

Student-Informed Teacher Training

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Imitation learning with a privileged teacher has proven effective for
2 learning complex control behaviors from high-dimensional inputs, such as im-
3 ages. In this framework, a teacher is trained with privileged task information,
4 while a student tries to predict the actions of the teacher with more limited obser-
5 vations, e.g., in a robot navigation task, the teacher might have access to distances
6 to nearby obstacles, while the student only receives visual observations of the
7 scene. However, privileged imitation learning faces a key challenge: the student
8 might be unable to imitate the teacher’s behavior due to partial observability. This
9 problem arises because the teacher is trained without considering if the student is
10 capable of imitating the learned behavior. To address this teacher-student asym-
11 metry, we propose a framework for joint training of the teacher and student poli-
12 cies, encouraging the teacher to learn behaviors that can be imitated by the student
13 despite the latter’s limited access to information and its partial observability.

14 1 Introduction

15 State-of-the-art (SotA) policy learning approaches often rely on imitation learning to accelerate
16 training [1, 2, 3, 4]. However, collecting expert demonstrations for imitation learning can be pro-
17 hibitively expensive, which has led to the development of the teacher-student framework. In this
18 framework, expert data is generated automatically by training a teacher policy using RL on priv-
19 ileged task information, benefiting from efficient simulation and a faster learning process. This
20 approach eliminates the need for the student to extensively explore the environment, which can be a
21 very challenging process when dealing with high-dimensional observations, such as images. How-
22 ever, privileged imitation learning can be hindered by information asymmetry between the teacher
23 and student, where the student receives less informative observations and struggles to imitate the
24 behavior of the teacher [5]. As a consequence of the information asymmetry, the teacher tends
25 to over-rely on its full observability of the environment without considering the more limited ob-
26 servation space of the student. This causes the teacher to provide target actions that the student
27 cannot infer from its observations, since the student lacks access to the same level of environmental
28 information. To tackle these challenges, we propose a teacher-student knowledge distillation frame-
29 work that encourages the teacher to learn behaviors that account for the capabilities of the student.
30 Specifically, the objective function of the teacher is extended by adding the upper bound of the stu-
31 dent performance within the imitation learning setting. This results in a reward term that penalizes
32 the teacher for visiting states where there is a significant action mismatch between the student and
33 teacher. Additionally, minimizing this upper bound leads to a second optimization term that directly
34 supervises the weights of the teacher network.

35 2 Student-Informed Teacher Training

36 As shown in [6, 7, 8], the performance gap between student and teacher is upper-bounded by the
37 action difference between both policies. Thus, minimizing the action difference under the state
38 distribution of the expert also minimizes the performance gap $J(\pi_T) - J(\pi)$. Instead of trying to
39 minimize the action difference by adjusting the student policy π_S , we propose to change the per-
40 spective and find a teacher policy π_T optimizing for the task reward while considering the alignment

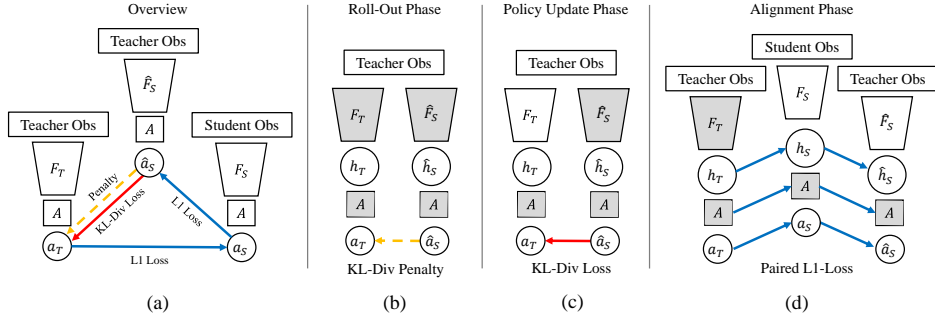


Figure 1: **Method Overview.** Our method leverages three networks (a), which are trained in three alternating phases: the roll-out phase (b), the policy update phase (c), and the alignment phase (d). The grey boxes represent networks frozen during the specific phase, i.e., the network parameters are not updated. (b) In the roll-out phase, the KL-Divergence between the proxy student \hat{F}_S and teacher F_T is used as a penalty term. (c) In addition to the policy gradient, the teacher encoder is updated by backpropagating through the KL-Divergence between the action distribution of the teacher and the proxy student. (d) Using student observations, the proxy student is aligned to the student F_S while the student is aligned to the teacher network using consistency losses.

41 between teacher and student. Thus, we want to find a teacher policy π_T that maximizes

$$\tilde{J}(\pi_T) = \mathbb{E}_{s \sim d_{\pi_T}, a \sim \pi_T(\cdot|s)}[r(s, a)] - \mathbb{E}_{s \sim d_{\pi_T}}[D_{KL}(\pi_T(\cdot|s), \pi_S(\cdot|s))] \quad (1)$$

$$= \mathbb{E}_{s \sim d_{\pi_T}, a \sim \pi_T(\cdot|s)}[r(s, a) - D_{KL}(\pi_T(\cdot|s), \pi_S(\cdot|s))] \quad (2)$$

$$\propto \mathbb{E}_{\tau \sim p_\theta}[R(\tau) - D_\theta(\tau)] = \int p_\theta(\tau)(R(\tau) - D_\theta(\tau))d\tau. \quad (3)$$

42 In the last step, the discounted state distribution is changed to the expectation over trajectories $\tau \sim$
 43 p_θ , which are induced by the expert policy π_T and represent the state and corresponding actions
 44 $\tau = \{s_0, a_0, s_1, a_1, \dots\}$. Additionally, we define the return $R(\tau) = \sum_{s_t, a_t \in \tau} \gamma^t r(s_t, a_t)$ and the
 45 sum of discounted KL-Divergences $D_\theta(\tau) = \sum_{s_t \in \tau} \gamma^t D_{KL}(\pi_T(\cdot|s_t), \pi_S(\cdot|s_t))$. We use subscript
 46 θ , to emphasize that the probability distribution over the trajectories p_θ and $D_\theta(\tau)$ is dependent on
 47 the parameter of the teacher network θ . Following the classical policy gradient to obtain the optimal
 48 policy, we take the gradient of Eq. 3 with respect to the teacher parameters θ

$$\nabla_\theta \tilde{J}(\pi_T) = \nabla_\theta \int p_\theta(\tau)(R(\tau) - D_\theta(\tau))d\tau \quad (4)$$

$$= \int \nabla_\theta p_\theta(\tau)R(\tau)d\tau - \int \nabla_\theta p_\theta(\tau)D_\theta(\tau)d\tau - \int p_\theta(\tau)\nabla_\theta D_\theta(\tau)d\tau \quad (5)$$

$$= \underbrace{\int \nabla_\theta p_\theta(\tau)(R(\tau) - D_\theta(\tau))d\tau}_{\text{Policy Gradient}} - \underbrace{\int p_\theta(\tau)\nabla_\theta D_\theta(\tau)d\tau}_{\text{KL-Div Gradient}}. \quad (6)$$

49 As can be observed, we end up with the standard policy gradient optimizing the task reward while
 50 also considering the teacher-student misalignment for each trajectory. This weighted KL-Divergence
 51 D_θ can be interpreted as a reward encouraging the teacher policy to visit states where the student
 52 and teacher are aligned and avoid states with a large misalignment. The second term contains the
 53 expectation of the gradient with respect to the teacher network over the expert states, which rep-
 54 represents a direct supervision on the teacher weights by enforcing the prediction of the same action
 55 distribution as the student.

56 3 Joint Learning Framework

57 Building on the formulation in Sec. 2, we propose a practical framework to tackle the teacher-
 58 student asymmetry. Following Eq. 6, we adapt the widely-used PPO algorithm [9] to train the
 59 teacher to learn behaviors that can be imitated by the student. At the same time, we train the student
 60 network to imitate the teacher based on a subset of the collected environment interactions containing

61 student and teacher observations. By using a subset of paired teacher and student data, we avoid the
62 (usually expensive) simulation of student observation for each time step the teacher interacts with
63 the environment. An overview of the proposed method is shown in Figure 1 a).

64 To implement the objective in Eq. 3 inside the teacher training, we implement two key components:
65 (i) a proxy student network taking as input teacher observations and (ii) a shared action decoder
66 network. Excluding the critic, our method consists of three different networks: the teacher network
67 F_T , the student network F_S , and the proxy student network \hat{F}_S , which all share the same action
68 decoder network A .

69 **Proxy Student Network** To compute the action difference used in the penalty term and to obtain the
70 KL-Div Gradient in Eq. 6, a forward pass through both the teacher and student networks is required
71 for each collected sample. However, simulating high-dimensional student observations, such as
72 images, is often computationally expensive, contradicting the initial goal of accelerating training.
73 To avoid this simulation overhead, we introduce a separate neural network \hat{F}_S that imitates the
74 current student policy based on the teacher observations. This allows us to approximate the actions
75 of the student at each expert state without additional simulation cost. The proxy student network is
76 trained during the alignment phase, where both student and teacher observations are available for a
77 subset of environment interactions. Our proposed framework consists of three alternating training
78 phases: (i) the classical policy roll-out, (ii) the policy update, and (iii) the alignment phase, see also
79 Figure 1 (b)-(c). The first two phases follow the standard on-policy training, while in the alignment
80 phase (iii), the student F_S is aligned to the teacher F_T , and the proxy student \hat{F}_S is aligned to the
81 student F_S . For both network alignments, paired student and teacher observations are used. In the
82 following, we provide more details about the specific training phases.

83 **Roll-out Phase** Our proposed framework introduces a minor modification to the roll-out phase of
84 the standard teacher training, specifically in the reward computation. In addition to the task reward,
85 we also add a penalty term computed based on the action difference between the teacher and the
86 proxy student. This penalty encourages the teacher to only visit states in which the student can pre-
87 dict the same actions as the teacher, thereby improving alignment between the student and teacher.
88 Furthermore, during each roll-out phase, we store a subset of expert states required in the alignment
89 phase. Depending on the simulation environment, this subset can be randomly selected states or a
90 fixed number of environments from which student observations are generated.

91 **Policy Update Phase** The gradient of the KL-Div term in Eq. 6 can be integrated into the policy up-
92 date of the teacher, during which the network weights are updated using the clipped policy gradient
93 of PPO. Since both the policy gradient and the KL-Divergence gradient are computed over the state
94 distribution of the teacher, we can use the teacher states inside the roll-out buffer to compute the
95 KL-Divergence between the action distributions of the teacher and the proxy student. This allows
96 us to update the network parameters of the teacher in a single backward pass through the combined
97 loss function.

98 **Alignment Phase** The alignment phase focuses on aligning the features across the encoders of the
99 teacher, proxy student, and student. This phase is the only one that requires paired teacher and
100 student observations, which are simulated from a subset of the teacher’s experiences during the roll-
101 out phase. We align the student encoder with the teacher encoder by computing the L1 loss between
102 their corresponding features and between the activations of the frozen shared action decoder. To
103 prevent the collapse of the model into predicting constant outputs, gradients are only backpropagated
104 to the student encoder. Similarly, the proxy student is aligned with the student using the L1 loss on
105 the encoded features, with gradients only backpropagated to the proxy student. The parameters of
106 the teacher remain unaffected during this phase and are only updated during the policy update phase.

107 4 Experiments

108 We evaluate our student-informed teacher training framework on the task of vision-based drawer
109 opening using a robot arm. We compare our method to multiple behavior cloning (BC) and DAgger
110 and ablate the benefits of the alignment penalty and loss terms in all three tasks.

111 **Setup** We adapt the publicly available *Omniverse Isaac Gym Reinforcement Learning Environments*

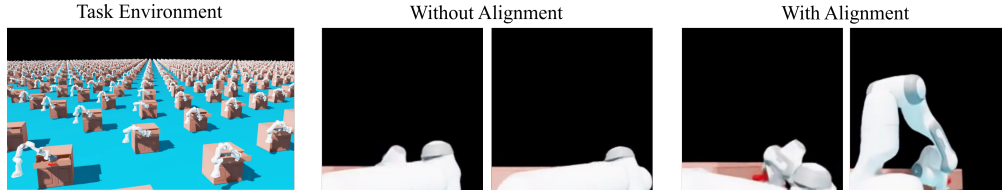


Figure 2: **Vision-Based Manipulation.** On the left, the task of opening a drawer with a robotic arm is visualized for all of the parallel environments. The two images in the center (Without Alignment) are sample images given to the student, which show the teacher behaviors trained without our alignment. Our approach with alignment leads to behaviors that the student can imitate more easily, i.e., the robot does not block the red drawer handle, as visualized in the two images on the right.

Methods	Success Rate
BC	0.27 ± 0.02
DAGger	0.31 ± 0.24
w/o Align (Ours)	0.60 ± 0.04
w Align (Ours)	0.77 ± 0.12

Table 1: **Manipulation Success Rates.** The mean and standard deviation of the success rate for the task of opening the drawer obtained from three different trainings runs.

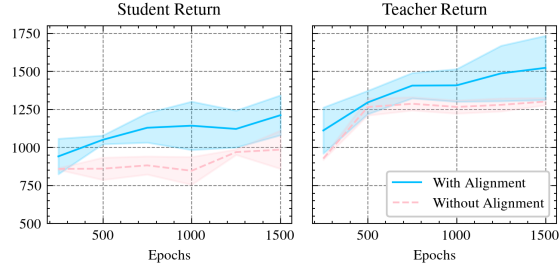


Figure 3: **Manipulation Returns.** The mean returns achieved by the student and teacher trained with and without alignment, averaged over three training runs.

112 for Isaac Sim repository, specifically modifying the cabinet-opening task. In this task, a Franka
 113 robot arm is trained to open a drawer that contains objects inside. The teacher receives the de-
 114 fault observations, which include the robot’s joint positions and velocities and the relative position
 115 between the gripper and the drawer handle, providing privileged information for effective manipu-
 116 lation. The student observations do not include relative distance information, and instead, an image
 117 is given together with the state of the robot arm as observation. Instead of an unobstructed top view,
 118 the viewpoint of the camera is selected closer to the robot arm, which enables self-occlusion. This
 119 camera setup is closer to real-world applications on mobile robots, such as humanoid robots, where
 120 certain arm configurations may obstruct cameras. We use a frozen DINOv2 [10] encoder to extract
 121 flattened image features, which are passed through a five-layer MLP before being fed into the shared
 122 action decoder, which comprises one layer. For the teacher and proxy student, we use a three-layer
 123 MLP to process the 1D observations.

124 Table 1 reports the success rates for our method (with and without alignment), a DAGger-trained
 125 student, and a BC student across three training runs. Our framework with and without alignment
 126 significantly outperforms students trained with DAGger and BC, with success rates of up to 0.77
 127 compared to 0.47 (DAGger) and 0.27 (BC). This improvement can be explained by the shared task
 128 decoder, which is trained by leveraging the training samples of the teacher. Our method with align-
 129 ment improves student success rates by 17% compared to the non-aligned framework while also
 130 consistently achieving higher returns (Figure 3). These results demonstrate that our framework
 131 helps the teacher learn better behaviors for the student. Student policies trained without alignment
 132 achieve non-zero success rates due to the small sampling interval of the cabinet position. This al-
 133 lows them to memorize behaviors without relying heavily on the images. All teachers, regardless of
 134 alignment, reach a 100% success rate. Interestingly, the return of the teacher trained with alignment
 135 is also constantly higher than without alignment. A possible explanation is that the teacher learns
 136 gripper movements optimized for robustness without trying multiple times to grab the handle, which
 137 is a difficult behavior for the student lacking relative pose information. As can be seen in Figure 2,
 138 our method leads to several different behaviors, which also explains the high variance of success
 139 rates and returns. With alignment, the teacher learns once to grab the handle from a top-down con-
 140 figuration while another teacher lowers its first two elements to make the red handle visible. In both
 141 cases, the red handle is visible right before the gripper touches it. This shows that our alignment
 142 leads to emerging teacher behaviors that consider the imitation difficulties of the student.

References

- 143
- 144 [1] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy:
145 Visuomotor policy learning via action diffusion. *Robotics: Science and Systems*, 2023.
- 146 [2] J. Luo, C. Xu, X. Geng, G. Feng, K. Fang, L. Tan, S. Schaal, and S. Levine. Multi-stage cable
147 routing through hierarchical imitation learning. *IEEE Transactions on Robotics*, 2024.
- 148 [3] O. M. Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna,
149 T. Kreiman, C. Xu, et al. Octo: An open-source generalist robot policy. *Robotics: Science and
150 Systems*, 2024.
- 151 [4] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine. Scaling cross-embodied learning:
152 One policy for manipulation, navigation, locomotion and aviation. *8th Conference on Robot
153 Learning*, 2024.
- 154 [5] H. Nguyen, A. Baisero, D. Wang, C. Amato, and R. Platt. Leveraging fully observable policies
155 for learning under partial observability. *arXiv preprint arXiv:2211.01991*, 2022.
- 156 [6] T. Xu, Z. Li, and Y. Yu. Error bounds of imitating policies and environments. In
157 H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in
158 Neural Information Processing Systems*, volume 33, pages 15737–15749. Curran Asso-
159 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/
160 2020/file/b5c01503041b70d41d80e3dbe31bbd8c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/b5c01503041b70d41d80e3dbe31bbd8c-Paper.pdf).
- 161 [7] U. Syed and R. E. Schapire. A reduction from apprenticeship learning to classification. *Ad-
162 vances in neural information processing systems*, 23, 2010.
- 163 [8] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured predic-
164 tion to no-regret online learning. In *Proceedings of the fourteenth international conference
165 on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Pro-
166 ceedings, 2011.
- 167 [9] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
168 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 169 [10] M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez,
170 D. Haziza, F. Massa, A. El-Nouby, R. Howes, P.-Y. Huang, H. Xu, V. Sharma, S.-W. Li,
171 W. Galuba, M. Rabbat, M. Assran, N. Ballas, G. Synnaeve, I. Misra, H. Jegou, J. Mairal,
172 P. Labatut, A. Joulin, and P. Bojanowski. Dinov2: Learning robust visual features without
173 supervision, 2023.