

---

# Towards Universal Visual Reward and Representation via Value-Implicit Pre-Training

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We introduce **Value-Implicit Pre-training (VIP)**, a self-supervised pre-trained  
2 visual representation capable of generating dense and smooth reward functions  
3 for unseen robotic tasks. VIP casts representation learning from human videos  
4 as an *offline goal-conditioned reinforcement learning* problem and derives a self-  
5 supervised *dual* goal-conditioned value-function objective that does not depend  
6 on actions, enabling pre-training on unlabeled human videos. Theoretically, VIP  
7 can be understood as a novel *implicit* time contrastive learning that makes for  
8 temporally smooth embedding that enables the value function to be implicitly  
9 defined via the embedding distance, which can be used as the reward function for  
10 any downstream task specified through goal images. Trained on large-scale Ego4D  
11 human videos and without any fine-tuning on task-specific robot data, VIP’s frozen  
12 representation can provide dense visual reward for an extensive set of simulated and  
13 **real-robot** tasks, enabling diverse reward-based policy learning methods, including  
14 visual trajectory optimization and online/offline RL, and significantly outperform  
15 all prior pre-trained representations. Notably, VIP can enable *few-shot* offline RL  
16 on a suite of real-world robot tasks with as few as 20 trajectories. Project website:  
17 <https://sites.google.com/view/rl-vip>

## 18 1 Value-Implicit Pre-Training

19 Due to space limit, we provide the full version of this section in Appendix D.

### 20 1.1 Foundation: Self-Supervised Value Learning from Human Videos

21 While human videos are out-of-domain data for robots, they are *in-domain* for learning a goal-  
22 conditioned human policy. Given that human videos naturally contain goal-directed behavior, one  
23 reasonable idea of utilizing offline human videos for representation learning is to solve an offline  
24 goal-conditioned RL problem over the space of human policies and then extract the learned visual  
25 representation. However, this idea is seemingly implausible because the offline human dataset does  
26 not come with any action labels that are typically required for *policy* learning. Our key insight is that,  
27 for a suitable choice of offline policy optimization problem, we can solve for the *dual* value learning  
28 problem that does not depend on any action label in the offline dataset. In particular, leveraging the  
29 idea of Fenchel duality (Rockafellar, 1970) from convex optimization, we have the following result:

30 **Proposition 1.1.** *Under assumption of deterministic transition dynamics, the dual optimization*  
31 *problem of (11) is*

$$\max_{\phi} \min_V \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [V(\phi(o); \phi(g))] + \log \mathbb{E}_{(o,o';g) \sim D} [\exp(r(o, g) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g)))] \right], \quad (1)$$

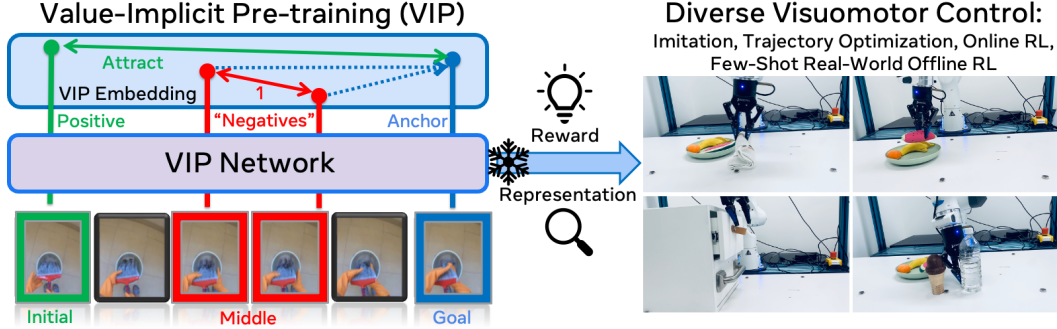


Figure 1: **Value-Implicit Pre-training (VIP)**. Pre-trained on large-scale, in-the-wild human videos, frozen VIP network can provide visual reward and representation for downstream robotics tasks and enable diverse visuomotor control strategies without any task-specific fine-tuning.

32 where  $\mu_0(o; g)$  is the goal-conditioned initial observation distribution, and  $D(o, o'; g)$  is the goal-  
 33 conditioned distribution of two consecutive observations in dataset  $D$ .

34 As shown, actions do not appear in the objective. Furthermore, since all expectations in (12) can be  
 35 sampled using the offline dataset, this dual value-function objective can be self-supervised with an  
 36 appropriate choice of reward function. In particular, since our goal is to acquire a value function that  
 37 extracts a general notion of goal-directed task progress from passive offline human videos, we set  
 38  $r(o, g) = \mathbb{I}(o == g) - 1$ , which we refer to as  $\tilde{\delta}_g(o)$  in shorthand. This reward provides a constant  
 39 negative reward when  $o$  is not the provided goal  $g$ , and does not require any task-specific engineering.  
 40 The resulting value function  $V(\phi(o); \phi(g))$  captures the discounted total number of steps required  
 41 to reach goal  $g$  from observation  $o$ , and will objective will encourage learning visual features  $\phi$  that  
 42 are amenable to predicting the discounted temporal distance between two frames in a human video  
 43 sequence. With enough size and diversity in the training dataset, we hypothesize that this value  
 44 function can generalize to completely unseen (robot) domains.

## 45 1.2 Analysis: Implicit Time Contrastive Learning

46 In this section, we show that (1) can be understood as a novel *implicit* temporal contrastive rep-  
 47 resentation learning that acquires temporally smooth embedding distance over video sequences,  
 48 underpinning VIP’s efficacy jointly as a visual representation and reward for downstream control.

49 Assuming that the optimal  $V^*$  is found in (1), with a few algebraic manipulation steps (see Appendix E  
 50 for a derivation), we can massage (13) into an expression that resembles the InfoNCE (Oord et al.,  
 51 2018) time contrastive learning (Sermanet et al., 2018) (see Appendix B.2 for a definition and  
 52 additional background) objective:

$$\min_{\phi} (1 - \gamma) \mathbb{E}_{p(g), \mu_0(o; g)} \left[ - \log \frac{e^{V^*(\phi(o); \phi(g))}}{\mathbb{E}_{D(o, o'; g)} [\exp(\tilde{\delta}_g(o) + \gamma V^*(\phi(o'); \phi(g)) - V^*(\phi(o), \phi(g)))]^{\frac{-1}{1-\gamma}}} \right] \quad (2)$$

53 In particular,  $p(g)$  can be thought of the distribution of “anchor” observations,  $\mu_0(s; g)$  the distribution  
 54 of “positives” samples, and  $D(o, o'; g)$  the distribution of “negatives” samples. Since the value  
 55 function encodes negative discounted temporal distance, due to the recursive nature of value temporal-  
 56 difference (TD), in order for the one-step TD error to be globally minimized along a video sequence,  
 57 observations that are temporally farther away from the goal will naturally be repelled farther away in  
 58 the representation space compared to observations that are nearby in time. Therefore, the repulsion  
 59 of the negative observations is an *implicit*, emergent property from the optimization of (2), instead of  
 60 an explicit constraint as in standard (time) contrastive learning. In Appendix D, we detail how this  
 61 implicit time contrast mechanism gives rise to a temporally smooth visual representation that makes  
 62 for effective zero-shot reward-specification.

### 63 1.3 Algorithm: Value-Implicit Pre-Training (VIP)

64 Recall that  $V^*$  is assumed to be known for the derivation in Section 1.2, but in practice, its analytical  
 65 form is rarely known. Now, given that  $V^*$  plays the role of a distance measure in our implicit  
 66 time contrastive learning framework, a simple and intuitive way to approximate  $V^*$  in practice is to  
 67 *implicitly* parameterize it to be a choice of distance measure. In this work, we choose the common  
 68 choice of the negative  $L_2$  distance used in prior work Sermanet et al. (2018); Nair et al. (2022):  
 69  $V^*(\phi(o), \phi(g)) := -\|\phi(o) - \phi(g)\|_2$ . Altogether, VIP training is illustrated in Alg. 2; it is simple  
 70 and its core training loop can be implemented in fewer than 10 lines of PyTorch code (Alg. 3).

---

#### Algorithm 1 Value-Implicit Pre-Training (VIP)

---

- 1: **Require:** Offline (human) videos  $D = \{(o_1^i, \dots, o_{i_h}^i)\}_{i=1}^N$ , visual architecture  $\phi$
  - 2: **for** number of training iterations **do**
  - 3:   Sample sub-trajectories  $\{o_t^i, \dots, o_k^i, o_{k+1}^i, \dots, o_T^i\}_{i=1}^B \sim D, t \in [1, i_h - 1], t \leq k < T, T \in (t, i_h], \forall i$
  - 4:    $\mathcal{L}(\phi) := \frac{1-\gamma}{B} \sum_{i=1}^B [\|\phi(o_t^i) - \phi(o_T^i)\|_2] + \log \frac{1}{B} \sum_{i=1}^B [\exp(\|\phi(o_k^i) - \phi(o_T^i)\|_2 - \delta_{o_t^i}(o_k^i) - \gamma \|\phi(o_{k+1}^i) - \phi(o_T^i)\|_2)]$
  - 5:   Update  $\phi$  using SGD:  $\phi \leftarrow \phi - \alpha_\phi \nabla \mathcal{L}(\phi)$
- 

## 71 2 Experiments

72 In this section, we demonstrate VIP’s effectiveness as both a pre-trained  
 73 visual reward and representation on three distinct reward-based policy  
 74 learning settings. Due to space limit,  
 75 we delve into results directly, and  
 76 all omitted experimental details are  
 77 contained in App. G; additional re-  
 78 sults and analysis are presented in  
 79 App.I. At a high level, VIP fixes the  
 80 visual architecture (ResNet50) and  
 81 pre-training dataset (Ego4D) as a  
 82 state-of-art pre-trained representation  
 83 R3M (Nair et al., 2022), differing pri-  
 84 marily in the training objective. We use FrankaKitchen (Gupta et al., 2019) for evaluation. Each task  
 85 is specified via only a goal image, requiring the pre-trained representations to provide embedding-  
 86 distance based reward (4) and visual encoding.  
 87  
 88

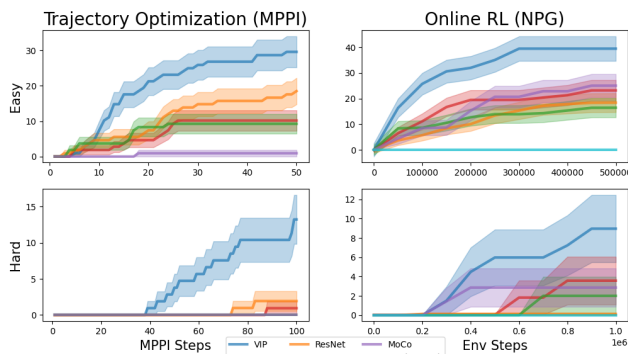


Figure 2: Visual traj. opt. and RL results (max success rate %).

### 89 2.1 Trajectory Optimization & Online Reinforcement Learning

90 We evaluate pre-trained representations’ capability as pure visual reward functions by using them  
 91 to directly synthesize a sequence of actions using a standard trajectory optimization algorithm. We  
 92 also evaluate online RL, which provides improved exploration but comes with the added challenge of  
 93 demanding the pre-trained representation to provide both the visual reward and representation for  
 94 learning a closed-loop policy. In Figure 2, we report each representation’s cumulative success rate  
 95 averaged over task configurations and random seeds (3 seeds \* 3 cameras \* 12 tasks = 108 runs).

96 Examining the MPPI results, we see that VIP is substantially better than all baselines in both Easy  
 97 and Hard settings, and is the only representation that makes non-trivial progress on the Hard setting.  
 98 These results demonstrate that VIP has superior capability as a pure visual reward function. In  
 99 Fig. 3, we couple VIP and the strongest baselines (R3M, Resnet)’s with increasingly powerful MPPI  
 100 optimizers (i.e., more trajectories per optimization step). As shown, while VIP steadily benefits from  
 101 stronger optimizers and can reach an average success rate of **44%**, baselines often do *worse* when  
 102 MPPI becomes more powerful, suggesting that their reward landscapes are filled with local minima  
 103 that do not correlate with task progress and are easily exploited by (stronger) optimizers.

104 Switching gear to online RL, VIP again achieves consistently superior performance, demonstrating  
 105 its joint effectiveness as visual reward and representation. VIP (Sparse)’s inability to solve any

Table 1: Real-robot offline RL results (success rate % averaged over 10 rollouts with standard deviation reported).

Environment	Pre-Trained				In-Domain		
	VIP-RWR	VIP-BC	R3M-RWR	R3M-BC	Scratch-BC	VIP-RWR	VIP-BC
CloseDrawer	100 ± 0	50 ± 50	80 ± 40	10 ± 30	30 ± 46	0 ± 0	0* ± 0
PushBottle	90 ± 30	50 ± 50	70 ± 46	50 ± 50	40 ± 48	0* ± 0	0* ± 0
PlaceMelon	60 ± 48	10 ± 30	0 ± 0	0 ± 0	0 ± 0	0* ± 0	0* ± 0
FoldTowel	90 ± 30	20 ± 40	0 ± 0	0 ± 0	0 ± 0	0* ± 0	0* ± 0

task indicates the necessity of dense reward in solving these challenging visual manipulation tasks. Whereas sparse reward still requires human engineering via installing additional sensors (Rajeswar et al., 2021; Singh et al., 2019) and faces exploration challenges (Nair et al., 2018), with VIP, the end-user has to provide only a goal image, and, without any additional state or reward instrumentation, can expect a significant improvement in performance.

## 2.2 Real-World Few-Shot Offline Reinforcement Learning

Finally, we demonstrate how VIP’s reward and representation can power a simple and practical system for real-world robot learning in the form of *few-shot* offline reinforcement learning, making offline RL simple, sample-efficient, and more effective than BC with almost no added complexity.

To this end, we consider a simple reward-weighted regression (RWR) (Peters & Schaal, 2007; Peng et al., 2019) approach, in which the reward and the encoder are provided by the pre-trained model  $\phi$ :

$$\mathcal{L}(\pi) = -\mathbb{E}_{D_{\text{task}}} [\exp(\tau \cdot R(o, o'; \phi, g)) \log \pi(a | \phi(o))], \quad (3)$$

where  $R$  is defined via (4) and  $\tau$  is the temperature scale. Compared to BC, which would be (3) with uniform weights to all transitions, RWR can focus policy learning on transitions that have high rewards (i.e., high task progress) under the deployed representation.

We introduce 4 tabletop manipulation tasks (see Figure 1 and Figure 10) requiring a real 7-DOF Franka robot to manipulate objects drawn from distinct categories of objects. For each task, we collect in-domain, task-specific offline data  $D_{\text{task}}$  of  $\sim 20$  demonstrations with randomized object initial placements for policy learning; we provide detailed task and experiment descriptions in Appendix H.

The average success rate (%) and standard deviation across 10 test rollouts are reported in Table 1. As shown, VIP-RWR improves upon VIP-BC on all tasks and provides substantial benefit in the harder tasks that are multi-stage in nature. In contrast, R3M-RWR, while able to improve R3M-BC on the simpler two tasks involving pushing an object, fails to make any progress on the harder tasks. The low performance of BC-based methods on the harder `PickPlaceMelon` and `FoldTowel` tasks indicates that in this low-data regime, regardless of the quality of visual representation, good reward information is necessary for task success. Finally, *in-domain* methods all fail in this low-data regime. Altogether, these results corroborate the necessity of pre-training in achieving real-world few-shot offline RL and highlight the unique effectiveness of VIP in realizing this goal.

## 3 Conclusion

We have proposed Value-Implicit Pre-training (VIP), a self-supervised value-based pre-training objective that is highly effective in providing both the visual reward and representation for downstream unseen robotics tasks. VIP is derived from first principles of dual reinforcement learning and admits an appealing connection to an implicit and more powerful formulation of time contrastive learning, which captures long-range temporal dependency and injects local temporal smoothness in the representation to make for effective zero-shot reward specification. Trained entirely on diverse, in-the-wild human videos, VIP demonstrates significant gains over state-of-art pre-trained representations on an extensive set of policy learning settings. Notably, VIP can enable simple and sample-efficient real-world offline RL with just handful of trajectories. Altogether, we believe that VIP makes an important contribution in both the algorithmic frontier of visual pre-training for RL and practical real-world robot learning.

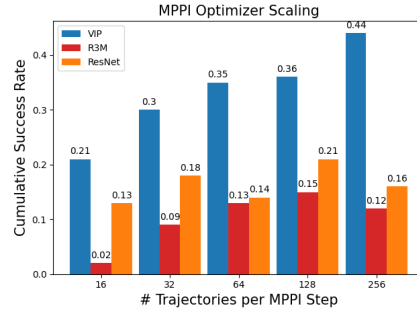


Figure 3: VIP benefits from scaling compute for downstream trajectory optimization.

148 **References**

- 149 Alekh Agarwal, Nan Jiang, and Sham M Kakade. Reinforcement learning: Theory and algorithms.  
150 2019.
- 151 Anish Agarwal, Abdullah Alomar, Varkey Alumootil, Devavrat Shah, Dennis Shen, Zhi Xu, and  
152 Cindy Yang. Persim: Data-efficient offline reinforcement learning with heterogeneous agents via  
153 personalized simulators. *Advances in Neural Information Processing Systems*, 34:18564–18576,  
154 2021.
- 155 Shikhar Bahl, Abhinav Gupta, and Deepak Pathak. Human-to-robot imitation in the wild. *arXiv*  
156 *preprint arXiv:2207.09450*, 2022.
- 157 Steven J Bradtke and Andrew G Barto. Linear least-squares algorithms for temporal difference  
158 learning. *Machine learning*, 22(1):33–57, 1996.
- 159 Yevgen Chebotar, Karol Hausman, Yao Lu, Ted Xiao, Dmitry Kalashnikov, Jake Varley, Alex Irpan,  
160 Benjamin Eysenbach, Ryan Julian, Chelsea Finn, et al. Actionable models: Unsupervised offline  
161 reinforcement learning of robotic skills. *arXiv preprint arXiv:2104.07749*, 2021.
- 162 Annie S Chen, Suraj Nair, and Chelsea Finn. Learning generalizable robotic reward functions from "  
163 in-the-wild" human videos. *arXiv preprint arXiv:2103.16817*, 2021.
- 164 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
165 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
166 pp. 248–255. Ieee, 2009.
- 167 Frederik Ebert, Yanlai Yang, Karl Schmeckpeper, Bernadette Bucher, Georgios Georgakis, Kostas  
168 Daniilidis, Chelsea Finn, and Sergey Levine. Bridge data: Boosting generalization of robotic skills  
169 with cross-domain datasets. *arXiv preprint arXiv:2109.13396*, 2021.
- 170 Benjamin Eysenbach, Tianjun Zhang, Ruslan Salakhutdinov, and Sergey Levine. Contrastive learning  
171 as goal-conditioned reinforcement learning. *arXiv preprint arXiv:2206.07568*, 2022.
- 172 Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzyńska, Susanne Westphal,  
173 Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, Florian Hoppe,  
174 Christian Thureau, Ingo Bax, and Roland Memisevic. The "something something" video database  
175 for learning and evaluating visual common sense, 2017.
- 176 Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy  
177 learning: Solving long-horizon tasks via imitation and reinforcement learning. *arXiv preprint*  
178 *arXiv:1910.11956*, 2019.
- 179 Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle  
180 for unnormalized statistical models. In *Proceedings of the thirteenth international conference on*  
181 *artificial intelligence and statistics*, pp. 297–304. JMLR Workshop and Conference Proceedings,  
182 2010.
- 183 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
184 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
185 pp. 770–778, 2016.
- 186 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
187 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on*  
188 *computer vision and pattern recognition*, pp. 9729–9738, 2020.
- 189 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
190 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*  
191 *Computer Vision and Pattern Recognition*, pp. 16000–16009, 2022.

- 192 Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine,  
193 and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In  
194 *Conference on Robot Learning*, pp. 991–1002. PMLR, 2022.
- 195 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
196 *arXiv:1412.6980*, 2014.
- 197 Aviral Kumar, Anikait Singh, Stephen Tian, Chelsea Finn, and Sergey Levine. A workflow for offline  
198 model-free robotic reinforcement learning. *arXiv preprint arXiv:2109.10813*, 2021.
- 199 Aviral Kumar, Joey Hong, Anikait Singh, and Sergey Levine. When should we prefer offline  
200 reinforcement learning over behavioral cloning? *arXiv preprint arXiv:2204.05618*, 2022.
- 201 Youngwoon Lee, Andrew Szot, Shao-Hua Sun, and Joseph J Lim. Generalizable imitation learning  
202 from observation via inferring goal proximity. In A. Beygelzimer, Y. Dauphin, P. Liang, and  
203 J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL  
204 <https://openreview.net/forum?id=lp9fo08AFoD>.
- 205 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,  
206 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 207 Yunfei Li, Tian Gao, Jiaqi Yang, Huazhe Xu, and Yi Wu. Phasic self-imitative reduction for sparse-  
208 reward goal-conditioned reinforcement learning. In *International Conference on Machine Learning*,  
209 pp. 12765–12781. PMLR, 2022.
- 210 Yecheng Jason Ma, Andrew Shen, Dinesh Jayaraman, and Osbert Bastani. Smoldice: Versatile offline  
211 imitation learning via state occupancy matching. *arXiv preprint arXiv:2202.02433*, 2022a.
- 212 Yecheng Jason Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. How far i’ll go: Offline goal-  
213 conditioned reinforcement learning via  $f$ -advantage regression. *arXiv preprint arXiv:2206.03023*,  
214 2022b.
- 215 Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao,  
216 John Emmons, Anchit Gupta, Emre Orbay, et al. Roboturk: A crowdsourcing platform for robotic  
217 skill learning through imitation. In *Conference on Robot Learning*, pp. 879–893. PMLR, 2018.
- 218 Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-  
219 Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline  
220 human demonstrations for robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.
- 221 Ofir Nachum and Bo Dai. Reinforcement learning via fenchel-rockafellar duality, 2020.
- 222 Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Algaedice:  
223 Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.
- 224 Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Over-  
225 coming exploration in reinforcement learning with demonstrations. In *2018 IEEE international*  
226 *conference on robotics and automation (ICRA)*, pp. 6292–6299. IEEE, 2018.
- 227 Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal  
228 visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- 229 Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations:  
230 Theory and application to reward shaping. In *Icml*, volume 99, pp. 278–287, 1999.
- 231 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive  
232 coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 233 Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, and Abhinav Gupta. The unsurprising  
234 effectiveness of pre-trained vision models for control. *arXiv preprint arXiv:2203.03580*, 2022.

- 235 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
236 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,  
237 high-performance deep learning library. *Advances in neural information processing systems*, 32,  
238 2019.
- 239 Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:  
240 Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 241 Jan Peters and Stefan Schaal. Reinforcement learning by reward-weighted regression for operational  
242 space control. In *Proceedings of the 24th international conference on Machine learning*, pp.  
243 745–750, 2007.
- 244 Sai Rajeswar, Cyril Ibrahim, Nitin Surya, Florian Golemo, David Vazquez, Aaron Courville, and  
245 Pedro O. Pinheiro. Haptics-based curiosity for sparse-reward tasks. In *5th Annual Conference on*  
246 *Robot Learning*, 2021. URL <https://openreview.net/forum?id=VfGkOELQ4LC>.
- 247 Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel  
248 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement  
249 learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- 250 R. Tyrrell Rockafellar. *Convex analysis*. Princeton Mathematical Series. Princeton University Press,  
251 Princeton, N. J., 1970.
- 252 Karl Schmeckpeper, Oleh Rybkin, Kostas Daniilidis, Sergey Levine, and Chelsea Finn. Reinforce-  
253 ment learning with videos: Combining offline observations with interaction. *arXiv preprint*  
254 *arXiv:2011.06507*, 2020.
- 255 Pierre Sermanet, Kelvin Xu, and Sergey Levine. Unsupervised perceptual rewards for imitation  
256 learning. *arXiv preprint arXiv:1612.06699*, 2016.
- 257 Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey  
258 Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In  
259 *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 1134–1141. IEEE,  
260 2018.
- 261 Rutav Shah and Vikash Kumar. Rrl: Resnet as representation for reinforcement learning. *arXiv*  
262 *preprint arXiv:2107.03380*, 2021.
- 263 Dandan Shan, Jiaqi Geng, Michelle Shu, and David Fouhey. Understanding human hands in contact  
264 at internet scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
265 *Recognition*, 2020.
- 266 Avi Singh, Larry Yang, Kristian Hartikainen, Chelsea Finn, and Sergey Levine. End-to-end robotic  
267 reinforcement learning without reward engineering. *arXiv preprint arXiv:1904.07854*, 2019.
- 268 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. 2018.
- 269 Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked visual pre-training for  
270 motor control. *arXiv preprint arXiv:2203.06173*, 2022.
- 271 Haoyu Xiong, Quanzhou Li, Yun-Chun Chen, Homanga Bharadhwaj, Samarth Sinha, and Animesh  
272 Garg. Learning by watching: Physical imitation of manipulation skills from human videos. In *2021*  
273 *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 7827–7834.  
274 IEEE, 2021.
- 275 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey  
276 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.  
277 In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.

- 278 Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Chelsea Finn, and Sergey Levine. How  
279 to leverage unlabeled data in offline reinforcement learning. *arXiv preprint arXiv:2202.01741*,  
280 2022.
- 281 Kevin Zakka, Andy Zeng, Pete Florence, Jonathan Tompson, Jeannette Bohg, and Debidatta Dwibedi.  
282 Xirl: Cross-embodiment inverse reinforcement learning. In *Conference on Robot Learning*, pp.  
283 537–546. PMLR, 2022.
- 284 Siyuan Zhang and Nan Jiang. Towards hyperparameter-free policy selection for offline reinforcement  
285 learning. *Advances in Neural Information Processing Systems*, 34:12864–12875, 2021.



286 **Part I**

287 **Appendix**

288 **Table of Contents**  
289

---

290	<b>A Problem Setting and Background</b>	<b>10</b>
291	A.1 Out-of-Domain Pre-Training Visual Representation . . . . .	10
292	A.2 Representation Evaluation . . . . .	10
293	<b>B Additional Background</b>	<b>10</b>
294	B.1 Goal-Conditioned Reinforcement Learning . . . . .	10
295	B.2 InfoNCE & Time Contrastive Learning. . . . .	11
296	<b>C Related Work</b>	<b>11</b>
297	<b>D Value-Implicit Pre-Training (Full-Version)</b>	<b>12</b>
298	D.1 Foundation: Self-Supervised Value Learning from Human Videos . . . . .	13
299	D.2 Analysis: Implicit Time Contrastive Learning . . . . .	13
300	D.3 Algorithm: Value-Implicit Pre-Training (VIP) . . . . .	14
301	<b>E Technical Derivations and Proofs</b>	<b>15</b>
302	E.1 Proof of Proposition D.1 . . . . .	15
303	E.2 VIP Implicit Time Contrast Learning Derivation . . . . .	16
304	E.3 VIP Implicit Repulsion . . . . .	16
305	<b>F VIP Training Details</b>	<b>17</b>
306	F.1 Dataset Processing and Sampling . . . . .	17
307	F.2 VIP Hyperparameters . . . . .	17
308	F.3 VIP Pytorch Pseudocode . . . . .	17
309	<b>G Simulation Experiment Details.</b>	<b>18</b>
310	G.1 FrankaKitchen Task Descriptions . . . . .	18
311	G.2 In-Domain Representation Probing . . . . .	18
312	G.3 Trajectory Optimization . . . . .	19
313	G.4 Reinforcement Learning . . . . .	21
314	<b>H Real-World Robot Experiment Details</b>	<b>21</b>
315	H.1 Task Descriptions . . . . .	21
316	H.2 Training and Evaluation Details . . . . .	22
317	H.3 Additional Analysis & Context . . . . .	22
318	H.4 Qualitative Analysis . . . . .	23
319	<b>I Additional Results</b>	<b>23</b>
320	I.1 Value-Based Pre-Training Ablation: Least-Square Temporal-Difference . . . . .	23
321	I.2 Visual Imitation Learning . . . . .	24
322	I.3 Embedding and True Rewards Correlation . . . . .	25
323	I.4 Embedding Distance Curves . . . . .	25
324	I.5 Embedding Distance Curve Bumps . . . . .	25
325	I.6 Embedding Reward Histograms (Real-Robot Dataset) . . . . .	29
326	I.7 Embedding Reward Histograms (Ego4D) . . . . .	29

327  
328  
329

## 330 A Problem Setting and Background

331 In this section, we describe our problem setting of out-of-domain pre-training and provide formalism  
 332 for downstream representation evaluation. Additional background on goal-conditioned reinforcement  
 333 learning and contrastive learning is included in Appendix B.

### 334 A.1 Out-of-Domain Pre-Training Visual Representation

335 We consider the problem setting of pre-training a frozen visual encoder for downstream control  
 336 tasks (Shah & Kumar, 2021; Parisi et al., 2022; Nair et al., 2022). More specifically, we have access  
 337 to a training set of video data  $D = \{v_i := (o_1^i, \dots, o_{i_h}^i)\}_{i=1}^N$ , where each  $o \in \mathbb{R}^{H \times W \times 3}$  is a raw RGB  
 338 image; note that this formalism also captures standard image datasets (e.g., ImageNet), if we take  
 339  $i_h = 1$  for all  $v_i$ . Like prior works, we assume  $D$  to be out-of-domain and does not include any robot  
 340 task or domain-specific data. A learning algorithm  $\mathcal{A}$  ingests this training data and outputs a visual  
 341 encoder  $\phi := \mathcal{A}(D) : \mathbb{R}^{H \times W \times 3} \rightarrow K$ , where  $K$  is the embedding dimension.

### 342 A.2 Representation Evaluation

343 Given a choice of representation  $\phi$ , every evaluation task can be instantiated as a Markov decision  
 344 process  $\mathcal{M}(\phi) := (\phi(O), A, R(o_t, o_{t+1}; \phi, g), T, \gamma, g)$ , in which the state space is the induced space  
 345 of observation embeddings, and the task is specified via a (set of) goal image(s)  $g$ . Specifically,  
 346 for a given transition tuple  $(o_t, o_{t+1})$ , we define the reward to be the goal-embedding distance  
 347 difference (Lee et al., 2021; Li et al., 2022):

$$R(o_t, o_{t+1}; \phi, \{g\}) := \mathcal{S}_\phi(o_{t+1}; g) - \mathcal{S}_\phi(o_t; g) := (1 - \gamma)\mathcal{S}_\phi(o_{t+1}; g) + (\gamma\mathcal{S}_\phi(o_{t+1}; g) - \mathcal{S}_\phi(o_t; g)), \quad (4)$$

348 where  $\mathcal{S}_\phi$  is a choice of distance function in the  $\phi$ -representation space; in this work, we set  
 349  $\mathcal{S}_\phi(o_t; g) := -\|\phi(o_t) - \phi(g)\|_2$ . This reward function can be interpreted as a raw embedding  
 350 distance reward with a reward shaping (Ng et al., 1999) term that encourages making progress  
 351 towards the goal. This preserves optimal policy but enables more efficient and robust policy learning.

352 Under this formalism, parameters of  $\phi$  are frozen during policy learning (it is considered a part of the  
 353 MDP), and we want to learn a policy  $\pi : \mathbb{R}^K \rightarrow A$  that outputs an action based on the embedded  
 354 observation  $a \sim \pi(\phi(o))$ .

## 355 B Additional Background

### 356 B.1 Goal-Conditioned Reinforcement Learning

357 This section is adapted from Ma et al. (2022b). We consider a goal-conditioned Markov decision  
 358 process from visual state space:  $\mathcal{M} = (O, A, G, r, T, \mu_0, \gamma)$  with state space  $O$ , action space  $A$ ,  
 359 reward  $r(o, g)$ , transition function  $o' \sim T(o, a)$ , the goal distribution  $p(g)$ , and the goal-conditioned  
 360 initial state distribution  $\mu_0(o; g)$ , and discount factor  $\gamma \in (0, 1]$ . We assume the state space  $O$  and  
 361 the goal space  $G$  to be defined over RGB images. The objective of goal-conditioned RL is to find a  
 362 goal-conditioned policy  $\pi : O \times G \rightarrow \Delta(A)$  that maximizes the discounted cumulative return:

$$J(\pi) := \mathbb{E}_{p(g), \mu_0(o; g), \pi(a_t | s_t, g), T(o_{t+1}, | o_t, a_t)} \left[ \sum_{t=0}^{\infty} \gamma^t r(o_t; g) \right] \quad (5)$$

363 The *goal-conditioned* state-action occupancy distribution  $d^\pi(o, a; g) : O \times A \times G \rightarrow [0, 1]$  of  $\pi$  is

$$d^\pi(o, a; g) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(o_t = o, a_t = a \mid o_0 \sim \mu_0(o; g), a_t \sim \pi(o_t; g), o_{t+1} \sim T(o_t, a_t)) \quad (6)$$

364 which captures the goal-conditioned visitation frequency of state-action pairs for policy  $\pi$ . The  
 365 state-occupancy distribution then marginalizes over actions:  $d^\pi(o; g) = \sum_a d^\pi(o, a; g)$ . Then, it

366 follows that  $\pi(a | o, g) = \frac{d^\pi(o, a; g)}{d^\pi(o; g)}$ . A state-action occupancy distribution must satisfy the *Bellman*  
 367 *flow constraint* in order for it to be an occupancy distribution for some stationary policy  $\pi$ :

$$\sum_a d(o, a; g) = (1 - \gamma)\mu_0(o; g) + \gamma \sum_{\tilde{o}, \tilde{a}} T(s | \tilde{o}, \tilde{a})d(\tilde{o}, \tilde{a}; g), \quad \forall o \in O, g \in G \quad (7)$$

368 We write  $d^\pi(o, g) = p(g)d^\pi(o; g)$  as the joint goal-state density induced by  $p(g)$  and the policy  $\pi$ .  
 369 Finally, given  $d^\pi$ , we can express the objective function (5) as  $J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{(o, g) \sim d^\pi(o, g)} [r(o, g)]$ .

## 370 B.2 InfoNCE & Time Contrastive Learning.

371 As VIP can be understood as a implicit and smooth time contrastive learning objective, we provide ad-  
 372 ditional background on the InfoNCE Oord et al. (2018) and time contrastive learning (TCN) (Sermanet  
 373 et al., 2018) objective to aid comparison in Section D.2.

374 InfoNCE is an unsupervised contrastive learning objective built on the noise contrastive estima-  
 375 tion (Gutmann & Hyvärinen, 2010) principle. In particular, given an ‘‘anchor’’ datum  $x$  (otherwise  
 376 known as context), and distribution of positives  $x_{\text{pos}}$  and negatives  $x_{\text{neg}}$ , the InfoNCE objective  
 377 optimizes

$$\min_{\phi} \mathbb{E}_{x_{\text{pos}}} \left[ -\log \frac{\mathcal{S}_{\phi}(x, x_{\text{pos}})}{\mathbb{E}_{x_{\text{neg}}} \mathcal{S}_{\phi}(x, x_{\text{neg}})} \right], \quad (8)$$

378 where  $\mathbb{E}_{x_{\text{neg}}}$  is often approximated with a fixed number of negatives in practice.

379 It is shown in Oord et al. (2018) that optimizing (8) is maximizing a lower bound on the mutual  
 380 information  $\mathcal{I}(x, x_{\text{pos}})$ , where, with slight abuse of notation,  $x$  and  $x_{\text{pos}}$  are interpreted as random  
 381 variables.

382 TCN is a contrastive learning objective that learns a representation that in timeseries data (e.g., video  
 383 trajectories). The original work (Sermanet et al., 2018) considers multi-view videos and perform  
 384 contrastive learning over frames in separate videos; in this work, we consider the single-view variant.  
 385 At a high level, TCN attracts representations of frames that are temporally close, while pushing apart  
 386 those of frames that are farther apart in time. More precisely, given three frames sampled from a  
 387 video sequence  $(o_{t_1}, o_{t_2}, o_{t_3})$ , where  $t_1 < t_2 < t_3$ , TCN would attract the representations of  $o_{t_1}$   
 388 and  $o_{t_2}$  and repel the representation of  $o_{t_3}$  from  $o_{t_1}$ . This idea can be formally expressed via the  
 389 following objective:

$$\min_{\phi} \mathbb{E}_{(o_{t_1}, o_{t_2} > t_1) \sim D} \left[ -\log \frac{\mathcal{S}_{\phi}(o_{t_1}; o_{t_2})}{\mathbb{E}_{o_{t_3} | t_3 > t_2 \sim D} [\mathcal{S}_{\phi}(o_{t_1}; o_{t_3})]} \right] \quad (9)$$

390 Given a ‘‘positive’’ window of  $K$  steps and a uniform distribution among valid positive samples, we  
 391 can write (9) as

$$\min_{\phi} \frac{1}{K} \sum_{k=1}^K \mathbb{E}_{(o_{t_1}, o_{t_1+k}) \sim D} \left[ -\log \frac{\mathcal{S}_{\phi}(o_{t_1}; o_{t_1+k})}{\mathbb{E}_{o_{t_3} | t_3 > t_1+k \sim D} [\mathcal{S}_{\phi}(o_{t_1}; o_{t_3})]} \right], \quad (10)$$

392 in which each term inside the expectation is a standalone InfoNCE objective tailored to observation  
 393 sequence data.

## 394 C Related Work

395 We review relevant literature on (1) Out-of-Domain Representation Pre-Training for Control, (2)  
 396 Perceptual Reward Learning from Human Videos, and (3) Goal-Conditioned RL as Representation  
 397 Learning.

398 **Out-of-Domain Representation Pre-Training for Control.** Bootstrapping visual control using  
 399 frozen representations learned pre-trained on out-of-domain non-robot data is a nascent field that has  
 400 seen fast progress over the past year. Shah & Kumar (2021) demonstrates that pre-trained ResNet (He  
 401 et al., 2016) representation on ImageNet (Deng et al., 2009) serves as effective visual backbone

402 for simulated dexterous manipulation RL tasks. Parisi et al. (2022) finds ResNet models trained  
403 with unsupervised objectives, such as momentum contrastive learning (MOCO) (He et al., 2020), to  
404 surpass supervised objectives (e.g. image classification) for both visual navigation and control tasks.  
405 Xiao et al. (2022) pre-trains visual representation on human video data (Goyal et al., 2017; Shan  
406 et al., 2020) using masked-autoencoding (He et al., 2022). Along this axis, the closest work to ours is  
407 Nair et al. (2022), which is also pre-trained on the Ego4D dataset and attempts to capture temporal  
408 information in the videos by using time-contrastive learning (Sermanet et al., 2018); it additionally  
409 leverages textual descriptions associated with the videos to encode semantic information. In contrast,  
410 our objective is fully self-supervised without dependence on textual annotations. Furthermore, VIP  
411 is the first to propose using a RL-based objective for out-of-domain pre-training and is capable of  
412 producing generalizable dense reward signals.

413 **Perceptual Reward Learning from Human Videos.** Human videos provide a rich natural source  
414 of reward and representation learning for robotic learning. Most prior works exploit the idea of  
415 learning an invariant representation between human and robot domains to transfer the demonstrated  
416 skills (Sermanet et al., 2016, 2018; Schmeckpeper et al., 2020; Chen et al., 2021; Xiong et al., 2021;  
417 Zakka et al., 2022; Bahl et al., 2022). However, training these representations require task-specific  
418 human *demonstration* videos paired with robot videos solving the same task, and cannot leverage the  
419 large amount of “in-the-wild” human videos readily available. As such, these methods require robot  
420 data for training, and learn rewards that are task-specific and do not generalize beyond the tasks they  
421 are trained on. In contrast, VIP do not make any assumption on the quality or the task-specificity of  
422 human videos and instead pre-trains an (implicit) value function that aims to capture task-agnostic  
423 goal-oriented progress, which can generalize to completely unseen robot domains and tasks.

424 **Goal-Conditioned RL as Representation Learning.** Our pre-training method is also related to the  
425 idea of treating goal-conditioned RL as representation learning. Chebotar et al. (2021) shows that a  
426 goal-conditioned Q-function trained with offline in-domain multi-task robot data learns an useful  
427 visual representation that can accelerate learning for a new downstream task in the same domain.  
428 Eysenbach et al. (2022) shows that goal-conditioned Q-learning with a particular choice of reward  
429 function can be understood as performing contrastive learning. In contrast, our theory introduces  
430 a new implicit time contrastive learning, and states that for *any* choice of reward function, the dual  
431 formulation of a regularized offline GCRL objective can be cast as implicit time contrast. This  
432 conceptual bridge also explains why VIP’s learned embedding distance is temporally smooth and can  
433 be used as an universal reward mechanism. Finally, whereas these two works are limited to training  
434 on in-domain data with robot action labels, VIP is able to leverage diverse out-of-domain human data  
435 for visual representation pre-training, overcoming the inherent limitation of robot data scarcity for  
436 in-domain training.

437 Our work is also closely related to Ma et al. (2022b), which first introduced the dual offline GCRL  
438 objective based on Fenchel duality (Rockafellar, 1970; Nachum & Dai, 2020; Ma et al., 2022a).  
439 Whereas Ma et al. (2022b) assumes access to the true state information and focuses on the offline  
440 GCRL setting using in-domain offline data with robot action labels, we extend the dual objective  
441 to enable out-of-domain, action-free pre-training from human videos. Our particular dual objective  
442 also admits a novel implicit time contrastive learning interpretation, which simplifies VIP’s practical  
443 implementation by letting the value function be implicitly defined instead of a deep neural network  
444 as in Ma et al. (2022b).

## 445 **D Value-Implicit Pre-Training (Full-Version)**

446 In this section, we demonstrate how a self-supervised value-function objective can be derived from  
447 computing the dual of an offline RL objective on passive human videos (Section D.1). Then, we  
448 show how this objective amounts to a novel implicit formulation of temporal contrastive learning  
449 (Section D.2), which naturally lends a temporally and locally smooth embedding favorable for  
450 downstream visual reward specification. Finally, we leverage this contrastive interpretation to  
451 instantiate a simple implementation (<10 lines of PyTorch code) of our dual value objective that does

452 not explicitly learn a value network (Section D.3), culminating in our final algorithm, Value-Implicit  
 453 Pre-training (VIP).

#### 454 D.1 Foundation: Self-Supervised Value Learning from Human Videos

455 While human videos are out-of-domain data for robots, they are *in-domain* for learning a goal-  
 456 conditioned policy  $\pi_H$  over human actions,  $a^H \sim \pi^H(\phi(o) \mid \phi(g))$ , for some human action space  
 457  $A^H$ . Therefore, given that human videos naturally contain goal-directed behavior, one reasonable idea  
 458 of utilizing offline human videos for representation learning is to solve an offline goal-conditioned  
 459 RL problem over the space of human policies and then extract the learned visual representation. To  
 460 this end, we consider the following KL-regularized offline RL objective (Nachum et al., 2019) for  
 461 some to-be-specified reward  $r(o, g)$ :

$$\max_{\pi_H, \phi} \mathbb{E}_{\pi_H} \left[ \sum_t \gamma^t r(o; g) \right] - (d^{\pi_H}(o, a^H; g) \| d^D(o, \tilde{a}^H; g)), \quad (11)$$

462 where  $d^{\pi_H}(o, a^H; g)$  is the distribution over observations and actions  $\pi_H$  visits conditioned on  $g$ .  
 463 Observe that a “dummy” action  $\tilde{a}$  is added to every transition  $(o_h^i, \tilde{a}_h^i, o_{h+1}^i)$  in the dataset  $D$  so that  
 464 the KL regularization is well-defined, and  $\tilde{a}_h^i$  can be thought of as the unobserved *true* human action  
 465 taken to transition from observation  $o_h^i$  to  $o_{h+1}^i$ . While this objective is mathematically sound and  
 466 encourages learning a conservative  $\pi^H$ , it is seemingly implausible because the offline dataset  $D^H$   
 467 does not come with any action labels nor can  $A^H$  be concretely defined in practice. However, what  
 468 this objective does provide is an elegant *dual* objective over a value function that does not depend on  
 469 any action label in the offline dataset. In particular, leveraging the idea of Fenchel duality (Rockafellar,  
 470 1970) from convex optimization, we have the following result:

471 **Proposition D.1.** *Under assumption of deterministic transition dynamics, the dual optimization*  
 472 *problem of (11) is*

$$473 \max_{\phi} \min_V \mathbb{E}_{p(g)} [(1 - \gamma) \mathbb{E}_{\mu_0(o; g)} [V(\phi(o); \phi(g))] + \log \mathbb{E}_{(o, o'; g) \sim D} [\exp(r(o, g) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g)))]], \quad (12)$$

474 where  $\mu_0(o; g)$  is the goal-conditioned initial observation distribution, and  $D(o, o'; g)$  is the goal-  
 475 conditioned distribution of two consecutive observations in dataset  $D$ .

476 As shown, actions do not appear in the objective. Furthermore, since all expectations in (12) can be  
 477 sampled using the offline dataset, this dual value-function objective can be self-supervised with an  
 478 appropriate choice of reward function. In particular, since our goal is to acquire a value function that  
 479 extracts a general notion of goal-directed task progress from passive offline human videos, we set  
 480  $r(o, g) = \mathbb{I}(o == g) - 1$ , which we refer to as  $\tilde{\delta}_g(o)$  in shorthand. This reward provides a constant  
 481 negative reward when  $o$  is not the provided goal  $g$ , and does not require any task-specific engineering.  
 482 The resulting value function  $V(\phi(o); \phi(g))$  captures the discounted total number of steps required to  
 483 reach goal  $g$  from observation  $o$ . Consequently, the overall objective will encourage learning visual  
 484 features  $\phi$  that are amenable to predicting the discounted temporal distance between two frames in a  
 485 human video sequence. With enough size and diversity in the training dataset, we hypothesize that  
 486 this value function can generalize to completely unseen (robot) domains and tasks.

#### 487 D.2 Analysis: Implicit Time Contrastive Learning

488 While (12) will learn some useful visual representation via temporal value function optimization,  
 489 in this section, we show that it can be understood as a novel *implicit* temporal contrastive learning  
 490 objective that acquires temporally smooth embedding distance over video sequences, underpinning  
 491 VIP’s efficacy jointly as a visual representation and reward for downstream control.

492 We begin by simplifying the expression in (12) by first assuming that the optimal  $V^*$  is found:

$$493 \min_{\phi} \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o; g)} [-V^*(\phi(o); \phi(g))] + \log \mathbb{E}_{D(o, o'; g)} \left[ \exp \left( \tilde{\delta}_g(o) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g)) \right) \right]^{-1} \right], \quad (13)$$

494 where we have also re-written the maximization problem as a minimization problem. Now, after  
 495 few algebraic manipulation steps (see App. E for a derivation), if we think of  $V^*(\phi(o); \phi(g))$  as a

496 *similarity metric* in the embedding space, then we can massage (13) into an expression that resembles  
 497 the InfoNCE (Oord et al., 2018) time contrastive learning (Sermanet et al., 2018) (see App. B.2 for a  
 498 definition and additional background) objective:

$$499 \min_{\phi} (1 - \gamma) \mathbb{E}_{p(g), \mu_0(o;g)} \left[ - \log \frac{e^{V^*(\phi(o); \phi(g))}}{\mathbb{E}_{D(o, o';g)} [\exp(\tilde{\delta}_g(o) + \gamma V^*(\phi(o'); \phi(g)) - V^*(\phi(o), \phi(g)))]^{\frac{-1}{1-\gamma}}} \right] \quad (14)$$

500 In particular,  $p(g)$  can be thought of the distribution of “anchor” observations,  $\mu_0(s; g)$  the distribution  
 501 of “positive” samples, and  $D(o, o'; g)$  the distribution of “negative” samples. Counter-intuitively and  
 502 in contrast to standard single-view time contrastive learning (TCN), in which the positive observations  
 503 are temporally closer to the anchor observation than the negatives, (14) has the positives to be as  
 504 temporally far away as possible, namely the initial frame in the the same video sequence, and the  
 505 negatives to be middle frames sampled in between. This departure is accompanied by the equally  
 506 intriguing deviation of the lack of explicit repulsion of the negatives from the anchor; instead, they  
 507 are simply encouraged to minimize the (exponentiated) one-step temporal-difference error in the  
 508 representation space (the denominator in (14)); see Fig. 1. Now, since the value function encodes  
 509 negative discounted temporal distance, due to the recursive nature of value temporal-difference (TD),  
 510 in order for the one-step TD error to be globally minimized along a video sequence, observations that  
 511 are temporally farther away from the goal will naturally be repelled farther away in the representation  
 512 space compared to observations that are nearby in time; in App. E.3, we formalize this intuition and  
 513 show that this repulsion always holds for optimal paths. Therefore, the repulsion of the negative  
 514 observations is an *implicit*, emergent property from the optimization of (14), instead of an explicit  
 515 constraint as in standard (time) contrastive learning.

516 Now, we dive into why this *implicit* time contrastive learning is desirable. First, the explicit  
 517 attraction of the initial and goal frames enables capturing *long-range* semantic temporal dependency  
 518 as two frames that meaningfully indicate the beginning and end of a task are made close  
 519 in the embedding space. This closeness is also  
 520 well-defined due to the one-step TD backup that  
 521 makes every embedding distance recursively defined  
 522 to be the discounted number of timesteps to the goal frame. **Combined with the implicit yet structured repulsion of intermediate frames,**

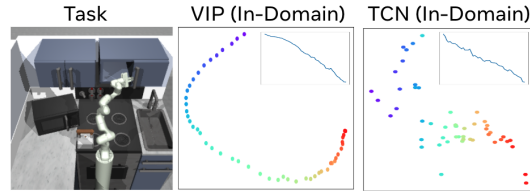


Figure 4: Learned 2D representation of a held-out task demonstration by VIP and TCN trained on task-specific in-domain data. The color gradient indicates trajectory time progression (purple for beginning, red for end). The inset plots are embedding distances to last frame.

523 this push-and-pull mechanism helps inducing a *temporally smooth* and consistent representation. In  
 524 particular, as we pass a video sequence in the training set through the trained representation, the em-  
 525 bedding should be structured such that two trends emerge: (1) neighboring frames are close-by in the  
 526 embedding space, (2) their distances to the last (goal) frame smoothly decrease due to the recursively  
 527 defined embedding distances. To validate this intuition, in Fig. 4, we provide a simple toy example  
 528 comparing implicit vs. standard time contrastive learning when trained on *in-domain, task-specific*  
 529 demonstrations; details are included in App. G.2. As shown, standard time contrastive learning only  
 530 enforces a coarse notion of temporal consistency and learns a non-locally smooth representation  
 531 that exhibits many local minima. In contrast, VIP learns a much better structured embedding that is  
 532 indeed temporally consistent and locally smooth. **As we will show, the prevalence of sharp “bumps”  
 533 in the embedding distance as in TCN can be easily exploited by the control algorithm, and VIP’s  
 534 ability to generate long-range temporally smooth embedding is the key ingredient for its effective  
 535 downstream zero-shot reward-specification.**

### 541 D.3 Algorithm: Value-Implicit Pre-Training (VIP)

542 The theoretical development in the previous two sections culminates in *Value Implicit Pre-Training*  
 543 (VIP), a simple value-based self-supervised pre-training objective, in which the value function is  
 544 implicitly represented via the learned embedding distance.

545 Recall that  $V^*$  is assumed to be known for the derivation in Section D.2, but in practice, its analytical  
 546 form is rarely known. Now, given that  $V^*$  plays the role of a distance measure in our implicit time  
 547 contrastive learning framework, a simple and practical way to approximate  $V^*$  is to simply set it to  
 548 be a choice of **similarity metric**, bypassing having to explicitly parameterize it as a neural network. In  
 549 this work, we choose the common choice of the negative  $L_2$  distance used in prior work Sermanet  
 550 et al. (2018); Nair et al. (2022):  $V^*(\phi(o), \phi(g)) := -\|\phi(o) - \phi(g)\|_2$ . Given this choice, our final  
 551 representation learning objective is as follows:

$$552 \quad \mathcal{L}(\phi) = \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [\|\phi(o) - \phi(g)\|_2] + \log \mathbb{E}_{(o,o';g) \sim D} \left[ \exp \left( \|\phi(o) - \phi(g)\|_2 - \tilde{\delta}_g(o) - \gamma \|\phi(o') - \phi(g)\|_2 \right) \right] \right], \quad (15)$$

553 in which we also absorb the exponent of the log-sum-exp term in 13 into the inner  $\exp(\cdot)$  term via  
 554 an Jensen’s inequality; we found this upper bound to be numerically more stable. To sample video  
 555 trajectories from  $D$ , because any sub-trajectory of a video is also a valid video sequence, VIP samples  
 556 these sub-trajectories and treats their initial and last frames as samples from the goal and initial-state  
 557 distributions (Step 3 in Alg. 2). Altogether, VIP training is illustrated in Alg. 2; it is simple and its  
 558 core training loop can be implemented in fewer than 10 lines of PyTorch code (Alg. 3 in App. F.3).

---

### Algorithm 2 Value-Implicit Pre-Training (VIP)

---

- 1: **Require:** Offline (human) videos  $D = \{(o_1^i, \dots, o_{h_i}^i)\}_{i=1}^N$ , visual architecture  $\phi$
  - 2: **for** number of training iterations **do**
  - 3:   Sample sub-trajectories  $\{o_t^i, \dots, o_k^i, o_{k+1}^i, \dots, o_T^i\}_{i=1}^B \sim D, t \in [1, h_i - 1], t \leq k < T, T \in (t, h_i], \forall i$
  - 4:    $\mathcal{L}(\phi) := \frac{1-\gamma}{B} \sum_{i=1}^B [\|\phi(o_t^i) - \phi(o_T^i)\|_2] + \log \frac{1}{B} \sum_{i=1}^B \left[ \exp \left( \|\phi(o_k^i) - \phi(o_T^i)\|_2 - \tilde{\delta}_{o_T^i}(o_k^i) - \gamma \|\phi(o_{k+1}^i) - \phi(o_T^i)\|_2 \right) \right]$
  - 5:   Update  $\phi$  using SGD:  $\phi \leftarrow \phi - \alpha_\phi \nabla \mathcal{L}(\phi)$
- 

## 559 E Technical Derivations and Proofs

### 560 E.1 Proof of Proposition D.1

561 We first reproduce Proposition D.1 for ease of reference:

562 **Proposition E.1.** *Under assumption of deterministic transition dynamics, the dual optimization*  
 563 *problem of*

$$\max_{\pi_H, \phi} \mathbb{E}_{\pi_H} \left[ \sum_t \gamma^t r(o; g) \right] - (d^{\pi_H}(o, a^H; g) \| d^D(o, \tilde{a}^H; g)), \quad (16)$$

564 is

$$\max_{\phi} \min_V \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [V(\phi(o); \phi(g))] + \log \mathbb{E}_{D(o,o';g)} [\exp(r(o, g) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g)))] \right], \quad (17)$$

565 where  $\mu_0(o; g)$  is the goal-conditioned initial observation distribution, and  $D(o, o'; g)$  is the goal-  
 566 conditioned distribution of two consecutive observations in dataset  $D$ .

567 *Proof.* We begin by rewriting (16) as an optimization problem over valid state-occupancy distribu-  
 568 tions. To this end, we have<sup>1</sup>

$$\begin{aligned} & \max_{\phi} \max_{d(\phi(o), a; \phi(g)) \geq 0} \mathbb{E}_{d(\phi(o), \phi(g))} [r(o; g)] - (d(\phi(o), a; \phi(g)) \| d^D(\phi(o), \tilde{a}; \phi(g))) \\ \text{(P)} \quad & \text{s.t.} \quad \sum_a d(\phi(o), a; \phi(g)) = (1 - \gamma) \mu_0(o; g) + \gamma \sum_{\tilde{o}, \tilde{a}} T(o | \tilde{o}, \tilde{a}) d(\phi(\tilde{o}), \tilde{a}; \phi(g)), \forall o \in O, g \in G \end{aligned} \quad (18)$$

569 Fixing a choice of  $\phi$ , the inner optimization problem operates over a  $\phi$ -induced state and goal space,  
 570 giving us (18). Then, applying Proposition 4.2 of Ma et al. (2022b) to the inner optimization problem,

---

<sup>1</sup>We omit the human action superscript  $H$  in this derivation.

571 we immediately obtain

$$(D) \quad \max_{\phi} \min_V \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [V(\phi(o); \phi(g))] \right. \\ \left. + \log \mathbb{E}_{d^D(\phi(o), a; \phi(g))} \left[ \exp(r(o, g) + \gamma \mathbb{E}_{T(o'|o, a)} [V(\phi(o'); \phi(g))] - V(\phi(o), \phi(g))) \right] \right] \quad (19)$$

572 Now, given our assumption that the transition dynamics is deterministic, we can replace the inner  
573 expectation  $\mathbb{E}_{T(o'|o, a)}$  with just the observed sample in the offline dataset and obtain:

$$\max_{\phi} \min_V \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [V(\phi(o); \phi(g))] \right. \\ \left. + \log \mathbb{E}_{d^D(\phi(o), \phi(o'); \phi(g))} \left[ \exp(r(o, g) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g))) \right] \right] \quad (20)$$

574 Finally, sampling embedded states from  $d^D(\phi(o), \phi(o'); \phi(g))$  is equivalent to sampling from  
575  $D(o, o'; g)$ , assuming there is no embedding collision (i.e.,  $\phi(o) \neq \phi(o'), \forall o \neq o'$ ), which can  
576 be satisfied by simply augmenting any  $\phi$  by concatenating the input to the end. Then, we have our  
577 desired expression:

$$\max_{\phi} \min_V \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [V(\phi(o); \phi(g))] + \log \mathbb{E}_{D(o, o'; g)} \left[ \exp(r(o, g) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g))) \right] \right] \quad (21)$$

578  $\square$

## 579 E.2 VIP Implicit Time Contrast Learning Derivation

580 This section provides all intermediate steps to go from (13) to (14). First, we have

$$\min_{\phi} \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [-V^*(\phi(o); \phi(g))] + \log \mathbb{E}_{D(o, o'; g)} \left[ \exp(\tilde{\delta}_g(o) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g))) \right]^{-1} \right]. \quad (22)$$

581 We can equivalently write this objective as

$$\min_{\phi} \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} [-\log e^{V^*(\phi(o); \phi(g))}] + \log \mathbb{E}_{D(o, o'; g)} \left[ \exp(\tilde{\delta}_g(o) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g))) \right]^{-1} \right]. \quad (23)$$

582 Then,

$$\min_{\phi} \mathbb{E}_{p(g)} \left[ (1 - \gamma) \mathbb{E}_{\mu_0(o;g)} \left[ -\log e^{V^*(\phi(o); \phi(g))} - \log \mathbb{E}_{D(o, o'; g)} \left[ \exp(\tilde{\delta}_g(o) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g))) \right]^{\frac{1}{1-\gamma}} \right] \right] \\ = \min_{\phi} (1 - \gamma) \mathbb{E}_{p(g), \mu_0(o;g)} \left[ \log \frac{e^{-V^*(\phi(o); \phi(g))}}{\mathbb{E}_{D(o, o'; g)} \left[ \exp(\tilde{\delta}_g(o) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g))) \right]^{\frac{1}{1-\gamma}}} \right] \quad (24)$$

583 This is (14) in the main text.

## 584 E.3 VIP Implicit Repulsion

585 In this section, we formalize the implicit repulsion property of VIP objective ((14)); in particular, we  
586 prove that under certain assumptions, it always holds for optimal paths.

587 **Proposition E.2.** *Suppose  $V^*(s; g) := -\|\phi(s) - \phi(g)\|_2$  for some  $\phi$ , under the assumption of*  
588 *deterministic dynamics (as in Proposition D.1), for any pair of consecutive states reached by the*  
589 *optimal policy,  $(s_t, s_{t+1}) \sim \pi^*$ , we have that*

$$\|\phi(s_t) - \phi(g)\|_2 > \|\phi(s_{t+1}) - \phi(g)\|_2, \quad (25)$$

590 *Proof.* First, we note that

$$V^*(s; g) = \max_a Q^*(s, a; g) \quad (26)$$

591 A proof can be found in Section 1.1.3 of Agarwal et al. (2019). Then, due to the Bellman optimality  
592 equation, we have that

$$Q^*(s, a; g) = r(s, g) + \gamma \mathbb{E}_{s' \sim T(s, a)} \max_{a'} Q^*(s', a'; g) \quad (27)$$



593 Given that the dynamics is deterministic and (26), we have that

$$Q^*(s, a; g) = r(s, g) + \gamma V^*(s'; g) \tag{28}$$

594 Now, for  $(s_t, a_t, s_{t+1}) \sim \pi^*$ , this further simplifies to

$$V^*(s_t; g) = r(s_t, g) + \gamma V^*(s_{t+1}; g) \tag{29}$$

595 Note that since  $V^*$  is also the optimal value function, given that  $r(s_t, g) = \mathbb{I}(s_t = g) - 1$ ,  $V^*(s_t; g)$   
 596 is the *negative* discounted distance of the shortest path between  $s_t$  and  $g$ . In particular, since  
 597  $V^*(g; g) = 0$  by construction, we have that  $V^*(s_t; g) = -\sum_{k=0}^K \gamma^k$  (this also clearly satisfies (29)),  
 598 where the shortest path (i.e., the path  $\pi^*$  takes) between  $s_t$  and  $g$  are  $K$  steps long. Now, giving that  
 599 we assume  $V^*(s_t; g)$  can be expressed as  $-\|\phi(s_t) - \phi(g)\|_2$  for some  $\phi$ , it immediately follows that  
 600

$$\|\phi(s_t) - \phi(g)\|_2 > \|\phi(s_{t+1}) - \phi(g)\|_2, \quad \forall (s_t, s_{t+1}) \sim \pi^* \tag{30}$$

601

□

602 The implication of this result is that at least along the trajectories generated by the optimal policy, the  
 603 representation will have monotonically decreasing and well-behaved embedding distances to the goal.  
 604 Now, since in practice, VIP is trained on goal-directed (human video) trajectories, which are near-  
 605 optimal for goal-reaching, we expect this smoothness result to be informative about VIP’s embedding  
 606 practical behavior and help formalize out intuition about the mechanism of implicit time contrastive  
 607 learning. As confirmed by our qualitative study in Section H.4, We highlight that VIP’s embedding is  
 608 indeed much smoother than other baselines along test trajectories on both Ego4D and on our real-robot  
 609 dataset. This smoothness along optimal paths makes it easier for the downstream control optimizer to  
 610 discover these paths, conferring VIP representation effective zero-shot reward-specification capability  
 611 that is not attained by any other comparison.

## 612 F VIP Training Details

### 613 F.1 Dataset Processing and Sampling

614 We use the exact same pre-processed Ego4D dataset as in R3M, in which long raw videos are first  
 615 processed into shorter videos consisting of 60-70 frames each. In total, there are approximately 72000  
 616 clips and 4.3 million frames in the dataset. Within a sampled batch, we first sample a set of videos,  
 617 and then sample a sub-trajectory from each video (Step 3 in Algorithm 2). In this formulation, each  
 618 sub-trajectory is treated as a video segment from the algorithm’s perspective; this can viewed as a  
 619 variant of trajectory data augmentation. As in R3M, we apply random crop at a video level within  
 620 a batch, so all frames from the same video sub-trajectory are cropped the same way. Then, each  
 621 raw observation is resized and center-cropped to have shape  $224 \times 224 \times 3$  before passed into the  
 622 visual encoder. Finally, as in standard contrastive learning and R3M, for each sampled sub-trajectory  
 623  $\{o_t^i, \dots, o_k^i, o_{k+1}^i, \dots, o_T^i\}$ , we also sample additional 3 negative samples  $(\tilde{o}_j, \tilde{o}_{j+1})$  from separate  
 624 video sequences to be included in the log-sum-exp term in  $\mathcal{L}(\phi)$ .

### 625 F.2 VIP Hyperparameters

626 Hyperparameters used can be found in Table 2.

### 627 F.3 VIP Pytorch Pseudocode

628 In this section, we present a pseudocode of VIP written in PyTorch (Paszke et al., 2019), Algorithm 3.  
 629 As shown, the main training loop can be as short as 10 lines of code.

Table 2: VIP Architecture &amp; Hyperparameters.

	Name	Value
Architecture	Visual Backbone	ResNet50 (He et al., 2016)
	FC Layer Output Dim	1024
Hyperparameters	Optimizer	Adam (Kingma & Ba, 2014)
	Learning rate	0.0001
	$L_1$ weight penalty	0.001
	$L_1$ weight penalty	0.001
	Mini-batch size	32
	Discount factor $\gamma$	0.98

**Algorithm 3** VIP PyTorch Pseudocode

```

# D: offline dataset
# phi: vision architecture

# training loop
for (o_0, o_t1, o_t2, g) in D:
    phi_g = phi(o_g)
    V_0 = - torch.linalg.norm(phi(o_0), phi_g)
    V_t1 = - torch.linalg.norm(phi(o_t1), phi_g)
    V_t2 = - torch.linalg.norm(phi(o_t2), phi_g)
    VIP_loss = (1-gamma)*-V_0.mean() + torch.logsumexp(V_t1+1-gamma*V_t2)
    optimizer.zero_grad()
    VIP_loss.backward()
    optimizer.step()

```

630 **G Simulation Experiment Details.**631 **G.1 FrankaKitchen Task Descriptions**

632 In this section, we describe the FrankaKitchen suite for our simulation experiments. We use 12 tasks  
633 from the v0.1 version<sup>2</sup> of the environment.

634 We use the environment default initial state as the initial state and frame for all tasks in the Hard  
635 setting. In the Easy setting, we use the 20th frame of a demonstration trajectory and its corresponding  
636 environment state as the initial frame and state. The goal frame for both settings is chosen to be the  
637 last frame of the same demonstration trajectory. The initial frames and goal frame for all 12 tasks and  
638 3 camera views are illustrated in Figure 5-6. In the Easy setting, the horizon for all tasks is 50 steps;  
639 in the Hard setting, the horizon is 100 steps. Note that using the 20th frame as the initial state is a  
640 crude way for initializing the robot, and for some tasks, this initialization makes the task substantially  
641 easier, whereas for others, the task is still considerably difficult. Furthermore, some tasks become  
642 naturally more difficult depending on camera viewpoints. For these reasons, it is worth noting that  
643 our experiment’s emphasis is on the *aggregate* behavior of pre-trained representations, instead of  
644 trying to solve any particular task as well as possible.

645 **G.2 In-Domain Representation Probing**

646 In this section, we describe the experiment we performed to generate the in-domain VIP vs. TCN  
647 comparison in Figure 4. We fit VIP and TCN representations using 100 demonstrations from the  
648 FrankaKitchen `sdoor_open` task (center view). For TCN, we use R3M’s implementation of the  
649 TCN loss without any modification; this also allows our findings in Figure 4 to extend to the main  
650 experiment section. The visual architecture is ResNet34, and the output dimension is 2, which enables  
651 us to directly visualize the learned embedding. Different from the out-of-domain version of VIP, we  
652 also do not perform weight penalty, trajectory-level random cropping data augmentation, or additional

<sup>2</sup>[https://github.com/vikashplus/mj\\_envs/tree/v0.1real/mj\\_envs/envs/relay\\_kitchen](https://github.com/vikashplus/mj_envs/tree/v0.1real/mj_envs/envs/relay_kitchen)



Figure 5: Initial frame (Easy), initial frame (Hard), and goal frame for all 12 tasks and 3 camera views in our FrankaKitchen suite.

653 negative sampling. Besides these choices, we use the same hyperparameters as in Table 2 and train  
 654 for 2000 batches.

### 655 G.3 Trajectory Optimization

656 We use a publicly available implementation of MPPI<sup>3</sup>, and make no modification to the algorithm or  
 657 the default hyperparameters. In particular, the planning horizon is 12 and 32 sequences of actions  
 658 are proposed per action step. Because the embedding reward ((4)) is the goal-embedding distance  
 659 difference, the score (i.e., sum of per-transition reward) of a proposed sequence of actions is equivalent  
 660 to the negative embedding distance (i.e.,  $S_\phi(\phi(o_T); \phi(g))$ ) at the last observation.

#### 661 G.3.1 Robot and Object Pose Error Analysis

662 In this section, we visualize the per-step robot and object pose  $L_2$  error with respect to the goal-image  
 663 poses. We report the non-cumulative curves (on the success rate as well) for more informative  
 664 analysis.

<sup>3</sup><https://github.com/aravindr93/trajopt/blob/master/trajopt/algos/mppi.py>



Figure 6: Initial frame (Easy), initial frame (Hard), and goal frame for all 12 tasks and 3 camera views in our FrankaKitchen suite.



Figure 7: Trajectory optimization results with pose errors.

Table 3: Real-world robotics tasks descriptions.

Environment	Object Type	Dataset	Success Criterion
CloseDrawer	Articulated Object	10 demos + 20 failures	the drawer is closed enough that the spring loads.
PushBottle	Transparent Object	20 demonstrations	the bottle is parallel to the goal line set by the icecream cone.
PlaceMelon	Soft Object	20 demonstrations	the watermelon toy is fully placed in the plate.
FoldTowel	Deformable Object	20 demonstrations	the bottom half of the towel is cleanly covered by the top half.



Figure 8: Real-robot setup.

#### 665 G.4 Reinforcement Learning

666 We use a publicly available implementation of NPG<sup>4</sup>, and make no modification to the algorithm or  
 667 the default hyperparameters. In the Easy (resp. Hard) setting, we train the policy until 500000 (resp.  
 668 1M) real environment steps are taken. For evaluation, we report the cumulative maximum success  
 669 rate on 50 test rollouts from each task configuration (50\*108=5400 total rollouts) every 10000 step.

## 670 H Real-World Robot Experiment Details

### 671 H.1 Task Descriptions

672 The robot learning environment is illustrated in Figure 8; a RealSense camera is mounted on the  
 673 right edge of the table, and we only use the RGB image stream without depth information for data  
 674 collection and policy learning.

675 We collect offline data  $D_{\text{task}}$  for each task via kinesthetic playback, and the object initial placement  
 676 is randomized for each trajectory. On the simplest CloseDrawer task, we combine 10 expert  
 677 demonstrations with 20 sub-optimal failure trajectories to increase learning difficulty. For the other  
 678 three tasks, we collect 20 expert demonstrations, which we found are difficult enough for learning  
 679 good policies. Each demonstration is 50-step long collected at 25Hz. The initial state for the robot is  
 680 fixed for each demonstration and test rollout, but the object initial position is randomized. The task  
 681 success is determined based on a visual criterion that we manually check for each test rollout. The  
 682 full task breakdown is described in Table 3.

683 Each task is specified via a set of goal images that are chosen to be the last frame of all demonstrations  
 684 for the task. Hence, the goal embedding used to compute the embedding reward ((4)) for each task is  
 685 the average over the embeddings of all goal frames.

686 The tasks (in their initial positions) using a separate high-resolution phone camera are visualized in  
 687 Figure 9. Sample demonstrations in the robot camera view are visualized in Figure 10.

<sup>4</sup>[https://github.com/aravindr93/mjrl/blob/master/mjrl/algos/npg\\_cg.py](https://github.com/aravindr93/mjrl/blob/master/mjrl/algos/npg_cg.py)

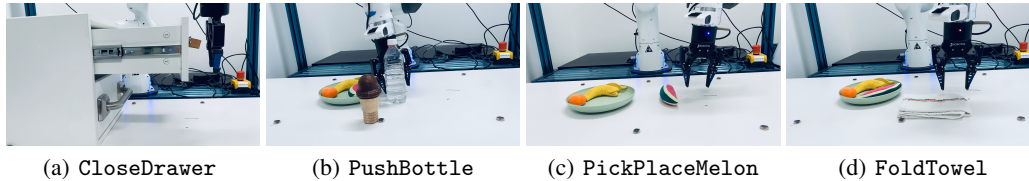


Figure 9: Side-view of real-robot tasks using a high-resolution smartphone camera.

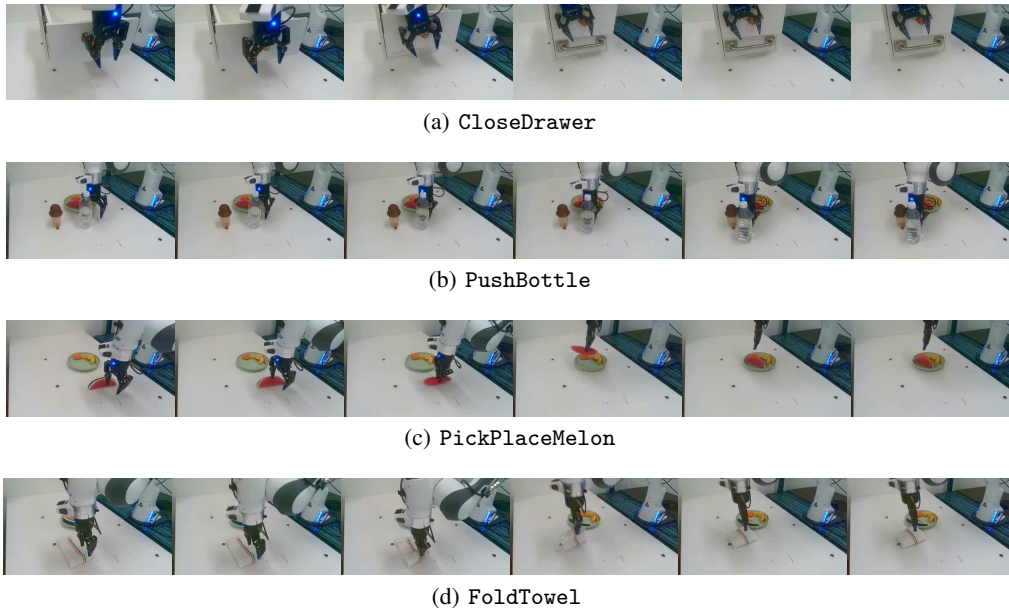


Figure 10: Real-robot task demonstrations (every 10th frame) in robot camera view. The first and last frames in each row are representative of initial and final goal observations for the respective task.

## 688 H.2 Training and Evaluation Details

689 The policy network is implemented as a 2-layer MLP with hidden sizes [256, 256]. As in R3M’s  
 690 real-world robot experiment setup, the policy takes in concatenated visual embedding of current  
 691 observation and robot’s proprioceptive state and outputs robot action. The policy is trained with a  
 692 learning rate of 0.001, and a batch size of 32 for 20000 steps.

693 For RWR’s temperature scale, we use  $\tau = 0.1$  for all tasks, except CloseDrawer where we find  
 694  $\tau = 1$  more effective for both VIP and R3M.

695 For policy evaluation, we use 10 test rollouts with objects randomly initialized to reflect the object  
 696 distribution in the expert demonstrations. The rollout horizon is 100 steps.

## 697 H.3 Additional Analysis & Context

698 **Offline RL vs. imitation learning for real-world robot learning.** Offline RL, though known  
 699 as the data-driven paradigm of RL (Levine et al., 2020), is not necessarily data *efficient* (Agarwal  
 700 et al., 2021), requiring hundreds of thousands of samples even in low-dimensional simulated tasks,  
 701 and requires a dense reward to operate most effectively (Mandlekar et al., 2021; Yu et al., 2022).  
 702 Furthermore, offline RL algorithms are significantly more difficult to implement and tune compared to  
 703 BC (Kumar et al., 2021; Zhang & Jiang, 2021). As such, the dominant paradigm of real-world robot  
 704 learning is still learning from demonstrations (Jang et al., 2022; Mandlekar et al., 2018; Ebert et al.,

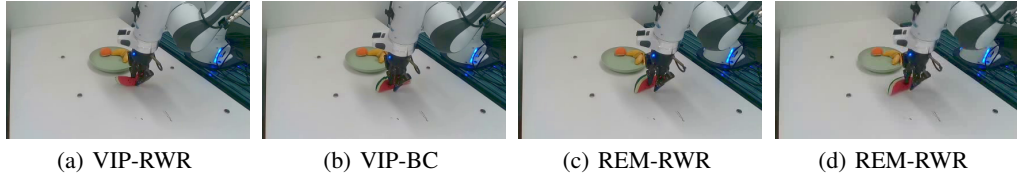


Figure 11: Comparison of failure trajectories on `PickPlaceMelon`. VIP-RWR is still able to reach the critical state of gripping watermelon, whereas baselines fail.

2021). With the advent of VIP-RWR, offline RL may finally be a practical approach for real-world robot learning at scale.

**Performance of R3M-BC.** Our R3M-BC, though able to solve some of the simpler tasks, appears to perform relatively worse than the original R3M-BC in Nair et al. (2022) on their real-world tasks. To account for this discrepancy, we note that our real-world experiment uses different software-hardware stacks and tasks from the original R3M real-world experiments, so the results are not directly comparable. For instance, camera placement, an important variable for real-world robot learning, is chosen differently in our experiment and that of R3M; in R3M, a different camera angle is selected for each task, whereas in our setup, the same camera view is used for all tasks. Furthermore, we emphasize that our focus is not the absolute performance of R3M-BC, but rather the relative improvement R3M-RWR provides on top of R3M-BC.

#### H.4 Qualitative Analysis

In this section, we study several interesting policy behaviors VIP-RWR acquires. Policy videos are included in our supplementary video.

**Robust key action execution.** VIP-RWR is able to execute key actions more robustly than the baselines; this suggests that its reward information helps it identify necessary actions. For example, as shown in Figure 11, on the `PickPlaceMelon` task, failed VIP-RWR rollouts at least have the gripper grasp onto the watermelon, whereas for other baselines, the failed rollouts do not have the watermelon between the gripper and often incorrectly push the watermelon to touch the plate’s outer edge, preventing pick-and-place behavior from being executed.

**Task re-attempt.** We observe that VIP-RWR often learns more robust policies that are able to perform recovery actions when the task is not solved on the first attempt. For instance, in both `CloseDrawer` and `FoldTowel`, there are trials where VIP-RWR fails to close the drawer all the way or pick up the towel edge right away; in either case, VIP-RWR is able to re-attempt and solves the task (see our supplementary video). This is a known advantage of offline RL over BC (Kumar et al., 2022; Levine et al., 2020); however, we only observe this behavior in VIP-RWR and not R3M-RWR, indicating that this advantage of offline RL is only realized when the reward information is sufficiently informative.

## I Additional Results

### I.1 Value-Based Pre-Training Ablation: Least-Square Temporal-Difference

While VIP is the first value-based pre-training approach and significantly outperforms all existing methods, we show that this effectiveness is also unique to VIP and not to training a value function. To this end, we show that a simpler value-based baseline does not perform as well. In particular, we consider Least-Square Temporal-Difference policy *evaluation* (**LSTD**) (Bradtke & Barto, 1996; Sutton & Barto, 2018) to assess the importance of the choice of value-training objective:

$$\min_{\phi} \mathbb{E}_{(o, o', g) \sim D} \left[ \left( \tilde{\delta}_g(o) + \gamma V(\phi(o'); \phi(g)) - V(\phi(o), \phi(g)) \right)^2 \right], \quad (31)$$

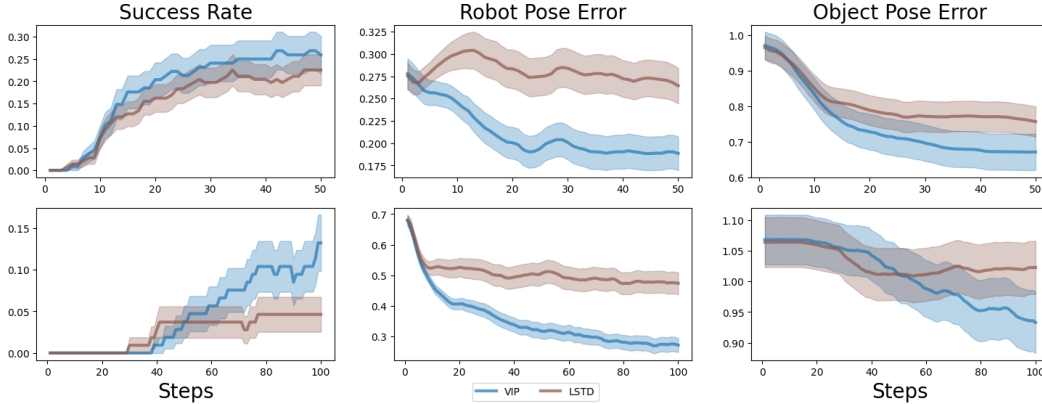


Figure 12: VIP vs. LSTD Trajectory Optimization Comparison.

Table 4: Visual Imitation Learning Results.

	<i>Self-Supervised</i>				<i>Supervised</i>		
	VIP (E)	LSTD (E)	R3M-Lang (E)	MOCO (I)	R3M (E)	ResNet50 (I)	CLIP (Internet)
Success Rate	<b>53.6</b>	51.5	51.2	45.0	<b>55.9</b>	41.8	44.3

740 in which we also parameterize  $V$  as the negative  $L_2$  embedding distance as in VIP. Given that human  
 741 videos are reasonably goal-directed, the value of the human behavioral policy computed via LSTD  
 742 should be a decent choice of reward; however, LSTD does not capture the long-range dependency  
 743 of initial to goal frames (first term in (12)), nor can it obtain a value function that outperforms that  
 744 of the behavioral policy. We train LSTD using the exact same setup as in VIP, differing in only the  
 745 training objective, and compare it against VIP in our trajectory optimization settings.

746 As shown in Fig. 12, interestingly, LSTD already works better than all prior baselines in the Easy  
 747 setting, indicating that value-based pre-training is indeed favorable for reward-specification. However,  
 748 its inability to capture long range temporal dependency as in VIP (the first term in VIP’s objective)  
 749 makes it far less effective on the Hard setting, which require extended smoothness in the reward  
 750 landscape to solve given the distance between the initial observation and the goal. These results  
 751 show that VIP’s superior reward specification comes precisely from its ability to capture both long-  
 752 range temporal dependencies and local temporal smoothness, two innate properties of its dual value  
 753 objective and the associated implicit time contrastive learning interpretation. To corroborate these  
 754 findings, we have also included LSTD in our qualitative reward curve and histogram analysis in  
 755 App. I.4, I.6, and I.7 and finds that VIP generates much smoother embedding than LSTD.

## 756 I.2 Visual Imitation Learning

757 One alternative hypothesis to VIP’s smoother embedding for its superior reward-specification capabil-  
 758 ity is that it learns a better visual representation, which then naturally enables a better visual reward  
 759 function. To investigate this hypothesis, we compare representations’ capability as a pure visual  
 760 encoder in a visual imitation learning setup. We follow the training and evaluation protocol of (Nair  
 761 et al., 2022) and consider 12 tasks combined from FrankaKitchen, MetaWorld (Yu et al., 2020), and  
 762 Adroit (Rajeswaran et al., 2017), 3 camera views for each task, and 3 demonstration dataset sizes,  
 763 and report the aggregate average maximum success rate achieved during training. **R3M-Lang** is the  
 764 publicly released R3M variant without supervised language training. The average success rates over  
 765 all tasks are shown in Table 4; the letter inside  $()$  stands for the pre-training dataset with  $E$  referring  
 766 to Ego4D and  $I$  Imagenet.

767 These results suggest that with current pre-training methods, the performance on visual imitation  
 768 learning may largely be a function of the pre-training dataset, as all methods trained on Ego4D, even  
 769 our simple baseline LSTD, performs comparably and are much better than the next best baseline



770 not trained on Ego4D. Conversely, this result also suggests that despite not being designed for this  
771 purely supervised learning setting, value-based approaches constitute a strong baseline, and VIP is  
772 in fact currently the state-of-art for self-supervised methods. While these results highlight that VIP  
773 is effective even as a pure visual encoder, a necessary requirement for joint effectiveness for visual  
774 reward and representation, it fails to explain why VIP is far superior to R3M in reward-based policy  
775 learning. As such, we conclude that studying representations’ capability as a pure visual encoder  
776 may not be sufficient for distinguishing representations that can additionally perform zero-shot  
777 reward-specification.

### 778 I.3 Embedding and True Rewards Correlation

779 In this section, we create scatterplots of embedding reward vs. true reward on the trajectories MPPI  
780 have generated to assess whether the embedding reward is correlated with the ground-truth dense  
781 reward. More specifically, for each transition in the MPPI trajectories in Figure 2, we plot its reward  
782 under the representation that was used to compute the reward for MPPI versus the true human-crafted  
783 reward computed using ground-truth state information. The dense reward in FrankaKitchen tasks  
784 is a weighted sum of (1) the negative object pose error, (2) the negative robot pose error, (3) bonus  
785 for robot approaching the object, and (4) bonus for object pose error being small. This dense reward  
786 is highly tuned and captures human intuition for how these tasks ought to be best solved. As such,  
787 high correlation indicates that the embedding is able to capture both intuitive robot-centric and  
788 object-centric task progress from visual observations. We only compare VIP and R3M here as a proxy  
789 for comparing our implicit time contrastive mechanism to the standard time contrastive learning.

790 The scatterplots over all tasks and camera views (Easy setting) are shown in Figure 13,14, and 15.  
791 VIP rewards exhibit much greater correlation with the ground-truth reward on its trajectories that  
792 do accomplish task, indicating that when VIP does solve a task, it is solving the task in a way that  
793 matches *human* intuition. This is made possible via large-scale value pre-training on diverse human  
794 videos, which enables VIP to extract a human notion of task-progress that transfers to robot tasks and  
795 domains. These results also suggest that VIP has the potential of *replacing* manual reward engineering,  
796 providing a data-driven solution to the grand challenge of reward engineering for manipulation tasks.  
797 However, VIP is not yet perfect in its current form. Both methods exhibit local minima where high  
798 embedding distances in fact map to lower true rewards; however, this phenomenon is much severe  
799 for R3M. On 8 out of 12 tasks, VIP at least has one camera view in which its rewards are highly  
800 correlated with the ground-truth rewards on its MPPI trajectories.

### 801 I.4 Embedding Distance Curves

802 In Figure 16, we present additional embedding distance curves for all methods on Ego4D and our  
803 real-robot offline RL datasets. For Ego4D, we randomly sample 4 videos of 50-frame long (see  
804 Appendix I.5 for how these short snippets are sampled), and for our robot dataset, we compute the  
805 embedding distance curves for the 4 sample demonstrations in Figure 10. As shown, on all tasks in  
806 the real-robot dataset, VIP is distinctively more smooth than any other representation. This pattern  
807 is less accentuated on Ego4D. This is because a randomly sampled 50-frame snippet from Ego4D  
808 may not coherently represent a task solved from beginning to completion, so an embedding distance  
809 curve is not inherently supposed to be smoothly declining. Nevertheless, VIP still exhibits more local  
810 smoothness in the embedding distance curves, and for the snippets that do solve a task (the first two  
811 videos), it stands out as the smoothest representation.

### 812 I.5 Embedding Distance Curve Bumps

813 In this section, we compute the fraction of negative embedding rewards (equivalently, positive  
814 slopes in embedding embedding distance curves) for each video sequence and average over all video  
815 sequences in a dataset. Each sequence in our robot dataset is of 50 frames, and we use each sequence  
816 without any further truncation. For Ego4D, video sequences are of variable length. For each long  
817 sequence of more than 50 frames, we use the first 50 frames. We do not include videos shorter than

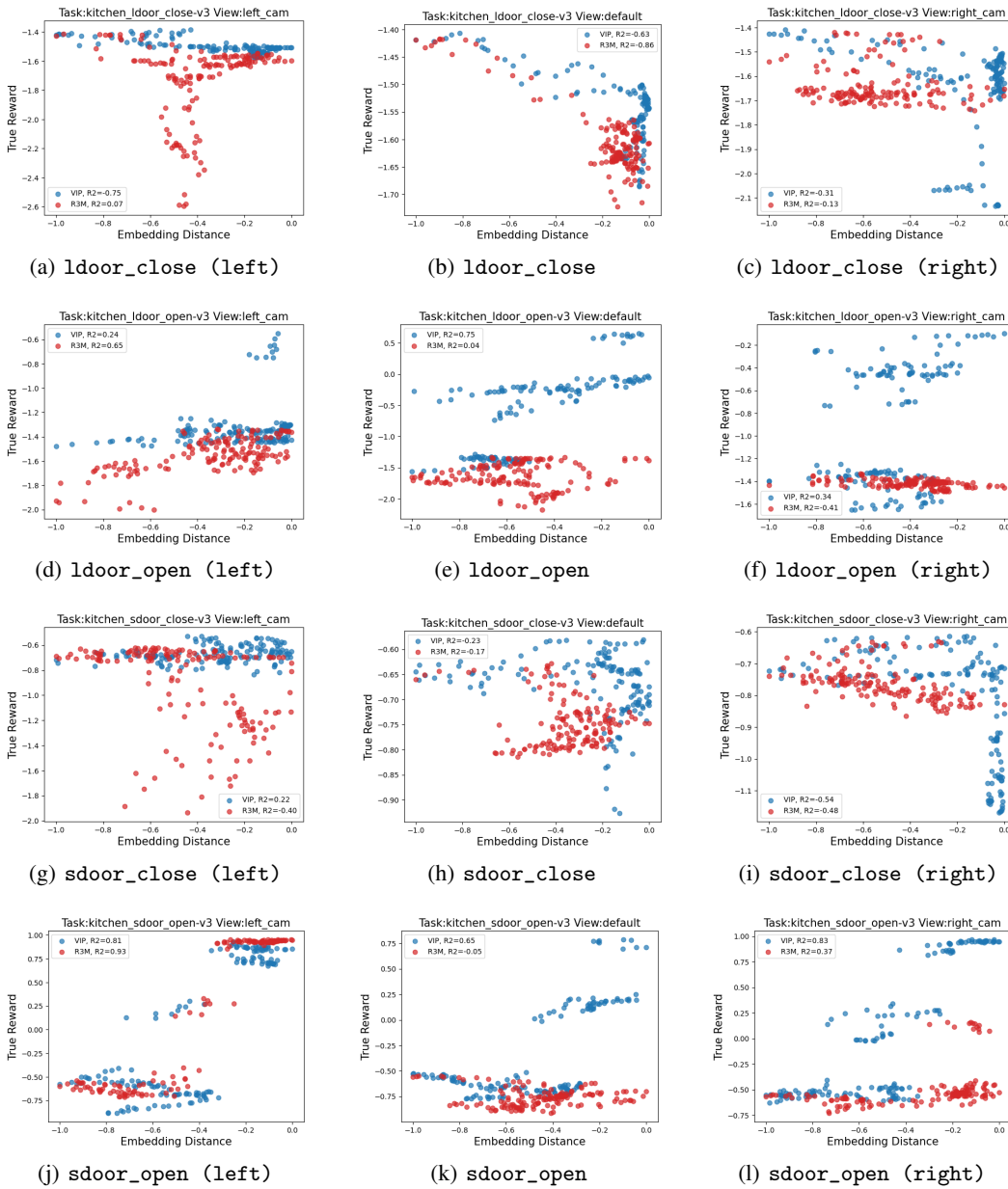


Figure 13: Embedding reward vs. ground-truth human-engineered reward correlation (VIP vs. R3M) part 1.

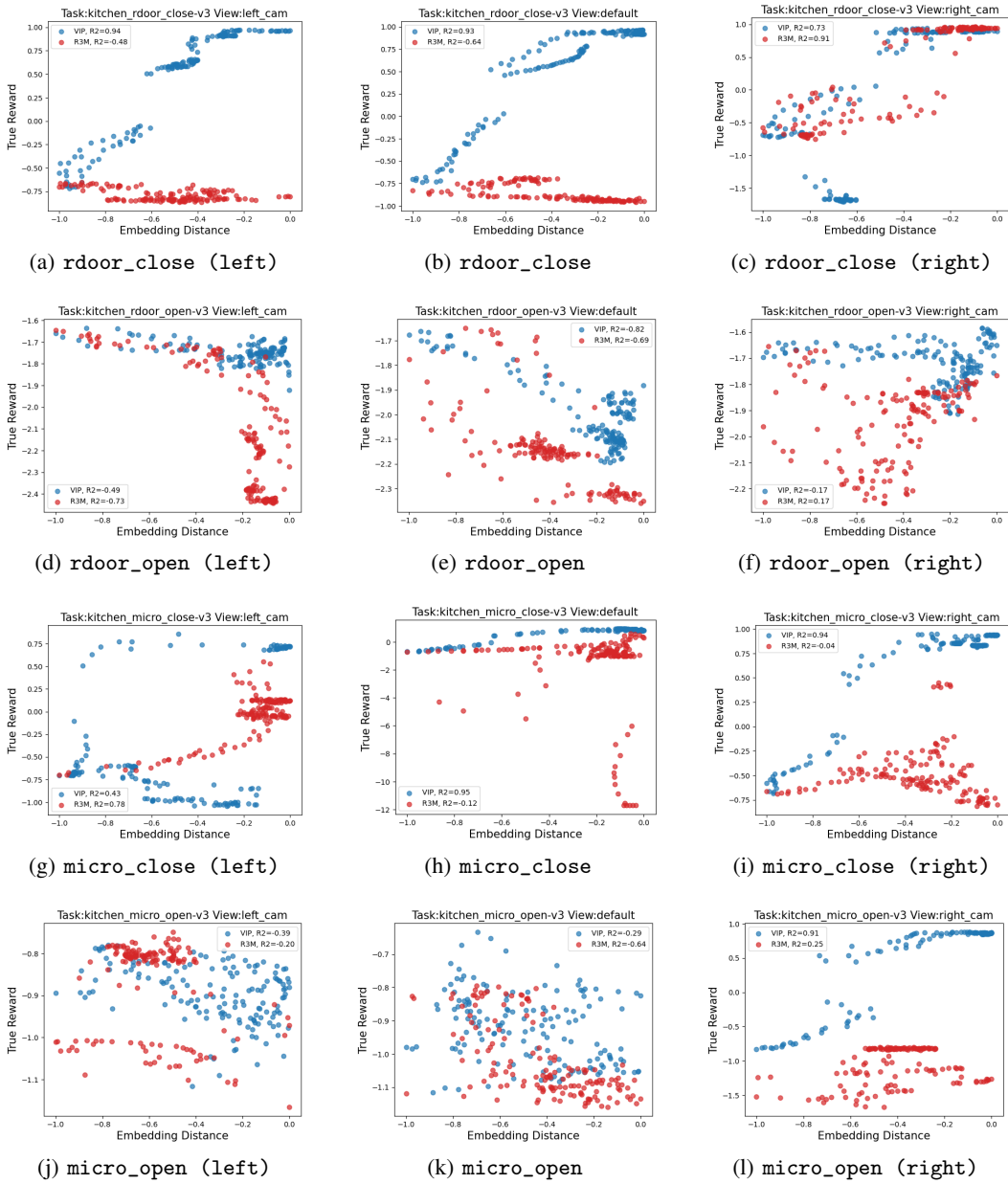


Figure 14: Embedding reward vs. ground-truth human-engineered reward correlation (VIP vs. R3M) part 2.

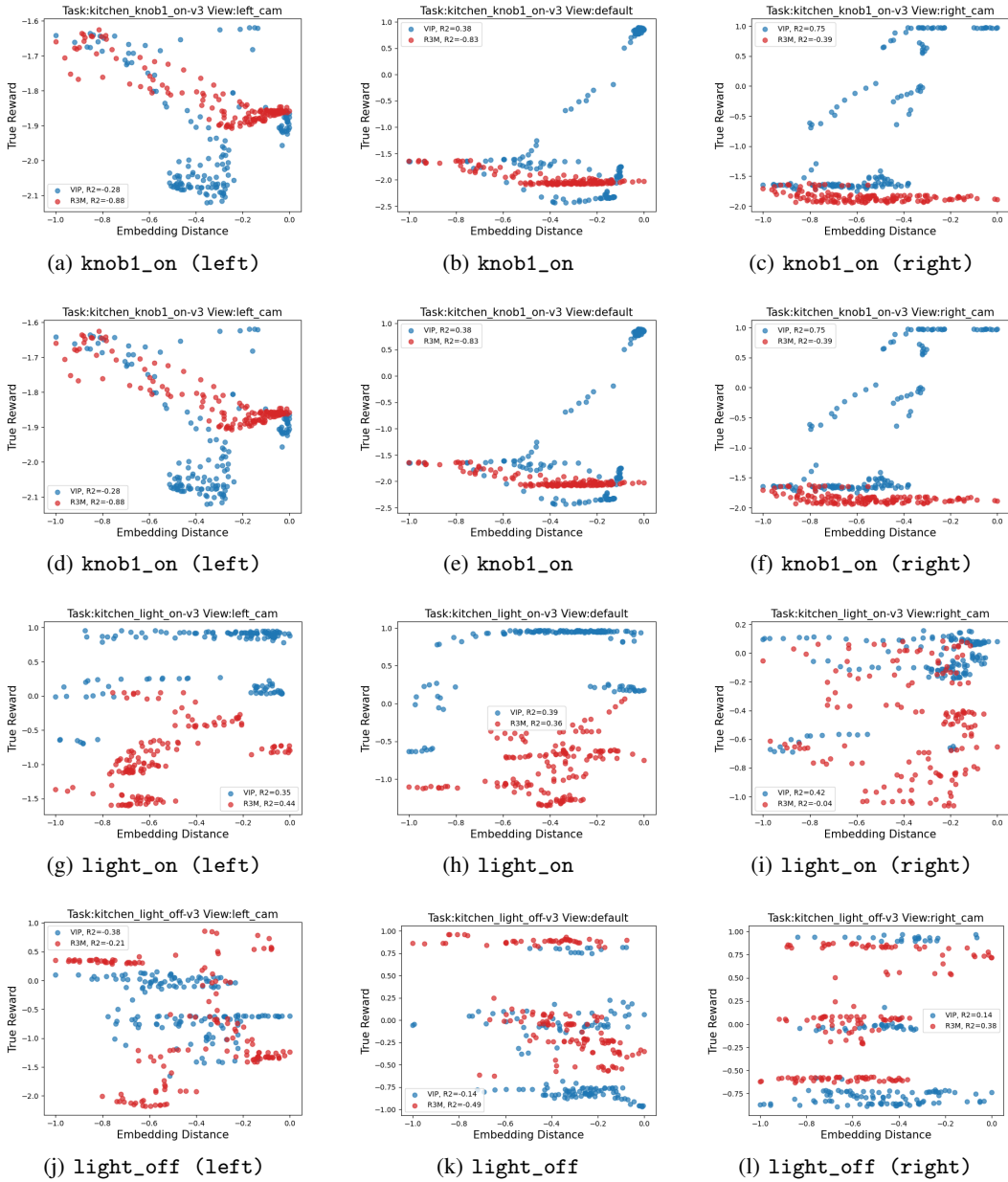
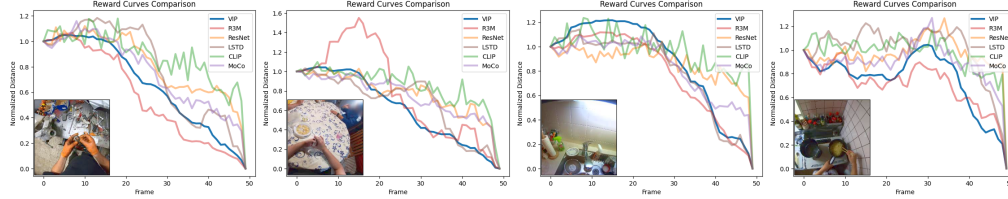
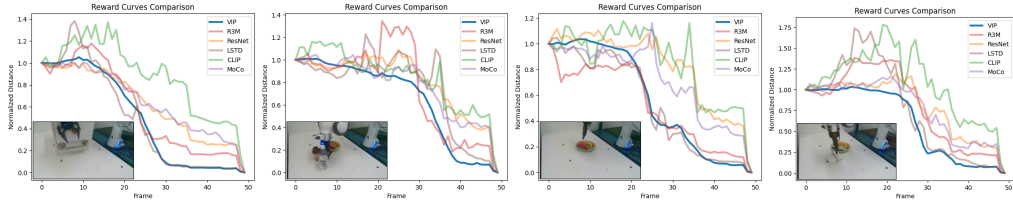


Figure 15: Embedding reward vs. ground-truth human-engineered reward correlation (VIP vs. R3M) part 3.



(a) Ego4D



(b) Real-robot dataset

Figure 16: Additional embedding distance curves on Ego4D and real-robot videos.

Table 5: Proportion of bumps in embedding distance curves.

Dataset	VIP (Ours)	R3M	ResNet50	MOCO	CLIP
Ego4D	$0.253 \pm 0.117$	$0.309 \pm 0.097$	$0.414 \pm 0.052$	$0.398 \pm 0.057$	$0.444 \pm 0.047$
In-House Robot Dataset	$0.243 \pm 0.066$	$0.323 \pm 0.076$	$0.366 \pm 0.046$	$0.380 \pm 0.052$	$0.438 \pm 0.046$

818 50 frames, in order to make the average fraction for each representation comparable between  
 819 the two distinct datasets. Note that for Ego4D, due to its in-the-wild nature, it is not guaranteed that a  
 820 50-frame segment represents one task being solved from beginning to completion, so there may be  
 821 naturally bumps in the embedding distance curve computed with respect to the last frame, as earlier  
 822 frames may not actually be progressing towards the last frame in a goal-directed manner. The full  
 823 results are shown in Table 5. VIP has fewest bumps in Ego4D videos, and this notion of smoothness  
 824 transfer to the robot dataset. Furthermore, since the robot videos are in fact visually simpler and each  
 825 video is guaranteed to be solving one task, the bump rate is actually *lower* despite the domain gap.  
 826 While this observation generally also holds true for other representations, it notably does not hold for  
 827 R3M, which is trained using standard time contrastive learning.

### 828 I.6 Embedding Reward Histograms (Real-Robot Dataset)

829 We present the reward histogram comparison against all baselines in Figure 17. The trend of VIP  
 830 having more small, positive rewards and fewer extreme rewards in either direction is consistent across  
 831 all comparisons.

### 832 I.7 Embedding Reward Histograms (Ego4D)

833 We present the reward histogram comparison against all baselines in Figure 18. The histograms are  
 834 computed using the same set of 50-frame Ego4D video snippets as in Appendix I.5. The y-axis is in  
 835 log-scale due to the large total count of Ego4D frames. As discussed, Ego4D video segments are  
 836 less regular than those in our real-robot dataset, and this irregularity contributes to all representations  
 837 having significantly more negative rewards compared to their histograms on the real-robot dataset.  
 838 Nevertheless, the relative difference ratio’s pattern is consistent, showing VIP having far more  
 839 rewards that lie in the first positive bin. Furthermore, VIP also has significantly fewer extreme  
 840 negative rewards compared to all baselines.

841

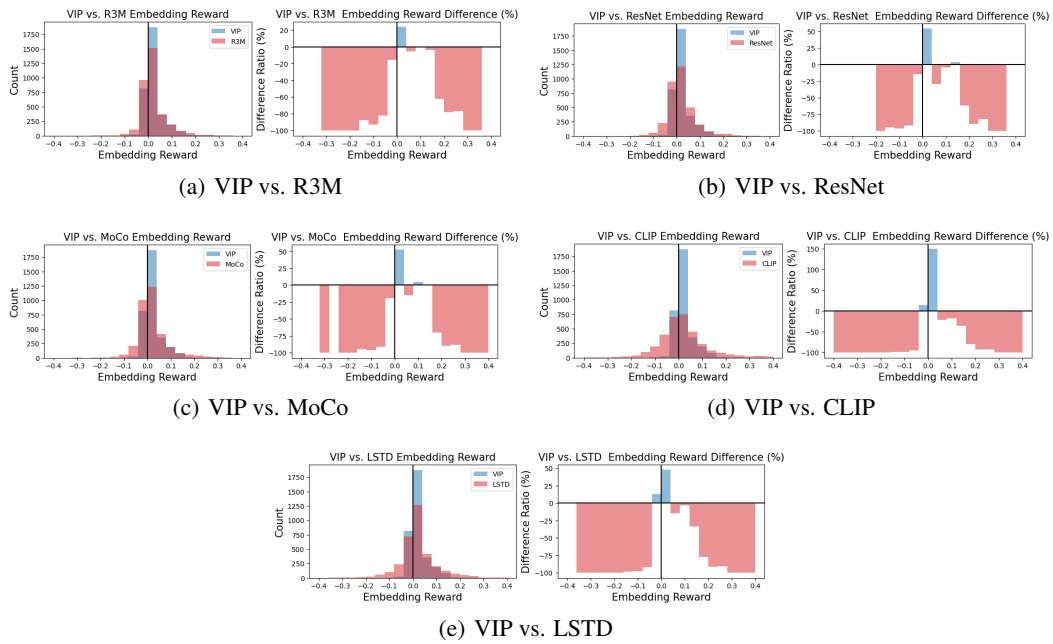


Figure 17: Embedding reward histogram comparison on real-robot dataset.

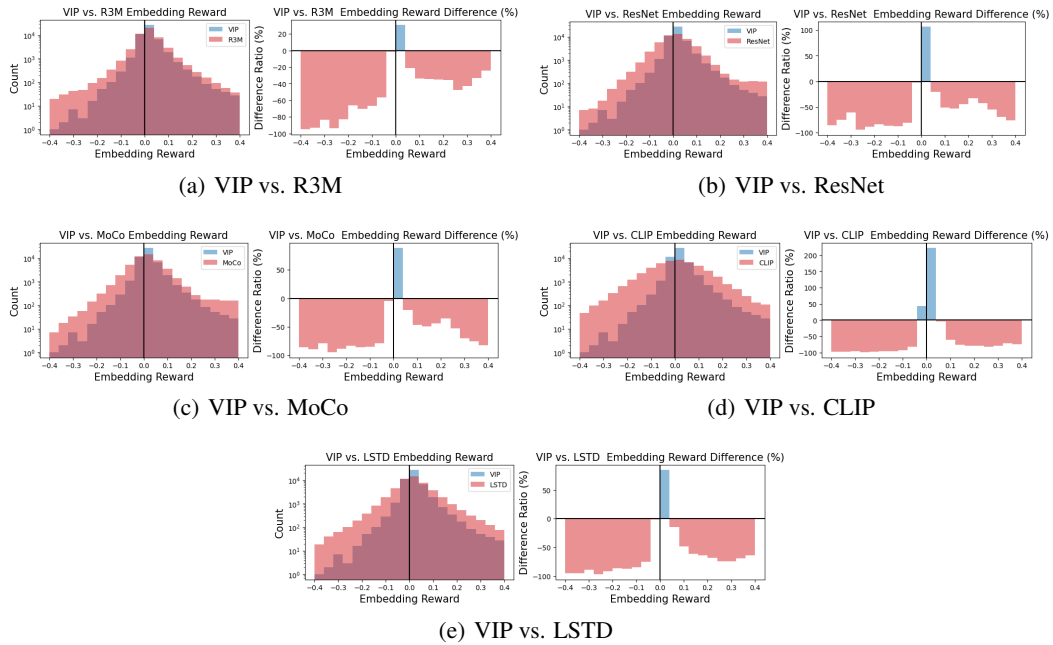


Figure 18: Embedding reward histogram comparison on Ego4D videos.