# There Is No Turning Back: A Self-Supervised Approach for Reversibility-Aware Reinforcement Learning

Nathan Grinsztajn, Johan Ferret,
Olivier Pietquin, Philippe Preux,
Matthieu Geist

*Reading notes*
Malik-Manel Hashim

# Agenda

- Overview
- Approach
- Reversibility and Reversibility Estimation
- Reversibility-Aware Reinforcement Learning
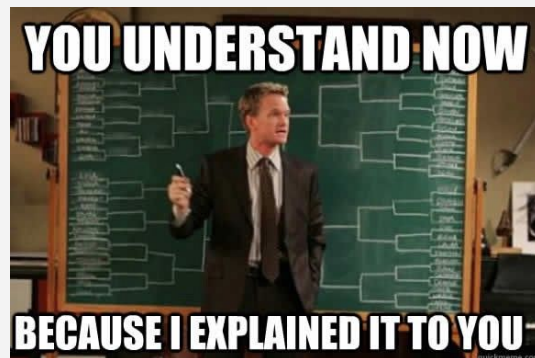- Experiments
- Summary
- Conclusion and Critique

# Agenda

- Overview
- ~~Approach~~                  1. *Idea*
- ~~Reversibility and Reversibility Estimation~~    2. *Math*
- ~~Reversibility-Aware Reinforcement Learning~~ 3. *How to RL*
- ~~Experiments~~                  4. *Tests and results*
- ~~Summary~~
- Conclusion and Critique



[2]

# Overview

- Knowing reversibility of an action = Knowing its potential risk
- A always before B $\Rightarrow$ A $\rightarrow$ B not reversible
  - Simple binary classification
- Can be used for exploration / control
- Performs great (on what they tested it on)



[3]

*(The Idea)*

# Approach

- Reversibility = Safety
- Approximation via temporal order
- Can be learned through a surrogate task



Is A → B reversible?

Yes, because B → A does not contradict the laws of physics!

Is A → B reversible?

Easy, since B → A is as likely as A → B!

*hard*      *simple(-ish)*

*(The Idea)*

# Approach

*(The Math)*

# **Reversibility and Reversibility Estimation**

Degree of reversibility within K steps

$$\phi_K(s, a) := \sup_{\pi} p_\pi(s \in \tau_{t+1:t+K+1} \mid s_t = s, a_t = a)$$

Degree of reversibility

$$\phi(s, a) := \sup_{\pi} p_\pi(s \in \tau_{t+1:\infty} \mid s_t = s, a_t = a)$$

*(The Math)*

# Reversibility and Reversibility Estimation

Degree of reversibility within K steps

$$\phi_{\pi,K}(s,a) := \sup_{\pi} p_\pi(s \in \tau_{t+1:t+K+1} \mid s_t = s, a_t = a)$$

Degree of reversibility

$$\phi_\pi(s,a) := \sup_{\pi} p_\pi(s \in \tau_{t+1:\infty} \mid s_t = s, a_t = a)$$

*(The Math)*

# **Reversibility and Reversibility Estimation**



[4]

- Reversibility can be predicted via *precedence*

  → *s* or *s' first on average*

Finite-horizon precedence estimator

$$\psi_{\pi,T}\left(s, s'\right) = \mathbb{E}_{\tau \sim \pi} \mathbb{E}_{s_t=s, s_{t'}=s'} \left[\mathbb{1}_{t'>t}\right]$$
$$t, t' < T$$

Emperical reversibility

$$\overline{\phi}_\pi(s, a) = \mathbb{E}_{s' \sim P(s,a)} \left[\psi_\pi\left(s', s\right)\right]$$

*(The Math)*

# **Reversibility and Reversibility Estimation**

Reversibility: $$\phi_\pi(s,a) := p_\pi(s \in \tau_{t+1:\infty} \mid s_t = s, a_t = a)$$

Empirical reversibility: $$\overline{\phi}_\pi(s,a) = \mathbb{E}_{s' \sim P(s,a)} \left[ \psi_\pi(s', s) \right]$$

Relation of reversibility and empirical reversibility

$$\overline{\phi}_\pi(s,a) \geq \frac{\phi_\pi(s,a)}{2}$$

*(The Math)*

# **Reversibility and Reversibility Estimation**

Reversibility: 
$$\phi_\pi(s, a) := p_\pi(s \in \tau_{t+1:\infty} \mid s_t = s, a_t = a)$$

[5]

Empirical reversibility: 
$$\overline{\phi}_\pi(s, a) = \mathbb{E}_{s' \sim P(s,a)} \left[ \psi_\pi(s', s) \right]$$

Relation of reversibility and empirical reversibility

$$\overline{\phi}_\pi(s, a) \geq \frac{\phi_\pi(s, a)}{2}$$



I'M AFRAID

I'M GONNA NEED PROOF

# Reversibility and Reversibility Estimation *(The Math)*



**A.3  Proofs of Theorem 1 and Theorem 2**

In the following, we prove simultaneously Theorem 1 and Theorem 2. We begin by two lemmas.

**Lemma 1.** *Given a trajectory $\tau$, we denote by $\#_T(s \to s')$ the number of pairs $(s, s')$ in $\tau_{1:T}$ such that $s$ appears before $s'$. We present a simple formula for $\psi(s', s)$ according to the structure of the state trajectory:*

$$\psi_{\pi, T}(s, s') = \frac{\mathbb{E}_{\tau \sim \pi}\left[\#_T(s \to s')\right]}{\mathbb{E}_{\tau \sim \pi}\left[\#_T(s \to s') + \#_T(s' \to s)\right]}.$$

*Proof.* In order to simplify the notations, we leave implicit the fact that indices are always sampled within $[0, T]$.

$$\psi_{\pi, T}(s, s') = \mathbb{E}_\pi \mathbb{E}_{t \neq t' | s_t = s, s_{t'} = s'}\left[\mathbb{1}_{t' > t}\right],$$
$$= \frac{\mathbb{E}_\pi \mathbb{E}_{t \neq t'}\left[\mathbb{1}_{t' > t}\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right]}{\mathbb{E}_\pi \mathbb{E}_{t \neq t'}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right]}.$$

Similarly, we have:

$$\mathbb{E}_\pi \mathbb{E}_{t' > t}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right] = \frac{\mathbb{E}_\pi \mathbb{E}_{t \neq t'}\left[\mathbb{1}_{t' > t}\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right]}{\mathbb{E}_{t \neq t'}\left[\mathbb{1}_{t' > t}\right]}.$$

Combining it with our previous equation:

$$\psi_{\pi, T}(s, s') = \frac{\mathbb{E}_\pi \mathbb{E}_{t' > t}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right] \mathbb{E}_{t \neq t'}\left[\mathbb{1}_{t' > t}\right]}{\mathbb{E}_\pi \mathbb{E}_{t \neq t'}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right]},$$
$$= \frac{1}{2}\frac{\mathbb{E}_\pi \mathbb{E}_{t' > t}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right]}{\mathbb{E}_\pi \mathbb{E}_{t \neq t'}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right]}.$$

Looking at the denominator, we can notice:

$$\mathbb{E}_\pi \mathbb{E}_{t \neq t'}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right] = \frac{1}{2}\mathbb{E}_\pi \mathbb{E}_{t < t'}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right] + \frac{1}{2}\mathbb{E}_\pi \mathbb{E}_{t' < t}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right],$$
$$= \frac{1}{2}\mathbb{E}_\pi \mathbb{E}_{t < t'}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'} + \mathbb{1}_{s_t = s'}\mathbb{1}_{s_{t'} = s}\right],$$

which comes from the fact that $t$ and $t'$ play a symmetrical role. Thus,

$$\psi_{\pi, T}(s, s') = \frac{\mathbb{E}_{\tau \sim \pi}\mathbb{E}_t \mathbb{E}_{t' > t}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right]}{\mathbb{E}_{\tau \sim \pi}\mathbb{E}_t \mathbb{E}_{t' > t}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'} + \mathbb{1}_{s_t = s'}\mathbb{1}_{s_{t'} = s}\right]}.$$

Since

$$\mathbb{E}_{\tau \sim \pi}\left[\#_T(s \to s')\right] = \sum_{i < j \leq T} \mathbb{1}_{s_i = s}\mathbb{1}_{s_j = s'},$$
$$= \binom{T}{2} \sum_{i < j \leq T} \frac{1}{\binom{T}{2}}\mathbb{1}_{s_i = s}\mathbb{1}_{s_j = s'},$$
$$= \binom{T}{2} \mathbb{E}_{\tau \sim \pi}\mathbb{E}_t \mathbb{E}_{t' > t}\left[\mathbb{1}_{s_t = s}\mathbb{1}_{s_{t'} = s'}\right],$$

16

*(The Math)*

# Reversibility and Reversibility Estimation

# Reversibility and Reversibility Estimation *(The Math)*

# Reversibility and Reve...

*(The Math)*

[Nathan et al., *NeurIPS* 2021]

# Reversi

*(The Math)*

[Nathan et al., *NeurIPS* 2021]

# Reversi

(The Math)

[1, p. 16ff]

*(The Math)*

# **Reversibility and Reversibility Estimation**

Reversibility:

$$\phi_\pi(s, a) := p_\pi(s \in \tau_{t+1:\infty} \mid s_t = s, a_t = a)$$

Empirical reversibility:

$$\overline{\phi}_\pi(s, a) = \mathbb{E}_{s' \sim P(s,a)} \left[ \psi_\pi(s', s) \right]$$

Relation of reversibility and empirical reversibility

$$\overline{\phi}_\pi(s, a) \geq \frac{\phi_\pi(s, a)}{2}$$



AND THEN I SAID

THE PROOF IS TRIVIAL

[6]

# Reversibility-Aware Exploration / Control *(How to RL)*

(How to RL)

# Reversibility-Aware Exploration / Control



(a) RAE penalizes irreversible transitions

(b) RAC hijacks irreversible actions

[1, p. 6]

# Reversibility-Aware Exploration / Control *(How to RL)*



[1, p. 5]

# Tests and results



Turf-Environment

[1, p. 8]

Sokoban

[7]

Cartpole

[8]

# Tests and results – Turf



(a) Initial state    (b) A trajectory    (c) PPO (500k)    (d) PPO+RAE (500k)

[1, p. 8]

- No irreversible actions
- Slower learning



[1, p. 9]

# Tests and results - Sokoban



[1, p. 8]

- Very challenging environment
- Sparse irreversible actions
- Better and more consistent performance with RAE

# **Tests and results - Cartpole**

Reward-free cartpole (RAE)

Cartpole+ (RAC)

Note: Color indicates
estimated reversibility values



Action 0 (move left)          Action 1 (move right)     [1, p. 9]

# **Summary**

- Safety = Reversibility
- Definition of reversibility via precedence
- Precedence classification as surrogate task
- Reversibility-Aware Exploration / Control
- Turf / Sakoban / Cartpole



[9]

# Conclusion and critique

- Simple representation of a complex task
- Extremely modular, can be used on any policy
- Requires (a lot) more testing



[10]

# Sources

- Paper:
  - [1] – Nathan Grinsztajn and Johan Ferret and Olivier Pietquin and Philippe Preux and Matthieu Geist, "There Is No Turning Back: A Self-Supervised Approach for Reversibility-Aware Reinforcement Learning", *NeurIPS 2021.*

- Images:
  - [2]–https://blog.ml6.eu/catching-the-ai-train-c0c496959999
  - [3]–https://medium.com/decktopus/15-memes-everyone-who-has-given-a-presentation-will-relate-to-e4946babfc6f
  - [4]–https://www.pinterest.de/pin/625718941963149650/
  - [5]–https://www.pinterest.es/pin/691935930226547050/
  - [6]–http://www.quickmeme.com/meme/3u2bs0
  - [7]–https://mobile.twitter.com/GoogleAI/status/1455973174319910915?cxt=HHwWhsCwner407QoAAAA
  - [8]–https://github.com/ganeshjha/Cartpole
  - [9]-https://imgflip.com/i/71ua6o
  - [10]-https://towardsdatascience.com/deep-learning-a-monty-hall-strategy-or-a-gentle-introduction-to-deep-q-learning-and-openai-gym-d66918ac5b26

# Questions?