

Supplementary Material

Guide Your Agent with Adaptive Multimodal Rewards

A Experiment Details

In this section, we describe the details for implementing Multimodal Reward Decision Transformer, and we provide our source code in the supplementary material.

Progen details We utilize a publicly available implementation³ to replicate the environments introduced by Di Langosco et al. [9]. We modify the simulator of the environments to render higher-resolution images to leverage pre-trained multimodal representations for both our method and baselines. In this particular setup, the observations obtained from the environment at each timestep t comprise an RGB image with dimensions of $256 \times 256 \times 3$ and a natural language instruction that delineates the desired goal. Throughout our experiments, we adhere to the *hard environment difficulty* as described in [7]. Maximum episode length for all tasks is 500. To gather expert demonstrations used for training data, we train PPG [8] agents on 500 training levels for 200M timesteps per task using hyperparameters provided in Cobbe et al. [8]. For evaluation purposes, we assess the test performance on 1000 different levels, encompassing previously unseen themes and goals that differ from those employed in training.

Architecture details Both *InstructRL* (Liu et al., 2022) and MRDT employ ViT-B/16 as the transformer-policy and pre-trained multimodal transformer encoder (M3AE; [12]) in all experiments, unless stated otherwise. Inspired by Gao et al. [11], we attach an additional 2-layer MLP to the end of a pre-trained multimodal transformer encoder and perform residual-style feature blending with the pre-trained features. In the training phase, we apply gradients only to the weight of these linear layers. Through empirical evaluation, we observe that this architecture yields superior performance in both our method and the baseline.

Training details We use $256 \times 256 \times 3$ RGB observations for training the return-conditioned policy. To stabilize training, we normalize multimodal returns following the method proposed by Chen et al. (2021), dividing them by 1000 in all experiments. We use the AdamW optimizer (Loshchilov et al., 2018) with a learning rate of 5×10^{-4} and weight decay 5×10^{-5} . A cosine decay schedule is utilized to adjust the training learning rate. In CoinRun experiments, data augmentation techniques such as color jitter and random rotation are applied to the RGB images o_t while maintaining alignment in the context. However, no augmentation is applied to RGB images in Maze I/II experiments. For scaling the return prediction loss in training the return-conditioned policy, we set $\lambda = 0.01$ in CoinRun experiments and $\lambda = 0.001$ in Maze I/II experiments. During the fine-tuning of the pre-trained multimodal encoder, a 2-layer MLP is attached to the end of both CLIP image and text encoders. Additionally, an extra 2-layer MLP is added as an action prediction layer for the IDM objective. The model is trained for 20 epochs, and the one with the lowest validation loss is used for generating multimodal rewards. To scale the IDM loss in fine-tuning CLIP, we employ $\beta = 1.5$ in CoinRun experiments and $\beta = 2.0$ in Maze I/II experiments.

Computation We use 24 CPU cores (Intel Xeon CPU @ 2.2GHz) and 2 GPUs (NVIDIA A100 40GB GPU) for training return-conditioned policy. The training of MRDT for 50 epochs takes approximately 4 hours for CoinRun experiments with the largest dataset size. For fine-tuning CLIP, we use 24 CPU cores (Intel Xeon CPU @ 2.2GHz) and 1 GPU (NVIDIA A100 40GB GPU), and it takes approximately 1.5 hours for Coinrun experiments.

³<https://github.com/JacobPfau/procgenAISC>

Hyperparameters We report the hyperparameters used in our experiments in Table 3.

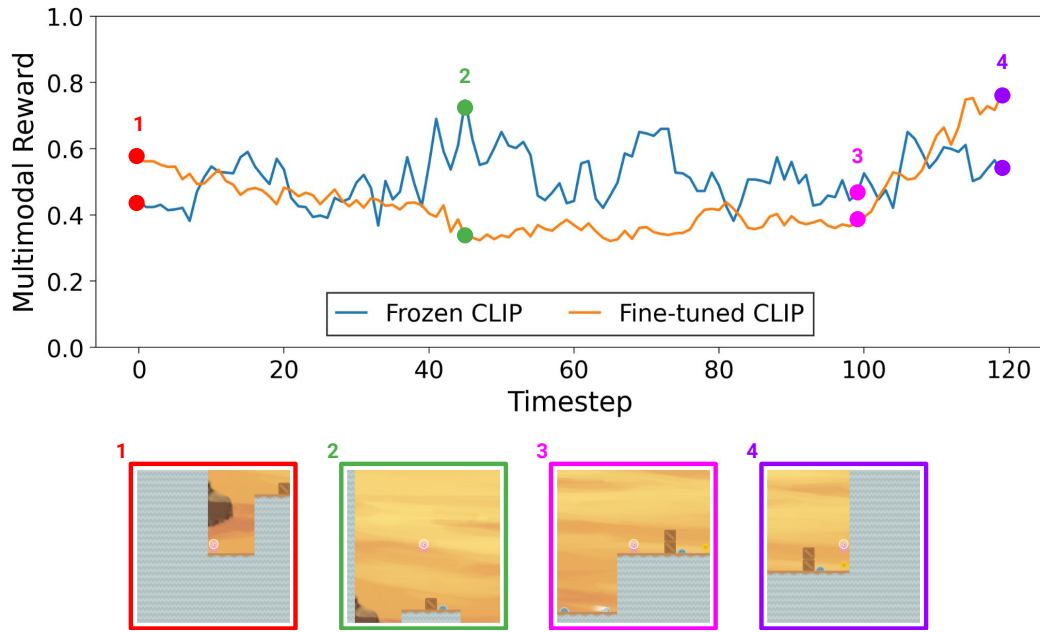
Table 3: Hyperparameters of Multimodal Reward Decision Transformer (MRDT). Unless specified, we use the same hyperparameters used in *InstructRL* [16].

Hyperparameter	Value
Policy batch size	64
Policy epochs	50
Policy context length	4
Policy learning rate	0.0005
Policy optimizer	AdamW [18]
Policy optimizer momentum	$\beta_1 = 0.9, \beta_2 = 0.999$
Policy weight decay	0.00005
Policy learning rate decay	Linear warmup and cosine decay (see code for details)
Policy context length	4
Policy transformer size	2 layers, 4 heads, 768 units
Fine-tuned CLIP batch size	64
Fine-tuned CLIP epochs	20
Fine-tuned CLIP learning rate	0.0001
Fine-tuned CLIP weight decay	0.001
Fine-tuned CLIP adapter layer size	2 layers, 1024 units
Fine-tuned CLIP optimizer	AdamW [18]
Fine-tuned CLIP optimizer momentum	$\beta_1 = 0.9, \beta_2 = 0.999$

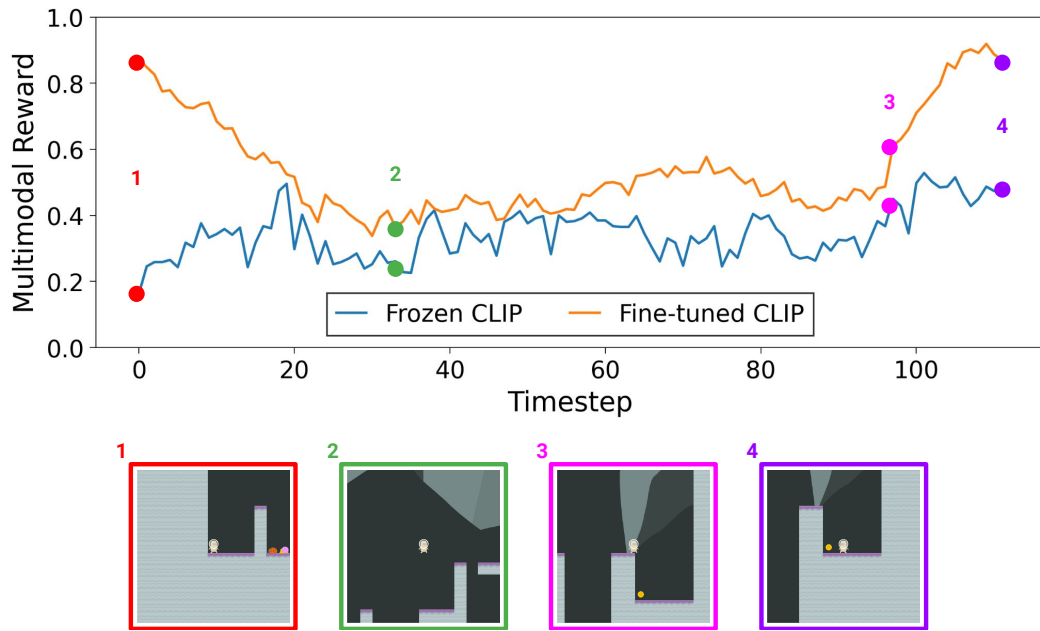
387

388 B Qualitative Results of Multimodal Rewards

389 In Figure 9, 10, 11, we present the curves of multimodal rewards for frozen/fine-tuned CLIP on the
 390 trajectories from training/held-out evaluation environments. We find that the multimodal reward
 391 exhibits an overall increasing trend as the agent approaches the goal in both frozen and fine-tuned
 392 CLIP, irrespective of the training and held-out evaluation environments. Furthermore, we observe
 393 that fine-tuned CLIP not only induces a reward that is temporally smoother in the intermediate stages
 394 compared to frozen CLIP (see Figure 9) but also demonstrates a steeper upward reward curve (see
 395 Figure 10, 11). These results support the claim that the quality of multimodal rewards from the
 396 fine-tuned CLIP outperforms those from the frozen CLIP (Section 4.2). Video examples of the
 397 trajectories are provided in the supplementary material.

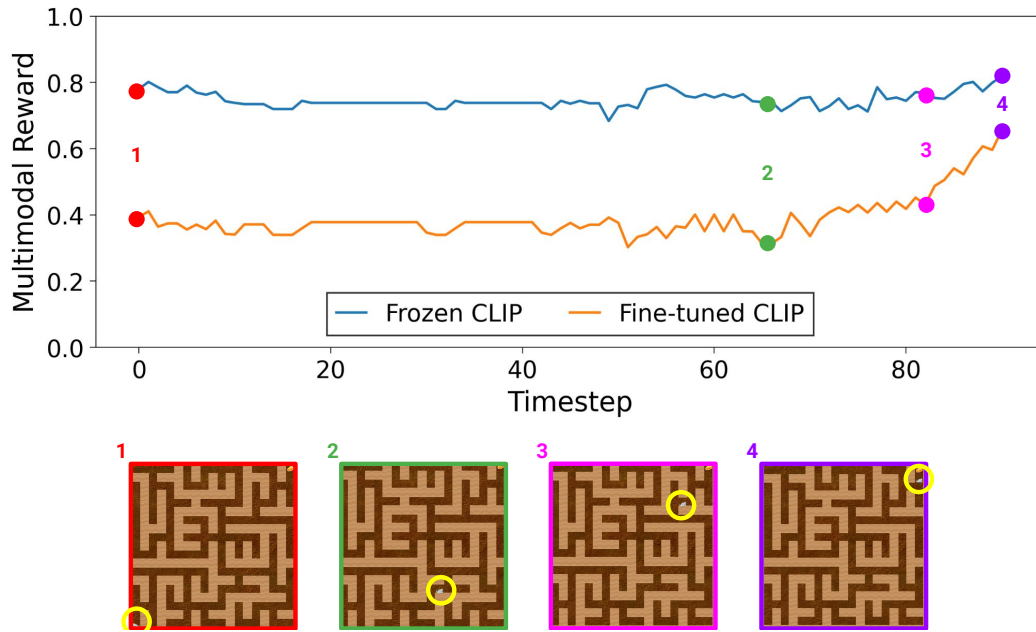


(a) Multimodal reward curve on the training environment.

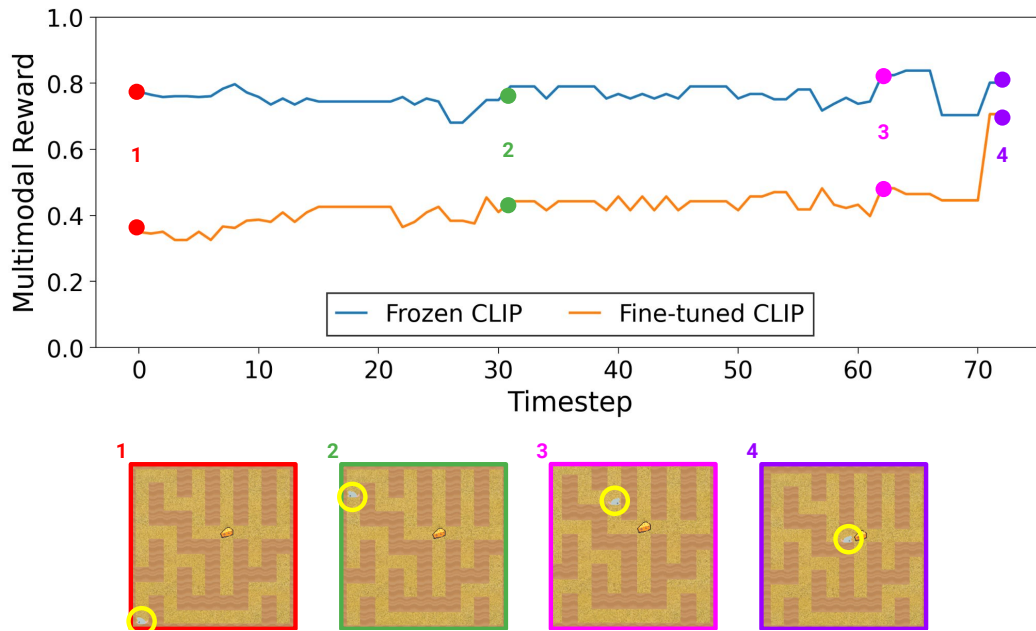


(b) Multimodal reward curve on the held-out evaluation environment.

Figure 9: Qualitative results of multimodal rewards in CoinRun environments.

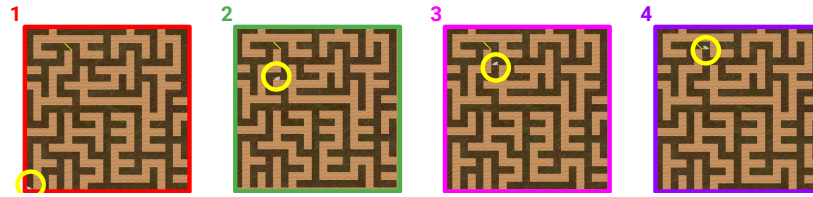
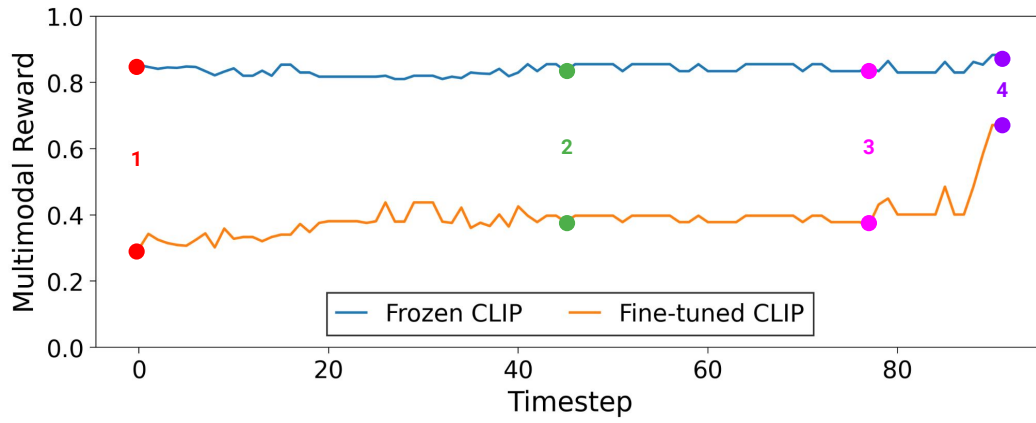


(a) Multimodal reward curve on the training environment.

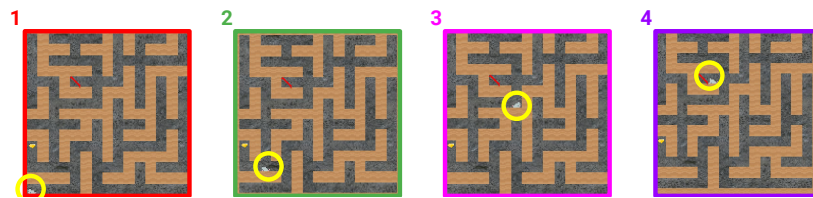
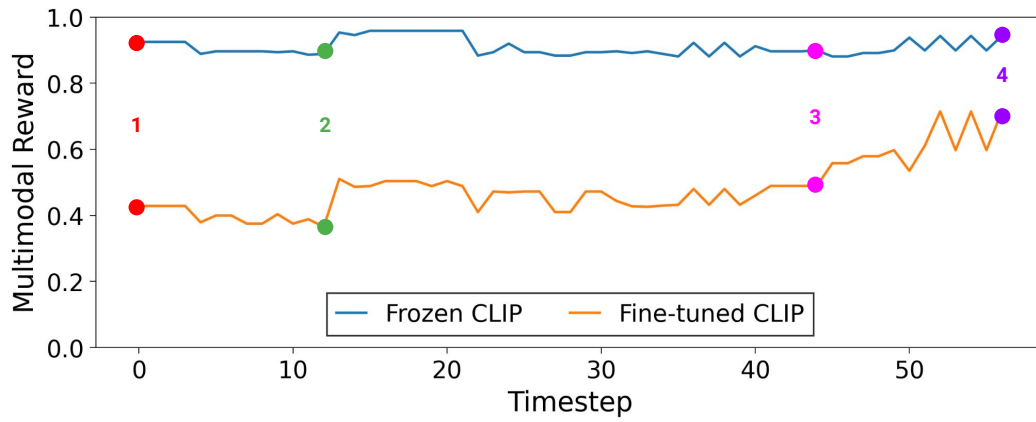


(b) Multimodal reward curve on the held-out evaluation environment.

Figure 10: Qualitative results of multimodal rewards in Maze I environments.



(a) Multimodal reward curve on the training environment.



(b) Multimodal reward curve on the held-out evaluation environment.

Figure 11: Qualitative results of multimodal rewards in Maze II environments.

Table 4: Expert-normalized scores on training/evaluation CoinRun environments investigating the effect of hyperparameter λ adjusting the scale of return prediction loss in training return-conditioned policy. The result shows the mean and standard variation averaged over three runs.

λ	\mathcal{L}_{FT}	Train (%)	Eval (%)
0.001	✗	89.93% \pm 3.94%	62.65% \pm 10.12%
	✓	85.42% \pm 1.80%	71.69% \pm 5.71%
0.01	✗	89.58% \pm 2.08%	63.32% \pm 2.01%
	✓	90.28% \pm 1.59%	72.36% \pm 3.48%
0.1	✗	87.15% \pm 2.62%	62.65% \pm 10.31%
	✓	85.76% \pm 3.18%	73.37% \pm 3.48%
1.0	✗	87.15% \pm 4.70%	61.64% \pm 6.38%
	✓	81.25% \pm 1.04%	73.37% \pm 2.66%

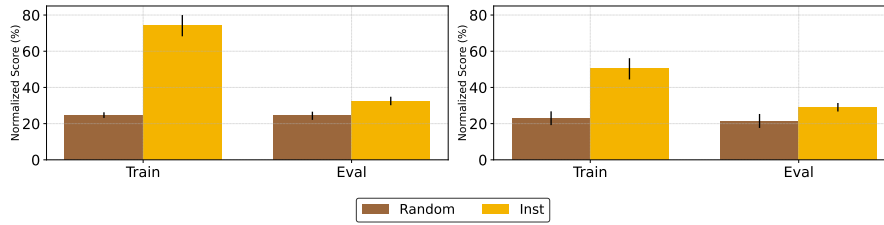


Figure 12: Expert-normalized scores on training/evaluation environments of MRDT trained using multimodal rewards generated with (i) instructive text (*i.e.*, Inst) and (ii) random text (*i.e.*, Random) in Maze I environments (left) and Maze II environments (right). The result shows the mean and standard deviation averaged over three runs.

Effect of scaling return prediction loss We investigate how the coefficient λ , which determines the weight of the return prediction loss in training return-conditioned policy, affects the performance of MRDT. To this end, we test various values of λ in CoinRun environments. Table 4 shows the performance of MRDT in training/evaluation environments with different λ . We find that performance is not significantly different according to the value of λ in the held-out evaluation environments. These results indicate that MRDT is robust to the choice of hyperparameter λ .

Extra ablation study on text instructions In Figure 12, we further investigate whether MRDT leverages adaptive signals from multimodal rewards in decision-making. We evaluate the quality of rewards generated with instructive text (*i.e.*, Inst) and random text (*i.e.*, Random) in Maze I/II environments. Specifically, we use a natural language instruction for each environment, as described in Section 4 for Inst, and "NeurIPS 2023 will be held again at the New Orleans Ernest N. Morial Convention Center" for Random. We find that using random text instructions results in a decline in performance in both training and evaluation environments. These findings align with the trend observed in Figure 8.

413 D Limitation and Future Work

414 One limitation of our work is that we currently rely on a single image-text pair to compute the
415 multimodal reward at every timestep t . Although our approach has shown effectiveness both quanti-
416 tatively and qualitatively, there are tasks where rewards depend on the history of past observations
417 (*i.e.*, non-Markovian) [2, 3, 15]. To address this limitation, it would be valuable to explore the
418 extension of our method to incorporate video-text pairs for calculating multimodal rewards. This
419 extension could involve generating multimodal rewards using pre-trained video-text multimodal
420 representations [17, 19, 24, 21], which presents an intriguing avenue for better generalization across
421 various goals in behavior learning. Another aspect to consider is that the tasks we have examined
422 so far are relatively simple, as they involve only a single condition for success. To tackle more
423 complex problems, we are interested in investigating approaches that leverage large language mod-
424 els [13, 14, 1, 10] in conjunction with our method. Finally, an interesting direction to explore would
425 be the utilization of multimodal rewards in combination with extrinsic rewards [22, 20, 5].

426 E Potential Negative Societal Impacts

427 We do not anticipate significant negative societal impacts in that our method is now limited to playing
428 simple simulation games. However, if our method is applied in real-world scenarios, privacy concerns
429 may arise considering that behavior cloning agents used in such applications, like autonomous driv-
430 ing [23] or real-time control [4, 10], require large amounts of data, which often contain controversial
431 information. Additionally, a behavior cloning policy presents a challenge as it imitates specified
432 demonstrations, potentially including undesirable actions. If some bad actions are included in expert
433 demonstrations (*e.g.*, behaviors that may be violent or harmful to the pedestrians are contained in
434 the training data for mobile manipulation tasks), the policy could have significant negative impacts
435 on users. To address this concern, future directions should focus on developing agents with safe
436 adaptation in addition to performance enhancement efforts.

References

- [1] Ahn, Michael, Brohan, Anthony, Brown, Noah, Chebotar, Yevgen, Cortes, Omar, David, Byron, Finn, Chelsea, Gopalakrishnan, Keerthana, Hausman, Karol, Herzog, Alex, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [2] Bacchus, Fahiem, Boutilier, Craig, and Grove, Adam. Rewarding behaviors. In *Proceedings of the National Conference on Artificial Intelligence*, pp. 1160–1167, 1996.
- [3] Bacchus, Fahiem, Boutilier, Craig, and Grove, Adam. Structured solution methods for non-markovian decision processes. In *AAAI/IAAI*, pp. 112–117. Citeseer, 1997.
- [4] Brohan, Anthony, Brown, Noah, Carbajal, Justice, Chebotar, Yevgen, Dabis, Joseph, Finn, Chelsea, Gopalakrishnan, Keerthana, Hausman, Karol, Herzog, Alex, Hsu, Jasmine, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- [5] Burda, Yuri, Edwards, Harrison, Storkey, Amos, and Klimov, Oleg. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019.
- [6] Chen, Lili, Lu, Kevin, Rajeswaran, Aravind, Lee, Kimin, Grover, Aditya, Laskin, Misha, Abbeel, Pieter, Srinivas, Aravind, and Mordatch, Igor. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- [7] Cobbe, Karl, Hesse, Chris, Hilton, Jacob, and Schulman, John. Leveraging procedural generation to benchmark reinforcement learning. In *International conference on machine learning*, pp. 2048–2056. PMLR, 2020.
- [8] Cobbe, Karl W, Hilton, Jacob, Klimov, Oleg, and Schulman, John. Phasic policy gradient. In *International Conference on Machine Learning*, pp. 2020–2027. PMLR, 2021.
- [9] Di Langosco, Lauro Langosco, Koch, Jack, Sharkey, Lee D, Pfau, Jacob, and Krueger, David. Goal misgeneralization in deep reinforcement learning. In *International Conference on Machine Learning*, pp. 12004–12019. PMLR, 2022.
- [10] Driess, Danny, Xia, Fei, Sajjadi, Mehdi SM, Lynch, Corey, Chowdhery, Aakanksha, Ichter, Brian, Wahid, Ayzaan, Tompson, Jonathan, Vuong, Quan, Yu, Tianhe, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- [11] Gao, Peng, Geng, Shijie, Zhang, Renrui, Ma, Teli, Fang, Rongyao, Zhang, Yongfeng, Li, Hongsheng, and Qiao, Yu. Clip-adapter: Better vision-language models with feature adapters. *arXiv preprint arXiv:2110.04544*, 2021.
- [12] Geng, Xinyang, Liu, Hao, Lee, Lisa, Schuurams, Dale, Levine, Sergey, and Abbeel, Pieter. Multi-modal masked autoencoders learn transferable representations. *arXiv preprint arXiv:2205.14204*, 2022.
- [13] Huang, Wenlong, Abbeel, Pieter, Pathak, Deepak, and Mordatch, Igor. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *arXiv preprint arXiv:2201.07207*, 2022.
- [14] Huang, Wenlong, Xia, Fei, Xiao, Ted, Chan, Harris, Liang, Jacky, Florence, Pete, Zeng, Andy, Tompson, Jonathan, Mordatch, Igor, Chebotar, Yevgen, Sermanet, Pierre, Jackson, Tomas, Brown, Noah, Luu, Linda, Levine, Sergey, Hausman, Karol, and brian ichter. Inner monologue: Embodied reasoning through planning with language models. In *6th Annual Conference on Robot Learning*, 2022. URL <https://openreview.net/forum?id=3R3Pz5i0tye>.
- [15] Kim, Changyeon, Park, Jongjin, Shin, Jinwoo, Lee, Honglak, Abbeel, Pieter, and Lee, Kimin. Preference transformer: Modeling human preferences using transformers for RL. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Peot1SFDX0>.

- 484 [16] Liu, Hao, Lee, Lisa, Lee, Kimin, and Abbeel, Pieter. Instruction-following agents with jointly
485 pre-trained vision-language models. *arXiv preprint arXiv:2210.13431*, 2022.
- 486 [17] Liu, Yuqi, Xiong, Pengfei, Xu, Luhui, Cao, Shengming, and Jin, Qin. Ts2-net: Token shift
487 and selection transformer for text-video retrieval. In *Computer Vision—ECCV 2022: 17th*
488 *European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIV*, pp. 319–
489 335. Springer, 2022.
- 490 [18] Loshchilov, Ilya and Hutter, Frank. Decoupled weight decay regularization. In *International*
491 *Conference on Learning Representations*, 2019. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Bkg6RiCqY7)
492 [id=Bkg6RiCqY7](https://openreview.net/forum?id=Bkg6RiCqY7).
- 493 [19] Luo, Huaishao, Ji, Lei, Zhong, Ming, Chen, Yang, Lei, Wen, Duan, Nan, and Li, Tianrui.
494 Clip4clip: An empirical study of clip for end to end video clip retrieval and captioning. *Neuro-*
495 *computing*, 508:293–304, 2022.
- 496 [20] Pathak, Deepak, Gandhi, Dhiraj, and Gupta, Abhinav. Self-supervised exploration via disagree-
497 ment. In *International Conference on Machine Learning*, 2019.
- 498 [21] Rasheed, Hanoona, khattak, Muhammad Uzair, Maaz, Muhammad, Khan, Salman, and Khan,
499 Fahad Shahbaz. Finetuned clip models are efficient video learners. In *The IEEE/CVF Conference*
500 *on Computer Vision and Pattern Recognition*, 2023.
- 501 [22] Seo, Younggyo, Chen, Lili, Shin, Jinwoo, Lee, Honglak, Abbeel, Pieter, and Lee, Kimin.
502 State entropy maximization with random encoders for efficient exploration. In *International*
503 *Conference on Machine Learning*, pp. 9443–9454. PMLR, 2021.
- 504 [23] Shah, Dhruv, Osiński, Błażej, Levine, Sergey, et al. Lm-nav: Robotic navigation with large
505 pre-trained models of language, vision, and action. In *Conference on Robot Learning*, pp.
506 492–504. PMLR, 2023.
- 507 [24] Wang, Yi, Li, Kunchang, Li, Yizhuo, He, Yinan, Huang, Bingkun, Zhao, Zhiyu, Zhang, Hongjie,
508 Xu, Jilan, Liu, Yi, Wang, Zun, et al. Internvideo: General video foundation models via generative
509 and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022.