

Subtask-Aware Visual Reward Learning from Segmented Demonstrations

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Reinforcement Learning (RL) agents have demonstrated their potential
2 across various robotic tasks. However, they still heavily rely on human-engineered
3 reward functions, requiring extensive trial-and-error and access to target behav-
4 ior information, often unavailable in real-world settings. This paper introduces
5 REDS: *REward learning from Demonstration with Segmentations*, a novel reward
6 learning framework that leverages action-free videos with minimal supervision.
7 Specifically, REDS employs video demonstrations segmented into subtasks from
8 diverse sources and treats these segments as ground-truth rewards. We train a
9 dense reward function conditioned on video segments and their corresponding
10 subtasks to ensure alignment with ground-truth reward signals by minimizing the
11 Equivalent-Policy Invariant Comparison distance. Additionally, we employ con-
12 trastive learning objectives to align video representations with subtasks, ensuring
13 precise subtask inference during online interactions. Our experiments show that
14 REDS significantly outperforms baseline methods on complex robotic manipula-
15 tion tasks in Meta-World and more challenging real-world tasks, such as furniture
16 assembly in FurnitureBench, with minimal human intervention.

17 1 Introduction

18 Reinforcement Learning (RL) has demonstrated significant potential for training autonomous agents
19 in various real-world robotic tasks, provided that appropriate reward functions are available [1, 2, 3,
20 4, 5]. However, reward engineering typically requires substantial trial-and-error [6, 7] and extensive
21 task knowledge, often necessitating specialized instrumentation (*e.g.*, motion trackers [8] or tactile
22 sensors [9]) or detailed information about target objects [10, 11, 12, 13, 14, 15], which are difficult
23 to obtain in real-world settings. Learning reward functions from action-free videos has emerged as a
24 promising alternative, as it avoids the need for detailed action annotations or precise target behavior
25 information, and video data can be easily collected from online sources [16, 17, 18]. Approaches in
26 this domain include learning discriminators between video demonstrations and policy rollouts [19,
27 20], training temporally aligned visual representations from large-scale video datasets [21, 22, 23,
28 24, 25] to estimate reward based on distance to a goal image, and using video prediction models to
29 generate reward signals [26, 27].

30 Despite this progress, existing methods often struggle with long-horizon, complex robotic tasks that
31 involve multiple subtasks. These approaches typically fail to provide context-aware reward signals,
32 relying only on a few consecutive frames or the final goal image without considering subsequent
33 subtasks. For example, in One Leg task (see Figure 2d) from FurnitureBench [28], prior methods
34 often overemphasize the reward for picking up the leg while neglecting crucial steps such as inserting
35 the leg into a hole and tightening it. Recent work [29] proposes a discriminator-based approach that
36 treats complex tasks as a sequence of subtasks. However, it assumes that the environment provides
37 explicit subtask identification, which often demands significant human intervention in real-world
38 scenarios. Moreover, discriminator-based methods are known to be prone to mode collapse [30,
39 31]. Consequently, designing an effective visual reward function for real-world, long-horizon tasks
40 remains an open problem.

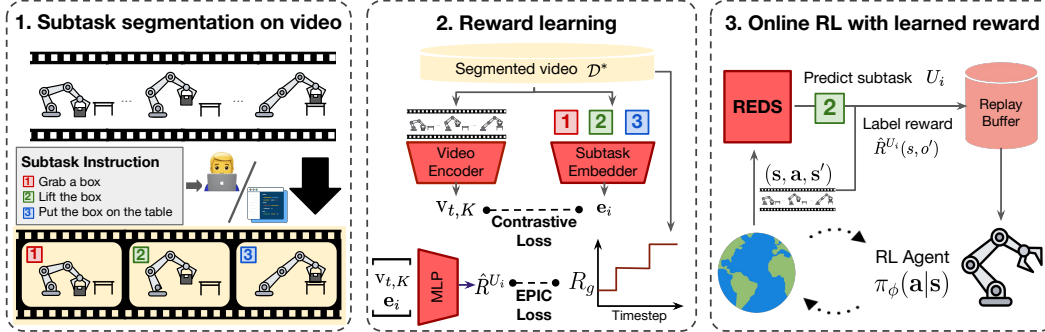


Figure 1: Illustration of REDS. Our main idea is to leverage expert demonstrations annotated with the ongoing subtask as the source of implicit reward signals (left). We train a reward model conditioned on video segments and corresponding subtasks with 1) contrastive loss to attract the video segments and corresponding subtask embeddings and 2) EPIC [32] loss to generate reward equivalent to subtask segmentations (middle). In online RL, REDS infers the ongoing subtask using only video segments at each timestep and computes the reward with that (right).

41 **Our approach** To address the aforementioned limitations, we propose a novel reward learning
 42 framework, REDS: *REward learning from Demonstration with Segmentations*, which infers subtask
 43 information from video segments and generates corresponding reward signals for each subtask. The
 44 key idea is to employ minimal supervision to produce appropriate reward signals for intermediate
 45 subtask completion. Specifically, REDS utilizes expert demonstrations, where subtasks are anno-
 46 tated at each timestep by various sources (e.g., human annotators, code snippets, vision-language
 47 models; see the left figure of Figure 1). These annotations serve as ground-truth rewards. For train-
 48 ing, we introduce a new objective function minimizing the Equivalent-Policy Invariant Comparison
 49 (EPIC) [32] between the learned reward function and the ground-truth rewards, guaranteeing a theo-
 50 retical upper bound on regret relative to the ground-truth reward function. Additionally, to correctly
 51 infer the ongoing subtask in online interactions, we adopt a contrastive learning objective to align
 52 video representations with task embeddings. In terms of architecture, our reward model is designed
 53 to capture temporal dependencies in video segments using transformers [33], leading to enhanced
 54 reward signal quality.

55 **Contributions** We present a novel visual reward learning framework REDS: *REward learning*
 56 *from Demonstration with Segmentations*, which can produce suitable reward signals aware of sub-
 57 tasks in long-horizon complex robotic manipulation tasks. We show that REDS significantly out-
 58 performs baselines in training RL agents for robotic manipulation tasks in Meta-world, and even
 59 surpasses dense reward functions in some tasks. Furthermore, we demonstrate that REDS can
 60 train real-world RL agents to perform long-horizon complex furniture assembly tasks from Fur-
 61 nitureBench.

62 2 REward learning from Demonstration with Segmentations

63 2.1 Intuition

64 The sparse reward function R provides feedback only on the overall success or failure of a task,
 65 which is insufficient for guiding the agent through intermediate states. To address this, drawing
 66 inspiration from previous work on long-range robotic manipulation tasks [28, 29], we decompose
 67 the task into m subtasks, denoted as $\mathcal{U} = \{U_1, \dots, U_m\}$, using domain knowledge. Each subtask U_i
 68 represents a distinct step in the task sequence, where i indicates its order. For task success, the agent
 69 must complete these subtasks in sequence, meaning subtask U_i must be finished before moving to
 70 U_{i+1} .¹ To further guide the agent, we provide text instructions $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^m$, describing how to
 71 solve each subtask. To obtain subtask segmentations, we map each observation o_t in the trajectory

¹For instance, Door Open can be divided into (i) reaching the door handle and (ii) pulling the door to the target position.

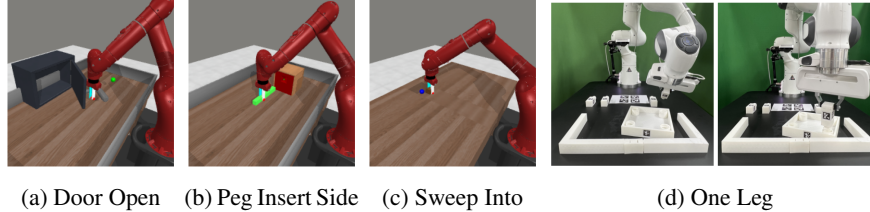


Figure 2: Examples of visual observations used in our experiments. We consider a variety of robotic manipulation tasks from Meta-world [12] and FurnitureBench [28].

72 $\tau = (o_0, \dots, o_T)$ to its corresponding subtask using a segmentation function $\psi : \mathcal{O} \rightarrow \mathcal{U}$ from
 73 various sources.

74 2.2 Reward Modeling

75 **Architecture** As mentioned in Section 1, previous reward learning methods generate rewards only
 76 by a single frame or consequent frames, not taking into account the order of subtasks. To resolve
 77 the issue, we propose a new reward predictor $\hat{R}^U = \hat{R}(s; U)$ conditioned on each subtask.

78 **Reward equivariance with subtask segmentation** Our key insight is that the subtask segmenta-
 79 tion function ψ can be thought of as the ground-truth reward function, providing implicit signals for
 80 solving intermediate tasks. To ensure our reward function induces the same set of optimal policies
 81 as ψ , we train to minimize EPIC [32] distance between our reward model \hat{R}_θ^U parameterized by θ
 82 and ψ for all subtasks.

83 **Progressive reward signal** However, minimizing EPIC with ψ alone can lead to overfitting and
 84 the inability to provide progressive signals within each subtask. To mitigate this issue, we propose an
 85 additional regularization term to enforce progressive reward signals. Inspired by previous work [34,
 86 35, 36], we view the reward function as a progress indicator for each subtask, and we regularize the
 87 reward function output to be higher in later states of expert demonstration.

88 **Aligning video representation with subtask embeddings** As the reward model lacks informa-
 89 tion about the ongoing subtasks in online interactions, it must infer the agent’s current subtask. To
 90 achieve this, we train the video representation to be closely aligned with the corresponding sub-
 91 task embedding by adopting a contrastive learning objective. The model can select the appropriate
 92 subtask embedding only by the video segment.

93 3 Experiments

94 3.1 Meta-world Experiments

95 **Setup** We first evaluate our method on 8 different visual robotic manipulation tasks from Meta-
 96 world [12]. As a backbone algorithm, we use DreamerV3 [37], a state-of-the-art visual model-based
 97 RL algorithm that learns from latent imaginary rollouts. For collecting subtask segmentations, we
 98 utilize a scripted teacher in simulation environments for scalability. Specifically, we use the pre-
 99 defined indicator for subtasks provided in the benchmark for all subtask segmentations (see Ap-
 100 pendix E for the list of subtasks and corresponding text instructions for each task). We do not use
 101 these indicators when training/evaluating RL agents. For training REDS, we first collect subtask
 102 segmentations from 50 expert demonstrations for initial training and train DreamerV3 agents for
 103 100K environment steps with the initial reward model to collect suboptimal trajectories, which is
 104 used for fine-tuning. In evaluation, we measure the success rate averaged over 10 episodes in every
 105 20K steps. Please refer to Appendix B for more details.

106 **Results** Figure 3 shows that REDS consistently improves the sample-efficiency of DreamerV3
 107 agents by outperforming all baselines. While baselines exhibit non-zero success rates in simple
 108 tasks like Faucet Close, their performance significantly deteriorates in more complex tasks, such as
 109 Peg Insert Side. On the other hand, our method maintains non-zero success rates across all tasks and
 110 even surpasses human-engineered reward functions in some tasks (e.g., Drawer Open, Push, Cof-

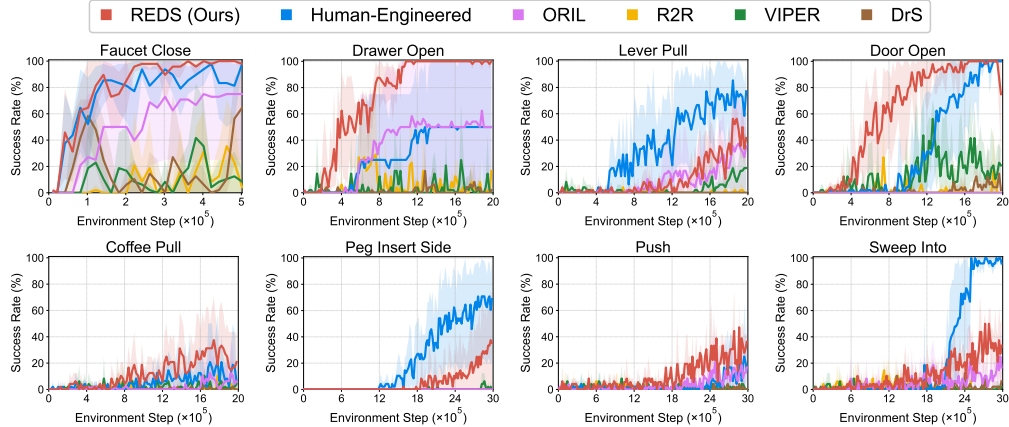


Figure 3: Learning curves of DreamerV3 [37] agents trained with different reward functions for solving eight robotic manipulation tasks from Meta-world [12], measured by success rate (%). The solid line and shaded regions represent the mean and stratified bootstrap interval across 4 runs.

111 fee Pull) without requiring task-specific reward engineering. These results show that REDS effectively
 112 generates appropriate rewards for solving intermediate tasks by leveraging subtask-segmented
 113 demonstrations.

114 3.2 FurnitureBench Experiments

115 **Setup** We further evaluate our method on real-world furniture assembly tasks from Furni-
 116 tureBench [28], specifically focusing on One Leg assembly. This task involves a sequence of com-
 117 plex subtasks such as picking up, inserting, and screwing (see Figure 2d). For training REDS, we
 118 use 300 expert demonstrations with subtask segmentations provided by FurnitureBench, along with
 119 an additional 200 rollouts from IQL [38] policy trained with expert demonstrations in a single train-
 120 ing iteration. To prevent misleading reward signals stemming from visual occlusions, we utilize
 121 visual observations from the front camera and wrist cameras in training REDS. For downstream RL,
 122 we first train offline RL agents using 300 expert demonstrations labeled with each reward model,
 123 followed by online fine-tuning to assess improvements. We provide more details in Appendix B.

124 **Results** As shown in Table 1, REDS achieves significant performance improvements through online
 125 fine-tuning, whereas the improvements from baselines are marginal. These results indicate that our
 126 method produces informative signals for solving a sequence of subtasks, while baselines either fail to
 127 provide context-aware signals or dense rewards for better exploration (see Appendix H for qualitative
 128 examples). Moreover, we note that our method outperforms the IQL trained with 500 expert demon-
 129 strations, achieving a score of 2.45 compared to 1.8 reported by FurnitureBench, despite using only 300
 130 expert demonstrations. Considering that REDS does not require additional human interventions beyond
 131 resetting the environment, these results highlight the potential to extend our approach to a wider
 132 range of real-world robotics tasks.

Table 1: Online fine-tuning results of IQL agents in One Leg from FurnitureBench. We report the initial performance after offline RL (left) and the performance after 150 episodes of online RL (right).

Method	# Expert Demos	Completed Subtasks (Offline → Online)
Sparse (Offline) [28]	500	1.8
VIPER	300	1.10 → 1.25
DrS	300	1.05 → 1.10
REDS (Ours)	300	1.10 → 2.45

138 4 Conclusion

139 We proposed REDS, a visual reward learning framework considering subtasks by utilizing subtask
 140 segmentation. Our main contribution is based on proposing a new reward model leveraging minimal
 141 domain knowledge as a ground-truth reward function. Our approach is generally applicable and
 142 does not require any additional instrumentations in online interactions. We believe REDS will sig-
 143 nificantly alleviate the burden of reward engineering and facilitate the application of RL to a broader
 144 range of real-world robotic tasks.

References

- [1] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 2016.
- [2] S. Gu, E. Holly, T. Lillicrap, and S. Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *IEEE International Conference on Robotics and Automation*, 2017.
- [3] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *International Journal of Robotics Research*, 2020.
- [4] L. Smith, I. Kostrikov, and S. Levine. A walk in the park: Learning to walk in 20 minutes with model-free reinforcement learning. In *Robotics: Science and Systems*, 2023.
- [5] A. Handa, A. Allshire, V. Makoviychuk, A. Petrenko, R. Singh, J. Liu, D. Makoviichuk, K. Van Wyk, A. Zhurkevich, B. Sundaralingam, et al. Dextreme: Transfer of agile in-hand manipulation from simulation to reality. In *IEEE International Conference on Robotics and Automation*, 2023.
- [6] S. Booth, W. B. Knox, J. Shah, S. Niekum, P. Stone, and A. Allievi. The perils of trial-and-error reward design: misdesign through overfitting and invalid task specifications. In *AAAI Conference on Artificial Intelligence*, 2023.
- [7] W. B. Knox, A. Allievi, H. Banzhaf, F. Schmitt, and P. Stone. Reward (mis) design for autonomous driving. *Artificial Intelligence*, 2023.
- [8] X. B. Peng, E. Coumans, T. Zhang, T.-W. Lee, J. Tan, and S. Levine. Learning agile robotic locomotion skills by imitating animals. *arXiv preprint arXiv:2004.00784*, 2020.
- [9] Y. Yuan, H. Che, Y. Qin, B. Huang, Z.-H. Yin, K.-W. Lee, Y. Wu, S.-C. Lim, and X. Wang. Robot synesthesia: In-hand manipulation with visuotactile sensing. *arXiv preprint arXiv:2312.01853*, 2023.
- [10] S. James, Z. Ma, D. R. Arrojo, and A. J. Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 2020.
- [11] Y. Zhu, J. Wong, A. Mandlekar, R. Martín-Martín, A. Joshi, S. Nasiriany, and Y. Zhu. robosuite: A modular simulation framework and benchmark for robot learning. *arXiv preprint arXiv:2009.12293*, 2020.
- [12] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, 2020.
- [13] T. Mu, Z. Ling, F. Xiang, D. C. Yang, X. Li, S. Tao, Z. Huang, Z. Jia, and H. Su. Maniskill: Generalizable manipulation skill benchmark with large-scale demonstrations. In *Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [14] J. Gu, F. Xiang, X. Li, Z. Ling, X. Liu, T. Mu, Y. Tang, S. Tao, X. Wei, Y. Yao, X. Yuan, P. Xie, Z. Huang, R. Chen, and H. Su. Maniskill2: A unified benchmark for generalizable manipulation skills. In *International Conference on Learning Representations*, 2023.
- [15] C. Sferrazza, D.-M. Huang, X. Lin, Y. Lee, and P. Abbeel. Humanoidbench: Simulated humanoid benchmark for whole-body locomotion and manipulation. *arXiv preprint arXiv:2403.10506*, 2024.
- [16] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.

- 189 [17] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola,
190 T. Green, T. Back, P. Natsev, et al. The kinetics human action video dataset. *arXiv preprint*
191 *arXiv:1705.06950*, 2017.
- 192 [18] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti,
193 J. Munro, T. Perrett, W. Price, and M. Wray. Scaling egocentric vision: The epic-kitchens
194 dataset. In *European Conference on Computer Vision*, 2018.
- 195 [19] A. S. Chen, S. Nair, and C. Finn. Learning generalizable robotic reward functions from” in-
196 the-wild” human videos. In *Robotics: Science and Systems*, 2021.
- 197 [20] D. Yang, D. Tjia, J. Berg, D. Damen, P. Agrawal, and A. Gupta. Rank2reward: Learning
198 shaped reward functions from passive video. In *IEEE International Conference on Robotics*
199 *and Automation*, 2024.
- 200 [21] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain. Time-
201 contrastive networks: Self-supervised learning from video. In *IEEE International Conference*
202 *on Robotics and Automation*, 2018.
- 203 [22] K. Zakka, A. Zeng, P. Florence, J. Tompson, J. Bohg, and D. Dwibedi. Xirl: Cross-embodiment
204 inverse reinforcement learning. In *Conference on Robot Learning*, 2021.
- 205 [23] S. Kumar, J. Zamora, N. Hansen, R. Jangir, and X. Wang. Graph inverse reinforcement learning
206 from diverse videos. In *Conference on Robot Learning*. PMLR, 2022.
- 207 [24] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang. VIP: Towards
208 universal visual reward and representation via value-implicit pre-training. In *International*
209 *Conference on Learning Representations*, 2023.
- 210 [25] Y. J. Ma, W. Liang, V. Som, V. Kumar, A. Zhang, O. Bastani, and D. Jayaraman. Liv:
211 Language-image representations and rewards for robotic control. In *International Conference*
212 *on Machine Learning*, 2023.
- 213 [26] A. Escontrela, A. Adeniji, W. Yan, A. Jain, X. B. Peng, K. Goldberg, Y. Lee, D. Hafner, and
214 P. Abbeel. Video prediction models as rewards for reinforcement learning. In *Conference on*
215 *Neural Information Processing Systems*, 2023.
- 216 [27] T. Huang, G. Jiang, Y. Ze, and H. Xu. Diffusion reward: Learning rewards via conditional
217 video diffusion. In *European Conference on Computer Vision*, 2024.
- 218 [28] M. Heo, Y. Lee, D. Lee, and J. J. Lim. Furniturebench: Reproducible real-world benchmark
219 for long-horizon complex manipulation. In *Robotics: Science and Systems*, 2023.
- 220 [29] T. Mu, M. Liu, and H. Su. Drs: Learning reusable dense rewards for multi-stage tasks. In
221 *International Conference on Learning Representations*, 2024.
- 222 [30] Z. Wang, J. S. Merel, S. E. Reed, N. de Freitas, G. Wayne, and N. Heess. Robust imitation of
223 diverse behaviors. In *Conference on Neural Information Processing Systems*, 2017.
- 224 [31] K. Zolna, S. Reed, A. Novikov, S. G. Colmenarejo, D. Budden, S. Cabi, M. Denil, N. de Freitas,
225 and Z. Wang. Task-relevant adversarial imitation learning. In *Conference on Robot Learning*.
226 PMLR, 2021.
- 227 [32] A. Gleave, M. D. Dennis, S. Legg, S. Russell, and J. Leike. Quantifying differences in reward
228 functions. In *International Conference on Learning Representations*, 2021.
- 229 [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polo-
230 sukhin. Attention is all you need. In *Conference on Neural Information Processing Systems*,
231 2017.

- 232 [34] Y. Lee, A. Szot, S.-H. Sun, and J. J. Lim. Generalizable imitation learning from observation
233 via inferring goal proximity. In *Conference on Neural Information Processing Systems*, 2021.
- 234 [35] K. Hartikainen, X. Geng, T. Haarnoja, and S. Levine. Dynamical distance learning for semi-
235 supervised and unsupervised skill discovery. In *International Conference on Learning Repre-
236 sentations*, 2020.
- 237 [36] Z. Wu, W. Lian, V. Unhelkar, M. Tomizuka, and S. Schaal. Learning dense rewards for contact-
238 rich manipulation tasks. In *IEEE International Conference on Robotics and Automation*, 2021.
- 239 [37] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap. Mastering diverse domains through world
240 models. *arXiv preprint arXiv:2301.04104*, 2023.
- 241 [38] I. Kostrikov, A. Nair, and S. Levine. Offline reinforcement learning with implicit q-learning.
242 In *International Conference on Learning Representations*, 2022.
- 243 [39] A. Pan, K. Bhatia, and J. Steinhardt. The effects of reward misspecification: Mapping and mit-
244 igating misaligned models. In *International Conference on Learning Representations*, 2022.
- 245 [40] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
246 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervi-
247 sion. In *International Conference on Machine Learning*. PMLR, 2021.
- 248 [41] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding
249 by generative pre-training. 2018.
- 250 [42] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Confer-
251 ence on Learning Representations*, 2019.
- 252 [43] D. Yarats, I. Kostrikov, and R. Fergus. Image augmentation is all you need: Regularizing deep
253 reinforcement learning from pixels. In *International Conference on Learning Representations*,
254 2021.
- 255 [44] D. Yarats, R. Fergus, A. Lazaric, and L. Pinto. Mastering visual continuous control: Improved
256 data-augmented reinforcement learning. In *International Conference on Machine Learning*,
257 2022.
- 258 [45] Y. Seo, D. Hafner, H. Liu, F. Liu, S. James, K. Lee, and P. Abbeel. Masked world models for
259 visual control. In *Conference on Robot Learning*. PMLR, 2022.
- 260 [46] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual represen-
261 tation for robot manipulation. In *Conference on Robot Learning*, 2022.
- 262 [47] P. J. Ball, L. Smith, I. Kostrikov, and S. Levine. Efficient online reinforcement learning with
263 offline data. In *International Conference on Machine Learning*, 2023.
- 264 [48] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*,
265 2016.
- 266 [49] K. Zolna, A. Novikov, K. Konyushkova, C. Gulcehre, Z. Wang, Y. Aytar, M. Denil, N. de Fre-
267 itas, and S. Reed. Offline learning from demonstrations and unlabeled experience. In *Confer-
268 ence on Neural Information Processing Systems*, 2020.
- 269 [50] D. Xu and M. Denil. Positive-unlabeled reward learning. In *Conference on Robot Learning*,
270 2021.
- 271 [51] W. Yan, Y. Zhang, P. Abbeel, and A. Srinivas. Videogpt: Video generation using vq-vae and
272 transformers. *arXiv preprint arXiv:2104.10157*, 2021.
- 273 [52] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet. Video diffusion
274 models. In *Conference on Neural Information Processing Systems*, 2022.

- 275 [53] A. Y. Ng and S. Russell. Algorithms for inverse reinforcement learning. In *International*
276 *Conference on Machine Learning*, 2000.
- 277 [54] P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *Inter-*
278 *national Conference on Machine Learning*, 2004.
- 279 [55] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey, et al. Maximum entropy inverse reinforce-
280 ment learning. In *AAAI Conference on Artificial Intelligence*, 2008.
- 281 [56] J. Ho and S. Ermon. Generative adversarial imitation learning. In *Conference on Neural*
282 *Information Processing Systems*, 2016.
- 283 [57] J. Fu, K. Luo, and S. Levine. Learning robust rewards with adversarial inverse reinforcement
284 learning. In *International Conference on Learning Representations*, 2018.
- 285 [58] K. Zolna, S. Reed, A. Novikov, S. G. Colmenarejo, D. Budden, S. Cabi, M. Denil, N. de Freitas,
286 and Z. Wang. Task-relevant adversarial imitation learning. In *Conference on Robot Learning*,
287 2021.
- 288 [59] B. Wulfe, L. M. Ellis, J. Mercat, R. T. McAllister, and A. Gaidon. Dynamics-aware comparison
289 of learned reward functions. In *International Conference on Learning Representations*, 2022.
- 290 [60] J. M. V. Skalse, L. Farnik, S. R. Motwani, E. Jenner, A. Gleave, and A. Abate. STARC:
291 A general framework for quantifying differences between reward functions. In *International*
292 *Conference on Learning Representations*, 2024.
- 293 [61] J. Rocamonde, V. Montesinos, E. Nava, E. Perez, and D. Lindner. Vision-language models are
294 zero-shot reward models for reinforcement learning. In *International Conference on Learning*
295 *Representations*, 2024.
- 296 [62] A. Adeniji, A. Xie, and P. Abbeel. Skill-based reinforcement learning with intrinsic reward
297 matching. *arXiv preprint arXiv:2210.07426*, 2022.
- 298 [63] X. Liang, K. Shu, K. Lee, and P. Abbeel. Reward uncertainty for exploration in preference-
299 based reinforcement learning. In *International Conference on Learning Representations*, 2022.
- 300 [64] S. James and A. J. Davison. Q-attention: Enabling efficient learning for vision-based robotic
301 manipulation. *IEEE Robotics and Automation Letters*, 2022.
- 302 [65] S. James, K. Wada, T. Laidlow, and A. J. Davison. Coarse-to-fine q-attention: Efficient learning
303 for visual robotic manipulation via discretisation. In *IEEE Conference on Computer Vision and*
304 *Pattern Recognition*, 2022.
- 305 [66] M. Shridhar, L. Manuelli, and D. Fox. Perceiver-actor: A multi-task transformer for robotic
306 manipulation. In *Conference on Robot Learning*, 2022.
- 307 [67] L. X. Shi, A. Sharma, T. Z. Zhao, and C. Finn. Waypoint-based imitation learning for robotic
308 manipulation. In *Conference on Robot Learning*, 2023.
- 309 [68] Z. Zhang, Y. Li, O. Bastani, A. Gupta, D. Jayaraman, Y. J. Ma, and L. Weihs. Universal visual
310 decomposer: Long-horizon manipulation made easy. In *IEEE International Conference on*
311 *Robotics and Automation*, 2024.
- 312 [69] L. Kou, F. Ni, Y. Zheng, J. Liu, Y. Yuan, Z. Dong, and H. Jianye. Kisa: A unified keyframe
313 identifier and skill annotator for long-horizon robotics demonstrations. In *International Con-*
314 *ference on Machine Learning*, 2024.
- 315 [70] A. Xie, L. Lee, T. Xiao, and C. Finn. Decomposing the generalization gap in imitation learning
316 for visual robotic manipulation. In *IEEE International Conference on Robotics and Automa-*
317 *tion*, 2024.

- 318 [71] A. Padalkar, A. Pooley, A. Jain, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai,
319 A. Singh, A. Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x mod-
320 els. *arXiv preprint arXiv:2310.08864*, 2023.
- 321 [72] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany,
322 M. K. Srirama, L. Y. Chen, K. Ellis, et al. Droid: A large-scale in-the-wild robot manipulation
323 dataset. *arXiv preprint arXiv:2403.12945*, 2024.
- 324 [73] S. Bahl, R. Mendonca, L. Chen, U. Jain, and D. Pathak. Affordances from human videos as
325 a versatile representation for robotics. In *IEEE Conference on Computer Vision and Pattern*
326 *Recognition*, 2023.
- 327 [74] C. Devin, P. Abbeel, T. Darrell, and S. Levine. Deep object-centric representations for gener-
328 alizable robot learning. In *IEEE International Conference on Robotics and Automation*, 2018.

329 A REDS Training and Inference

330 **Inference** For each transition (s_t, a, s_{t+1}) at timestep t , we compute the reward using s_{t+1} . To
331 infer ongoing subtasks in REDS, we first encode the visual observations from executed actions and
332 the history of previous observations using a pre-trained visual encoder and a causal transformer.
333 The transformer’s final output, \mathbf{v}_t , is used to predict the subtask. REDS selects the subtask index \hat{i}
334 by choosing the subtask embedding $e_{\hat{i}}$ resulting in the highest cosine similarity with \mathbf{v}_t . The final
335 reward is computed as $\hat{R}_{\theta}(s_{t-1}, o_t; U_{\hat{i}})$. Please refer to Appendix B for more details.

336 **Training** We outline the training procedure for REDS. First, we collect subtask segmentations
337 from expert demonstrations, creating a dataset \mathcal{D}^0 , and use it to train the initial reward model,
338 M^0 . However, reward models trained solely on expert data are susceptible to reward misspecifica-
339 tion [39]. To address this, we iteratively collect suboptimal demonstrations and fine-tune the reward
340 model using expert and suboptimal data. Unlike expert demonstrations, suboptimal demonstrations
341 cover a broader range of states and more diverse observations, making manual segmentation labor-
342 intensive and error-prone. To reduce the burden on human annotators, we develop an automatic
343 subtask inference procedure, avoiding the need for manual segmentation.

344 Before the iterative process, we compute similarity scores for all states in the expert demonstrations
345 using the initial reward model M^0 . For each subtask U_i , we calculate a threshold T_{U_i} based on
346 the similarity scores between the expert states and the corresponding instructions, ensuring T_{U_i}
347 represents the minimum similarity required for successful subtask completion. In each iteration
348 $i \in \{1, \dots, n\}$, we proceed as follows:

- 349 • Step 1 (Suboptimal data collection): We train an RL agent using the reward model M^i and collect
350 suboptimal demonstrations $\mathcal{D}_{\text{replay}}^i$ from the agent’s replay buffer.
- 351 • Step 2 (Subtask inference for suboptimal data): For each timestep in the suboptimal trajectory,
352 we infer the subtask index \hat{i} using the same procedure as in inference and compute $\text{sim}(\mathbf{v}_t, e_{\hat{i}})$.
353 If the similarity falls below the threshold T_{U_i} at any timestep, we mark the subtask as failed and
354 assign the remaining timesteps to that subtask.
- 355 • Step 3 (Fine-tuning): We fine-tune the reward model M^{i-1} using the combined dataset $\mathcal{D}^i =$
356 $\mathcal{D}^i \cup \mathcal{D}_{\text{replay}}^i$ to obtain M^i .

357 We use the final reward model M^n for downstream RL training.

358 B Experiment Details

359 **Training and inference details** We used the open-source pre-trained CLIP [40] with ViT-B/16
360 architecture to encode images and subtask instructions for all experiments. We adopt a GPT [41]
361 architecture with 3 layers and 8 heads for the causal transformer. To canonicalize our reward func-
362 tions, we use the same \mathcal{D} for both coverage distribution \mathcal{D}_C and potential shaping distribution \mathcal{D}_S ,
363 and we estimate the expectation over state distributions using a sample-based average over 8 ad-
364 ditional samples from \mathcal{D} per sample. All models are trained with AdamW [42] optimizer with a
365 learning rate of 1×10^{-4} and a mini-batch size of 32. To ensure robustness against visual changes,
366 we apply data augmentations, including random shifting [43, 44] and color jittering. For optimiza-
367 tion, we train REDS with AdamW [42] optimizer with a learning rate of 1×10^{-4} , weight decay of
368 2×10^{-2} , and a cosine decay schedule for adjusting the training learning rate. We apply a warm-up
369 scheduling for the initial 500 gradient steps starting from a learning rate of 0. Note that the param-
370 eters for CLIP visual/text encoders have not been updated. For training downstream RL agents, we
371 normalize the reward by dividing it by the maximum value observed in the expert demonstrations.
372 We report the hyperparameters used in our experiments in Table 2. For both coverage distribution
373 \mathcal{D}_C and potential shaping distribution \mathcal{D}_S , we use the same dataset with subtask segmentations \mathcal{D}^i ,
374 unlike prior work dealing with arbitrarily random distributions because of the absence of subtask
375 segmentations. To canonicalize our reward functions, we estimate the expectation over state dis-
376 tributions using a sample-based average over 8 additional samples from \mathcal{D} per sample. To prevent

377 false positive cases in predicting subtasks in online interactions, we add margins to similarity scores
 378 inversely proportional to the subtasks in online interactions. Specifically, we infer the subtask \hat{i} as
 379 follows:

$$\hat{i} = \operatorname{argmax}_{i \in \{1, \dots, k\}} (\operatorname{sim}(\mathbf{v}_t, \mathbf{e}_i) + \eta * (k - i)), \quad (1)$$

380 where η is a hyperparameter for the margin between subtasks.

Table 2: Hyperparameters of REDS used in our experiments.

Hyperparameter	Value
Batch size	32 (Meta-world, RLBench), 8 (FurnitureBench)
Training steps	5000
Learning rate	0.0001
Optimizer	AdamW [42]
Optimizer momentum	$\beta_1 = 0.9, \beta_2 = 0.999$
Weight decay	0.02
Learning rate decay	Linear warmup and cosine decay
Warmup steps	500
Context length	4
Causal transformer size	3 layers, 8 heads, 512 units
EPIC canonical samples	8
ϵ for progressive reward signal	0.05
η for inferring subtasks	0.01 (Meta-world, RLBench), 0.05 (FurnitureBench)
number of training iterations n	2 (Meta-world, RLBench), 1 (FurnitureBench)

381 **Meta-world experiments** We use visual observations of $64 \times 64 \times 3$. To consistently use a single
 382 camera viewpoint over all tasks, we use the modified version of the `corner2` viewpoint as suggested
 383 by Seo et al. [45]. Expert demonstrations for each task are collected using scripted policies publicly
 384 released in the benchmark. We use an action repeat of 2 to accelerate training and set the maximal
 385 episode length as 250 for all Meta-world tasks. For downstream RL, we use the implementation
 386 of DreamerV3 from VIPER². We report the hyperparameters of DreamerV3 agents used in our
 387 experiments in Table 3. Unless otherwise specified, we use the same set of hyperparameters as
 388 VIPER.

389 **RLBench experiments** For training both reward models and downstream RL agents, we utilize
 390 $64 \times 64 \times 3$ RGB observations from the front camera and wrist camera. For downstream RL, we
 391 don't use any expert demonstrations and we use the same set of hyperparameters as VIPER.

392 **FurnitureBench experiments** We use the implementation of IQL from FurnitureBench³ for our
 393 experiments. We utilize $224 \times 224 \times 3$ RGB observations from the front camera and wrist cameras,
 394 along with proprioceptive states, to represent the current state. We encode each image with pre-
 395 trained R3M [46] for visual observations. Following [38], we first run offline RL for 1M gradient
 396 steps, then continue training while collecting environment interaction data, adding it to the replay
 397 buffer, and repeating this process for 150 episodes. Before online fine-tuning, we pre-fill the re-
 398 play buffer with 10 rollouts from the pre-trained IQL policy. We adopt techniques from RLPD [47]
 399 for efficient offline-to-online RL training. Specifically, we sample 50% of the data from the replay
 400 buffer and the remaining 50% from the offline data buffer containing 300 expert demonstrations. We
 401 also apply LayerNorm [48] in the critic/value network of the IQL agent to prevent catastrophic over-
 402 estimation. We list the hyperparameters used in our experiments in Table 4. For training REDS, we
 403 collect subtask segmentations for suboptimal demonstrations using the automatic subtask inference
 404 procedure described in Appendix A, and we manually modified some subtask segmentations with
 405 false negatives to guarantee stable performance. For baselines, we compare against VIPER and DrS.
 406 We emphasize that our method enables fully autonomous training in online RL sessions, in contrast

²https://github.com/Alescontrela/viper_rl

³https://github.com/clvrai/furniture-bench/tree/main/implicit_q_learning

Table 3: Hyperparameters of DreamerV3 [26] used in Meta-world experiments.

Hyperparameter	Value
General	
Replay Capacity (FIFO)	5×10^5
Start learning (prefill)	5000
MLP size	2×512
World Model	
RSSM size	512
Base CNN channels	32
Codes per latent	32

Table 4: Hyperparameters of IQL [38] used in FurnitureBench experiments.

Hyperparameter	Value
Learning rate	3×10^{-4}
Batch size	256
Policy # hidden units	(512, 256, 256)
Critic/value # hidden units	(512, 256, 256)
Image encoder	R3M [46]
Discount factor (γ)	0.996
Expectile (τ)	0.8
Inverse Temperature (β)	10.0

407 to DrS, which relies on a subtask indicator provided by humans. In our DrS experiments, subtasks
 408 were manually identified by a human. We measure the average number of completed subtasks over
 409 20 rollouts for evaluation.

410 **Computation** We use 24 Intel Xeon CPU @ 2.2GHz CPU cores and 4 NVIDIA RTX 3090 GPUs
 411 for training our reward model, which takes about 1.5 hours in Meta-world and 3 hours in Furni-
 412 tureBench due to high-resolution visual observations from multiple views. For training DreamerV3
 413 agents in Meta-world, we use 24 Intel Xeon CPU @ 2.2GHz CPU cores and a single NVIDIA RTX
 414 3090 GPU, which takes approximately 4 hours over 500K environment steps. For training IQL
 415 agents in FurnitureBench, we use 24 Intel Xeon CPU @ 2.2GHz CPU cores and a single NVIDIA
 416 RTX 3090 GPU, taking approximately 2 hours for 1M gradient steps in offline RL and 4.5 hours
 417 over 150 episodes of environment interactions in online RL.

418 C REDS Architecture Details

419 We encode visual observations with a pre-trained CLIP [40] ViT-B/16 visual encoder, utilizing all
 420 representations from the sequence of patches. We adopt 1D learnable parameters with the same size
 421 for positional embedding, and we add these parameters to 2D fixed sin-cos embeddings and add them
 422 to features. To encode temporal dependencies in visual observations, we use a GPT [41] architecture
 423 with 3 layers and 8 heads. In FurnitureBench, we use a sequence of images from both the front
 424 camera and wrist camera as input. Given $s_t^{\text{front}}/s_t^{\text{wrist}}$ from the front/wrist camera, we concatenate
 425 visual observations to $[o_{t-K-1}^{\text{front}}, o_{t-K-1}^{\text{wrist}}, \dots, o_t^{\text{front}}, o_t^{\text{wrist}}]$, add positional embeddings, 2D fixed sin-
 426 cos embeddings, and additional 1D learnable parameters for each viewpoint for effectively utilizing
 427 images from multiple cameras. We then pass the features to the transformer layer, the same as the
 428 model with a single image. The subtask embedder and final reward predictor are implemented as
 429 2-layer MLPs.

430 D Baseline Details

431 We consider the following baselines: (1) human-engineered reward functions provided in the bench-
 432 mark, (2) ORIL [49], an adversarial imitation learning (AIL) method trained only with offline
 433 demonstrations, (3) Rank2Reward (R2R) [20], an AIL method which trains a discriminator weighted
 434 with temporal ranking of video frames to reflect task progress, (4) VIPER [26], a reward model uti-
 435 lizing likelihood from a pre-trained video prediction model as a reward signal, and (5) DrS [29], an
 436 AIL method that assumes subtask information from the environment and trains a separate discrimi-
 437 nator for each subtask.

438 **ORIL [49]** For implementing ORIL with visual observations, we use the CNN architecture from
 439 Yarats et al. [43] to encode image observations. For training data, we use the same set of demon-

440 strations as for training REDS. Since our training data are divided into success and failure demon-
441 strations, we do not use positive-unlabeled learning [50] in our experiments. For robustness against
442 visual changes, we apply the same augmentation techniques used for training REDS.

443 **Rank2Reward (R2R) [20]** To ensure compatibility with backbone RL algorithms [37, 38] imple-
444 mented in JAX, we reimplement the reward model with JAX following the official implementation of
445 Rank2Reward⁴ and use the same hyperparameters. We first pre-train the ranking network using the
446 same expert demonstrations as REDS, and we then train a discriminator for the expert demonstration
447 and policy rollouts, weighted by the output from the pre-trained ranking network. For training effi-
448 ciency, we use the CNN architecture from Yarats et al. [43] for encoding visual observations instead
449 of R3M [46], finding no significant difference when we use the pre-trained visual representations
450 like R3M, but with much slower training in online RL. We observe that our R2R implementation
451 with DreamerV3 in JAX outperforms the original version implemented with DrQ-V2 [44] agents.

452 **DrS [29]** Similar to R2R, we reimplement DrS with JAX following the official implementation of
453 DrS⁵, and use the same set of hyperparameters for reward learning. As the original DrS implemen-
454 tation is based on a state-based environment, we switch the backbone RL algorithm from SAC to
455 DrQ-V2 [44] and apply the augmentation technique in the reward learning phase for processing vi-
456 sual observations efficiently. To report the RL performance, we use the learned dense reward model
457 to train new RL agents. In FurnitureBench experiment, we train the reward model with the same
458 expert/failure demonstrations as in Section 3.2, without online interaction, to avoid unsafe behaviors
459 and a significant increase in training time from online interactions.

460 **VIPER [26]** We use the official implementation of VIPER⁶ for our experiments. Given the
461 similarities among robotic manipulation tasks, we use the same set of hyperparameters as in RL-
462 Bench [10] experiments to train VQ-GAN and VideoGPT. We train 100K steps, choosing the check-
463 point with the minimum validation loss. In FurnitureBench experiment, we use images from the
464 front camera, resized to $64 \times 64 \times 3$, and set the exploration objective β as 0.

465 E Task Descriptions

466 In this section, we list the subtasks and corresponding text instructions for each task in Table 5. For
467 Meta-world tasks, we provide the code snippet used to determine the success of each subtask (Please
468 refer to the Meta-world [12] for more details). For the FurnitureBench One Leg task, we outline the
469 criteria used by human experts to assess the success of each subtask based on the metric defined in
470 FurnitureBench [28].

471 F Related Work

472 **Reward learning from videos** Learning from observations without expert actions has been a
473 promising research area because it does not require extensive instrumentation and allows for the
474 easy collection of vast amounts of video from online sources. Notably, several studies have pro-
475 posed methods for learning rewards directly from videos and using the signal to train RL agents.
476 Previous work has been focused on learning a reward function by aligning video representations in
477 temporal order [21, 22, 23] while others train a reward function for expressing the progress of the
478 agent towards the goal [35, 34, 20]. Most recent work [26] inspired by the success of video gener-
479 ative models [51, 52] utilizes the likelihood of pre-trained video prediction models as a reward. To
480 effectively utilize video for long-horizon tasks, we propose a new reward model conditioned both
481 on video segments and corresponding subtasks trained with subtask segmentations.

⁴<https://github.com/dxyang/rank2reward>

⁵<https://github.com/tongzhoumu/DrS>

⁶https://github.com/Alescontrela/viper_rl

Table 5: A list of subtasks and language description for each subtask used in our experiments.

Task	Subtask	Success condition	Language description
Meta-world Faucet Close	1	$\text{object_grasped} \leq 0.9$	a robot arm reaching the faucet handle.
	2	$\text{target_to_obj} \leq 0.07$	a robot arm rotating the faucet handle to the right.
Meta-world Drawer Open	1	$\text{gripper_error} \leq 0.03$	a robot arm grabbing the drawer handle.
	2	$\text{handle_error} \leq 0.03$	a robot arm opening a drawer to the green target point.
	3	$\text{handle_error} \leq 0.03$	a robot arm holding the drawer handle near the green target point after opening.
Meta-world Lever Pull	1	$\text{ready_to_lift} > 0.9$	a robot arm touching the lever.
	2	$\text{lever_error} \leq \text{np.pi}/24$	a robot arm pulling up the lever to the red target point.
Meta-world Door Open	1	$\text{reward_ready} \geq 1.0$	a robot arm grabbing the door handle.
	2	$\text{abs}(\text{obs}[4] - \text{self_target_pos}[0]) \leq 0.08$	a robot arm opening a door to the green target point.
	3	$\text{abs}(\text{obs}[4] - \text{self_target_pos}[0]) \leq 0.08$	a robot arm holding the door handle near the green target point after opening.
Meta-world Coffee Pull	1	$\text{tcp_to_obj} < 0.04 \wedge \text{tcp_open} > 0$	a robot arm grabbing the coffee cup.
	2	$\text{obj_to_target} \leq 0.07$	a robot arm moving the coffee cup to the green target point.
	3	$\text{obj_to_target} \leq 0.07$	a robot arm holding the cup near the green target point.
Meta-world Peg Insert Side	1	$\text{tcp_to_obj} < 0.02 \wedge \text{tcp_open} > 0$	a robot arm grabbing the green peg.
	2	$\text{obj}[2] - 0.1 > \text{self.obj_init_pos}[2]$	a robot arm lifting the green peg from the floor.
	3	$\text{obj_to_target} \leq 0.07$	a robot arm inserting the green peg to the hole of the red box.
	4	$\text{obj_to_target} \leq 0.07$	a robot arm holding the green peg after inserting.
Meta-world Push	1	$\text{tcp_to_obj} \leq 0.03$	a robot arm grabbing the red cube.
	2	$\text{target_to_obj} \leq 0.05$	a robot arm pushing the grabbed red cube to the green target point.
	3	$\text{target_to_obj} \leq 0.05$	a robot arm holding the grabbed red cube near the green target point.
Meta-world Sweep Into	1	$\text{self.touching_main_object} > 0 \wedge \text{tcp_opened} > 0$	a robot arm grabbing the red cube.
	2	$\text{target_to_obj} \leq 0.05$	a robot arm sweeping the grabbed red cube to the blue target point.
	3	$\text{target_to_obj} \leq 0.05$	a robot arm holding the grabbed red cube near the blue target point.
Meta-world Door Close	1	$\text{in_place} == 1.0$	a robot arm grabbing the door handle.
	2	$\text{obj_to_target} \leq 0.08$	a robot arm closing a door to the green target point.
	3	$\text{obj_to_target} \leq 0.08$	a robot arm holding the door handle near the green target point after closing.
Meta-world Window Close	1	$\text{tcp_to_obj} \leq 0.05$	a robot arm grabbing the window handle.
	2	$\text{target_to_obj} \leq 0.05$	a robot arm closing a window from left to right.
	3	$\text{target_to_obj} \leq 0.05$	a robot arm holding the window handle after closing.
FurnitureBench One Leg	1	robot gripper tips make contact with one surface of the tabletop.	a robot arm picking up the white tabletop.
	2	nearest corner of the tabletop is placed close to the right edge of the obstacle.	a robot arm pushing the white tabletop to the front right corner.
	3	robot gripper securely grasps a leg of the table and lifts it.	a robot arm picking up the white leg.
	4	leg is inserted into one of the screw holes of the tabletop, and the robot releases the gripper.	a robot arm inserting the white leg into screw hole.
	5	leg is fully assembled to the tabletop.	a robot arm screwing the white leg until tightly lifted.
	6	leg is fully assembled to the tabletop.	a robot arm holding the white leg in place.
RLBench Take Umbrella Out of Umbrella Stand	1	$\text{GraspedCondition}(\text{self.robot.gripper, self.umbrella}).\text{condition}_{\text{net}}[0]$	a robot arm grasping the umbrella.
	2	$\text{DetectedCondition}(\text{self.umbrella, self.success.sensor, negated} = \text{True}).\text{condition}_{\text{net}}[0]$	a robot arm taking the grasped umbrella out of the umbrella stand.
	3	$\text{DetectedCondition}(\text{self.umbrella, self.success.sensor, negated} = \text{True}).\text{condition}_{\text{net}}[0]$	a robot arm holding the umbrella on the umbrella stand.

482 **Inverse reinforcement learning** Designing an informative reward function remains a long-
483 standing challenge for training RL agents. To achieve this, Inverse Reinforcement Learning (IRL)
484 [53, 54, 55] aims to estimate the underlying reward function from expert demonstrations. Adversarial
485 imitation learning (AIL) approaches [56, 57, 49, 58, 29] address this by training a discriminator
486 network to discriminate transitions from expert data or policy rollouts and using the output from the
487 discriminator as a reward for training agents with RL. The most similar work to ours is DrS [29],
488 which also utilizes subtask information of the multi-stage task. While DrS assumes that the infor-
489 mation on ongoing subtasks can be obtained from the environment during online interaction, our
490 method has no such assumption, so it can be applied in more general cases when the segmenting of
491 the subtask is hard in automatic ways (e.g., [28]).

492 **Quantifying differences between reward functions** Previous work has explored methods for
493 measuring the difference between reward functions without relying on policy optimization proce-
494 dures [32, 59, 60]. In particular, Gleave et al. [32] introduced the EPIC distance, a pseudometric
495 invariant to equivalent classes of reward functions. Subsequent work [61, 62, 63] has employed
496 EPIC to assess the quality of reward functions. In this paper, we take a different approach by using
497 EPIC distance as an optimization objective. While Adeniji et al. [62] also utilizes EPIC distance for
498 optimizing intrinsic reward functions in skill discovery, our method applies EPIC distance to train
499 dense reward functions for long-horizon tasks, serving as a direct reward signal for RL training.

500 **Segmenting demonstrations for long-horizon manipulation tasks** Several approaches have
501 been proposed to decompose long-horizon demonstrations into multiple subgoals to prevent error
502 accumulation and provide intermediate signals for agent training. These include extracting key
503 points from proprioceptive states [64, 65, 66, 67], employing greedy heuristics on off-the-shelf vi-
504 sual representations pre-trained with robotic data [68], and learning additional modules on top of
505 pre-trained visual-language models to align with keyframes [69]. Our work builds on these efforts
506 by leveraging subtask segmentations but focuses on developing a reward learning framework that
507 explicitly incorporates subtask decomposition to generate suitable reward signals for intermediate
508 tasks. Additionally, we further demonstrate that our model generalizes effectively to unseen tasks
509 and robot embodiments.

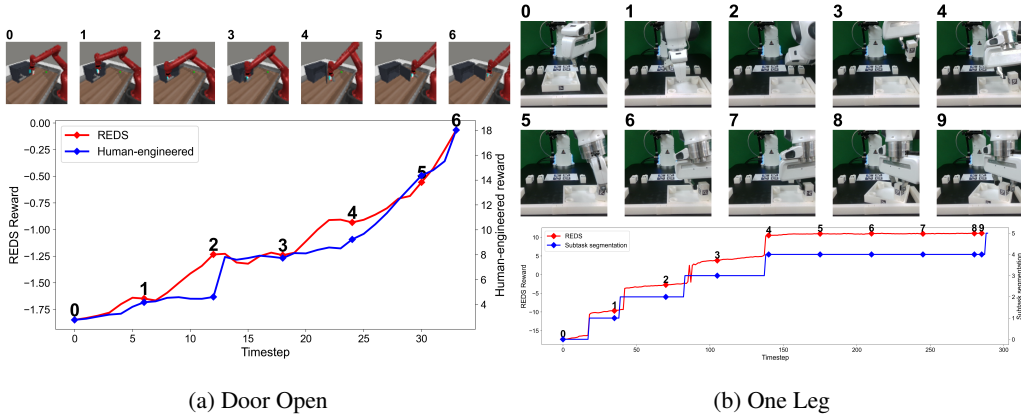


Figure 4: Qualitative results of REDS in Door Open in Meta-world [12] and One Leg from FurnitureBench [28]. We observe that REDS produces suitable reward signals aligned with ground-truth reward functions by predicting ongoing subtasks effectively and providing progressive reward signals.

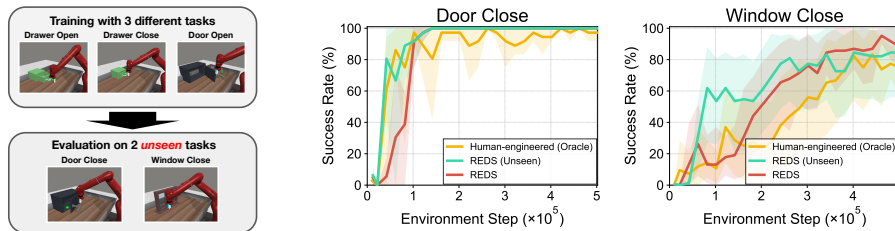


Figure 5: We train REDS with 3 different tasks from Meta-world [12] and use this model to train RL agents in 2 unseen tasks (left). We present learning curves on Door Close (center) and Window Close (right), as measured by success rate (%). The solid line and shaded regions represent the mean and stratified bootstrap interval across 4 runs.

510 G Additional Experimental Results

511 G.1 Alignment with Ground-truth Rewards

512 **EPIC measurement** To quantitatively
 513 validate the alignment of our method with
 514 ground-truth reward functions, we mea-
 515 sure the EPIC distance with a set of unseen
 516 demonstrations during training. Specifi-
 517 cally, we use rollouts from the reference
 518 policy trained with expert demonstrations
 519 for state distribution. In Table 6, we
 520 observe that REDS exhibits significantly
 521 lower EPIC distance than baselines across
 522 all tasks. Particularly, the difference between REDS and baselines is more pronounced in complex
 523 tasks like One Leg. This result consistently supports the empirical findings from previous sections.

Table 6: EPIC [32] distance (lower is better) between learned reward functions and hand-engineered reward functions (Meta-world) / subtask segmentations (FurnitureBench) in unseen data.

Task	VIPER	R2R	ORIL	REDS (Ours)
Meta-world Door Open	0.5934	0.5649	0.7071	0.4913
Meta-world Push	0.6144	0.6838	0.7073	0.5381
Meta-world Peg Insert Side	0.5974	0.5806	0.6989	0.4674
Meta-world Sweep Into	0.6248	0.6413	0.7001	0.4673
FurnitureBench One Leg	0.7035	0.6001	0.7014	0.0713

524 **Qualitative analysis** We provide the graph of computed rewards from REDS in Figure 4. We
 525 observe that REDS can induce suitable reward signals aligned with ground-truth reward functions.
 526 For example, REDS provides subtask-aware signals in transition states (e.g., between 2 and 3, and
 527 between 4 and 5) and generates progressive reward signals throughout each subtask. Please refer to
 528 Appendix H for the extensive comparison between REDS and baselines.

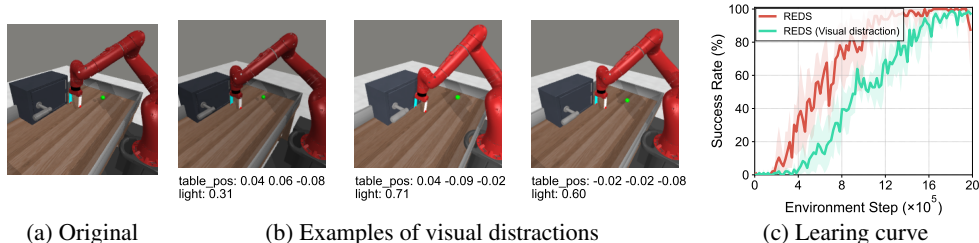


Figure 6: We provide visual observations from (a) the original environment and (b) unseen environments with visual distractions used in our experiments in Section G.2 .

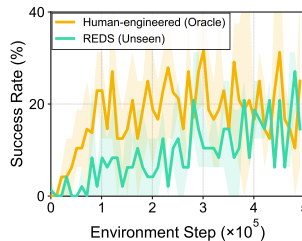
529 G.2 Generalization Capabilities

530 **Transfer to unseen tasks** REDS can be applied as a reward function in unseen tasks. To validate
 531 this, we conduct additional experiments by training REDS with segmentation data from 3 tasks
 532 (Door Open, Drawer Open/Close) and using the reward model to train RL agents in two unseen
 533 tasks. In Door Close, we aim to validate that REDS can provide informative signals for a new task
 534 involving a previously seen object and behaviors. In Window Close, we aim to determine whether
 535 REDS can provide suitable reward signals for familiar behaviors (closing) with an unseen object
 536 (window). In evaluation, we change the text instruction following the target object (as shown in
 537 Table E), and we do not fine-tune the reward model. Figure 5 shows that REDS provides effective
 538 reward signals on unseen tasks and achieves comparable or even better RL performance than REDS
 539 trained on the target task. This result demonstrates that REDS can be applied to RL training in
 540 unseen tasks that share properties with training tasks.

541 **Robustness to visual distractions** To prove the robust performance of REDS against visual dis-
 542 tractions, we train RL agents with our reward model in new Meta-world environments incorporating
 543 visual distractions, such as varying light and table positions following [70] (see Figure 6b). Note
 544 that the reward model was trained using demonstrations only from the original environment. As
 545 Figure 6c shows, REDS can generate robust reward signals despite visual distractions and train RL
 546 agents to solve the task effectively.

547 **Transfer to unseen embodiments** Since our framework lever-
 548 ages only action-free video data, we hypothesize that transferring to other robot embod-
 549 iments with similar DoFs is feasible. To support this claim, we train REDS with demonstrations of the Franka
 550 Panda Arm and then compute the reward of an unseen demon-
 551 stration of the Sawyer Arm in Take Umbrella Out of Stand from
 552 RL Bench [10]. Figure 8 shows that REDS generates informative
 553 reward signals even with the unseen embodiment. For instance,
 554 REDS can capture the behavior of taking the umbrella out of the
 555 stand, as indicated by the increased reward signals between 6 and
 556 7. Additionally, Figure 7 shows that REDS trained only with the
 557 Panda Arm can be used to train downstream RL agents in the environment with the Sawyer Arm.
 558

Figure 7: Learning curve for DreamerV3 agents in environments of the Sawyer Arm.



559 G.3 Ablation Studies

560 **Effect of training objectives** We investigate the effect of each training objective in Figure 9a.
 561 Specifically, we compare REDS with 1) a baseline trained with regression to subtask segmentation
 562 instead of EPIC loss $\mathcal{L}_{\text{EPIC}}$, 2) a baseline that utilizes only video representations without subtask
 563 embeddings, and 3) a baseline trained without the regularization loss \mathcal{L}_{reg} . We observe that RL
 564 performance significantly degrades without each component, implying that our losses synergistically
 565 improve reward quality.

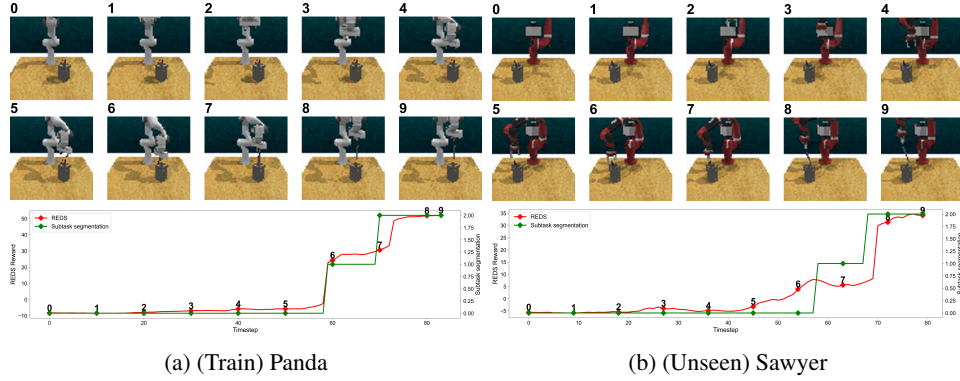


Figure 8: Qualitative results of REDS with different robot embodiments. REDS was trained using demonstrations from the Panda Arm and evaluated on an unseen demonstration from the Sawyer Arm in Take Umbrella Out of Stand from RLbench [10]. We visualize several frames above the graph and mark them with a diamond symbol.

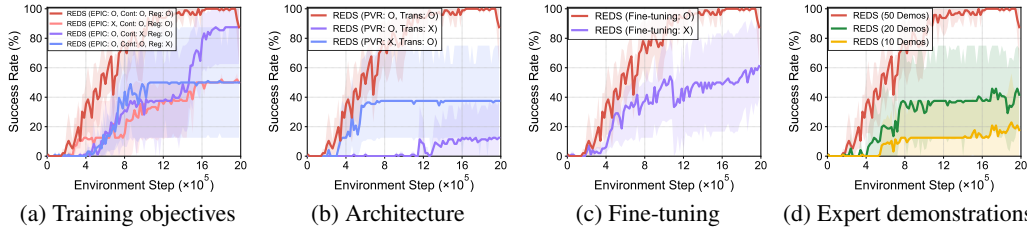


Figure 9: Learning curves for two Meta-world [12] robotic manipulation tasks, measured by success rate (%), to examine the effects of (a) training objectives, (b) architecture, (c) fine-tuning, and (d) the number of expert demonstrations. The solid line and shaded regions show the mean and stratified bootstrap interval across 8 runs.

566 **Effect of architecture** To verify the design choice, we compare REDS with 1) a baseline using
 567 a CNN for encoding images instead of pre-trained visual representations (PVR) and 2) a baseline
 568 simply concatenating pre-trained visual representations without a causal transformer. Figure 9b
 569 shows that both baselines show worse performance compared to ours. Notably, detaching a causal
 570 transformer significantly degrades RL performance, implying that temporal information is essential
 571 for providing suitable reward signals in robotic manipulation.

572 **Effect of fine-tuning** In Figure 9c, we compare REDS trained only with the expert demonstrations
 573 in the initial phase to REDS fine-tuned with additional suboptimal demonstrations as described
 574 in Appendix A. REDS shows improved RL performance when trained with additional suboptimal
 575 demonstrations, indicating that the coverage of state distribution impacts the reward quality. Further
 576 investigation on how to efficiently collect suboptimal demonstrations to enhance the performance of
 577 learned reward function is a promising future direction.

578 **Effect of the number of expert demonstrations** We investigate the effect of the number of expert
 579 demonstrations by measuring the RL performance of DreamerV3 agents with REDS trained with
 580 different numbers of expert demonstrations in 2 tasks (Door Open, Drawer Open) from Meta-world.
 581 Figure 9d shows that the agents' RL performance positively correlates with the number of expert
 582 demonstrations trained for reward learning.

583 **H Extended Qualitative Analysis**

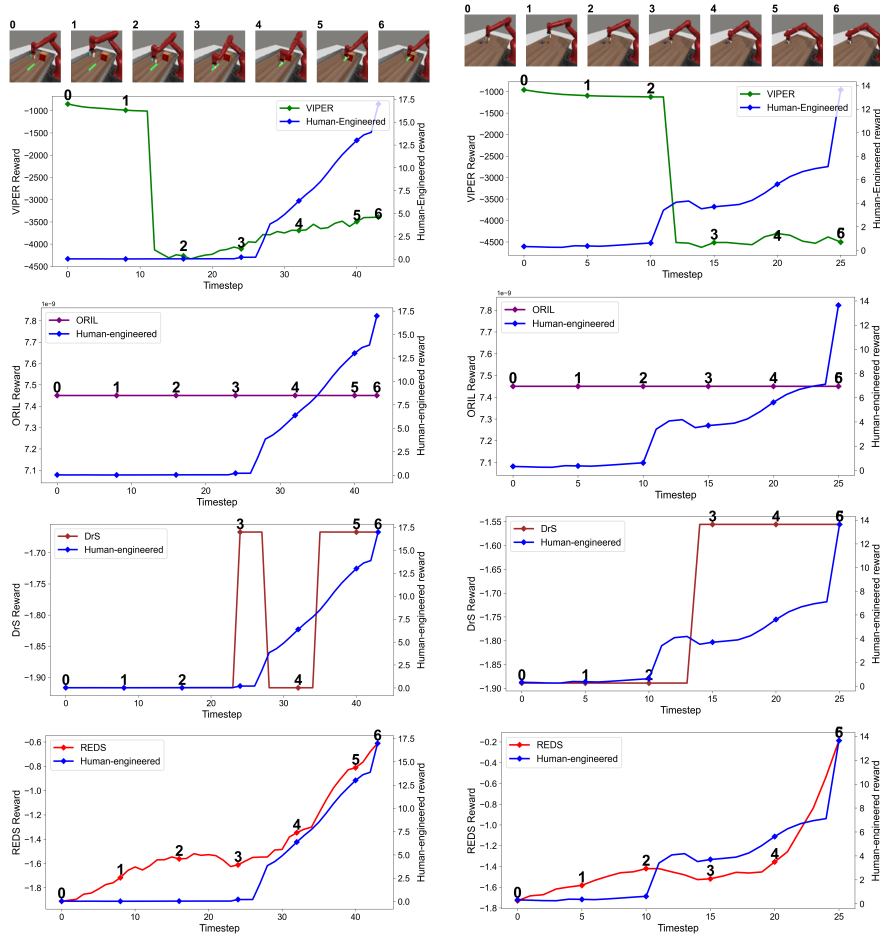


Figure 10: Qualitative results of VIPER [26], ORIL [49], DrS [29], and REDS (Ours) in Peg Insert Side (left), and Sweep Into (right) from Meta-world [12]. We visualize several frames above the graph and mark them with a diamond symbol.

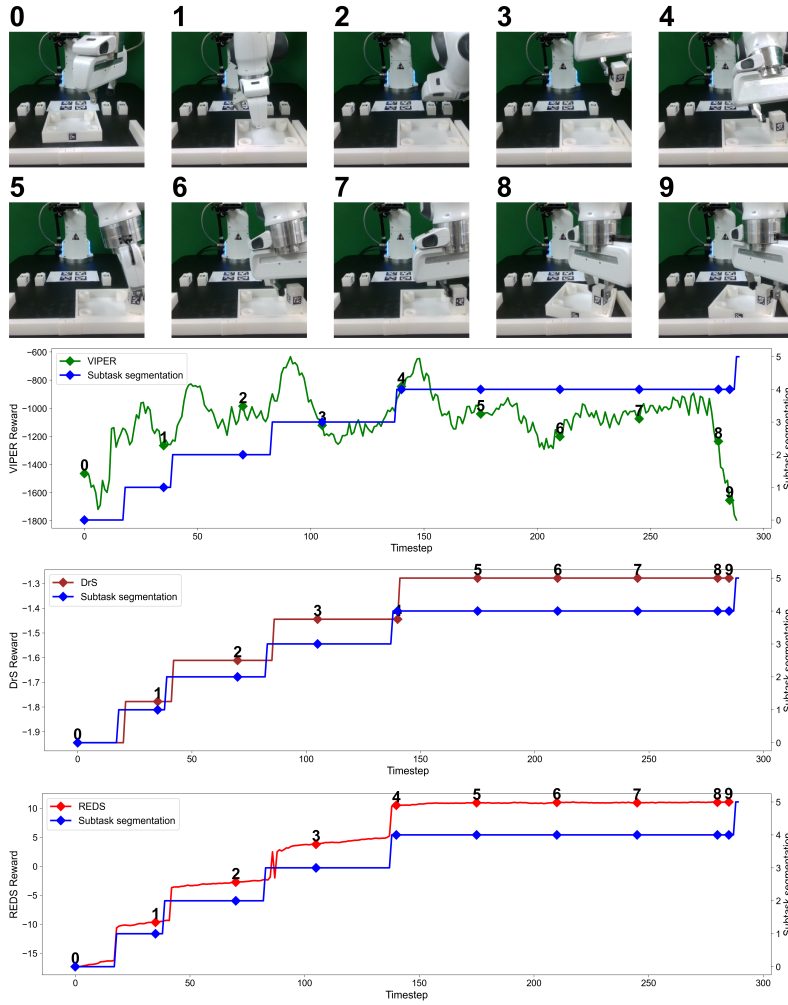


Figure 11: Qualitative results of VIPER [26], DrS [29], and REDS (Ours) in One Leg from FurnitureBench [28]. We visualize several frames above the graph and mark them with a diamond symbol. VIPER, which does not utilize subtask information, failed to produce suitable rewards for transitioning between subtasks. While DrS uses ground-truth subtask information from the environment, it produces sparse reward signals within each subtask. In contrast, REDS provides subtask-aware signals in transition states (e.g., between 2 and 3, and 4 and 5) and generates dense reward signals throughout each subtask.

584 **I Limitation and Future Directions**

585 One limitation of our work is the reliance on pre-trained representations trained with natural im-
586 age/text data for encoding videos and subtasks. Although REDS proves its effectiveness in various
587 robotic manipulation tasks, we observe that REDS struggles to distinguish subtle changes. We be-
588 lieve that the quality of rewards can be further improved by utilizing 1) larger models pre-trained
589 with large-scale data with diverse robotic tasks [71, 72] and 2) representations trained with objec-
590 tives considering affordances [73] or object-centric methods [74]. Moreover, the number of expert
591 demonstrations and the number of iteration for fine-tuning REDS are determined by empirical trials.
592 Investigating how to efficiently collect failure demonstrations to mitigate reward misspecification is
593 an interesting future direction.

594 **J Ethics Statement**

595 Video demonstrations and subtask segmentations used in the experiments were sourced from pub-
596 licly available benchmarks (Meta-world, RLBench, FurnitureBench), ensuring no personal or sen-
597 sitive information is involved. Potential risks could arise when training and deploying RL agents
598 directly in real-world scenarios, particularly in human-robot interactions. Ensuring the safety and
599 reliability of these agents before deployment is essential to prevent harm.

600 **K Reproducibility Statement**

601 For the reproducibility of REDS, we have provided a detailed explanation of implementation details
602 and experimental setups in Appendix A, Section 3, and Appendix B. In addition, to further facilitate
603 the reproduction, we attach the source code used in our experiments in the supplementary materials.