# Action and Perception as Divergence Minimization

Danijar Hafner[12], Pedro Ortega[3], Jimmy Ba[2], Thomas Parr[4], Karl Friston[4], Nicolas Heess[3]

[1]Google Brain   [2]University of Toronto   [3]DeepMind   [4]University College London

## 1 Overview

Action and Perception as Divergence Minimization (APD)...

1. Map of **possible agent objectives** that correspond to different latent variables, target factorizations, and divergence measures
2. **Unified perspective** on representation learning, infogain exploration, and empowerment
3. Representation learning should be **paired** with infogain exploration for a temporally consistent objective
4. **World models** as path toward adaptive infomax agents while making task rewards optional
5. **Future objectives** should be derived from a joint divergence to facilitate comparison and make target explicit

## 2 Agents with Latent Variables

Perception



**Agent Beliefs**
actions, objects, rules, etc

Action

**Input Sequence**
images, proprioception, etc

Parameterize belief (incl actions) by $\phi$

## 3 Joint KL Minimization

Formulate agent objective as bringing its current **actual distribution** toward a **target distribution**:

actions, state estimates, parameters, skills

replay buffer

on-policy data or planning

$$\min_\phi \mathrm{KL}\left[\quad \| \quad\right] \quad \text{many different options}$$

**Actual distribution**    **Target distribution**
$$p_\phi(x,z) \qquad \tau(x,z)$$

$x$ lifetime trajectory of inputs
$z$ set of agent latents
$\phi$ parameters of agent beliefs   $p_\phi(z\,|\,x_<)$ over repr and actions

## 4 Target Dependencies

**Factorized Targets**
$$\tau(x)\tau(z)$$

Inputs and latents have zero mutual information under the target

Agent **minimizes** mutual information in the actual dist

Examples
- MaxEnt RL uses reward factor and action prior to solve the task while keeping actions as random as possible

**Expressive Targets**
$$\tau(x\,|\,z)\tau(z)$$

Target knows or learns depen. between inputs and latents

Agent **maximizes** the mutual information in the actual dist

Examples
- World models learn representations that are informative of past inputs
- Reverse predictors learn skills that maximally influence future inputs

## 5 Information Bounds

Minimizing joint KL to an expressive target...
- realizes the **preferences** expressed by the target
- maximizes variational bound on the **mutual information** between inputs and latents
- bound is tighter the better the target can express dependencies

simplicity    accuracy    input entropy

$$\underbrace{\mathrm{KL}\left[p_\phi(x,z)\,\|\,\tau(x,z)\right]}_{\text{joint divergence}} = \underbrace{\mathrm{E}\,\mathrm{KL}\left[p_\phi(z\,|\,x)\,\|\,\tau(z)\right]}_{\text{realizing latent preferences}} - \underbrace{\mathrm{E}\left[\ln\tau(x\,|\,z) - \ln p_\phi(x)\right]}_{\text{information bound}}$$

control    information gain

$$\underbrace{\mathrm{KL}\left[p_\phi(x,z)\,\|\,\tau(x,z)\right]}_{\text{joint divergence}} = \underbrace{\mathrm{E}\,\mathrm{KL}\left[p_\phi(x\,|\,z)\,\|\,\tau(x)\right]}_{\text{realizing input preferences}} - \underbrace{\mathrm{E}\left[\ln\tau(z\,|\,x) - \ln p_\phi(z)\right]}_{\text{information bound}}$$

## 6 Past and Future

Agents with expressive targets...
- infer latent representations that are **informative** of past inputs
- explore future inputs that are **informative** of the representations

$$\mathrm{I}[x;z] = \mathrm{I}[x_<;z] + \mathrm{I}[x_>;z\,|\,x_<]$$



Actual distribution $p_\phi$    Target distribution $\tau$   (see eq 6 in the paper)

## 7 Types of Latents

| Latent Variable | Past Infomax | Future Infomax |
| --- | --- | --- |
| Actions | N/A<br>past actions are observed | Empowerment & MaxEnt RL<br>VIM, ACIE, EPC, SQL, SAC |
| Skills | N/A<br>past skills are observed | Skill Discovery<br>VIC, SNN, DIAYN, VALOR |
| State Estimates | State Estimation<br>VAE, DVBF, SOLAR, PlaNet | State Information Gain<br>NDIGO, DVBF-LM |
| Dynamics Parameters | System Identification<br>PETS, Bayesian PlaNet | Dynamics Information Gain<br>VIME, MAX, Plan2Explore |
| Policy Parameters | Belief over Policies<br>BootDQN, Bayesian DQN | Policy Information Gain<br>BootDQN, Bayesian DQN |

## 8 Action Perception Cycle

Action and perception optimize the same objective but receive and affect different variables.

Under a unified target distribution...
- perception makes the **agent beliefs consistent with the world**
- actions make the **world consistent with the agent beliefs**

action

$$\min_\phi \mathrm{KL}\left[\underbrace{p_\phi(z\,|\,x)}_{\text{beliefs}}\underbrace{p_\phi(x)}_{\text{inputs}}\,\|\,\underbrace{\tau(x,z)}_{\text{target}}\right]$$

perception

## 9 Niche Seeking

Minimizing a joint divergence also brings the marginals together

$$\mathrm{KL}\left[p_\phi(x,z)\,\|\,\tau(x,z)\right] \geq \mathrm{KL}\left[p_\phi(x)\,\|\,\tau(x)\right]$$

The marginal target distribution over inputs is the marginal likelihood

The agent thus converges to an **ecological niche**...
- See inputs propto how well the agent can learn to predict them
- That is large because of the information gain exploration
- That it can inhabit despite external perturbations

The agent thus seeks out a large niche that it can inhabit and understand

Models that assign high prob to **more trajectories** lead to **larger niches**