# Google DeepMind
# Curiosity in Hindsight:
# Intrinsic Exploration in Stochastic Environments

Daniel Jarrett, Corentin Tallec, Florent Altché, Thomas Mesnard, Rémi Munos, Michal Valko

## Introduction

How to explore the world when external rewards are sparse or absent?

- **Curiosity-Driven Exploration**: Prioritize exploring (and learning) what is not yet understood, where "understanding" := ability to **predict** outcome.

- **Problem of Stochasticity**: Predictive error-based exploration agents are often "stuck" around **high-entropy** elements in the state-action space.

- **Example (Noisy TV Problem)**: Since the agent never knows what the next state is, they will stand in front of it, hoping to learn what is never learnable.

**Key Idea**: Disentangle (irreducible) "**noise**" from (reducible) "**novelty**", and only use novelty to guide exploration. How?

- Learn **representations of the future** capturing precisely the **unpredictable** aspects of each outcome (no more, no less).

- Use this as additional input when making predictions, such that intrinsic rewards only reflect the **predictable** aspects of world dynamics.

## Two Models + Two Objectives

Learn a **Reconstructor** $f_\eta$ for $f(X_t, A_t, Z_{t+1})$ and a **Generator** $p_\theta$ for $p_\eta(Z_{t+1} | X_t, A_t, X_{t+1})$.



$Z_{t+1}$ should capture *at least* all aspects that are unpredictable (so we <u>don't</u> reward the agent for *irreducible* error).

**Objective 1 (Reconstruction)**

$$\mathcal{L}_\eta^{\text{rec.}}(x_t, a_t, z_{t+1}, x_{t+1}) := \left\| x_{t+1} - f_\eta(x_t, a_t, z_{t+1}) \right\|_2^2$$

$$\mathcal{R}_{\theta,\eta}^{\text{rec.}}(x_t, a_t) := \mathbb{E}_{\substack{X_{t+1} \sim \tau(\cdot | x_t, a_t) \\ Z_{t+1} \sim p_\theta(\cdot | x_t, a_t, X_{t+1})}} \mathcal{L}_\eta^{\text{rec.}}(x_t, a_t, z_{t+1}, x_{t+1})$$

**Objective 2 (Invariance)**

$$\mathcal{L}_\theta^{\text{inv.}}(x_t, a_t, z_{t+1}) := \log \frac{p_\theta(z_{t+1} | x_t, a_t)}{p_\theta(z_{t+1})}$$

$$\mathcal{R}_\theta^{\text{inv.}}(x_t, a_t) := \mathbb{E}_{\substack{X_{t+1} \sim \tau(\cdot | x_t, a_t) \\ Z_{t+1} \sim p_\theta(\cdot | x_t, a_t, X_{t+1})}} \mathcal{L}_\theta^{\text{inv.}}(x_t, a_t, z_{t+1})$$

$Z_{t+1}$ should capture *at most* all aspects that are unpredictable (so we *do* reward the agent for *reducible* error).

## Implementation

Ask a **Critic** to $g_\nu$ to maximize a contrastive loss.

**Objective 3 (Contrastive Learning)**

$$\mathcal{L}_{\theta,\nu}^{K,\text{con.}}(x_t, a_t, z_{t+1}) := \mathbb{E} \quad \log \frac{e^{g_\nu(x_t, a_t, z_{t+1})}}{\frac{1}{K}\left( e^{g_\nu(x_t, a_t, z_{t+1})} + \sum_{i=1}^{K-1} e^{g_\nu(x_t, a_t, Z_{t+1}^i)} \right)}$$

$$(X_t^1, \ldots, X_t^{K-1}) \sim \prod_{i=1}^{K-1} \rho_\pi$$
$$(A_t^1, \ldots, A_t^{K-1}) \sim \prod_{i=1}^{K-1} \pi(\cdot | X_t^i)$$
$$(X_{t+1}^1, \ldots, X_{t+1}^{K-1}) \sim \prod_{i=1}^{K-1} \tau(\cdot | X_t^i, A_t^i)$$
$$(Z_{t+1}^1, \ldots, Z_{t+1}^{K-1}) \sim \prod_{i=1}^{K-1} p_\theta(\cdot | X_t^i, A_t^i, X_{t+1}^i)$$

$$\mathcal{R}_{\theta,\nu}^{K,\text{con.}}(x_t, a_t) := \mathbb{E}_{\substack{X_{t+1} \sim \tau(\cdot | x_t, a_t) \\ Z_{t+1} \sim p_\theta(\cdot | x_t, a_t, X_{t+1})}} \mathcal{L}_{\theta,\nu}^{K,\text{con.}}(x_t, a_t, z_{t+1})$$

**In Practice**: (1) batch size $K < \infty$, (2) $\nu$ is not fully optimized, (3) $\lambda$ is a hyperparameter.

The intrinsic reward is now:

$$\mathcal{R}_{\theta,\eta,\nu}^K(x_t, a_t) := \frac{1}{\lambda}\mathcal{R}_{\theta,\eta}^{\text{rec.}}(x_t, a_t) + \mathcal{R}_{\theta,\nu}^{K,\text{con.}}(x_t, a_t)$$

and the agent performs:

$$\begin{array}{cc} \text{(policy)} & \text{(model)} \\ \underset{\pi}{\text{maximize}} & \underset{\theta,\nu}{\min}\,\underset{\nu}{\max} \end{array} \mathbb{E}_{\substack{X_t \sim \rho_\pi \\ A_t \sim \pi(\cdot | X_t)}} \mathcal{R}_{\theta,\eta,\nu}^K(X_t, A_t)$$

Overall, simple *drop-in modification* on top of any choice of curiosity-driven exploration.

## Curiosity

**Definition 1 (Curiosity)**: Define the *intrinsic reward*
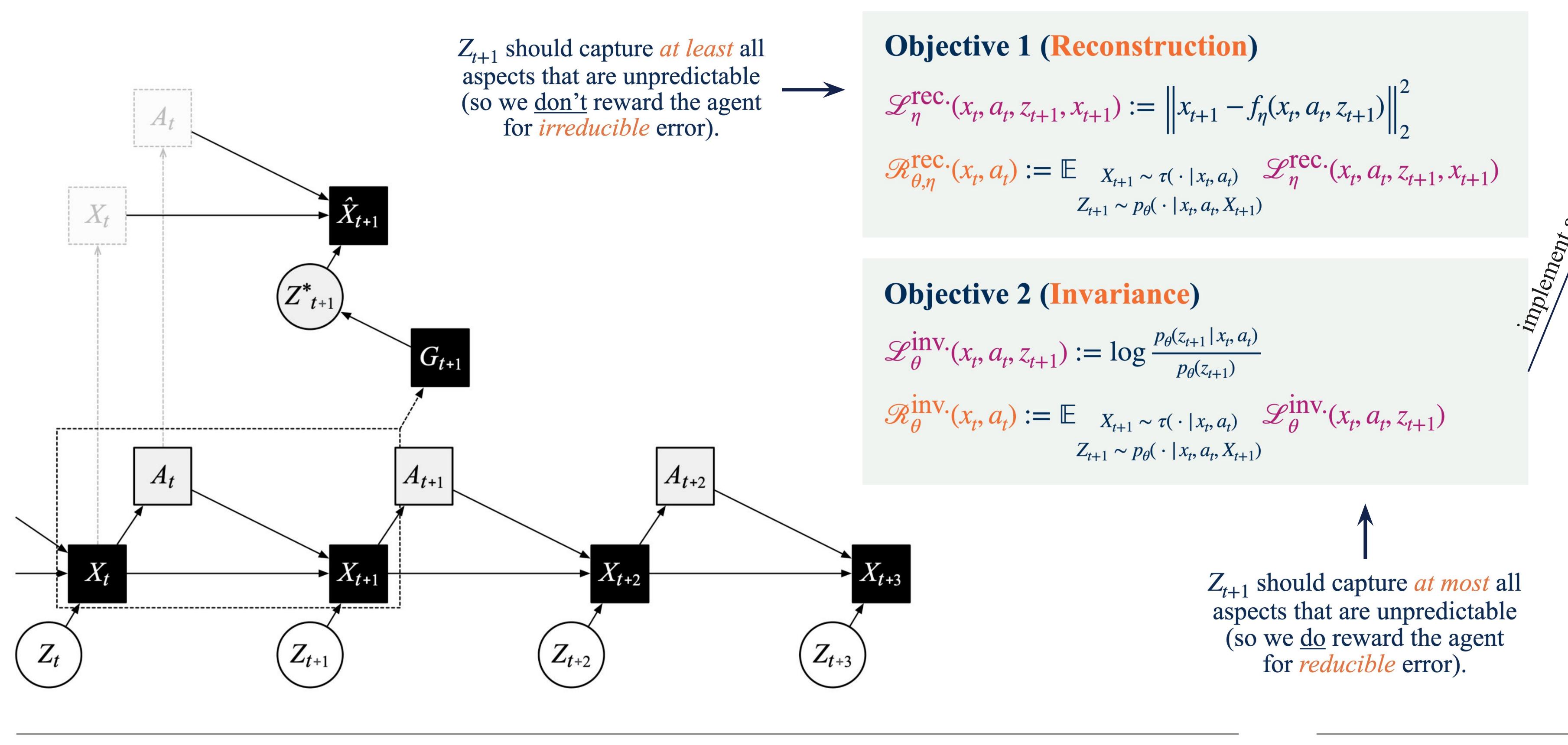
$$\mathcal{R}_\eta(x_t, a_t) := -\mathbb{E}_{X_{t+1} \sim \tau(\cdot | x_t, a_t)} \log \tau_\eta(X_{t+1} | x_t, a_t)$$

The agent performs

$$\begin{array}{cc} \text{(policy)} & \text{(model)} \\ \underset{\pi}{\text{maximize}} & \underset{\underline{\eta}}{\min} \end{array} \mathbb{E}_{\substack{X_t \sim \rho_\pi \\ A_t \sim \pi(\cdot | X_t)}} \mathcal{R}_\eta(X_t, A_t)$$

## Curiosity in Hindsight

**Definition 2 (Curiosity in Hindsight)**: Define the *hindsight intrinsic reward*

$$\mathcal{R}_{\theta,\eta}(x_t, a_t) := \frac{1}{\lambda}\mathcal{R}_{\theta,\eta}^{\text{rec.}}(x_t, a_t) + \mathcal{R}_\theta^{\text{inv.}}(x_t, a_t)$$

The agent performs

$$\begin{array}{cc} \text{(policy)} & \text{(model)} \\ \underset{\pi}{\text{maximize}} & \underset{\underline{\theta,\eta}}{\min} \end{array} \mathbb{E}_{\substack{X_t \sim \rho_\pi \\ A_t \sim \pi(\cdot | X_t)}} \mathcal{R}_{\theta,\eta}(x_t, a_t)$$

## Illustrative Example



Montezuma's Revenge



... with Sticky Actions