

# Learning to Remember the Past by Learning to Predict the Future

PhD Defense - April 7th, 2025

Gaspard Lambrechts

# A Matter of Perception







# What should the robot do?













-2 7 M

















# **Partial observability**

• The agent does not observe the full state of the environment.



• Instead, its perception is limited to a partial observation.



# Past and future

• We need to remember the state of things from the past.



• We need to learn the impact of our actions on the future.



# **Intelligent behaviors**



# Learning intelligent behaviors



# Learning curves



# **Sufficient memory**



# **Sufficient memory**



# Reinforcement Learning under Partial Observability

# Partially observable Markov decision process



A **POMDP** is described by a model  $\mathcal{P} = (\mathcal{S}, \mathcal{A}, \mathcal{O}, T, R, O, P, \gamma)$ .

- States  $s_t \in \mathcal{S}$ ,
- Actions  $a_t \in \mathcal{A}$ ,
- Observations  $o_t \in \mathcal{O}$ ,
- Discount  $\gamma \in [0, 1)$ ,

- Transition  $T(s_{t+1} \mid s_t, a_t)$ ,
- Reward  $r_t \sim R(\cdot \mid s_t, a_t)$ ,
- Perception  $O(o_t \,|\, s_t),$
- Initialisation  $P(s_0)$ .

# **History-dependent policies**



The **history** at time t is  $h_t = (o_0, a_0, ..., o_t) \in \mathcal{H}.$ 

#### Definition 1: History-dependent policy.

A history-dependent policy  $\eta \in \mathbf{H} = \mathcal{H} \to \Delta(\mathcal{A})$  is a mapping from histories to distributions over the actions, with density  $\eta(a \mid h)$ .

# **History-dependent policies**



# Optimal control under partial observability



The problem of RL in POMDP is to find an optimal history-dependent policy,

$$\eta^* \in \operatorname*{argmax}_{\eta \in \mathbf{H}} \underbrace{\mathbb{E}^{\eta} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \right]}_{J(\eta)},$$

from samples  $(o_0, a_0, r_0, ..., o_{t-1}, a_{t-1}, r_{t-1}, o_t)$ .

# Q-function of a policy

Definition 2: Q-function of a policy.

The **Q-function of a policy**  $\eta \in \mathbf{H}$ , or **critic**, is defined as,

$$Q^{\eta}(h,a) = \mathbb{E}^{\eta} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \; \middle| \; H_0 = h, A_0 = a \right].$$



 $\mathbf{NB}: J(\eta) = \mathbb{E}[Q^{\eta}(H_0, A_0)].$ 



Definition 3: Q-function.

The optimal Q-function, or simply Q-function, is defined as,

$$Q(h,a) = \max_{\eta \in \mathcal{H}} Q^{\eta}(h,a) \big( = Q^{\eta^*}(h,a) \big).$$



 $\mathbf{NB}: \operatorname{Supp}(\eta^*(\cdot \mid h)) \subseteq \operatorname{argmax}_{a \in \mathcal{A}} Q(h, a).$ 

# History-dependent reinforcement learning

We use function approximators for the policy or Q-function estimation.

• Sliding window, transformer, recurrent networks, etc.



Fig. 2: Q-function approximator.

# History-dependent reinforcement learning

We use function approximators for the policy or Q-function estimation.

• Sliding window, transformer, recurrent networks, etc.



Fig. 2: Q-function approximator.
# History-dependent reinforcement learning

We use **function approximators** for the **policy** or **Q-function** estimation.

• Sliding window, transformer, recurrent networks, etc.



Fig. 2: Q-function approximator.

# **Reinforcement learning techniques**

In this thesis, we considered three main reinforcement learning techniques:

1. **Policy-gradient** (and actor-critic) methods:

 $\nabla_{\psi} J \big( \eta_{\psi} \big) \approx \mathbb{E}^{d^{\eta_{\psi}}} \big[ Q_{\varphi}^{\eta_{\psi}}(H, A) \nabla_{\psi} \log \eta_{\psi}(A \,|\, H) \big]$ 

2. Value-based methods (Q-learning):

$$Q_{\varphi}(h,a) \approx \mathbb{E} \Big[ R + \gamma \max_{a' \in \mathcal{A}} Q_{\varphi}(H',a') \, \Big| \, H = h, A = a \Big]$$

3. Model-based methods (supervised learning):

 $q_{\theta}(r,o'\,|\,h,a)\approx \Pr(r,o'\,|\,h,a)$ 

# Outline

# Outline

- I. Learning and Remembering
  - Learning for Remembering
  - Remembering for Learning
- II. Leveraging Additional Information
  - Sufficiency through Additional Information
  - Learning Faster with Additional Information
- III. Entangling Predictions and Decisions
  - Rolling the Dice First
  - Just Looking at the Dice

# Outline

- I. Learning and Remembering
  - Learning for Remembering
  - Remembering for Learning
- II. Leveraging Additional Information
  - Sufficiency through Additional Information
  - Learning Faster with Additional Information
- III. Entangling Predictions and Decisions
  - Rolling the Dice First
  - Just Looking at the Dice

# I. Learning and Remembering

### I.1. Learning for Remembering

Published in Transactions on Machine Learning Research (08/2022)

### Recurrent networks, hidden states and beliefs in partially observable environments

Gaspard Lambrechts Montefiore Institute, University of Liège

Adrien Bolland Montefiore Institute, University of Liège gaspard. lambrechts Guliege. be adrien. bolland Guliege. be

dernstØulieae.be

Damien Ernst Montefiore Institute, University of Liège LTCI, Telecom Paris, Institut Polytechnique de Paris

Reviewed on OpenReview: https://openreview.net/forum?id=dkHfV3uB2l

#### Abstract

Reinforcement learning aims to learn optimal policies from interaction with environments whose dynamics are unknown. Many methods rely on the approximation of a value function to derive near-optimal policies. In partially observable environments, these functions depend on the complete sequence of observations and past actions, called the history. In this work, we show empirically that recurrent neural networks trained to approximate such value functions internally filter the posterior probability distribution of the current state given the history, called the belief. More precisely, we show that, as a recurrent neural network learns the Q-function, its hidden states become more and more correlated with the beliefs of state variables that are relevant to optimal control. This correlation is measured through their mutual information. In addition, we show that the expected return of an agent increases with the ability of its recurrent architecture to reach a high mutual information between its hidden states and the beliefs. Finally, we show that the mutual information between the hidden states and the beliefs of variables that are irrelevant for optimal control decreases through the learning process. In summary, this work shows that in its hidden states, a recurrent neural network approximating the Q-function of a partially observable environment reproduces a sufficient statistic from the history that is correlated to the relevant part of the belief for taking optimal actions.

#### 1 Introduction

Latest advances in reinforcement learning (RL) rely heavily on the ability to approximate a value function (i.e., state or state-action value function). Modern RL algorithms have been shown to be able to produce amproximation of the value functions of Markow devices approximation processor which high could use the value of the value realized of the value of the value functions of the value functions of the value functions of the value functions.

# I.1. Learning for Remembering



# **Belief sufficiency**

**Definition 4:** Belief of a history. The belief b = f(h) of a history  $h \in \mathcal{H}$  is defined as,  $b(s) = \Pr(s \mid h).$ 

Theorem 1: Belief recurrence.

$$f(h') = u(f(h), a, o').$$

Theorem 2: Belief sufficiency.

$$Q(h,a) = Q'(f(h),a).$$



 $\Rightarrow$  If the belief is known, we can discard the history. But it is usually not known.

## **Recurrent reinforcement learning**

In **recurrent Q-learning**, a recurrent approximator learns the Q-function.

- The RNN state  $z = f_{\theta}(h)$  is a statistic of the history (memory).



 $\Rightarrow$  Should RNN states encode the belief?

## **Recurrent reinforcement learning**

We showed that **beliefs emerge** in RNN states during model-free recurrent Q-learning (LSTM, GRU, BRC, NBRC, MGU).



Fig. 5: Return and mutual information throughout training.

We estimate  $\hat{I}(\theta) \approx I(z, b)$  under stationary distribution  $p^{\eta_{\theta}}(h)$ .

# **Recurrent reinforcement learning**

The belief of irrelevant state variable is not encoded in RNN states (GRU).



Fig. 6: Return and mutual information (relevant or irrelevant belief).

We estimate  $\hat{I}(\theta) \approx I(z, b)$  under stationary distribution  $p^{\eta_{\theta}}(h)$ .

The statistic of the history encodes the belief of relevant state variables.

### I.2. Remembering for Learning

#### Neural Networks 166 (2023) 645-665



Contents lists available at ScienceDirect

Neural Networks



#### journal homepage: www.elsevier.com/locate/neunet

#### Warming up recurrent neural networks to maximise reachable multistability greatly improves learning



#### Gaspard Lambrechts<sup>a,\*,1</sup>, Florent De Geeter<sup>a,1</sup>, Nicolas Vecoven<sup>a,1</sup>, Damien Ernst<sup>a,b</sup>, Guillaume Drion<sup>a</sup>

<sup>a</sup> Monteflore Institute, University of Liège, 10 allée de la découverte, Liège, 4000, Belgium <sup>b</sup> LTCL Telecom Paris, Institut Polytechnique de Paris, 19 place Marguerite Perey, Palaiseau, 91120, France

#### ARTICLE INFO

#### ABSTRACT

Article history: Received 27 July 2022 Received in revised form 12 June 2023 Accepted 14 July 2023 Available online 7 August 2023

Repriverds: Recurrent neural network Multistability Initialisation procedure Long-term memory Warmup Long time dependencies Training recurrent neural networks is known to be difficult when time dependencies become long. In this work, we show that most standard cells only have one stable equilibrium at initialisation, and that learning on tasks with long time dependencies generally occurs once the number of network stable equilibria increases; a property known as multistability. Multistability is often not easily attained by initially monostable networks, making learning of long time dependencies between inputs and outputs difficult. This insight leads to the design of a novel way to initialise any recurrent cell connectivity through a procedure called "warmup" to improve its capability to learn arbitrarily long time dependencies. This initialisation procedure is designed to maximise network reachable multistability, i.e., the number of equilibria within the network that can be reached through relevant input trajectories, in few gradient steps. We show on several information restitution, sequence classification, and reinforcement learning benchmarks that warming up greatly improves learning speed and performance, for multiple recurrent cells, but sometimes impedes precision. We therefore introduce a double-laver architecture initialised with a partial warmup that is shown to greatly improve learning of long time dependencies while maintaining high levels of precision. This approach provides a general framework for improving learning abilities of any recurrent cell when long time dependencies are present. We also show empirically that other initialisation and pretraining procedures from the literature implicitly foster reachable multistability of recurrent cells.

© 2023 Elsevier Ltd. All rights reserved.

#### 1. Introduction

Despite their performances and videspread use, recurrent neural networks (NNs) are known to be blackbox models with extremely complex internal dynamics. A growing body of work has focused on understanding the internal dynamics of trained contain. Canzult. & Sarsillo, 2019; Sarsillo & Bankz, 2013) perduding invaluable Sarsillo, 2019; Sarsillo & Bankz, 2013) privileng invaluable intuition into the RMN prediction process. This viewpoint has already been used to understand the difficulties. FRNNs to capture longer time dependencies (Bengio, Frascon, 2013). This line of work has argued that locating such fixed points efficiently could provide insights into RNN vhammics and input-output properties. Here, we build upon this line of work by studying the impact of the number of reachable fixed points in an RNN on the ability to learn long time dependencies. Moreover, we highlight how maximising the number of reachable fixed points at initialisation can improve RNN learning, in particular in the presence of arbitrarily long dependencies.

More precisely, we introduce a fast-to-compute measure of the multistability of a network called variability amongst attractors (VAA). This measure gives the number of reachable attractors

# I.2. Remembering for Learning



# **II. Leveraging Additional Information**

### **II.1. Sufficiency through Additional Information**

RLJ | RLC 2024

#### Informed POMDP: Leveraging Additional Information in Model-Based RL

Gaspard Lambrechts gaspard.lambrechts@uliege.be Montefiore Institute, University of Liège Adrien Bolland adrien.bolland@uliege.be Montefiore Institute, University of Liège

#### Damien Ernst

dernst@uliege.be Montefiore Institute, University of Liège LTCI, Télécom Paris, Institut Polytechnique de Paris

#### Abstract

In this work, we generalize the problem of learning through interaction in a POMDP by accounting for eventual additional information available at training time. First, we introduce the informed POMDP, a new learning paradigm offering a clear distinction between the information at training and the observation at execution. Next, we propose an objective that leverages this information for learning a sufficient statistic of the history for the optimal courtor. We then adapt this informed objective to learn a world model able to sample latent trajectories. Finally, we empirically word model in the Dreawer algorithm. These reasting and the simplicity of the proposed adaptation advocate for a systematic consideration ad eventual additional information when learning in a POMDP using model-based RL.

#### 1 Introduction

Reinforcement learning (RL) aims to learn to act optimally through interaction with neuroimmetrs whose dynamics are unknown. A mapping challenge in this field is partial observability, where only a partial observability of the Markovian state of the environment is is available for taking action a taking action at our formation of a partial bole parabol Markov desion process (POMDP). In this context, an optimal policy  $\eta(a|b)$  generally depends on the history h of all observations and previous actions, which grows the incord with the relative theoretically possible to find a statistic f(b) of the history h that is updated recurrently and that summarizes all relevant control. Formally, a statistic f(b) a statistic f(b) a statistic f(b) as statistic to  $f(b) = \eta(d)$ , and each time a action a is taken and a new observation of is received, with h' = (h, a, c). And statistic f(b) is sufficient for the output observables an output and point g(a(b)). And statistic f(b) is sufficient for the output output when the excitation action  $f(b) = \eta(d) = \eta(d)$ .

### **II.1. Sufficiency through Additional Information**



# A story of partial observability



## **Classical POMDP**



The problem of RL in POMDP is to find an optimal history-dependent policy

$$\eta^* \in \operatorname*{argmax}_{\eta \in \mathcal{H}} \mathbb{E}^{\eta} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \right]$$

from samples  $(o_0, a_0, r_0, ..., o_t)$ .

28/49

### **Informed POMDP**



The problem of RL in POMDP is to find an optimal history-dependent policy

$$\eta^* \in \operatorname*{argmax}_{\eta \in \mathcal{H}} \mathbb{E}^{\eta} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \right]$$

from samples  $(\mathbf{i_0}, o_0, a_0, r_0, ..., \mathbf{i_t}, o_t)$ .

# Sufficiency for optimal control



The history *h* is compressed into a statistic *z* by a function *f*.  $\Rightarrow$  It should summarize all relevant information to act optimally.

Definition 5: Sufficiency for optimal control.

A statistic  $f:\mathcal{H}\to\mathcal{Z}$  is sufficient for optimal control if, and only if,  $\max_g J(g\circ f)=\max_\eta J(\eta).$ 

For example, the **belief** is a sufficient statistic.

# Sufficiency in an informed POMDP

Theorem 3: Sufficiency of recurrent predictive statistics.

- A statistic  $f:\mathcal{H}\to\mathcal{Z}$  is sufficient for optimal control if it is,
- (i) **recurrent**: f(h') = u(f(h), a, o'),
- (ii) predictive:  $p(r, i' \mid h, a) = p(r, i' \mid f(h), a)$ .

It motivates the **informed world model**  $q(r, i' \mid f(h), a)$  objective:

$$\max_{f,q} \mathop{\mathbb{E}}_{p(r,i'\mid h,a)} q(r,i'\mid f(h),a).$$



**Fig. 8**: Statistic z = f(h) of the history *h* encoding the transition distribution.

# Informed world model

The Informed Dreamer uses a variational recurrent neural network.

- Prior  $\hat{e} \sim q^e(\cdot \mid z, a)$
- Information  $\hat{\imath} \sim q^i(\cdot \,|\, z, \hat{e})$ 
  - + Instead of observation  $\hat{o} \sim q^o(\cdot \,|\, z, \hat{e})$
- Reward  $\hat{r} \sim q^r(\cdot \,|\, z, \hat{e})$
- Encoder  $e \sim q^e(\cdot \,|\, z, a, o')$
- Update z' = u(z, a, e)



### Learning in imagination

From the learned world model, we can learn a latent policy  $g_{\varphi}: \mathcal{Z} \to \Delta(\mathcal{A}).$ 



**Fig. 9**: Imagining trajectories with the Informed Dreamer and latent policy  $g_{\omega}$ .

### **Informed Dreamer**



Fig. 10: Varying Mountain Hike

### **Informed Dreamer**



Fig. 11: Velocity DeepMind Control

### **Informed Dreamer**



Fig. 12: Pop Gym

Informed world models provide faster policy learning than observational ones.

### **II.2. Learning Faster with Additional Information**

#### A Theoretical Justification for Asymmetric Actor-Critic Algorithms

Gaspard Lambrechts<sup>1\*</sup> Damien Ernst<sup>1</sup> Aditya Mahajan<sup>2</sup>

#### Abstract

In reinforcement learning for partially observable environments, many successful algorithms were devolped within the asymmetric learning pardigm. This pardigm leverages additional state information available at training time for infester learning. Although the proposed learning objectives are usually theoretically sound, these methods still lack a theoretical justication for their potential benefits. We propose such a justification for asymmetric actor-relification infestion convergence analysis but is setting. The resulting information to the same state of the resulting information are not remaining from aliasing in the agent state.

#### 1. Introduction

Reinforcement learning (RL) is an appealing framework for solving decision making problems, notably because it makes very few assumptions about the problem at hand. In its purest form, the promise of an RL algorithm is to learn an optimal behavior from interaction with an environment whose dynamics are unknown. More formally, an RL algorithm aims at learning a policy (i.e., a mapping from observations to actions) in order to maximize a reward signal from samples obtained by interacting with an environment. While RL has offered empirical successes for a plethora of challenging problems ranging from games to robotics (Mnih et al., 2015; Schrittwieser et al., 2020; Levine et al., 2015; Akkaya et al., 2019), most of these achievements have assumed full observability. A more realistic assumption is partial observability, where only a partial observation of the state of the environment is available for history-dependent policies, usually using a recurrent neural network to process historise (Blaker, 2001). Wierstra et al., 2010; Hausknecht & Stone, 2015; Heess et al., 2015; Zhang et al., 2016; Zhu et al., 2017; Chren ted diffuelty of learning effective history-dependent policies, various auxiliary representation learning objective bash bee hen proposed to compress the history into useful representations (Eql et al., 2018; Buessing et al., 2018; Guote et al., 2018; Group et al., 2019; Han et al., 2019; Hanf et al., 2019; Guote et al., 2020; Lee et al., 2020; Shummain et al., 2022; Ni et al., 2023; Lee et al., 2020; Shummain et al., 2022; Ni et al., 2023;

While these methods are theoretically able to learn optimal history-dependent policies, they learn solely the partial state observations, which can be too restrictive. Indeed, assuming the same partial observability at training time and execution time can be too pessimistic for many environments, notably for those that are simulated. This motivated the asymmetric learning paradigm, where additional state information available at training time is leveraged during the learning process of the history-dependent policy. Although the optimal policies obtained by asymmetric learning are theoretically equivalent to those learned by symmetric learning, the promise of asymmetric learning is to improve the convergence speed towards a near-optimal policy. Early approaches notably proposed to imitate a privileged policy conditioned on the state (Choudhury et al., 2018), or to use an asymmetric critic conditioned on the state (Pinto et al., 2018). These heuristic methods initially lacked a theoretical framework, and a recent line of work has focused on proposing theoretically grounded asymmetric learning objectives. First, imitation learning of a privileged policy was known to be suboptimal, and it was addressed by constraining the privileged policy so that its imitation results in an optimal policy for the partially observable environment (Warrington et al., 2021), Similarly, asymmetric actor-critic approaches were proven to provide biased gradients, and an

### **II.2. Learning Faster with Additional Information**



### **II.2. Learning Faster with Additional Information**



### Asymmetric actor-critic algorithm

Actor-critic algorithms are policy-gradient methods with a critic  $Q_{\varphi}^{\eta_{\psi}} \approx Q^{\eta_{\psi}}$ .

- The critic is **only used** for estimating the policy-gradient.
- It can be **informed** with additional information:  $Q(h, a) \rightarrow Q(h, i, a)$ .



 $\Rightarrow$  Very effective, but no theoretical justification for its benefits.

# Asymmetric actor-critic algorithm

- State-informed:
  - We study the case where i = s.
- Fixed statistic:
  - Fixed update  $z' \sim U(\cdot \mid z, a, o')$ , and policy  $a \sim \pi(\cdot \mid z)$ .
- Finite state Q-functions:
  - Asymmetric  $\mathcal{Q}^{\pi}(s, z, a)$  and symmetric  $Q^{\pi}(z, a)$ .
- Linear approximations:
  - $\bullet \ \hat{Q}^{\pi}_{\beta}(s,z,a) = \langle \beta, \varphi(s,z,a) \rangle \text{ and } \hat{Q}^{\pi}_{\beta}(z,a) = \langle \beta, \chi(z,a) \rangle.$
  - $\bullet \ \pi_{\theta}(a \,|\, z) \propto \exp(\langle \theta, \psi(z,a) \rangle).$

### Actor-critic algorithm

Algorithm 1: Asymmetric and symmetric actor-critic.

- 1. Initialize policy parameters  $\psi_0$ .
- 2. For t = 1...T
  - 1. Estimate  $\hat{\mathcal{Q}}^{\pi}_{\varphi} \approx \mathcal{Q}^{\pi_{\psi}}$  or  $\hat{\mathcal{Q}}^{\pi}_{\chi} \approx \mathcal{Q}^{\pi_{\psi}}$  (**TD learning**).
  - 2. Estimate  $g_{t-1} \approx \nabla_{\psi} J(\pi_{\psi_{t-1}})$  using  $\mathcal{Q}_{\varphi}$  or  $Q_{\chi}$  (NPG estimation).
  - 3. Update policy  $\psi_t = \psi_{t-1} + \eta g_{t-1}$ .
- 3. Return  $\pi_{\psi_T}$

From the belief  $b(s | h) = \Pr(s | h)$  and approximate belief  $\hat{b}(s | z) = \Pr(s | z)$ , we introduce a **measure of the aliasing** of the agent state z.

Aliasing measure.

$$\varepsilon_{\rm alias} \propto \mathbb{E} \big[ \big\| b(\cdot \,|\, h) - \hat{b}(\cdot \,|\, z) \big\| \big].$$

### Finite-time bound for the critics

**Theorem 4:** Finite-time bound for asymmetric and symmetric Q-functions. For any  $\pi \in \Pi_{\mathcal{M}}$ , and any  $m \in \mathbb{N}$ , we have for TD learning with  $\alpha = \frac{1}{K}$ ,

$$\begin{split} &\sqrt{\mathbb{E}\Big[\left\|\mathcal{Q}^{\pi}-\overline{\mathcal{Q}}^{\pi}\right\|_{d^{\pi}}^{2}\Big]} \leq \varepsilon_{\mathrm{td}}+\varepsilon_{\mathrm{app}}+\varepsilon_{\mathrm{shift}} \\ &\sqrt{\mathbb{E}\Big[\left\|\mathcal{Q}^{\pi}-\overline{\mathcal{Q}}^{\pi}\right\|_{d^{\pi}}^{2}\Big]} \leq \varepsilon_{\mathrm{td}}+\varepsilon_{\mathrm{app}}+\varepsilon_{\mathrm{shift}}+\varepsilon_{\mathrm{alif}} \end{split}$$

$$\begin{split} \varepsilon_{\mathrm{td}} &= \sqrt{\frac{4B^2 + \left(\frac{1}{1-\gamma} + 2B\right)^2}{2\sqrt{K}(1-\gamma^m)}} \\ \varepsilon_{\mathrm{app}} &= \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_{\varphi}^B} \|f - Q^{\pi}\|_{d^{\pi}} \\ \varepsilon_{\mathrm{shift}} &= \left(B + \frac{1}{1-\gamma}\right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m^{\pi} \otimes \pi - d^{\pi} \otimes \pi\|_{\mathrm{TV}}} \\ \varepsilon_{\mathrm{alias}} &= \frac{2}{1-\gamma} \left\| \mathbb{E}^{\pi} \left[ \sum_{k=0}^{\infty} \gamma^{km} \left\| \hat{b}_{km} - b_{km} \right\|_{\mathrm{TV}} \right| Z_0 = \cdot, A_0 = \cdot \right] \right\|_{d^{\pi}}. \end{split}$$
### Finite-time bound for the actors

**Theorem 5:** Finite-time bound for asymmetric and symmetric NAC. For any  $(\mathcal{Z}, U)$ , we have for NAC with  $\alpha = \frac{1}{K}$ ,  $\zeta = \frac{B\sqrt{1-\gamma}}{\sqrt{2N}}$ ,  $\eta = \frac{1}{\sqrt{T}}$ ,  $(1-\gamma)\min_{0 \le t < T} \mathbb{E}[J(\pi^*) - J(\pi_t)] \le \varepsilon_{\text{nac}} + \varepsilon_{\text{actor}} + \varepsilon_{\text{inf}} + \varepsilon_{\text{grad}} + \frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t}$ ,

$$\begin{split} \varepsilon_{\rm nac} &= \frac{B^2 + 2 \log |A|}{2 \sqrt{T}} \\ \varepsilon_{\rm actor} &= \overline{C}_{\infty} \sqrt{\frac{(2 - \gamma)B}{(1 - \gamma)\sqrt{N}}} \\ \varepsilon_{\rm inf,asym} &= 0 \quad \varepsilon_{\rm inf,sym} = 2 \mathbb{E}^{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k \left\| \hat{b}_k - b_k \right\|_{\rm TV} \right] \\ \varepsilon_{\rm grad,asym} &= 2 \overline{C}_{\infty} \sup_{0 \leq t < T} \sqrt{\min_w \mathcal{L}_t(w)} \quad \varepsilon_{\rm grad,sym} = 2 \overline{C}_{\infty} \sup_{0 \leq t < T} \sqrt{\min_w L_t(w)} \\ \varepsilon_{\rm critic,asym}^{\pi_t} &= 2 \overline{C}_{\infty} \sqrt{6} \big( \varepsilon_{\rm td} + \varepsilon_{\rm app} + \varepsilon_{\rm shift} \big) \quad \varepsilon_{\rm critic,sym}^{\pi_t} = 2 \overline{C}_{\infty} \sqrt{6} \big( \varepsilon_{\rm td} + \varepsilon_{\rm app} + \varepsilon_{\rm shift} \big) \end{split}$$

Asymmetric learning is insensitive to aliasing in the statistic of the history.

# III. Entangling Predictions and Decisions

#### **III.1. Rolling the Dice First**

ICML 2024 Next Generation of Sequence Modeling Architectures Workshop

#### Parallelizing Autoregressive Generation with Variational State Space Models

Gaspard Lambrechts\* Yann Claes\* Pierre Geurts Damien Ernst Montefiore Institute, University of Liège GASPARD.LAMBRECHTS @ ULIEGE.BE Y.CLAES @ ULIEGE.BE P.GEURTS @ ULIEGE.BE DERNST @ ULIEGE.BE

#### Abstract

Attention-based models such as Transformers and recurrent models like state space models (SSMs) have emerged as successful methods for antoregressive sequence modeling. Although both enables parallel training, none enable parallel generation due to their autoregressive ness. We propose the variational SSM (VSSMs) a variational autorescoler (VAS) where both the encoder and decoder and SSMs. Since sampling the latent variables and decoding them with the SSM can be parallelized, both training and gSM exercision can be conducted in parallel. Moreover, the decoder and decoder are sequenced and the sequence of the structure of the sequence of the structure and the sequence framily, we propose the autoreanguage generation task. Interestingly, the autoregression VSM still enables parallel generation. We highlight on top problems (MMST, CIFAR) the empirical gains in speed-up and show that it competes with relational models in terms of generation quality (Transformer, Namhas SSM).

Keywords: Parallel, Autoregressive, Generation, VAE, SSM, VSSM

#### 1. Introduction

Sequence modeling tasks, namely time-series forecasting and text generation, have gained in populativ and various types of architectures were dosigned to tackle such problems. Transformers were proven effective [17, 19], yet they nonetheless reprocess the complete sequence at each timestep, making generation less efficient. Recurrent neural networks (RNNs) [3, 8] update a hidden state based on new impacts at each timestep, enabling efficient generation. SSMs [9-11, 18], a recently introduced class of RNNs, enable parallel training thanks to their linear recurrence. Alternatively, sevend works adapt VAEs for sequenting modeling. Some architectures integrate Transformers [13, 14] and enable parallel training, although little work [5] proposes models that can be conditioned on parial relatizations (e.g., prompts). Conversely, variational RNNs (VRNNs) [4] loose parallelizabi-

# **III.1. Rolling the Dice First**











#### **III.2.** Just Looking at the Dice



ames, remension cu sociaes cuisis, inagua: Done cu phrus. Quisque venction, una esci utilizes autôta; pede lorem egesta dui, et consullis elit erat sel nulla. Done cu hcuts. Curabiur et nuac. Alignan dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonumny in, fermentum faucibus, egestas vel, odio, Lambrechts et al., 2024).

Etian eusimod. Fasce facilisis lacinia dui. Suspendisse potenti. In ni erat, cursus id, nonumny sect. ullancorper eges, sapien. Praesear pretum, magan ie neferind egestas, pede pede pretimi lorem, quis consecteture tortor aspien facilisis magan. Mauris quis magan varius nulla scelerisque imperdiet. Aliquam non quan. Allequam portitori quan a lacus. Praesearti vel arcu ut toric cursus voltari. In vitae pede quis diam bibendum placerat. Fasce elementum convalin ence. Sed dolor orei, scelerisque ac, dapibus nec, utiricies ut. mi. Diais nec dui quis los aguittis commodo.

Nulla malesanda portitior diam. Done felis cerat, congue non, voltapat at, inicidunt tristique, libers. Vivoums vierrar fermentum felis. Dones communy pelleturospece aute. Phasello solipsicing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, piacent a, molestie nec, los. Macenna lactina, Nun ipsun ligada, diefend at, accentum ance, saccepia a, pineum. Morbi blaudit ligada fengida mugaa. Nune emin. Prasent enismod mune en purus. Dance biberdum quam in tellus. Nullan carnas pulvinar lectus. Dance et mi. Nam volpatate mette se enim. Weitsbum pelletetsper felis en massa.

## **III.2.** Just Looking at the Dice



# A Matter of Abstractions

## Conclusion

Instead of the history MDP, we should consider the **structure of the solution**, focusing on **recurrent approximators** that are **predictive of the future**.



- I. Good initial memory learning to focus on the relevant belief,
- II. Predictive of the state, or can be augmented with state at training,
- III. Learned as **predictive latent** of the future, sampled and decoded in parallel.

## **Future perspectives**

We developed techniques to **learn an optimal behavior** for any **given task**.  $\rightarrow$  Could we instead **learn a behavior** that optimally **generalize to all tasks**?

Generalization is equivalent to optimal control under partial observability!

