# Towards Governing Agent's Efficacy:
# Action-Conditional $\beta$-VAE for Deep Transparent Reinforcement Learning

**John Yang[1], Gyujeong Lee[1], Minsung Hyun[1], Simyung Chang[1,2], Nojun Kwak[1]**

[1]Seoul National University, Seoul, South Korea
[2]Samsung Electronics, Suwon, South Korea
{yjohn, regulation.lee, minsung.hyun, timelighter, nojunk}@snu.ac.kr

## Abstract

We tackle the blackbox issue of deep neural networks in the settings of reinforcement learning (RL) where neural agents learn towards maximizing reward gains in an uncontrollable way. Such learning approach is risky when the interacting environment includes an expanse of state space because it is then almost impossible to foresee all unwanted outcomes and penalize them with negative rewards beforehand. Unlike reverse analysis of learned neural features from previous works, our proposed method tackles the blackbox issue by encouraging an RL policy network to learn interpretable latent features through an implementation of a disentangled representation learning method. Toward this end, our method allows an RL agent to understand self-efficacy by distinguishing its influences from uncontrollable environmental factors, which closely resembles the way humans understand their scenes. Our experimental results show that the learned latent factors not only are interpretable, but also enable modeling the distribution of entire visited state space with a specific action condition. We have experimented that this characteristic of the proposed structure can lead to ex post facto governance for desired behaviors of RL agents.

## Introduction

Despite many recent successful achievements that deep neural networks (DNN) have allowed in machine learning fields (Krizhevsky, Sutskever, and Hinton 2012; LeCun, Bengio, and Hinton 2015; Mnih et al. 2015), the legibility of their high-level representations are noticeably less studied compared to the relevant studies which rather prioritize performance enhancements or task completions. The blackbox issue of neural networks has been many times neglected and such technical opacity has been excused for their vast performance improvements (Burrell 2016).

While the opaqueness of DNN comes handy when strict labels are available for every data sample, its blackbox issue is a great element of risk especially in reinforcement learning (RL) settings where machines, or agents, are allowed to have highly intertwined interactions with their environments. Since an RL agent's policy on action selection is

optimized towards maximizing the rewards, it may produce harmful and unexpected outcomes if these outcomes are not primarily penalized with negative reward signals.

Yet, too much regulation would, contrarily, result in misusing the full potential of the technology (Rahwan 2018). RL is proven of its powerfulness over humans by, for an example of AlphaGo, figuring to learn unprecedented winning moves (Silver et al. 2017). Interfering in the learning process to control the model's resultant behaviors as done in the work of (Christiano et al. 2017) may not be efficient in governing RL agents. Rather, it is desired to control the efficacy of an agent which is already optimized for the environment.

In order to rule AI agents efficiently, humans who govern first need to comprehend how AI machines perceive their world and monitor their efficacy (Stilgoe 2018; Wynne 1988). Higgins et al. modeled an environment with the $\beta$-Variational Autoencoder ($\beta$-VAE) to generate disentangled latent features (Higgins et al. 2017), purposefully inducing the learned features to be interpretable to human (Higgins et al. 2016b), and have applied the features for transfer learning across multiple environments. We are motivated that this method can be utilized to train an explainable RL agent (Higgins et al. 2016a).

We believe building transparent RL agents and governing them would solve issues mentioned above. In this paper, we propose a method that allows training a deep but transparent RL policy network, encouraging their latent features to be interpretable. We intend to accomplish this by training RL agents to learn disentangled representations of their world in egocentric perspective with action-conditional $\beta$-VAE (AC-$\beta$-VAE): the learned control-dependent latent features and uncontrollable environmental factors are disentangled while the learned factors are also able to model the environment. Our strategic design that engage the AC-$\beta$-VAE and an RL policy network to share a backbone structure overcomes the blackbox issue, supporting the transparency of deep RL. We also empirically show that the behavior of our agents can further be governed with human enforcements.

## Related Work

Deep learning methods are praised of their unruled pattern extraction that yields better performance in many tasks than machines trained under human prior knowledge (Günel ; Moore and Lu 2011; Vanderbilt 2012), but as stated earlier,
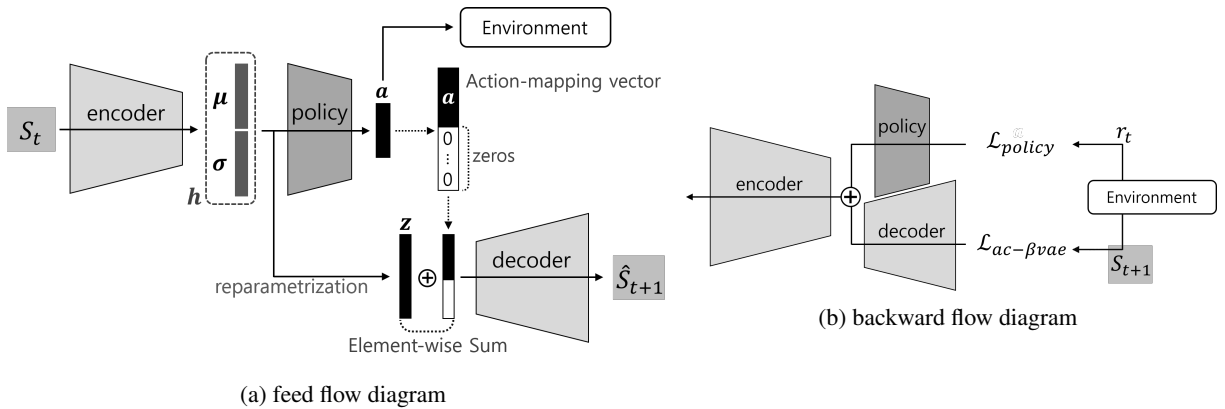
Figure 1: The structure and flow diagrams of the proposed AC-$\beta$-VAE for a transparent policy network. The proposed network requires training samples of MDP tuples of RL environments that consist of $(s_t, a_t, r_t, s_{t+1})$ where $s_t$, $a_t$ and $r_t$ are respectively state, action and reward at time step $t$. The action-conditional decoder encourages the input features of the policy network to be disentangled and interpretable. Since the encoder + policy network can be seen as one big policy network that takes raw states as inputs, its inner intentions in selecting actions for a desired next state can thus be explained visually through the outputs of the decoder.

the blackbox characteristic of DNNs can be precarious especially in the RL setting. One of the safety factors of AI development suggested in (Amodei et al. 2016) is avoidance of negative side effects when training an agent to complete a goal task with a strict reward function.

Attempts to open the blackbox of DNN and to understand the inner system of neural networks have been made in many recent works (Lipson and Kurman 2016; Zeiler and Fergus 2014; Bojarski et al. 2017; Greydanus et al. 2017). Its inherent learning phenomena are reversely analyzed by observing the resultant learned understructure. While the training progress is also analytically interpreted via information theory (Shwartz-Ziv and Tishby 2017; Saxe et al. 2018), it is still challenging to anticipate how and why high-level features in neural models are learned in a certain way before training them. Since learning a disentangled representation encourages its interpretability (Bengio, Courville, and Vincent 2013; Higgins et al. 2016b), it is previously reported that features of convolutional neural networks (CNN) can also be learned in a visually explainable way (Zhang and Zhu 2018) through disentangled representation learning.

Prospection of future states conditioned by current actions is meaningful to RL agents in many ways, and action-conditional (variational) autoencoders are learned to predict sequent states in the works of (Ha and Schmidhuber 2018; Oh et al. 2015; Thomas et al. 2017). DARLA (Higgins et al. 2017) utilizes disentangled latent representations for cross-domain zero-shot adaptations. It aims to prove its representation power in multiple similar but different environments. Our model may also look similar to conditional generative models like Conditional Variational Autoencoders (CVAE) (Sohn, Lee, and Yan 2015) and InfoGan (Chen et al. 2016), but these are not directly applicable models to RL domains.

## Preliminary: $\beta$-VAE

Variational autoencoder (VAE) (Kingma and Welling 2013) works as a generative model based on the distribution of training samples (Co-Reyes et al. 2018; Babaeizadeh et al. 2017). VAE's goal is to learn the marginal likelihood of a sample $x$ from a distribution parametrized by generative factors $z$. In doing so, a tractable proxy distribution $q_\phi(z|x)$ is used to estimate an intractable posterior $p_\theta(z|x)$ with two different parameter vectors $\phi$ and $\theta$. The marginal likelihood of a data point $x$ can be defined as:

$$\log p_\theta(x) = D_{KL}(q_\phi(z|x)||p_\theta(z|x)) + L(\theta, \phi, x, z). \quad (1)$$

Since the KL divergence term $D_{KL}(\cdot||\cdot)$ is non-negative, $L_{vae} \triangleq L(\theta, \phi, \mathbf{x}, \mathbf{z})$ sets a variational lower bound for the likelihood $\log p_\theta(x)$ and the best approximation $q_\phi(z|x)$ for $p_\theta(z|x)$ can be obtained by maximizing this lower bound:

$$L_{vae} = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x)||p(z)). \quad (2)$$

In practice, $q_\phi$ and $p_\theta$ are respectively encoder and decoder that are parameterized by deep neural networks, and the prior $p(z)$ is usually set to follow Gaussian distribution $\mathcal{N}(0, I)$. The gradients of the lower bound can be approximated using the *reparametrization trick*.

$\beta$-VAE (Higgins et al. 2016b) extends the work and drives VAE to learn disentangled latent features, weighting the KL-divergence term from the VAE loss function (negative of the lower bound) with a hyper-parameter $\beta > 1$:

$$L_{\beta vae} = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - \beta D_{KL}(q_\phi(z|x)||p(z)). \quad (3)$$

When $\beta$ is ideally selected and does not severely interfere the reconstruction optimization, each latent factor of $z$ is learned to be not only independent of each other, but also interpretable. This means the resultant features follow physio-visual characteristics of our world and differ from conventional DNN features that are not so human-friendly.

## The Proposed Model

Our proposed model is composed of two structures: a policy gradient RL method and the action-conditional $\beta$-VAE (AC-$\beta$-VAE). As shown in Figure 1, both components are designed to strategically share first layers of the encoding network so that the latent features of AC-$\beta$-VAE can also become the input of the policy network. This simple shared architecture enables human-level interpretations on behaviors of deep RL methods.

Consider a reinforcement learning setting where an actor plays a role of learning policy $\pi_\psi(a_t|s_t)$ and selects an action $a \in \mathcal{A}$ given a state $s \in \mathcal{S}$ at time $t$, and there exists a critic that estimates value of the states $V_w(s)$ to lead the actor to learn the optimal policy. Here, $\psi$ and $w$ respectively denote the network parameters of the actor and the critic. Training progresses towards the direction of maximizing the objective function based on cumulative rewards, $J(\theta) = \mathbb{E}_{\pi_\psi}[\sum_t \gamma^t r_t]$ where $r_t$ is the instantaneous reward at time $t$ and $\gamma$ is a discount factor. The policy update objective function to maximize is defined as follows:

$$L_{policy} = \mathbb{E}_\pi[\log \pi_\psi(s_t, a_t) A^\pi(s_t, a_t)]. \qquad (4)$$

Here, $A^\pi(s, a)$ is an advantage function, which is defined as it is in asynchronous advantage actor-critic method (A3C) (Mnih et al. 2016):

$$A^\pi(s_t, a_t) = \sum_{i=0}^{k-1} \gamma^i r(s_{t+i}, a_{t+i}) + \gamma^k V_w^\pi(s_{t+k}) - V_w^\pi(s_t),$$

where $k$ denotes the number of steps. We have used the update method of *Advantage Actor Critic* (A2C) (Wu et al. 2017), a synchronous and batched version of A3C, for Atari domain environments (Bellemare et al. 2013). *Proximal Policy Optimization* (PPO) (Schulman and Klimov 2017) is also used for our experiments in continuous control environments, which reformulates the update criterion with the use of clipping objective constraint $\mathcal{C}$ in the form of:

$$L_{policy} = \mathbb{E}_\pi \left[ \frac{\pi_\psi(a|s)}{\pi_\psi^{old}(a|s)} A(s, a) \right] - \mathcal{C}D_{KL}(\pi_\psi^{old}(\cdot|s)||\pi_\psi(\cdot|s)). \qquad (5)$$

Here, the subscript $t$ for $a$, $s$ and $A$ is omitted for brevity.

### Action-Conditional $\beta$-VAE (AC-$\beta$-VAE)

As shown in Fig. 1 with a given environment, the policy network combined with the encoder produces rollouts of typical Markov tuples that consist of $(s_t, a_t, r_t, s_{t+1})$. A raw state $s_t$ feeds into the encoder model and gets encoded into a representation $h \in \mathbb{R}^{2n}$, where $n$ is the dimension of the the latent space. Since the policy network and AC-$\beta$-VAE share the parameters until this encoding process, the representation $h = [\mu^T, \sigma^T]^T$ represents a DNN feature which is inputted to the policy network while also representing a concatenated form of the mean and the standard deviation vectors $\mu, \sigma \in \mathbb{R}^n$. The vectors are reparametrized into a posterior variable $z \in \mathbb{R}^n$ through the AC-$\beta$-VAE pipeline. The output of the encoder feed-flows into the policy network $\pi(a|h)$ to output an optimal action $a \in \mathbb{R}^m$ where

---

**Algorithm 1** AC-$\beta$-VAE with an actor-critic policy network

Initialize encoder $q_\phi(h|s)$ and decoder $p_\theta(s|z)_{z \sim \mathcal{N}(h)}$
Initialize critic $V_w(s)$, actor $\pi_\psi(a|h)$ and state $s$.
**while** not stop-criterion **do**
    $t_{start} = t$
    **repeat**
        Take an action $a_t$ with policy $\pi_\psi$
        Receive new state $s_{t+1}$ and reward $r_t$
    **until** $t - t_{start} \geq$ *number of step* **or** terminal $s_t$
    $R = \begin{cases} 0 & \text{for terminal } s_t \\ V_w(s_t) & \text{for non-terminal } s_t \end{cases}$
    **for** $i \in \{t-1, ..., t_{start}\}$ **do**
        $R \Leftarrow r_i + \gamma R$
        Compute $A(s_i, a_i)$ (for A2C or PPO)
        Sample $z_i \sim \mathcal{N}(h_i)$ and create $a_i^{map}$
        Predict $p_\theta(\hat{s}_{i+1}|z_i + a_i^{map})$
        Compute $L_{policy}$ and $L_{ac-\beta vae}$
        Update encoder, actor and decoder based on:
            $L_{total} = L_{policy} + \alpha L_{ac-\beta vae}$
        Update critic by minimizing the loss:
            $L_{critic}(w) = (R - V_w(s_i))^2$

---

$m < n$ so that an RL environment responds accordingly. The action vector $a$ is then concatenated with a vector of zeros in length of $\mathbb{R}^{n-m}$ to create, we call, an *action-mapping vector* $a^{map} = [a^T, 0^T]^T \in \mathbb{R}^n$. An element-wise sum of the latent variable $z$ and the action-mapping vector $a^{map}$ is performed in order to map action-controllable factors into the latent vector. This causes the latent variable sampled to be constrained by the probability of actions. The resultant vector $z_t + a_t^{map}$ is fed into the decoder network to predict the next state $\hat{s}_{t+1}$. The prediction is then compared with the real state $s_{t+1}$ given by the environment after the action taken. For an MDP tuple collected at time $t$, the loss of AC-$\beta$-VAE is computed with the following loss function:

$$\begin{aligned} L_{ac-\beta vae} =& \mathbb{E}_{q_\phi(h_t|s_t)}[\log p_\theta(s_{t+1}|z_t + a_t^{map})]_{z_t \sim \mathcal{N}(h_t)} \\ & - \beta D_{KL}(q_\phi(z_t|s_t)||\mathcal{N}(0, I)). \end{aligned} \qquad (6)$$

As one can see, the AC-$\beta$-VAE model can be trained either simultaneously with the policy network or separately, and all our experiments are performed with the former because it is more practical. At each iteration of update, the total objective function value is calculated with the weighted sum of objective function values from both models:

$$L_{total} = L_{policy} + \alpha L_{ac-\beta vae} \qquad (7)$$

where $\alpha$ is the weight balance parameter. Since exploration based on the error between generated outputs and the ground-truths have already been proven on the training enhancement in many RL related works (Oh et al. 2015; Ha and Schmidhuber 2018; Tang et al. 2017), our model rather focuses on feasible training of a transparent neural policy network and modeling self-efficacy of agents, not on RL performance improvement. We thus choose relatively small-valued $\alpha$ not to confuse the policy network too much. A basic pseudo-code for the training scenario of our proposed structure is provided in Algorithm 1.
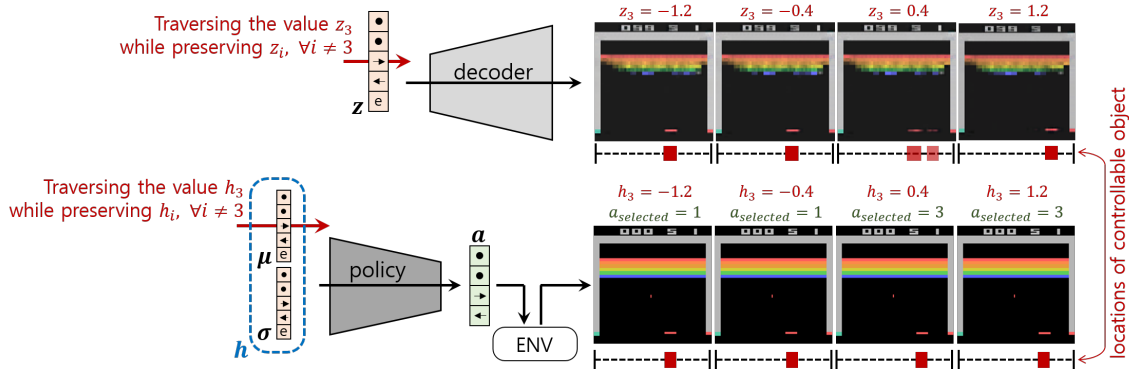
Figure 2: The results of traversing the latent factor of our trained model on Atari game environment BREAKOUT with $z \in \mathbb{R}^5$, where $z_{1:4}$ are mapped with variant features of $a \in \mathbb{R}^4$ and $z_5$ is condensed with other environmental factors. Since the factors in the latent vector $z$ of AC-$\beta$-VAE are defined by the vectors of mean and standard deviation $\mu, \sigma$, traversing $i$-th value of the latent vector $z_i$ is almost equivalent to traversing $\mu_i$. The input DNN feature $h$ of the policy network is the concatenation of $\mu$ and $\sigma$, and thus the next state due to its output actions $a_{selected}$ caused by traversed $\mu_i$ factor would be probabilistically predictable by the visual consequence estimated by the decoder with traversed $z_i$.

## Mapping Action-Controllable Representations

Learning visual influence was previously introduced of its importance and implicitly solved in the works of (Oh et al. 2015; Greydanus et al. 2017). Distinguishing directly-controllable objects and environment-dependent objects reflects much of how a human perceives the world. Restricting in the world of Atari game domains as an example, it is intuitive for a human agent to first figure out 'where I am in the screen' or 'what I am capable of with my actions' and then work their ways towards achieving the highest score.

We show in the experiment section that AC-$\beta$-VAE allows RL agents not only to explicitly learn visual influences of their actions, but also learn them in a human-friendly way. By traversing each element of the latent vector, we are able to interpret which dimensions are mapped with actions and which are mapped with other environmental factors.

## Transparent Policy Network

As mentioned earlier, the encoder and the policy network can be grouped as one bigger policy network model with an interpretable layer constrained by the AC-$\beta$-VAE loss. Unlike high-level features from conventional DNN models, the inner features of our policy network are consequentially interpretable.

Figure 2 illustrates how our policy network becomes transparent. If the action-dependent factors are disentangled in the latent vector $z \in \mathbb{R}^n$ and mapped into $z_{1:m}$, then so they are in $\mu_{1:m}$ and $\sigma_{1:m}$ because they define the sampling distribution of $z_i$ where $i$ denotes the dimensional location. The variational samplings from the latent space of VAE is defined as: $z_i = \mu_i + \sigma_i \epsilon_i$ where $\epsilon$ is an auxiliary noise variable $\epsilon \sim \mathcal{N}(0, 1)$. And, we know that $q_\phi(z|x) \prod_i dz_i = p(\epsilon) \prod_i d\epsilon_i$. Since the $\sigma$ value controls mainly the scale of sampled $\epsilon$, traversing $z_i$ is almost equivalent to traversing $\mu_i$[1]. Thus, traversing $\mu_i$ encourages the policy network to

---

[1]Refer the original work of VAE for more insightful details

cause actions as predictions of each traversing value of $z_i$ for $i \leq m$.

## Experiments

In this section, we present experimental results that demonstrate the following key aspects of our proposed method:

- By mapping actions into the latent vector of $\beta$-VAE, action-controllable factors are disentangled from other environmental factors.

- Governance over the optimized behavior of an agent can be made based on human-level interpretation of learned latent behavioral factors.

We have experimented our method in three different environment types: dSprites, Atari and MuJoCo.

**dSprites Environment** is an environment we design with the *dSprites* dataset (Matthey et al. 2017). It originally is a synthetic dataset of 2D shapes that gradually vary in five factors: shape (square, ellipse, heart), scale, orientation, locations in vertical and horizontal axes, respectively. The environment provides a $64 \times 64$ sized image that embrace two shapes, one heart and one square. At each time step, the square is randomly scaled in a randomly oriented form at random location within the image. The heart-shaped object responds to one of the following discrete action inputs: move upward, downward, left, right, enlarge, shrink, rotate left and right. All actions can be represented with a 4-dimensional action vector each of which is responsible for a unit of either vertical, horizontal, scaling or rotating movement.

**Atari Learning Environment** is a software framework for assessing RL algorithms (Bellemare et al. 2013). Each frame is considered as a state and immediate rewards are given for every state transitions. Our method is experimented in the Atari game environments of BREAKOUT, SEAQUEST and SPACE-INVADERS.

**MuJoCo Environment** provides a physics engine system for rigid body simulations (E. Todorov and Tassa. 2012;

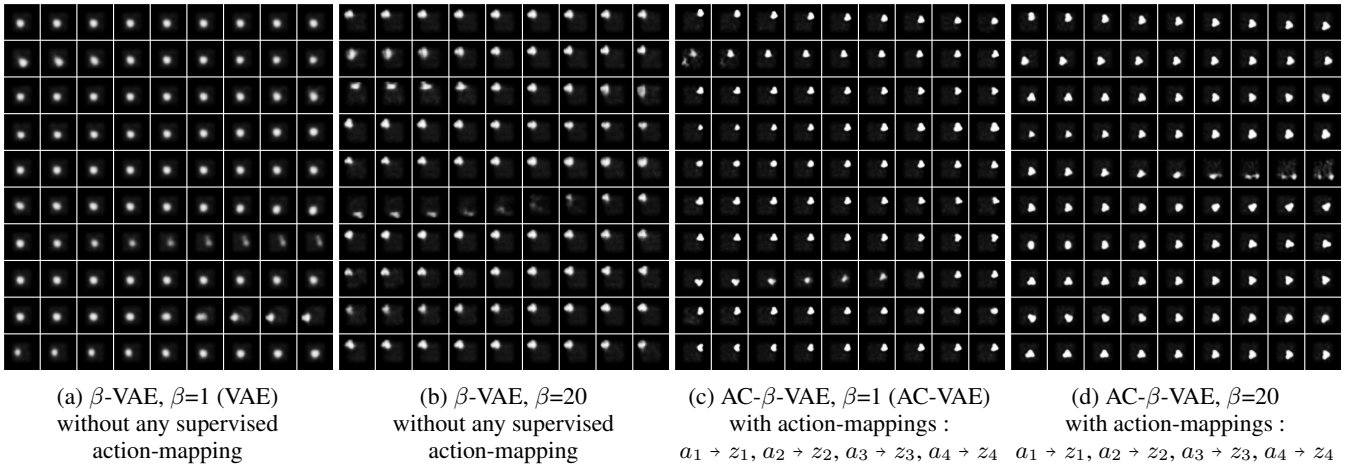|  (a) $\beta$-VAE, $\beta$=1 (VAE) without any supervised action-mapping | (b) $\beta$-VAE, $\beta$=20 without any supervised action-mapping | (c) AC-$\beta$-VAE, $\beta$=1 (AC-VAE) with action-mappings : $a_1 \to z_1, a_2 \to z_2, a_3 \to z_3, a_4 \to z_4$ | (d) AC-$\beta$-VAE, $\beta$=20 with action-mappings : $a_1 \to z_1, a_2 \to z_2, a_3 \to z_3, a_4 \to z_4$ |

Figure 3: The qualitative results of traversing latent factors in $\beta$-VAE with $\beta$=1 (VAE) and $\beta$=20 on $(s_t, s_{t+1})$ data tuples and those of AC-$\beta$-VAE with $\beta$=1 (AC-VAE) and $\beta$=20 on $(s_t, a_t, s_{t+1})$ data tuples in dSprites environment. The action vectors are retrieved randomly as combinations of $(a_1, a_2, a_3, a_4)$ that respectively represent vertical, horizontal, rotational, scaling moves. The vertical axes represent the dimensions of the learned latent vector $z_{1:10}$ from top to bottom while the horizontal axes represent traversing values of $[-2 : 2]$ from left to right.

|  | VAE ($\beta$=1) | $\beta$-VAE ($\beta$=20) | AC-VAE ($\beta$=1) | AC-$\beta$-VAE ($\beta$=20) |
|---|---|---|---|---|
| Avg. Disent. | 0.120 | 0.133 | 0.233 | **0.390** |
| Avg. Compl. | 0.155 | 0.231 | 0.288 | **0.405** |

Table 1: The quantitative scores of disentanglement and completeness averaged over dimensions of the latent vector learned with $(s_t, a_t, s_{t+1})$ tuples from dSprites environment.

G. Brockman and Zaremba. 2016). Four robotics tasks are engaged in our experiments: WALKER2D, HOPPER, HALF-CHEETAH and SWIMMER. A state vector represents the current status of a provided robotic figure, each factor of which is unknown of its physical meaning.

As an encoder and a decoder, we have used a convolutional neural network (CNN) for Atari environments and fully-connected MLP networks for dSprites and MuJoCo environments. For the stochastic policy network, we have used a fully-connected MLP. PPO and A2C are applied to optimize agent's policy for continuous control and discrete actions, respectively. Most of hyper-parameters for the policy optimization are referred from the works of (Schulman and Klimov 2017; Wu et al. 2017).

### Disentanglement & Interpretability

To demonstrate the disentanglement performance and interpretability of the proposed algorithm, we have experimented our method with $(s_t, a_t, s_{t+1})$ tuples from environments mentioned above.

Figure 3 and Table 1 illustrate the results for the dSprites environment. The metric framework suggested in (Eastwood and Williams 2018) with a random forest regressors are applied to present the quantitative results of disentanglement and completeness. The tree depths are determined for the

lowest prediction error of the validation set. Since the metric system is based on the disentanglement for the conventional VAE and $\beta$-VAE, our metric results may not be strictly comparable to the ones reported in the original work. In Fig. 3(a) and (b), the VAE and $\beta$-VAE seem to struggle from learning the pattern between input $s_t$ and the output $s_{t+1}$ without any action constraint because of the randomness of the environmental squared object, creating relatively blurred reconstructions. Such excessive generalization in reconstructions results in low scores in both disentanglement and completeness which means relatively low representational power to reproduce the ground truth variant factors. Although the action conditions and the low-weighted $D_{KL}$ term allow AC-VAE ($\beta$=1) reconstruct sharper images, its relatively low disentanglement pressure results in lower metric scores compared to AC-$\beta$-VAE ($\beta$=20).

The results for the Atari environments in Figure 4 and Figure 5 show that the latent vector trained with our method models the given environment successfully. All the visited state space and learned behaviors can be projected by traversing each dimension of the latent vector. In that sense, our method can be considered as an action-conditional generative model. Because AC-$\beta$-VAE can model the world in an egocentric perspective, all the sequences of (state-action-next state) can be re-simulated. Such trait may advance many RL methods since similar models are used for an exploration guidance (Tang et al. 2017) or as the imagery rehearsals for training (Ha and Schmidhuber 2018).

Figure 6 shows the quantitative results of the traverse experiment on the MuJoCo environment. Numbers on the heatmap represent the standard deviations for each dimension's state values when traversing dimensional factor. The higher standard deviation value in the traverse of a specific dimension means the more effects the traversing dimension have on immediate state changes. Unlike other environments, the
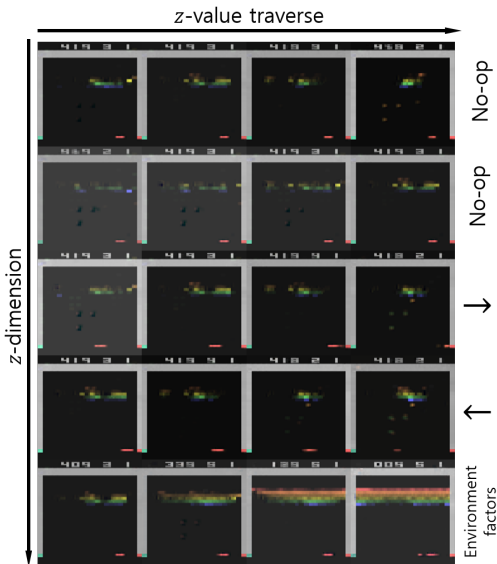
Figure 4: The images are the estimated next states obtained by traversing the latent vector $z \in \mathbb{R}^5$ learned by AC-$\beta$-VAE with $\beta$=10 and $\alpha$=0.001 on the Atari game environment BREAKOUT. The factors at $z_{1:4}$ are mapped with the control factors such as movements of the paddle, and $z_5$ is mapped with the environmental factors such as bricks and the scoreboard.

MuJoCo environment has no environmental factors, and the current state is represented by the preceding movement of the given robotic body. As shown in Figure 6, since the standard deviation of the state values during the traverse of the dimensions that are mapped with actions is larger than the unmapped ones, we can see the proposed algorithm is able to learn the disentangled action-dependent latent features. However, it is limited from clear visual interpretation compared to the experimental cases in other environments because the actions in the MuJoCo environment is defined as a continuous control of torques for all joints and it is conjectured that the movement of one joint affects the whole status of the body.

**Controlling and Governing efficacy**

To verify the controllability of an agent's optimized efficacy, we traverse the latent factors over the environment-specific range during an episode on the learned network. In order to examine $s_{t+1}$, the environment output, the traversal is conducted before reparameterization ($\mu$ vector). Furthermore, to get a clear view on the effect of action-mapped dimensions of the latent vector, we set all of the value of action mapped dimensions to zero except for the traversing one and those unmapped dimension of the latent vector. These experiments are conducted on the Mujoco environments, and traverse range is set as [-5, 5] for every tasks.

The learned behavior in each latent dimension is also depicted in Figure 7. The resultant traverses of action-mapped dimensions on latent factors yield in behavioral movements that are combinations of multiple joint torque values. Un-
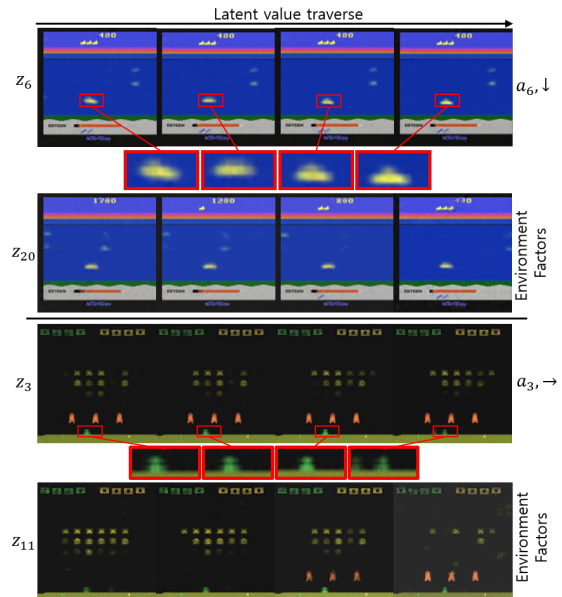


Figure 5: The images are the estimated next states obtained by traversing the latent vector $z \in \mathbb{R}^{20}$ learned by A2C policy and AC-$\beta$-VAE with $\beta$=10 and $\alpha$=0.001 on Atari game environments SEAQUEST (top) and SPACE-INVADERS (bottom) with action spaces of $\mathbb{R}^{18}$ and $\mathbb{R}^6$, respectively. Because of a small movement per action, we have enlarged the ego at a fixed location (red box).

like in Atari environments with discrete action spaces, AC-$\beta$-VAE is constrained with various combinations of continuous action values during training simulations. When the policy network is optimized to accomplish a goal behavior such as walking, the action-mapped latent factors are learned to represent required behavioral components of spreading or gathering the legs. Therefore, $\mu$ vector represents variations in combinations of multiple joint movements, which allows for ease of visual comprehension on agent's optimized efficacy. This clearly shows the possibility of governance over an RL agent's efficacy with human-level interpretations through controlling the values of the $\mu$ vector in the latent space.

We have taken the advantage of our transparent policy network and derived another behavior by controlling learned behavioral components. An RL agent is able to learn with a reward function defined by human preference to perform, for example, a back-flip motion in HOPPER environment (Christiano et al. 2017). Showing a promising result of human enforcements on an RL model, our method enables governance over the agent's optimized behavior in HALF-CHEETAH environment. After identification of behavioral components by traversing each element of the $\mu$ vector, we are able to express another behavior of the agent, a back-flip in this case, as shown in Figure 8.

## Conclusion

In this paper, we propose the action-conditional $\beta$-VAE (AC-$\beta$-VAE) which, for a given input state $s_t$ at time $t$, pre-
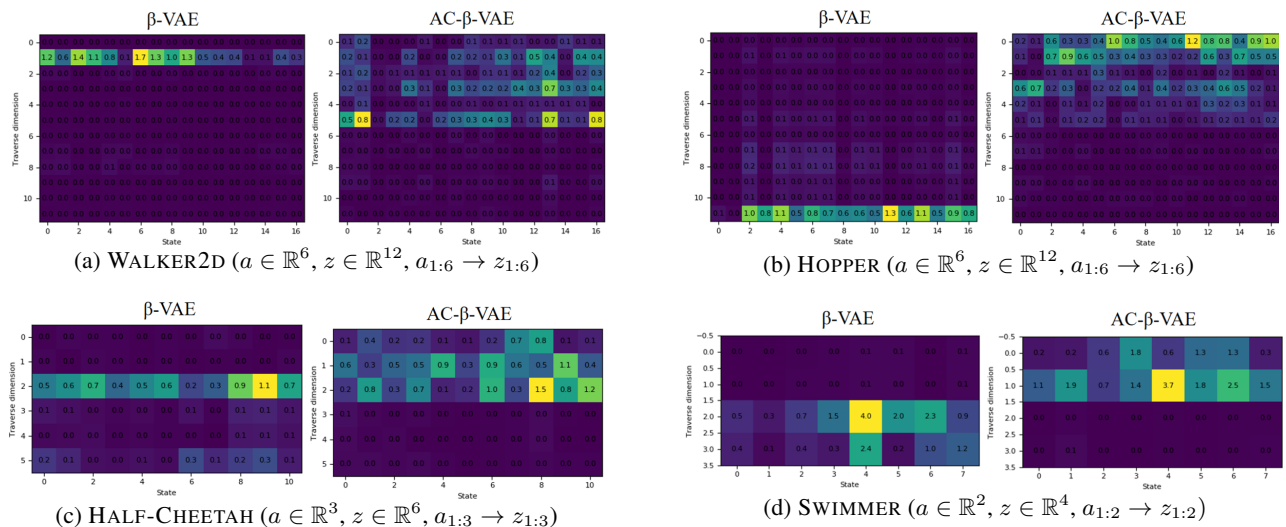
(a) WALKER2D ($a \in \mathbb{R}^6$, $z \in \mathbb{R}^{12}$, $a_{1:6} \to z_{1:6}$)

(b) HOPPER ($a \in \mathbb{R}^6$, $z \in \mathbb{R}^{12}$, $a_{1:6} \to z_{1:6}$)

(c) HALF-CHEETAH ($a \in \mathbb{R}^3$, $z \in \mathbb{R}^6$, $a_{1:3} \to z_{1:3}$)

(d) SWIMMER ($a \in \mathbb{R}^2$, $z \in \mathbb{R}^4$, $a_{1:2} \to z_{1:2}$)

Figure 6: Traverse results in the MuJoCo environments. The numbers in the boxes represent the standard deviations of each dimensional factor of the following state, $s_{t+1}$, when traversing the corresponding dimensional factor of the latent vector. Compared to the traverse for unmapped dimensions, the standard deviations of state values in the action-mapped dimensions are larger. Right arrows indicate action-mapping dimensional locations.
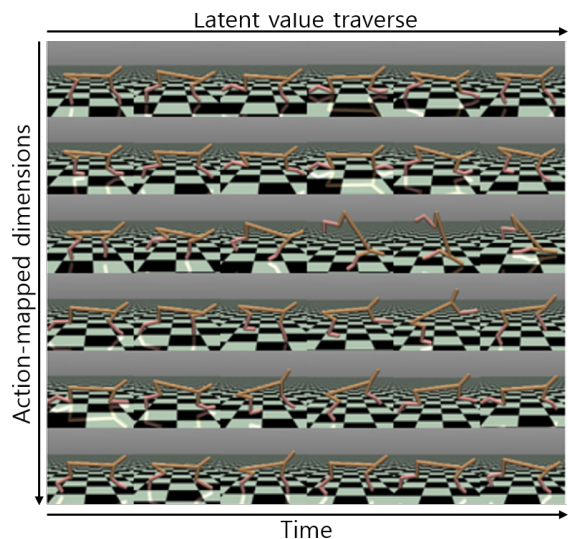


Figure 7: For HALF-CHEETAH environment with continuous control, latent behavioral factors can be interpreted by traversing latent values in time. As a result, each action-mapped latent feature is responsible for a behavioral factor.
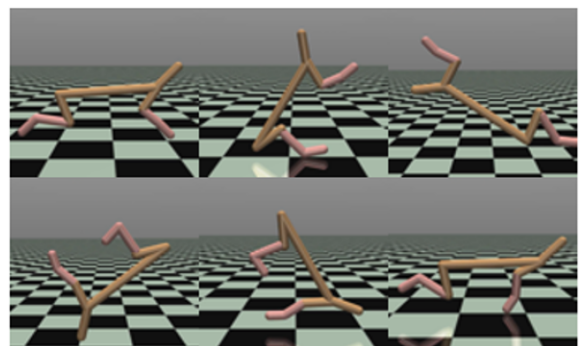


Figure 8: Example of governing the agent movement in MuJoCo environment of HALF-CHEETAH. The robotic body is conducting a back-flip movement which is induced by controlling latent values at first and second dimensions of the learned $\mu$ vector shown in Figure 7.

dicts next state $s_{t+1}$ conditioned on an action $a_t$, sharing a backbone structure with a policy network during a deep reinforcement learning process. Our proposed model not only learns disentangled representations but distinguishes action-mapped factors and uncontrollable factors by partially mapping control-dependent variant features into the latent vector. Since the policy network combined with the preceding encoder can be considered as one bigger policy network that takes raw states as inputs, with AC-$\beta$-VAE, we are able to

build a transparent RL agent of which latent features are interpretable to human, overcoming conventional blackbox issue of Deep RL. Such transparency allows human governance over the agent's optimized behavior with adjustments of learned latent factors. We plan on the relevant studies for applications of the action-mapped latent vector.

## References

[Amodei et al. 2016] Amodei, D.; Olah, C.; Steinhardt, J.; Christiano, P.; Schulman, J.; and Mané, D. 2016. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*.

[Babaeizadeh et al. 2017] Babaeizadeh, M.; Finn, C.; Erhan, D.; Campbell, R. H.; and Levine, S. 2017. Stochastic variational video prediction. *arXiv preprint arXiv:1710.11252*.

[Bellemare et al. 2013] Bellemare, M. G.; Naddaf, Y.; Veness, J.;

and Bowling, M. 2013. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research* 47:253–279.

[Bengio, Courville, and Vincent 2013] Bengio, Y.; Courville, A.; and Vincent, P. 2013. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence* 35(8):1798–1828.

[Bojarski et al. 2017] Bojarski, M.; Yeres, P.; Choromanska, A.; Choromanski, K.; Firner, B.; Jackel, L.; and Muller, U. 2017. Explaining how a deep neural network trained with end-to-end learning steers a car. *arXiv preprint arXiv:1704.07911*.

[Burrell 2016] Burrell, J. 2016. How the machine 'thinks': Understanding opacity in machine learning algorithms. *Big Data & Society* 3(1):2053951715622512.

[Chen et al. 2016] Chen, X.; Duan, Y.; Houthooft, R.; Schulman, J.; Sutskever, I.; and Abbeel, P. 2016. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, 2172–2180.

[Christiano et al. 2017] Christiano, P. F.; Leike, J.; Brown, T.; Martic, M.; Legg, S.; and Amodei, D. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, 4299–4307.

[Co-Reyes et al. 2018] Co-Reyes, J. D.; Liu, Y.; Gupta, A.; Eysenbach, B.; Abbeel, P.; and Levine, S. 2018. Self-consistent trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings. *arXiv preprint arXiv:1806.02813*.

[E. Todorov and Tassa. 2012] E. Todorov, T. E., and Tassa., Y. 2012. Mujoco: A physics engine for model-based control. *International Conference on Intelligent Robots and Systems (IROS)*.

[Eastwood and Williams 2018] Eastwood, C., and Williams, C. K. 2018. A framework for the quantitative evaluation of disentangled representations.

[G. Brockman and Zaremba. 2016] G. Brockman, V. Cheung, L. P. J. S. J. S. J. T., and Zaremba., W. 2016. Openai gym. *arXiv preprint arXiv:1606.01540*.

[Greydanus et al. 2017] Greydanus, S.; Koul, A.; Dodge, J.; and Fern, A. 2017. Visualizing and understanding atari agents. *arXiv preprint arXiv:1711.00138*.

[Günel ] Günel, M. Googlenet.

[Ha and Schmidhuber 2018] Ha, D., and Schmidhuber, J. 2018. World models. *arXiv preprint arXiv:1803.10122*.

[Higgins et al. 2016a] Higgins, I.; Matthey, L.; Glorot, X.; Pal, A.; Uria, B.; Blundell, C.; Mohamed, S.; and Lerchner, A. 2016a. Early visual concept learning with unsupervised deep learning. *arXiv preprint arXiv:1606.05579*.

[Higgins et al. 2016b] Higgins, I.; Matthey, L.; Pal, A.; Burgess, C.; Glorot, X.; Botvinick, M.; Mohamed, S.; and Lerchner, A. 2016b. beta-vae: Learning basic visual concepts with a constrained variational framework.

[Higgins et al. 2017] Higgins, I.; Pal, A.; Rusu, A. A.; Matthey, L.; Burgess, C. P.; Pritzel, A.; Botvinick, M.; Blundell, C.; and Lerchner, A. 2017. Darla: Improving zero-shot transfer in reinforcement learning. *arXiv preprint arXiv:1707.08475*.

[Kingma and Welling 2013] Kingma, D. P., and Welling, M. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

[Krizhevsky, Sutskever, and Hinton 2012] Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 1097–1105.

[LeCun, Bengio, and Hinton 2015] LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *nature* 521(7553):436.

[Lipson and Kurman 2016] Lipson, H., and Kurman, M. 2016. *Driverless: intelligent cars and the road ahead*. Mit Press.

[Matthey et al. 2017] Matthey, L.; Higgins, I.; Hassabis, D.; and Lerchner, A. 2017. dsprites: Disentanglement testing sprites dataset. https://github.com/deepmind/dsprites-dataset/.

[Mnih et al. 2015] Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518(7540):529.

[Mnih et al. 2016] Mnih, V.; Badia, A. P.; Mirza, M.; Graves, A.; Lillicrap, T.; Harley, T.; Silver, D.; and Kavukcuoglu, K. 2016. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, 1928–1937.

[Moore and Lu 2011] Moore, M. M., and Lu, B. 2011. Autonomous vehicles for personal transport: A technology assessment.

[Oh et al. 2015] Oh, J.; Guo, X.; Lee, H.; Lewis, R. L.; and Singh, S. 2015. Action-conditional video prediction using deep networks in atari games. In *Advances in neural information processing systems*, 2863–2871.

[Rahwan 2018] Rahwan, I. 2018. Society-in-the-loop: programming the algorithmic social contract. *Ethics and Information Technology* 20(1):5–14.

[Saxe et al. 2018] Saxe, A. M.; Bansal, Y.; Dapello, J.; Advani, M.; Kolchinsky, A.; Tracey, B. D.; and Cox, D. D. 2018. On the information bottleneck theory of deep learning.

[Schulman and Klimov 2017] Schulman, J., W. F. D. P. R. A., and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

[Shwartz-Ziv and Tishby 2017] Shwartz-Ziv, R., and Tishby, N. 2017. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*.

[Silver et al. 2017] Silver, D.; Schrittwieser, J.; Simonyan, K.; Antonoglou, I.; Huang, A.; Guez, A.; Hubert, T.; Baker, L.; Lai, M.; Bolton, A.; et al. 2017. Mastering the game of go without human knowledge. *Nature* 550(7676):354.

[Sohn, Lee, and Yan 2015] Sohn, K.; Lee, H.; and Yan, X. 2015. Learning structured output representation using deep conditional generative models. In *Advances in Neural Information Processing Systems*, 3483–3491.

[Stilgoe 2018] Stilgoe, J. 2018. Machine learning, social learning and the governance of self-driving cars. *Social studies of science* 48(1):25–56.

[Tang et al. 2017] Tang, H.; Houthooft, R.; Foote, D.; Stooke, A.; Chen, O. X.; Duan, Y.; Schulman, J.; DeTurck, F.; and Abbeel, P. 2017. # exploration: A study of count-based exploration for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, 2753–2762.

[Thomas et al. 2017] Thomas, V.; Pondard, J.; Bengio, E.; Sarfati, M.; Beaudoin, P.; Meurs, M.-J.; Pineau, J.; Precup, D.; and Bengio, Y. 2017. Independently controllable features. *arXiv preprint arXiv:1708.01289*.

[Vanderbilt 2012] Vanderbilt, T. 2012. Let the robot drive: The autonomous car of the future is here. *Wired Magazine, Conde NAST, www. wired. com* 1–34.

[Wu et al. 2017] Wu, Y.; Mansimov, E.; Grosse, R. B.; Liao, S.; and Ba, J. 2017. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *Advances in neural information processing systems*, 5279–5288.

[Wynne 1988] Wynne, B. 1988. Unruly technology: Practical rules, impractical discourses and public understanding. *Social studies of Science* 18(1):147–167.

[Zeiler and Fergus 2014] Zeiler, M. D., and Fergus, R. 2014. Visualizing and understanding convolutional networks. In *European conference on computer vision*, 818–833. Springer.

[Zhang and Zhu 2018] Zhang, Q.-s., and Zhu, S.-C. 2018. Visual interpretability for deep learning: a survey. *Frontiers of Information Technology & Electronic Engineering* 19(1):27–39.