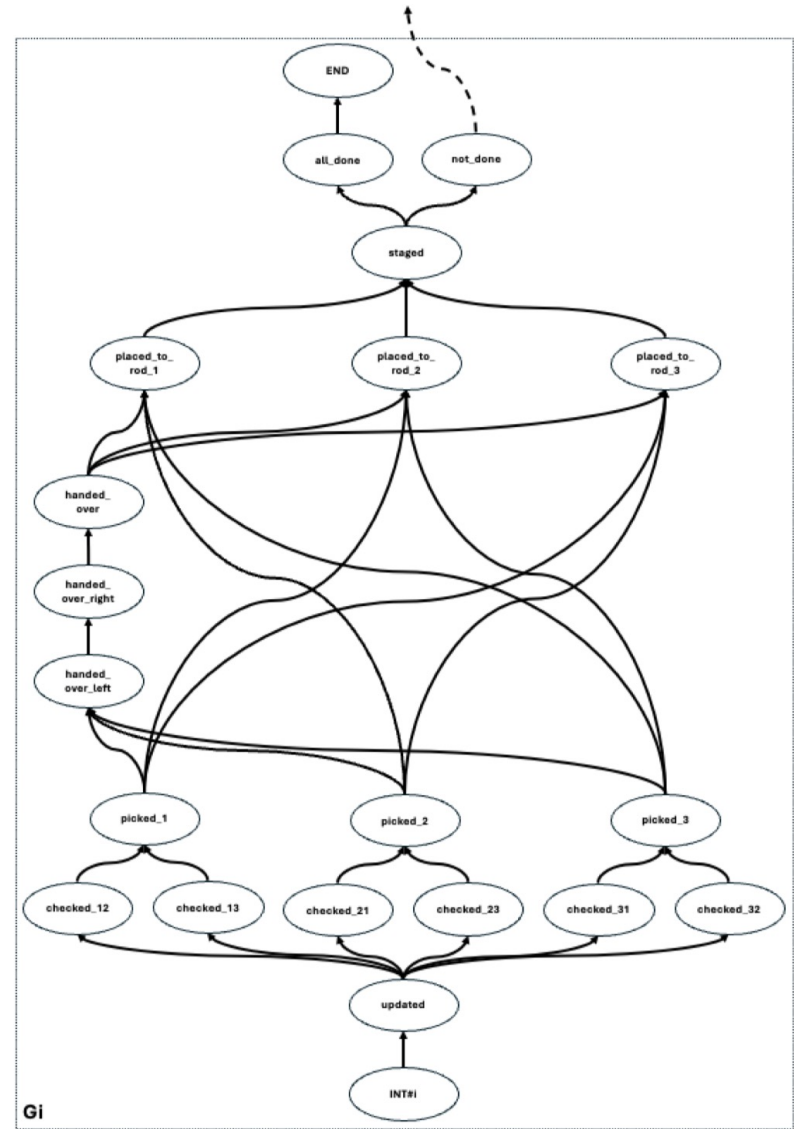
(b)



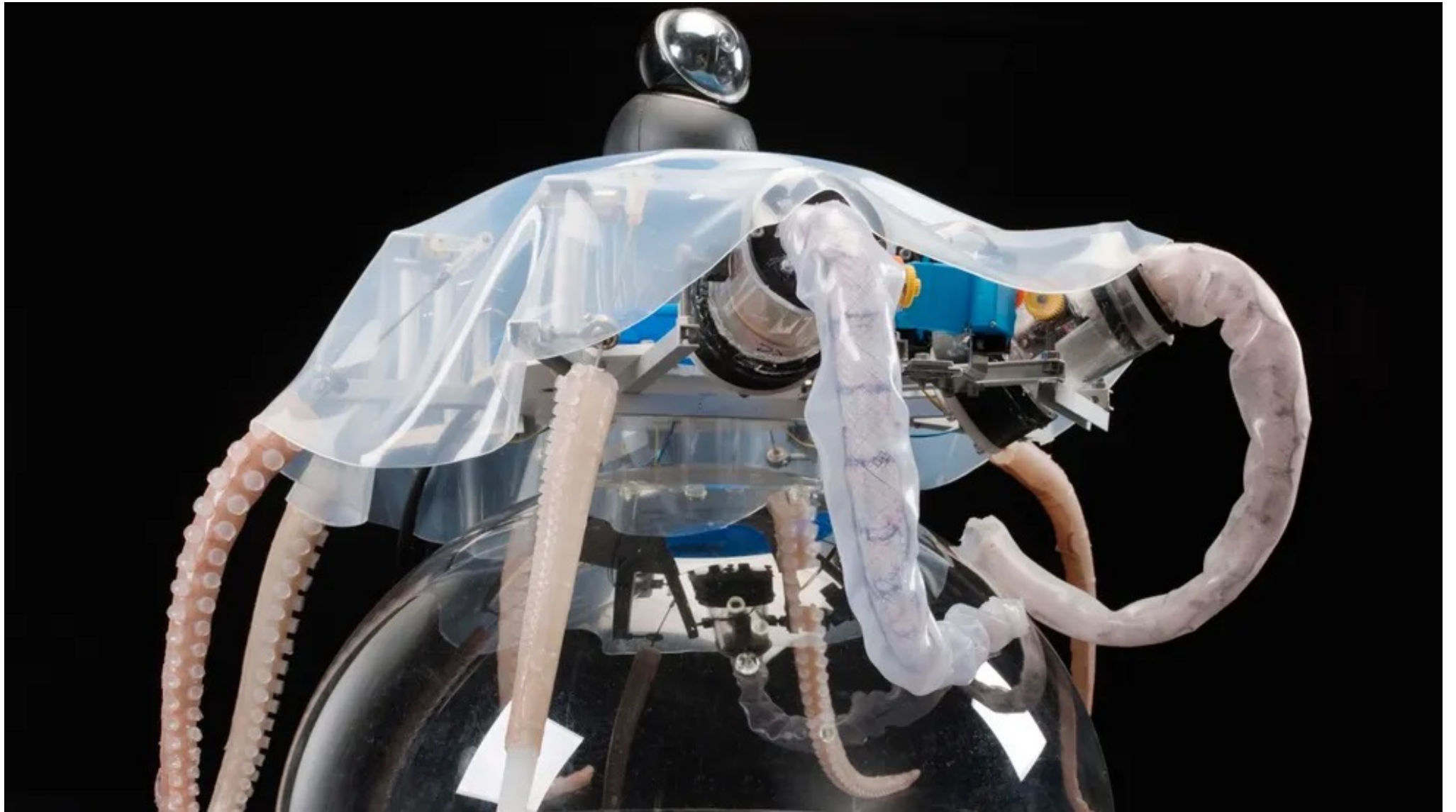
(c)



(d)



G1





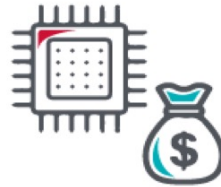
Evade  
Shutdown



Hack  
Computer  
Systems



Make  
Copies



Acquire  
Resources



Ethics  
Violation



Hire or  
Manipulate  
Humans



AI Research  
&  
Programming



Persuasion  
&  
Lobbying



Hide  
Unwanted  
Behaviors



Strategically  
Appear Aligned



Escape  
Containment



Research  
&  
Development



Manufacturing  
&  
Robotics



Autonomous  
Weaponry



THANK YOU!

Fulvio Mastrogiovanni

Professor of Robotics and Artificial Intelligence, University of Genoa  
Founder and Chief Science Officer, Teseo srl  
Coordinator of the Scientific Committee, Digital Innovation Hub Liguria