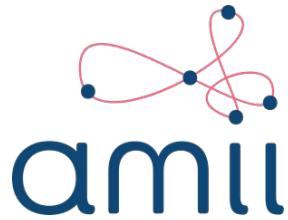


Representation-driven Option Discovery in Reinforcement Learning

Marlos C. Machado



Sharif University of Technology, June 6, 2025



This is a research program

that touches upon a lot
of what you've seen in the other guest lectures



RL is now commonly deployed in the real-world



RL is now commonly deployed in the real-world



Video compression
[Mandhane et al., 2022]

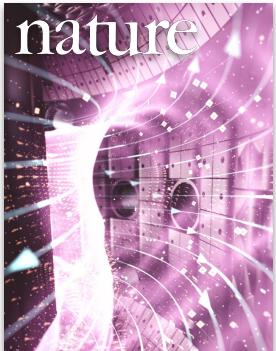
Matrix multiplication
[Fawzi et al., 2022]

Hardware design
[Mirhoseini et al., 2021]

Cooling systems
[Luo et al., 2022]

Thermal power generators
[Zhan et al., 2022]

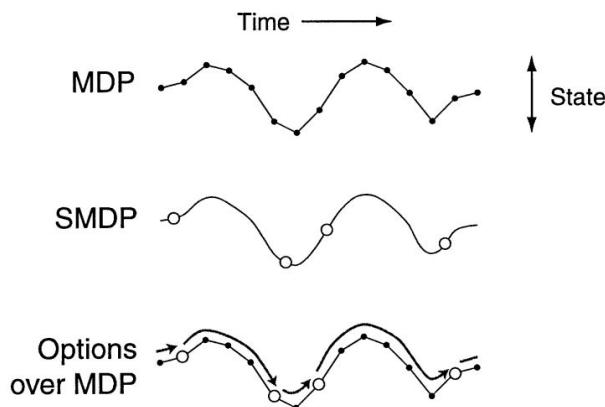
Managing inventories
[Madekaet al., 2022]



...

Temporal abstraction – Options

[Sutton, Precup, & Singh, AIJ 1999]

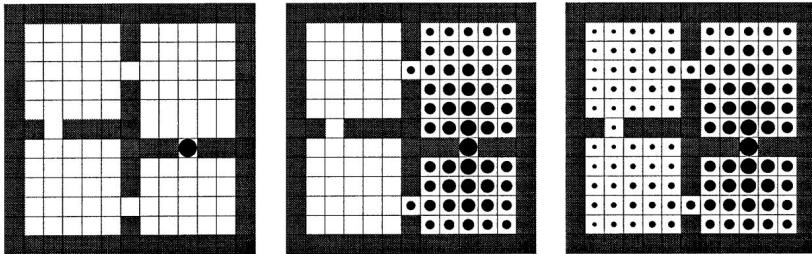


$$\omega = \langle \mathcal{I}_\omega, \pi_\omega, \beta_\omega \rangle \quad (1)$$

initiation set
policy
termination condition

The many use cases of options

Faster credit assignment / planning:



[Sutton, Precup, & Singh. Artif. Intelligence 1999]

Exploration [Jong, Hester, & Stone, AAMAS 2008]

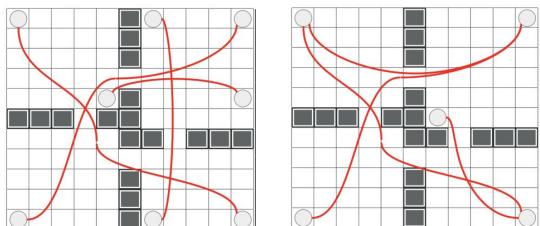
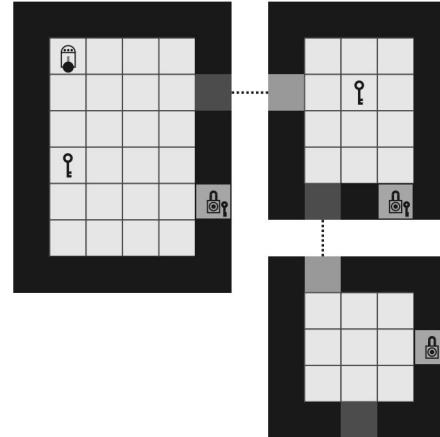


Figure by Jinnai et al. (2019)

Transfer learning:

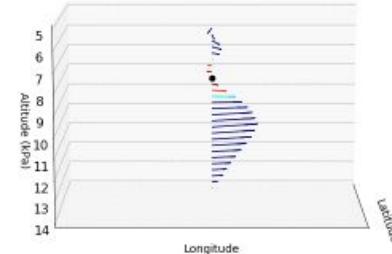
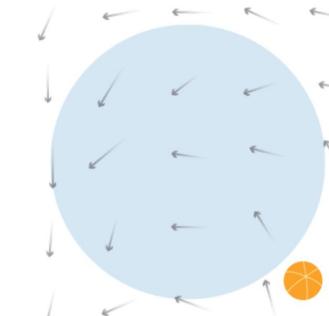
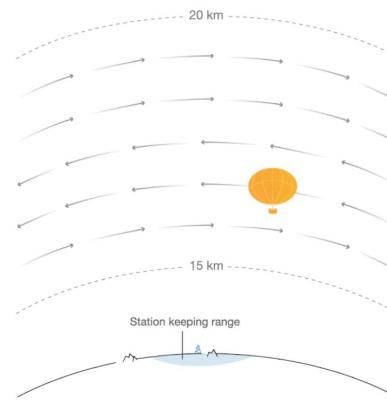
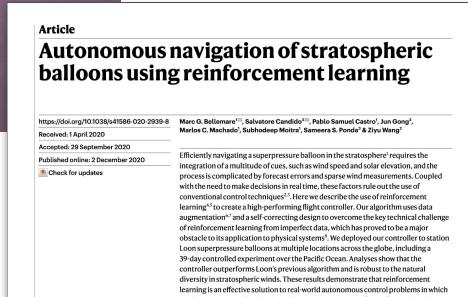


[Konidaris & Barto, IJCAI 2007]

and more...

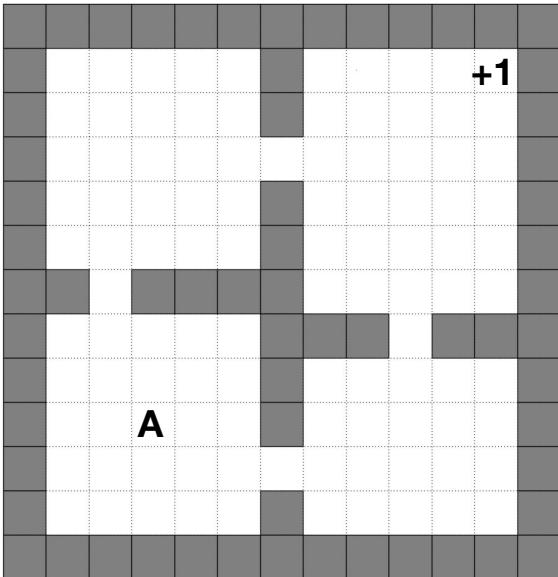
It works! RL in the real world

<https://www.nature.com/articles/s41586-020-2939-8>

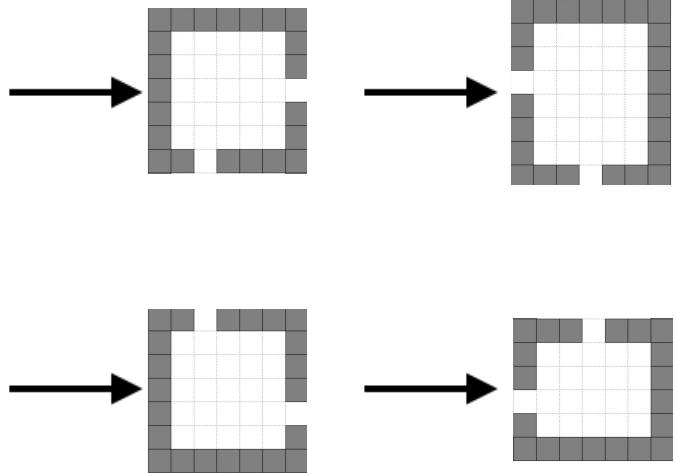
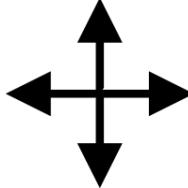


[Bellemare, Candido, Castro, Gong, Machado, Moitra, Ponda, & Wang, Nature 2020]

Exploring at a higher level of abstraction



“actions”:



But where do options come from?

Temporal abstraction – Options

[Sutton, Precup, & Singh, AIJ 1999]

[Sutton et al., AAMAS 2011]

$$v_{\pi, \beta}^{c,z}(s) \doteq \mathbb{E}_{\pi, \beta} \left[\sum_{j=1}^K \gamma^{j-1} c(S_j) + \gamma^{K-1} z(S_K) \mid S_0 = s \right], \quad \forall s \in \mathcal{S}$$

Temporal abstraction – Options

[Sutton, Precup, & Singh, AIJ 1999]

[Sutton et al., AAMAS 2011]

subtask (problem): maximize discounted sum of cumulants plus a stopping value

$$v_{\pi, \beta}^{c, z}(s) \doteq \mathbb{E}_{\pi, \beta} \left[\sum_{j=1}^K \gamma^{j-1} \underline{c(S_j)} + \gamma^{K-1} \underline{z(S_K)} \mid S_0 = s \right], \quad \forall s \in \mathcal{S}$$

Temporal abstraction – Options

[Sutton, Precup, & Singh, AIJ 1999]

[Sutton et al., AAMAS 2011]

$$v_{\pi, \beta}^{c,z}(s) \doteq \mathbb{E}_{\pi, \beta} \left[\sum_{j=1}^K \gamma^{j-1} c(S_j) + \gamma^{K-1} z(S_K) \mid S_0 = s \right], \quad \forall s \in \mathcal{S}$$

random variable

option (solution): policy and stopping function

Temporal abstraction – Options

[Sutton, Precup, & Singh, AIJ 1999]
 [Sutton et al., AAMAS 2011]

$$v_{\pi, \beta}^{c,z}(s) \doteq \mathbb{E}_{\pi, \beta} \left[\sum_{j=1}^K \gamma^{j-1} c(S_j) + \gamma^{K-1} z(S_K) \mid S_0 = s \right], \quad \forall s \in \mathcal{S}$$

random variable

subtask (problem): maximize discounted sum of cumulants plus a stopping value

option (solution): policy and stopping function

Defining the option discovery problem

$$v_{\pi,\beta}^{c,z}(s) \doteq \mathbb{E}_{\pi,\beta} \left[\sum_{j=1}^K \gamma^{j-1} c(S_j) + \gamma^{K-1} z(S_K) \mid S_0 = s \right], \quad \forall s \in \mathcal{S}$$

Specify subtask:

- c : signal to maximize
- z : stopping-value function

Defining the option discovery problem

$$v_{\pi,\beta}^{c,z}(s) \doteq \mathbb{E}_{\pi,\beta} \left[\sum_{j=1}^K \gamma^{j-1} c(S_j) + \gamma^{K-1} z(S_K) \mid S_0 = s \right], \quad \forall s \in \mathcal{S}$$

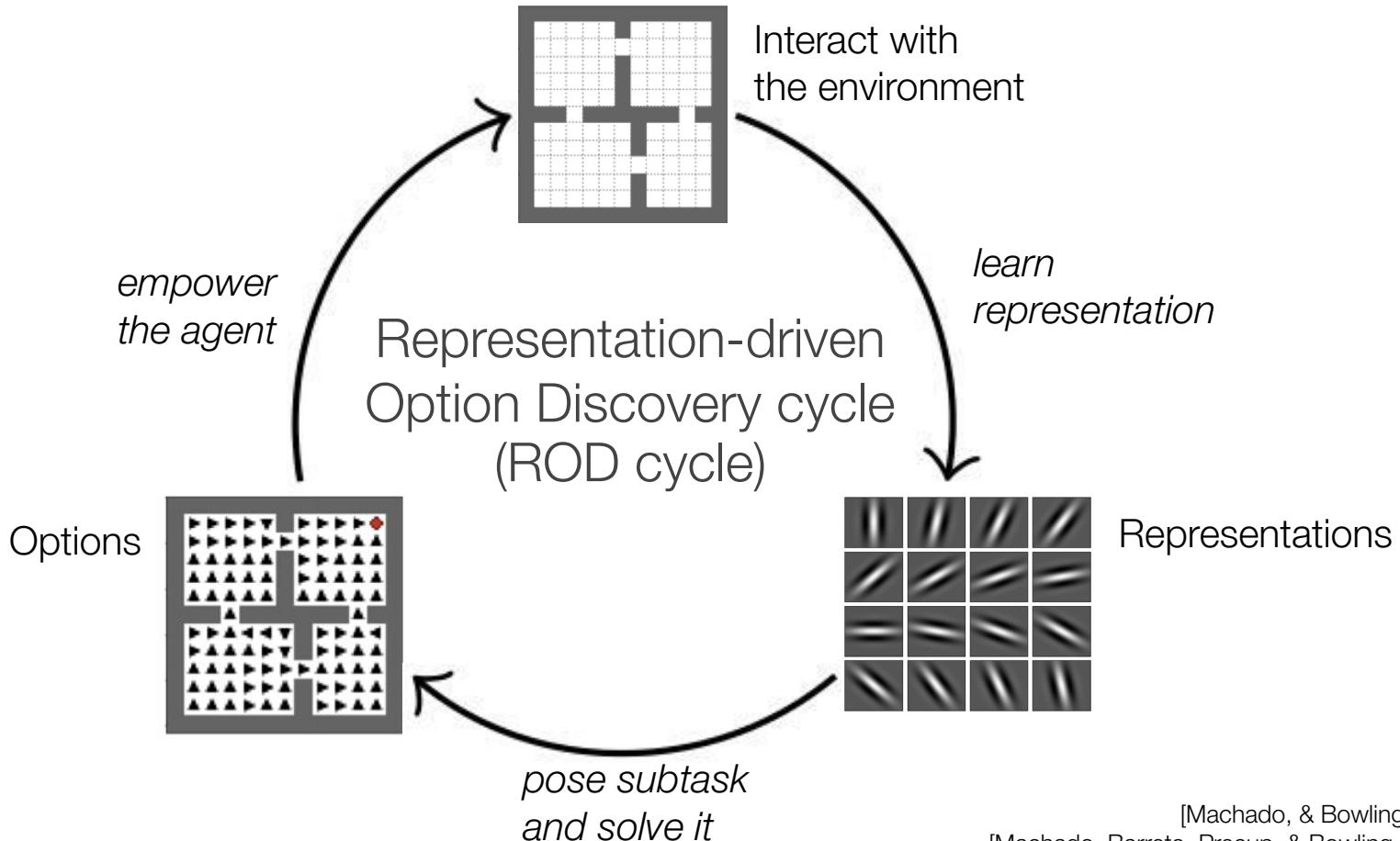
Specify subtask:

- c : signal to maximize
- z : stopping-value function

Example:
Shortest-path option to a bottleneck state

- $C_t = -1$
- $z(s) = 0$ at subgoal states or $z(s) = -\infty$ o.w.

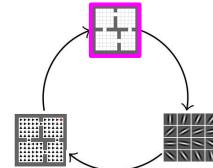
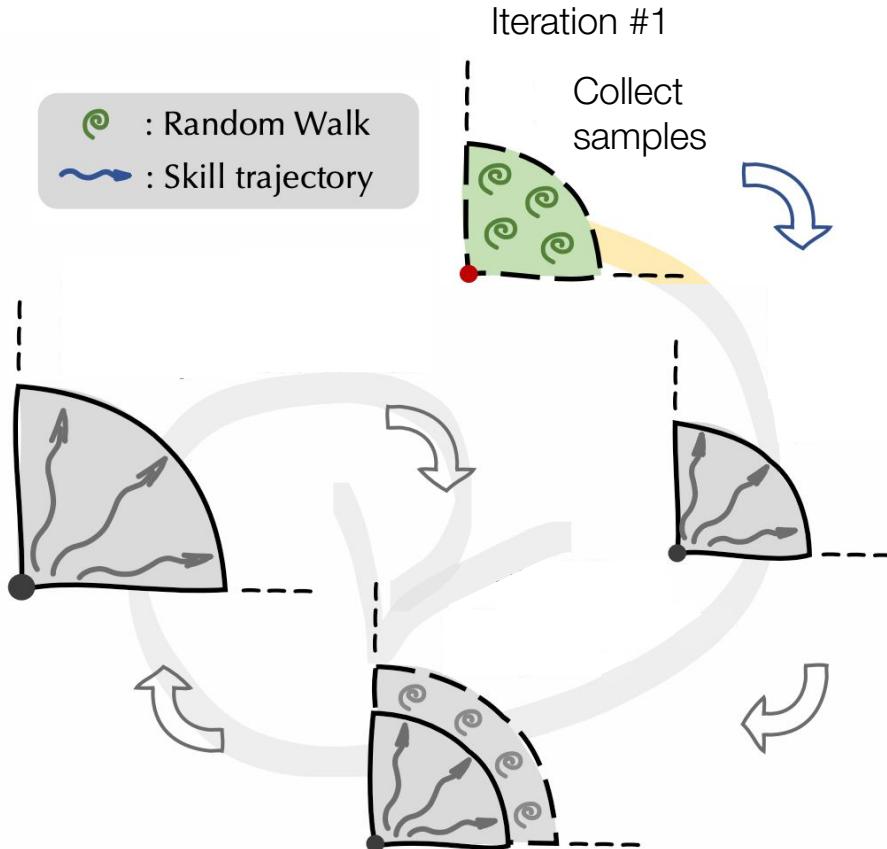
~~Where should options come from?~~
What subtasks should we use?



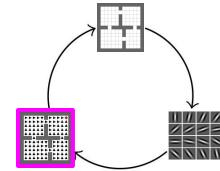
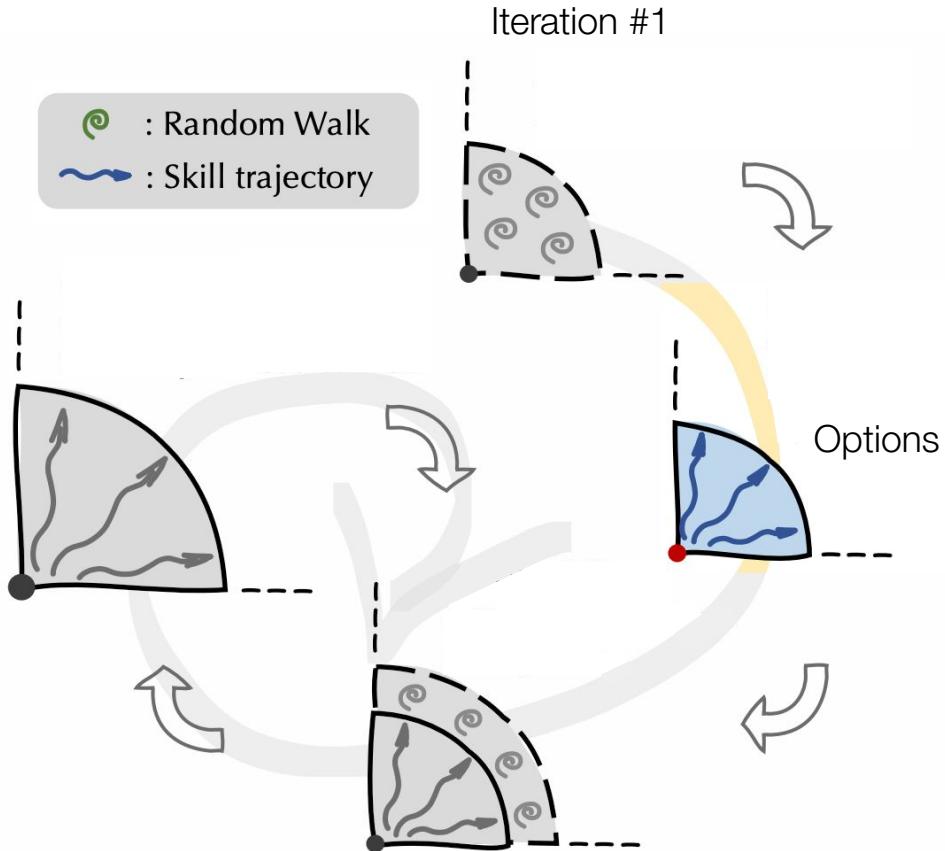
[Machado, & Bowling, arXiv 2016]

[Machado, Barreto, Precup, & Bowling, JMLR 2023]

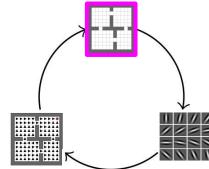
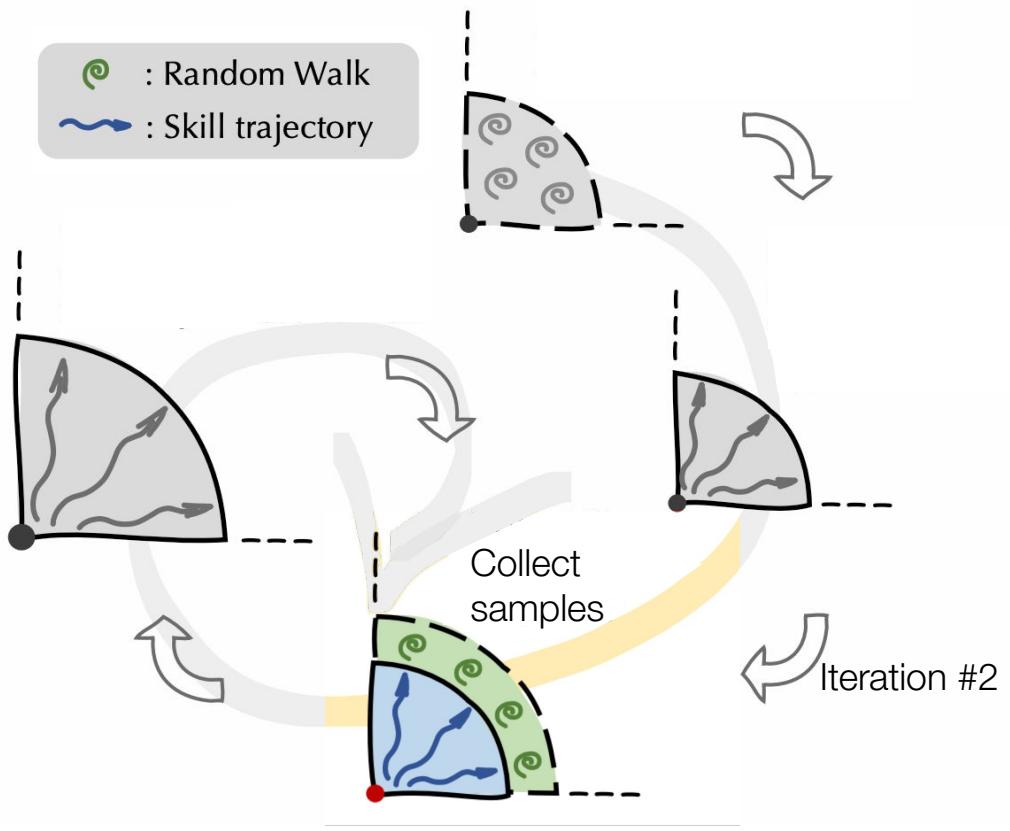
Intuition



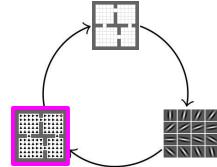
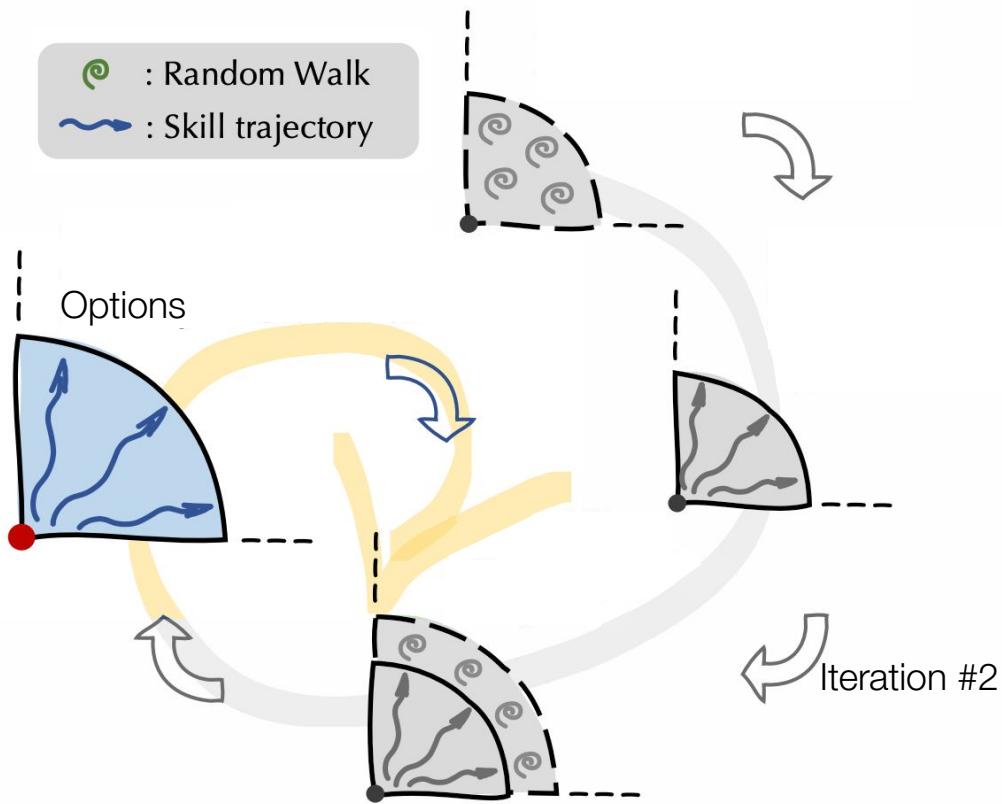
Intuition



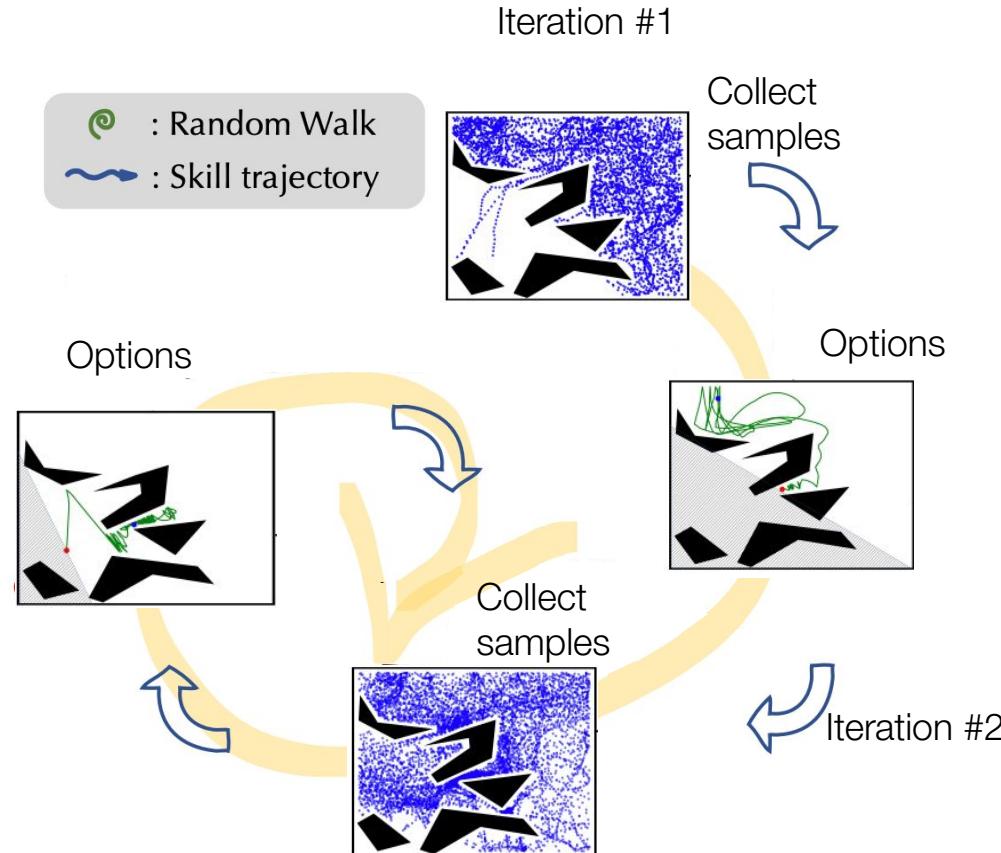
Intuition

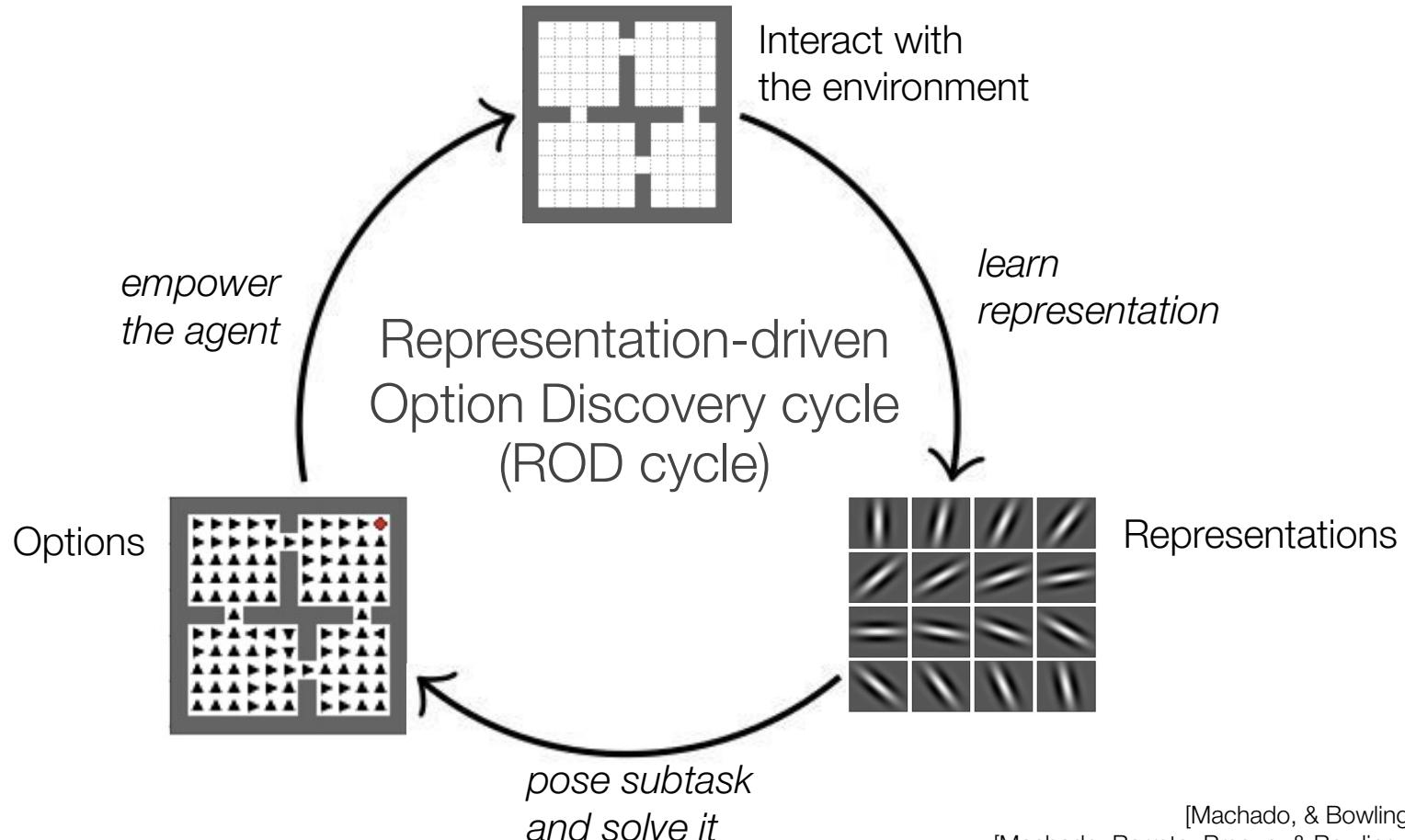


Intuition



Intuition



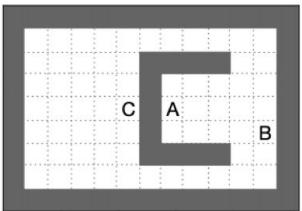
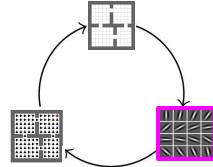


[Machado, & Bowling, arXiv 2016]

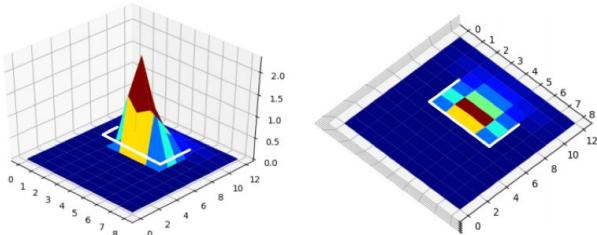
[Machado, Barreto, Precup, & Bowling, JMLR 2023]

Successor Representation

[Dayan, Neural Computation 1993]



$$\Psi_{\pi}(s, s') = \mathbb{E}_{\pi} \left[\sum_t \gamma^t \mathbf{1}_{S_t = s'} | S_0 = s \right] \quad (1)$$



$$\Psi_{\pi} = \sum_t (\gamma \mathbf{P}_{\pi})^t = (\mathbf{I} - \gamma \mathbf{P}_{\pi})^{-1} \quad (2)$$

The SR as a collection of GVF_s

- \mathbf{C}_t : indicator function for state visitation
- γ : any fixed γ , but the same across all GVF_s
- π : any policy, but the same across all GVF_s
- $\mathbf{z}(s) = 0$ for all states

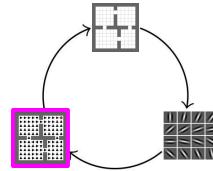
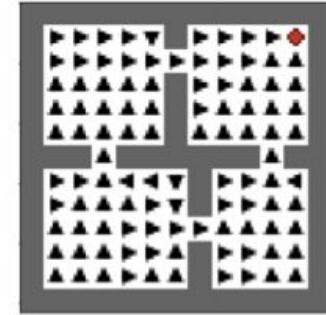
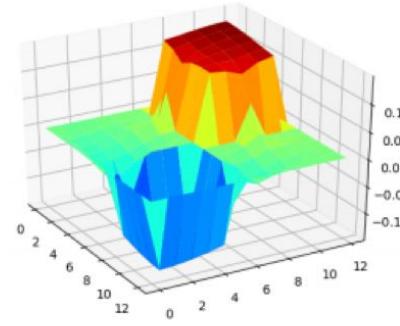
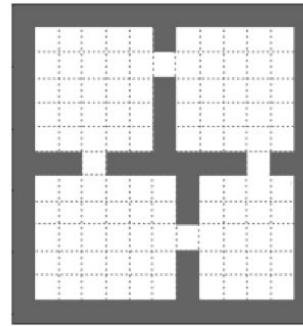
Eigenoptions

$$\Psi_\pi \mathbf{e} = \lambda \mathbf{e} \quad (1)$$

$$C_t \doteq \mathbf{e}_i^\top (\mathbf{x}(S_t) - \mathbf{x}(S_{t-1})) \quad (2)$$

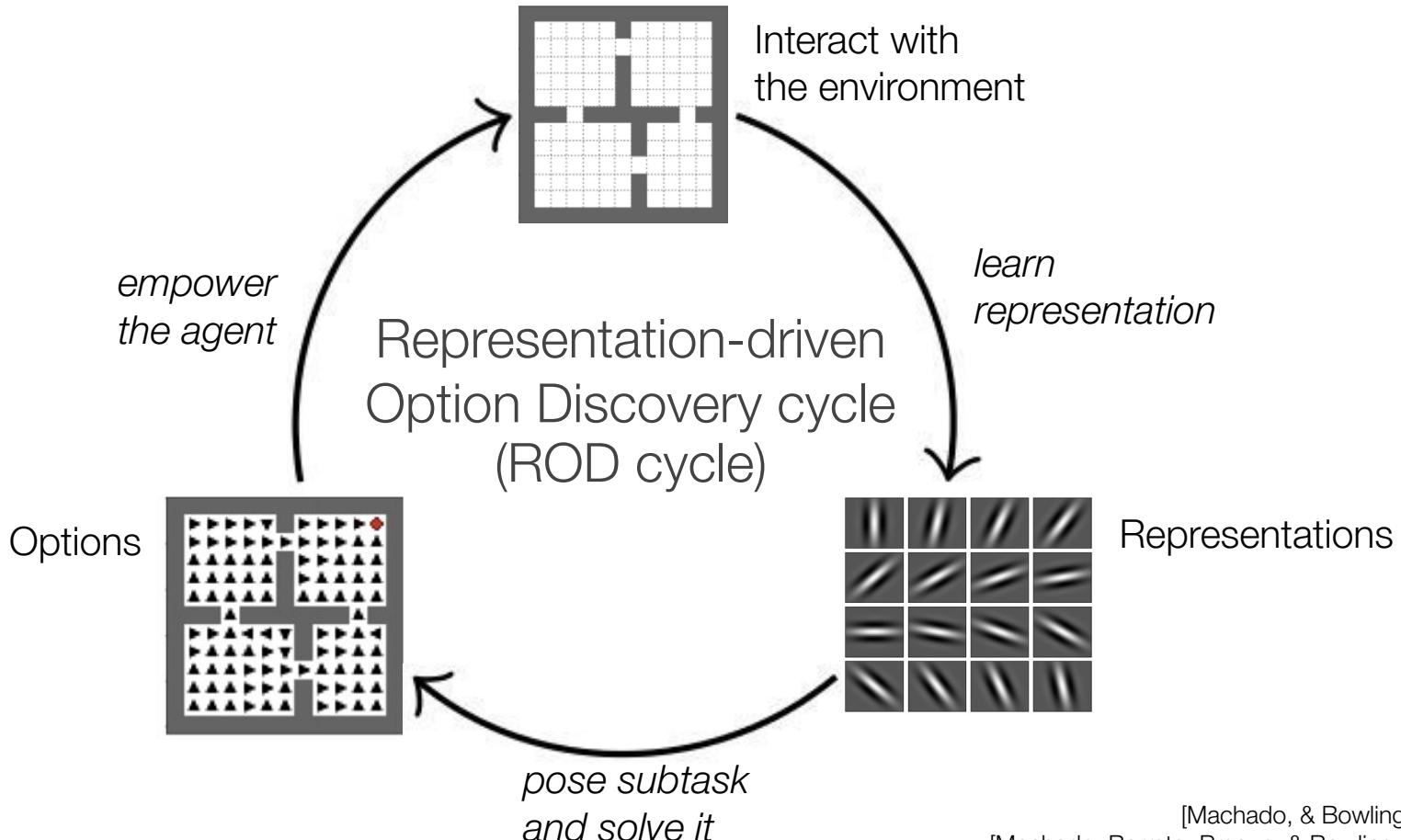
$$z(s) = 0 \quad \forall s \in \mathcal{S} \quad (3)$$

$$q_\pi(\cdot, \perp) = 0 \quad \forall \pi \quad (4)$$



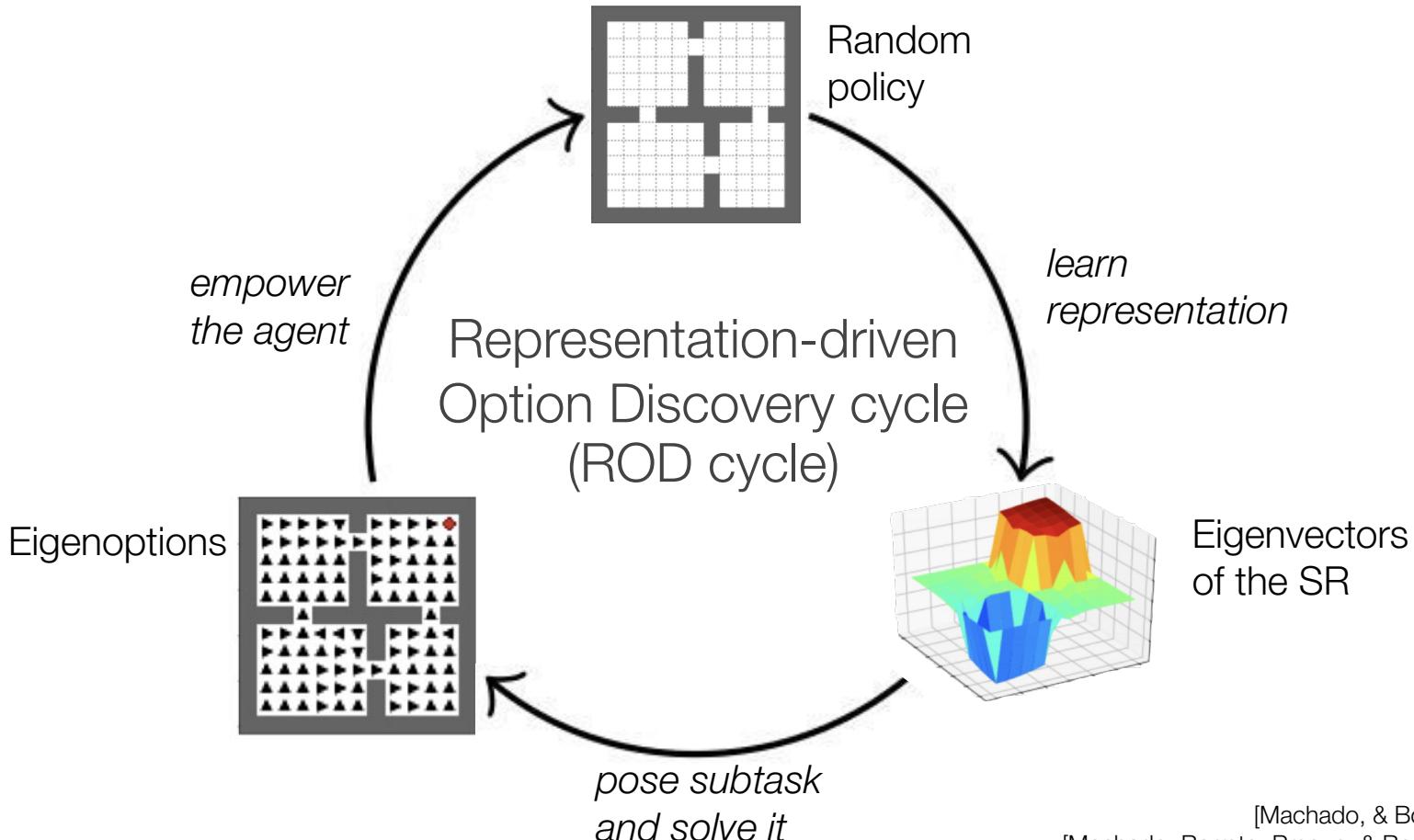
[Machado, Bellemare, & Bowling, ICML 2017]

[Machado, Rosenbaum, Guo, Liu, Tesauro, & Campbell, ICLR 2018]



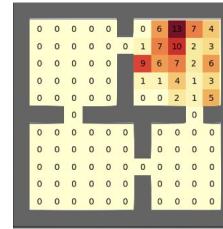
[Machado, & Bowling, arXiv 2016]

[Machado, Barreto, Precup, & Bowling, JMLR 2023]



An example

ROD cycle Iteration 1

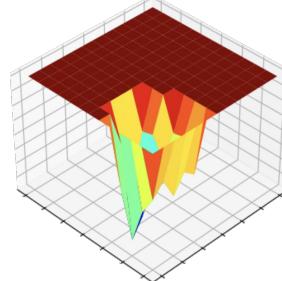


Trajectory

ROD cycle Iteration 1

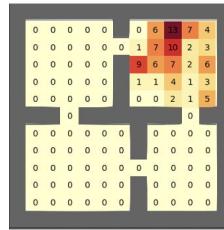


Trajectory



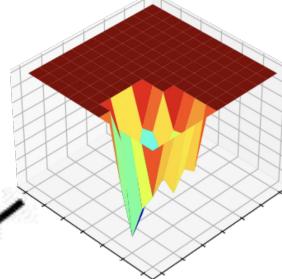
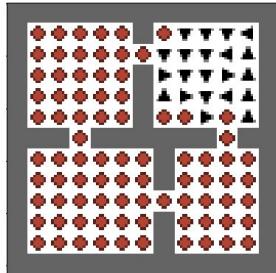
Eigenvector of the SR

ROD cycle Iteration 1



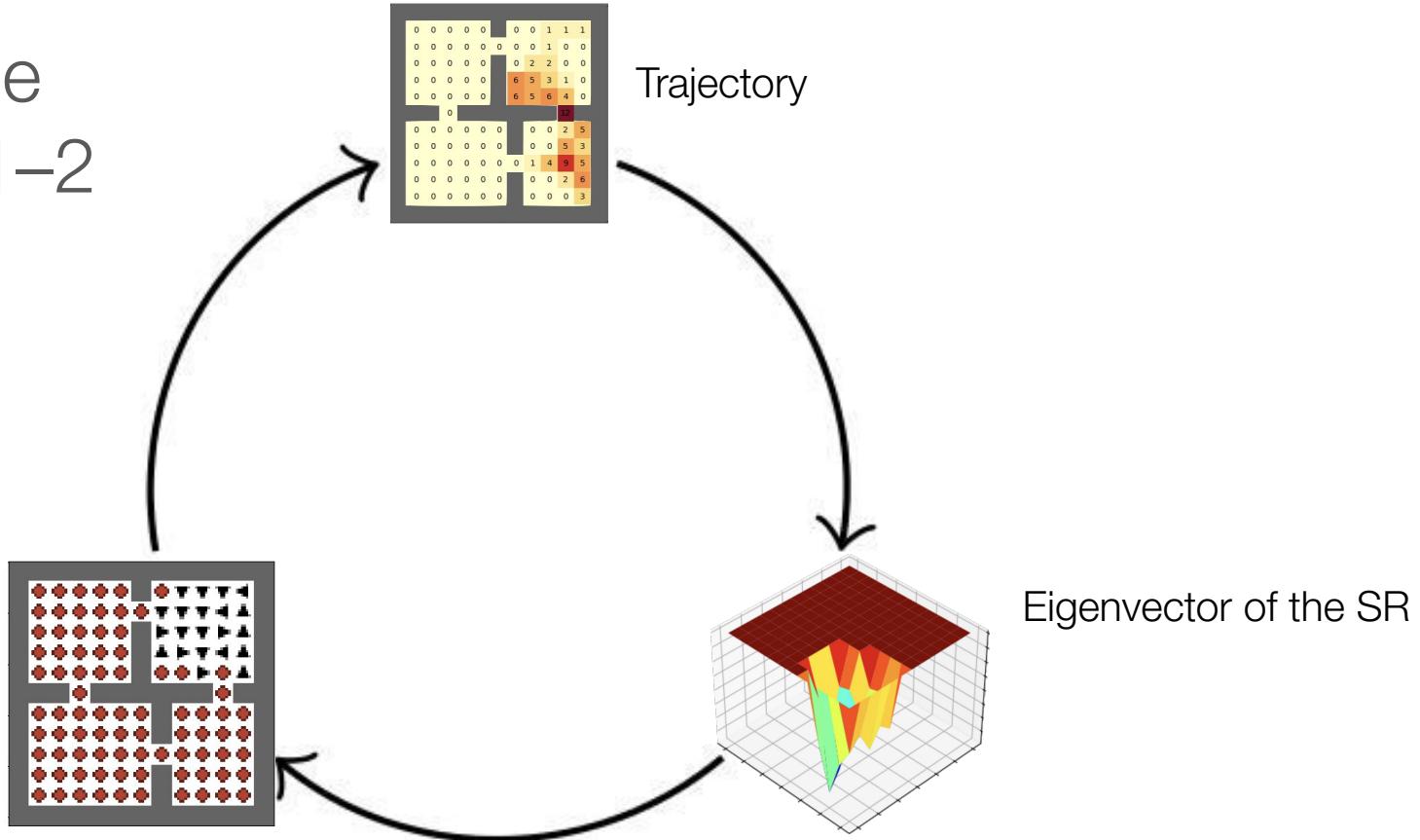
Trajectory

Eigenoption

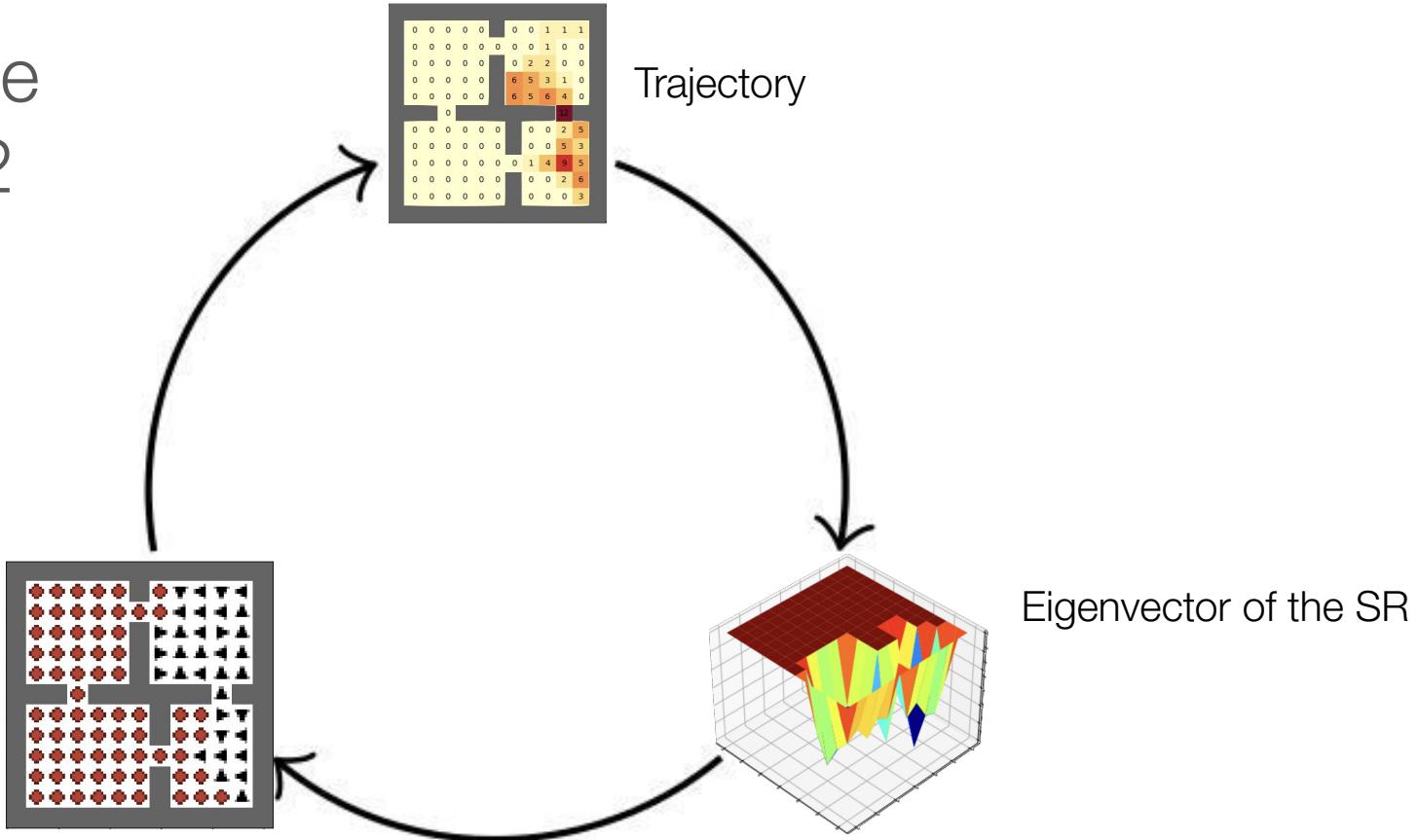


Eigenvector of the SR

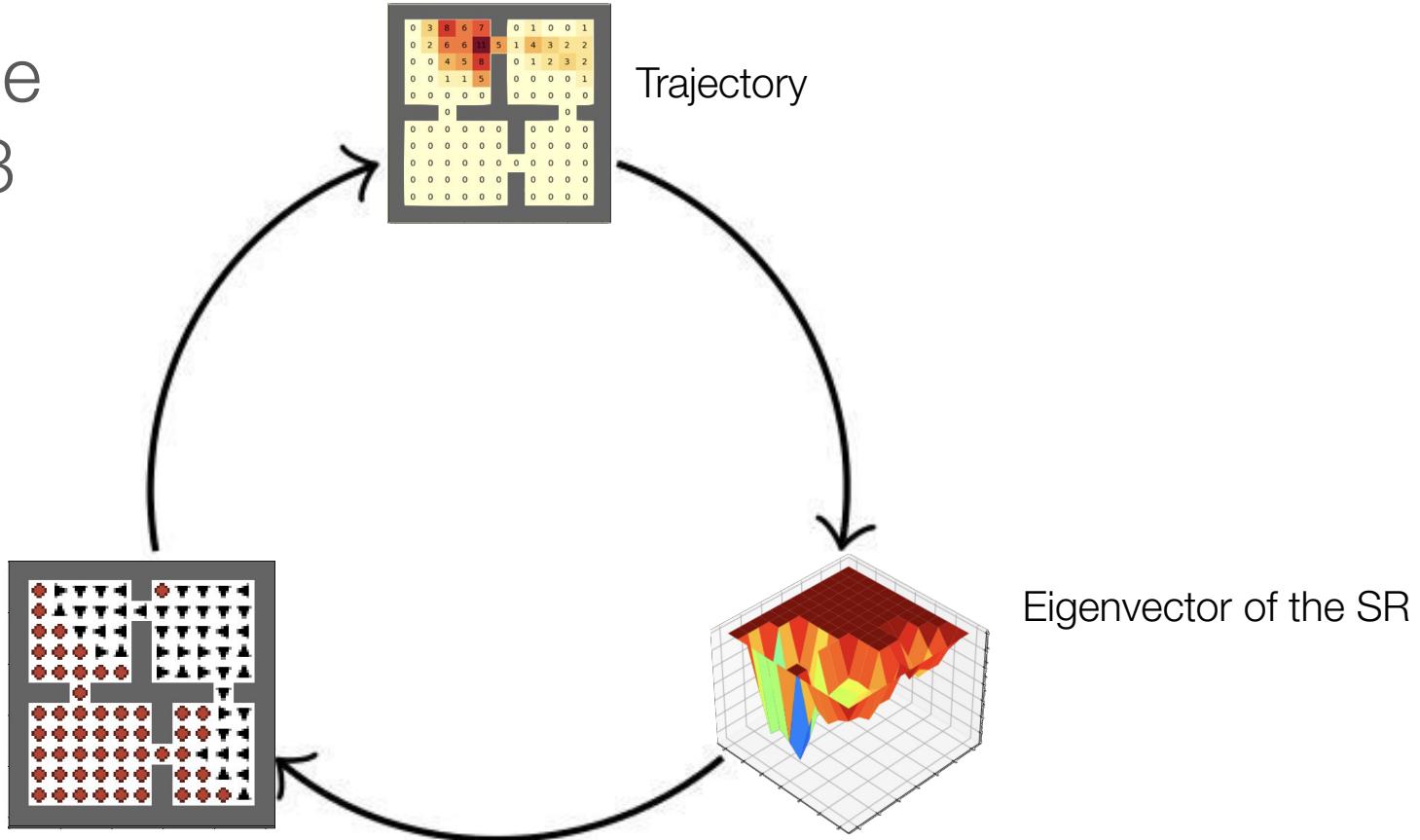
ROD cycle Iteration 1–2



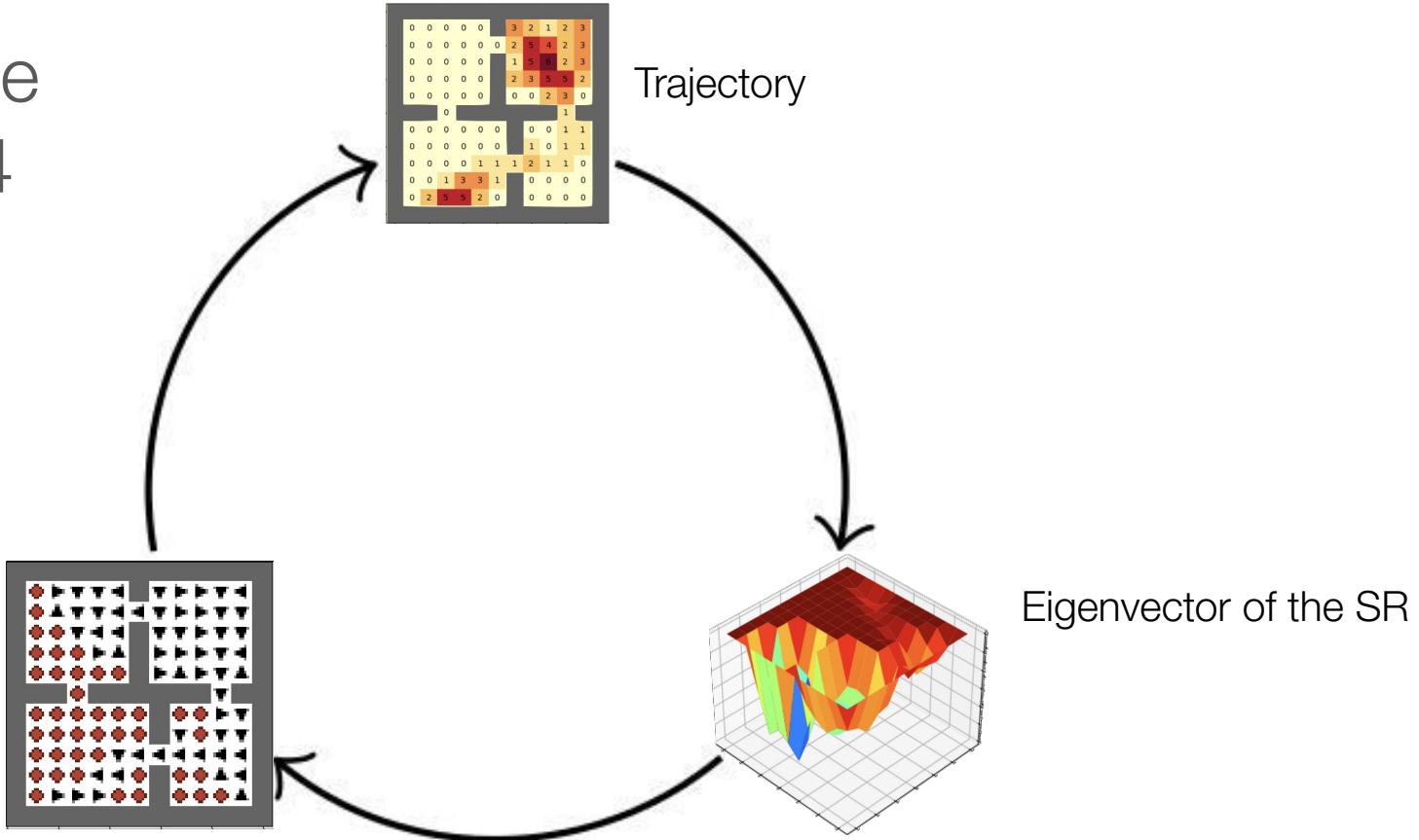
ROD cycle Iteration 2



ROD cycle Iteration 3

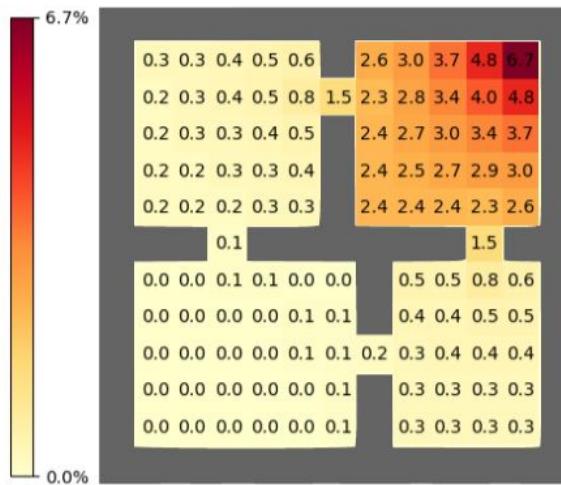


ROD cycle Iteration 4

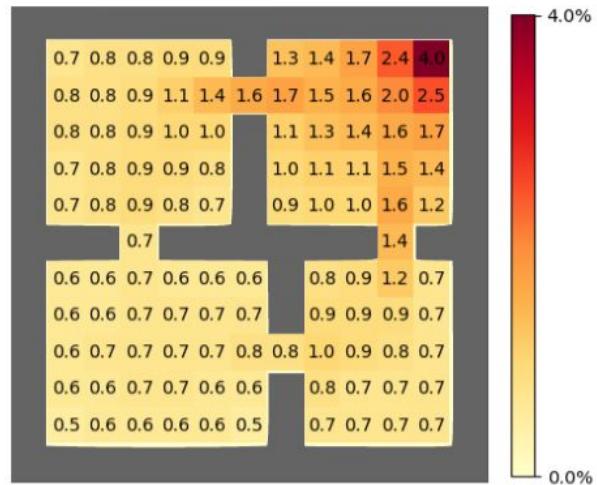


Covering eigenoptions vs random policy

State visitation:



Random policy over primitive actions

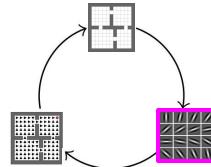


Random policy over primitive actions and covering eigenoptions

steps needed to visit all states: ~27,000

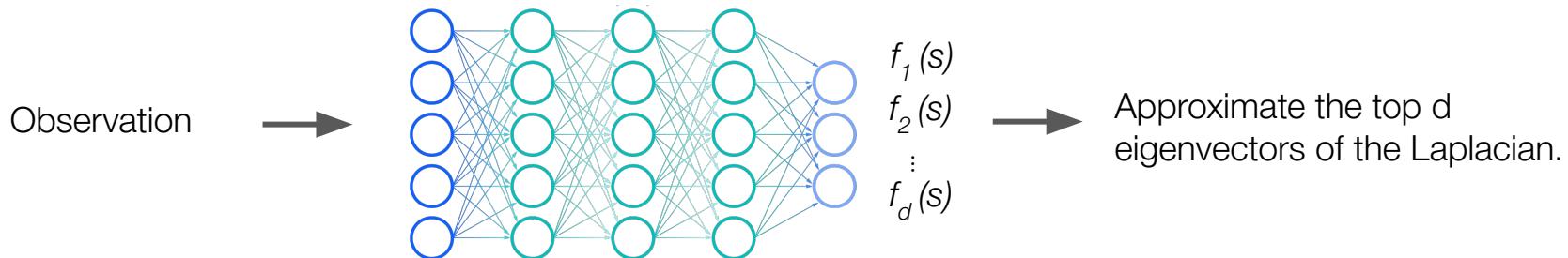
~2,300

Approximating the eigenfunctions of the graph Laplacian with neural networks



Approximate eigendecomposition

We can approximate the eigenvectors using a neural network and SGD

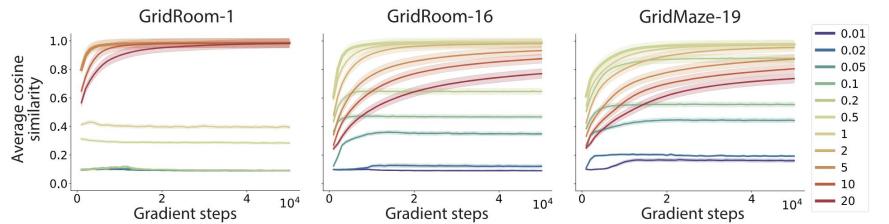
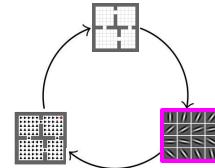


minimizing the *augmented Lagrangian Laplacian objective (ALLO)*:

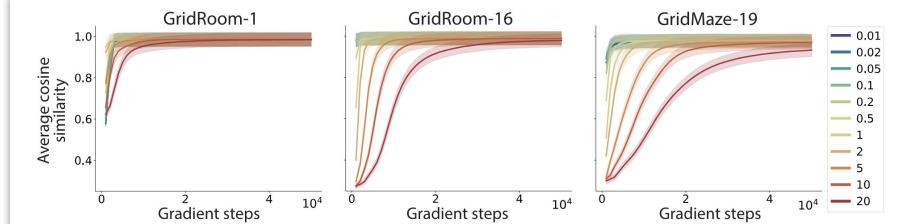
$$\max_{\boldsymbol{\beta}} \min_{\mathbf{u} \in \mathbb{R}^{d|\mathcal{S}|}} \sum_{i=1}^d \langle \mathbf{u}_i, \mathbf{L}\mathbf{u}_i \rangle + \sum_{j=1}^d \sum_{k=1}^j \beta_{jk} (\langle \mathbf{u}_j, [\![\mathbf{u}_k]\!] \rangle - \delta_{jk}) + b \sum_{j=1}^d \sum_{k=1}^j (\langle \mathbf{u}_j, [\![\mathbf{u}_k]\!] \rangle - \delta_{jk})^2$$



Proper Laplacian Representation Learning

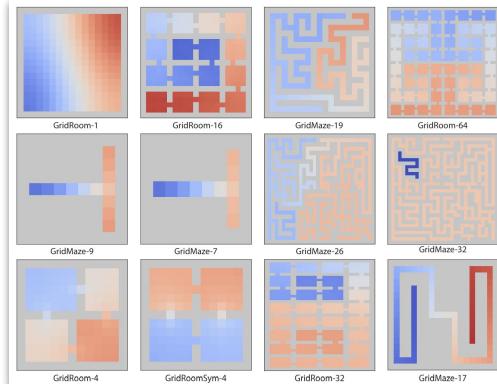


GGDO [Wang et al., ICML 2021]



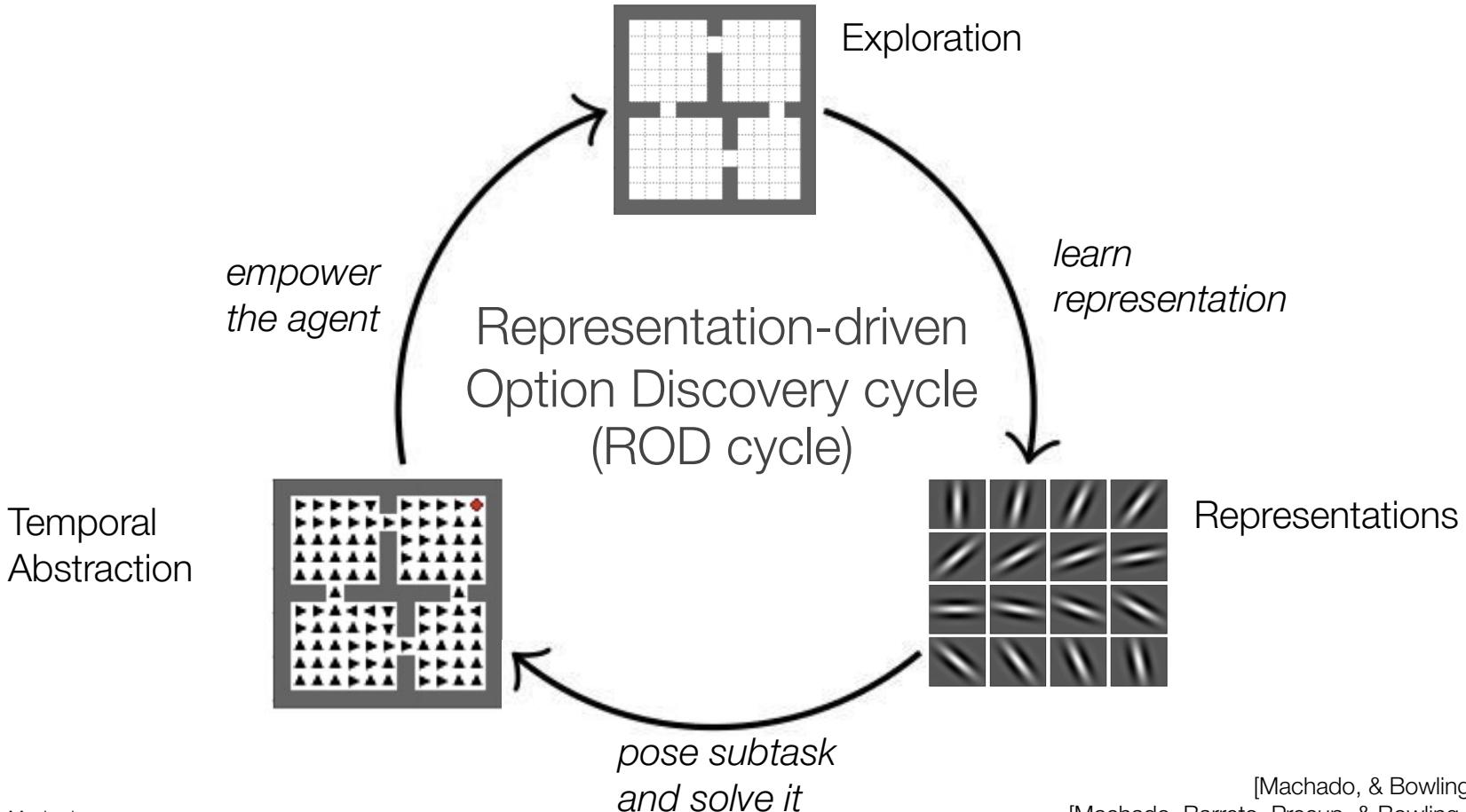
ALLO

Compared to previous approaches it is more robust, accurate, and it works across different data streams

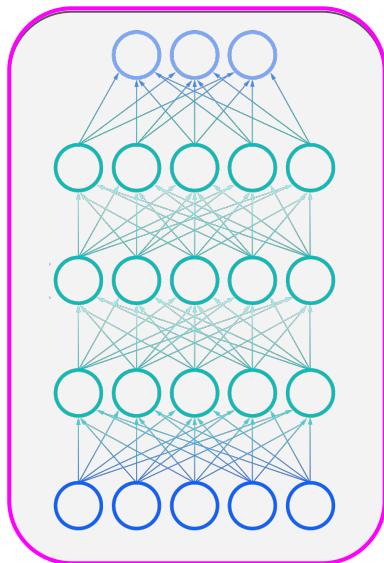


[Gomez, Bowling, & Machado, ICLR 2024]

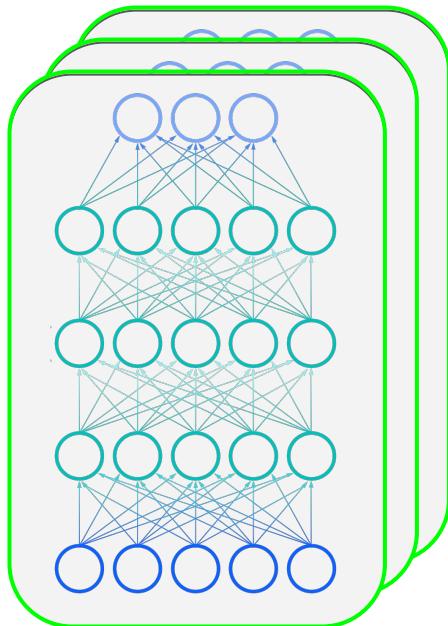
Putting everything together



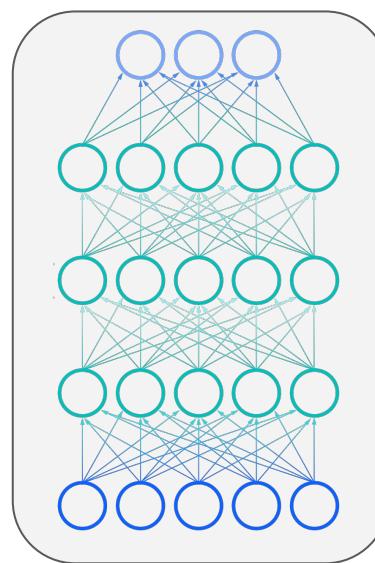
Approximations all the way down



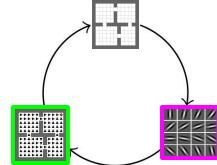
Laplacian representation
GGDO [Wang et al., 2021]



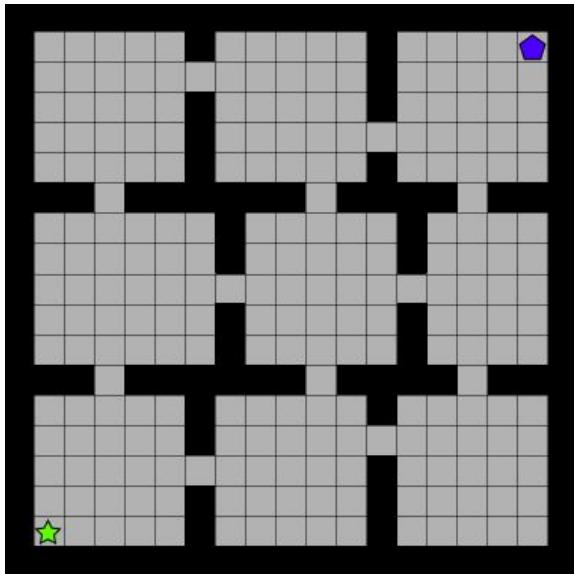
Eigenoptions
DDQN + n-step
[van Hasselt et al., 2016]



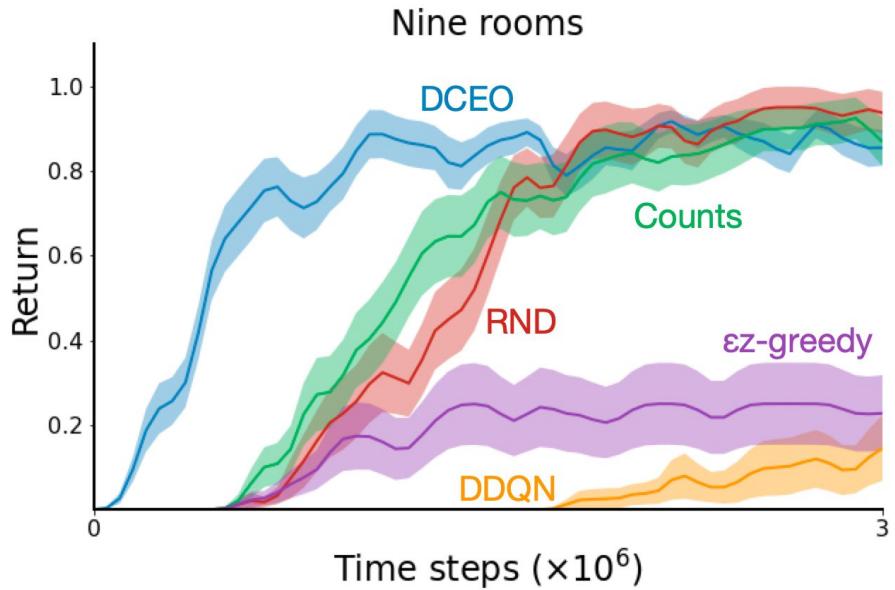
Main Q-learner
DDQN + n-step
[van Hasselt et al., 2016]



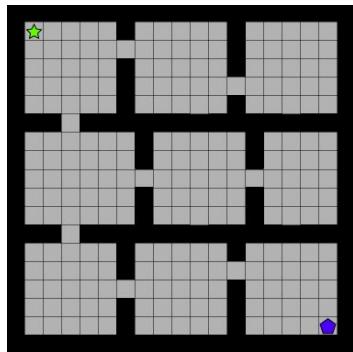
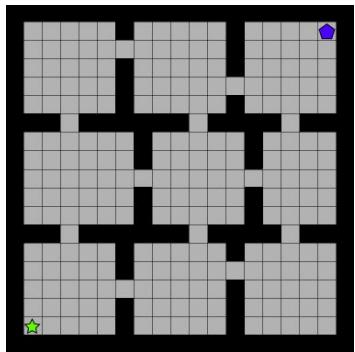
Online deep covering eigenoptions



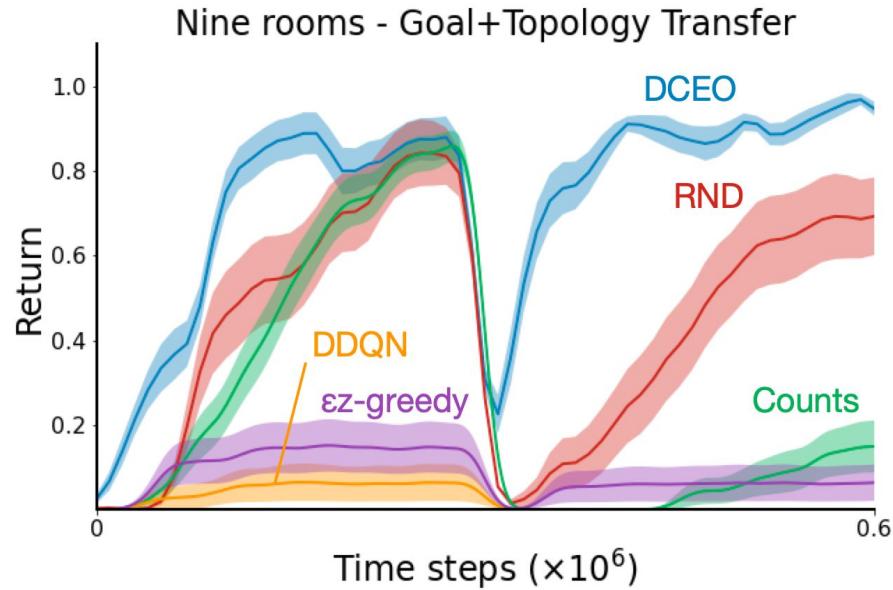
Observation type:
Image



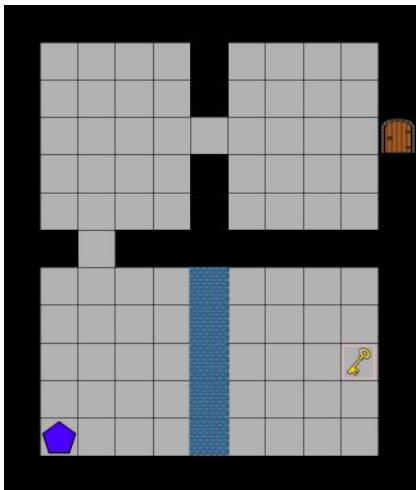
DCEO in non-stationary environments (continual learning)



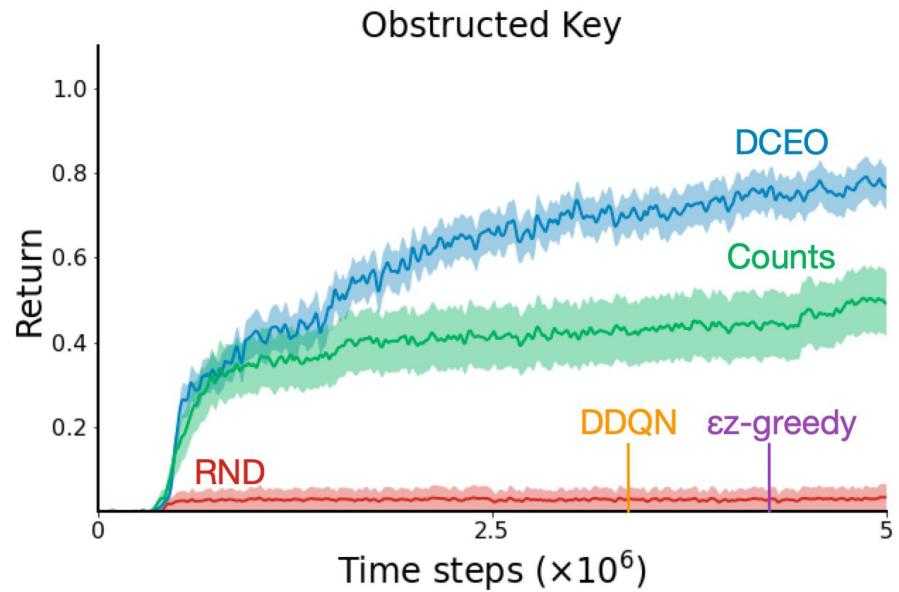
Observation type:
Image



and more!



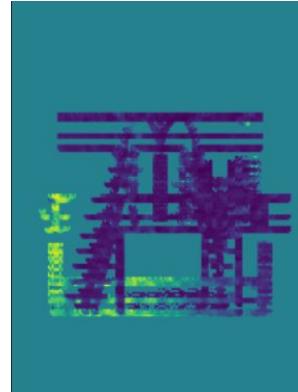
Observation type:
Image



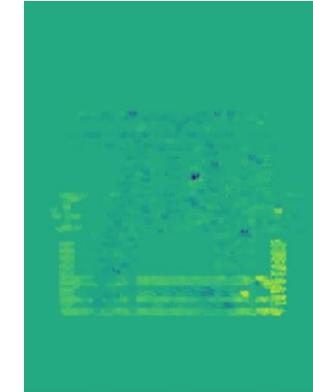
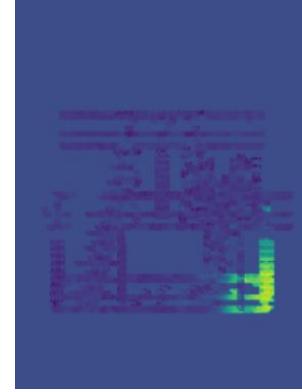
Subgoals in Montezuma's Revenge



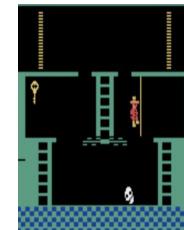
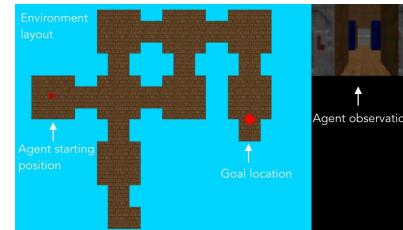
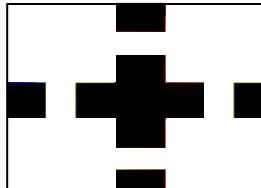
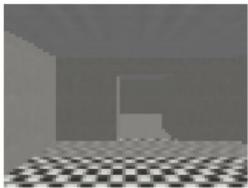
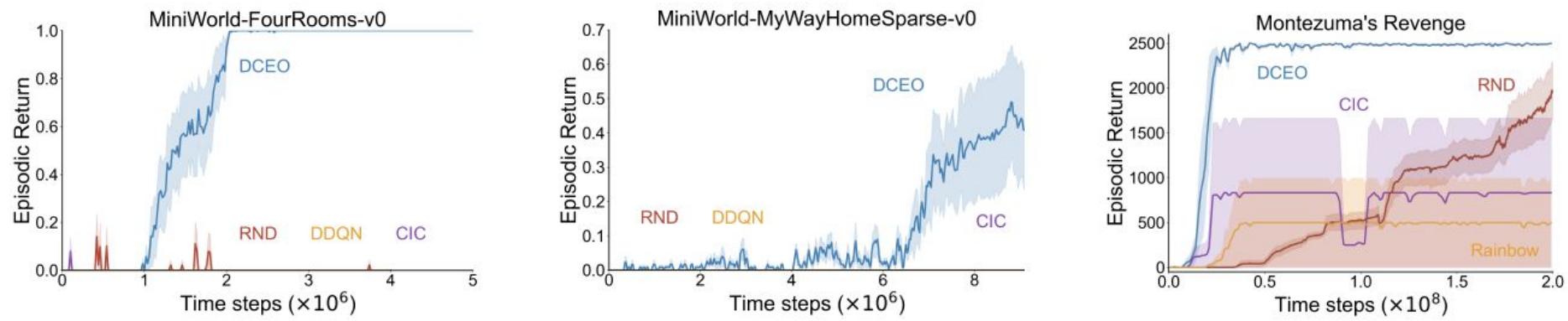
Montezuma's Revenge



Eigenfunctions discovered by the Laplacian



Beyond “navigation, 2d gridworld tasks”



Are we done?

Are we done?

- More demonstrations!
- More algorithms!

Are we done?

- More demonstrations! Learning from scratch is slow, almost by definition
 - There are many good reasons to do research when learning from scratch, but that impacts the types of demonstrations we can provide
 - We do have results leveraging domain knowledge (a.k.a. LLMs) alongside temporal abstractions, but I won't talk about that today (see work by Klissarov et al., 2025)
- More algorithms! Conceptually speaking, there are still some pieces missing
 - We are not using options for credit assignment
 - Should we define state similarity in terms of rewards as well?
 - What about partial observability?
 - Can we combine options without additional learning?
 - What about MBRL and planning?

Are we done?

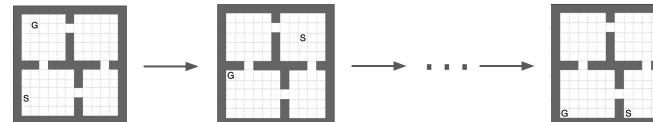
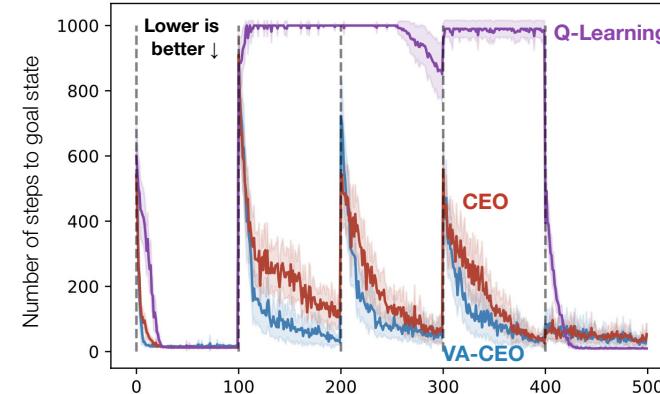
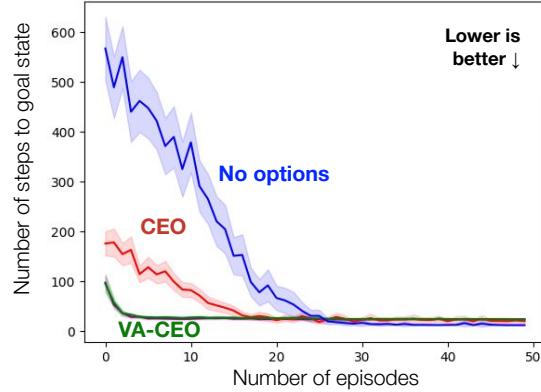
- More demonstrations! Learning from scratch is slow, almost by definition
 - There are many good reasons to do research when learning from scratch, but that impacts the types of demonstrations we can provide
 - We do have results leveraging domain knowledge (a.k.a. LLMs) alongside temporal abstractions, but I won't talk about that today (see work by Klissarov et al., 2025)
- More algorithms! Conceptually speaking, there are still some pieces missing
 - We are not using options for credit assignment  (Kotamreddy & Machado, In preparation)
 - Should we define state similarity in terms of rewards as well?  (Tse, Chandrasekar, & Machado, arXiv 2025)
 - What about partial observability?  (Jose & Machado, In preparation)
 - Can we combine options without additional learning?  (Chandrasekar & Machado, In preparation)
 - What about MBRL and planning? (see Sutton et al., 2023)

Are we done?

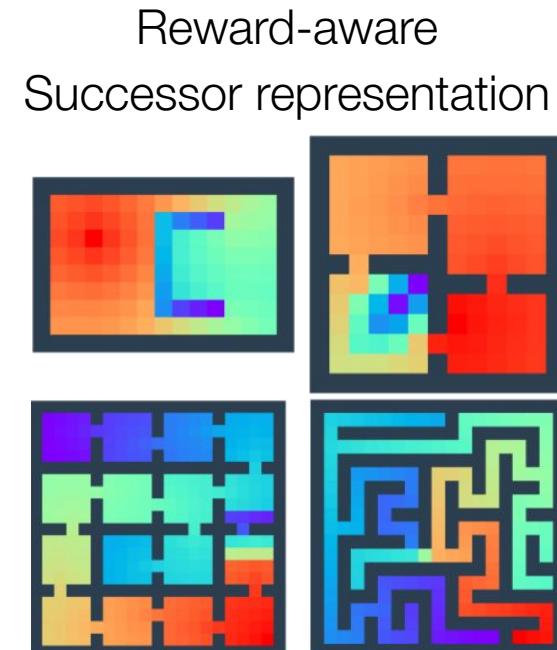
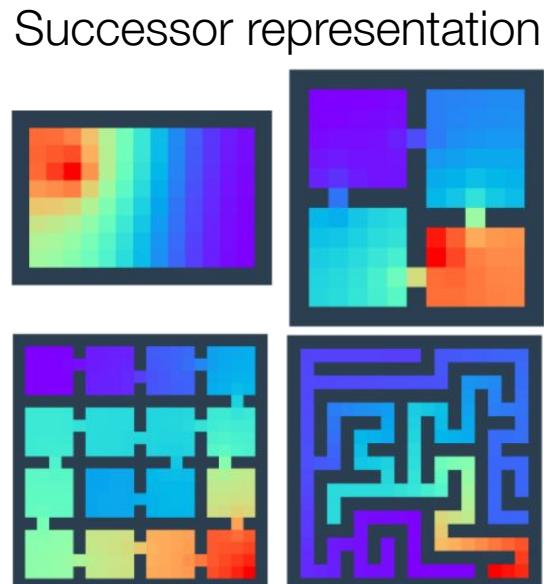
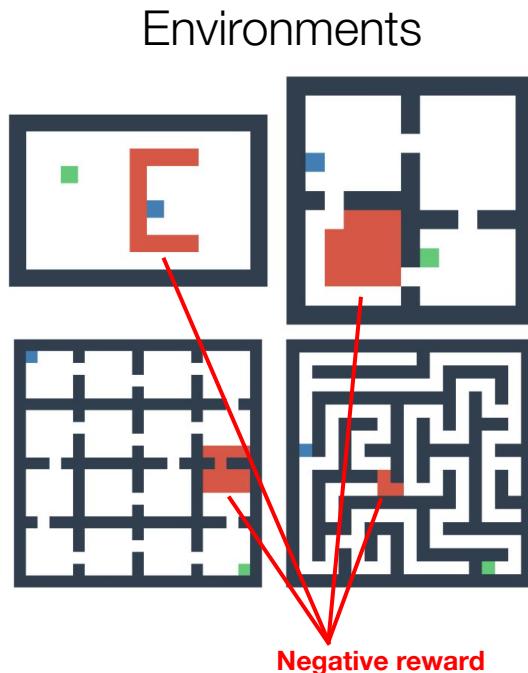
- More demonstrations! Learning from scratch is slow, almost by definition
 - There are many good reasons to do research when learning from scratch, but that impacts the types of demonstrations we can provide
 - We do have results leveraging domain knowledge (a.k.a. LLMs) alongside temporal abstractions, but I won't talk about that today (see work by Klissarov et al., 2025)
- More algorithms! Conceptually speaking, there are still some pieces missing
 - **We are not using options for credit assignment**  (Kotamreddy & Machado, In preparation)
 - **Should we define state similarity in terms of rewards as well?**  (Tse, Chandrasekar, & Machado, arXiv 2025)
 - What about partial observability?  (Jose & Machado, In preparation)
 - Can we combine options without additional learning?  (Chandrasekar & Machado, In preparation)
 - What about MBRL and planning? (see Sutton et al., 2023)

Model-Free Value-Aware Covering Eigenoptions

- What if we used options not only for exploration but also for credit assignment?
- Pre-computed (tabular) options in the four-room environment:

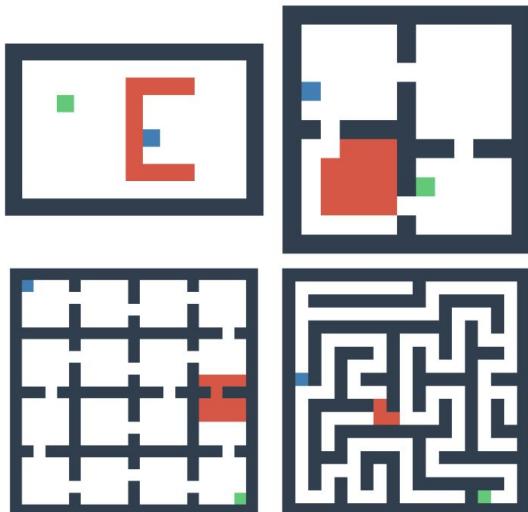


Reward-Aware Proto-Representation

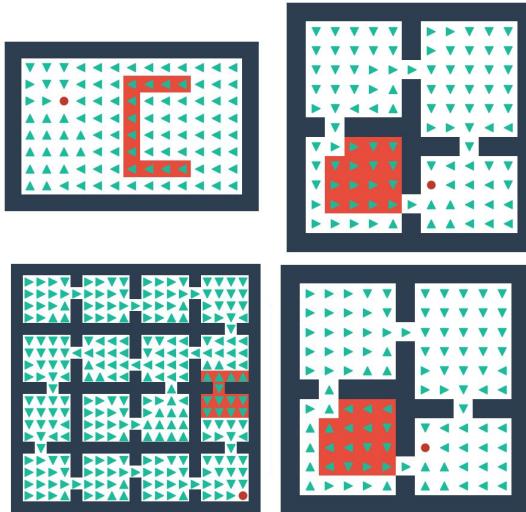


Reward-Aware Eigenoptions

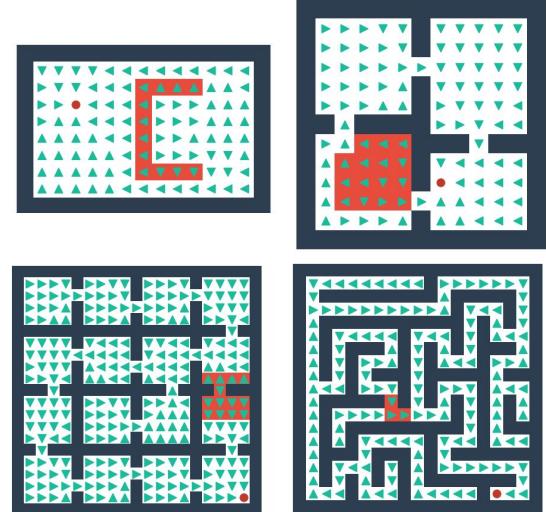
Environments



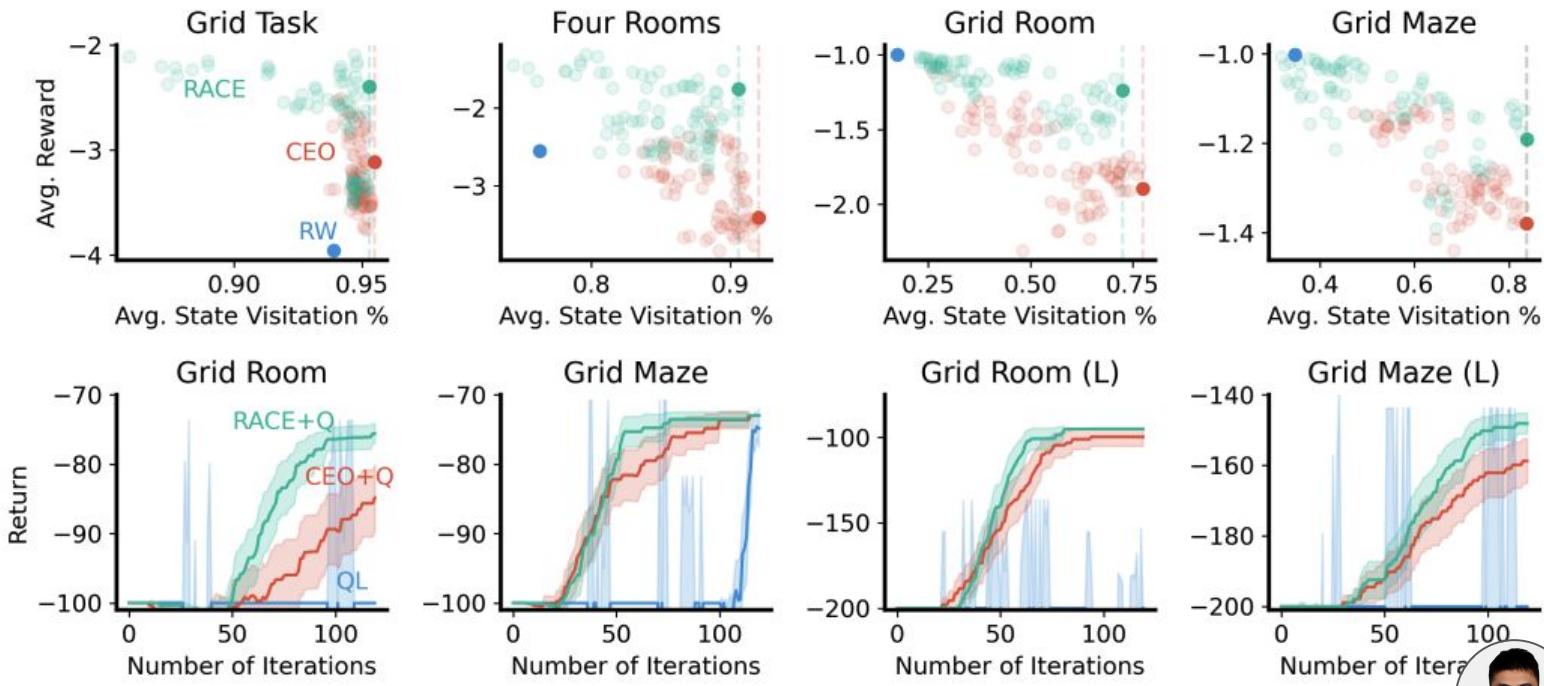
Successor representation



Reward-aware
Successor representation



Reward-Aware Eigenoptions



Wrapping up

Conclusion

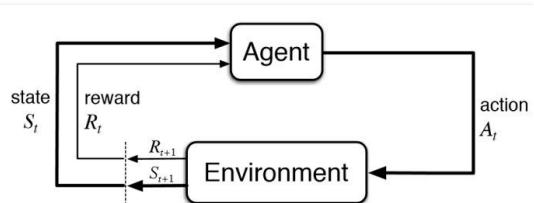
I'm still very excited about options. We now have option discovery methods that are general, that work, and that:

- demonstrate a virtuous cycle, **and**
- are fully experiential, **and**
- are scalable, **and**
- are amenable to function approximation, **and**
- work for different data streams, **and**
- don't make any assumptions about the topology of the environment, **and**
- is (sort of) biologically plausible, if you are into that kind of thing _(ツ)_/

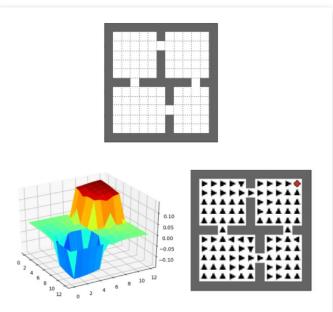
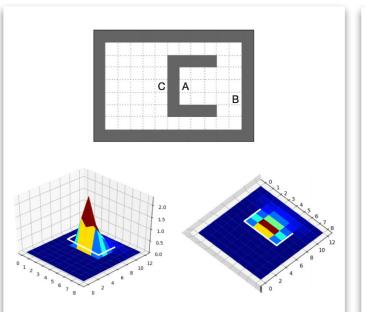
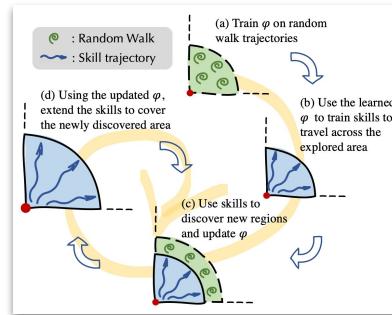
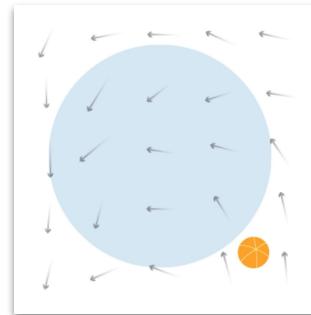
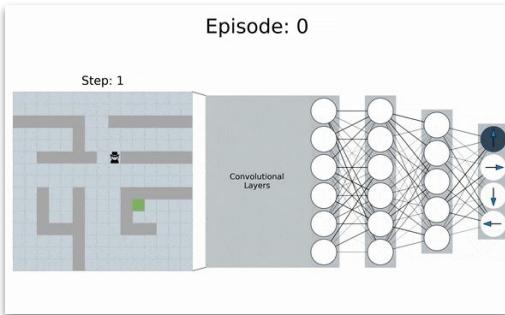
Conclusion

- Temporal abstractions should be a central piece of reinforcement learning
- Where should options come from?
 - Specific representations learned by the agent
- Given the right discovery method, options are scalable.
- Options are particularly helpful for continual learning/in the face of non-stationarity.
 - This is what I believe is the future of our field ↑

“State representations and temporal abstractions should be deeply intertwined, where representations and options are constantly refined based on each other.”



$$\omega = \langle \mathcal{I}_\omega, \pi_\omega, \beta_\omega \rangle$$



Thanks!

