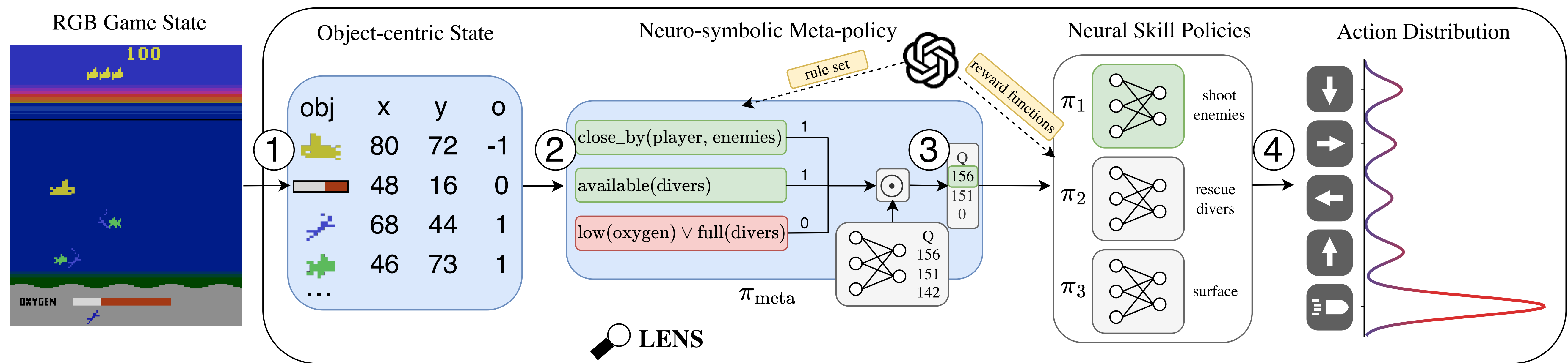


Interpretable Reinforcement Learning via Meta-Policy Guidance

Raban Emunds¹, Jannis Blüml^{1,2},
Quentin Delfosse^{1,3}, Kristian Kersting^{1,2,4,5}



Improve RL interpretability by combining symbolic meta-policies with neural skills.



Motivation

RGB Image		Object Centric	
Neural	Hierarchical	Decision Tree	Logic Rule Set
✗ black box	✗ black box	✓ transparent	✓ transparent
✗ monolithic	✓ modular	✗ monolithic	✗ monolithic
✗ learns shortcuts	✗ entangled skills	✗ excessive size	✗ excessive size

- Traditional RL approaches often exhibit **misaligned behavior** that is **difficult to identify or correct** without interpretability.^[1,2]
- But: Existing interpretable RL methods deployed on atomic actions quickly become **excessively complex**.
- Instead: Employ interpretable RL on **abstract skills** based on **object-centric** input and learn skills with **LLM-generated rewards**.

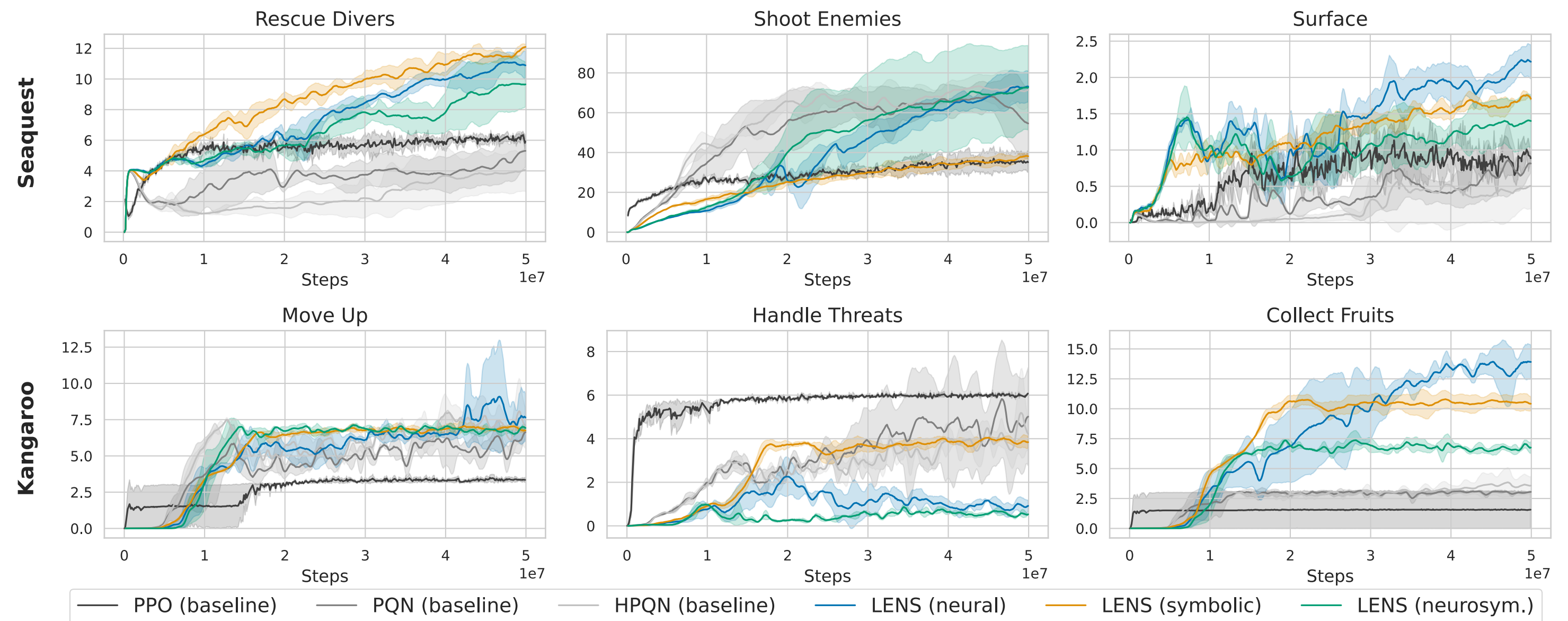
Logically Enhanced Neural Skills

- (1) Extract Objects + Attributes** : First, transform image state into object-centric (OC) state using existing methods (e.g. [3, 4, 5]).
- (2) Filter Relevant Skills**: Remove skills not applicable in the current situation using LLM-generated filters (based on OC-input).
- (3) Maximize Meta-Q-Values**: Select the neural skill to be executed by maximizing the Q-values of the meta policy with remaining skills.
- (4) Execute Neural Skill**: Obtain the next action by maximizing the skill-specific Q-value function.
- Three meta-policy variations: **neural**, **symbolic**, **neuro-symbolic**

```
def meta_policy(st: state):  
    if enemy_close(st.enemies,  
st.player):  
        return fight_enemies()  
    elif is_available(st.divers):  
        return rescue_divers()  
    elif is_low(st.oxygen):  
        return surface()  
    elif all_collected(st.divers):  
        return surface()  
    return rescue_divers()  
  
def meta_policy_rules(st: state):  
    fight_enemies = False  
    rescue_divers = True  
    surface = False  
    if enemy_close(st.enemies,  
st.player):  
        fight_enemies = True  
    if oxygen_low(st.oxygen):  
        surface = True  
    return [fight_enemies,  
rescue_divers, surface]
```

Examples of a **symbolic meta-policy** (left) and rule-set for the **neuro-symbolic meta-policy** (right) in Seaquest.

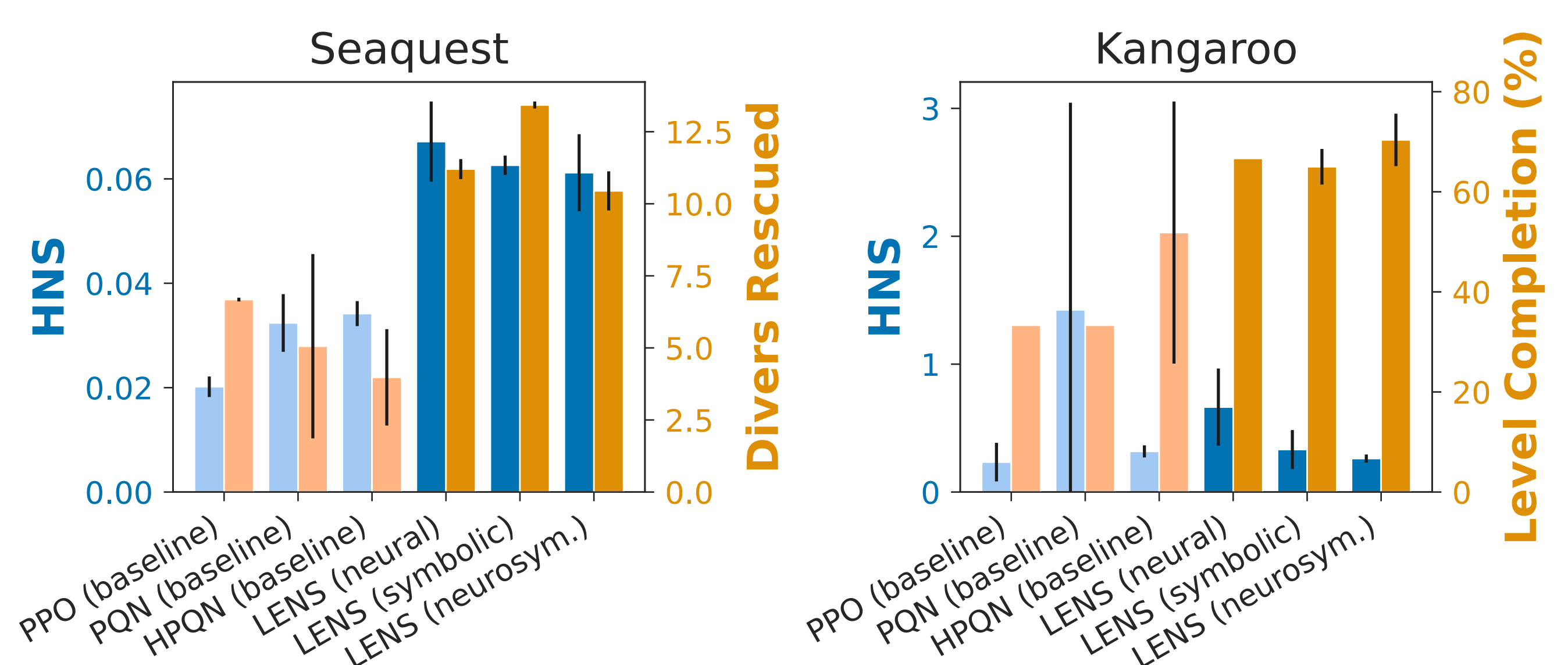
Results



1. LENS learns disentangled Skills jointly from off-policy data.

	enemy_closeby fight_enemies: 156.10 diver_available rescue_divers: 151.42 oxygen_low surface: 0.0 all_collected
	enemy_closeby fight_enemies: 273.32 diver_available rescue_divers: 277.47 oxygen_low surface: 0.0 all_collected
	enemy_closeby fight_enemies: 395.02 diver_available rescue_divers: 0.0 oxygen_low surface: 405.52 all_collected

2. LENS produces interpretable yet flexible high-level plans. Decisions in ambiguous situations are made intuitively.



3. LENS is competitive and better aligned to actual environment goals.

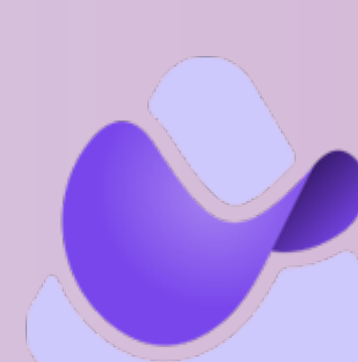
- [1] Rudin. "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead." (2019)
[2] Delfosse et al. "Interpretable Concept Bottlenecks to Align Reinforcement Learning Agents" (2024)
[3] Li et al. "Object-sensitive Deep Reinforcement Learning" (2017)
[4] Locatello et al. "Object-Centric Learning with Slot Attention." (2020)
[5] Delfosse et al. "Boosting Object Representation Learning via Motion and Object Continuity" (2023)



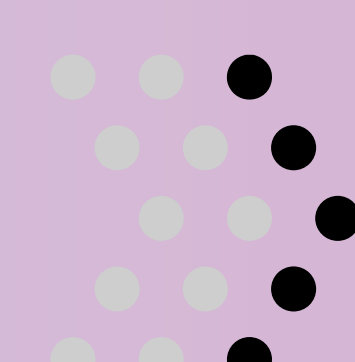
Raban Emunds Jannis Blüml Quentin Delfosse



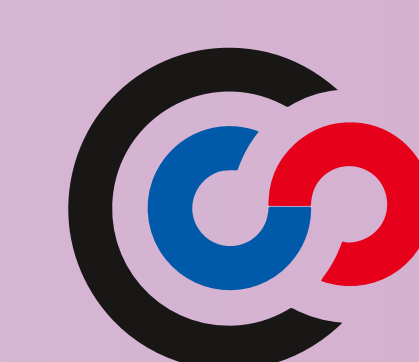
¹AIML Lab
TU Darmstadt



²hessian.AI



³ATHENE



⁴TUDa Centre for
Cognitive Science



⁵German Research
Center for AI