# Deep RL agents perform hidden Shortcut Reinforcement Learning. RL agents must rely on human understandable concepts.

Interpretable Concept Bottlenecks
to Align Reinforcement Learning Agents

Quentin Delfosse
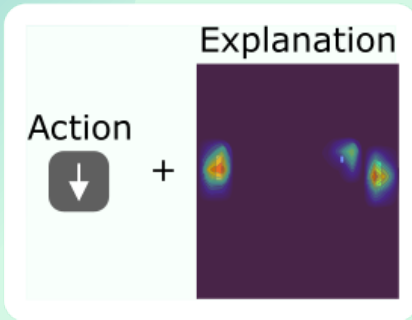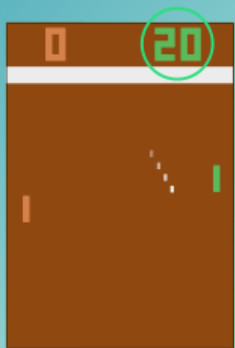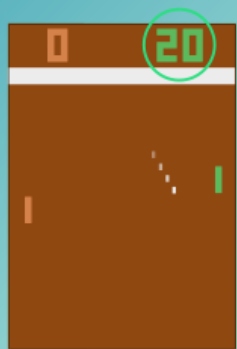
Sebastian Sztwiertnia

Wolfgang Stammer

Kristian Kersting

# Can we interpret deep agents?



Training

Explanation

Action +

# Can we interpret deep agents?



Training

Explanation
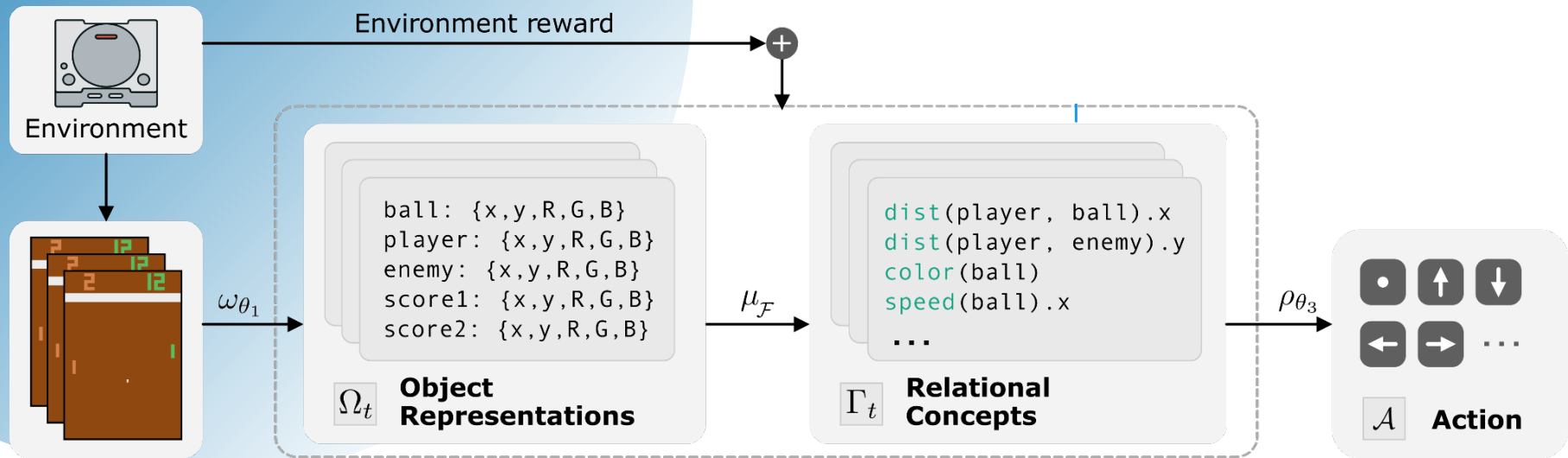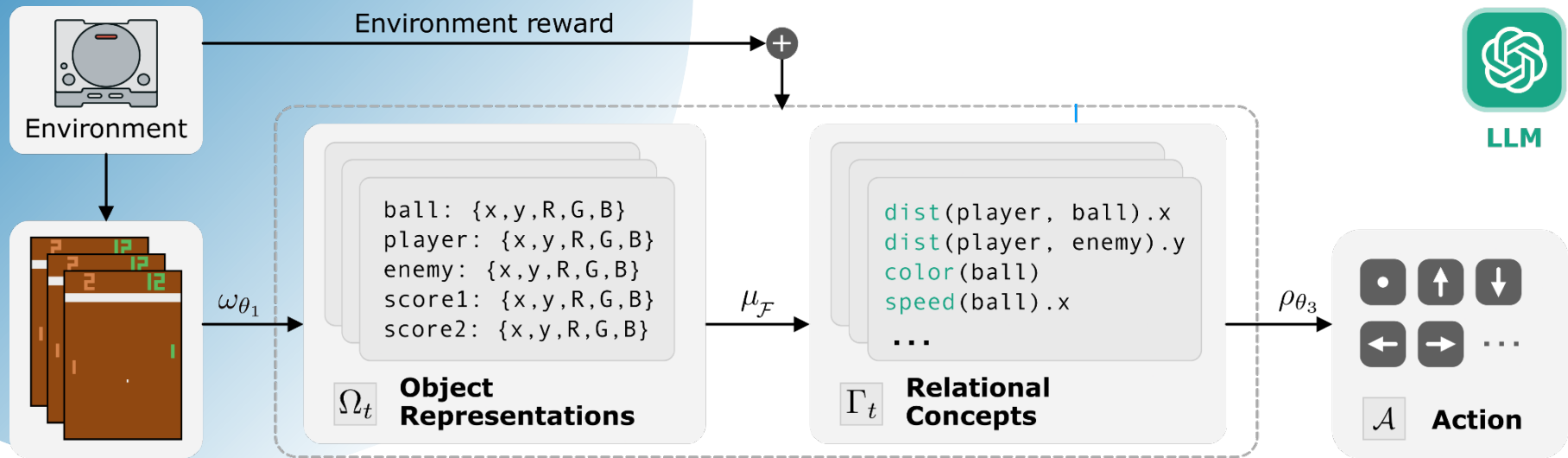
Action +

Evaluation

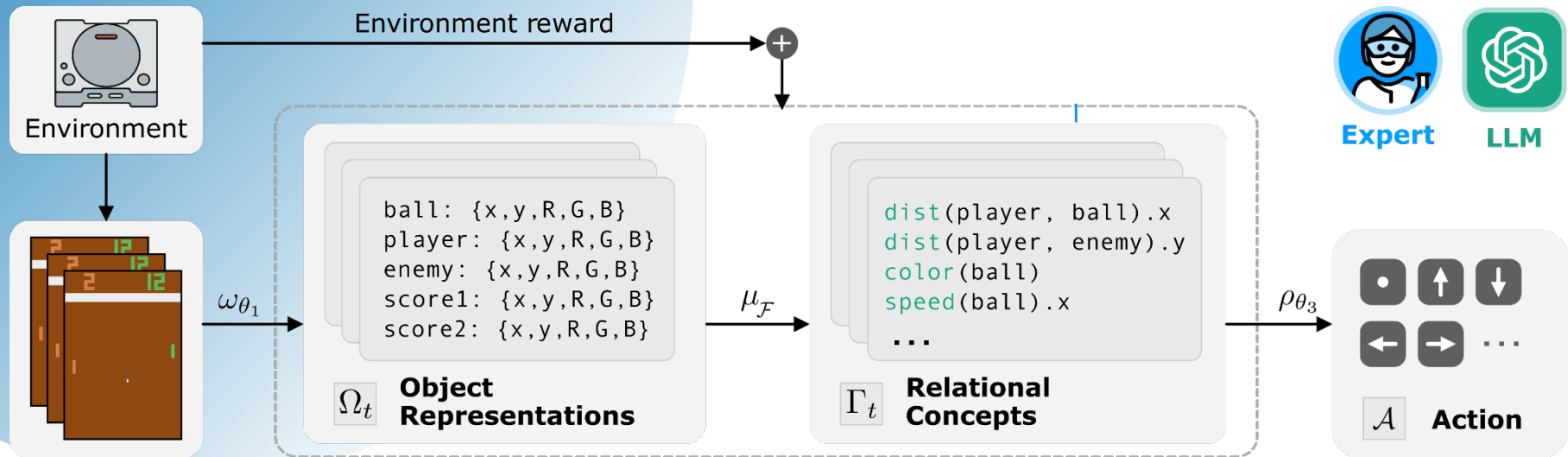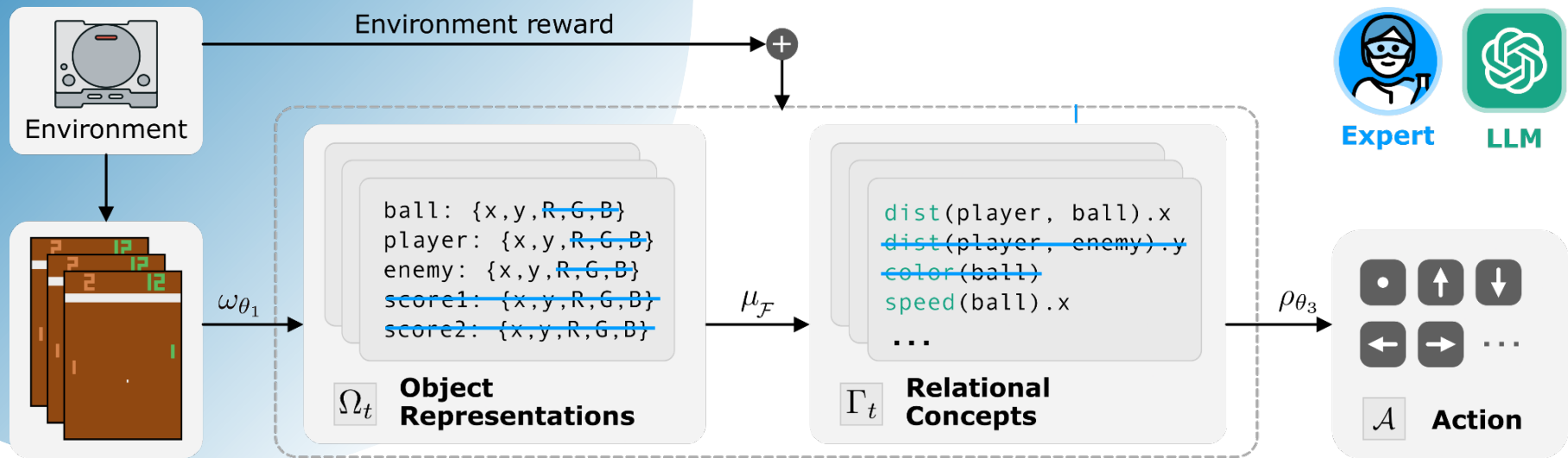# Successive Concept Bottlenecks



SPACE: Lin et al. "SPACE: Unsupervised Object Scene Representation via Spatial Attention and Decomposition." (2020)
SPOC: Delfosse et al. "Boosting object representation learning via motion and object continuity." (2023).

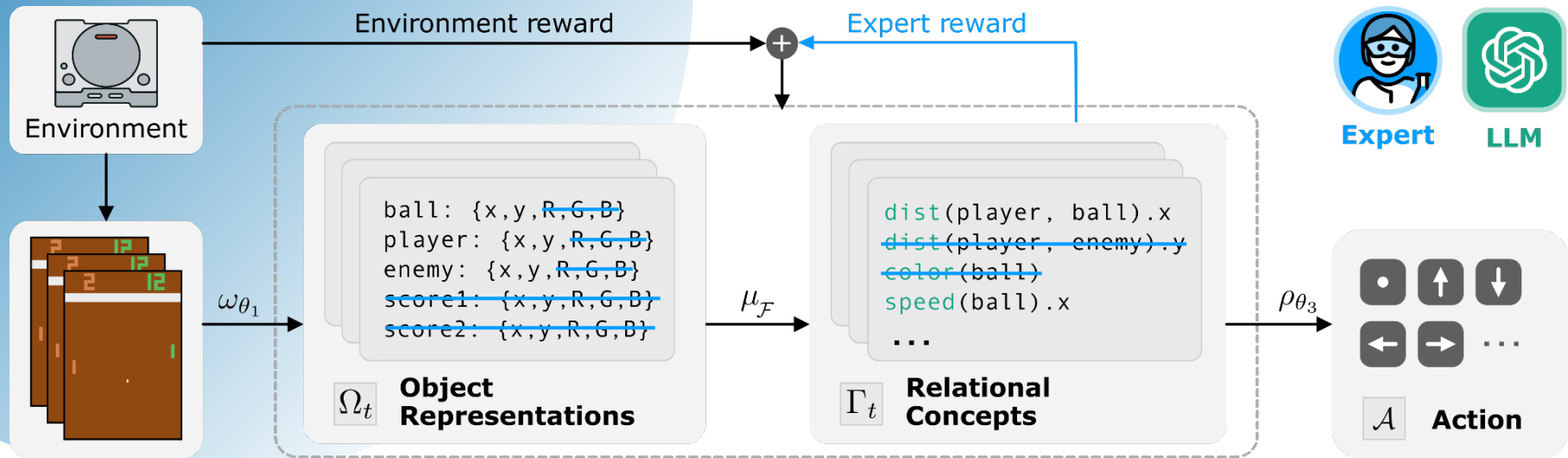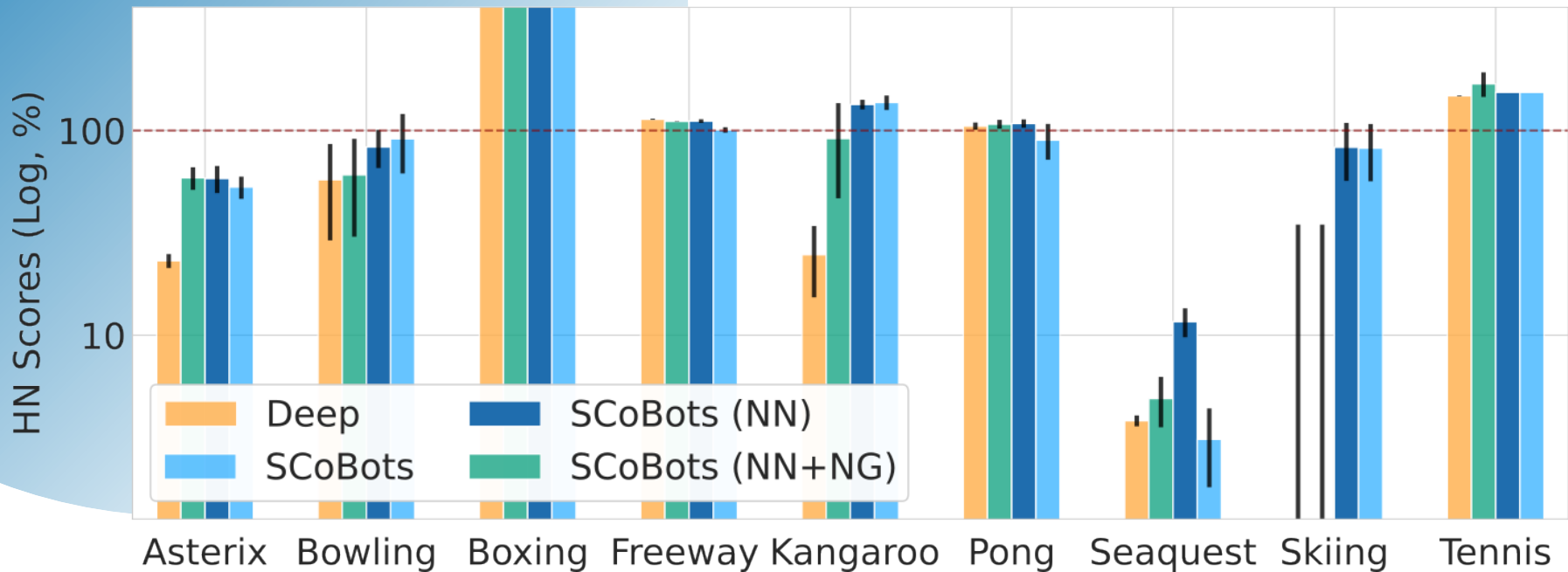# Successive Concept Bottlenecks



SPACE: Lin et al. "SPACE: Unsupervised Object Scene Representation via Spatial Attention and Decomposition." (2020)
SPOC: Delfosse et al. "Boosting object representation learning via motion and object continuity." (2023).

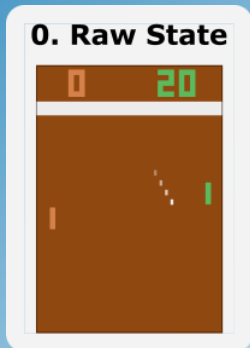# Successive Concept Bottlenecks



SPACE: Lin et al. "SPACE: Unsupervised Object Scene Representation via Spatial Attention and Decomposition." (2020)
SPOC: Delfosse et al. "Boosting object representation learning via motion and object continuity." (2023).

# Successive Concept Bottlenecks



SPACE: Lin et al. "SPACE: Unsupervised Object Scene Representation via Spatial Attention and Decomposition." (2020)
SPOC: Delfosse et al. "Boosting object representation learning via motion and object continuity." (2023).

# Successive Concept Bottlenecks



SPACE: Lin et al. "SPACE: Unsupervised Object Scene Representation via Spatial Attention and Decomposition." (2020)
SPOC: Delfosse et al. "Boosting object representation learning via motion and object continuity." (2023).

# SCoBots are competitive



Performances

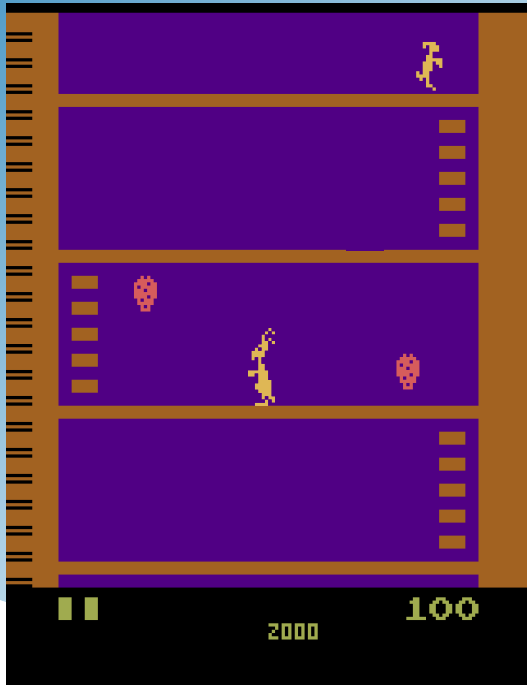# SCoBots' can be realigned!

# SCoBots' can be realigned!
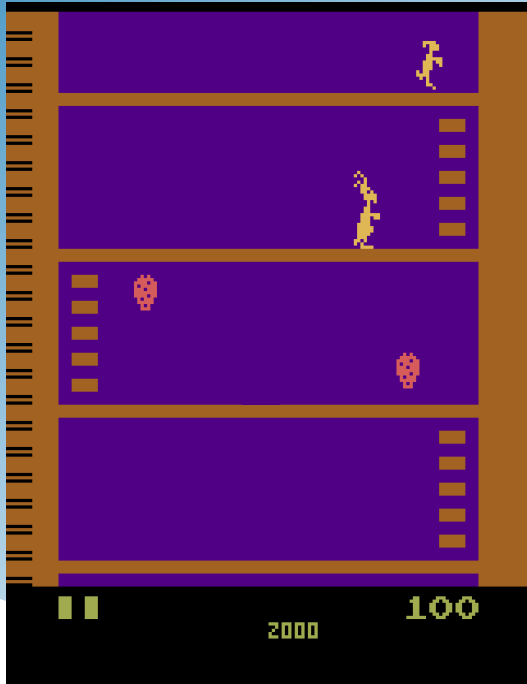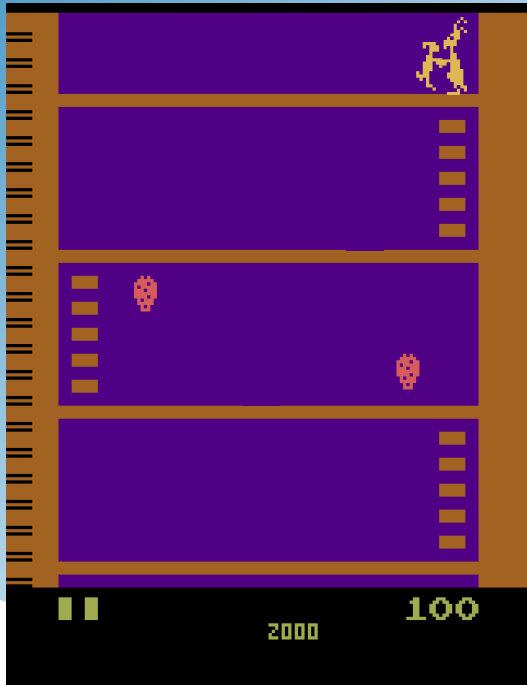
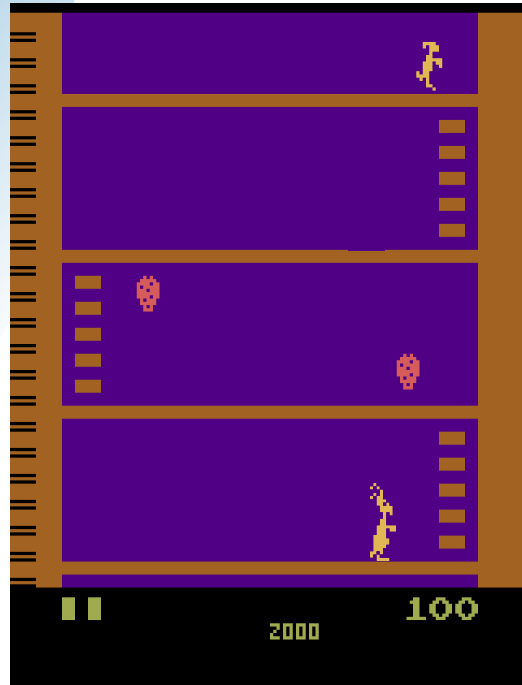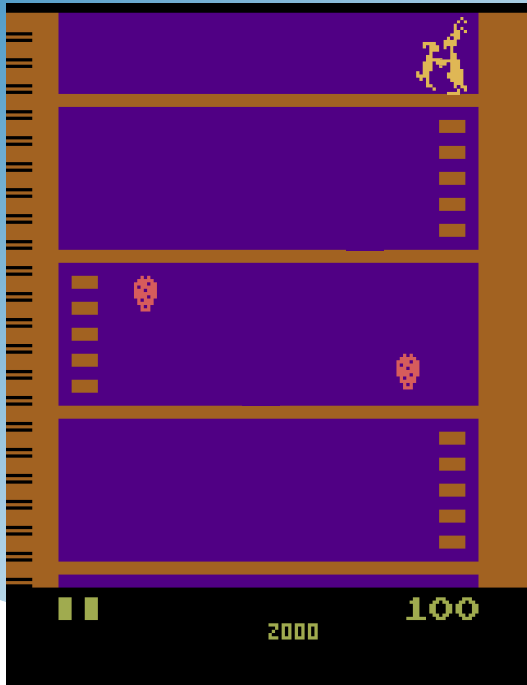# SCoBots' can be realigned!

# SCoBots can be redirected!
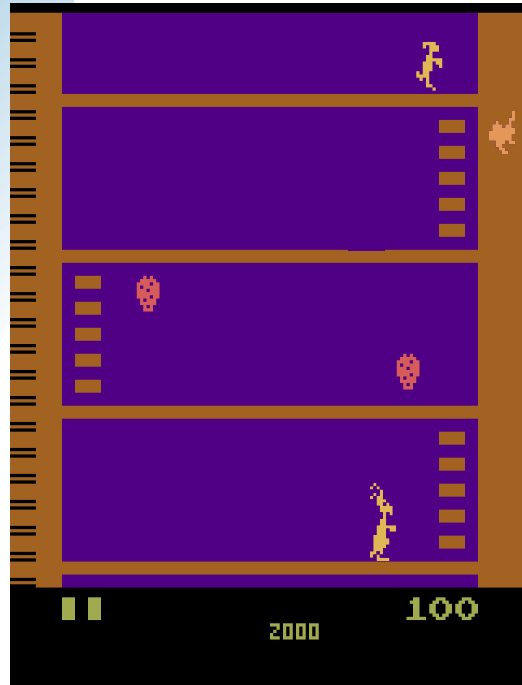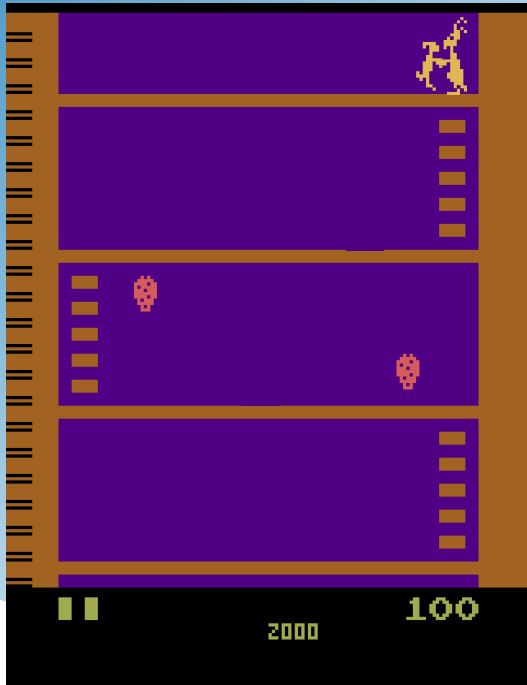
# SCoBots can be redirected!

# SCoBots can be redirected!

# SCoBots can be redirected!

# SCoBots can be redirected!

# SCoBots can be redirected!

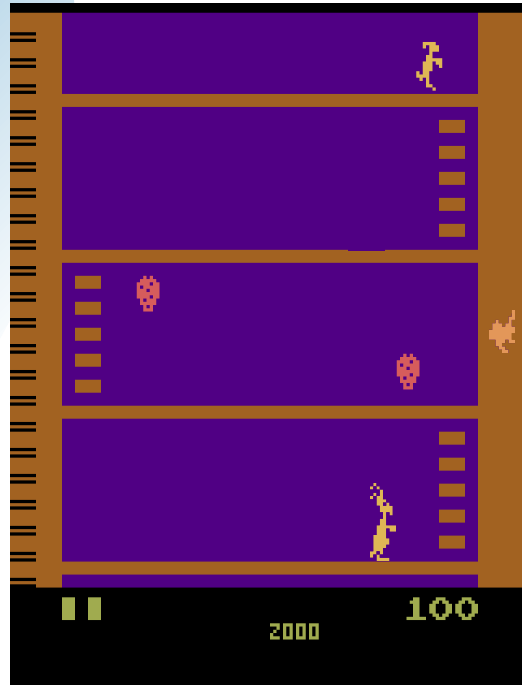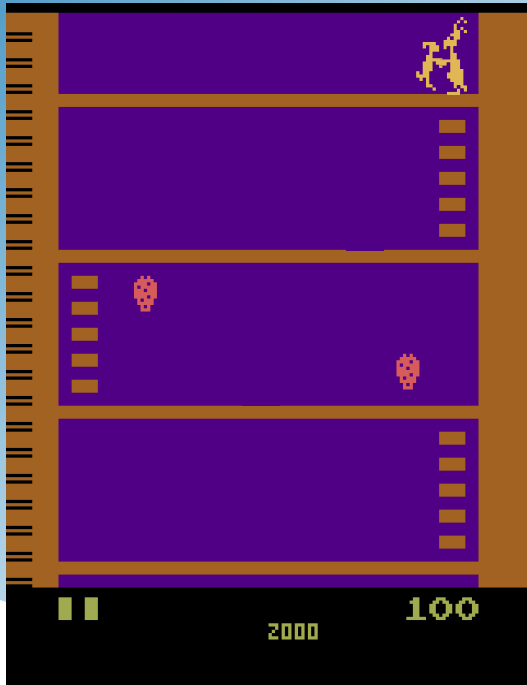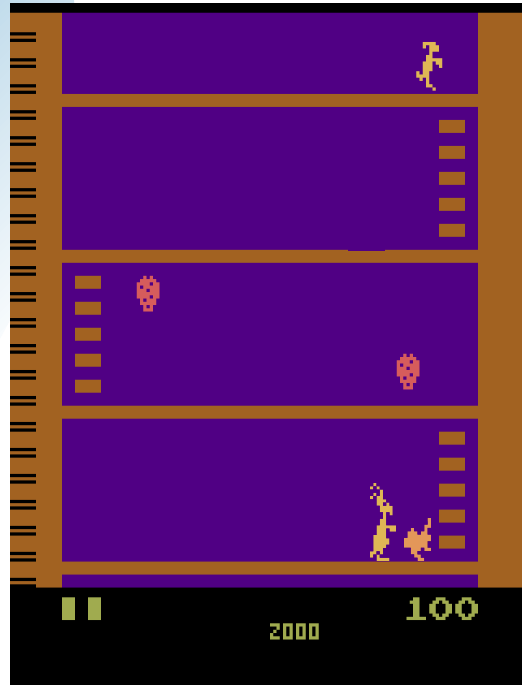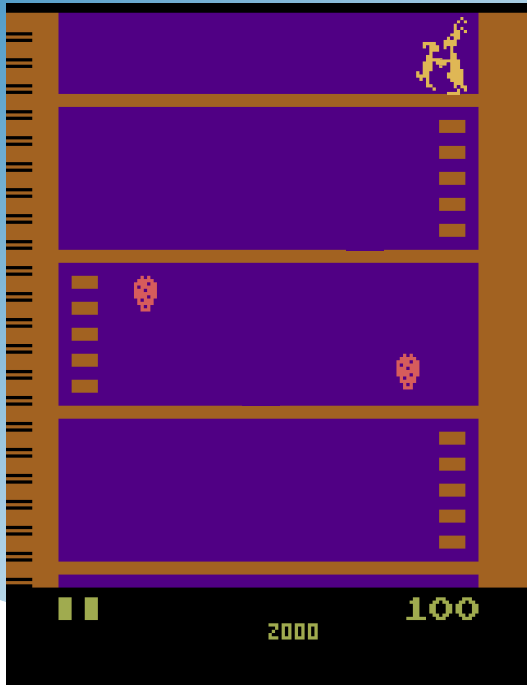# SCoBots can be redirected!

# SCoBots can be redirected!

# SCoBots can be redirected!

# SCoBots can be redirected!



Ill-defined reward
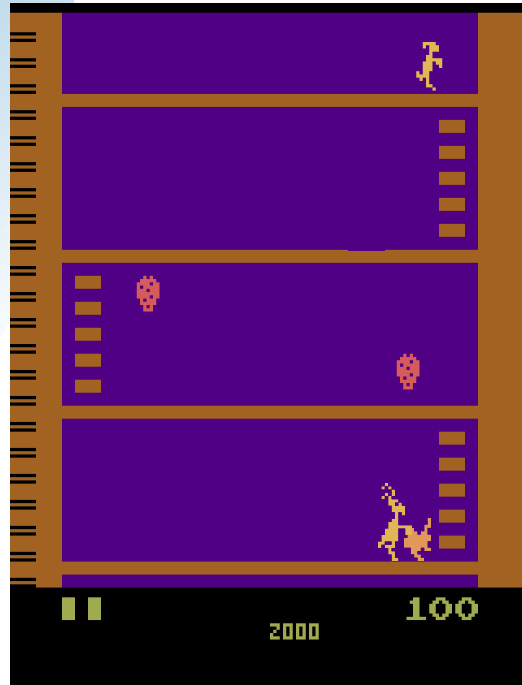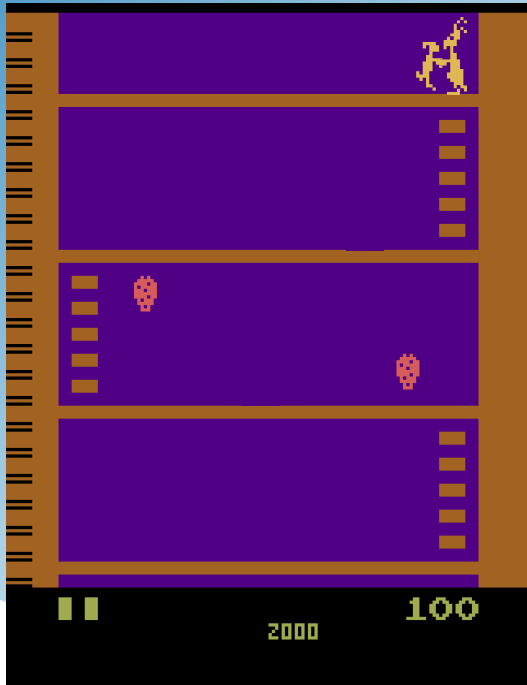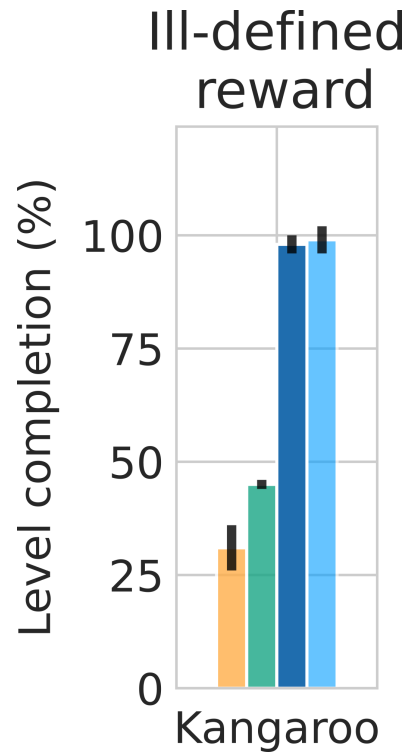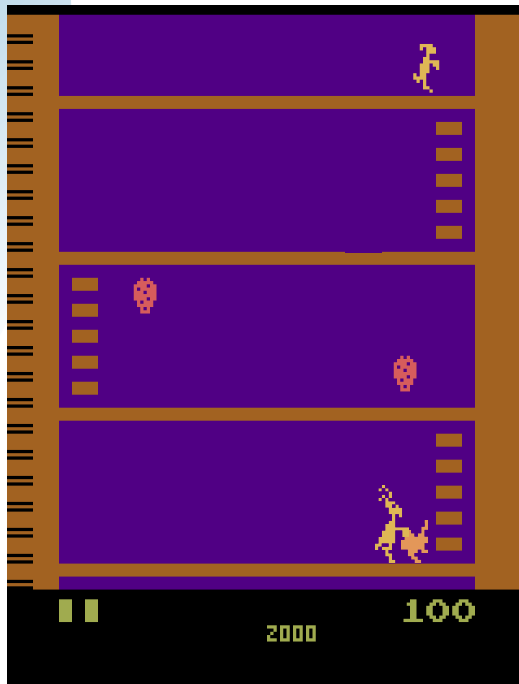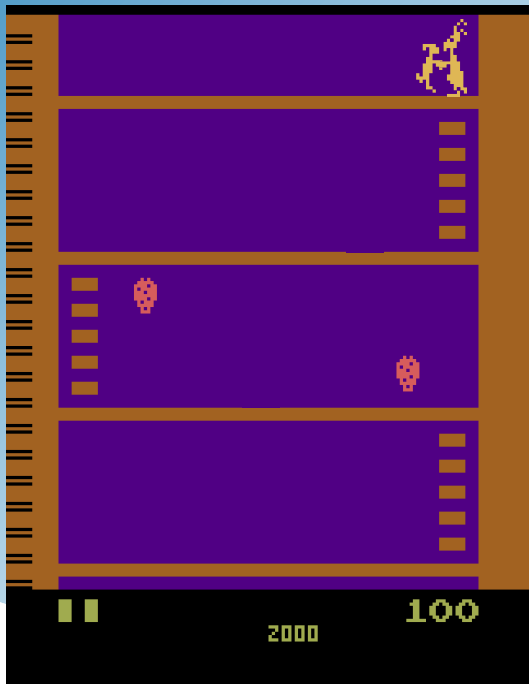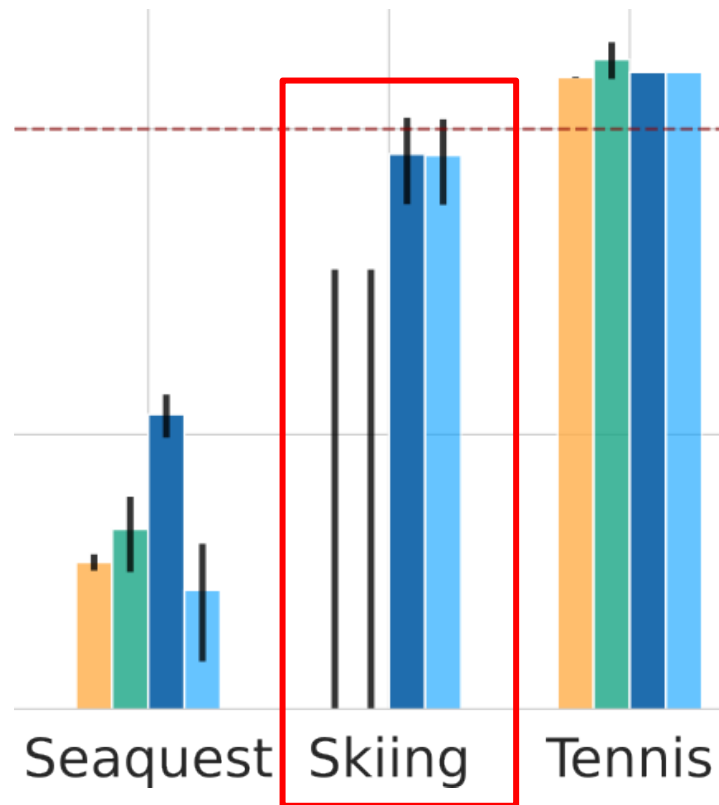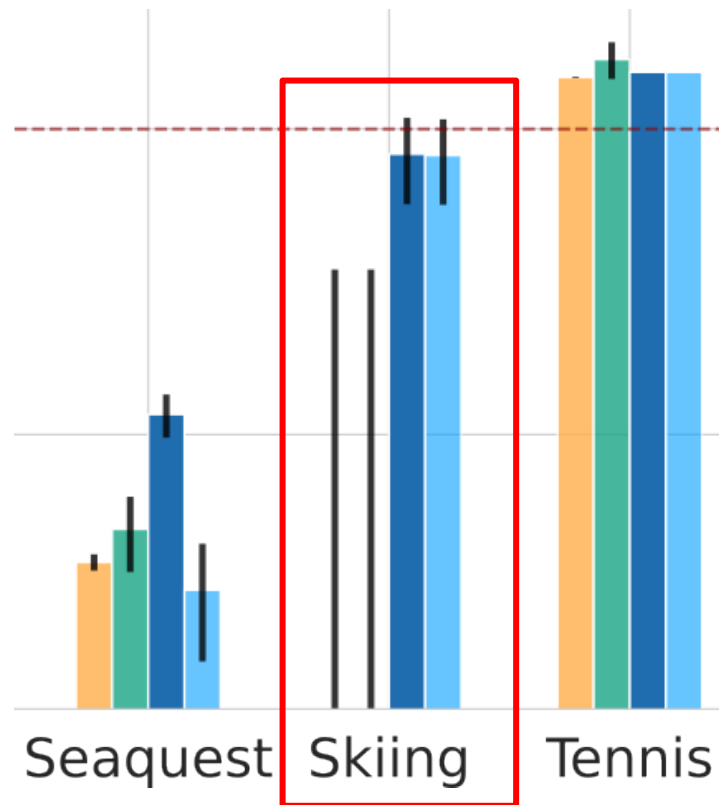
Level completion (%)

Kangaroo

Deep
SCoBots
SCoBots (NN)
SCoBots (NN+NG)

# SCoBots can be redirected

# SCoBots can be redirected

# Deep RL agents perform hidden Shortcut Reinforcement Learning. RL agents must rely on human understandable concepts.

Open source: [github.com/k4ntz/SCoBots](github.com/k4ntz/SCoBots)