# Algorithmic category theory for reinforcement learning

Tanner Duve, Theo M. Tyburn, Andrea Abeje-Stine, Onkar KaleAdjoint School 2025Mentor: Georgios BakirtzisTA: Michail Savvas

Tue 3 Jun 2025

## Reinforcement learning (RL) solves problems through reward signals

Training can happen without having a complete model of the task.



The video game Dota 2

# But RL faces serious limitations in solving complex problems

Some of the limitations of RL arise from:

- reward sparsity,
- huge action spaces  $\rightarrow$  sample inefficiency,
- lack of stability  $\rightarrow$  poor safety guarantees,
- poor generalization, and
- poor interpretability.



Typical framing of an RL scenario

Compositionality is being used in RL through two main paradigms:

#### 1. Hierarchical RL

models sequential tasks that can be split into independent sequential subtasks.

#### 2. Functional RL

models parallel tasks that can be split into interacting subtasks.



# A robotic task decomposed hierarchically and functionally

Two main problems in RL are

- **Scalability**: RL algorithms can have too complicated state and action spaces to keep track of.
- **Generalization**: RL agents may overfit to their training environments and perform poorly on new environments.

Compositionality can be used to solve both

- scalability through more efficient scaling methods and
- generalization by learning composable problems

## Multi-agents scenarios suffer from both scaling and generalization issues

Multi-agent scenarios are generally modeled by taking the products of MDPs:  $M_1 \times \cdots \times M_n$ .



Problem: The action space  $A_1 \times \cdots \times A_n$  rapidly blows up.

#### We want to reduce the action space to something like this



An Markov Decision Process (MDP)  $\mathcal{M} = (S, A, \psi, T, R)$  is composed of the data:

- *S* State space (measurable)
- A Action space
- $\psi: \mathbf{A} \rightarrow \mathbf{S}$  Action source function
- $T: A \rightarrow \mathcal{P}_S$  Transition probability function
- $R: A \rightarrow \mathbb{R}$  Reward function



# MDPs model sequential decision making under uncertainty







The category **MDP** has pushouts and partial pullbacks.

subtasks	injective morphisms		
safety	punctured MDP		
symmetry	quotient MDP		
sequential tasks	pushouts		
parallel tasks	products		

Through the category of **MDP** we studied

- A "box product" for MDP to potentially improve scalability
- A "twisted MDP" which is a combination of linear temporal logic specification for requirements (safety and liveness) and MDPs for behavior to potentially improve generalization
- Comonads to carry side information around in the category MDP

# Injective maps in MDP model subprocess MDPs



 $\mathcal{M}_1 \stackrel{(f,g)}{\longleftrightarrow} \mathcal{M}_2$ 

**Safe grid world:** navigate to a goal state while avoiding obstacle states.



# Punctured MDPs enforce safety conditions

Let  $\mathcal{M} = (S, A, \psi, T)$  be an MDP,  $\mathbb{O} \subset S$ .

• Punctured MDP:

$$\mathcal{M}^{\circ} = (S^{\circ}, A^{\circ}, \psi^{\circ}, T^{\circ}) \text{ where:} \\ 1. S^{\circ} = S \setminus \mathbb{O}. \\ 2. A^{\circ} = A \setminus (\psi^{-1}(\mathbb{O}) \cup \text{supp}_{T}(\mathbb{O})). \\ 3. \psi^{\circ} = \psi|_{A^{\circ}}. \\ 4. T^{\circ} = T|_{A^{\circ}}. \end{cases}$$

•  $\mathcal{M}^{\circ}$  is the canonical maximal subprocess of  $S^{\circ}$ .



# Sequential tasks can be modelled using pushouts of zig-zag diagrams



# Sequential tasks can be modelled using pushouts of zig-zag diagrams



#### Parallel tasks can be modeled using products

To model *n* agents acting in parallel, we can take the cartesian product of MDPs:  $\mathcal{M}_{\times} \cdots \times \mathcal{M}_{n}$ .



Problem: The size of the action space  $A_1 \times \cdots \times A_n$  is  $|A_1| \cdots |A_n|$ .

#### We want to reduce the action space to something like this

Can we define a product  $M_1 \Box \ldots \Box M_n$  where the action space has a size comparable to  $|A_1| + \cdots + |A_n|$ ?



#### Stationary actions are the building blocks of the box product



The box product is a generalization of the box product for graphs inspired by the funny tensor product for categories.

#### Theorem

Given MDPs  $M_1$  and  $M_2$  we get the following pushout diagram:



we call  $\mathcal{M}_1 \Box \mathcal{M}_2$  the box product of MDPs  $\mathcal{M}_1$  and  $\mathcal{M}_2$ .

- 1. State space  $S = S_1 \times S_2$
- 2. Action space is given by

$$A = \overline{A_1} imes (A_2)_0 \amalg (A_1)_0 imes \overline{A_2} / \sim$$

where  $\sim$  identifies redundant stationary states. For large action spaces:

$$|S_1| \cdot |A_2| + |S_2| \cdot |A_1| \ll |A_1| \cdot |A_2|$$

#### The size of action space grows much slower for the box product



# We wrote a python library to conduct experiments

#### https://codeberg.org/bakirtzis/ctmdp

C Explore About FAQ Help Don	te		A Register	
💂 bakirtzis / ctmdp				
<> Code ⊙ Issues 11 Pull requests -∿- Act				
Markov decision processes, categorically				
🕲 28 commits	🗜 1 branch	🛇 0 tags	<b>Ө 86</b> Ків	
🐉 main 👻 🕅 Find a file		HTTPS https://coo	deberg.org/bakirtzis/ctmdp.git	ወ
Théo Tyburn 51e28da924 bp_experiments: adapt to refactor yesterda				
🖿 ctmdp				
🖿 examples	bp_experiments: adapt to refactor			
🗅 .gitignore				
C README.org				
T README.org				

- Getting started
- Evamples

#### Motivation.

In functional programming, **comonads** represent computations that use context. In RL, **side information** is additional knowledge **not** encoded in the MDP

### Hypothesis.

Comonads in the category **MDP** can formally model side information

### Examples of side info we constructed comonads for:

- Safety constraints (e.g. unsafe region avoidance via twisted MDP)
- Symmetry structure (e.g. group actions, orbit-aware generalization)
- History context (e.g. fixed-length memory)

The approach: use category theory to come up with concrete algorithms.

The hope: the category **MDP** (or an extession of it) is the right framework to develop a theory of decomposition of RL agents that is stable under generalization.

The goal: provide RL practitioners with tools for:

- reward design,
- reusability and generalization,
- training efficiency via action space reduction.

#### 1. Problems we tried solving

- $\circ$  generalization
- scaling
- 2. Solution: use compositional methods
  - Sequential composition: pushouts of MDPS
  - Functional composition: box products
  - Composing side information: comonads
- 3. Further things we want to investigate
  - Extensions of span and cospan categories
  - Introducing error terms into MDP morphisms
  - MDPs as coalgebras
  - (Probabilistic) Linear Temporal Logic
  - Partially Observable MDPs

- Categorical semantics of compositional reinforcement learning, Georgios Bakirtzis, Michail Savvas, and Ufuk Topcu
- Structure in Deep Reinforcement Learning: A Survey and Open Problems, Aditya Mohan, Amy Zhang, Marius Lindauer
- Categorical Semantics of motion planning, Georgios Bakirtzis and Ufuk Topcu

#### THANK YOU FOR YOUR ATTENTION!