
Complex Skill Acquisition through Simple Skill Imitation Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Humans often think of complex tasks as combinations of simpler subtasks in order
2 to learn those complex tasks more efficiently. For example, a backflip could
3 be considered a combination of four subskills: jumping, tucking knees, rolling
4 backwards, and thrusting arms downwards. Motivated by this line of reasoning,
5 we propose a new algorithm that trains neural network policies on simple, easy-
6 to-learn skills in order to cultivate latent spaces that accelerate imitation learning
7 of complex, hard-to-learn skills. We focus on the case in which the complex task
8 comprises a *concurrent* (and possibly *sequential*) combination of the simpler sub-
9 tasks, and therefore our algorithm can be seen as a novel approach to *concurrent*
10 *hierarchical imitation learning*. We evaluate our algorithm on difficult tasks in a
11 high-dimensional environment and see that it consistently outperforms a state-of-
12 the-art baseline in training speed and overall performance.

13 1 Introduction

14 Humans have the power to reason about complex tasks as combinations of simpler, interpretable
15 subtasks. There are many hierarchical reinforcement learning approaches designed to handle tasks
16 comprised of sequential subtasks [14, 8], but what if a task is made up of *concurrent* subtasks?
17 For example, someone who wants to learn to do a backflip may consider it to be combination of
18 sequential *and* concurrent subtasks: jumping, tucking knees, rolling backwards, and thrusting arms
19 downwards. Little focus has been given to designing algorithms that decompose complex tasks
20 into distinct concurrent subtasks. Even less effort has been put into finding decompositions that
21 are made up of independent yet interpretable concurrent subtasks, even though analogous approaches
22 have been effective on many challenging artificial intelligence problems [3, 2].

23 We propose a new generative model for encoding and generating arbitrarily complex trajectories. We
24 augment the VAE objective used in [15] in order to induce latent space structure that captures the
25 relationship between a behavior and the subskills that comprise this behavior in a disentangled and
26 interpretable way. We evaluate both the original and modified objectives on a moderately complex
27 imitation learning problem, in which agents are trained to perform a behavior after being trained on
28 subskills that qualitatively comprise that behavior.

29 2 Embedding and reconstructing trajectories

30 We use a conditional variational autoencoder (CVAE) [13, 7] to learn a semantically-meaningful
31 low-dimensional embedding space that can (1) help an agent learn new behaviors more quickly, (2)
32 be sampled from to generate behaviors, (3) and shed light on high-level factors of variation (e.g.
33 subskills) that comprise complex behaviors.

34 Illustrated by Figure 1, our CVAE has a bi-directional LSTM (BiLSTM) [6, 12] state-sequence
35 encoder $q_\phi(z|s_{1:T})$, an attention module [1, 17] that maps the BiLSTM output to values that

36 parametrize the distribution from which the latent (i.e. trajectory) embedding z is sampled, a
 37 conditional WaveNet [10] state decoder $\mathcal{P}_\psi(s_{t+1}|s_t, z)$, which serves as a *dynamics model*, and
 38 a multi-layer perceptron (MLP) action decoder $\pi_\theta(a_t|s_t, z)$, which serves as a *policy* whose outputs
 39 parametrize the normal distribution from which a_t is sampled. The bidirectional-LSTM captures
 40 sequential information over the states of the trajectories, and the conditional WaveNet allows for
 41 exact density modeling of the possibly multi-modal dynamics.

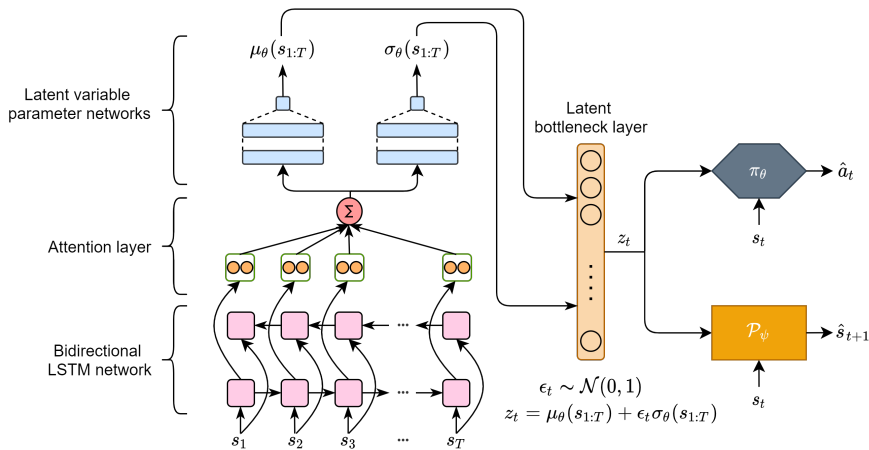


Figure 1: The conditional VAE we use to encode and generate trajectories.

42 We can train this CVAE by maximizing the following objective

$$\begin{aligned}
 \mathcal{L}(\theta, \phi, \psi; \tau^i) = \mathbb{E}_{z \sim q_\phi(z|s_{1:T_i}^i)} & \left[\sum_{t=1}^{T_i} \log \pi_\theta(a_t^i | s_t^i, z) + \log \mathcal{P}_\psi(s_{t+1}^i | s_t^i, z) \right] \\
 & + D_{KL}(q_\phi(z|s_{1:T_i}^i) \| p(z)). \quad (1)
 \end{aligned}$$

43 In Section 3 we will modify this objective in order to encourage the latent space to capture semanti-
 44 cally meaningful relationships between a behavior and the subskills that comprise this behavior.

45 3 Shaping the latent (i.e. trajectory embedding) space

46 Some skills can be seen as approximate combinations of certain subskills. Training a VAE to embed
 47 and reconstruct demonstrations of these skills and subskills using (1) would generally result in an
 48 embedding space with no clear relationship between skill and subskill embedding, especially if the
 49 dimensionality of the latent space is large or the number of demonstrated behaviors is small.

50 Motivated by semantically meaningful latent representations found in other work [9], we aim to
 51 induce a latent space structure so that a behavior embedding is the sum of its subskill embed-
 52 dings. Concretely, if z_A is a backflip embedding and z_a, z_b, z_c, z_d are embeddings corresponding
 53 to jumping, tucking knees, rolling backwards, and thrusting arms downwards, we want to have
 54 $z_A = z_a + z_b + z_c + z_d$. An example of such latent space restructuring is shown in Figure 2.

55 However, the VAE models probability distributions, so enforcing equality between one instance of
 56 a behavior and one instance of its subskills is insufficient. Instead, we want the random variables
 57 (RVs) representing the embeddings of the subskills to relate to the RV representing the embedding
 58 of the behavior comprised of those subskills. Another way to do this is to relate the subskill embed-
 59 ding RVs with the RV representing the trajectory generated by decoder networks \mathcal{P}_ψ and π_θ when
 60 conditioned on an embedding of the behavior.

61 Suppose τ_A is a behavior comprised of M subskills $\{\tau_{(1)}, \tau_{(2)}, \dots, \tau_{(M)}\}$. Let $\tilde{\tau}_A =$
 62 $(s_1, a_1, s_2, a_2, \dots, s_T, a_T)$ represent the trajectory generated from an embedding corresponding to
 63 τ_A . Define $V = z_1 + z_2 + \dots + z_M$, where $z_i \sim q_\phi(z|s_{(i)}, 1:T_{(i)})$. To train the encoder $q_\phi(z|s_{1:T})$,
 64 state decoder $\mathcal{P}_\psi(s_t|s_{t-1}, z)$, and action decoder $\pi_\theta(a_t|s_t, z)$ simultaneously, we aim to maximize

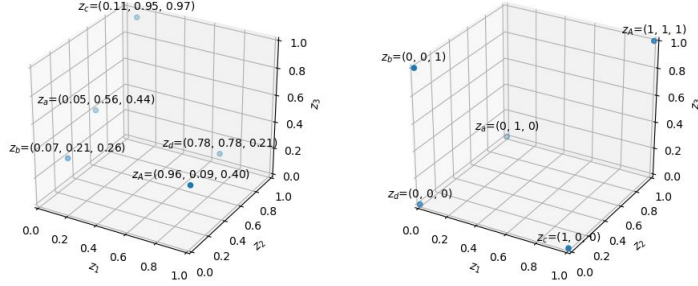


Figure 2: An example of latent space restructuring. *Left*: original latent space. *Right*: Hypothetical latent space induced by our approach (created intentionally for illustrative purposes).

65 the mutual information between V and $\tilde{\tau}$, which can be expressed as

$$\begin{aligned}
 I(V; \tilde{\tau}) &= H(V) - H(V|\tilde{\tau}) \\
 &= -\mathbb{E}_{V \sim p(V)} [\log p(V)] - \mathbb{E}_{V \sim p(V|\tilde{\tau})} [\log p(V|\tilde{\tau})]. \quad (2)
 \end{aligned}$$

66 If the latent variable prior distribution $p(z_i)$ is Gaussian, $H(V)$ is easy to compute, with an analytical
 67 solution under minor assumptions. We describe how to evaluate $H(V)$ in Appendix B.

68 3.1 Lower bounding mutual information through variational inference

69 However, we don't have access to the true posterior distribution $p(V|\tilde{\tau})$. We instead introduce a
 70 distribution $Q(V|\tilde{\tau})$ as a variational approximation to $p(V|\tilde{\tau})$ to get $L_I(\tilde{\tau}, Q)$, a variational lower
 71 bound of $I(V; \tilde{\tau})$

$$\begin{aligned}
 L_I(\tilde{\tau}, Q) &= \mathbb{E}_{V \sim p(V), \tau \sim \tilde{\tau}|V} [\log Q(V|\tau)] + H(V) \\
 &= \mathbb{E}_{\tau \sim \tilde{\tau}} [\mathbb{E}_{V \sim p(V|\tau)} [\log Q(V|\tau)]] + H(V) \\
 &\leq I(V; \tilde{\tau})
 \end{aligned}$$

72 in an approach similar to that of [3].

73 However, unlike in [3], $Q(V|\tilde{\tau})$ is *not* the same as $q(z|s_{1:T})$, the distribution approximated by
 74 the encoder network in our CVAE. Furthermore, even though embedding variables z_1, z_2, \dots, z_M
 75 are independent, they are *not* conditionally independent given $\tilde{\tau}$. Therefore, we *cannot* simply
 76 replace $Q(V|\tilde{\tau})$ with $\sum_{i=1}^M q(z_i|\tilde{\tau})$ and would instead need to again use variational inference to find
 77 $Q(V|\tilde{\tau})$, which would require training an additional VAE.

78 3.2 Lower bounding mutual information without variational inference

79 We derive a simpler lower bound to $I(V; \tilde{\tau})$ that allows us to circumvent the time and memory costs
 80 associated with training a VAE to model $Q(V|\tilde{\tau})$. We show the main result (3) here, and provide
 81 our derivation of this result in Appendix A.

$$I(V; \tilde{\tau}) \gtrsim -\mathbb{E}_{V \sim p(V)} [\log p(V)] + \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^M \log q_\phi(z_{n,i}|\tilde{\tau}) \quad (3)$$

82 By maximizing the lower bound in (3), we (approximately) maximize $I(V; \tilde{\tau})$
 83

84 3.3 Regularization with variational approximation

85 To encourage a semantically meaningful relationship between a behavior embedding and this behav-
 86 ior's subskill embeddings, we regularize the objective in (1) with $L_I(\tilde{\tau}, Q_\alpha)$ to get

$$\begin{aligned}
 \mathcal{L}(\theta, \phi, \psi; \tau^i) &= \mathbb{E}_{z \sim q_\phi(z|s_{1:T_i}^i)} \left[\sum_{t=1}^{T_i} \log \pi_\theta(a_t^i | s_t^i, z) + \log \mathcal{P}_\psi(s_{t+1}^i | s_t^i, z) \right] \\
 &\quad + D_{KL}(q_\phi(z|s_{1:T_i}^i) \| p(z)) + \lambda L_I(\tilde{\tau}, Q_\alpha), \quad (4)
 \end{aligned}$$

87 where $\lambda > 0$ is a hyperparameter that controls the trade-off between original objective and degree
 88 of shaping the latent space.

89 **3.4 Regularization without variational approximation**

90 If we want to avoid performing potentially expensive variational inference, we can use (6), the result
 91 we derived earlier in place of $L_I(\tilde{\tau}, Q)$,

$$\mathcal{L}(\theta, \phi, \psi; \tau^i) = \mathbb{E}_{z \sim q_\phi(z|s_{1:T_i}^i)} \left[\sum_{t=1}^{T_i} \log \pi_\theta(a_t^i | s_t^i, z) + \log \mathcal{P}_\psi(s_{t+1}^i | s_t^i, z) \right] \\
 + D_{KL}(q_\phi(z|s_{1:T_i}^i) \parallel p(z)) + \lambda \left(-\mathbb{E}_{V \sim p(V)} [\log p(V)] + \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^M \log q_\phi(z_{n,i} | \tilde{\tau}) \right). \quad (5)$$

92 As shown in Appendix B, the inner expectation in (5) can be evaluated analytically if the latent
 93 variables $\{z_i\}_{i=1}^M$ are independent and normally distributed—the standard case with VAEs.

94 **4 Experiments and results**

95 We evaluate our approach on a 197-dimensional state and 34-dimensional action space humanoid
 96 simulated in Bullet [4]. We use policies that were pre-trained by [11] to perform *kick*, *spin*, and
 97 *jump*, as subskills that qualitatively comprise the behavior *spin kick*. We also take a similar ap-
 98 proach for the behavior *backflip*. We train three sets of five VAEs on the subskills: one set optimizes
 99 for the original VAE objective (1), another set optimizes for the objective regularized by the varia-
 100 tional approximation (4), and the third set optimizes for the objective regularized without variational
 101 inference. To compare the proposed approach with the original, we evaluate the training process of
 102 each set of VAEs by considering the similarity between the generated trajectories and the pre-trained
 103 *spin kick* and *backflip* policy demonstrations. Results of the mean squared error (MSE) between the
 104 generated and demonstration states averaged over 5 different random seeds are shown in Figure 3.

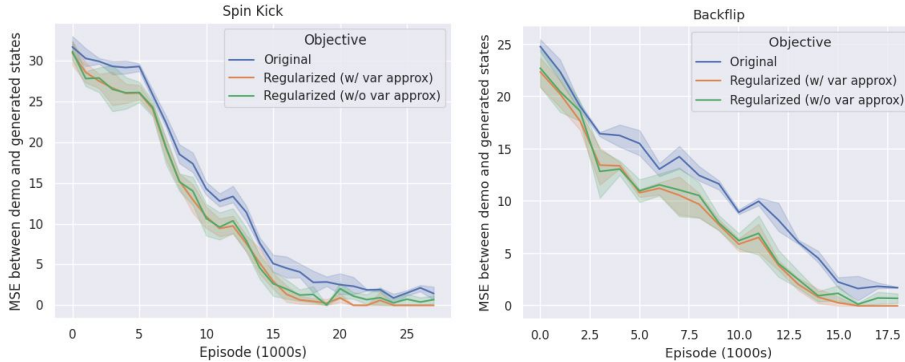


Figure 3: MSE (lower is better) between demonstration states and generated states on the Deep-
 Mimic *spin kick* and *backflip* tasks averaged over 5 different random seeds. **Regularized denotes
 our approaches (4), (5), and Original denotes the state-of-the-art baseline (1).**

105 We see that our proposed approaches attain better overall performance and train faster than the base-
 106 line algorithm. This suggests that we can bootstrap learning of difficult tasks by training agents on
 107 simpler, related subtasks while inclining their representations toward certain hierarchical structures.

108 **5 Discussion and future work**

109 We explored the idea of inducing certain latent structure through the maximization of mutual in-
 110 formation between generated behaviors and embeddings of the subskills that qualitatively comprise
 111 those behaviors, which, to the best of our knowledge, has not yet been investigated. Though our al-
 112 gorithm outperformed the state-of-the-art baseline, there is much room for future work. The CVAE
 113 could be replaced with a β -CVAE [5] to control disentanglement of z . The proposed approach could
 114 be evaluated on behaviors and subskills that more strictly adhere to concurrent relationship desired.
 115 A larger number of behaviors, such as those put forth by [16], could be trained at once, both to
 116 constrain the latent space and to enrich the pool of subskills from which to train on and inspect
 117 relationships between. The non-variational mutual information approximation could be compared
 118 to the variational one in order to quantify accuracy. Interpolations within the convex hull of subskill
 119 embeddings could be used to fine-tune known behaviors or generate completely new behaviors.

References

- 120
- 121 [1] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align
122 and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- 123 [2] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner.
124 Understanding disentangling in β -vae. *arXiv preprint arXiv:1804.03599*, 2018.
- 125 [3] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Inter-
126 pretable representation learning by information maximizing generative adversarial nets. In
127 *Advances in neural information processing systems*, pages 2172–2180, 2016.
- 128 [4] E. Coumans. *Bullet*, 2015.
- 129 [5] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Ler-
130 chner. beta-vae: Learning basic visual concepts with a constrained variational framework. *Iclr*,
131 2(5):6, 2017.
- 132 [6] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–
133 1780, 1997.
- 134 [7] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint*
135 *arXiv:1312.6114*, 2013.
- 136 [8] G. Konidaris and A. G. Barto. Building portable options: Skill transfer in reinforcement learn-
137 ing. In *IJCAI*, volume 7, 2007.
- 138 [9] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of
139 words and phrases and their compositionality. In *Advances in neural information processing*
140 *systems*, pages 3111–3119, 2013.
- 141 [10] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner,
142 A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint*
143 *arXiv:1609.03499*, 2016.
- 144 [11] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne. Deepmimic: Example-guided deep re-
145 inforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*,
146 37(4):1–14, 2018.
- 147 [12] M. Schuster and K. K. Paliwal. Bidirectional recurrent neural networks. *IEEE transactions on*
148 *Signal Processing*, 45(11):2673–2681, 1997.
- 149 [13] K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep conditional
150 generative models. In *Advances in neural information processing systems*, pages 3483–3491,
151 2015.
- 152 [14] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for tem-
153 poral abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- 154 [15] Z. Wang, J. S. Merel, S. E. Reed, N. de Freitas, G. Wayne, and N. Heess. Robust imitation of
155 diverse behaviors. In *Advances in Neural Information Processing Systems*, pages 5320–5329,
156 2017.
- 157 [16] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A
158 benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on*
159 *Robot Learning*, pages 1094–1100, 2020.
- 160 [17] P. Zhou, W. Shi, J. Tian, Z. Qi, B. Li, H. Hao, and B. Xu. Attention-based bidirectional long
161 short-term memory networks for relation classification. In *Proceedings of the 54th Annual*
162 *Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages
163 207–212, 2016.

164 **A Derivation of mutual information lower bound without variational**
 165 **approximation**

166 For clarity in the following derivation, let $V_p = \sum_{i=p}^M z_i$. Then we have

$$\begin{aligned}
 H(V|\tilde{\tau}) &= H(V_1|\tilde{\tau}) \\
 &= H(z_1 + z_2 + \cdots + z_M|\tilde{\tau}) \\
 &= H(z_1|\tilde{\tau}) + H(z_1 + z_2 + \cdots + z_M|z_1, \tilde{\tau}) - H(z_1|z_1 + z_2 + \cdots + z_M, \tilde{\tau}) \\
 &= H(z_1|\tilde{\tau}) + H(z_2 + z_3 + \cdots + z_M|z_1, \tilde{\tau}) - H(z_1|z_1 + z_2 + \cdots + z_M, \tilde{\tau}) \\
 &\leq H(z_1|\tilde{\tau}) + H(z_2 + z_3 + \cdots + z_M|\tilde{\tau}) - H(z_1|z_1 + z_2 + \cdots + z_M, \tilde{\tau}) \\
 &= H(z_1|\tilde{\tau}) + H(V_2|\tilde{\tau}) - H(z_1|V_1, \tilde{\tau})
 \end{aligned}$$

167 By rolling out $H(V_p|\tilde{\tau})$ recursively for $p = 1, 2, 3, \dots, M - 1$, we get

$$\begin{aligned}
 H(V|\tilde{\tau}) &\leq \sum_{i=1}^M [H(z_i|\tilde{\tau}) - H(z_i|V_i, \tilde{\tau})] \\
 &\leq \sum_{i=1}^M H(z_i|\tilde{\tau}) \\
 &= \sum_{i=1}^M -\mathbb{E}_{z_i \sim p(z_i|\tilde{\tau})} [\log p(z_i|\tilde{\tau})] \\
 &\approx \sum_{i=1}^M -\mathbb{E}_{z_i \sim q_\phi(z_i|\tilde{\tau})} [\log q_\phi(z_i|\tilde{\tau})]
 \end{aligned}$$

168 if $p(z|\tilde{\tau}) \approx q_\phi(z|\tilde{\tau})$. Plugging this result into (2) allows us to lower bound $I(V; \tilde{\tau})$ as follows,

$$\begin{aligned}
 I(V; \tilde{\tau}) &\geq -\mathbb{E}_{V \sim p(V)} [\log p(V)] + \sum_{i=1}^M \mathbb{E}_{z_i \sim p(z_i|\tilde{\tau})} [\log p(z_i|\tilde{\tau})] \\
 &\approx -\mathbb{E}_{V \sim p(V)} [\log p(V)] + \sum_{i=1}^M \mathbb{E}_{z_i \sim q_\phi(z_i|\tilde{\tau})} [\log q_\phi(z_i|\tilde{\tau})],
 \end{aligned}$$

169 and we can obtain an unbiased estimate of the second term by sampling $z_i \sim q_\phi(z_i|\tilde{\tau})$ to get

$$I(V; \tilde{\tau}) \gtrsim -\mathbb{E}_{V \sim p(V)} [\log p(V)] + \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^M \log q_\phi(z_{n,i}|\tilde{\tau}), \quad (6)$$

170 where $x \gtrsim y$ denotes that x is approximately greater than or equal to y .

171 **B Evaluating entropy of sum of subskill embeddings**

172 Computing the entropy for an arbitrary distribution may be difficult, but by setting X to be a Gaussian RV—the standard choice for VAE encoders—the entropy $H(X)$ has the simple, closed-form
173 expression
174

$$H(X) = \frac{1}{2}(1 + \ln(2\pi\sigma_X^2)),$$

175 where σ_X is the standard deviation of X . We choose $q_\phi(z|s_{1:T})$ to parametrize a Gaussian dis-
176 tribution and assume that state sequences from different subskills are sufficiently unrelated so that
177 they can be considered statistically independent. This is generally a safe assumption because even
178 minor differences in subskills will tend to place trajectories corresponding to different skills in very
179 different locations within the trajectory space. It follows that V is the sum of Gaussian RVs and has
180 the simple form

$$V \sim \mathcal{N}(\mu_{z_a} + \mu_{z_b} + \dots + \mu_{z_M}, \sigma_{z_a}^2 + \sigma_{z_b}^2 + \dots + \sigma_{z_M}^2),$$

181 and the entropy of V is

$$H(V) = \frac{1}{2}(1 + \ln(2\pi(\sigma_{z_a}^2 + \sigma_{z_b}^2 + \dots + \sigma_{z_M}^2))). \quad (7)$$