

# Inria TAU Seminar

## Learning-based Control on Dynamical Systems

---

R. Hosseinkhan Boucher<sup>1</sup>, S. Douka<sup>1</sup>, O. Semeraro<sup>1</sup>, L. Mathelin<sup>1</sup>

Laboratoire Interdisciplinaire des Sciences du Numérique, Université Paris-Saclay, CNRS

**Doctoral School:** *Sciences et technologies de l'information et de la communication (STIC)*

Granted by the Agence Nationale de la Recherche (ANR) under projet ANR-21-CE46-0008 Reinforcement Learning as Optimal control for Shear Flows (REASON)



# **Dynamical Systems Control: Challenges**

---

# Challenges in Dynamical Systems Control

## Optimal Control Problem

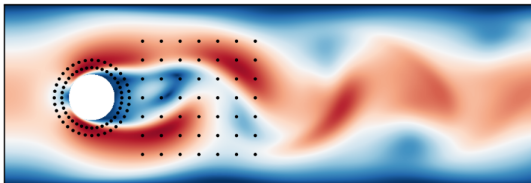
**Dynamics:**  $\partial_t x(z, t) = P[x, u](z, t)$

**Objective:**  $\min_u J(u) = \int_0^T c(x(t), u(t)) dt$

## Example

$P$  is the Navier-Stokes operator

Energy criterion:  $c(x, u) = \|x\|^2 + \|u\|^2$



Cylinder flow drag reduction. Partial observation through sensors.

## Challenges<sup>1</sup>

- Partial observability (PO) and delays
- Controllability
- Sampling complexity
- Robustness
- High dimensional hidden state space  $\mathcal{X}$
- Extremely large degrees of freedom (sensor placement, actuators, amplitude, optimization problem). No benchmark

## Rigorously

- Control problem with **continuous** time and **infinite** state space (Relaxed Stochastic Control)

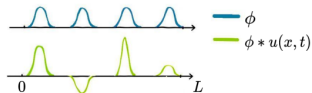
<sup>1</sup>J. Viquerat et al. "A review on deep reinforcement learning for fluid mechanics: An update", AIP Publishing (2022)

# Controlled Kuramoto-Sivashinsky (KS)<sup>1,2</sup>

**Controlled KS:**  $\partial_t x(z, t) + x(z, t) \partial_x x(z, t) = -\partial_x^2 x(z, t) - \partial_x^4 x(z, t) + \langle \phi, \mathbf{u} \rangle(z, t)$

$$x(z + L, t) = x(z, t) \text{ and } (z, t) \in [0, L] \times [0, T]$$

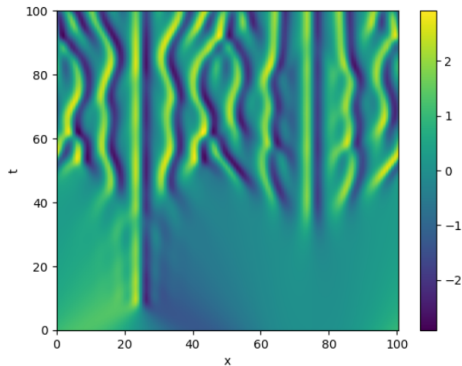
**Control term:**  $\langle \phi, \mathbf{u} \rangle = \sum_{i=1}^r u_i f_{\mathcal{N}}(\mu_i, \sigma^2)$



$\phi$  define a given gaussian mixture,  $\mathbf{u}$  is **unknown**

## Properties

- Spatio-temporal chaos, 4th order non-linear
- Equilibria, relative equilibria, symmetries
- 4 equilibria  $x_e^0(z) = 0$ ,  $x_e^1(z)$ ,  $x_e^2(z)$ ,  $x_e^3(z)$



Evolution of the Kuramoto-Sivashinsky equation with  $L = 100$

<sup>1</sup>Y. Kuramoto. "Diffusion-Induced Chaos in Reaction Systems", *Progress of Theoretical Physics Supplement* (1978)

<sup>2</sup>G.I. Sivashinsky. "Nonlinear analysis of hydrodynamic instability in laminar flames—I. Derivation of basic equations", *Acta Astronautica* (1977)

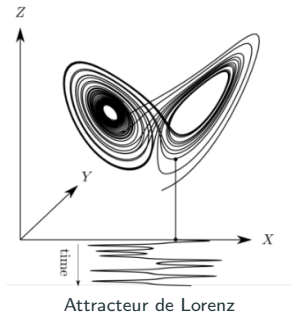
# Controlled Lorenz<sup>1</sup>

**Controlled Lorenz:** 
$$\begin{cases} \partial_t x_1 = \sigma(x_2 - x_1) + u_1 \\ \partial_t x_2 = x_1(\rho - x_3) - x_2 + u_2 \\ \partial_t x_3 = x_1 x_2 - \beta x_3 + u_3 \end{cases}$$

**Control Term:**  $u = (u_1, u_2, u_3)$

## Properties

- Chaos, instabilities, symmetries
- Equilibria  $x_e^0, x_e^1, x_e^2$



<sup>1</sup>T. L. Vincent, J. Yu. "Control of a chaotic system", *Dynamics and Control* (1991)

# Partially Observable Markov Decision Process (POMDP)

## Dynamics

$$\partial_t x(z, t) = P[x, u](z, t), \quad x(\cdot, t) \in \mathbb{L}^2(\mathcal{X}) \text{ and } u(\cdot, t) \in \mathbb{L}^2(\mathcal{U}) \text{ for any } t \in [0, T]$$

## Spatial Discretisation

$$\mathbb{L}^2(\mathcal{X}) \simeq \mathcal{X}^{d_x} \quad \mathbb{L}^2(\mathcal{U}) \simeq \mathcal{U}^{d_u}$$

## Temporal Discretisation

$$[0, T] \simeq (k\delta)_{0 \leq k \leq n}$$

Continuous operator  $\longrightarrow$  Discrete<sup>1</sup> operator:  $x_{k+1} = P(x_k, u_k)$ ,  $x_k \in \mathcal{X}^{d_x}$ ,  $u_k \in \mathcal{U}^{d_u}$

---

<sup>1</sup>The same notations (operator, time horizon etc.) as the continuous time framework will be used for the discrete time framework.

# Partially Observable Markov Decision Process (POMDP)

## Dynamics

$$\partial_t x(z, t) = P[x, u](z, t), \quad x(\cdot, t) \in \mathbb{L}^2(\mathcal{X}) \text{ and } u(\cdot, t) \in \mathbb{L}^2(\mathcal{U}) \text{ for any } t \in [0, T]$$

## Spatial Discretisation

$$\mathbb{L}^2(\mathcal{X}) \simeq \mathcal{X}^{d_x} \quad \mathbb{L}^2(\mathcal{U}) \simeq \mathcal{U}^{d_u}$$

## Temporal Discretisation

$$[0, T] \simeq (k\delta)_{0 \leq k \leq n}$$

Continuous operator  $\longrightarrow$  Discrete<sup>1</sup> operator:  $x_{k+1} = P(x_k, u_k)$ ,  $x_k \in \mathcal{X}^{d_x}$ ,  $u_k \in \mathcal{U}^{d_u}$

## Generalisation: Partially Observable Markov Decision Process (POMDP)

$$\begin{aligned} X_{k+1} &= P(X_k, U_k, \eta_k) & \eta_k &\sim \mathcal{N}(0, \sigma_\eta^2 I_d) \\ Y_{k+1} &= Q(X_k) + \epsilon_k & \epsilon_k &\sim \mathcal{N}(0, \sigma_\epsilon^2 I_d) \end{aligned} \tag{1}$$

with  $X_0 \sim \mathcal{N}(x_e, \sigma_e^2 I_d)$ .

---

<sup>1</sup>The same notations (operator, time horizon etc.) as the continuous time framework will be used for the discrete time framework.

# Modeling as a Markov Decision Process (MDP)

State space  $\mathcal{X}$ , control space  $\mathcal{U}$ , observation space  $\mathcal{Y}$

## Random Dynamics

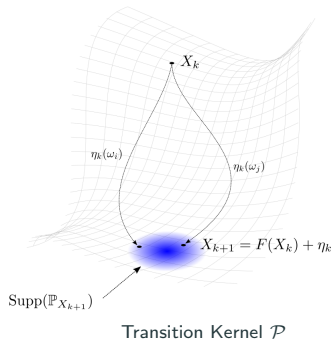
$\mathcal{P}(dx_{k+1} \mid (x_k, u_k)) \rightarrow$  probability on  $\mathcal{X}$  given  $(x_k, u_k) \in \mathcal{X} \times \mathcal{U}$

## Random Observation

$\mathcal{Q}(dy_k \mid x_k) \rightarrow$  probability on  $\mathcal{Y}$  given  $x_k \in \mathcal{X}$

## Random Control

$\pi(du_k \mid y_k) \rightarrow$  probability on  $\mathcal{U}$  given  $y_k \in \mathcal{Y}$





# Modeling as a Markov Decision Process (MDP)

State space  $\mathcal{X}$ , control space  $\mathcal{U}$ , observation space  $\mathcal{Y}$

## Random Dynamics

$\mathcal{P}(dx_{k+1} | (x_k, u_k)) \rightarrow$  probability on  $\mathcal{X}$  given  $(x_k, u_k) \in \mathcal{X} \times \mathcal{U}$

## Random Observation

$\mathcal{Q}(dy_k | x_k) \rightarrow$  probability on  $\mathcal{Y}$  given  $x_k \in \mathcal{X}$

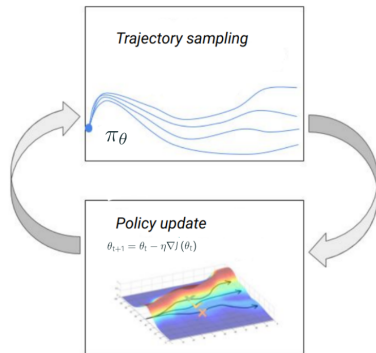
## Random Control

$\pi(du_k | y_k) \rightarrow$  probability on  $\mathcal{U}$  given  $y_k \in \mathcal{Y}$

## Controlled Hidden Markov Chain

$$P^\pi(dx_0 du_0 dy_0 dx_1 du_1 \dots dx_T) = P_{X_0}(dx_0) \mathcal{Q}(dy_0 | x_0) \pi(du_0 | y_0) \mathcal{P}(dx_1 | x_0, u_0)$$

$$\mathcal{Q}(dy_1 | x_1) \pi(du_1 | y_1) \dots \pi(du_{T-1} | y_{T-1}) \mathcal{P}(dx_T | x_{T-1}, u_{T-1})$$



Policy gradient iterations to solve  
 $\arg \min_{\pi} \mathbb{E}^\pi [\sum_{k=0}^T c(X_k, U_k)]$

## **Maximum Entropy: Noise Robustness**

---

# Robustness: Maximum Entropy and Flat Minima

## Maximum Entropy in Reinforcement Learning

$$\arg \min_{\pi} \mathbb{E}^{\pi} \left[ \sum_{k=0}^T \gamma^k \|X_k\|^2 - \alpha \mathcal{H}(\pi(du \mid X_k)) \right], \quad \alpha > 0, \quad \mathcal{H} : \text{entropy}$$

## Observations

- Better exploration
- Robustness
- Flat minima and optimisation regularity (recent work: Ahmed et al. ICLR (2019)<sup>1</sup>)

---

<sup>1</sup>A. Ahmed et al. "Understanding Flat Minima in Neural Networks", *ICLR* (2019)

# Robustness: Maximum Entropy and Flat Minima

## Maximum Entropy in Reinforcement Learning

$$\arg \min_{\pi} \mathbb{E}^{\pi} \left[ \sum_{k=0}^T \gamma^k \|X_k\|^2 - \alpha \mathcal{H}(\pi(du \mid X_k)) \right], \quad \alpha > 0, \quad \mathcal{H} : \text{entropy}$$

## Observations

- Better exploration
- Robustness
- Flat minima and optimisation regularity (recent work: Ahmed et al. ICLR (2019)<sup>1</sup>)

## Questions:

Why does entropy improve robustness? Why does entropy regularise the optimisation landscape?

## Objective

Understanding robustness-entropy-regularity synergy

## Hypothesis

Entropy  $\longrightarrow$  Policy Complexity

---

<sup>1</sup>A. Ahmed et al. "Understanding Flat Minima in Neural Networks", *ICLR* (2019)

## Partial Observability

$$\begin{aligned} X_{k+1} &= P(X_k, U_k, \eta_k) \\ Y_{k+1} &= Q(X_k) + \epsilon_k \quad \epsilon_k \sim \mathcal{N}(0, \sigma_\epsilon^2 I_d) \end{aligned} \tag{2}$$

## Notation

When  $\epsilon \equiv 0 \longrightarrow P^\pi$

When  $\epsilon \not\equiv 0 \longrightarrow P^{\pi, \epsilon}$

# Excess Risk Under Noise

## Partial Observability

$$\begin{aligned} X_{k+1} &= P(X_k, U_k, \eta_k) \\ Y_{k+1} &= Q(X_k) + \epsilon_k \quad \epsilon_k \sim \mathcal{N}(0, \sigma_\epsilon^2 I_d) \end{aligned} \tag{2}$$

## Notation

When  $\epsilon \equiv 0 \longrightarrow P^\pi$

When  $\epsilon \not\equiv 0 \longrightarrow P^{\pi, \epsilon}$

## Rate of Excess Risk Under Noise

$$\dot{\mathcal{R}}^\pi = \frac{J^{\pi, \epsilon} - J^\pi}{J^\pi} \tag{3}$$

with  $J^{\pi, \epsilon} = \mathbb{E}^{\pi, \epsilon} \left[ \sum_{k=0}^T \gamma^k \|X_k\|^2 \right]$

# Training with different temperature levels $\alpha$

## Objective

$$\pi_{\alpha}^* = \arg \min_{\pi} \mathbb{E}^{\pi} \left[ \sum_{k=0}^T \gamma^k \|X_k\|^2 - \alpha \mathcal{H}(\pi(du \mid X_k)) \right], \quad \alpha > 0$$

**Initial condition**  $X_0 \sim \mathcal{N}(x_e^2, \sigma^2)$  and  $\pi_{\theta}(\cdot \mid X_k) \sim \mathcal{N}_{d_{\mathcal{U}}}(\mu_{\theta}(X_k), \theta_{\sigma_{\pi}} I_{d_{\mathcal{U}}})$

**Goal** control  $x_k \rightarrow x_e^0 = 0$

# Training with different temperature levels $\alpha$

## Objective

$$\pi_{\alpha}^* = \arg \min_{\pi} \mathbb{E}^{\pi} \left[ \sum_{k=0}^T \gamma^k \|X_k\|^2 - \alpha \mathcal{H}(\pi(du | X_k)) \right], \quad \alpha > 0$$

**Initial condition**  $X_0 \sim \mathcal{N}(x_e^2, \sigma^2)$  and  $\pi_{\theta}(\cdot | X_k) \sim \mathcal{N}_{d_{\mathcal{U}}}(\mu_{\theta}(X_k), \theta_{\sigma_{\pi}} I_{d_{\mathcal{U}}})$

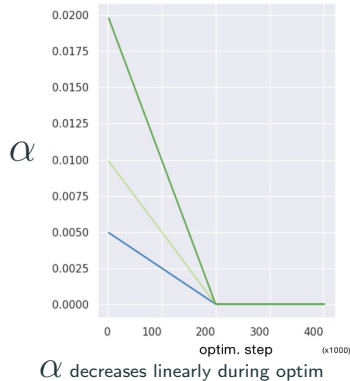
**Goal** control  $x_k \rightarrow x_e^0 = 0$

## Hypothesis

With  $\alpha > 0$  the policies  $\pi_{\alpha}^*$  are more robust than  $\pi_{\alpha=0}^*$

## Experimental Plan

- Fix 5 entropy levels  $\alpha$
- 10 trainings for each  $\alpha$  for 2m of iterations with the system
- $\alpha$  decreases linearly
- Study of the regularity of  $\pi_{\alpha}^*$  and its robustness





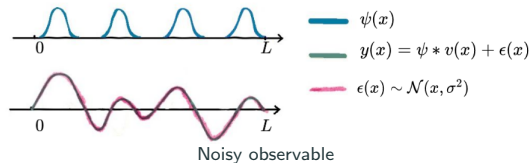
# Evaluation of the policy with noisy observation

## Hypothesis

$\epsilon \nearrow \longrightarrow J^{\pi^*, \epsilon} \nearrow$  (noise impacts perf)  
 $\alpha > 0 \longrightarrow \mathring{\mathcal{R}}^{\pi, \alpha} \searrow$  (robustness)

## Experimental Plan

- *Test*  $\pi_\alpha^*$  with *different noise levels*  $\epsilon$  on  $Y$
- Compare  $J^{\pi^*, \epsilon}$  according to  $J^{\pi^*}$  i.e.  $\mathring{\mathcal{R}}^\pi = \frac{J^{\pi^*, \epsilon} - J^{\pi^*}}{J^{\pi^*}}$



with  $J^{\pi^*} = \mathbb{E}^{\pi^*} \left[ \sum_{k=0}^T \gamma^k \|X_k\|^2 \right]$   
and  $J^{\pi^*, \epsilon}$  same quantity evaluated  
with **noisy observables**

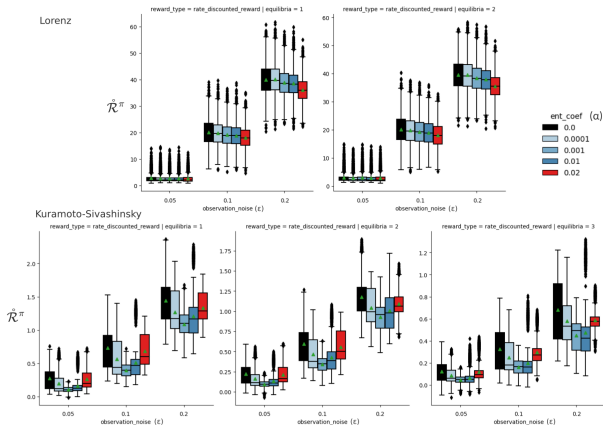
# Observation noise robustness by Maximum Entropy

## Experiment

- Evaluate 10 models  $\theta_{\alpha}^*$  for each value of  $\alpha$
- **Total** : 50 models  $\theta_{\alpha_i}^*$
- $\forall \theta_{\alpha_i}^*$  evaluate 200 trajectories until  $T$

## Results

- Noise  $\epsilon$  increases globally the cost  $J^{\pi^*}$
- **KS and Lorenz**:  $\alpha = 0$  noise sensitive
- **KS**:  $\alpha_{\max}$  noise sensitive



Variation  $\frac{J^{\pi^*}_{\alpha, \epsilon} - J^{\pi^*}}{J^{\pi^*}}$ . Each bar block : noise intensity  $\epsilon$ .  
 Colors:  $\alpha = 0$  (black),  $\alpha > 0$  (blue),  $\alpha_{\max}$  (red)

# Complexity measures<sup>1</sup>

## Complexity Measure

$$\mathcal{M}: \pi \in \Pi \rightarrow \mathbb{R}_+$$

$\mathcal{M}(\pi)$  measures the **complexity** of the model  $\pi$

## Robustness Measure

$$\mathring{\mathcal{R}}^\pi \leq f(\mathcal{M}(\pi))$$

where  $f$  is an increasing function

## Objective

Identify proper complexity measures for robustness

---

<sup>1</sup>B. Neyshabur et al. "Exploring Generalization in Deep Learning" *NIPS* (2017)

# Complexity Measure: Lipschitz Upper Bound

## Lipshitz Bound

$$\pi_{\theta}(\cdot|X_k) \sim \mathcal{N}_{d_{\mathcal{U}}}(\mu_{\theta}(X_k), \theta_{\sigma_{\pi}} I_{d_{\mathcal{U}}})$$

$$\text{If } \mu_{\theta}(x) = (\sigma_l \circ \sigma_{l-1} \circ \dots \circ \sigma_1)(x),$$

$$\text{Lips}(\mu_{\theta}) \leq \prod_{i=1}^l \text{Lips}(\sigma_i) = \prod_{i=1}^l \|\theta_i\|,$$

where  $\theta_i$  weight matrix  $i$ .

# Complexity Measure: Lipschitz Upper Bound

## Lipshitz Bound

$$\pi_{\theta}(\cdot | X_k) \sim \mathcal{N}_{d_{\mathcal{U}}}(\mu_{\theta}(X_k), \theta_{\sigma_{\pi}} I_{d_{\mathcal{U}}})$$

$$\text{If } \mu_{\theta}(x) = (\sigma_I \circ \sigma_{I-1} \circ \dots \circ \sigma_1)(x),$$

$$\text{Lips}(\mu_{\theta}) \leq \prod_{i=1}^I \text{Lips}(\sigma_i) = \prod_{i=1}^I \|\theta_i\|,$$

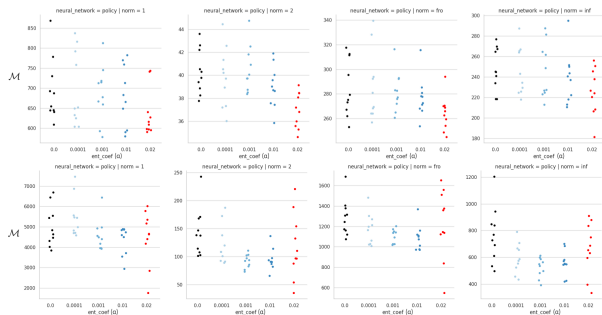
where  $\theta_i$  weight matrix  $i$ .

## Lipshitz Complexity Measure

- $\mathcal{M}(\pi_{\theta}) = \prod_{i=1}^I \|\theta_i\|$

## Result

Low  $\mathcal{M}(\pi_{\theta}^{\alpha})$  corresponds to low  $\hat{\mathcal{R}}^{\pi}$



$$\mathcal{M}(\pi_{\theta}^{\alpha}) = \prod_{i=1}^I \|\theta_i^{\alpha}\|$$

Colors:  $\alpha = 0$ ,  $\alpha > 0$ ,  $\alpha_{\max}$

Top: Lorenz, Bottom: KS

# Complexity Measure: Conditional Fisher Information Trace

## Hessian and Fisher Information

$\alpha > 0 \longrightarrow$  Flat Minima (already observed)<sup>1</sup>

$$\nabla_{\theta}^2 J^{\pi_{\theta}} = \mathbb{E}^{\pi_{\theta}} \left[ \sum_{h,i,j=0}^T c(X_h, U_h) \left( \nabla_{\theta} \log \pi_{\theta}(U_i | X_j) \nabla_{\theta} \log \pi_{\theta}(U_j | X_j)^T + \nabla_{\theta}^2 [\log \pi_{\theta}(U_i | X_j)] \right) \right]$$

$$\text{Fisher Information: } \mathcal{I}(\theta) = -\mathbb{E}^{X \sim \rho, U \sim \pi_{\theta}(\cdot|X)} [\nabla_{\theta}^2 \log \pi_{\theta}(U|X)]$$

# Complexity Measure: Conditional Fisher Information Trace

## Hessian and Fisher Information

$\alpha > 0 \rightarrow$  Flat Minima (already observed)<sup>1</sup>

$$\nabla_{\theta}^2 J^{\pi_{\theta}} = \mathbb{E}^{\pi_{\theta}} \left[ \sum_{h,i,j=0}^T c(X_h, U_h) \left( \nabla_{\theta} \log \pi_{\theta}(U_i | X_j) \nabla_{\theta} \log \pi_{\theta}(U_j | X_j)^T + \nabla_{\theta}^2 [\log \pi_{\theta}(U_i | X_j)] \right) \right]$$

$$\text{Fisher Information: } \mathcal{I}(\theta) = -\mathbb{E}^{X \sim \rho, U \sim \pi_{\theta}(\cdot|X)} \left[ \nabla_{\theta}^2 \log \pi_{\theta}(U|X) \right]$$

## Fisher Information Complexity Measure

- $\mathcal{M}(\pi_{\theta}) = -\text{Tr} \left( \mathbb{E}^{X \sim \rho^{\pi_{\theta}}, U \sim \pi_{\theta}(\cdot|X)} \left[ \nabla_{\theta}^2 \log \pi_{\theta}(U|X) \right] \right)$

# Complexity Measure: Conditional Fisher Information Trace

## Hessian and Fisher Information

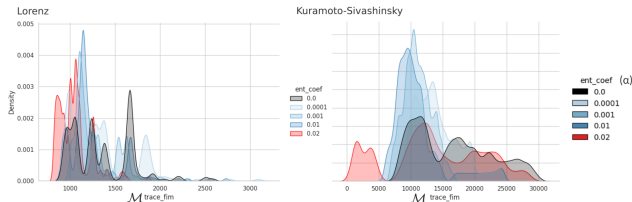
$\alpha > 0 \rightarrow$  Flat Minima (already observed)<sup>1</sup>

$$\nabla_{\theta}^2 J^{\pi_{\theta}} = \mathbb{E}^{\pi_{\theta}} \left[ \sum_{h,i,j=0}^T c(X_h, U_h) \left( \nabla_{\theta} \log \pi_{\theta}(U_i | X_j) \nabla_{\theta} \log \pi_{\theta}(U_j | X_j)^T + \nabla_{\theta}^2 [\log \pi_{\theta}(U_i | X_j)] \right) \right]$$

Fisher Information:  $\mathcal{I}(\theta) = -\mathbb{E}^{X \sim \rho, U \sim \pi_{\theta}(\cdot|X)} [\nabla_{\theta}^2 \log \pi_{\theta}(U|X)]$

## Fisher Information Complexity Measure

- $\mathcal{M}(\pi_{\theta}) = -Tr \left( \mathbb{E}^{X \sim \rho^{\pi_{\theta}}, U \sim \pi_{\theta}(\cdot|X)} [\nabla_{\theta}^2 \log \pi_{\theta}(U|X)] \right)$



## Result

Low  $\mathcal{M}(\pi_{\theta}^{\alpha})$  corresponds to low  $\mathring{\mathcal{R}}^{\pi}$

$$\mathcal{M}(\pi_{\theta}^{\alpha}) = -Tr(\mathbb{E}^{X \sim \rho^{\pi_{\theta}}, U \sim \pi_{\theta}(\cdot|X)} [\nabla_{\theta}^2 \log \pi_{\theta}(U|X)])$$

Colors:  $\alpha = 0$ ,  $\alpha > 0$ ,  $\alpha_{\max}$



## Hypothesis

**Entropy**  $\longrightarrow$  **Flat Minimum**    Already observed in (Ahmed et al. 2019)

**Flat Minimum**  $\longleftrightarrow$  **Robustness**  $\longleftrightarrow$  **Fisher Information of  $\theta_\pi$**     ✓

## Robustness of the results

- For  $\alpha_{\max}$  we lose robustness because we no longer solve the same objective
- Lorenz (fully observable) does not discriminate policies (because deterministic solution?)

## Perspectives

Optimization scheme based on flat minima or Fisher Information

# **Learning Based Control: Sampling Strategies**

---

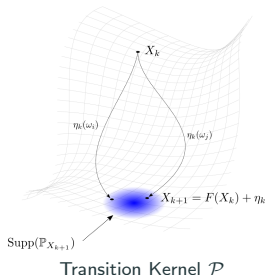
# Gaussian Process Modeling

## Controlled Hidden Markov Chain

$$P^\pi(dx_0 du_0 dx_1 du_1 \dots dx_T) = P_{X_0}(dx_0) \mathcal{Q}(dy_0 | x_0) \pi(x_0, du_0) \mathcal{P}(dx_1 | x_0, u_0) \\ \mathcal{Q}(dy_1 | x_1) \pi(x_1, du_1) \dots \pi(x_{T-1}, du_{T-1}) \mathcal{P}(dx_T | x_{T-1}, u_{T-1})$$

## Learning Dynamics with Gaussian Process

$$\hat{\mathcal{P}}_{\mathcal{D}}(\cdot, (x, u)) \sim \mathcal{N}(\mu_{(x,u)}, k_{(x,u), (x,u)} | \mathcal{D}) \quad (4)$$



<sup>1</sup>C. E. Rasmussen et al. "Gaussian Processes in Reinforcement Learning" *NIPS* (2003)

# Dynamics Approximation with Model Predictive Control

## Model Predictive Control

$$\pi^{\text{MPC}}(x) = u_0^* \quad (5)$$

$$s.t. \quad (u_0^*, \dots, u_{T^{\text{MPC}}}^*) = \arg \min_{(u_0, \dots, u_{T^{\text{MPC}}})} \mathbb{E}^{(u_0, \dots, u_{T^{\text{MPC}}})} \left[ \sum_{k=0}^{T^{\text{MPC}}} c(X_k, u_k) \mid X_0 = x \right] \quad (6)$$

## New Problem

- Sampling budget  $\rightarrow n$
- Collect  $\mathcal{D}_n$  such that  $\hat{\mathcal{P}}_{\mathcal{D}} \simeq \mathcal{P}$
- Data  $\mathcal{D}_n = \{(x_0, u_0), \dots, (x_n, u_n)\}$  is collected online along an observed dynamics
- How to sample next data point?

---

<sup>1</sup>R. Y. Rubinstein et al. The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization, Monte-Carlo Simulation, and Machine Learning *Springer* (2004)

# Entropy Map

How to quantify the uncertainty on  $P_{X_{k+1}}$ ?

Infinitesimal volume element of  $\mathcal{X} \rightarrow dx$

**Uncertainty on  $dx$**

$$I(dx) = \log\left(\frac{1}{P_{X_{k+1}}(dx)}\right)$$

# Entropy Map

How to quantify the uncertainty on  $P_{X_{k+1}}$ ?

Infinitesimal volume element of  $\mathcal{X} \rightarrow dx$

**Uncertainty on  $dx$**

$$I(dx) = \log\left(\frac{1}{P_{X_{k+1}}(dx)}\right)$$

**Entropy (average uncertainty)**

$$\mathcal{H}(P_{X_{k+1}}) = \int_{\mathbb{R}} \log \frac{1}{P_{X_{k+1}}(dx)} P_{X_{k+1}}(dx)$$

# Entropy Map

How to quantify the uncertainty on  $P_{X_{k+1}}$ ?

Infinitesimal volume element of  $\mathcal{X} \rightarrow dx$

## Uncertainty on $dx$

$$l(dx) = \log\left(\frac{1}{P_{X_{k+1}}(dx)}\right)$$

## Entropy (average uncertainty)

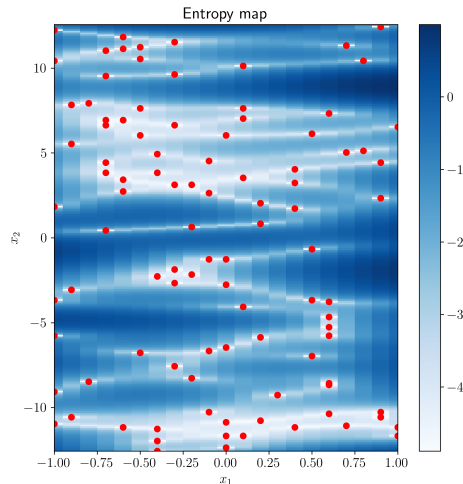
$$\mathcal{H}(P_{X_{k+1}}) = \int_{\mathbb{R}} \log \frac{1}{P_{X_{k+1}}(dx)} P_{X_{k+1}}(dx)$$

## Gaussian Entropy

$$P_{X_{k+1}}(dx) = f_{\mathcal{N}}(\mu_{(x,u)}, k_{(x,u)}, (x,u))(x)dx$$

$$\mathcal{H}(P_{X_{k+1}}) = \frac{1}{2} \log(2\pi e k_{(x,u)}, (x,u))$$

In the Gaussian case, the variance  $k$  characterise the entropy



# Expected Information Gain

Process *trajectory*  $\longrightarrow H_T = (X_0, U_0, \dots, U_T, X_T)$

Optimal trajectory under  $\hat{\mathcal{P}}_{\mathcal{D}} \longrightarrow H_T^*$

## Expected Information Gain

$$\text{EIG}_n(x, u) = \mathcal{H} \left[ \hat{H}_T^* \mid \mathcal{D}_n \right] - \mathbb{E}_{P_{X_{n+1} \mid \mathcal{D}_n, X_n=x, U_n=u}} \left[ \mathcal{H} \left[ \hat{H}_T^* \mid \mathcal{D}_n, X_n = x, U_n = u, X_{n+1} \right] \right] \quad (7)$$

EIG  $\longrightarrow$  conditional mutual information between  $\hat{H}_T^*$  and  $X_{n+1}$  given  $\mathcal{D}_n$

---

<sup>1</sup>V. Mehta et al. "An Experimental Design Perspective on Model-Based Reinforcement Learning" *ICLR* (2022)



# Expected Information Gain

Process *trajectory*  $\longrightarrow H_T = (X_0, U_0, \dots, U_T, X_T)$

Optimal trajectory under  $\hat{\mathcal{P}}_{\mathcal{D}} \longrightarrow H_T^*$

## Expected Information Gain

$$\text{EIG}_n(x, u) = \mathcal{H} \left[ \hat{H}_T^* \mid \mathcal{D}_n \right] - \mathbb{E}_{P_{X_{n+1} \mid \mathcal{D}_n, X_n=x, U_n=u}} \left[ \mathcal{H} \left[ \hat{H}_T^* \mid \mathcal{D}_n, X_n=x, U_n=u, X_{n+1} \right] \right] \quad (7)$$

EIG  $\longrightarrow$  conditional mutual information between  $\hat{H}_T^*$  and  $X_{n+1}$  given  $\mathcal{D}_n$

## By symmetry of cond. MI

$$\text{EIG}_n(x, u) = \mathcal{H} [X_{n+1} \mid \mathcal{D}_n, X_n=x, U_n=u] - \mathbb{E}_{P_{\hat{H}_T^* \mid \mathcal{D}_n}} \left[ \mathcal{H} [X_{n+1} \mid \mathcal{D}_n, X_n=x, U_n=u, \hat{H}_T^*] \right] \quad (8)$$

---

<sup>1</sup>V. Mehta et al. "An Experimental Design Perspective on Model-Based Reinforcement Learning" *ICLR* (2022)

# Randomized Decision Epochs: Temporal Abstraction with options<sup>1</sup>

## New Problem

- Data  $\mathcal{D}_n = \{(x_0, u_0), \dots, (x_n, u_n)\}$  is collected online along an observed dynamics
- When to sample next data point?

---

<sup>1</sup>R. S. Sutton et al. "Between MDPs and semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning", *NIPS* (1999)

# Randomized Decision Epochs: Temporal Abstraction with options<sup>1</sup>

## New Problem

- Data  $\mathcal{D}_n = \{(x_0, u_0), \dots, (x_n, u_n)\}$  is collected online along an observed dynamics
- When to sample next data point?

## Hypothesis

Wait for auto-decorrelation of  $(X_{n+1}, U_{n+1})$  from  $\mathcal{D}_n$

---

<sup>1</sup>R. S. Sutton et al. "Between MDPs and semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning", *NIPS* (1999)

# Randomized Decision Epochs: Temporal Abstraction with options<sup>1</sup>

## New Problem

- Data  $\mathcal{D}_n = \{(x_0, u_0), \dots, (x_n, u_n)\}$  is collected online along an observed dynamics
- When to sample next data point?

## Hypothesis

Wait for auto-decorrelation of  $(X_{n+1}, U_{n+1})$  from  $\mathcal{D}_n$

Random decision epochs  $\rightarrow (\kappa_j)_{j \in \mathbb{N}}$

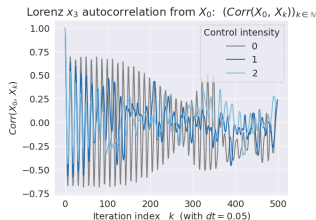
Semi-Markov Decision Process  $\rightarrow (X_{\kappa_j})_{j \in \mathbb{N}}$

$\mathcal{P}^{\text{SMDP}}(dx' | (x, (u, t))) = P(X_{k+t} | X_k = x, U_{k:k+t-1} = u)$

Time-delay  $\rightarrow t \in \mathcal{T}$

Constant control between  $\kappa_j$  and  $\kappa_{j+1}$

New action space  $\rightarrow \mathcal{U} \times \mathcal{T}$



$(\text{Cov}(X_0, X_k))_{k \in \mathbb{N}}$  for the controlled Lorenz system  $x_3$  component under multiple control intensities.

<sup>1</sup>R. S. Sutton et al. "Between MDPs and semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning", *NIPS* (1999)

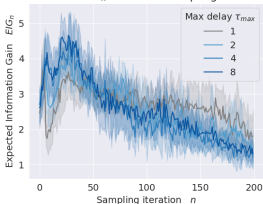
# Semi-Markovian Expected Information Gain

## New Information Gain

$$\text{EIG}_n^{\text{SM-TIP}}(x, (u, t)) =$$

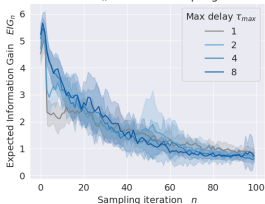
$$\mathcal{H}[X_{\kappa_n+t+1} \mid \mathcal{D}_n, X_{\kappa_n+t} = x, U_{\kappa_n+t} = u, \kappa_n] - \mathbb{E}_{P_{\hat{H}_T^* | \mathcal{D}_n}} \left[ \mathcal{H}[X_{\kappa_n+t+1} \mid \mathcal{D}_n, X_{\kappa_n+t} = x, U_{\kappa_n+t} = u, \hat{H}_T^*, \kappa_n] \right] \quad (9)$$

Evolution of  $\text{EIG}_n$  over the sampling iterations  $n$



(a) Lorenz

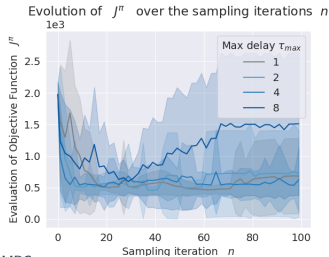
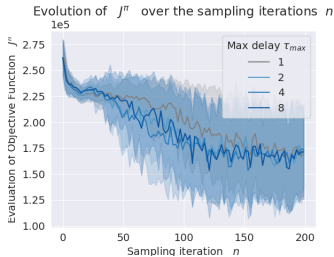
Evolution of  $\text{EIG}_n$  over the sampling iterations  $n$



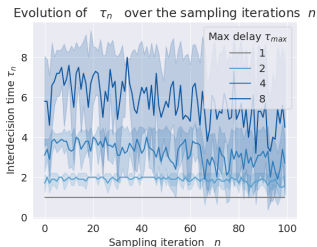
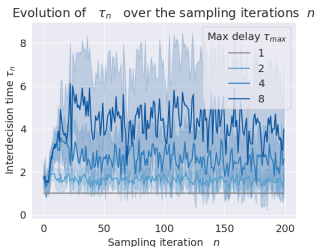
(b) Pendulum

Expected Information Gain  $\text{EIG}_n^{\text{SM-TIP}}$  over  $n$ .

# Objective Function and Inter-Decision Time



Objective function  $J^{\hat{\pi}^{MPC}}$  during training.



(a) Lorenz

(b) Pendulum

## Hypothesis

**Temporal Abstraction**  $\longrightarrow$  Information (fixed sampling budget  $n$ ) ✓

## Robustness of the results

- Evaluation fairness when one model can explore the dynamics further than the other
- How to characterise the  $\delta_t$  threshold where temporal abstraction is beneficial?

## Perspectives

Correct bootstrapping error in the SMDP