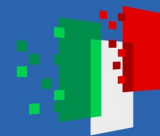




Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



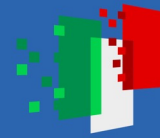
Future
Artificial
Intelligence
Research

Integrazione di Informazioni Fisiche in LLM Multimodali per la Manipolazione Robotica

Giulia Pasquale,
Istituto Italiano di
Tecnologia

24/09/2024, Napoli





Progresso di Large Language Models e Vision-Language Models

LLM e VLM permettono interazione naturale con le macchine

Come estendere questo progresso alla robotica?

Un solo VLM può essere istruito per svariati task visivi:

Embodied machines



“describe the image”

“what is close to the pear?”

“what time is it?”

VLM

“a kitchen counter with a hand cutting a peach, ...”

“a grey jar with wooden lid”

“11:50”

Manipolazione di oggetti

perception



+ proprioception

“cut the peach”

“put the pear in the bowl”

“add salt to the pasta”

Multimodal Model

action

Progresso di Large Language Models e Vision-Language Models

LLM e VLM permettono interazione naturale con le macchine

Un solo VLM può essere istruito per svariati task visivi:



“describe the image”

“what is close to the pear?”

“what time is it?”

VLM

“a kitchen counter with a hand cutting a peach, ...”

“a grey jar with wooden lid”

“11:50”

Obiettivi a lungo termine:

- Apprendimento di task di manipolazione
- “Grounding” di comandi espressi in linguaggio naturale

Manipolazione di oggetti

perception



+ proprioception

“cut the peach”

“put the pear in the bowl”

“add salt to the pasta”

Multimodal Model

action



Progresso di Large Language Models e Vision-Language Models

LLM e VLM permettono interazione naturale con le macchine

Un solo VLM può essere istruito per svariati task visivi:



“describe the image”

“what is close to the pear?”

“what time is it?”

VLM

“a kitchen counter with a hand cutting a peach, ...”

“a grey jar with wooden lid”

“11:50”

Obiettivi a lungo termine:

- **Apprendimento di task di manipolazione**
- “Grounding” di comandi espressi in linguaggio naturale

Manipolazione di oggetti

perception



+ proprioception

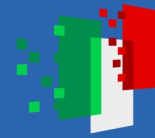
“cut the peach”

“put the pear in the bowl”

“add salt to the pasta”

Multimodal Model

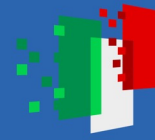
action



Apprendimento di task di manipolazione

(In Robotica) esistono due paradigmi per l'apprendimento di task di manipolazione:

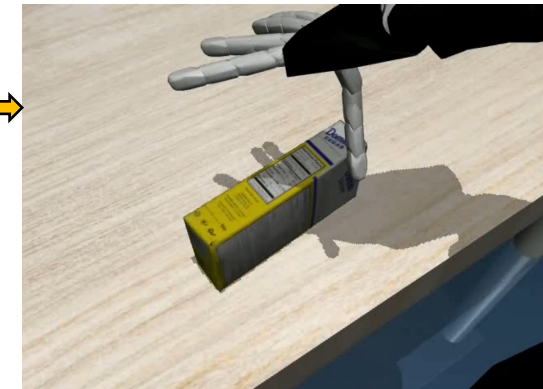
- *Reinforcement Learning*
 - Può essere indispensabile
 - Spesso parte in simulazione (la fase iniziale è critica)
- *Imitation Learning*
 - Può velocizzare la fase iniziale dell'apprendimento
 - Necessita di dimostrazioni (ed è limitato dalle stesse)



Accelerazione di metodi di Reinforcement Learning (RL)

(In Robotica) esistono due paradigmi per l'apprendimento di task di manipolazione:

- **Reinforcement Learning**
 - Può essere indispensabile
 - Spesso parte in simulazione (la fase iniziale è critica)
- **Imitation Learning**
 - Può velocizzare la fase iniziale dell'apprendimento
 - Necessita di dimostrazioni (ed è limitato dalle stesse)



Una parte iniziale del lavoro si è focalizzata sull'accelerazione di metodi di RL per prese con più dita:

- [1] Sfruttando modelli disponibili per la pianificazione della presa (*grasp planners*)**
- [2] Colmando il divario tra simulazione e realtà mediante *residual RL***

[1] Ceola, F., Maiettini, E., Rosasco, L., and Natale, L., A Grasp Pose is All You Need: Learning Multi-fingered Grasping with Deep Reinforcement Learning from Vision and Touch, in IEEE/RSJ International Conference on Intelligent Robots and Systems, Detroit, MI, USA, 2023.

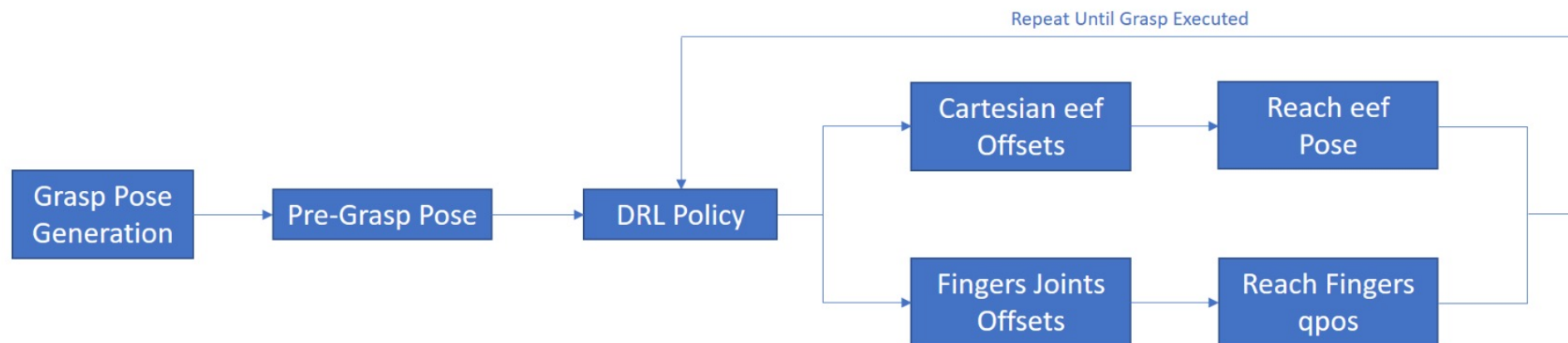
[2] Ceola, F., Rosasco, L., and Natale, L., RESPECT: Speeding-up Multi-fingered Grasping with Residual Reinforcement Learning, IEEE Robotics & Automation Letters, vol. 9, no. 4, 2024.

Pres a più dita con RL sfruttando grasp planner disponibili

G-PAYN: A Grasp Pose is All You Need

1. Il manipolatore raggiunge una posa di **pre-grasp** calcolata a partire da un **grasp planner disponibile**
2. La **policy di RL** utilizza **RGB + tatto + propriocezione** per **raffinare** gli offset cartesiani e dei giunti e completare la presa

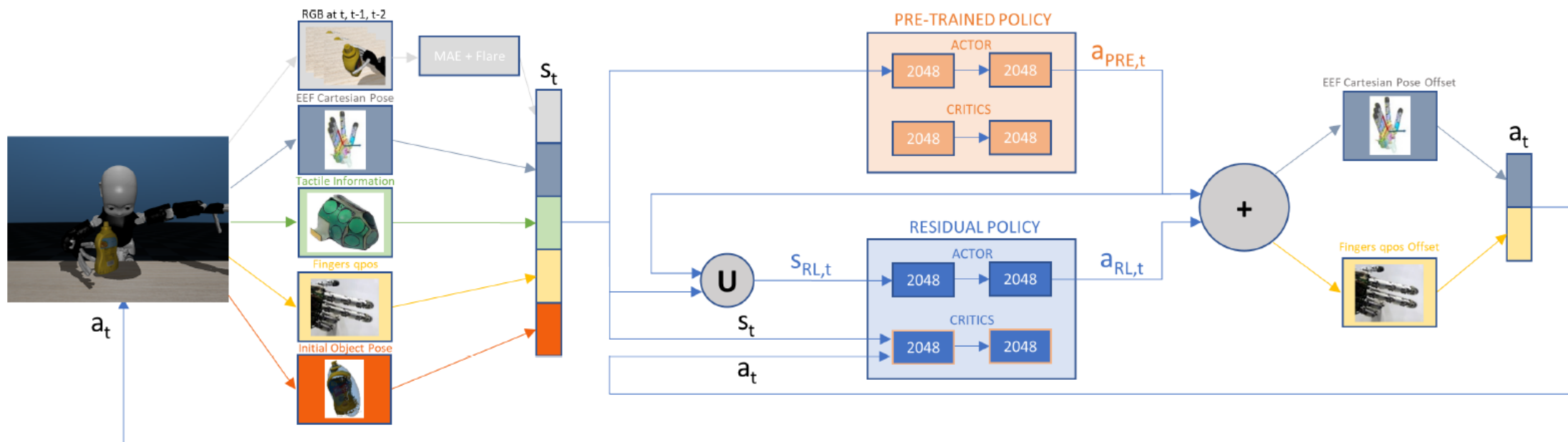
L'apprendimento della policy è accelerato da **dimostrazioni calcolate automaticamente da un grasp planner disponibile**.



Colmare il divario tra simulazione e realtà mediante residual RL

RESPRECT: RESidual learning with PREtrained CriTics

1. Policy pre-addestrata su molti oggetti in simulazione
2. Apprendimento di una policy residuale per la presa di un oggetto nuovo in una frazione del tempo (~5× più veloce)

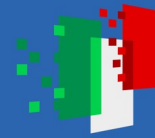




Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca

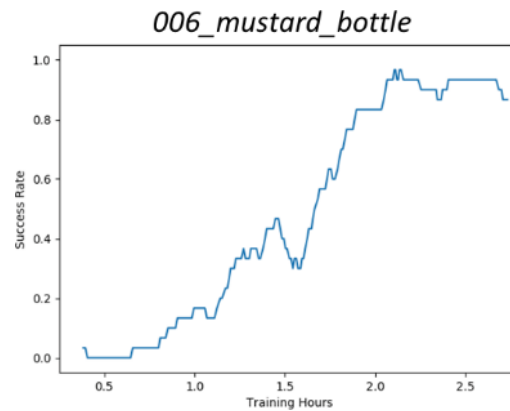
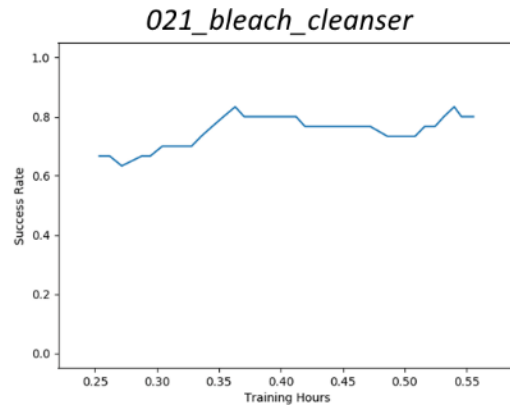


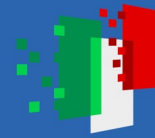
Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



Future
Artificial
Intelligence
Research

Colmare il divario tra simulazione e realtà mediante residual RL

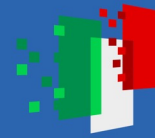




Apprendimento di task di manipolazione

(In Robotica) esistono due paradigmi per l'apprendimento di task di manipolazione:

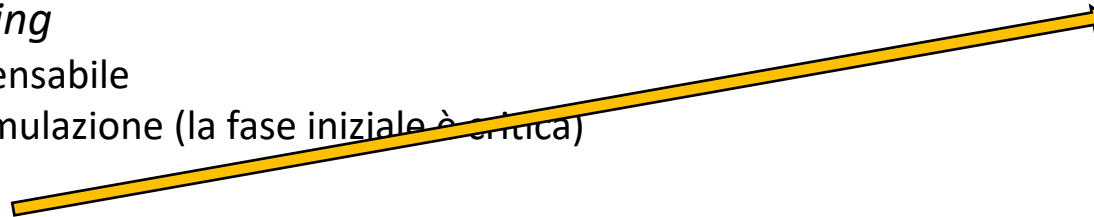
- *Reinforcement Learning*
 - Può essere indispensabile
 - Spesso parte in simulazione (la fase iniziale è critica)
- *Imitation Learning*
 - Può velocizzare la fase iniziale dell'apprendimento
 - Necessita di dimostrazioni (ed è limitato dalle stesse)



Acquisizione dati e metodi per Imitation Learning

(In Robotica) esistono due paradigmi per l'apprendimento di task di manipolazione:

- *Reinforcement Learning*
 - Può essere indispensabile
 - Spesso parte in simulazione (la fase iniziale è critica)
- **Imitation Learning**
 - Può velocizzare la fase iniziale dell'apprendimento
 - Necessita di dimostrazioni (ed è limitato dalle stesse)



Dimostrazioni da **retargeting del movimento**:

- teleoperazione del robot
 - ✓ registrazione interazioni fisiche
 - ! precise ma necessitano del setup
- video di persone (o altri robot)
 - ! come predire le forze di contatto?
 - ✓ meno precise ma disponibili

Una parte di lavoro attuale studia l'apprendimento da dimostrazioni:

- [1] **retargeting della posa della mano su manipolatori con N-DoFs**
- [2] **strategie di campionamento per Diffusion Policies**

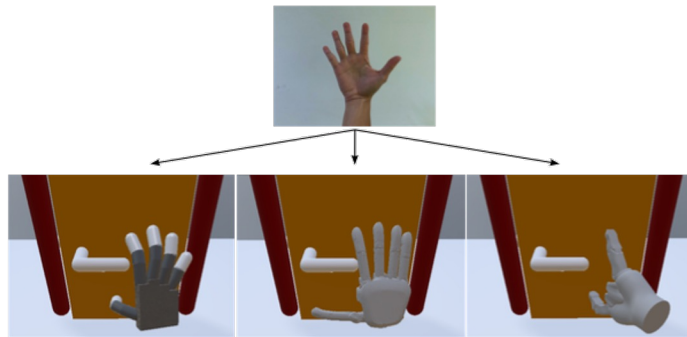
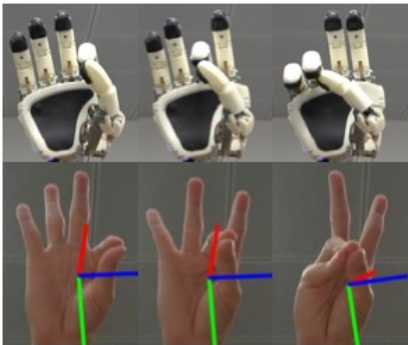
[1] Puang, E. Y., Ceola, F., Pasquale, G., and Natale, L., Hand Pose Retargeting on Manipulators with N-DoFs, submitted.

[2] Rosasco, A., Ceola, F., Pasquale G., and Natale, L., D-BTS: A Sampling Strategy for Diffusion-Binned Trajectory Selection, submitted.



Retargeting della posa della mano

1. Posa della mano espressa da un insieme di ancore
2. Composizione di metodi per dimensionality reduction
→ Sinergie comuni tra manipolatori diversi

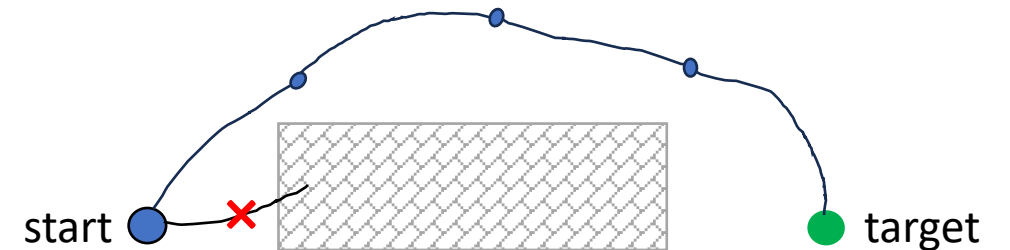


Strategie di campionamento per Diffusion Policies

Diffusion Policy [1] è un approccio che applica metodi di generazione di immagini alla generazione di azioni.

La stocasticità del processo di denoising può generare azioni particolarmente fallimentari.

→ Studiamo strategie di campionamento per evitarle.



Rosasco, A., Ceola, F., Pasquale G., and Natale, L., D-BTS: A Sampling Strategy for Diffusion-Binned Trajectory Selection, submitted.

[1] Chi, C., Xu, Z., Feng, S., Cousineau, E., Du, Y., Burchfiel, B., Tedrake, R., and Song, S., Diffusion Policy: Visuomotor Policy Learning via Action Diffusion, Proceedings of Robotics: Science and Systems (RSS), 2023, & The International Journal of Robotics Research, 2024.



Progresso di Large Language Models e Vision-Language Models

LLM e VLM permettono interazione naturale con le macchine

Un solo VLM può essere istruito per svariati task visivi:



“describe the image”

“what is close to the pear?”

“what time is it?”

VLM

“a kitchen counter with a hand cutting a peach, ...”

“a grey jar with wooden lid”

“11:50”

Obiettivi a lungo termine:

- Apprendimento di task di manipolazione
- “Grounding” di comandi espressi in linguaggio naturale

Manipolazione di oggetti

perception



+ proprioception

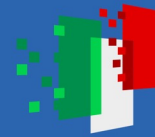
“cut the peach”

“put the pear in the bowl”

“add salt to the pasta”

Multimodal Model

action



Grounding di comandi espressi in linguaggio naturale

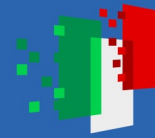
Task comunicato mediante una istruzione di alto livello espressa in linguaggio naturale:

1. **riflessione + pianificazione di alto livello:** per scomporre il task in sotto-tasks
2. **policy di controllo a basso livello:** esecuzione di ogni sotto-task



Fig. 3: **Sub-tasks decomposition** of a *Place the bowl on the plate and the cup in the bowl matching the color* sequence.

Task
Clean the pan.
Cook the capsicum and place it on a plate.
Cook the vegetables.
Dry the plate.
Hide the teddy bear in the red bowl.
Match the cups with the appropriate bowls.
Place the bowl on the plate and the cup in the bowl matching the color.
Place the bowls on the appropriate plates.
Prepare two cups of tea.
Put a highlighter on each book.
Put the ball in the red pot.
Roll the dices in the bowl.
Serve the vegetables in different plates.
Set the table.
Sort the balls from left to right in order of size.
Stack green blocks.
Stack the bowls.
Stack the cups.
Throw away the rubbish paper.
Water the potted plant and put the can on the plate.



Grounding di comandi espressi in linguaggio naturale

Task comunicato mediante una istruzione di alto livello espressa in linguaggio naturale:

1. **riflessione + pianificazione di alto livello:** per scomporre il task in sotto-tasks
2. **policy di controllo a basso livello:** esecuzione di ogni sotto-task

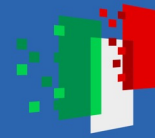


Fig. 3: **Sub-tasks decomposition** of a *Place the bowl on the plate and the cup in the bowl matching the color* sequence.

I dataset per task a ``orizzonte lungo`` soffrono le limitazioni seguenti:

- disponibili solo in ambiente simulato
- oggetti semplici in condizioni su tavolo relativamente semplici

Task
Clean the pan.
Cook the capsicum and place it on a plate.
Cook the vegetables.
Dry the plate.
Hide the teddy bear in the red bowl.
Match the cups with the appropriate bowls.
Place the bowl on the plate and the cup in the bowl matching the color.
Place the bowls on the appropriate plates.
Prepare two cups of tea.
Put a highlighter on each book.
Put the ball in the red pot.
Roll the dices in the bowl.
Serve the vegetables in different plates.
Set the table.
Sort the balls from left to right in order of size.
Stack green blocks.
Stack the bowls.
Stack the cups.
Throw away the rubbish paper.
Water the potted plant and put the can on the plate.



Grounding di comandi espressi in linguaggio naturale

LHManip: A Dataset for Long-Horizon Language-Grounded Manipulation Tasks

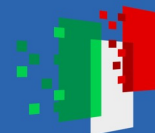
- 20 task con un totale di 33 **oggetti quotidiani**
- ambiente da tavolo **molto disordinato**
- 10 dimostrazioni per task (oggetti e configurazioni diverse)
- acquisito mediante **teleoperazione**: *visione + propriocezione + azione*



Task
Clean the pan.
Cook the capsicum and place it on a plate.
Cook the vegetables.
Dry the plate.
Hide the teddy bear in the red bowl.
Match the cups with the appropriate bowls.
Place the bowl on the plate and the cup in the bowl matching the color.
Place the bowls on the appropriate plates.
Prepare two cups of tea.
Put a highlighter on each book.
Put the ball in the red pot.
Roll the dices in the bowl.
Serve the vegetables in different plates.
Set the table.
Sort the balls from left to right in order of size.
Stack green blocks.
Stack the bowls.
Stack the cups.
Throw away the rubbish paper.
Water the potted plant and put the can on the plate.



Fig. 4: **Tasks variations**: we consider different plate-bowl colors for the *Place the bowls on the appropriate plates* task (left) and different plates for the *Dry the plate* task (right).



Conclusioni e direzioni future

Apprendimento di task di manipolazione

Accelerazione di RL

con
dimostrazioni
automatiche
[IROS 2023]

colmando
divario
sim-2-real
[RAL 2024]

Dati e metodi per IL

retargeting
posa
della mano
[submitted]

campionamento
metodi basati su
Diffusion Policies
[submitted]

Grounding comandi in linguaggio naturale

Dati

Dataset
in teleoperazione
setting realistico
[RSS workshop & ICRA 2024]

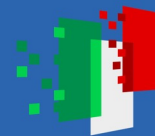
- Apprendimento con combinazione di RL e IL
- IL: Utilizzo di dimostrazioni disponibili da video di persone o altri robot
- Grounding di comandi in linguaggio naturale: Metodi
- Integrazione di sensori tattili



Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



Future
Artificial
Intelligence
Research

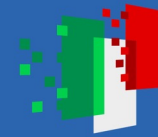
Grazie!



Finanziato
dall'Unione europea
NextGenerationEU



Ministero
dell'Università
e della Ricerca



Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA



Future
Artificial
Intelligence
Research

Riferimenti

Accelerazione di RL

F. Ceola, E. Maiettini, L. Rosasco and L. Natale, "A Grasp Pose is All You Need: Learning Multi-Fingered Grasping with Deep Reinforcement Learning from Vision and Touch," 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), Detroit, MI, USA, 2023, pp. 2985-2992, doi: 10.1109/IROS55552.2023.10341776.

F. Ceola, L. Rosasco and L. Natale, "RESPRECT: Speeding-up Multi-Fingered Grasping With Residual Reinforcement Learning," in IEEE Robotics and Automation Letters, vol. 9, no. 4, pp. 3045-3052, April 2024, doi: 10.1109/LRA.2024.3363532.

Dati e metodi per IL

Puang, E. Y., Ceola, F., Pasquale, G., and Natale, L., Hand Pose Retargeting on Manipulators with N-DoFs, submitted.

Rosasco, A., Ceola, F., Pasquale G., and Natale, L., D-BTS: A Sampling Strategy for Diffusion-Binned Trajectory Selection, submitted.

Grounding comandi in linguaggio naturale

Ceola F., Natale L., Sunderhauf N., Rana K., "LHManip: A Dataset for Long-Horizon Language-Grounded Manipulation Tasks in Cluttered Tabletop Environments", RSS Workshop on Mechanisms for Mapping Human Input to Robots From Robot Learning to Shared Control/Autonomy

Percezione tattile

G. M. Caddeo, A. Marcani, P. D. Alfano, L. Rosasco, and L. Natale, "Sim2Real Bilevel Adaptation for Object Surface Classification using Vision-Based Tactile Sensors", 2024 IEEE/RSJ International Conference on Robotics and Automation (ICRA), Yokohama, JP, 2024.