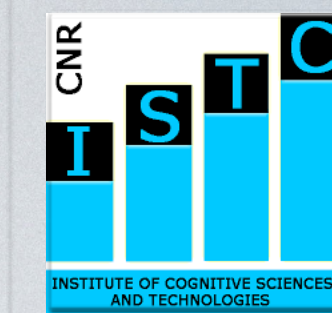




UNIVERSITÀ  
DEGLI STUDI  
DI PALERMO

**dj** dipartimento  
di ingegneria  
unipa



# ENDOWING ROBOTS WITH SELF-MODELING ABILITIES FOR TRUSTFUL HUMAN-ROBOT INTERACTIONS

*Cristiano Castelfranchi, Antonio Chella, Rino Falcone, Francesco Lanza and Valeria Seidita*



20<sup>TH</sup> WORKSHOP  
“FROM OBJECTS TO  
AGENTS”



# OBJECTIVE

---

***How to model and develop teammate robots performing trustful interaction with humans?***

- Modeling and representing robot's knowledge
- The robot has decide and act in an autonomous fashion
- The robot has be self-adaptive



# KEY IDEAS

- Triggering the decision process by means of attributing mental state to itself and to the others
- Integrating self-modeling and trust
- Employing BDI paradigm and Jason
  - extending BDI reasoning cycle for including self-modeling and justification



# HUMAN-ROBOT TEAMING INTERACTION

---

- Two main situations:
  - Known and unchanging environment
  - Partially known and changing environment



# HUMAN-ROBOT TEAMING INTERACTION

---

- Known and unchanging environment
  - Each teammate knows everything before starting
  - Use knowledge for performing task and plans

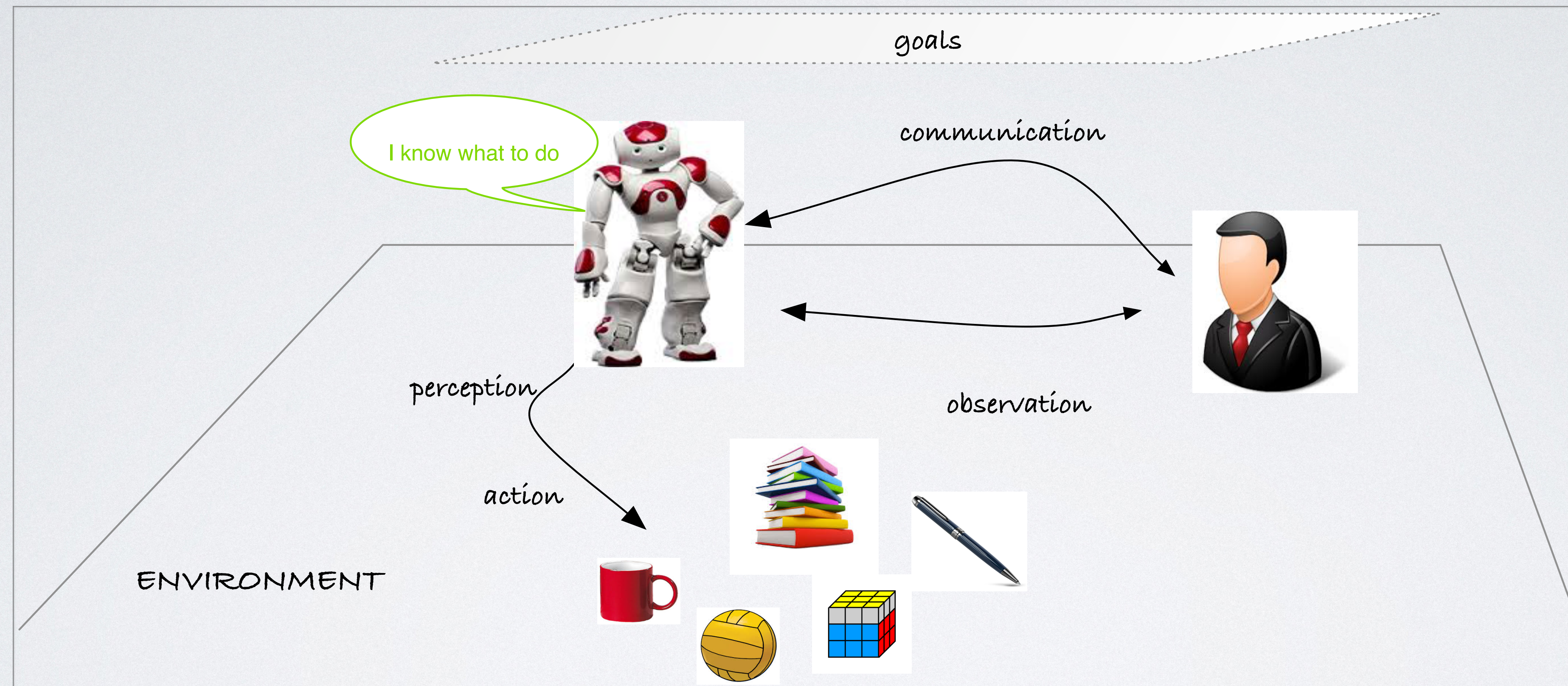


# HUMAN-ROBOT TEAMING INTERACTION

- Partially known and changing environment
  - Each teammate interacts mainly for:
    - Enhancing knowledge on the environment and on himself
    - Acquiring knowledge on what to do
- Each teammate is:
  - aware of his own limitation and capabilities
  - establish a level of confidence in the other



# HUMAN-ROBOT TEAMING INTERACTION





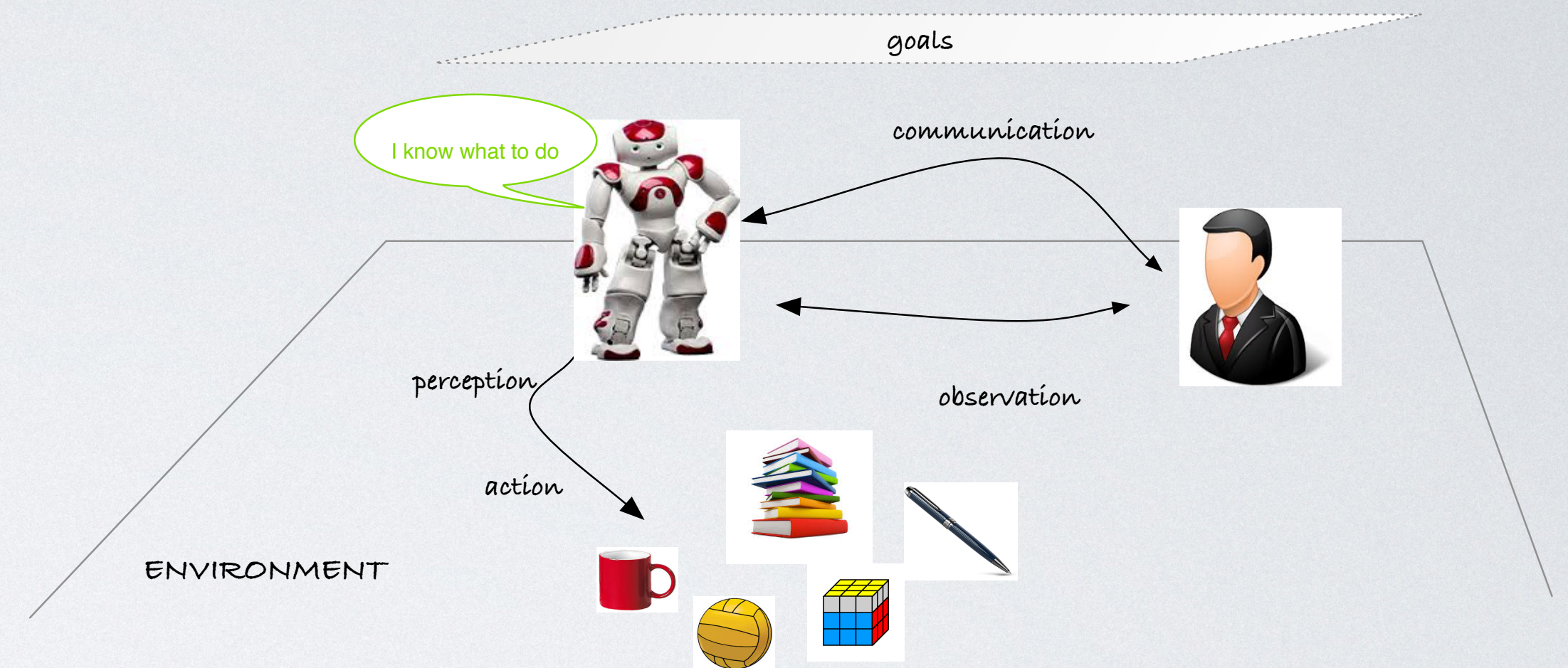
# HUMAN-ROBOT TEAMING INTERACTION

- Complex systems
  - where requirements are identified at runtime
    - changing environment conditions
    - presence of interacting users
  - global behavior emerges at runtime
- Need for exhibiting adaptation

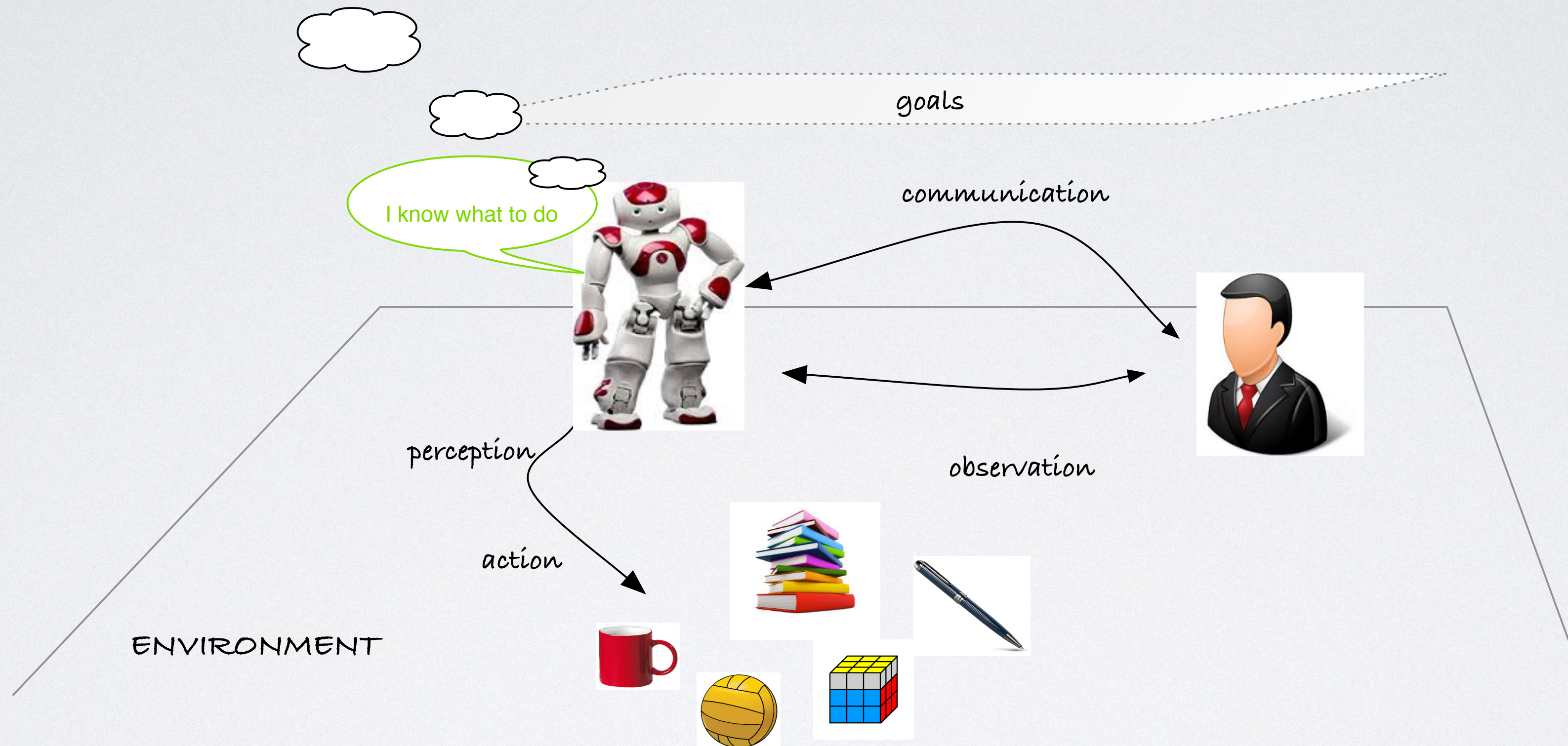
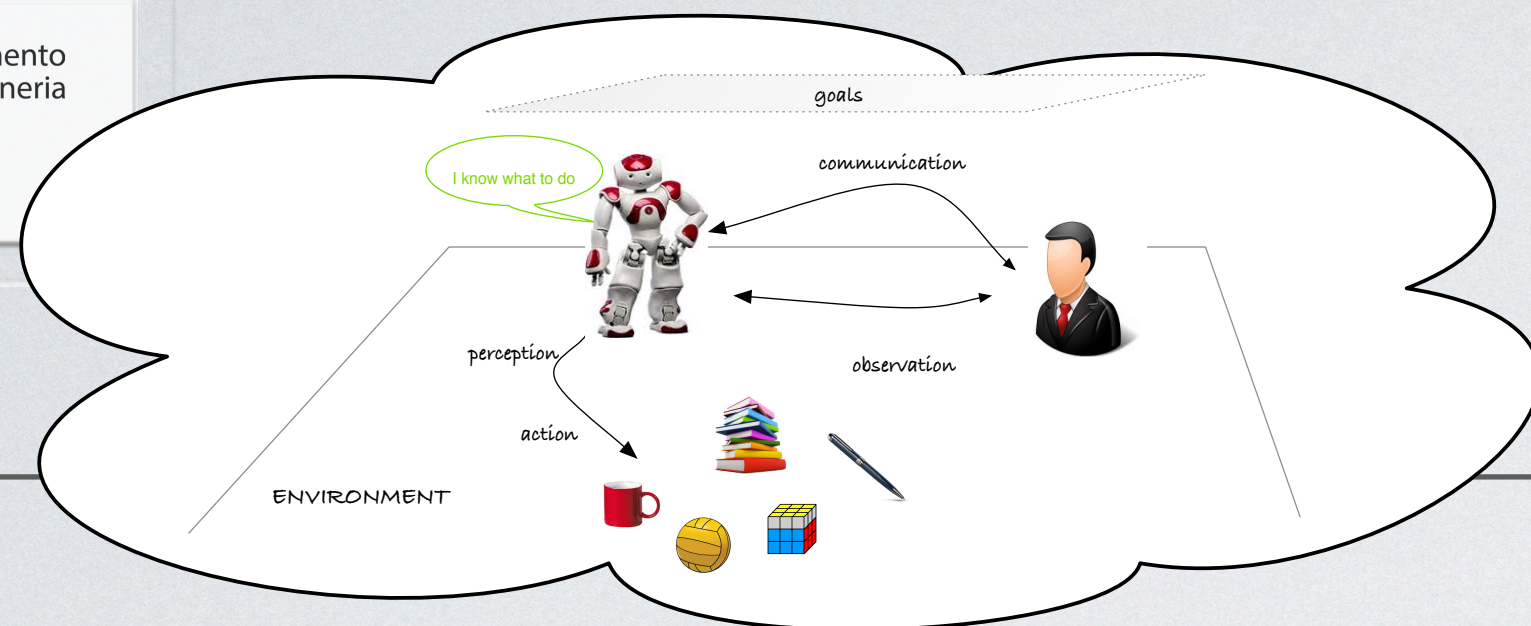


# A NEW PERSPECTIVE

- Changing perspective
- Robot is part of the environment
  - it senses itself as part of the environment along with all the other objects/agents









# CHALLENGES

---

- Environment is not known a-priori
  - plans, tasks and all that is necessary for acting and deciding cannot be established at design time
- Equipping the robot with the ability **(at runtime)** to select the best action to perform
  - Knowledge acquisition
  - Decision-making process
- Several ingredients trigger the decision process:
  - goals, capabilities, mental states, emotions, **trust**...in general: awareness



# SELF-MODELING AND TRUST

---

- The role of self-modeling and trust -> triggering the decision process
- The robot si able to create a model of the self
  - It is able to select an action on the base of what it knows about itself
  - it is able to decide which action to adopt



# IMPLEMENTING SELF-MODELING ABILITIES

---

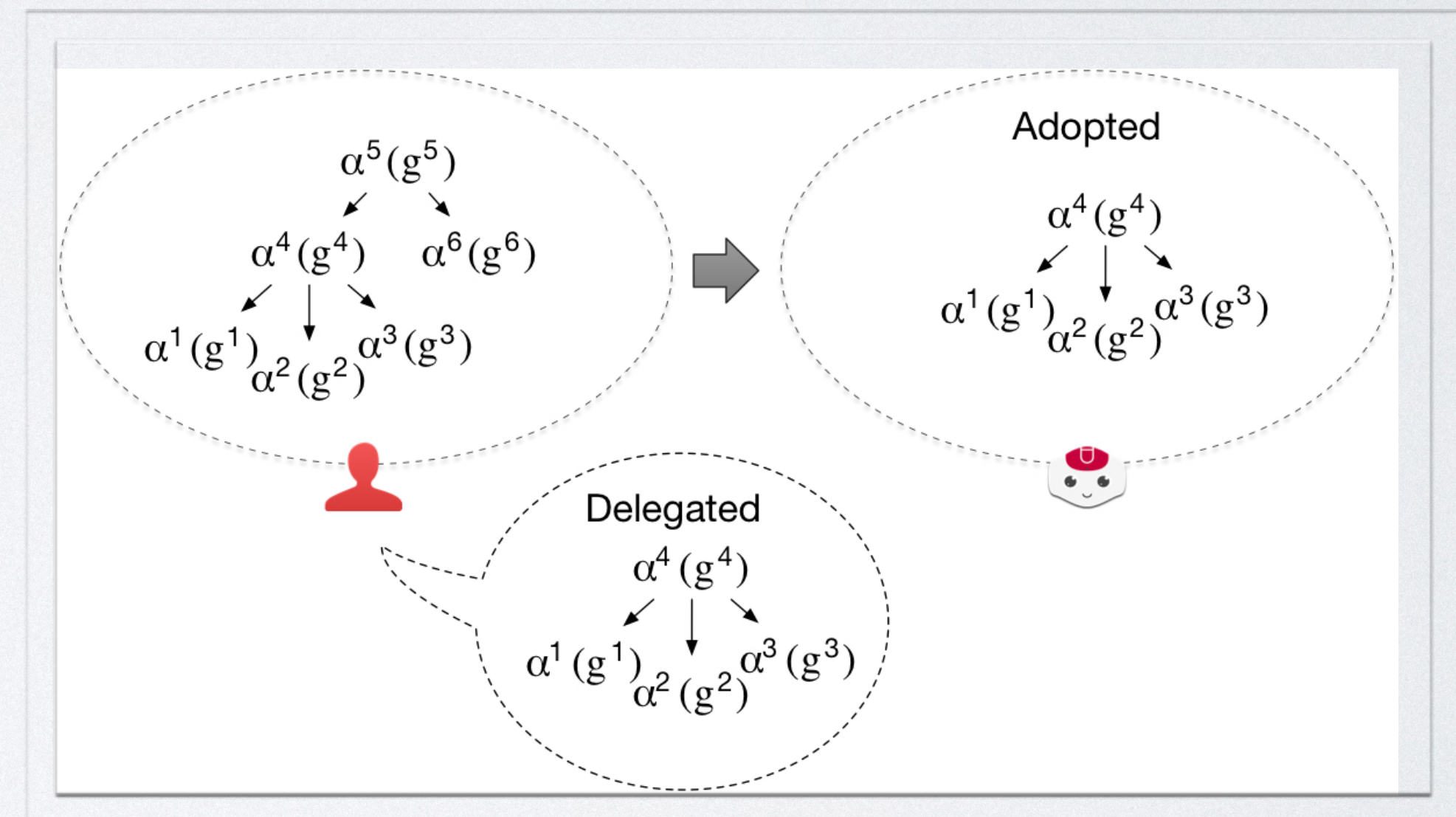
- Exploiting BDI practical reasoning and the trust model
- We extend the deliberation process



# THE TRUST MODEL

- The trust of a trustor agent in a trustee agent for a specific context to perform acts to realize the outcome result.

$$TRUST(X \ Y \ C \ \tau \ g_X)$$





# THE BDI REASONING CYCLE

Robot has:  
to commit or delegate actions and  
reason on its own action

```
1.   $B \leftarrow B_0$ ;          /*  $B_0$  are initial beliefs */
2.   $I \leftarrow I_0$ ;        /*  $I_0$  are initial intentions */
3.  while true do
4.    get next percept  $\rho$  via sensors;
5.     $B \leftarrow brf(B, \rho)$ ;
6.     $D \leftarrow options(B, I)$ ;
7.     $I \leftarrow filter(B, D, I)$ ;
8.     $\pi \leftarrow plan(B, I, Ac)$ ; /*  $Ac$  is the set of actions */
9.    while not ( $empty(\pi)$  or  $succeeded(I, B)$  or  $impossible(I, B)$ ) do
10.      $\alpha \leftarrow$  first element of  $\pi$ ;
11.      $execute(\alpha)$ ;
12.      $\pi \leftarrow$  tail of  $\pi$ ;
13.     observe environment to get next percept  $\rho$ ;
14.      $B \leftarrow brf(B, \rho)$ ;
15.     if  $reconsider(I, B)$  then
16.        $D \leftarrow options(B, I)$ ;
17.        $I \leftarrow filter(B, D, I)$ ;
18.     end-if
19.     if not  $sound(\pi, I, B)$  then
20.        $\pi \leftarrow plan(B, I, Ac)$ 
21.     end-if
22.   end-while
23. end-while
```



# MAPPING TRUST TO BDI

- Making beliefs explicit
- Breaking down actions and results  $\longrightarrow$  plans and sub-results

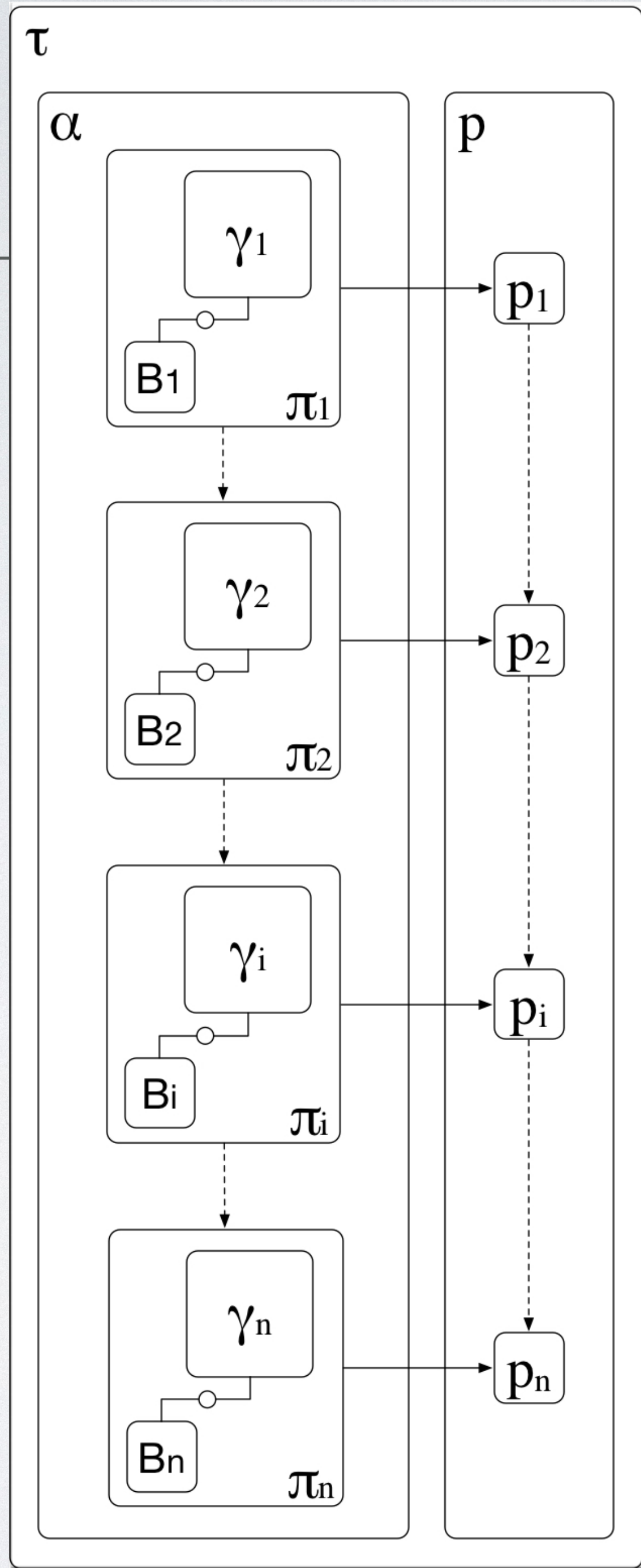
$$TRUST(X \ Y \ C \ \tau \ g_X)$$

where  $\tau = (\alpha, p)$  and  $g_X \equiv p$ ;

$$\tau = (\alpha, p) \quad \text{where} \quad \alpha = \bigcup_{i=1}^n \pi_i \quad \text{and} \quad p = \bigcup_{i=1}^n p_i$$

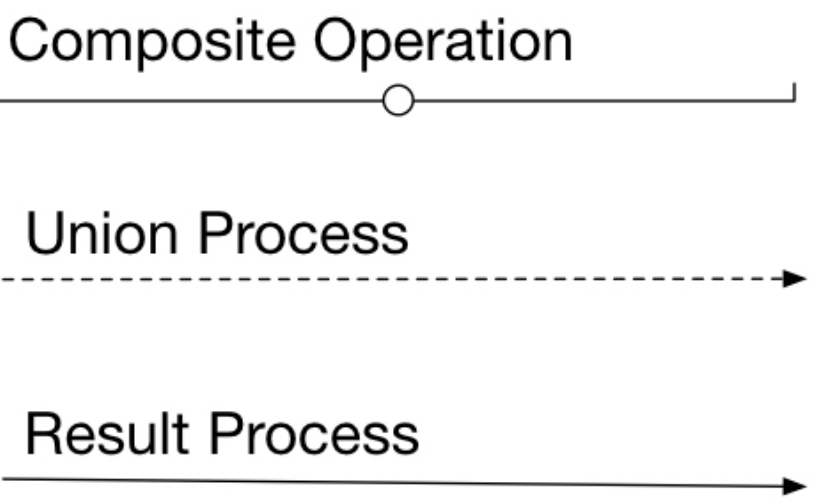
$$\pi_i = \gamma_i \circ B_i \Rightarrow \alpha = \bigcup_{i=1}^n (\gamma_i \circ B_i)$$





### Legend

$\tau$	causal process
$\alpha$	union of couple plan-kb
$p$	union of sub-goal
$\pi_i$	couple plan-kb
$\gamma_i$	sub-action
$B_i$	partial knowledge base
$p_i$	sub-goal





```

1.   $B \leftarrow B_0;$       /*  $B_0$  are initial beliefs */
2.   $I \leftarrow I_0;$       /*  $I_0$  are initial intentions */
3.  while true do
4.    get next percept  $\rho$  via sensors;
5.     $B \leftarrow brf(B, \rho);$ 
6.     $D \leftarrow options(B, I);$ 
7.     $I \leftarrow filter(B, D, I);$ 
8.     $\pi \leftarrow plan(B, I, A_c);$  /*  $A_c$  is the set of actions */
9.    while not ( $empty(\pi)$  or  $succeeded(I, B)$  or  $impossible(I, B)$ ) do
10.      $\alpha \leftarrow$  first element of  $\pi;$ 
11.      $execute(\alpha);$ 
12.      $\pi \leftarrow$  tail of  $\pi;$ 
13.     observe environment to get next percept  $\rho;$ 
14.      $B \leftarrow brf(B, \rho);$ 
15.     if  $reconsider(I, B)$  then
16.        $D \leftarrow options(B, I);$ 
17.        $I \leftarrow filter(B, D, I);$ 
18.     end-if
19.     if not  $sound(\pi, I, B)$  then
20.        $\pi \leftarrow plan(B, I, A_c)$ 
21.     end-if
22.   end-while
23. end-while

```

$A_c \leftarrow action(B_{\alpha_i}, Cap)$

$evaluate(\alpha_i);$

$J \leftarrow justify(\alpha_i, B_{\alpha_i});$

**Self-Modeling  
Trustful interaction**



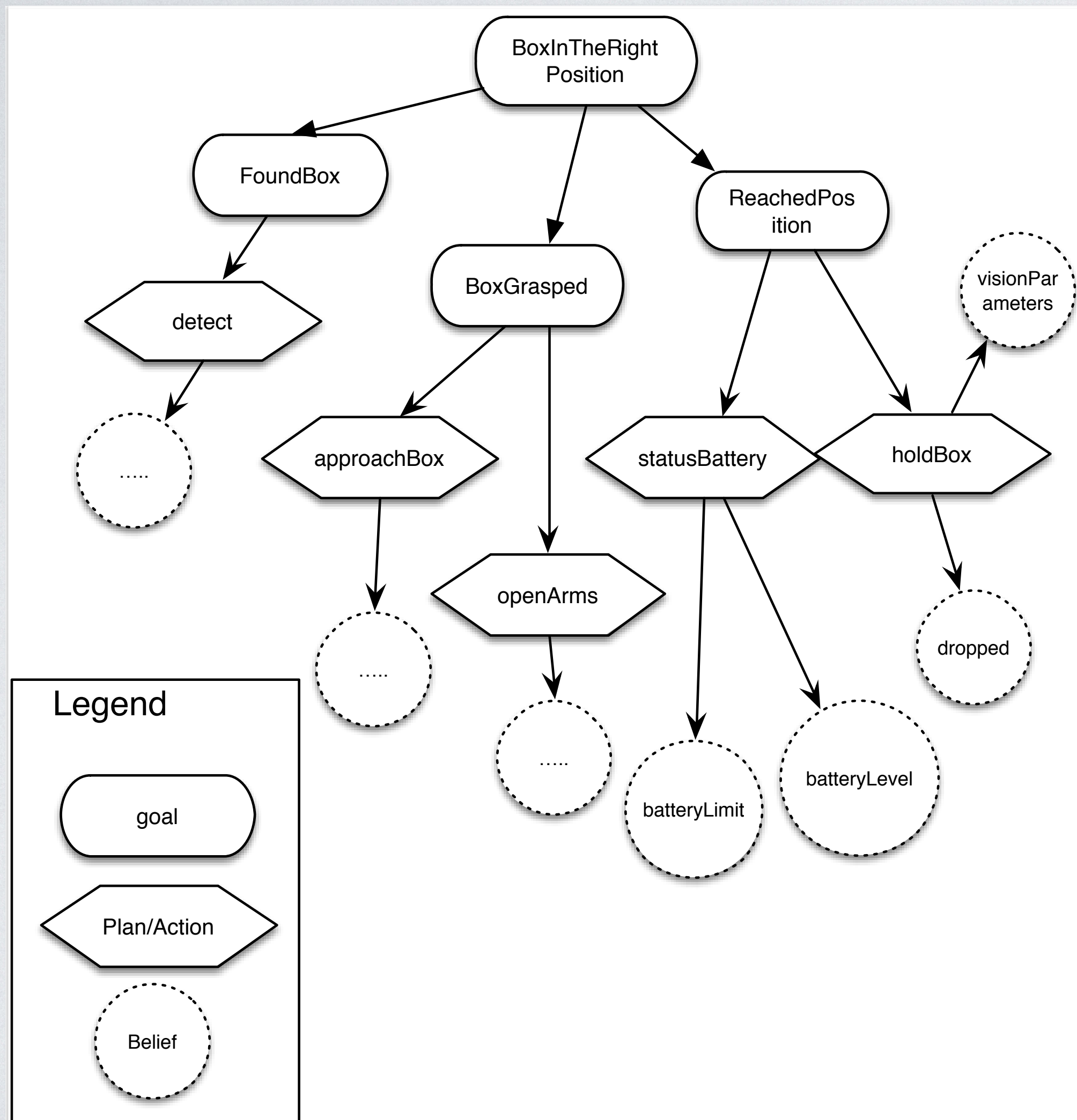
# THE ROBOT IN ACTION

---

- Jason Agent
- CArtAgO artifact
- CArtAgO @Operation
- a reference model of the environment
  - all the internal elements of the agent/robot as part of the environment



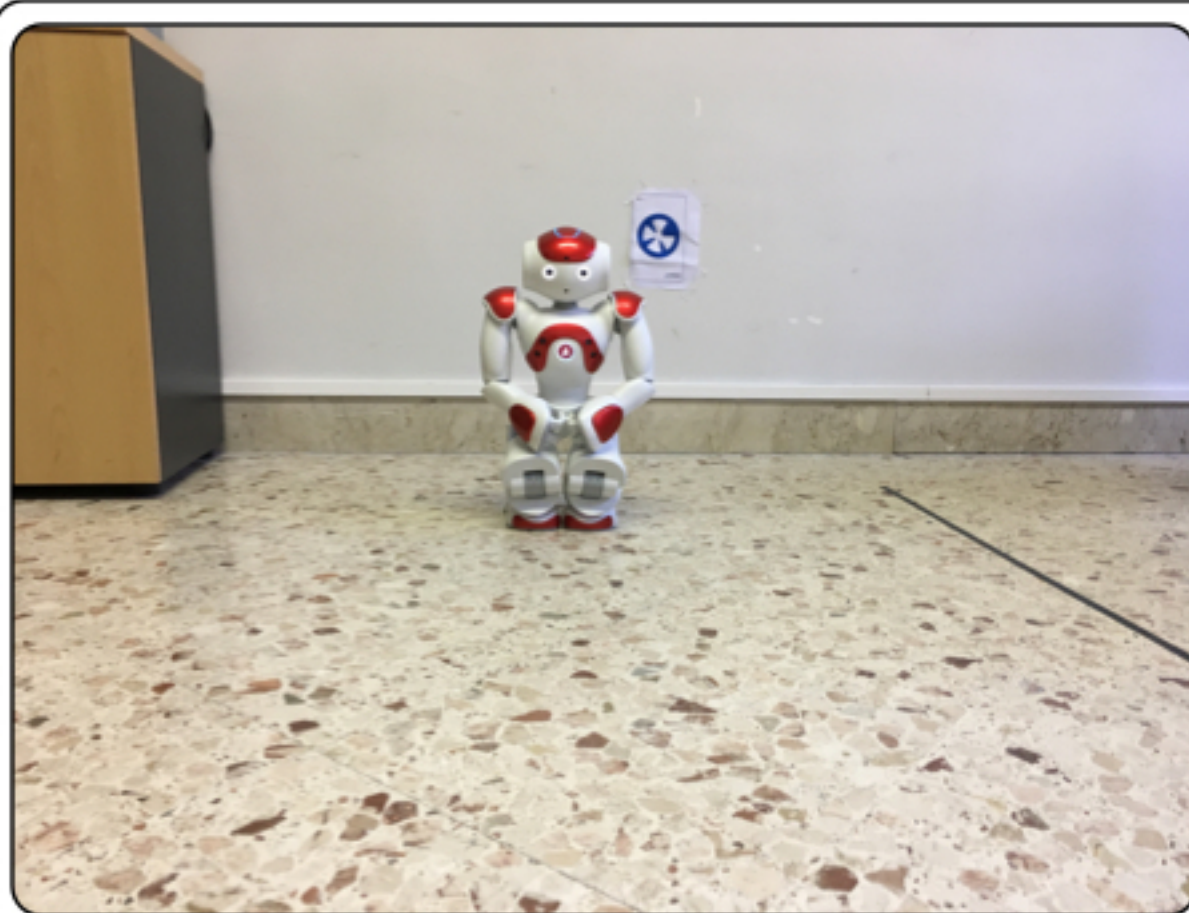
# THE ROBOT IN ACTION



- 1  $+\text{!ReachedPosition: true} \leftarrow \text{goAhead; holdBox. } [\tau];$
- 2  $+\text{!goAhead: batteryLimit}(X) \ \& \ \text{batteryLevel}(Y) \ \& \ Y < X \leftarrow \text{say}(\text{"My battery is exhaust. Please let me charge."}). [\gamma_1^+];$
- 3  $+\text{!goAhead: batteryLimit}(X) \ \& \ \text{batteryLevel}(Y) \ \& \ Y \geq X \leftarrow \text{execActions. } [\gamma_1^-];$
- 4  $B_1: \text{batteryLimit, batteryLevel};$
- 5  $+\text{!holdBox: dropped}(X) \ \& \ \text{visionParameters}(Y) \ \& \ X == \text{false} \leftarrow \text{execAct}(Y). [\gamma_2^+];$
- 6  $+\text{!holdBox: dropped}(X) \ \& \ \text{visionParameters}(Y) \ \& \ X == \text{true} \leftarrow \text{say}(\text{"The box is dropped."}). [\gamma_2^-];$
- 7  $B_2: \text{dropped, visionParameters};$



# THE ROBOT IN ACTION



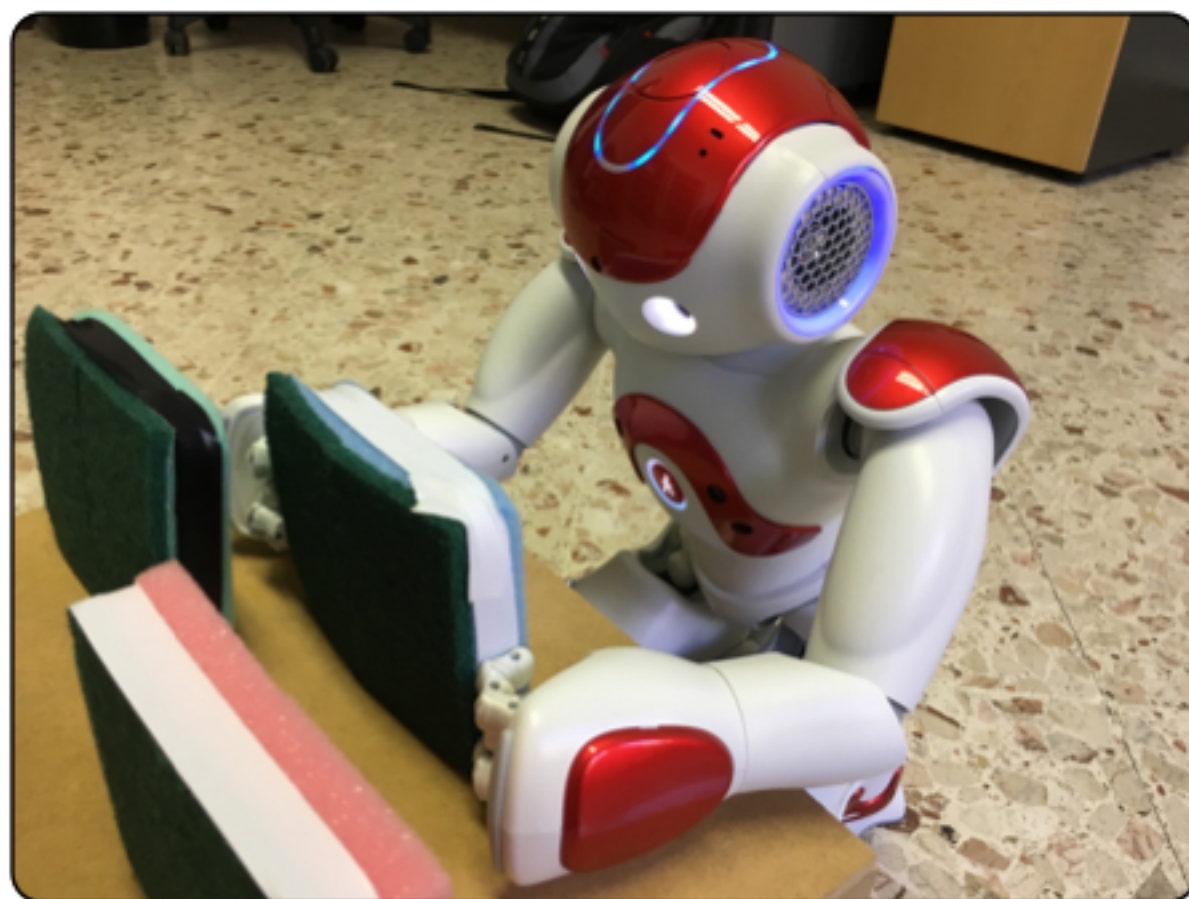
Initial Position



Boxes Assets



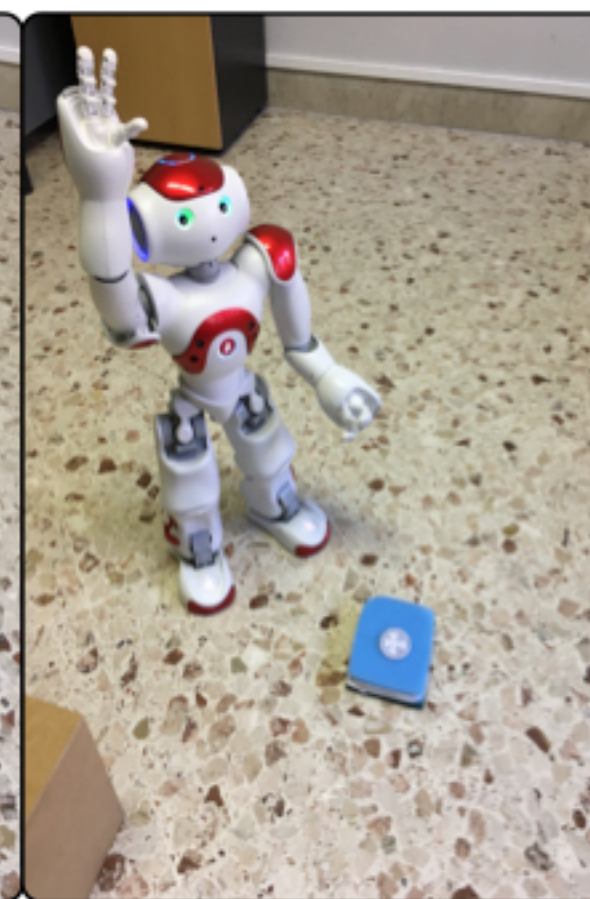
Detecting Phase



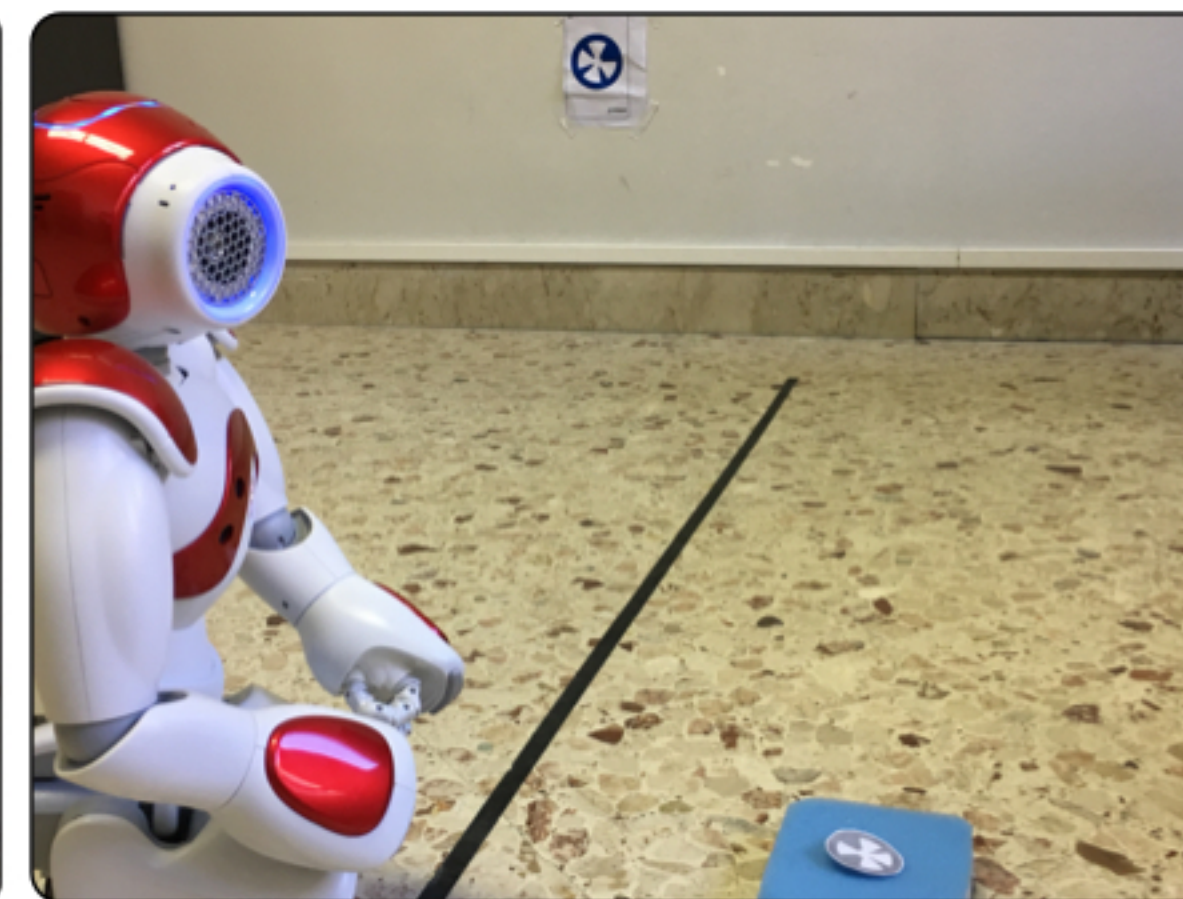
Taking Box



Box Dropped



Justification Phase



Box In The Right Position



# CONCLUSIONS AND REMARKS

---

- Equipping the robot with self-modeling abilities
- Integrating Trust model with BDI deliberation process
- Exploiting JASON and CArtAgO
  - natively support BDI theory and have a well-established counterpart for actions, plans, knowledge implementation



# CONCLUSIONS AND REMARKS

- Trust as a first element for triggering the decision process
- The integration in the BDI cycle —> two main elements of interaction
  - self-modeling
  - trust level in the interaction
- In the future:
  - implementing the other levels of adoption/delegation by Falcone&Castelfranchi
  - adding organization (MOISE)
  - adding other elements for triggering the decision process —> theory of mind

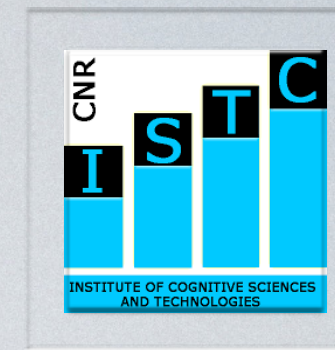




UNIVERSITÀ  
DEGLI STUDI  
DI PALERMO



dipartimento  
di ingegneria  
unipa



Thanks for your attention!

**[valeria.seidita@unipa.it](mailto:valeria.seidita@unipa.it)**