Visual Adversarial Imitation Learning using Variational Models

Anonymous Author(s) Affiliation Address email

Abstract

Reward function specification, which requires considerable human effort and itera-1 2 tion, remains a major impediment for learning behaviors through deep reinforce-3 ment learning. In contrast, providing visual demonstrations of desired behaviors presents an easier and more natural way to teach agents. We consider a setting 4 where an agent is provided a fixed dataset of visual demonstrations illustrating how 5 to perform a task, and must learn to solve the task using the provided demonstra-6 tions and unsupervised environment interactions. This setting presents a number of 7 8 challenges including representation learning for visual observations, sample complexity due to high dimensional spaces, and learning instability due to the lack of a 9 fixed reward or learning signal. Towards addressing these challenges, we develop a 10 variational model-based adversarial imitation learning (V-MAIL) algorithm. The 11 model-based approach provides a strong signal for representation learning, enables 12 sample efficiency, and improves the stability of adversarial training by enabling on-13 policy learning. Through experiments involving several vision-based locomotion 14 and manipulation tasks, we find that V-MAIL learns successful visuomotor policies 15 in a sample-efficient manner, has better stability compared to prior work, and also 16 achieves higher asymptotic performance. We further find that by transferring the 17 learned models, V-MAIL can learn new tasks from visual demonstrations without 18 any additional environment interactions. All results including videos can be found 19 online at https://sites.google.com/view/variational-mail. 20

21 **1 Introduction**

The ability of reinforcement learning (RL) agents to autonomously learn by interacting with the 22 environment presents a promising approach for learning diverse skills. However, reward specification 23 has remained a major challenge in the deployment of RL in practical settings [2, 9, 37]. The ability 24 to imitate humans or other expert trajectories allows us to avoid the reward specification problem, 25 while also circumventing challenges related to exploration in RL. Visual demonstrations are also 26 a more natural way to teach robots various tasks and skills in real-world applications. However, 27 this setting is also fraught with a number of technical challenges including representation learning 28 for visual observations, sample complexity due to the high dimensional observation spaces, and 29 learning instability [35, 25, 32] due to lack of a stationary learning signal. We aim to overcome these 30 challenges and to develop an algorithm that can learn from limited demonstration data as well as 31 scale to high-dimensional observation and action spaces often encountered in robotics applications. 32

Behaviour cloning (BC) is a classic algorithm to imitate expert demonstrations [34], which uses supervised learning to greedily match the expert behaviour at demonstrated expert states. Due to environment stochasticity, covariate shift, and policy approximation error, the agent may drift away from the expert state distribution and ultimately fail to mimic the demonstrator [40]. While a wide initial state distribution [41] or the ability to interactively query the expert policy [40] can



Figure 1: Left: the variational dynamics model, which enables joint representation learning from visual inputs and a latent space dynamics model, and the discriminator which is trained to distinguish latent states of expert demonstrations from that of policy rollouts. Dashed lines represent inference and solid lines represent the generative model. **Right**: the policy training, which uses the discriminator as the reward function, so that the policy induces a latent state visitation distribution that is indistinguishable from that of the expert. The learned policy network is composed with the image encoder from the variational model to recover a visuomotor policy.

circumvent these difficulties, such conditions require additional supervision and are difficult to meet 38 in practical applications. An alternate line of work based on inverse RL [13, 14] and adversarial 39 imitation learning [22, 12] aims to not only match actions at demonstrated states, but also the long 40 term visitation distribution [16]. These approaches explicitly train a GAN-based classifier [17] to 41 distinguish the visitation distribution of the agent from the expert, and use it as a reward signal 42 for training the agent with RL. While these methods have achieved substantial improvement over 43 behaviour cloning without additional expert supervision, they are difficult to deploy in realistic 44 scenarios, primarily due to three reasons: (1) the objective requires on-policy data collection leading 45 to high sample complexity; (2) the non-stationarity reward function changes as the RL agent learns; 46 and (3) high-dimensional observation spaces require representation learning and exacerbate the 47 optimization challenges. 48

Our main contribution in this work is the development of a new algorithm, variational model-based 49 adversarial imitation learning (V-MAIL), which aims to overcome each of the aforementioned chal-50 lenges within a single framework. As illustrated in Figure 1, V-MAIL trains a variational latent-space 51 dynamics model and a discriminator that provides a learning reward signal by distinguishing latent 52 rollouts of the agent from the expert. The key insight of our approach is that variational models can 53 address these challenges simultaneously by (a) making it possible to collect on-policy roll-outs inside 54 the model without environment interaction, leading to an efficient and stable optimization process 55 and (b) providing a rich auxiliary objective for efficiently learning compact state representations 56 and which regularizes the discriminator. Furthermore, the variational model also allows V-MAIL to 57 perform zero-shot transfer to new imitation learning tasks. By generating on-policy rollouts within 58 the model, and training the discriminator using these rollouts along with demonstrations of a new 59 task, V-MAIL can learn policies for new tasks without any additional environment interactions. 60

Through experiments on a collection of vision-based locomotion and manipulation tasks, we find that 61 V-MAIL can learn successful visuomotor control policies through imitation learning. In particular, 62 V-MAIL exhibits stable and near-monotonic learning, is highly sample efficient, and asymptotically 63 matches the expert level performance on most tasks. In contrast, prior algorithms exhibit unstable 64 learning and poor asymptotic performance, often achieving less that 20% of expert level performance. 65 We further show the ability to transfer our models to novel task and acquire qualitatively new behaviors 66 using only a few demonstrations and no additional environment interactions. To our knowledge this is 67 the first approach to use variational model-based training for zero-shot or few-shot imitation learning. 68

69 2 Related Work

70 Here, we review the relevant literature on imitation learning and image-based RL.

Imitation Learning. Recent model-free imitation learning can be categorized as either adversarial
 or non-adversarial. Adversarial methods inspired by GANs [17] train an explicit classifier between

expert and policy behaviour and optimize the agent in a two-player minimax game. GAIL [22]
and AIRL [14] are two such algorithms; however they often have poor sample efficiency due to the
requirement of on-policy rollouts in the environment. To address sample efficiency issues, off-policy
variants such as DAC [27] and SAM [5] have been developed, however they suffer from an objective
mismatch when using off-policy data [28], often resulting in learning instability [6].

An alternate line of research attempts to forego adversarial training: SQIL [39] frames the problem 78 as regularized behaviour cloning and trains an off-policy algorithm with rewards of 1 for expert 79 trajectories and 0 for policy ones. RCE [10] uses a very similar approach, but derives it as maximizing 80 probability of task success, which they show is equivalent to minimizing the Hellinger distance 81 between the policy occupation distribution and a particular target distribution. ValueDICE [28] uses 82 the same key result for iterative distribution matching as RCE in conjunction with the Donsker-83 Varadhan representation to obtain an off-policy distribution matching algorithm. In Swamy et al. 84 [42] the authors derive distribution matching as a bound on policy under-performance, similar to our 85 analysis in Section 4.1 and propose a practical non-adversarial algorithm AdVIL, however in reported 86 experiments it does not outperform behaviour cloning. A few papers have considered model-based 87 imitation learning as well: Baram et al. [3] is an adversarial algorithm conceptually similar to our 88 approach, but only focuses on low-dimensional state-based tasks and train the discriminator using 89 off-policy replay buffer, which does not allow it to generalize to new tasks. Related to our method is 90 Finn et al. [13] which uses a similar reward learning in combination with a locally linear dynamics 91 model, which leads to trajectory centric algorithms and the inability to transfer the model to new 92 tasks. Das et al. [8] considers a similar setting for inverse RL using a simplified parameterization 93 94 of the cost function. In this work we develop end-to-end model for adversarial imitation learning in high-dimensional POMDPs and generalization to novel tasks without hand-designed features. 95

Reinforcement Learning From Images with Variational Models. Reinforcement learning from images is an inherently difficult task, since the agent needs to learn meaningful visual representations to support policy learning. A recent line of research [15, 19, 30, 20, 36] train a variational model of the image-based environment as an auxiliary task, either for representation learning only [15, 30] or for additionally generating on-policy data by rolling out the model [20]. Our method builds upon these ideas, but unlike these prior works, considers the problem of learning from visual demonstrations without access to rewards.

103 3 Preliminaries

We consider the problem setting of learning in partially observed Markov decision processes (POMDPs), which can be described with the tuple: $\mathcal{M} = (S, \mathcal{A}, \mathcal{X}, \mathcal{R}, \mathcal{T}, \mathcal{U}, \gamma)$, where $s \in S$ is the state space, $a \in \mathcal{A}$ is the action space, $x \in \mathcal{X}$ is the observation space and $r = \mathcal{R}(s, a)$ is a reward function. The state evolution is Markovian and governed by the dynamics as $s' \sim \mathcal{T}(\cdot|s, a)$. Finally, the observations are generated through the observation model $x \sim \mathcal{U}(\cdot|s)$. The widely studied Markov decision process (MDP) is a special case of this 7-tuple where the underlying state is directly observed in the observation model.

In this work, we study imitation learning in unknown POMDPs. Thus, we do not have access to the underlying dynamics, the true state representation of the POMDP, or the reward function. In place of the rewards, the agent is provided with a fixed set of expert demonstrations collected by executing an expert policy π^E , which we assume is optimal under the unknown reward function. The agent can interact with the environment and must learn a policy $\pi(a_t|x_{\leq t})$ that mimics the expert.

116 3.1 Imitation learning as divergence minimization

In line with prior work, we interpret imitation learning as a divergence minimization problem [22, 16, 24]. For simplicity of exposition, we consider the MDP case in this section, and discuss POMDP extensions in Section 4.2. Let $\rho_{\mathcal{M}}^{\pi}(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} P(s_{t} = s, a_{t} = a)$ be the discounted state-action visitation distribution of a policy π in MDP \mathcal{M} . Then, a divergence minimization objective for imitation learning corresponds to

$$\min_{\tau} \mathbb{D}(\rho_{\mathcal{M}}^{\pi}, \rho_{\mathcal{M}}^{E}), \tag{1}$$

where $\rho_{\mathcal{M}}^E$ is the discounted visitation distribution of the expert policy π^E , and \mathbb{D} is a divergence measure between probability distributions such as KL-divergence, Jensen-Shannon divergence, or a generic f-divergence. To see why this is a reasonable objective, let $J(\pi, \mathcal{M})$ denote the expected value of a policy π in \mathcal{M} . Inverse RL [46, 22, 12] interprets the expert as the optimal policy under some unknown reward function. With respect to this unknown reward function, the sub-optimality of any policy π can be bounded as:

$$\left|J(\pi^{E}, \mathcal{M}) - J(\pi, \mathcal{M})\right| \leq \frac{R_{\max}}{1 - \gamma} \mathbb{D}_{TV}(\rho_{\mathcal{M}}^{\pi}, \rho_{\mathcal{M}}^{E}),$$

since the policy performance is $J(\pi, \mathcal{M}) = \mathbb{E}_{(s,a) \sim \rho_{\mathcal{M}}^{\pi}} [r(s, a)]$. We use \mathbb{D}_{TV} to denote total variation distance. Since various divergence measures are related to the total variation distance, optimizing the divergence between visitation distributions in state space amounts to optimizing a bound on the policy sub-optimality.

132 3.2 Generative Adversarial Imitation Learning (GAIL)

With the divergence minimization viewpoint, any standard generative modeling technique including 133 density estimation, VAEs, GANs etc. can in principle be used to minimize Eq. 1. However, in 134 practice, use of certain generative modeling techniques can be difficult. A standard density estimation 135 technique would involve directly parameterizing $\rho_{\mathcal{M}}^{\pi}$, say through auto-regressive flows, and learning 136 the density model. However, a policy that induces the learned visitation distribution in \mathcal{M} is not 137 guaranteed to exist and may prove hard to recover. Similar challenges prevent the direct application of 138 a VAE based generative model as well. In contrast, GANs allow for a policy based parameterization, 139 since it only requires the ability to sample from the generative model and does not require the 140 likelihood. This approach was followed in GAIL, leading to the optimization 141

$$\max_{\pi} \min_{D_{\psi}} \mathbb{E}_{(\boldsymbol{s},\boldsymbol{a}) \sim \rho_{\mathcal{M}}^{E}} \left[-\log D_{\psi}(\boldsymbol{s},\boldsymbol{a}) \right] + \mathbb{E}_{(\boldsymbol{s},\boldsymbol{a}) \sim \rho_{\mathcal{M}}^{\pi}} \left[-\log \left(1 - D_{\psi}(\boldsymbol{s},\boldsymbol{a}) \right) \right], \quad (2)$$

where D_{ψ} is a discriminative classifier used to distinguish between samples from the expert distri-142 bution and the policy generated distribution. Results from Goodfellow et al. [17] and Ho & Ermon 143 [22] suggest that the learning objective in Eq. 2 corresponds to the divergence minimization objective 144 in Eq. 1 with Jensen-Shannon divergence. In order to estimate the second expectation in Eq. 2 we 145 require on-policy samples from π , which is data-inefficient. Adversarial off-policy algorithms, such 146 as [27, 5] replace the expectation under the policy distribution with expectation under the current 147 replay buffer distribution, which allows for off-policy training, but no longer guarantee that the policy 148 marginal distribution will match the expert. 149

150 4 Variational Model-Based Adversarial Imitation Learning

Generative modeling in the context of imitation learning poses unique challenges. Improving the 151 generative distribution (policy in our case) requires samples from $\rho_{\mathcal{M}}^{\pi}$, which requires rolling out 152 π in the environment. Furthermore, the complex optimization landscape of a saddle point problem 153 requires many iterations of learning, each of which requires on-policy rollouts. This is unlike typical 154 generative modeling applications where generating samples from the generator is cheap and does 155 not require any environment interactions. To overcome this challenge, we present a model-based 156 imitation learning algorithm. For conceptual clarity and ease of exposition, we will first present 157 our conceptual algorithm in the MDP setting in Section 4.1. Subsequently, we will extend this 158 algorithm to the POMDP case in Section 4.2. Finally, we present a practical version of our algorithm 159 in Section 4.3. 160

161 4.1 Model-Based Adversarial Imitation Learning

Model-based algorithms for RL and IL involve learning an approximate dynamics model $\hat{\mathcal{T}}$ using environment interactions. The learned dynamics model can be used to construct an approximate MDP $\widehat{\mathcal{M}}$. In our context of imitation learning, learning a dynamics model allows us to generate samples from $\widehat{\mathcal{M}}$ as a surrogate for samples from \mathcal{M} , leading to the objective:

$$\min_{\pi} \mathbb{D}(\rho_{\widehat{\mathcal{M}}}^{\pi}, \rho_{\mathcal{M}}^{E}), \tag{3}$$

which can serve as a good proxy to Eq. 1 as long as the model approximation is accurate. In particular, with an α -approximate dynamics model given by $\mathbb{D}_{TV}(\widehat{\mathcal{T}}(s, a), \mathcal{T}(s, a)) \leq \alpha \ \forall (s, a)$, we can

bound the policy suboptimality with respect to the expert as:

$$\left| J(\pi^{E}, \mathcal{M}) - J(\pi, \mathcal{M}) \right| \leq \frac{R_{\max}}{1 - \gamma} \mathbb{D}_{TV}(\rho_{\widehat{\mathcal{M}}}^{\pi}, \rho_{\mathcal{M}}^{E}) + \frac{\alpha \cdot R_{\max}}{(1 - \gamma)^{2}}.$$
(4)

¹⁶⁹ Thus, the divergence minimization in Eq. 3 serves as an approximate bound on the sub-optimality

with a bias that is proportional to the model error. Thus, we ultimately propose to solve the following

171 saddle point optimization problem:

$$\max_{\pi} \min_{D_{\psi}} \mathbb{E}_{(\boldsymbol{s},\boldsymbol{a}) \sim \rho_{\mathcal{M}}^{E}} \left[-\log D_{\psi}(\boldsymbol{s},\boldsymbol{a}) \right] + \mathbb{E}_{(\boldsymbol{s},\boldsymbol{a}) \sim \rho_{\widehat{\mathcal{M}}}^{\pi}} \left[-\log \left(1 - D_{\psi}(\boldsymbol{s},\boldsymbol{a}) \right) \right], \quad (5)$$

which requires generating on-policy samples only from the learned model $\hat{\mathcal{M}}$. We can interleave policy learning according to Eq. 5 with performing policy rollouts in the real environment to iteratively improve the model. Provided the policy is updated sufficiently slowly, Rajeswaran et al. [38] show that such interleaved policy and model learning corresponds to a stable and convergent algorithm, while being highly sample efficient.

177 4.2 Extension to POMDPs

In POMDPs, the underlying state is not directly observed, and thus cannot be directly used by the policy. In this case, we typically use the notion of *belief state*, which is defined to be the filtering distribution $P(s_t|h_t)$, where we denote history with $h_t := (x_{\le t}, a_{< t})$. By using the historical information, the belief state provides more information about the current state, and can enable the learning of better policies. However, learning and maintaining an explicit distribution over states can be difficult. Thus, we consider learning a latent representation of the history $z_t = q(h_t)$, so that $P(s_t|h_t) \approx P(s_t|z_t)$.

To develop an algorithm for the POMDP setting, we first make the key observation that imitation learning in POMDPs can be reduced to divergence minimization in the latent belief state representation. To formalize this, we introduce the following theorem. A formal version of the theorem and proof are provided in the appendix.

Theorem 1. (Divergence bound in latent space; Informal) Consider a POMDP \mathcal{M} , and let z_t be a latent space representation of the history and belief state such that $P(s_t|x_{\leq t}, a_{< t}) = P(s_t|z_t)$. Let D_f be a generic f-divergence. Then the following inequalities hold:

$$D_f(\rho_{\mathcal{M}}^{\pi}(\boldsymbol{x},\boldsymbol{a})||\rho_{\mathcal{M}}^{E}(\boldsymbol{x},\boldsymbol{a})) \leq D_f(\rho_{\mathcal{M}}^{\pi}(\boldsymbol{s},\boldsymbol{a})||\rho_{\mathcal{M}}^{E}(\boldsymbol{s},\boldsymbol{a})) \leq D_f(\rho_{\mathcal{M}}^{\pi}(\boldsymbol{z},\boldsymbol{a})||\rho_{\mathcal{M}}^{E}(\boldsymbol{z},\boldsymbol{a}))$$

Theorem 1 suggests that the divergence of visitation distributions in the latent space represents 189 an upper bound of the divergence in the state and observation spaces. This is particularly useful, 190 since we do not have access to the ground-truth states of the POMDP and matching the expert 191 marginal distribution in the high-dimensional observation space (such as images) could be difficult. 192 Furthermore, based on the results in Section 3.1, minimizing the state divergence results in minimizing 193 a bound on policy sub-optimality as well. These results provide a direct way to extend the results 194 from Section 4.1 to the POMDP setting. If we can learn an encoder $z_t = q(x_{\le t}, a_{< t})$ that captures 195 sufficient statistics of the history, and a latent state space dynamics model $z_{t+1} \sim \hat{T}(\cdot | z_t, a_t)$, then we can learn the policy by extending Eq. 5 to the induced MDP in the latent space as: 196 197

$$\max_{\pi} \min_{D_{\psi}} \mathbb{E}_{(\boldsymbol{z},\boldsymbol{a}) \sim \rho_{\mathcal{M}}^{E}(\boldsymbol{z},\boldsymbol{a})} \left[-\log D_{\psi}(\boldsymbol{z},\boldsymbol{a}) \right] + \mathbb{E}_{(\boldsymbol{z},\boldsymbol{a}) \sim \rho_{\mathcal{M}}^{\pi}(\boldsymbol{z},\boldsymbol{a})} \left[-\log \left(1 - D_{\psi}(\boldsymbol{z},\boldsymbol{a}) \right) \right].$$
(6)

¹⁹⁸ Once learned, the policy can be composed with the encoder for deployment in the POMDP.

199 4.3 Practical Algorithm with Variational Models

The divergence bound of Theorem 1 allows us to develop a practical algorithm if we can learn a good belief state representation. Towards that end we turn to the theory of deep Bayesian filters [23] and begin with the likelihood:

$$\log P(\boldsymbol{x}_{1:T} | \boldsymbol{a}_{1:T}) = \log \int \prod_{t=1}^{T} U(\boldsymbol{x}_t | \boldsymbol{s}_t) \mathcal{T}(\boldsymbol{s}_t | \boldsymbol{a}_{t-1}, \boldsymbol{s}_{t-1}) d\boldsymbol{s}_{1:T}$$

We can introduce the belief distribution $q(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}, \mathbf{a}_{1:T-1}) = \prod_{t=1}^{T} q(\mathbf{z}_t|\mathbf{x}_t, \mathbf{z}_{t-1}, \mathbf{a}_{t-1})$, which considers only model classes that satisfy the the sufficient statistics requirement. Using the introduced belief distribution as the variational distribution, we derive the evidence lower bound (ELBO) [4, 26]:

Algorithm 1 V-	MAIL: `	Variational	Model-	Based A	dversarial	Imitation 1	Learning
----------------	---------	-------------	--------	---------	------------	-------------	----------

1: **Require**: Expert demos \mathcal{B}_E , environment buffer \mathcal{B}_{π} .

- 2: Randomly initialize variational model $\{q_{\theta}, \hat{\mathcal{T}}_{\theta}\}$, policy π_{ψ} and discriminator D_{ψ}
- 3: for number of iterations do
- 4: // Environment Data Collection
- 5: for timestep t = 1 : T do
- Estimate latent state from the belief distribution $z_t \sim q_{\theta}(\cdot | x_t, z_{t-1}, a_{t-1})$ 6:
- 7: Sample action $\boldsymbol{a}_t \sim \pi_{\psi}(\boldsymbol{a}_t | \boldsymbol{z}_t)$
- 8: Step environment and get observation x_{t+1}
- 9: Add data $\{x_{1:T}, a_{1:T-1}\}$ to policy replay buffer \mathcal{B}_{π}
- 10: for number of training iterations do
- 11: // Dynamics Learning
- 12: Sample a batch of trajectories $\{x_{1:T}, a_{1:T-1}\}$ from the joint buffer $\mathcal{B}_E \cup \mathcal{B}_{\pi}$
- 13: Optimize the variational model $\{q_{\theta}, \mathcal{T}_{\theta}\}$ using Equation 7
- 14: // Adversarial Policy Learning
- Sample trajectories from expert buffer $\{x_{1:T}^E, a_{1:T-1}^E\} \sim \mathcal{B}_E$ 15:
- Infer expert latent states $\boldsymbol{z}_{1:T}^E \sim q_{\theta}(\cdot | \boldsymbol{x}_{1:T}^E, \boldsymbol{a}_{1:T-1}^E)$ using the belief model q_{θ} 16:
- Generate latent rollouts $\boldsymbol{z}_{1:H}^{\pi_{\psi}}$ using the policy π_{ψ} from the forward model $\widehat{\mathcal{T}}_{\theta}$ Update the discriminator D_{ψ} with data $\boldsymbol{z}_{1:T}^{E}, \boldsymbol{z}_{1:H}^{\pi_{\psi}}$ using Equation 6 17:
- 18:
- Update the policy π_{ψ} to improve the value function in Equation 8 19:

$$\log P(\boldsymbol{x}_{1:T}|\boldsymbol{a}_{1:T}) \geq \mathbb{E}_{q(\boldsymbol{z}_{1:T}|\boldsymbol{x}_{1:T},\boldsymbol{a}_{1:T-1})} \left[\log \prod_{t=1}^{T} \mathcal{U}(\boldsymbol{x}_{t}|\boldsymbol{z}_{t}) \frac{\mathcal{T}(\boldsymbol{z}_{t}|\boldsymbol{a}_{t-1},\boldsymbol{z}_{t-1})}{q(\boldsymbol{z}_{t}|\boldsymbol{x}_{t},\boldsymbol{z}_{t-1},\boldsymbol{a}_{t-1})} \right]$$

To estimate the expectation, we can use sequential sampling from the belief distribution $z_t \sim$ 206 $q(|\mathbf{z}_t, \mathbf{z}_{t-1}, \mathbf{a}_{t-1}), t = 1 : T$ and the reparameterization trick [26]. This ultimately leads to the 207 empirical variational model training objective: 208

$$\max_{\theta} \widehat{\mathbb{E}}_{q_{\theta}} \left[\sum_{t=1}^{T} \underbrace{\log \widehat{\mathcal{U}}_{\theta}(\boldsymbol{x}_{t} | \boldsymbol{z}_{t})}_{\text{reconstruction}} - \underbrace{\mathbb{D}_{KL}(q_{\theta}(\boldsymbol{z}_{t} | \boldsymbol{x}_{t}, \boldsymbol{z}_{t-1}, \boldsymbol{a}_{t-1}) || \widehat{\mathcal{T}}_{\theta}(\boldsymbol{z}_{t} | \boldsymbol{z}_{t-1}, \boldsymbol{a}_{t-1}))}_{\text{forward model}} \right].$$
(7)

That is, we jointly train a belief representation q_{θ} and a Markovian dynamics model $\hat{\mathcal{T}}$, which allows 209 us to optimize Eq. 5 in our learned belief space. A number of recent works have considered similar 210 models [44, 45, 30, 15, 19, 20]. We base our architectural choice on the recurrent state space model 211 [19, 20], as it has shown strong performance in RL tasks from images. In principle, any on-policy RL 212 algorithm can be used to train the policy using Eq. 6. In our setup, the RL objective is a differentiable 213 function of the policy, model, and discriminator parameters. Based on this, we setup a K step value 214 expansion objective [11, 7] given below, and use it for policy learning. 215

$$V_{\theta,\psi}^{K}(\boldsymbol{z}_{t}) = \mathbb{E}_{\pi_{\psi},\widehat{T}_{\theta}} \left[\sum_{\tau=t}^{t+K-1} \gamma^{\tau-t} \log D_{\psi}(\boldsymbol{z}_{\tau}^{\pi_{\psi}}, \boldsymbol{a}_{\tau}^{\pi_{\psi}}) + \gamma^{K} V_{\psi}(\boldsymbol{z}_{t+K}^{\pi_{\psi}}) \right]$$
(8)

Finally, we train the discriminator D_{ψ} using Eq. 5 with on-policy rollots from the model $\hat{\mathcal{T}}$. Our full 216 approach is outlined in Algorithm 1. 217

Zero-Shot Transfer to New Imitation Tasks 4.4 218

Our model-based approach is well suited to the problem of zero-shot transfer to new imitation learning 219 tasks, i.e. transferring to a new task using a modest number of demonstrations and no additional 220 samples collected in the environment. In particular, we assume a set of source tasks $\{\mathcal{T}^i\}$, each 221 with a buffer of expert demonstrations \mathcal{B}_{E}^{i} . Each source task corresponds to a different POMDP 222 with different underlying rewards, but shared dynamics. The underlying state space may also change 223

Algorithm 2 Zero-Shot Transfer with V-MAIL

- 1: **Require**: Expert demos \mathcal{B}_E^i for each source task, expert demos \mathcal{B}_E for target task 2: Randomly initialize policy π_{ψ} , and discriminator D_{ψ}
- 3: Train Alg 1 on source tasks, yielding shared model $\{q_{\theta}, \widehat{\mathcal{T}}_{\theta}\}$ and aggregated replay buffer \mathcal{B}_{π}
- 4: for number of training iterations do
- // Dynamics Fine-Tuning using Expert Trajectories 5:
- Update the variational model $\{q_{\theta}, \widehat{\mathcal{T}}_{\theta}\}$ using Equation 7 with data from $\mathcal{B}_E \cup \mathcal{B}_{\pi}$ 6:
- 7: // Adversarial Policy Learning
- 8: Update discriminator D_{ψ} and policy π_{ψ} with Equations 6 and 8.

across tasks, but the dynamics and observation model are shared across tasks. During training, the 224 agent can interact with each source environment and collect additional data. At test time, we're 225

introduced with a new target task T with corresponding expert demonstrations \mathcal{B}_E and the goal is to 226

obtain a policy that achieves high reward without additional interaction with the environment. 227

Our key observation is that we can optimize Eq. 6 under our model and still obtain an upper bound 228 on policy sub-optimality via Eq. 4. Furthermore, the sub-optimality is bound by the accuracy of our 229 model over the marginal state-action distribution of the target task expert. Specifically, we first train 230 on all of the source tasks using Algorithm 1, training a single shared variational model across the 231 tasks. By fine-tuning that model on data that includes the target task expert demonstrations our hope 232 is that we can get an accurate model and thus a high-quality policy. Similarly to Algorithm 1, we 233 then train a discriminator and policy for the target task using only model rollouts. This approach is 234 outlined in Algorithm 2. 235

5 **Experiments** 236

In our experiments, we aim to answer several questions: (1) can V-MAIL successfully scale to 237 environments with image observations, (2) how does V-MAIL compare to state of the art model-free 238 imitation approaches, (3) can V-MAIL solve realistic manipulation tasks and environments with 239 complex physical interactions, and (4) can V-MAIL enable zero-shot transfer to new tasks? All 240 experiments were carried out on a single Titan RTX GPU using an internal cluster for about 1000 241 GPU hours. 242

5.1 Single-Task Experiments 243

Comparisons. To answer question (2), we choose to compare V-MAIL to model-free adversarial 244 and non-adversarial imitation learning methods. For the former, we choose DAC [27] as a representa-245 tive approach, which we equip with DrQ data augmentation for greater performance on vision-based 246 tasks. For the latter, we consider SQIL [39], also equipped with DrQ training. We refer to each 247 approach with data augmentation as DA-DAC and DA-SQIL respectively. Both of these methods 248 are off-policy algorithms, which we expect to be considerably more sample efficient than on-policy 249 methods like GAIL [22] and AIRL [14]. 250

Environments and Demonstration Data. To answer the above questions, we consider the five 251 visual control environments illustrated in Figure 2. We first evaluate our method on the visual 252 Cheetah and visual Walker tasks from the DeepMind Control Suite [43]. Following SQIL [39] we 253 also consider the classic Car Racing environment, which is difficult to solve even with ground-truth 254 rewards. In addition, we benchmark our method on a custom D'Claw environment from the Robel 255



Figure 2: Illustration of the environments used in our experiments: Cheetah, Walker, Car Racing, D'Claw, and Baoding Balls. In all environments, the agent has access only to the RGB image frames as observations, except with additional access to proprioception in the Baoding Balls environment.



Figure 3: Learning curves showing ground truth reward versus number of environment steps for V-MAIL (ours), prior model-free imitation learning approaches, and behavior cloning on five visual imitation tasks. We find that V-MAIL consistently outperforms prior methods in terms of sample efficiency, final performance, and stability, particularly for the first four environments where V-MAIL reaches near-expert performance. In the most challenging visual Baoding Balls task, which is notably difficult even with ground-truth state, only V-MAIL is able to make some progress, but all methods struggle. Confidence intervals are shown with 1 SD over 3 runs.

suite [1], entirely from images without proprioception. This makes the task challenging due to 256 a complex action dynamics, contact dynamics, and occlusions from the robot fingers. Our final 257 258 environment is the Baoding balls task from Nagabandi et al. [33]. This is an extremely challenging task for policy learning, even in the state-based case. All tasks are from raw RGB images, while the 259 Baoding balls task additionally includes robot proprioception. All methods receive access to use 10 260 expert demonstrations, with the exception of the Baoding environment, which uses 25 demonstrations. 261 The demonstrations for the DeepMind Control and D'Claw tasks are generated using a policy trained 262 with SAC [18], the expert data for the Car Racing environment is generated using Dreamer [20], and 263 the demonstrations for the Baoding task is generated using PDDM [33] from low-dimensional states. 264 Additional details on the experimental set-up are provided in the appendix. 265

Results. Experiment results are shown in Figure 3. To answer questions (1) and (2), we compare V-266 MAIL to DA-SQIL and DA-DAC on the Cheetah and Walker tasks. We find that V-MAIL efficiently 267 and reliably solves both tasks; in contrast, the model-free methods initially outperform V-MAIL, 268 but their performance has high variance across random seeds and exhibits significant instability. 269 Such stability issues have also been observed by Swamy et al. [42], which provides some theoretical 270 explanation in the case of SQIL and the suggestion of early stopping as a mitigation technique. In 271 the case of DAC, the reasons for instability are less clear. Motivated by instability we observed in 272 the critic loss for DA-DAC, we experimented with a number of mitigation strategies in an attempt 273 274 to improve DA-DAC, including constraining the discriminator, varying the buffer and batch sizes, 275 and separating the convolutional encoders of the discriminator and the actor/critic; however, these techniques didn't fully prevented the degradation in performance. 276

On the Car Racing environment, we find that DA-SQIL and DA-DAC can reach or outperform behavior cloning, but struggle to reach expert-level performance. In contrast, V-MAIL stably and reliably achieves near-expert performance in about 200k environment steps. Note that Reddy et al. [39] report expert-level performance on this task, but in an easier setting with double the number of expert demonstrations available (20 vs. 10). Given that tracks are randomly generated per episode demanding significant generalization, it is not surprising that the problem becomes considerably more difficult with only 10 demonstrations.

Finally, to answer question (3), we consider the D'Claw and Baoding Balls tasks. In the D'Claw environment, SQIL fails to make progress, while DA-DAC makes significant progress initially but quickly degrades. V-MAIL solves the task in less than 100k environment steps. In the most challenging visual Baoding Balls problem, involving a 26-dimensional control space, V-MAIL is the only algorithm to reach any success.

289 5.2 Transfer Experiments

Transfer Scenarios. To evaluate V-MAIL's ability to learn new imitation tasks in a zero-shot way 290 (i.e. without any additional environment samples) we deploy Algorithm 2 on two domains: in a 291 locomotion experiment we train on the Walker Stand and Walker Run (target speed greater than 8) 292 tasks and and evaluate transfer to the Walker Walk (target speed between 2 and 4) task from the 293 DeepMind Control suite. In a manipulation scenario, we use a set of custom D'Claw Screw tasks from 294 the Robel suite [1]. We train our model on the 3-prong tasks with clockwise and counter-clockwise 295 rotation, as well as the 4-prong task with counter-clockwise rotation and evaluate transfer to the 296 4-prong task with clockwise rotation. 297

Comparisons. To our knowledge, no prior work has considered this zero-shot transfer scenario 298 previously. Thus, we devise several points of comparison. First, we compare to directly applying 299 the policy learned in the most related source task to the target task. This tests whether the target 300 task demands qualitatively distinct behavior. Second, we compare to an offline version of DAC, 301 augmented with the CQL approach [29], where samples collected from the source task are used to 302 update the policy, with the target task demonstrations used to learn the reward. Finally, we also 303 compare to behavior cloning on the target task demonstrations (without leveraging any source task 304 data), and an oracle that performs V-MAIL on the target task directly. 305

Results. Our results are shown in Table 1. 306 Policy transfer performs poorly, suggesting 307 that the target task indeed requires quali-308 tatively different behaviour from the few 309 training tasks available. Further, behavior 310 cloning on the target demonstrations is not 311 sufficient to learn the task. Offline DAC 312 also shows poor performance. Finally, we 313 see that V-MAIL almost matches the per-314 formance of the agent explicitly trained on 315 task, indicating the learned model and the 316 algorithm for training within that model can 317 be used not just for efficient visual imitation 318 learning, but also for zero-shot transfer to 319 320 new tasks.

Method	Walker Walk	Claw Rotate
Offline DAC	8.8%	-0.7%
Behavior cloning	26.8%	8.3%
Policy transfer	21.3%	5.6%
V-MAIL (ours)	92.7%	97.9%
Target task IL (oracle)	98.2%	102.3%

Table 1: Performance on zero-shot transfer to a new imitation learning task as percent of expert return. Each method is provided with 10 demonstrations of the target task, and zero additional samples in the environment. V-MAIL can solve the target tasks within its learned model without any additional samples, while model-free transfer learning approaches fail.

321 6 Conclusion

In this work we presented V-MAIL, a model-based imitation learning algorithm that works from high-dimensional image observations. V-MAIL learns a model of the environment, which serves a strong supervision signal for visual representation learning, as well as allowing us to train an imitation learning algorithm on-policy, without sacrificing sample efficiency. V-MAIL achieves better asymptotic returns, is more stable, and matches the sample efficiency of off-policy model-free approaches. We also find that by training a policy using only model rollouts, our approach is a strong procedure for zero-shot transfer to novel imitation learning tasks.

Future Work. We believe this work opens the door for many potential developments. One direction is to use recent developments in variational models to train our procedure using only expert observations without access to expert actions, which is an even more realistic scenario. This setup is quite difficult for model-free approaches, since expert actions usually serve as a strong supervision. Another direction is to use on-policy model based rollouts to efficiently train other algorithms that inherently require on-policy data, such as multi-modal imitation [31, 21]. Finally, the experiments suggest that this algorithm is efficient enough to be applied to real robots, an interesting direction for future work.

Limitations. Although successful in domains with complex dynamics, crucially our approach relies on variational models with compact, single-level, latent state spaces. It is possible that this model class could not have the capacity to represent complex realistic scenes such as large-scale cluttered environments, cloth and deformable object dynamics, realistic city scenes or home environments, which would limit real world applications.

Negative Societal Impacts. We do not anticipate any negative societal impacts that are unique to this paper compared to prior imitation learning works.

343 **References**

- [1] Michael Ahn, Henry Zhu, Kristian Hartikainen, Hugo Ponte, Abhishek Gupta, Sergey Levine,
 and Vikash Kumar. ROBEL: RObotics BEnchmarks for Learning with low-cost robots. In
 Conference on Robot Learning (CoRL), 2019.
- [2] Dario Amodei, Chris Olah, J. Steinhardt, Paul F. Christiano, John Schulman, and Dan Mané.
 Concrete problems in ai safety. *ArXiv*, abs/1606.06565, 2016.
- [3] Nir Baram, Oron Anschel, and Shie Mannor. Model-based adversarial imitation learning.
 Conference on Neural Information Processing Systems, 2016.
- [4] David M. Blei, A. Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112:859 – 877, 2016.
- [5] Lionel Blondé and Alexandros Kalousis. Sample-efficient imitation learning via generative
 adversarial nets. *AISTATS*, 2019.
- [6] Lionel Blondé, Pablo Strasser, and Alexandros Kalousis. Lipschitzness is all you need to tame
 off-policy generative adversarial imitation learning, 2020.
- [7] Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample efficient reinforcement learning with stochastic ensemble value expansion. *Conference on Neural Information Processing Systems*, 2019.
- [8] Neha Das, Sarah Bechtle, Todor Davchev, Dinesh Jayaraman, Akshara Rai, and Franziska Meier.
 Model-based inverse reinforcement learning from visual demonstrations. *Conference on Robot Learning*, 2020.
- [9] Tom Everitt and Marcus Hutter. Reward tampering problems and solutions in reinforcement
 learning: A causal influence diagram perspective. *ArXiv*, abs/1908.04734, 2019.
- [10] Benjamin Eysenbach, Sergey Levine, and Ruslan Salakhutdinov. Replacing rewards with
 examples: Example-based policy search via recursive classification, 2021.
- [11] Vladimir Feinberg, Alvin Wan, Ion Stoica, Michael I. Jordan, Joseph E. Gonzalez, and Sergey
 Levine. Model-based value estimation for efficient model-free reinforcement learning. *International Conference on Machine Learning*, 2018.
- [12] Chelsea Finn, Paul Christiano, Pieter Abbeel, and Sergey Levine. A connection between
 generative adversarial networks, inverse reinforcement learning, and energy-based models.
 ArXiv Preprint, 2016.
- [13] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal
 control via policy optimization. In *International conference on machine learning*, pp. 49–58.
 PMLR, 2016.
- ³⁷⁶ [14] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse ³⁷⁷ reinforcement learning. *International Conference on Learning Representations*, 2018.
- [15] Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G. Bellemare. Deep mdp: Learning continuous latent space models for representation learning. *International Conference on Machine Learning*, 2019.
- [16] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimiza tion perspective on imitation learning methods. *Conference on Robot Learning*, 2019.
- [17] Ian J. Goodfellow, Jean Pouget-Abadie, M. Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
 Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [18] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
 maximum entropy deep reinforcement learning with a stochastic actor. *International Conference on Machine Learning*, 2018.
- [19] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and
 James Davidson. Learning latent dynamics for planning from pixels. *International Conference on Machine Learning*, 2019.
- [20] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control:
 Learning behaviors by latent imagination. *International Conference on Learning Representations*, 2020.

- [21] Karol Hausman, Yevgen Chebotar, Stefan Schaal, Gaurav Sukhatme, and Joseph Lim. Multi modal imitation learning from unstructured demonstrations using generative adversarial nets,
 2017.
- [22] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Conference on Neural Information Processing Systems*, 2016.
- [23] Maximilian Karl, Maximilian Soelch, Justin Bayer, and Patrick van der Smagt. Deep variational
 bayes filters: Unsupervised learning of state space models from raw data. *International Conference on Machine Learning*, 2017.
- Liyiming Ke, Matt Barnes, W. Sun, Gilwoo Lee, Sanjiban Choudhury, and S. Srinivasa. Imitation
 learning as f-divergence minimization. *ArXiv*, abs/1905.12888, 2019.
- [25] Khimya Khetarpal, Matthew Riemer, I. Rish, and Doina Precup. Towards continual reinforce ment learning: A review and perspectives. *ArXiv*, abs/2012.13490, 2020.
- ⁴⁰⁶ [26] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.
- [27] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan
 Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in
 adversarial imitation learning. *International Conference on Learning Representations*, 2019.
- [28] Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribu tion matching. *International Conference on Learning Representations*, 2020.
- [29] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for
 offline reinforcement learning. *Conference on Neural Information Processing Systems*, 2020.
- [30] Alex X. Lee, Anusha Nagabandi, Pieter Abbeel, and Sergey Levine. Stochastic latent actor-critic:
 Deep reinforcement learning with a latent variable model. *Conference on Neural Information Processing Systems*, 2020.
- [31] Yunzhu Li, Jiaming Song, and Stefano Ermon. Infogail: Interpretable imitation learning from
 visual demonstrations. *Conference on Neural Information Processing Systems*, 2017.
- [32] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, P. Abbeel, and Igor Mordatch. Multi-agent
 actor-critic for mixed cooperative-competitive environments. In *NIPS*, 2017.
- [33] Anusha Nagabandi, Kurt Konolige, Sergey Levine, and Vikash Kumar. Deep dynamics models
 for learning dexterous manipulation. *Conference on Robot Learning*, 2019.
- [34] Dean A Pomerleau. Alvinn: an autonomous land vehicle in a neural network. In *Proceedings* of the 1st International Conference on Neural Information Processing Systems, pp. 305–313,
 1988.
- [35] Rémy Portelas, Cédric Colas, Lilian Weng, Katja Hofmann, and Pierre-Yves Oudeyer. Auto matic curriculum learning for deep rl: A short survey. *ArXiv*, abs/2003.04664, 2020.
- [36] Rafael Rafailov, Tianhe Yu, Aravind Rajeswaran, and Chelsea Finn. Offline reinforcement
 learning from images with latent space models. *arXiv preprint arXiv:2012.11547*, 2020.
- [37] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
 Todorov, and Sergey Levine. Learning Complex Dexterous Manipulation with Deep Reinforce ment Learning and Demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*,
 2018.
- [38] Aravind Rajeswaran, Igor Mordatch, and Vikash Kumar. A Game Theoretic Framework for
 Model-Based Reinforcement Learning. In *ICML*, 2020.
- [39] Siddharth Reddy, Anca D. Dragan, and Sergey Levine. Sqil: Imitation learning via reinforcement
 learning with sparse rewards. *International Conference on Learning Representations*, 2020.
- [40] Stephane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. A reduction of imitation learning
 and structured prediction to no-regret online learning. *AISTATS*, 2011.
- [41] Jonathan Spencer, Sanjiban Choudhury, Arun Venkatraman, Brian Ziebart, and J. Andrew
 Bagnell. Feedback in imitation learning: The three regimes of covariate shift. *ArXiv Preprint*, 2021.
- [42] Gokul Swamy, Sanjiban Choudhury, Zhiwei Steven Wu, and J. Andrew Bagnell. Of moments
 and matching: Trade-offs and treatments in imitation learning. 2021.

- [43] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David
 Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin
 Riedmiller. Deepmind control suite, 2018.
- [44] Manuel Watter, Jost Tobias Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed
 to control: A locally linear latent dynamics model for control from raw images, 2015.
- ⁴⁵⁰ [45] Marvin Zhang, Sharad Vikram, Laura Smith, Pieter Abbeel, Matthew J. Johnson, and Sergey
- Levine. Solar: Deep structured representations for model-based reinforcement learning. *International Conference on Machine Learning*, 2019.
- [46] Brian D. Ziebart, Andrew L. Maas, J. Bagnell, and A. Dey. Maximum entropy inverse reinforce ment learning. In *AAAI*, 2008.

455 Checklist

456	1.	For al	ll authors
457		(a) I	Do the main claims made in the abstract and introduction accurately reflect the paper's
458		С	contributions and scope? [Yes]
459		(b) I	Did you describe the limitations of your work? [Yes] Limitations are discussed in the
460		С	conclusion.
461		(c) I	Did you discuss any potential negative societal impacts of your work? [Yes] This is
462			addressed in the conclusion.
463		(d) H	Have you read the ethics review guidelines and ensured that your paper conforms to
404	2	ייייי נ	
465	Ζ.	II you	
466		(a) I	Did you state the full set of assumptions of all theoretical results? [Yes] The com-
467		p	plete assumptions and full proofs are included as an Appendix with the supplemental
400		(b) I	Did you include complete proofs of all theoretical results? [Ves] The complete assump-
470		(b) 1 t	ions and full proofs are included as an Appendix with the supplemental materials.
471	3.	If you	ran experiments
472		(a) I	Did you include the code, data, and instructions needed to reproduce the main experi-
473		n	nental results (either in the supplemental material or as a URL)? [Yes] Our codebase
474		i	s available with the supplemental materials.
475		(b) I	Did you specify all the training details (e.g., data splits, hyperparameters, how they were
476		С	chosen)? [Yes] Training details are available as an appendix with the supplemental
477		n () r	naterials.
478 479		(c) I n	nents multiple times)? [Yes] Consult Figure 3.
480		(d) I	Did you include the total amount of compute and the type of resources used (e.g., type
481		С	of GPUs, internal cluster, or cloud provider)? [Yes] Refer to Experiments section.
482	4.	If you	are using existing assets (e.g., code, data, models) or curating/releasing new assets
483		(a) I	If your work uses existing assets, did you cite the creators? [Yes] We build our
484		C ·	codebases on top of several publicly available developments, more details are included
485		1	n the supplement materials.
486		(b) I	Did you mention the license of the assets? [N/A]
487 488		(c) 1 (Did you include any new assets either in the supplemental material or as a URL? [Yes] Our codebase is included with the supplement materials.
489		(d) I	Did you discuss whether and how consent was obtained from people whose data you're
490		u	using/curating? [N/A]
491		(e) I	Did you discuss whether the data you are using/curating contains personally identifiable
492		i	nformation or offensive content? [N/A]
493	5.	If you	used crowdsourcing or conducted research with human subjects
494		(a) I	Did you include the full text of instructions given to participants and screenshots, if
495		a	applicable? [N/A]
496		(b) I	Did you describe any potential participant risks, with links to Institutional Review
497		ł	Board (IKB) approvals, if applicable? [N/A]
498 499		(c) I s	Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? $[N/A]$