VIVA: VIDEO-TRAINED VALUE FUNCTIONS FOR GUIDING ONLINE RL FROM DIVERSE DATA

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027 028 029

030

Paper under double-blind review

ABSTRACT

Online reinforcement learning (RL) with sparse rewards poses a challenge partly because of the lack of feedback on states leading to the goal. Furthermore, expert offline data with reward signal is rarely available to provide this feedback and bootstrap online learning. How can we guide online agents to the right solution without this on-task data? Reward shaping offers a solution by providing fine-grained signal to nudge the policy towards the optimal solution. However, reward shaping often requires domain knowledge to hand-engineer heuristics for a specific goal. To enable more general and inexpensive guidance, we propose and analyze a data-driven methodology that automatically guides RL by learning from widely available video data such as Internet recordings, off-task demonstrations, task failures, and undirected environment interaction. By learning a model of optimal goal-conditioned value from diverse passive data, we open the floor to scaling up and using a wide variety of data sources to model general goal-reaching behaviors relevant to guiding online RL. Specifically, we use intent-conditioned value functions to learn from diverse video and incorporate these goal-conditioned values into the reward. Our experiments show that video-trained value functions work well with a variety of data sources, exhibit positive transfer from human video pre-training, can generalize to unseen goals, and scale with dataset size.

1 INTRODUCTION

Many sequential decision-making tasks are naturally defined with a sparse reward, meaning the agent only receives positive signal when the goal has been achieved. Unfortunately, these sparse reward tasks are especially challenging in reinforcement learning (RL) (Sutton, 2018) since they provide no signal at intermediate states, effectively requiring exhaustive search. Practitioners often resort to collecting task-relevant prior data (Pomerleau, 1988) or hand-designing task-relevant dense reward functions (Mataric, 1994). However, manually collecting this high-quality data or defining a task-specific reward is time-intensive and not general.

To solve this problem in RL, we should guide the search procedure online towards the desired goal. This dictates the usage of some general prior informing the agent what states lead to others to direct it to the goal. Humans make use of extensive prior knowledge when attempting to accomplish tasks: for example, we know that finding a mug generally requires us to try looking in cabinets and that opening them requires interacting with the handle. We posit that this prior can in fact be learned with task-agnostic environment data and general manipulation videos to develop a sense of "how the world works." This data is easily collected in the environment or mined from the web, respectively.

To leverage both of these data types, we choose to learn from *video*, enabling the use of a myriad of datasets without needing embodiment-specific actions or task-specific rewards. The nature of video data availability on the web also allows for training on other environments. We hypothesize that learning models from various video sources will expand the data support, enabling generalization and successful goal-reaching guidance. Crucially, we elect to represent our prior as a *goalconditioned state-value function* V(s, g), that for any image s and desired target g, estimates the temporal distance between the two states. Learning this type of model easily plugs into online RL by penalizing the predicted distance from the goal. Also, using a value-learning approach allows ingesting suboptimal reaching data, further relaxing our requirements for the training data, as opposed to other behavioral prior methods (Escontrela et al., 2023). Lastly, goal-conditioned value



Figure 1: Left: ViVa uses samples from internet-scale video to learn a value-function that encodes goal-reaching priors. Middle: ViVa finetunes on robotics-relevant data to bring the value function 072 into the domain of the tasks we wish to solve. Right: During online RL, we freeze the value function and augment the extrinsic reward with a guidance signal that captures temporal distances. We choose 074 to include the robotics-relevant interaction data in our online pipeline to assist exploration.

077 functions naturally extend to multiple tasks by flexible goal specification. In essence, we desire a 078 simply learned function that scales with widely available off-task and off-environment video data to 079 generally inform the agent about useful states leading to the goal.

We can instantiate our method by pre-training an Intent-conditioned Value Function (ICVF) (Ghosh 081 et al., 2023) on Internet-scale egocentric interaction data in various settings (Ego4D) (Grauman et al., 2022). We use this training to develop strong visual features as well as a priors over object 083 manipulation and interaction outcomes. We then finetune this ICVF on environment-specific, yet 084 task-agnostic data, to specialize the function for the setting of interest. During online RL, we provide 085 the temporal distance estimates to the agent in the form of a reward penalty.

We observe that reformulating online sparse RL problems with Video-trained Value functions (ViVa) 087 shows a number of benefits. Firstly, we see generalization to new goals unseen in prior data in the 088 Antmaze environment (Fu et al., 2021), a simple state-based control setting. Secondly, we also see 089 improvement in performance by training on off-task data in a visual robotic simulator, RoboVerse 090 (Singh et al., 2020). Thirdly, we see that pre-training on Ego4D can significantly improve perfor-091 mance but is not sufficient to solve online RL alone, necessitating some environment finetuning. 092 Lastly, we see that ViVa improves online performance as data scale increases and can enable solving 093 complex robotic tasks on Franka Kitchen (Gupta et al., 2019), another robotic simulator.

094 095

071

073

075 076

2 **RELATED WORK**

096 097

098 Solving sparse online RL problems is a difficult challenge due to the lack of reward feedback. One 099 way to make it easier is to better explore the environment to more reliably reach the goal state and begin backing up rewards. They range from simple noisy behaviors (Haarnoja et al., 2018b) to 100 structured behavioral priors (Ecoffet et al., 2021; Bharadhwaj et al., 2021; Kearns & Singh, 2002; 101 Brafman & Tennenholtz, 2003). Some methods utilize intrinsic bonuses (Schmidhuber, 2010) to 102 minimize uncertainty (Kolter & Ng, 2009; Pathak et al., 2019; Houthooft et al., 2017; Still & Pre-103 cup, 2012) or to seek novelty (Burda et al., 2018; Pathak et al., 2017; Ostrovski et al., 2017; Tang 104 et al., 2017; Bellemare et al., 2016). Unfortunately, these methods break down in complex visual 105 environments and intricate robotic control settings due to the large state and action space. 106

To narrow this search, a prior is desirable to inform the agent of what states or actions to explore 107 more. One way to do this is to inject domain knowledge into the reward function, guiding it to the goal. This family of approaches, known as reward shaping, can accelerate learning the optimal policy (Ng et al., 1999; Mataric, 1994; Hu et al., 2020; Devlin & Kudenko, 2012; Wiewiora, 2003). However, hand-crafting these rewards does not generally scale to many tasks and is often over-designed for one domain (Jiang et al., 2020; Mahmood et al., 2018; Haarnoja et al., 2018a; Malysheva & Kudenko, 2018; Hussein et al., 2017; Brys et al., 2015). A more ideal way to have a general prior is to learn it from a wide range of available data. In this work, we explore the effect of various video data sources in providing robotics-relevant dynamics information to downstream RL.

115 Many methods elect to use this cheap video data to learn a rich image representation through re-116 construction objectives (Xiao et al., 2022), constrastive learning (Nair et al., 2022), value-functions 117 (Bhateja et al., 2023), or predictive objectives (Shah & Kumar, 2021). There is also a family of 118 approaches that model videos through inferring latent actions from states, and use environmentspecific action-labelled data to map these latent actions to real actions (Ye et al., 2024; Edwards 119 et al., 2019; Schmidt & Jiang, 2024; Bruce et al., 2024). Bhateja et al. (2023) propose V-PTR which 120 is particularly similar to our approach but only utilizes the trained ICVF encoders as a pre-trained 121 representation for offline RL. Our method aims to use a distance-function rather than a pre-trained 122 encoder to directly guide a goal-conditioned online RL agent. This resembles temporal distance 123 learning methods (Pong et al., 2020; Mezghani et al., 2023) such as Dynamical Distance Learning 124 (DDL) (Hartikainen et al., 2020) where policy-conditioned distance learning and online RL for dis-125 tance minimization is alternated. However, DDL uses distances for unsupervised skill discovery and 126 preference-learning, and importantly do not extend to internet-scale interaction data. 127

The most similar approach is Value-Implicit Pretraining (VIP) (Ma et al., 2023) whereby a internet-128 scale video-trained representation function induces a distance to shape the reward. Our method 129 differs from VIP in a few different ways. First, our method explicitly uses temporal-difference 130 learning as opposed to time-contrastive learning, as done in VIP. Second, we explicitly focus on 131 downstream online RL rather than direct imitation or smooth trajectory optimization. Third, we 132 present a bi-level pre-training procedure to not only take advantage of task-agnostic human video, 133 but also environmental interaction data. We therefore identify with other offline-to-online methods 134 (Xie et al., 2022; Lee et al., 2021; Agarwal et al., 2022; Zheng et al., 2023; Andrychowicz et al., 135 2018; Li et al., 2023) whereas VIP compares to other pre-trained representation distances. These offline-to-online methods often assume action access though which limits the scope of usable data. 136 Our method's access to environmental interaction data dictates comparison to RLPD (Ball et al., 137 2023), a method which runs online RL and mixes training batches with offline data samples, as well 138 as JSRL (Uchendu et al., 2023), a method which condenses offline data into a policy to assist online 139 exploration. 140

141 142

143

3 PRELIMINARIES

144 Let S be the state space and A be the action space. We consider a sparse-reward Markov Decision 145 Process (MDP), M defined by a tuple (S, A, P, r, γ) where P(s'|s, a) is the transition dynamics and 146 γ is the discount factor. We additionally consider a goal specified by a goal state set G. The reward 147 r(s) is the set inclusion indicator $r(s) = \mathbb{1}[s \in G]$. The objective in this setting is learn a policy π 148 that maximizes the expected return $\mathbb{E}_{a \sim \pi(s_t), s_{t+1} \sim P(.|s_t, a)}[\sum_{t=0}^{\infty} \gamma^t r(s_t)]$ where the expectation is 149 taken over the policy π and the environment dynamics.

For our experiments, we assume access to a video dataset of human egocentric interactions, \mathcal{D}_{video} , and a dataset of environment-specific interaction, \mathcal{D}_{env} . \mathcal{D}_{video} contains data out of the desired domain and does not use the same embodiment as used for the target MDP \mathcal{M} . \mathcal{D}_{env} is environmentspecific data that contains actions, uses the embodiment of interest, but is either agnostic to the actual task at hand, or does not contain any successful trajectories due to the expensive nature of positive data trajectories.

156

4 VIVA : VIDEO-TRAINED VALUE FUNCTIONS

157 158

Our proposed solution for the sparse online RL case when faced with a lack of demonstrations is to develop a prior that guides the agent towards a valid goal, $g \in \mathcal{G}$. We elect to learn a value function V(s,g) to give the value of any given state, s, in the context of the task of reaching the state g optimally. As detailed in 4.1, we can train this value function to directly represent the temporal 



Figure 2: Left: A visualization of trajectories from the corrupted dataset shown in green. Middle: The learned ICVF values across all states with the goal at the red star. **Right:** The optimal dense reward (i.e. L2 distance) for all states with the goal at the red star.

distance from s to g, thus giving a simple reward penalty. This allows us to create a guided reward which has an injected prior towards the goal of choice.

$$\hat{r}(s,a) = r(s,a) + V(s,g) \tag{1}$$

179 4.1 VALUE-FUNCTION GUIDANCE

180 Our desired function is V(s,g), which generally yields a higher value for states closer to g on the 181 optimal path from s to g. Since we aim to learn this model from action and reward free video data, 182 we elect to model an Intent-conditioned Value Function, $ICVF(s, s^+, g)$, which is fully trainable 183 from this passive data. The ICVF models the unnormalized likelihood of reaching some outcome 184 state, s^+ , when starting in state s and acting optimally to reach some goal state g, otherwise 185 known as the "intent". To precisely define the ICVF, we denote $r_g : s \mapsto r$ as a reward function 186 corresponding to reaching any goal state. The optimal policy, $\pi_{r_g}^*$, induces a state-transition which 187 can define the value function based on the following expectation:

$$P_{g}(s_{t+1}|s_{t}) = P^{\pi_{r_{g}}}(s_{t+1}|s_{t})$$

$$r_{g}(s) = \mathbb{1}[s = g] - 1$$

$$ICVF(s, s^{+}, g) = \mathbb{E}_{s_{0}=s, s_{t+1} \sim P_{g}(.|s_{t})} \sum_{t=0}^{\infty} \gamma^{t} r_{s^{+}}(s_{t}).$$
(2)

By applying a scalar shift of -1 to our sparse reward, the reward-to-go is equivalent to the negative discounted number of timesteps to reach the goal. This negative temporal distance is well-suited to be used as an additive reward penalty. Furthermore, if we use g as not only the goal, but also the outcome, s^+ , we can model the negated time to reach g from s if the agent were to act optimally towards g thereafter. This is what we are looking for and can let us define our desired value function:

$$V(s,g) = ICVF(s,g,g) = \mathbb{E}_{s_0=s,s_{t+1}\sim P_g(.|s_t)} \sum_{t=0}^{\infty} \gamma^t r_g(s_t)$$
$$\hat{r}(s,a) = r(s,a) + ICVF(s,g,g).$$
(3)

We incorporate \hat{r} , our guided reward, into the online RL system. This allows the agent to apply knowledge of state-goal relationships contained in the learned ICVF. We note that usage of a potential-based instrinsic reward could be used for provable policy invariance as shown by Ng et al. (1999), but we observe higher variance returns which could destabilize training shown in Appendix A.3.

208 4.2 VALUE-FUNCTION TRAINING

We model the ICVF as a monolithic neural network, $V_{\theta}(s, s^+, g)$. This differs from the original multilinear formulation, $\phi_{\theta}(s)^T T_{\theta}(g) \psi_{\theta}(s^+)$, since we found a monolithic architecture to produce higher-quality value functions as shown in Figure 11 in the Appendix. When working with image states, we elect to feed in learnable latent representations of the inputs to the value function. We detail the training procedure below.

Given a video dataset of image sequences, \mathcal{D} , we first sample a starting frame and neighboring frame (s, s') from the same trajectory. Second, we sample some outcome s^+ from the future of the same



Figure 3: All plots detail the mean evaluation return computed over 10 evaluation episodes. Left: Online RL for pick-and-place on COG as we scale to more and more on-task data. The rows below show example off-task successful trajectories with the WidowX robot from the drawer_prior and blocked_drawer datasets. Right: Online RL for pick-and-place on COG when including Ego4D pretraining and off-task data sources. The rows below are a failure and a success from the prior dataset.

trajectory, and we sample a goal, g, in an identical way to s^+ . We additionally follow Ghosh et al. (2023) in sometimes sampling identical images or random images for s^+ and q for better training. After retrieving a sample, we minimize the temporal-difference (TD) error, in Equation 4. Inspired by Kostrikov et al. (2021a), we use the expectile regression framework with an advantage heuristic, shown in Equation 5, to relax any maximization operators. This expectile biases the objective to more strongly weight samples (s, s') that are approaching q under our current model of value.

$$\min_{\theta} |\alpha - \mathbb{1}(A \le 0)| * (V_{\theta}(s, s^+, g) - \mathbb{1}(s = s^+) - \gamma V_{\theta}(s', s^+, g))^2$$
(4)

$$A = \mathbb{1}(s = g) + \gamma V_{\theta}(s', g, g) - V_{\theta}(s, g, g)$$
(5)

Essentially, if transitioning to s' while conditioned on g is advantageous under our current value 255 estimates, we assume that the transition is implicitly running the optimal action to reach g. This 256 allows us to update our value function without a maximum operation across actions. As a result, we just minimize the one-step TD error which is equivalent to regressing our value estimate of 258 $V_{\theta}(s, s^+, g)$ towards $\mathbb{1}(s = s^+) + \gamma V_{\theta}(s', s^+, g)$. We use the expectile, α , to decide how hard or soft this assumption is, with $\alpha = 0.5$ equating all samples to be equal weight, and $\alpha = 1$ forcing only 260 using positive advantage samples for updates. As shown by Kostrikov et al. (2021a), this converges in the limit as α approaches 1 262

4.3 SYSTEM OVERVIEW

236

237

238

239

240

241 242 243

244

245

246

247

248 249 250

251

253 254

257

259

261

263

264

265 Video pre-training Using the training process described in 4.2, we first train an ICVF on Ego4D, 266 or \mathcal{D}_{video} . Ego4D is a dataset of first-person camera video from hundreds of participants across 267 many diverse scenes. This video data contains humans doing daily-life activities such as laundry, lawn-mowing, sports, gardening, and more. Approximately 3000 hours of video data is included and 268 we reshape to 128×128 and apply a random crop augmentation further detailed in Appendix A.1. 269 As detailed earlier, we utilize a -1 reward shift for the self-supervised reward targets to ensure the 270 value to-go matches a temporal distance as desired. We elect to sample future outcomes and goals 271 from the same trajectory 80% of the time and use a 10% chance for both choosing random goals or 272 goals equal to current sampled state. Lastly, we choose an expectile of 0.9 which ensures backups 273 are biased to occur stronger for transitions where the advantage heuristic is positive. This expectile 274 allows for the convergence guarantees in optimal value function learning as the expectile approaches 1 shown in Kostrikov et al. (2021a). We utilize ResNetv2 (He et al., 2016) on JAX (Bradbury et al., 275 2018) as our neural architecture and functional paradigm for this video pre-training. We encode the 276 three input images, (s, s^+, g) with the ResNet before passing them into an ensemble of two 2-layer MLPs for min-Q learning. 278

279 Environment fine-tuning Secondly, we use 280 available environment data, \mathcal{D}_{env} , to finetune the ICVF. The finetuning is done exactly the 281 same way as pre-training but with environment 282 video data. This finetuning brings the model 283 into the domain of the RL task and can help 284 to develop setting-specific features relevant to 285 tasks in the environment. We hypothesize usage of \mathcal{D}_{video} will develop general visual features 287 and fusion between the input and goal images. 288 Furthermore, it can learn priors about the cause-289 and-effect of manipulation. Alternatively, the 290 fine-tuning on \mathcal{D}_{env} will help to develop task-291 specific features and the visual dynamics of the target environment. 292

Guided online RL After our value function is trained, we lastly run online RL and utilize the available environmental data, \mathcal{D}_{env} , as prior



Figure 4: The online evaluation return in AntMaze when training ViVa with corrupted data. As seen, learning a value-function prior for online RL provides a more generalizable reward model when offline rewarded data is absent. Learning a behavioral prior also works in this setting.

296 data. Specifically, every batch update for online RL includes 50% sampled online data from the re-297 play buffer and 50% sampled offline data from the prior dataset. The addition of the prior data into our RL training assists online exploration by providing offline trajectories to backup across and ex-298 plore which may not be otherwise explored online. We run experiments with this system set up to 299 study and analyze generalization characteristics, performance, and data scaling properties. We pro-300 vide details of our chosen algorithms Soft Actor-Critic (Haarnoja et al., 2018b) and its DrQ variant 301 (Kostrikov et al., 2021b) in the Appendix. As shown by Haarnoja et al. (2018b), soft policy iteration 302 is shown to converge. 303

304 305

306

5 Results

We analyze ViVa through different lenses to understand the benefits of video-trained value functions for downstream online RL. Specifically, we seek to study the **generalization** capabilities of ViVa in providing effective guidance for tasks it has not been provided data for, the **performance** of ViVa in difficult control tasks, and the **scaling** properties of ViVa as more diverse data is incorporated in greater quantities.

312

313 5.1 BASELINES

314 We also choose to compare to other methods which take advantage of offline data to determine 315 whether video-trained value functions are an effective mathematical object for representing a prior 316 for online RL. We firstly compare against Reinforcement Learning with Prior Data (RLPD), a 317 method which simply includes offline prior data in the update batches exactly as we do in ViVa 318 . Importantly, RLPD only uses extrinsic reward signals and does not pretrain or finetune a value for 319 relabeling offline data as ViVa does. We also compare against Jump-start Reinforcement Learning 320 (JSRL) which learns a behavioral prior policy from offline data and then runs online RL by execut-321 ing the learned prior policy for N random steps and then giving control to the agent's policy until termination. This method aims to condense prior experience into a policy for improving exploration 322 towards desired goals. For our experiments, we train an imitative policy from offline data and use 323 that as the behavioral prior for JSRL. Lastly, we use vanilla DreamerV2, a competitive world-model



Figure 5: All plots detail the mean evaluation return computed over 10 evaluation episodes. Left: Online RL for the Hinge Cabinet task in FrankaKitchen. The bottom row is an image trajectory of a demonstration of opening the hinge cabinet. **Right:** Online RL for the Sliding Cabinet task in FrankaKitchen. The bottom row is an image trajectory of a demonstration of opening the sliding cabinet.

approach for online RL (Hafner et al., 2022) that uses latent imagination of rollouts for training. We use these baselines to study the importance of explicitly learning a value function from prior data as opposed to directly including it in the replay buffer or learning a behavioral prior from it. We also compare against vanilla Soft-Actor Critic (SAC) and ablate Ego4D pre-training from ViVa to further analyze our method.

5.2 EXPERIMENTS

341

342

343

344

345 346 347

348

349

350

351

352 353

354

355 Corrupted AntMaze We first use the D4RL AntMaze (Fu et al., 2021) environment to visually 356 analyze the robustness to states seen outside of the training distribution. Environment and training 357 details are further expanded upon in Appendix A.3. We modify the D4RL diverse prior dataset 358 which includes the training transitions of a random start goal-reaching policy. Importantly, we corrupt this dataset by removing all trajectories containing points near the goal-region as shown in 359 Figure 2. We train a 3-layer Multilayer Perceptron with 512 units each using Equation 4 as the 360 training objective and display the learned value function, after 45 minutes of training on a Tesla 361 V100 16GB GPU, in Figure 2. Evidently, we observe generalization to the goal region when it has 362 not been seen during value training. The benefits of this generalization can be seen when running 363 downstream online RL are shown in Figure 4. We conclude that learning a simple ICVF network 364 on offline data is enough to develop a prior that generalizes to the unseen goal and prevents the 365 failures that RLPD exhibits on sparse-reward tasks when the offline dataset doesn't contain the 366 goal. Our comparison shows that JSRL exhibits similar extrapolation to new goals in the space of 367 expert actions as opposed to values. However, this similar extrapolation ability seems to fail when 368 introduced to more complex visual environments. In these settings ViVa is able to take advantage of Ego4D pre-training whereas JSRL cannot. 369

370 RoboVerse off-task transfer We use the RoboVerse (Singh et al., 2020) simulator (COG) to test 371 whether ViVa can generalize to new tasks never seen before in a visual domain, rather than state-372 based. This simulator has a variety of settings and accompanying datasets using a 6 DoF WidowX-373 250 robot on a desk. We choose to evaluate on a pick-and-place task to move a randomly placed 374 object into a tray – this task has two datasets of interest, one of 10K grasping attempts (with around 375 40% success), known as the prior dataset, and one task-specific dataset of 5K placing attempts (with around 90% success) labeled with rewards. For our experimental setup, we choose 376 to exclude the task-specific dataset to emphasize the absence of positive demonstration data. 377 During training, we combine the prior data with various other off-task sets which contain inter389

390

391

392

393

394 395



Figure 6: Left: An ablation study of including Ego4D pretraining or not across different environment finetuning data availability levels. At 0%, we are evaluating a random value function and the zero-shot performance of an Ego4D trained value function. **Right:** Each row is a randomly sampled trajectory from the Ego4D training showing car washing, cooking, and construction (from top to bottom).

actions with an open drawer, a closed drawer, and an obstructed drawer. All these datasets contain 397 no rewards as they do not match the desired reward we are using. We train ViVa for 9 hours on a 398 v4 TPU with these datasets and use it to guide online RL for pick-and-place. We leave details of 399 the training and environment in the Appendix. The results in Figure 3 show that ViVa succeeds in 400 solving the task whereas plainly sampling the available data offline fails. This is since the offline 401 data includes no rewards so RLPD fails to benefit from offline batch updates. On the other hand, the 402 imitative prior that JSRL uses does not explore the right areas which slows down learning. Interest-403 ingly, this experiment shows that ViVa is able to take advantage of diverse off-task environmental 404 and video data to inform goal-reaching. To concretize this conclusion, we ablate these data sources to show that this is what enables guidance to unseen goals in the Figure 9 in Appendix A.4. 405

To more deeply understand how the trained ICVF behaves on out-of-distribution examples, we also plot value curves over trajectories of failure demonstrations and unseen successes in Figure 7. As shown, ViVa provides guidance towards unseen goals, resembling how a control value function trained on positive demonstrations does. Similar to the AntMaze experiment, ViVa provides generalization to unseen goals and assists downstream online RL while also taking advantage of Internetscale video data. We analyze the usage of how ViVa behaves as more data is available and how useful this Internet-scale pre-training is in the experiments below.

413 **RoboVerse scaling law** In this experiment, we seek to assess whether a greater amount of task-414 relevant data has a positive effect on the downstream RL guidance. We train ViVa on a varying 415 amount of data (from the task-specific and prior datasets) for 5 hours on a v4 TPU 416 and then examine the online performance. We elect to include only the prior dataset in the online 417 RL phase since including the task-specific data would significantly simplify the problem. As shown in Figure 3, we can see there is a strong performance increase as data scales upwards. This 418 shows that ViVa benefits from the diversity and coverage of its training data and has positive scaling 419 behavior. 420

421 RoboVerse pre-training ablation We run the same analysis in our scaling law, but we remove 422 environment-agnostic, task-agnostic video to assess the direct impact of Ego4D pretraining. In 423 Figure 6 we can see a diminishing yet positive return from pre-training the value function on internetscale Ego4D video. Specifically, in the low-data regime, we observe a 2x increase in performance 424 when including Internet-scale video. This demonstrates the effective transfer to online RL from 425 including videos of interaction data supporting our initial hypothesis of developing a goal-reaching 426 prior for guidance. Although, the poor zero-shot performance depicts the importance of fine-tuning 427 on environmental data given the significant domain shift from Ego4D to RoboVerse. 428

Franka Kitchen Lastly, to evaluate on a more difficult robotic benchmark, we run ViVa on the
 FrankaKitchen (Fu et al., 2021; Gupta et al., 2019) environment which simulates a 9-DoF Franka
 robot tasked to interact with different objects in a kitchen. We use datasets of ~ 1K failure trajectories when attempting to interact with the Hinge Cabinet and Sliding Cabinet as the environmental



Figure 7: The top shows the image trajectory being evaluated and the bottom is the corresponding value function plot. Left: Value across a successful trajectory conditioned on a picking goal. Middle: Value across a failure trajectory conditioned on a picking goal. Right: Value on an unseen placing trajectory with an unseen placing goal. The blue is our model generalizing, and the orange reference is an optimal value function learned on the placing task.

interaction data to finetune the video-pretrained value function. Finetuning is run for 2.5 hours on a v4 TPU. Results of each task with and without video pretraining are shown in Figure 5. We observe that RLPD works well due to the negative reward shift that encourages the agent to be near the terminal states of the prior data. JSRL works poorly yet still succeeds on some seeds since imitating the interaction failures allow for exploring near the right area. Evidently, the inclusion of Ego4D pretraining tends to improve sample-efficiency.

6 DISCUSSION

444

445

446

447 448 449

450

451

452

453

454

455 456

457

In this paper, we proposed a method for transferring goal-reaching priors found in video data to downstream online RL problems by learning an intent-conditioned value function. This method can enable sparse-reward task solving, generalization to new goals, and positive transfer between tasks. Our analysis of ViVa illustrates the importance of using value function pre-training on video data as opposed to other methods of utilizing prior data. Our scaling experiments show that this is due to the wide support that this method can take advantage of, namely from the availability and generality of video data as well as the lack of assumptions for value learning.

Comparisons with JSRL depict the superiority of value functions as a representation of prior data as
opposed to classical imitative policies. We hypothesize this is because value learning uses a method
akin to shortest-path finding within data to discover an underlying temporal structure as opposed
to naively matching the next action. Furthermore, direct imitative policies prevent support from
action-free data sources such as Ego4D. However, latent-imitation methods could be explored to
leverage actionless datasets. Regardless, the ViVa paradigm should therefore provide insight to RL
practitioners looking to harness extra data and ameliorate the absence of rewarded prior data.

472 Limitations and future work We note that a limitation of value functions is the weak zero-shot ex-473 trapolation ability when far out of domain. This can be seen through the poor 0% scale performance 474 shown in Figure 6 which is presumably because RoboVerse is significantly different than Ego4D. 475 But when there is fine-tuning involved (shown at scales larger than 0% in Figure 6), this Ego4D 476 pretraining helps, offering 2x performance boost in the low-data regime. These results make it clear 477 that pre-training offers a way to make fine-tuning more effective, but cannot work on its own as it'd need task-relevant data. A direction of future work would be to find ways to encode more explicit 478 forms of abstraction in the value function in order to extrapolate deeply when given only off-domain 479 pre-training data (such as Ego4D). This would help to improve pure zero-shot performance when 480 given no environment data. 481

We also notice that ViVa utilizes some action-labelled robotic data for fine-tuning which is assumed
to be exploratory or somewhat relevant to the downstream task. An exciting future direction would
be to pair ViVa without fine-tuning with an exploration algorithm online to run the fine-tuning during the online RL phase, thus simplifying the training pipeline by removing a separate fine-tuning
phase. This method would also allow for resolving value errors through collecting counterfactual

examples since these errors can be detrimental to performance by forcing the agent into states that
are erroneously near the goal. This way, ViVa could even be used to form a curriculum based on
state-values or uncertainties in state-values to guide exploration in harder problems. Lastly, a natural extension includes utilizing language goals for the intent-conditioned value function, harnessing
multi-modal features, and extending into the real-world.

491 492 493

500 501

502

527

528

529

530

7 REPRODUCIBILITY STATEMENT

We include important training parameters in our system overview, in Section 4.3. These include image shapes, augmentation choices, reward shift, and hyperparameters that control the data sampling for training. In Section 5, we include domain specific parameters as well as the datasets used for fine-tuning ViVa . Lastly, we include an Appendix with fine-grained training details, datasets, and code-bases used. The Appendix expands upon details mentioned in the main paper and gives parameters for exactly reproducing the models we have trained on the code repositories we used.

References

- Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G. Bellemare.
 Reincarnating reinforcement learning: Reusing prior computation to accelerate progress, 2022.
 URL https://arxiv.org/abs/2206.01626.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob
 McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay, 2018.
- Philip J. Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning with offline data, 2023.
- Gabriel Barth-Maron, Matthew W. Hoffman, David Budden, Will Dabney, Dan Horgan, Dhruva TB,
 Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed distributional deterministic
 policy gradients, 2018. URL https://arxiv.org/abs/1804.08617.
- Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation, 2016.
- Homanga Bharadhwaj, Animesh Garg, and Florian Shkurti. Leaf: Latent exploration along the frontier, 2021.
- 519
 520
 521
 522
 523
 524
 525
 525
 526
 527
 527
 528
 529
 529
 520
 520
 521
 521
 522
 522
 523
 524
 525
 525
 526
 527
 527
 528
 529
 529
 529
 520
 520
 520
 521
 522
 522
 523
 524
 525
 525
 526
 526
 527
 528
 529
 529
 520
 520
 520
 520
 521
 522
 522
 522
 523
 524
 525
 525
 526
 526
 527
 528
 529
 529
 520
 520
 520
 521
 522
 522
 522
 522
 523
 524
 525
 526
 526
 527
 528
 529
 529
 520
 520
 520
 521
 521
 522
 522
 522
 523
 524
 525
 526
 527
 528
 529
 529
 520
 520
 520
 521
 521
 522
 521
 522
 522
 523
 524
 525
 526
 526
 527
 528
 528
 529
 529
 520
 520
 520
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
 Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http:
 //github.com/google/jax.
 - Ronen I. Brafman and Moshe Tennenholtz. R-max a general polynomial time algorithm for nearoptimal reinforcement learning. 3(null):213–231, mar 2003. ISSN 1532-4435. doi: 10.1162/ 153244303765208377. URL https://doi.org/10.1162/153244303765208377.
- Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, Yusuf Aytar, Sarah Bechtle, Feryal Behbahani, Stephanie Chan, Nicolas Heess, Lucy Gonzalez, Simon Osindero, Sherjil Ozair, Scott Reed, Jingwei Zhang, Konrad Zolna, Jeff Clune, Nando de Freitas, Satinder Singh, and Tim Rocktäschel. Genie: Generative interactive environments, 2024. URL https://arxiv.org/abs/2402.15391.
- Tim Brys, Anna Harutyunyan, Halit Bener Suay, S. Chernova, Matthew E. Taylor, and Ann Nowé. Reinforcement learning from demonstration through shaping. In *International Joint Conference on Artificial Intelligence*, 2015. URL https://api.semanticscholar.org/ CorpusID:1557568.

540

541

distillation, 2018. 542 Sam Michael Devlin and Daniel Kudenko. Dynamic potential-based reward shaping. In 11th In-543 ternational Conference on Autonomous Agents and Multiagent Systems (AAMAS 2012), pp. 433-544 440. IFAAMAS, 2012. 546 Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. Go-explore: a 547 new approach for hard-exploration problems, 2021. 548 Ashley D. Edwards, Himanshu Sahni, Yannick Schroecker, and Charles L. Isbell. Imitating latent 549 policies from observation, 2019. URL https://arxiv.org/abs/1805.07914. 550 551 Alejandro Escontrela, Ademi Adeniji, Wilson Yan, Ajay Jain, Xue Bin Peng, Ken Goldberg, Young-552 woon Lee, Danijar Hafner, and Pieter Abbeel. Video prediction models as rewards for reinforce-553 ment learning, 2023. URL https://arxiv.org/abs/2305.14343. 554 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning, 2021. URL https://arxiv.org/abs/2004.07219. 556 Dibya Ghosh, Chethan Anand Bhateja, and Sergey Levine. Reinforcement learning from passive 558 data via latent intentions. In International Conference on Machine Learning, pp. 11321–11339. 559 PMLR, 2023. 560 561 Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Gird-562 har, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, 563 Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano 565 Fragomeni, Qichen Fu, Abrham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang, 566 Yifei Huang, Wenqi Jia, Weslie Khoo, Jachym Kolar, Satwik Kottur, Anurag Kumar, Federico 567 Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan 568 Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanova, 569 Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, 570 Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbelaez, David Cran-571 dall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna 572 Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, 573 Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, 574 Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the world in 3,000 hours of egocentric video, 2022. URL https://arxiv.org/abs/2110.07058. 575 576 Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy 577 learning: Solving long-horizon tasks via imitation and reinforcement learning. arXiv preprint 578 arXiv:1910.11956, 2019. 579 Tuomas Haarnoja, Vitchyr Pong, Aurick Zhou, Murtaza Dalal, Pieter Abbeel, and Sergey Levine. 580 Composable deep reinforcement learning for robotic manipulation, 2018a. 581 582 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy 583 maximum entropy deep reinforcement learning with a stochastic actor, 2018b. 584 585 Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with dis-586 crete world models, 2022. URL https://arxiv.org/abs/2010.02193. Kristian Hartikainen, Xinyang Geng, Tuomas Haarnoja, and Sergey Levine. Dynamical distance 588 learning for semi-supervised and unsupervised skill discovery, 2020. 589 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual 591 networks, 2016. URL https://arxiv.org/abs/1603.05027. 592 Rein Houthooft, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. Vime: Variational information maximizing exploration, 2017.

Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network

- 594 Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng 595 Wu, and Changjie Fan. Learning to utilize shaping rewards: A new approach of reward 596 shaping. ArXiv, abs/2011.02669, 2020. URL https://api.semanticscholar.org/ 597 CorpusID:226254296. 598 Ahmed Hussein, Eyad Elyan, Mohamed Medhat Gaber, and Chrisina Jayne. Deep reward shaping from demonstrations. 2017 International Joint Conference on Neural Networks (IJCNN), pp. 600 510-517, 2017. URL https://api.semanticscholar.org/CorpusID:1744332. 601 602 Yuqian Jiang, Suda Bharadwaj, Bo Wu, Rishi Shah, Ufuk Topcu, and Peter Stone. Temporal-603 logic-based reward shaping for continuing reinforcement learning tasks. In AAAI Conference on 604 Artificial Intelligence, 2020. URL https://api.semanticscholar.org/CorpusID: 605 231845313. 606 607 Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. Machine learning, 49:209-232, 2002. 608 609 J. Zico Kolter and Andrew Y. Ng. Near-bayesian exploration in polynomial time. In Proceedings 610 of the 26th Annual International Conference on Machine Learning, ICML '09, pp. 513–520, 611 New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605585161. doi: 612 10.1145/1553374.1553441. URL https://doi.org/10.1145/1553374.1553441. 613 614 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-615 learning, 2021a. 616 617 Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels, 2021b. URL https://arxiv.org/abs/2004. 618 13649. 619 620 Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online 621 reinforcement learning via balanced replay and pessimistic q-ensemble, 2021. URL https: 622 //arxiv.org/abs/2107.00591. 623 624 Qiyang Li, Jason Zhang, Dibya Ghosh, Amy Zhang, and Sergey Levine. Accelerating exploration 625 with unlabeled prior data, 2023. URL https://arxiv.org/abs/2311.05067. 626 Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy 627 Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training, 628 2023. 629 630 A. Rupam Mahmood, Dmytro Korenkevych, Brent J. Komer, and James Bergstra. Setting up a 631 reinforcement learning task with a real-world robot, 2018. 632 633 Aleksandra Malysheva and Daniel Kudenko. Learning to run with reward shaping from video data. 634 2018. URL https://api.semanticscholar.org/CorpusID:53512677. 635 Maja J Mataric. Reward functions for accelerated learning. In William W. Cohen and 636 Haym Hirsh (eds.), Machine Learning Proceedings 1994, pp. 181-189. Morgan Kauf-637 mann, San Francisco (CA), 1994. ISBN 978-1-55860-335-6. doi: https://doi.org/10.1016/ 638 B978-1-55860-335-6.50030-1. URL https://www.sciencedirect.com/science/ 639 article/pii/B9781558603356500301. 640 641 Lina Mezghani, Sainbayar Sukhbaatar, Piotr Bojanowski, Alessandro Lazaric, and Alahari Karteek. 642 Learning goal-conditioned policies offline with self-supervised reward shaping. In Conference 643 on Robot Learning, 2023. URL https://api.semanticscholar.org/CorpusID: 644 253985364. 645 Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A uni-646 versal visual representation for robot manipulation, 2022. URL https://arxiv.org/abs/ 647
 - 12

2203.12601.

648 649 650 651	Andrew Y. Ng, Daishi Harada, and Stuart J. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In <i>Proceedings of the Sixteenth International Confer-</i> <i>ence on Machine Learning</i> , ICML '99, pp. 278–287, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc. ISBN 1558606122.
652 653 654	Georg Ostrovski, Marc G. Bellemare, Aaron van den Oord, and Remi Munos. Count-based explo- ration with neural density models, 2017.
655 656	Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction, 2017.
657 658	Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement, 2019.
659 660 661	Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. Advances in neural information processing systems, 1, 1988.
662 663	Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models: Model- free deep rl for model-based control, 2020.
664 665 666 667	Jürgen Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). <i>IEEE Transactions on Autonomous Mental Development</i> , 2(3):230–247, 2010. doi: 10.1109/TAMD. 2010.2056368.
668 669	Dominik Schmidt and Minqi Jiang. Learning to act without actions, 2024. URL https: //arxiv.org/abs/2312.10812.
670 671	Rutav Shah and Vikash Kumar. Rrl: Resnet as representation for reinforcement learning, 2021. URL https://arxiv.org/abs/2107.03380.
672 673 674 675	Avi Singh, Albert Yu, Jonathan Yang, Jesse Zhang, Aviral Kumar, and Sergey Levine. Cog: Connecting new skills to past experience with offline reinforcement learning, 2020. URL https://arxiv.org/abs/2010.14500.
676 677	Susanne Still and Doina Precup. An information-theoretic approach to curiosity-driven reinforce- ment learning. <i>Theory in Biosciences</i> , 131:139–148, 2012.
678	Richard S Sutton. Reinforcement learning: An introduction. A Bradford Book, 2018.
679 680 681 682	Haoran Tang, Rein Houthooft, Davis Foote, Adam Stooke, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. exploration: A study of count-based exploration for deep reinforcement learning, 2017.
683 684 685	Ikechukwu Uchendu, Ted Xiao, Yao Lu, Banghua Zhu, Mengyuan Yan, Joséphine Simon, Matthew Bennice, Chuyuan Fu, Cong Ma, Jiantao Jiao, Sergey Levine, and Karol Hausman. Jump-start reinforcement learning, 2023. URL https://arxiv.org/abs/2204.02372.
686 687 688	E. Wiewiora. Potential-based shaping and q-value initialization are equivalent. <i>Journal of Artificial Intelligence Research</i> , 19:205–208, September 2003. ISSN 1076-9757. doi: 10.1613/jair.1190. URL http://dx.doi.org/10.1613/jair.1190.
689 690 691	Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked visual pre-training for motor control, 2022. URL https://arxiv.org/abs/2203.06173.
692 693 694	Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, and Yu Bai. Policy finetuning: Bridging sample-efficient offline and online reinforcement learning, 2022. URL https://arxiv.org/abs/2106.04895.
695 696 697 698	Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Sejune Joo, Jianwei Yang, Baolin Peng, Ajay Man- dlekar, Reuben Tan, Yu-Wei Chao, Bill Yuchen Lin, Lars Liden, Kimin Lee, Jianfeng Gao, Luke Zettlemoyer, Dieter Fox, and Minjoon Seo. Latent action pretraining from videos, 2024. URL https://arxiv.org/abs/2410.11758.
700 701	Han Zheng, Xufang Luo, Pengfei Wei, Xuan Song, Dongsheng Li, and Jing Jiang. Adaptive pol- icy learning for offline-to-online reinforcement learning, 2023. URL https://arxiv.org/ abs/2303.07693.



Hyperparameter	Value
$p_{random goal}$	0.1
$p_{trajgoal}$	0.8
$p_{currgoal}$	0.1
reward_scale	1.0
reward_shift	-1.0
$p_{same a o a l}$	0.5
intent_sametraj	True
Encoder	ResNet-v2
MLP Hidden Dims	[256, 256]
Value Ensemble Size	2
Optimizer Learning Rate	6e-5
Optimizer Epsilon	0.00015
Discount	0.98
Expectile	0.9
Target Update Rate	0.005
Batch Size	64

Figure 8: ICVF Ego4D value loss over training.

Table 1: ICVF Ego4D Training Settings. We include parameters from the ICVF public code base to control the image sampling mechanism.

A APPENDIX

726 A.1 VIVA TRAINING

728 We pre-train ViVa on the Ego4D video dataset. We use the public ICVF codebase and use settings 729 shown in Table 1. We preprocess the video dataset by shaping to 256×256 , center cropping the 730 middle 224×224 , then resizing it to 128×128 . The ICVF itself is structured with an encoder 731 which converts the state, future outcome, and goal into embeddings. For the encoder, we utilize the 732 26-layer ResNet-v2. The training loss is displayed in Figure 8. We train with 1 v4 TPU for 1.5 days.

Once embedded, we concatenate the latents and pass them into an ensemble of 2 Multilayer Perceptrons, each with LayerNorm and to produce the value estimate. We train the ICVF for 1 million steps. For RoboVerse and Franka Kitchen, we apply the same exact training process, but on a the fine-tuning dataset. Antmaze doesn't utilize pretraining and functions on states, so it has a different setup. For our final experiments, we swept across checkpoints to identify strong value functions to run online RL with.

A.2 ONLINE RL

When running online RL with a trained ICVF, we formulate our reward as:

$$\tilde{r}(s,a) = r(s,a) + ICVF(s,g,g) \tag{6}$$

However, we experimented with a different approach where

$$\tilde{r}(s,a,s') = r(s,a) + (\gamma \Phi_g(s') - \Phi_g(s)) \tag{7}$$

$$\Phi_g(s) = ICVF(s, g, g) \tag{8}$$

which follows the potential-based reward shaping strategy as formulated by Ng et al. (1999). They show a learned Q-function under the proposed reward transformation is:

$$\tilde{Q}_g^*(s,a) = Q_g^*(s,a) - \Phi_g(s) \tag{9}$$



Figure 9: Left: An experiment ablating out off-task data and Ego-4D pre-training. As seen, offtask data and off-environment pre-training are sigificant for performance boost. **Right:** Comparison between V-PTR and ViVa on the COG pick-and-place task showing how ViVa's design of reward guidance trumps simple representation transfer on COG.



Figure 10: Left: Comparison in AntMaze between using value potential or pure value as the reward augmentation. Middle: Comparison in AntMaze between agents given access to extrinsic reward labels or not. Right: Comparison between agent given prior data access (RLPD) and agents given 0 prior data (SAC)

which is evidently invariant of actions and thus admits the same optimal policy. Although this is
favorable in theory, in practice we observed no changes in results except variance in policy rollout
returns which could destabilize training. This was tested in Antmaze by training an ICVF on the full
Antmaze dataset and utilizing the potential-based shaping reward versus the simple value-guided
reward as shown in Figure 10.

794 A.3 ANTMAZE

795

767

768

769

770 771

780

781

782

783

796 Value training Our first experiment involves the AntMaze environment specified in the D4RL 797 experiment suite. It is build upon Mujoco and controls an 8 DoF ant with 4 legs through a maze. It 798 starts in the bottom left and is tasked to reach the top right using a sparse reward. In practice, we 799 do not utilize any ICVF Ego4D pretraining since we run this experiment in a state-based fashion. 800 The state is 29-dimensional including positions, velocities, angles, and angular velocities. We use a different ICVF setup for training too. Specifically, we utilize a discount of 0.999, a learning rate 801 of 3e-4 and an epsilon of 1e-8. We use a 3 layer, 512 unit MLP with LayerNorm as the value 802 function. We experimented with using the original multilinear formulation proposed by Ghosh et al. 803 (2023) but noticed early collapse during training as well as noisy values, shown in Figure 11. This 804 motivated our choice to use a single, monolithic neural architecture to represent value. 805

806

RL training We run on the RLPD public codebase and detail RLPD hyperparameters in Table 2.
 RLPD simply runs the Soft Actor-Critic algorithm but adds offline sampling and some extra design choices as detailed in their paper. We edit every update batch reward by adding the ICVF value for the current state conditioned on the goal times 0.001. We use 5 seeds for all baseline experiments.



Figure 11: These plots show trained value function curves across time on a trajectory from s to m to e, representing start, middle, and end, respectively. We denote x as representing the state in the trajectory at a given timestep. Each row is a comparison of values when setting the goal-conditioning to start, middle, or end. As depicted, the monolithic values much more smoothly express distance from the start, middle, or end as we move across the trajectory.

839	Hyperparameter	Value
840	CNN Features	(32, 64, 128, 256)
841	CNN Filters	(3, 3, 3, 3)
842	CNN Strides	(2, 2, 2, 2)
843	CNN Padding	"VALID"
844	CNN Latent Dimension	50
845	Update-to-Data Ratio	1
846	Offline Ratio	0.5
9/7	Start Training	5000
047	Backup Entropy	True
848	Hidden Dims	(256, 256)
849	Batch Size	256
850	Q Ensemble Size	2
851	Temperature LR	3e-4
852	Init Temperature	0.1
853	Actor LR	3e-4
854	Critic LR	3e-4
855	Discount	0.99
856	Tau	0.005
857	Critic Layer Norm	True
	Horizon	40

Hyperparameter	Value
Update-to-Data Ratio	20
Offline Ratio	0.5
Start Training	5000
Backup Entropy	False
Hidden Dims	(256, 256, 256)
Q Ensemble Size	1
Temperature LR	3e-4
Init Temperature	1.0
Actor LR	3e-4
Critic LR	3e-4
Discount	0.99
Tau	0.005
Critic Layer Norm	True
Horizon	1000

Figure 13: RLPD Settings for Antmaze

Figure 12: RLPD Settings for COG RoboVerse and FrankaKitchen

864	Method	AntMaze Corrupt	COG Pick-Place	Franka Hinge	Franka Slide
865	ViVa	N/A	8.73 ± 13.48	30.38 ± 32.07	62.64 ± 12.85
866	ViVa (No Ego4D)	0.89 ± 0.14	6.42 ± 11.25	14.91 ± 25.20	54.75 ± 14.98
867	JSRL	0.95 ± 0.04	0 ± 0	0 ± 0	1.68 ± 3.76
868	DreamerV2	N/A	0 ± 0	0 ± 0	15.11 ± 22.36
869	RLPD	0 ± 0	0 ± 0	21.06 ± 23.64	54.38 ± 19.71
870	SAC	0 ± 0	0.01 ± 0.03	0 ± 0	0.44 ± 0.98
871	ViVa	N/A	16.71 ± 16.73	47.56 ± 33.64	75.61 ± 11.46
872	ViVa (No Ego4D)	0.9 ± 0.13	8.42 ± 14.58	47.15 ± 33.36	74.61 ± 15.06
072	JSRL	0.98 ± 0.01	0 ± 0	8.7 ± 23.00	3.06 ± 6.84
073	DreamerV2	N/A	0 ± 0	0 ± 0	25.32 ± 30.71
874	RLPD	0 ± 0	0 ± 0	41.704 ± 33.40	72.47 ± 21.44
875	SAC	0 ± 0	0.02 ± 0.03	0 ± 0	0.80 ± 1.79
876					

Table 2: Experimental Suite Results. The top set of results are at the halfway point through the online RL training process. The bottom rows are metrics at the final step.

878 879 880

881

877

A.4 COG ROBOVERSE

882 We use the RoboVerse simulator, publicly located here, which simulates a WidowX 883 We use the datasets created in the COG paper which is pubrobot through PyBullet. 884 We run experiments on the pick-and-place task which sparsely relicly located here. 885 wards the agent for picking up a target object randomly placed on a table and placing 886 For ICVF fine-tuning, we utilize a number of data combinations it into a silver tray. 887 for different experiments detailed in the paper, but select from the main group of COG datasets: pickplace_prior, pickplace_task, DrawerOpenGrasp, drawer_task, closed_drawer_prior, blocked_drawer_1_prior, and blocked_drawer_2_prior. 889 We only include pickplace_task in the scaling law and elect to remove it for all other ex-890 periments. 891

892 During the online RL phase, we adopt the same RLPD system but we use the DrQ regularization 893 methods for image-based RL. Specifically, we utilize the D4PG (Barth-Maron et al., 2018) visual encoder. We attach experimental hyperparameters in Table 12 and use 8 seeds each. We additionally 894 compare our method to V-PTR, a similar method which uses the trained representations from the 895 ICVF rather than the actual value network outputs. V-PTR uses the trained encoders to map the 896 observations into an embedding space for the policy network to learn on. Since our method uses 897 the notion of distance itself and more actively enforces this signal directly into the reward, we 898 hypothesize it'd be more directly useful for sparse reward RL. Our comparison results are in Figure 899 9 which motivate our decision to use values directly. We additionally ablate off-task data and Ego4D 900 pre-training to show the effect of each data source in Figure 9. 901

901 902 903

A.5 FRANKA KITCHEN

904 Our final experiment uses the Franka Kitchen environment available on D4RL here which simu-905 lates a 9-DoF Franka Robot in a kitchen environment. We control the robot in joint velocity mode 906 clipped between -1 and 1 rad/s. The 9 degrees of freedom are 7 joints and 2 fingers of the gripper. We analyze two tasks which are opening the sliding cabinet and opening the hinge cabinet. 907 These tasks are specified with a sparse reward. As mentioned in the paper, the datasets we used 908 contain failed interactions with the target objects. We collect this data by controlling the robot with 909 expert demonstration actions with added Gaussian noise. We then filter out all successes from this 910 data to form our dataset. The hinge failures dataset contains 1013 trajectories, whereas the sliding 911 door dataset contains 630 trajectories. These trajectories are 50 steps each. The RLPD settings for 912 FrankaKitchen are the exact same as for RoboVerse but we use a horizon of 50 steps rather than 40, 913 and we run 6 seeds per baseline.

- 914
- 915
- 916