*11th International Conference on Machine Learning, Optimization and Data Science*

# Deep Active Inference Agents
## for Delayed and Long-Horizon Environments
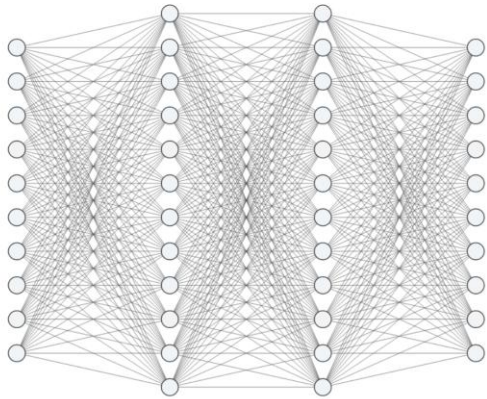
*Yavar Taheri Yeganeh [1], Mohsen Jafari [2], Andrea Matta [1]*

*yavar.taheri@polimi.it, jafari@soe.rutgers.edu, andrea.matta@polimi.it*

*[1] Politecnico di Milano, [2] Rutgers University*

22.09.2025 | Yavar Taheri Yeganeh
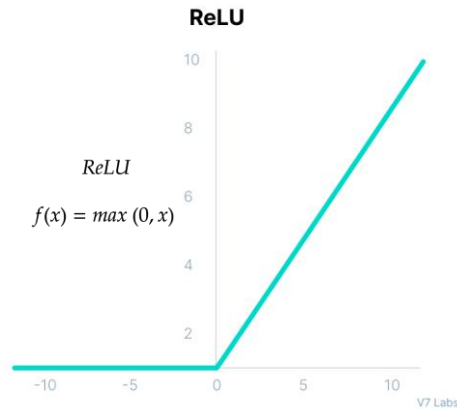
# Neuroscience in Machine Learning and AI

## Arterial Neural Networks



Modern Machine Learning
Universal Function
Approximators

Inspired by
Biological Neural Networks

## ReLU Activation Function



Deep Learning

Inspired by neural firings in the brain "*Neurons either fire or do not fire (binary activation)*"
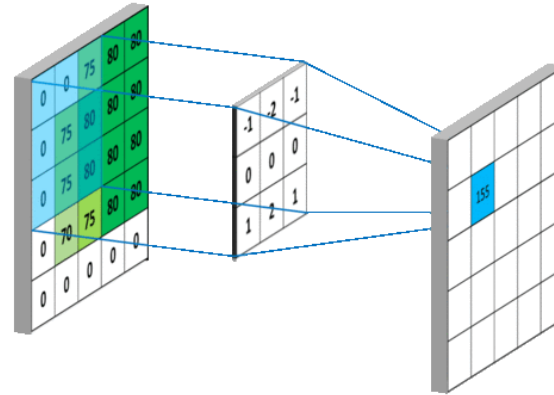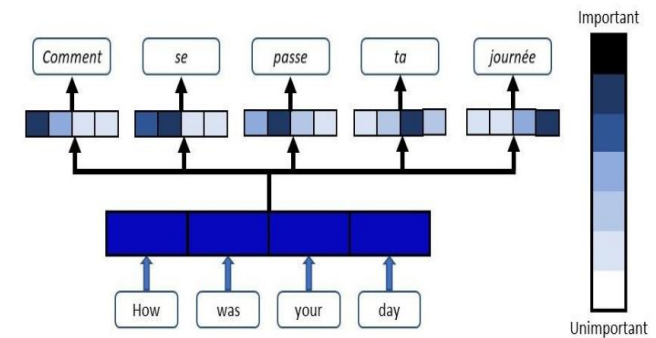
## Convolutional Neural Networks



Image Processing
Computer Vision

Inspired by the visual processing in the brain (the hierarchical structure of the visual cortex )
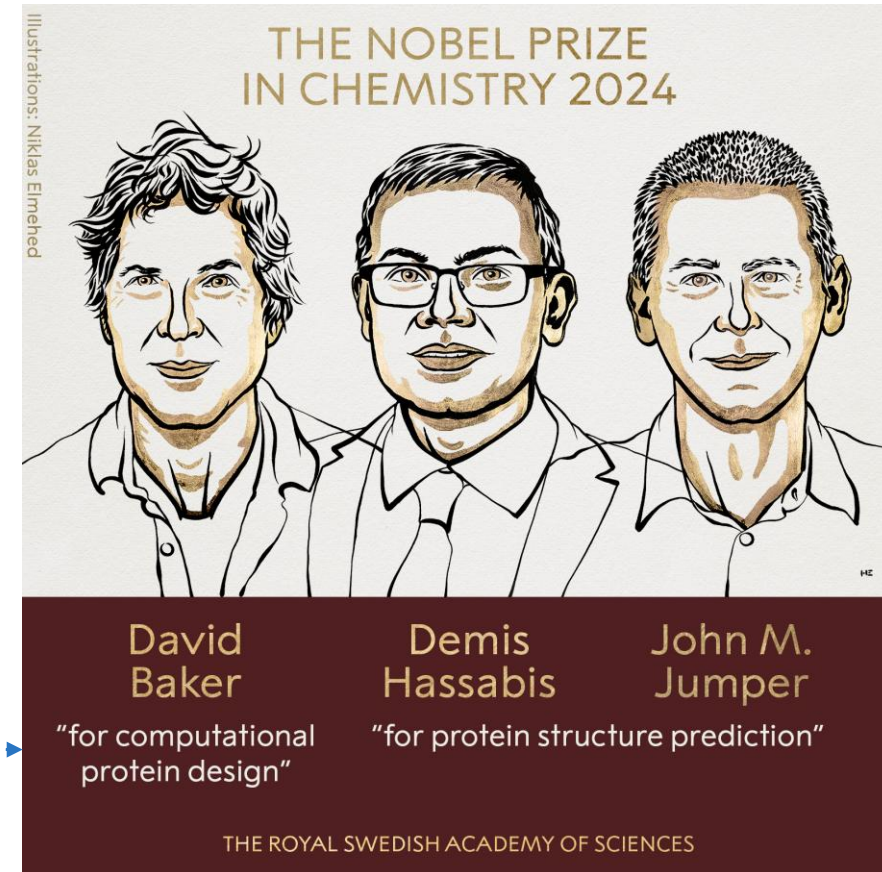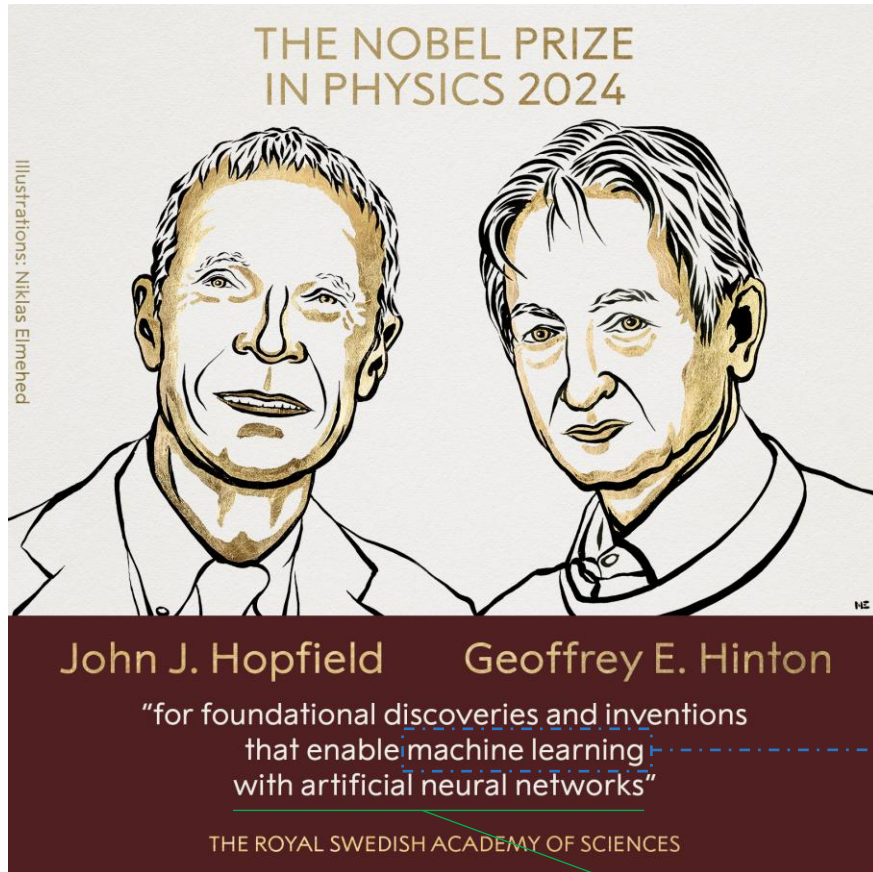
## Attention Mechanism



Transformers
LLMs (e..g, ChatGPT)

Connections to modulations of neural activities in the brain (selectively focus on specific stimuli)

# Neuroscience in Machine Learning and AI



Notably impacted by concepts from Neuroscience (and with Mathematics and Physics)

# Neuroscience: Active Inference Trend



*Karl Friston*

# Active Inference: Concept

Theory of Biological Brains

Based on Predictive Coding: Both are grounded on Generative Model and Predictions

Perception and Learning are driven by Prediction Errors.

Observations are influenced by actions: **Observation (action)**

Predictive Coding → Active Inference



Free Energy Principle: Minimizing surprise $(-\ln p(observation|model)$ or *Free Energy*$)$

Active Inference unifies perception, learning, and action:

- Variational Free Energy: Perception & Learning Model

- Expected Free Energy: Action

# Active Inference: Robotics



Robot Control [2]



Navigation of Autonomous Vehicles [3]

[2] Lanillos et al., 2021
[3] Huang et al., 2024

# Active Inference: Against Model-Free (DAIMC [4])



**Algorithm 1** DAIMC on-policy training

1: **for** $t = 1, 2, \ldots,$ max iterations **do**
2:     Randomize environment and sample a new observation $\tilde{o}_t$.
3:     Run planner and compute prior policy $P(a_t)$.
4:     Compute $Q_{\phi_s}(s_t)$ using $\tilde{o}_t$.
5:     Compute $Q_{\phi_a}(a_t)$ using a sampled state $\tilde{s}_t \sim Q_{\phi_s}(s_t)$.
6:     Compute $D_t = D_{\mathrm{KL}}\left[Q_{\phi_a}(a_t) \parallel P(a_t)\right]$.
7:     Apply a gradient step on $\phi_a$ using $D_t$ as loss.
8:     Compute $\omega_{t+1}$ from Eq. (10) using $D_t$.
9:     Apply action $\tilde{a}_t \sim P(a_t)$ to the environment and sample a new observation $\tilde{o}_{t+1}$.
10:     Compute $\mu, \sigma$ from $P_{\theta_s}(s_{t+1}|\tilde{s}_t, \tilde{a}_t)$.
11:     Compute $Q_{\phi_s}(s_{t+1})$ using $\tilde{o}_{t+1}$.
12:     Apply a gradient step on $\theta_s$ using $D_{\mathrm{KL}}\left[Q_{\phi_s}(s_{t+1}) \parallel \mathcal{N}(\mu, \sigma^2/\omega_t)\right]$.
13:     Apply a gradient step on $\phi_s$, $\theta_o$ using $-\mathbb{E}_{Q(s_{t+1})}\left[\log P_{\theta_o}(o_{t+1}|s_{t+1})\right] + D_{\mathrm{KL}}\left[Q_{\phi_s}(s_{t+1}) \parallel \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2/\omega_t)\right]$.
14: **end for**

[4] Fountas et al., 2020

# Active Inference: Core

The core of the formalism is the Generative Model.

Realization of this formalism is through Deep Learning: Generative Modeling and Inferences



$\tilde{o}_t$ .............................. $Q_{\phi_s}(s_t)$ .................. $P_{\theta_s}(s_{t+1}|\tilde{s}_t, \tilde{a}_t)$ ......... $P_{\theta_o}(o_{t+1}|\tilde{s}_{t+1})$

State (i.e., representation) space of the model

**Encoder**        **Transition**        **Decoder**

$\tilde{o}_t$        $\longrightarrow$        $\tilde{o}_{t+1}$

$\tilde{a}_t \, , \, \Delta_t$

Variational Auto Encoder

**POLITECNICO MILANO 1863** | DIPARTIMENTO DI MECCANICA

# World-Modelling and Model-Based RL

## Model-Based RL

- Use a model for making decisions
- Engaging methods for policy improvement and planning

## World Models

- Build an internal predictive/generative model of the world (often task agnostic)
- Focuses on modelling and representations, either latent or ambient
- Learning policy (or planning) in model imaginations (roll outs)

**General agents contain world models**

Jonathan Richens[1]  David Abel[1]  Alexis Bellot[1]  Tom Everitt[1]

**Abstract**

Are world models a necessary ingredient for flexible, goal-directed behaviour, or is model-free learning sufficient? We provide a formal answer to this question, showing that any agent capable of generalizing to multi-step goal-directed tasks must have learned a predictive model of its environment. We show that this model can be extracted from the agent's policy, and that increasing the agents performance or the complexity of the goals it can achieve requires learning increasingly accurate world models. This has a number of consequences: from developing safe and general agents, to bounding agent capabilities in complex environments, and providing new algorithms for learning an agent's world model.

Figure 1: Our result complements previous insights from planning and inverse RL. While planning uses a world model and a goal to determine a policy, and IRL and inverse planning use an agent's policy and a world model to identify its goal, our result uses an agent's policy and its goal to identify a world model



(a) Learn dynamics from experience   (b) Learn behavior in imagination   (c) Act in the environment

*Sequential rollouts*

$$
\text{RSSM} \begin{cases}
\text{Sequence model:} & h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1}) \\
\text{Encoder:} & z_t \sim q_\phi(z_t \mid h_t, x_t) \\
\text{Dynamics predictor:} & \hat{z}_t \sim p_\phi(\hat{z}_t \mid h_t) \\
\text{Reward predictor:} & \hat{r}_t \sim p_\phi(\hat{r}_t \mid h_t, z_t) \\
\text{Continue predictor:} & \hat{c}_t \sim p_\phi(\hat{c}_t \mid h_t, z_t) \\
\text{Decoder:} & \hat{x}_t \sim p_\phi(\hat{x}_t \mid h_t, z_t)
\end{cases}
$$

*Dreamer* learns policy with Actor-Critic in model imaginations [5]

[5] Hafner et al., 2023.

**POLITECNICO MILANO 1863** | DIPARTIMENTO DI MECCANICA

# Motivation

|  | Concept/Formalism | Architeture/Algorithm | KPI Improvement |

Machine Learning

**Theoretical Development**　　　**Model Development**　　　**Application Development**

Promises of Deep Learning: Generative Models:

- *Bayesian and Probabilistic*
- *Representation Learning for both latent and ambient spaces*
- *Utilizes the observatory (cheap and raw) data*

Promises of Active Inference (AIF):

- *Bayesian Framework: Integrates perception, learning, and action*
- *Model-Driven*
- *Adaptability*
- *Intrinsic Objectives: Exploration to resolve uncertainty*
- *End-to-End Action: Doesn't require high-quality reward signals*

Promises of Word-Model RL Agents:

- *Planning via imagination*
- *Generalization & Sample-efficient*
- *Reduced reliance on rewards*

Problem:

*AIF agents struggle in delayed and long-horizon environments due to reliance on:*
- *Immediate predictions.*
- *Exhaustive planning.*

Train a Policy within the Generative World-Model via Active Inference
Conencting: DL, AIF, RL

**POLITECNICO MILANO 1863 | DIPARTIMENTO DI MECCANICA**

CONNECTS

# Overshooting Trick: Integrating Policy



Figure 3: The agent's architecture and generative framework resemble that of a VAE. The line on top represents the agent simulating the future and making a prediction, while on the bottom, the agent receives a new observation after $\Delta_t$ of taking an action, $\tilde{a}_t$.



Figure 4: The agent's architecture and generative framework with an an impact/update/planning horizon, h (e.g., 100).



Selecting an impact/update/planning horizon: h (e.g., 300)

# Formalism: Deep Active Inference (DAIF)

Model Learning (Calibration): Variational Free Energy

$$\theta^* = \arg \min_{\theta} \left( \mathbb{E}_{Q_\phi(s_t, a_t)} \left[ \log Q_\phi(s_t, a_t) - \log P_\theta(o_t, s_t, a_t) \right] \right)$$

Decision (Planning & Action): Expected Free Energy

$$G(\pi, \tau) = -\mathbb{E}_{Q(\theta|\pi)Q(s_\tau|\theta,\pi)Q(o_\tau|s_\tau,\theta,\pi)} \left[ \log P(o_\tau|\pi) \right] \quad \text{①}$$
$$+ \mathbb{E}_{Q(\theta|\pi)} \left[ \mathbb{E}_{Q(o_\tau|\theta,\pi)} H(s_\tau|o_\tau, \pi) - H(s_\tau|\pi) \right] \quad \text{②}$$
$$+ \mathbb{E}_{Q(\theta|\pi)Q(s_\tau|\theta,\pi)} H(o_\tau|s_\tau, \theta, \pi) - \mathbb{E}_{Q(s_\tau|\pi)} H(o_\tau|s_\tau, \pi) \quad \text{③}$$

1) Extrinsic Value 2) State Epistemic Uncertainty 3) Parameter Epistemic Uncertainty

Policy Optimization via Gradients

1. Integrate the encoded policy parameters $\hat{\pi}(\phi_a)$ in the model and EFE: $G_\theta(\tilde{o}, \phi_a)$

2. Take the gradient of EFE w.r.t parameters: $\phi_a \leftarrow \phi_a - \alpha \nabla_{\phi_a} \mathbb{E}\left[ G(\phi_a) \right]$



$$\Pi : \mathcal{Q}_{\phi_a} \to \hat{\pi}$$

POLITECNICO MILANO 1863 | DIPARTIMENTO DI MECCANICA

CONNECTS

# DAIF: Architecture

- Actor directly interact with the environment.
- Actor parameters are encoded and integrated into transition.
- Model predicts H steps into future in a single rollout.
- Gradient of policy parameters are taken via EFE of the prediction of the H steps into future.



$\tilde{o}_t \ldots \ldots \ldots \ldots \ldots \ldots Q_{\phi_s}(s_t) \ldots \ldots \ldots \ldots \ldots \ldots \ldots P_{\theta_s}(s_{t+H}|\tilde{s}_t, \hat{\pi}) \ldots \ldots \ldots P_{\theta_o}(o_{t+H}|\tilde{s}_{t+H})$

Encoder $\phi_s$　　Transition $\theta_s$　　Decoder $\theta_o$

$\hat{\pi}(\phi_a)$　　$\nabla_{\phi_a}\mathbb{E}[G(\phi_a)]$

Actor $\phi_a$

$\tilde{o}_t \ldots \ldots \ldots \ldots Q_{\phi_a}(a_t|\tilde{o}_t)$

Environment

# DAIF: Algorithm

- Actor directly interact with the environment.
- Actor parameters are encoded and integrated into transition.
- Model predicts H steps into future in a single rollout.
- Gradient of policy parameters are taken via EFE of the prediction of the H steps into future.

**Algorithm 1** Deep AIF Agent Training (per epoch)

1: Initialize $\theta = \{\theta_s, \theta_o\}$, $\phi = \{\phi_s, \phi_a\}$, $\mathcal{M}$
2: Randomly initialize $E$
3: **for** $n = 1, 2, ..., N$ **do**
   ▷ ENVIRONMENT INTERACTION
4:　$\hat{\pi}_t \leftarrow \Pi(\mathcal{Q}_{\phi_a})$
5:　**for** $\tau = t + 1, t + 2, ..., t + H$ **do**
6:　　Sample a new observation $\tilde{o}_\tau$ from $E$
7:　　Apply $\tilde{a}_\tau \sim Q_{\phi_a}(a_\tau | \tilde{o}_\tau)$ to $E$
8:　　Sample a new observation $\tilde{o}_{\tau+1}$ from $E$
9:　$\mathcal{M} \leftarrow \mathcal{M} \cup \{(\tilde{o}_t, \hat{\pi}_t, \tilde{o}_{t+H})\}$
   ▷ MODEL LEARNING
10:　$\{(\tilde{o}_{t'}, \hat{\pi}_{t'}, \tilde{o}_{t'+H})\}^{B_1} \sim \mathcal{M}$
11:　**for** $t' = 1, 2, ..., B_1$ **do**
12:　　**run** Model$(\tilde{o}_{t'}, \hat{\pi}_{t'}, \tilde{o}_{t'+H})$
13:　　$\mathcal{L}_s \leftarrow \mathcal{L}_s + D_{\text{KL}} \left[ Q_{\phi_s}(s_{t'+H}) \| \mathcal{N}(\mu, \sigma^2) \right]$
14:　　$\mathcal{L}_o \leftarrow \mathcal{L}_o - \mathbb{E}_{Q(s_{t'+H})} \left[ \log P_{\theta_o}(o_{t'+H} | \tilde{s}_{t'+H}) \right]$
15:　　$\mathcal{L}_o \leftarrow \mathcal{L}_o + \beta * D_{\text{KL}} \left[ Q_{\phi_s}(s_{t'+H}) \| \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2) \right]$
16:　$\theta_s \leftarrow \theta_s - \xi \nabla_{\theta_s} \mathbb{E}[\mathcal{L}_s(\theta_s)]$
17:　$\phi_s \leftarrow \phi_s - \gamma \nabla_{\phi_s} \mathbb{E}[\mathcal{L}_s(\phi_o)]$
18:　$\theta_o \leftarrow \theta_o - \eta \nabla_{\theta_o} \mathbb{E}[\mathcal{L}_o(\theta_o)]$
   ▷ POLICY OPTIMIZATION
19:　$\{\tilde{o}_\tau\}^{B_2} \sim \mathcal{M}$
20:　**for** $\tau = 1, 2, ..., B_2$ **do**
21:　　Compute $Q_{\phi_s}(s_\tau)$ using $\tilde{o}_\tau$
22:　　Sample $\tilde{s}_\tau \sim Q_{\phi_s}(s_\tau)$
23:　　**for** $s = 1, 2, ..., S_1$ **do**
24:　　　Compute $\mu, \sigma \leftarrow P_{\theta_s}(s_{\tau+H} | \tilde{s}_\tau, \hat{\pi}_t)$
25:　　　Sample $\tilde{s}_{\tau+H} \sim \mathcal{N}(\mu, \sigma^2)$
26:　　　Compute $P_{\theta_o}(o_{\tau+H} | \tilde{s}_{\tau+H})$
27:　　　Compute $Q_{\phi_s}(\tilde{s}_{\tau+H})$ using $\tilde{o}_{\tau+H}$
28:　　　Compute $\mu', \sigma' \leftarrow Q_{\phi_s}(\tilde{s}_{\tau+H})$
29:　　　$G \leftarrow G - \log \Psi \left[ P_{\theta_o}(o_{\tau+H} | \tilde{s}_{\tau+H}) \right]$
30:　　　$G \leftarrow G + [H(\mu', \sigma') - H(\mu, \sigma)]$
31:　　　**for** $s = 1, 2, ..., S_2$ **do**
32:　　　　Sample $\tilde{s}_{\tau+H} \sim P_{\theta_s}(s_{\tau+H} | \tilde{s}_\tau, \hat{\pi}_\tau)$ ▷ Re-computed with dropout.
33:　　　　Compute $\mu'', \sigma'' \leftarrow P_{\theta_o}(o_{\tau+H} | \tilde{s}_{\tau+H})$
34:　　　　Sample $\tilde{s}_{\tau+H} \sim \mathcal{N}(\mu, \sigma^2)$
35:　　　　Compute $\mu''', \sigma''' \leftarrow P_{\theta_o}(o_{\tau+H} | \tilde{s}_{\tau+H})$
36:　　　　$G \leftarrow G + [H(\mu'', \sigma'') - H(\mu''', \sigma''')]$
37:　$\phi_a \leftarrow \phi_a - \alpha \nabla_{\phi_a} \mathbb{E}[G(\phi_a)]$

**Agent components:**
　Model:
　　Encoder $Q_{\phi_s}$.
　　Transition $P_{\theta_s}$.
　　Decoder $P_{\theta_o}$.
　Actor $Q_{\phi_a}$.
　Actor mapping $\Pi$.
　Preference mapping $\Psi$.

**Other components:**
　Environment $E$.
　Experience Memory $\mathcal{M}$.

**Hyperparameters:**
　Iterations $N$.
　Beta $\beta$.
　Horizon $H$.
　Batch size $B_1$, $B_2$.
　Sample size $S_1$, $S_2$.
　Learning rate $\xi, \gamma, \eta, \alpha$.

**Run** Model$(\tilde{o}_i, \hat{\pi}, \tilde{o}_{i+H})$:
　Compute $Q_{\phi_s}(s_i)$ using $\tilde{o}_i$
　Sample $\tilde{s}_i \sim Q_{\phi_s}(s_i)$
　Compute $\mu, \sigma \leftarrow P_{\theta_s}(s_{i+H} | \tilde{s}_i, \hat{\pi})$
　Compute $Q_{\phi_s}(\tilde{s}_{i+H})$ using $\tilde{o}_{i+H}$
　Compute $\mu', \sigma' \leftarrow Q_{\phi_s}(\tilde{s}_{i+H})$
　Sample $\tilde{s}_{i+H} \sim \mathcal{N}(\mu, \sigma^2)$
　Compute $P_{\theta_o}(o_{i+H} | \tilde{s}_{i+H})$

**POLITECNICO** MILANO 1863 | DIPARTIMENTO DI MECCANICA

# Benchmark: AIF vs Realistic Industrial Environment

## Application

Energy-Efficient Control of simulated workstations within automotive manufacturing system composed of parallel, identical machines.

## Challenges

Stochastic - Delayed - Long-horizon - Multi-Modal Observation

Requires extensive planning with horizon of one-shift (~ 3000 actions).





Active Inference ←——————————————→ Reinforcement Learning

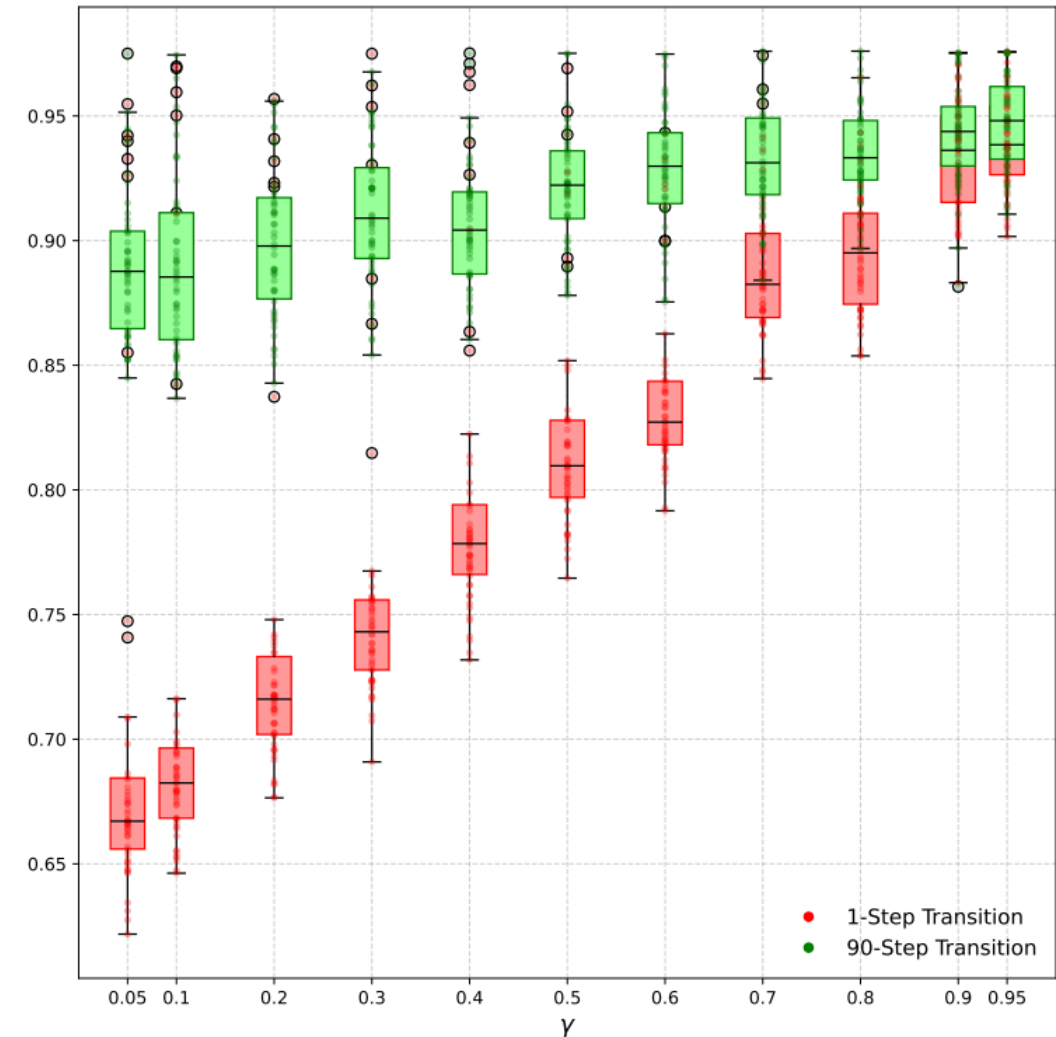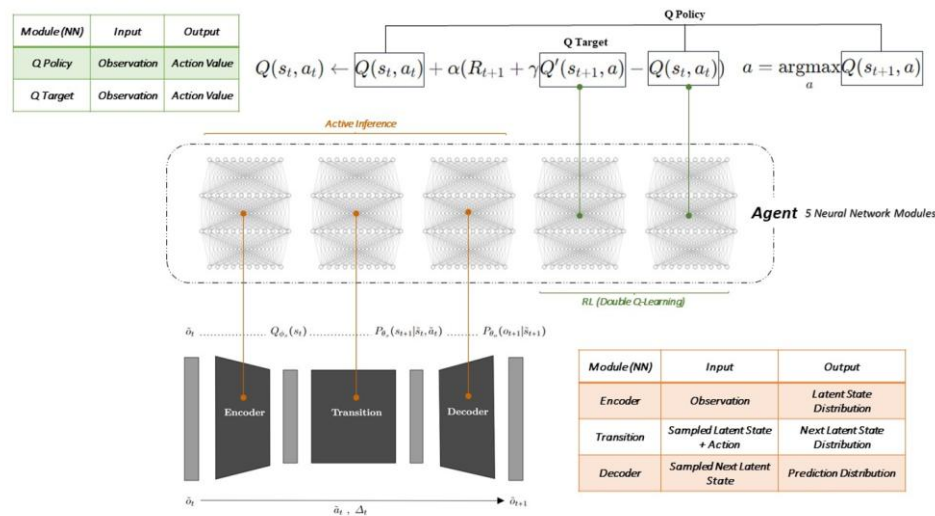# Results: DAIF vs Realistic Industrial Environment

## Application

Energy-Efficient Control of simulated workstations within automotive manufacturing system composed of parallel, identical machines.

## Challenges

Stochastic - Delayed - Long-horizon - Multi-Modal Observation

Requires extensive planning with horizon of one-shift (~ 3000 actions).

## Improving Energy-Efficiency:

Controlling the number of active machines to minimize idle time with negligible production loss.



| Agent($\phi$) | Production Loss [%] | EN Saving [%] |
|---|---|---|
| DQN (0.93) | $4.82 \pm 0.34$ | $10.87 \pm 0.76$ |
| DQN (0.94) | $3.34 \pm 0.23$ | $9.92 \pm 0.69$ |
| **DAIF** | $\mathbf{2.59 \pm 0.16}$ | $\mathbf{12.49 \pm 0.04}$ |
| DQN (0.95) | $1.27 \pm 0.05$ | $7.00 \pm 0.07$ |
| DQN (0.96) | $1.27 \pm 0.09$ | $7.62 \pm 0.12$ |
| DQN (0.97) | $1.20 \pm 0.05$ | $7.72 \pm 0.10$ |
| DQN (0.98) | $0.54 \pm 0.04$ | $2.72 \pm 0.19$ |
| DQN (0.99) | $0.40 \pm 0.03$ | $2.46 \pm 0.01$ |

*Production loss versus energy-saving (EN) across reward parameters $\phi$ of DQN agents (best model-free RL) against the DAIF agent .*

*Performance of the agents versus overshooting horizon H.*

POLITECNICO MILANO 1863 | DIPARTIMENTO DI MECCANICA

CONNECTS

# DAIF: Promises



Derived from a Bayesian-grounded framework (heuristic-free)

Natively scalable to both discrete and continuous action spaces

A unified explore–exploit gradient

Effective encoding of long-horizon dynamics

Capturing distinct forms of epistemic uncertainty, providing an intrinsic drive for model refinement

Reliance on observations (cheap, raw data), instead of often expensive, engineered reward signals

A degree of interpretability

Inherent adaptability in non-stationary settings via AIF

Minimal computational cost during both inference (model-free) and planning (EFE calculated in a single H-step forward pass)
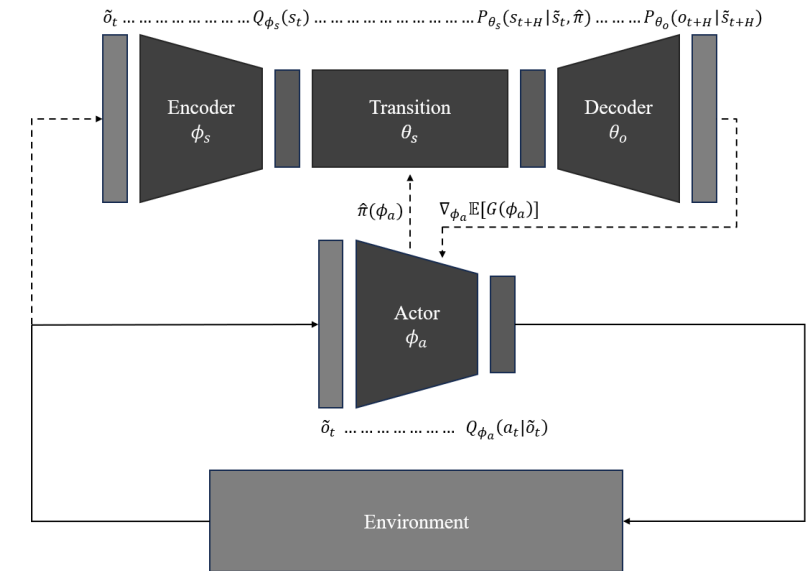
# DAIF: Potential

Scope: Data-Driven Decision Making or Optimization:

Sequential decision making under uncertainty (POMDP)

- Control / Planning / Search
- Reasoning / Inference / Search
- Generation / Synthesis

Domains:

- Robotics
- Production Control
- Healthcare
- Drug Discovery / Protein Design
- …

Settings:

Delayed and Long-Horizon Environments

- Stochastic & Partial Observability
- Delayed/Sparse/Expensive Reward (i.e., Scarce high-quality data)
- Long Horizon



$\tilde{o}_t$ ..................... $Q_{\phi_a}(a_t|\tilde{o}_t)$

Agent $\phi_a$

~ $\tilde{a}_t$

Task Environment



+500 Goal

+10 Misleading Goal Post

-4 +2.67

Two sparse rewards [1]

[1] Chakraborty et al., 2022

POLITECNICO MILANO 1863 | DIPARTIMENTO DI MECCANICA

# Future Work

DAIF is a concept: Generative world model with multi-step latent transitions + differentiable policy; backpropagates EFE through long horizons without tree search; scales to all spaces while keeping the explore–exploit balance.

Next steps

Generative Model (DL): Try diffusion/flow-matching world models instead of VAE.

Sample efficiency (DL): Still needs experience collection every H steps; aim to reduce data needs.

    Sequence aggregation: Replace slow recurrent unrolling with set-based pooling of H embeddings (with simple positional encodings) or operator-learning for resolution-invariant aggregation.

Optimization (RL): Explore actor–critic training and add regularization to stabilize/variance-reduce EFE gradients.

EFE Methods (AIF): Improving distinct term estimations (e.g., sampling).

Adaptation (AIF): Focus on rapid learning in non-stationary environments.

Expanding Experimental: In addition to more domains, looking a range of settings.

Reasoning Model: Exploring potential for reasoning tasks.

Big picture: Bridges generative world-modelling with active inference and RL—compact, end-to-end probabilistic agents

**POLITECNICO MILANO 1863** | DIPARTIMENTO DI MECCANICA

Hi CONNECTS

# References

Yeganeh, Yavar Taheri, Mohsen Jafari, and Andrea Matta. "Deep Active Inference Agents for Delayed and Long-Horizon Environments." *arXiv preprint arXiv:2505.19867* (2025).



*Paper*



*Code*

[1] Chakraborty, Souradip, et al. "Dealing with sparse rewards in continuous control robotics via heavy-tailed policies." arXiv preprint arXiv:2206.05652 (2022).

[2] Lanillos, P., Meo, C., Pezzato, C., Meera, A. A., Baioumy, M., Ohata, W., ... & Tani, J. (2021). Active inference in robotics and artificial agents: Survey and challenges. arXiv preprint arXiv:2112.01871.

[3] Huang, Y., Li, Y., Matta, A., & Jafari, M. (2024). Navigating Autonomous Vehicle on Unmarked Roads with Diffusion-Based Motion Prediction and Active Inference. arXiv preprint arXiv:2406.00211.

[4] Fountas, Z., Sajid, N., Mediano, P., & Friston, K. (2020). Deep active inference agents using Monte-Carlo methods. Advances in neural information processing systems, 33, 11662-11675.

[5] Hafner, D., Pasukonis, J., Ba, J., & Lillicrap, T. (2023). Mastering diverse domains through world models. arXiv preprint arXiv:2301.04104.

**POLITECNICO** | DIPARTIMENTO
MILANO 1863 | DI MECCANICA

Hi CONNECTS