

Supplementary Material

1 Formulations for Algorithm

This section serves as an introduction for the underlying formulation of the algorithm. Since we use a meta-learning algorithm to train the programmatic policy, the training procedure consists of an outer loop and an inner loop. We use Reptile as the meta-learning framework to update the parameter in outer loop. In inner loop, we use Proximal Policy Optimization (PPO) as the basic reinforcement learning (RL) algorithm to learn "fast" parameter.

In our method, we use PPO-Clip, which utilizes a clipping mechanism in the objective function to remove incentives that may prompt a new policy towards diverging greatly from an existing policy. The PPO-Clip algorithm does not incorporate a KL-divergence term in its objective function. The PPO-Clip objective is:

$$\begin{aligned}\theta_{k+1} &= \arg \max_{\theta} J(\theta) \\ &= \arg \max_{\theta} \frac{1}{|\mathcal{D}_k| T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \min \left(\frac{\mathcal{P}_{\theta}(a_t | s_t)}{\mathcal{P}_{\theta_k}(a_t | s_t)} A^{\mathcal{P}_{\theta_k}}(s_t, a_t), g(\epsilon, A^{\mathcal{P}_{\theta_k}}(s_t, a_t)) \right)\end{aligned}$$

, where

$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \\ (1 - \epsilon)A & A < 0 \end{cases}.$$

To illustrate, \mathcal{D}_k represents a set of trajectories by running the policy \mathcal{P}_{θ_k} in the environment. The $A^{\mathcal{P}_{\theta_k}}$ is the advantage estimation based on the value function and the hyperparameter ϵ corresponds to how far away the new policy can go from the old while still profiting the objective. The θ is the parameter for the programmatic policy \mathcal{P}_{θ} , which indeed are two parameters $\theta = (\mathcal{W}, \varphi)$. The \mathcal{W} represents the policy architecture as well as the φ is the parameter for policy actions. These two parameters are jointly optimized by PPO. The train process is a bilevel iterative optimization process. At each iteration k of training, we carry out two steps. During the first step, we optimize the φ while preserving the frozen architecture parameter weights \mathcal{W} :

$$\varphi_{k+1} = \arg \max_{\varphi} J(\mathcal{W}_k, \varphi)$$

Next, during the second step, we optimize the \mathcal{W} by fixing φ :

$$\mathcal{W}_{k+1} = \arg \max_{\mathcal{W}} J(\mathcal{W}, \varphi_k)$$

The bilevel training steps alternate throughout the training iterations until the reward converges. Then we choose the optimal \mathcal{W} to freezing the architecture of the programmatic policy. Finally, we train the parameters of the selected architecture, continuously until the parameter values learned by means of RL converge to $\hat{\theta}$.

In the outer loop, the "fast" parameter is used for updating the "slow" parameter. Multiple gradient descents are applied on each task in order to obtain the corresponding parameter value $\hat{\theta}_i$. Afterwards, the difference vector between the parameters of each task and the main task is calculated as the update direction:

$$\theta_{i+1} \leftarrow \beta(\theta_i - \hat{\theta}_i)$$

, where β is the meta learning rate. By repeating this process iteratively, the global initialized parameters are finally obtained. Intuitively, by using the gradient of a single task's parameter as the rough direction of gradient descent and decreasing the total loss function of the training tasks roughly but steadily, one can often obtain a good initialization parameter.

2 Implementation Details

In this section, we present the details of implementation, which cover the input and reward structures of each environment, architecture design and training hyperparameters.

2.1 Environment Details

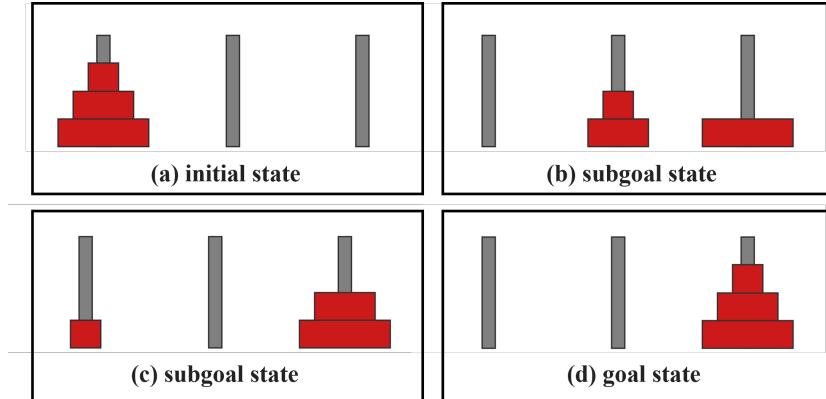


Figure 1: Details for sub-goals in hanoi

Hanoi In the Tower of Hanoi environment, there are three adjacent pillars labeled A, B, C, and several different-sized discs stacked in a pyramid shape from bottom to top on pillar A. The objective is to move all the discs one by one to pillar B, but at no time

can a larger disc be placed on top of a smaller one on the same pillar. Each action is to move a disc from the current pillar to another pillar.

In the Tower of Hanoi environment, the variation is the number and size of discs. For different environment variations, the goal of the task is to move all the discs to another pillar. Based on the number of discs on each pillar, as well as the location of the discs, based on their size, we encode the input of this environment into a vector with a dimension of 1×9 . Regardless of the number of discs, we can only set 6 actions, i.e., A to B, A to C, B to A, B to C, C to A, and C to B.

The Tower of Hanoi environment is a discrete environment with low input and output dimensions. However, completing such a task requires the agent to sequentially complete a series of actions, and a small number of mistakes in the action sequence can result in very low rewards for the whole task, or even the inability to complete it. Thus, we set some intermediate rewards for hanoi environment to assist the agent in learning the policy more easily. As shown in Figure 1, for a Tower of Hanoi environment with three discs, we set (b) and (c) as sub-goals, the environment returns a small reward when the agent accomplishes this for the first time. When the agent moves a disc once, the reward value is -1. When the agent performs an invalid action, such as placing a large disc on top of a small disc, the environment will judge that the action has failed, keep the environment state unchanged, and give a penalty of reward value -2 to the agent. When the agent completes the task, i.e. placing the discs in the state shown in (d), the environment gives a large reward value.

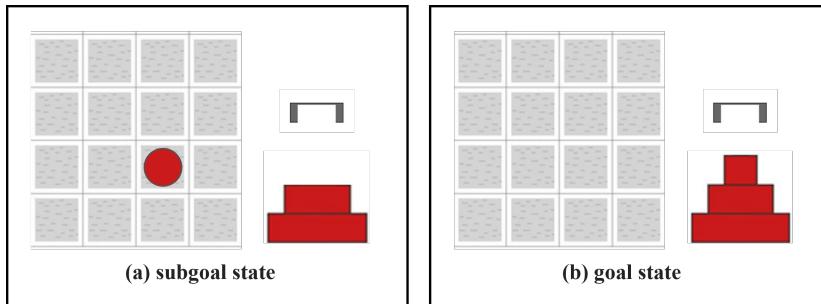


Figure 2: Details for sub-goals in stacking

Stacking In the Stacking environment, there are several differently-sized plates on the table, and a gripper that can move above the table. The goal is to collect them one by one using the gripper, putting them together in descending order of size.

In the Stacking environment, the variations of the environment are the number, size, and location of the plates. To simplify, we use a grid to represent the table. The gripper can move on the grid in up, down, left, and right directions, and can grab a plate. In this environment, we encode the table status, the stacked plates and the location of the gripper as the environment state, with a dimension of 1×35 . The gripper can perform actions of moving up, down, left, and right and selecting a plate to stack and the action is a 5-dimension vector.

The Stacking environment is similar to the Tower of Hanoi environment, as the agent needs to control the gripper to select plates in a specific order to complete the task. To reduce the sparseness of rewards, we also set intermediate rewards. For example, as shown in (a) of Figure 2, when the agent first completes stacking two plates, it receives a small reward. Each action taken by the gripper results in a reward value of -1. When the gripper performs an incorrect action, such as moving out of the table or trying to stack plates of wrong size, the environment will judge that the action has failed. Then the environment state keeps unchanged, and give a penalty of -2 to the agent. When the agent completes the task, i.e. stacking the plates in the state shown in (b), the environment gives a positive final reward.

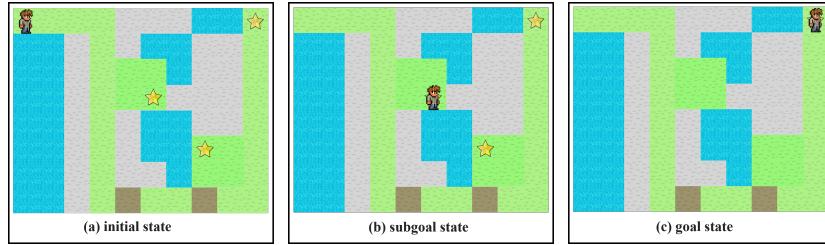


Figure 3: Details for sub-goals in hiking

Hiking In the Hiking environment, the characters on the map represent the agent and the stars on the map are the targets that the agents need to collect. The character can move on the ground, with green and gray blocks indicating walkable ground and blue water blocks representing obstacles that the character cannot pass through. The goal in this environment is to make the character move step by step and collect all the stars on the map.

The variations in this environment lies in the the number and position of the stars. The environment is also discretized into a grid. The entire state of the map is encoded as the state of the environment, with the dimension of 6×7 . The character, ground, and water on the map are all marked with different symbols. The character's actions are up, down, left, and right, so the dimension of the action vector for the environment is 1×4 .

In the Hiking environment, we also set up intermediate rewards, as shown in (b) of Figure 3, where the environment gives the agent a intermmediate reward when the character moves to a star and the star disappears from the map. Every time the character take an action in this environment, the environment returns a reward of -1. When the character moves out of the map boundary or into the water, the environment judges the action as a failure and keeps the entire state unchanged, while giving a punishment of reward value -2 to the agent. The final state is shown as in (c).

PandaReach, PandaPush, PandaSlide, PandaStack PandaReach, PandaPush, PandaSlide and PandaStack are designed for robot manipulation, as in Figure 4. A Panda robot, equiped with a robot arm, is placed on a control console with a series of manipulable

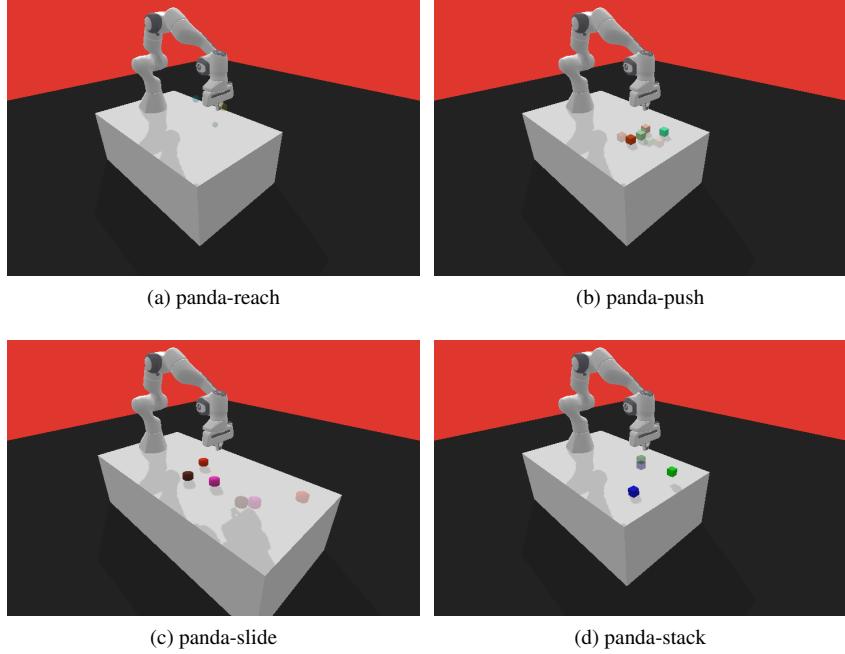


Figure 4: Benchmarks for continuous tasks

objects. The robot is required to manipulate the end-effector to finish the specified tasks. For example, in the environment of PandaPush, the robot needs to take actions to move the blocks to the target positions with the corresponding color. Once all the blocks are placed in the target positions, the task is completed.

The variations in this environment are derived from the number and position of the manipulable objects. Since the four environments have similar scenario, we consider the observation states uniformly. Through the integration of the position and velocity of the end-effector with the position, rotation, velocity, angular velocity of the manipulable objects, we encode the observation state as a 1×32 vector. The action of the robot is a 4-dimension vector, which contains a 3D coordinate for the end-effector and a finger control.

We set up intermediate rewards to improve the training. When one of the objects is placed to the target position, a positive return is provided. Each move of the robot will obtain a penalty return from the environments.

2.2 Hyperparameters

Except the programmatic policy, all the actor networks and value networks are constructed with a 3-layer linear networks. The input layer is used to process the input of environment state. The hidden size of the middle layer is 64. The dimension of the output layer is determined by the action space of the environment. For continuous tasks,

we use Gaussian distribution to construct the output layer. Following hyperparameters are used to train PPO algorithm.

- max train steps $2e6$.
- replay buffer of size 2048.
- mini-batch size 64.
- discount factor 0.99.
- Adam optimizer; actor learning rate $3e - 4$; critic learning rate $3e - 4$.
- clip parameter ϵ 0.2.
- Entropy coefficient 0.01.
- Max gradient norm 0.5.
- KL-Divergence limit 0.03.

Following hyperparameters are used to train SAC algorithm.

- max train steps $2e6$.
- replay buffer of size 512.
- mini-batch size 64.
- discount factor 0.99.
- tau 0.005.
- Adam optimizer; actor learning rate $2e - 4$; critic learning rate $2e - 4$.
- clip parameter ϵ 0.2.
- entropy coefficient 0.01.
- max gradient norm 0.5.
- KL-divergence limit 0.03.

Following hyperparameters are used to train ReptilePPO algorithm.

- max train steps $5e6$.
- meta iterations 250.
- inner iterations 1.
- max train steps $2e6$.
- replay buffer of size 2048.

- mini-batch size 64.
- discount factor 0.99.
- Adam optimizer; actor learning rate 2e-2; critic learning rate 2e-2.
- Adam optimizer; actor learning rate 1e-5; critic learning rate 1e-5.
- clip parameter ϵ 0.2.
- entropy coefficient 0.01.
- max gradient norm 0.5.
- KL-divergence limit 0.03.

The components of the programmatic policy are all linear cells. Following hyperparameters are used to train ReptilePRL algorithm as well as PRL.

- depth 6.
- max train steps 5e6.
- meta iterations 250.
- inner iterations 1.
- max train steps 5e6.
- replay buffer of size 2048.
- batch size
- mini-batch size 64.
- discount factor 0.99.
- Adam optimizer; actor learning rate 2e-2; critic learning rate 2e-2.
- Adam optimizer; actor learning rate 1e-4; critic learning rate 1e-4.
- clip parameter ϵ 0.2.
- entropy coefficient 0.01.
- max gradient norm 0.5.
- KL-divergence limit 0.03.

In our method, compared with ReptilePRL algorithm, some of the components are GRU layers. And in the updating process, the state of the RNN will be preserved across different episodes. Besides, due to the use of RNN, we have slightly changed the structure of the replay buffer. Following hyperparameters are used to train our method.

- depth 6.

- max train steps 5e6.
- meta iterations 250.
- inner iterations 1.
- max train steps 2e6.
- replay buffer of size 500.
- batch size 16
- mini-batch size 2.
- Discount factor 0.99.
- Adam optimizer; actor learning rate 2e-2; critic learning rate 2e-2.
- Adam optimizer; actor learning rate 4e-4; critic learning rate 4e-4.
- clip parameter ϵ 0.2.
- entropy coefficient 0.01.
- max gradient norm 0.5.
- KL-divergence limit 0.03.

3 Comprehensive Experimental Results

This section aims to present the comprehensive experimental results. In table 1 and table 2, we evaluated the PPO, SAC, HIRO, PRL, ReptilePPO, ReptilePRL and Our approach in seven benchmarks—Hanoi, Hiking, Stacking for task planning and PandaReach, PandaPush, PandaSlide and PandaStack for robot manipulation. The number of steps for the agents is utilized as the performance metric. A smaller number of action steps indicates better performance. The performance is measured by mean number of steps an agent takes to achieve goals, along with standard deviations. The best score is marked with bold font.

To conclude, our method basically showcases the best performance across the benchmarks. In some tasks, the best scores are achieved by ReptilePPO or ReptilePRL. It means that the policies trained by a meta-learning algorithm are capable of generalizing to different tasks in one task domain. In addition, when evaluated in the simple situations, such as Hiking1, Stacking1 or PandaReach1, the policies trained by SAC, HIRO or PRL are evaluated as the best scores. However, our method performs better in more complex environments according to the table. This indicates that our method has the ability to synthesize policies that show strong generalization in complex tasks. Even when confronted with unseen tasks from the same task domain, the policy learned by our method achieves high scores. This is most likely because the policy can capture reusable subroutines and states patterns.

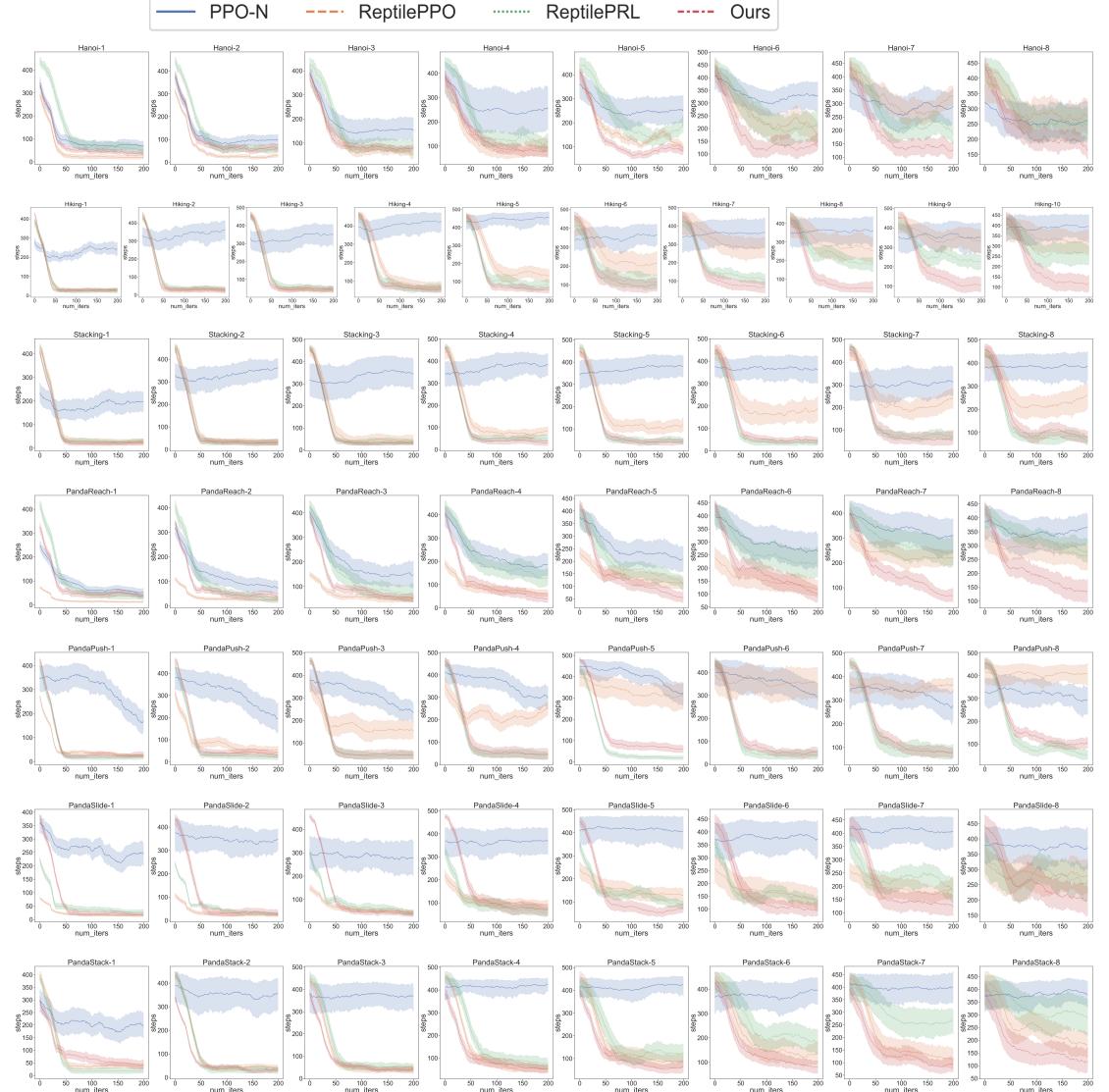


Figure 5: Comparison of the training curve for **PPO-N**, **ReptilePPO**, **ReptilePRL** and **Ours** in the seven benchmarks. Results are averaged over 10 random seeds.

The comprehensive results of the ablation study are illustrated as Figure 5. The vertical axis represents the number of steps required for the agent to achieve its objective, while the horizontal axis denotes the number of iterations involving agent-environment interaction episodes.

The policies learned by PPO-N struggle to achieving convergence in all benchmarks which indicates the necessity of a meta-learning algorithm. Overall, ReptilePRL

and ReptilePPO perform well under different benchmarks respectively, according to Figure 5. Overall, our approach performs better than ReptilePRL and ReptilePPO in most benchmarks. The more complex the environment, the better our methods perform. In unseen tasks, the ability to learn and reuse subroutines becomes very important. Our proposed idea of the combination of a meta-learning framework with a designed programmatic policy and the utilization of RNN blocks in policy architecture, is very useful.

test envs	PPO	SAC	HIRO	PRL	PPO	SAC	HIRO	PRL	PPO	SAC	HIRO	PRL
train envs: Hanoi-1												
Hanoi-1	212	176	284	223	53	82	8	7	27	37	19	217
Hanoi-2	331	271	394	102	137	24	24	100	119	53	53	384
Hanoi-3	444	329	481	165	137	81	110	187	241	95	95	423
Hanoi-4	466	413	491	228	227	153	222	247	267	153	153	483
Hanoi-5	491	424	499	487	271	284	231	294	290	257	223	493
Hanoi-6	493	402	500	345	288	319	355	342	299	297	251	499
Hanoi-7	499	427	499	497	379	311	383	412	386	387	311	500
Hanoi-8	500	370	499	500	412	345	381	434	399	343	366	500
train envs: Hanoi-2												
Hanoi-1	16	20	180	53	15	26	25	26	25	25	25	217
Hanoi-2	58	73	315	102	24	40	35	51	51	51	51	384
Hanoi-3	116	139	374	165	52	66	51	41	41	41	41	423
Hanoi-4	217	211	395	228	98	102	102	102	102	102	102	483
Hanoi-5	252	260	434	271	150	128	128	128	128	128	128	493
Hanoi-6	308	342	466	345	174	160	113	113	113	113	113	499
Hanoi-7	349	360	467	379	250	258	151	151	151	151	151	500
Hanoi-8	419	345	478	412	256	283	187	187	187	187	187	500
train envs: Hanoi-4												
Hanoi-1	16	20	180	53	15	26	25	26	25	25	25	217
Hanoi-2	58	73	315	102	24	40	35	51	51	51	51	384
Hanoi-3	116	139	374	165	52	66	51	41	41	41	41	423
Hanoi-4	217	211	395	228	98	102	102	102	102	102	102	483
Hanoi-5	252	260	434	271	150	128	128	128	128	128	128	493
Hanoi-6	308	342	466	345	174	160	113	113	113	113	113	499
Hanoi-7	349	360	467	379	250	258	151	151	151	151	151	500
Hanoi-8	419	345	478	412	256	283	187	187	187	187	187	500
train envs: Hanoi-1234												
Hanoi-1	16	20	180	53	15	26	25	26	25	25	25	217
Hanoi-2	58	73	315	102	24	40	35	51	51	51	51	384
Hanoi-3	116	139	374	165	52	66	51	41	41	41	41	423
Hanoi-4	217	211	395	228	98	102	102	102	102	102	102	483
Hanoi-5	252	260	434	271	150	128	128	128	128	128	128	493
Hanoi-6	308	342	466	345	174	160	113	113	113	113	113	499
Hanoi-7	349	360	467	379	250	258	151	151	151	151	151	500
Hanoi-8	419	345	478	412	256	283	187	187	187	187	187	500
train envs: Hiking-1												
Hiking-1	21.07 ± 13.92	25.33 ± 15.34	21.49 ± 14.06	21.52 ± 12.92	12.29 ± 7.85	14.03 ± 8.28	24.24 ± 14.61	72.43 ± 30.98	34.39 ± 20.58	34.74 ± 20.56	34.90 ± 21.36	13.92 ± 10.65
Hiking-2	22.35 ± 77.96	24.0 ± 23.0	73.67 ± 29.5	30 ± 82.92	38.26 ± 18.38	21.41 ± 15.20	52.78 ± 26.40	139.54 ± 53.61	121.17 ± 54.23	138.73 ± 53.38	105.56 ± 51.98	34.83 ± 24.54
Hiking-3	380.65 ± 81.55	311.43 ± 83.87	306.07 ± 86.98	357.43 ± 74.99	137.72 ± 49.15	137.28 ± 55.24	132.88 ± 58.79	197.81 ± 73.00	210.65 ± 61.66	220.81 ± 80.06	202.00 ± 57.02	117.88 ± 49.32
Hiking-4	374.77 ± 86.81	381.91 ± 85.01	319.57 ± 88.40	401.71 ± 79.83	227.94 ± 80.10	214.57 ± 76.21	227.00 ± 75.32	257.18 ± 83.75	252.40 ± 86.02	273.58 ± 89.51	213.98 ± 108.51	77.97 ± 108.51
Hiking-5	373.78 ± 85.43	338.88 ± 96.64	322.52 ± 94.46	402.69 ± 79.18	280.89 ± 89.26	275.07 ± 93.61	249.22 ± 79.19	306.33 ± 91.74	328.44 ± 66.49	368.17 ± 89.64	323.16 ± 61.49	267.31 ± 69.90
Hiking-6	367.13 ± 91.95	409.93 ± 81.68	409.23 ± 66.29	380.54 ± 85.24	360.78 ± 84.82	393.33 ± 70.38	321.71 ± 90.82	285.38 ± 80.83	321.71 ± 90.61	301.36 ± 96.33	319.04 ± 88.74	316.29 ± 90.69
Hiking-7	364.89 ± 92.28	360.70 ± 87.07	342.16 ± 93.41	362.11 ± 93.44	322.85 ± 94.04	342.52 ± 92.80	323.37 ± 92.62	328.89 ± 90.11	305.85 ± 91.81	370.23 ± 73.04	290.12 ± 91.32	372.27 ± 77.96
Hiking-8	353.80 ± 97.65	385.44 ± 89.61	379.94 ± 89.63	379.63 ± 88.51	404.62 ± 78.48	350.28 ± 92.84	342.77 ± 95.77	312.27 ± 89.26	368.32 ± 88.97	342.81 ± 92.32	348.86 ± 91.92	300.00 ± 91.92
Hiking-9	347.96 ± 97.67	382.71 ± 84.60	388.65 ± 87.41	379.67 ± 86.09	372.46 ± 85.51	364.27 ± 70.20	310.40 ± 86.14	195.27 ± 70.20	203.27 ± 90.25	197.81 ± 76.28	369.17 ± 83.43	320.51 ± 96.15
Hiking-10	355.69 ± 95.04	360.78 ± 95.82	376.14 ± 84.47	376.78 ± 90.42	304.82 ± 103.10	385.12 ± 77.07	352.51 ± 94.55	419.59 ± 69.93	373.66 ± 63.78	414.39 ± 73.58	378.64 ± 62.93	355.74 ± 61.30
train envs: Hiking-2												
Hiking-1	50.28 ± 24.61	43.22 ± 23.07	91.03 ± 34.63	22.65 ± 12.91	17.48 ± 11.56	23.05 ± 15.74	21.43 ± 15.19	17.43 ± 15.19	23.05 ± 15.74	21.43 ± 15.19	23.05 ± 15.74	19.37 ± 10.65
Hiking-2	98.18 ± 51.29	119.42 ± 55.93	93.88 ± 34.48	52.75 ± 24.98	31.99 ± 20.57	25.66 ± 17.45	31.14 ± 20.12	17.45 ± 17.45	31.14 ± 20.12	31.14 ± 20.12	31.14 ± 20.12	31.14 ± 20.12
Hiking-3	190.83 ± 73.43	152.45 ± 72.71	73.79 ± 13.62	21.56 ± 5.51	13.95 ± 7.19	47.21 ± 18.84	32.18 ± 16.36	18.84 ± 18.84	32.18 ± 16.36	32.18 ± 16.36	32.18 ± 16.36	32.18 ± 16.36
Hiking-4	253.99 ± 78.08	197.51 ± 85.01	209.76 ± 67.47	186.93 ± 79.02	107.11 ± 48.85	53.05 ± 29.52	33.42 ± 22.08	33.42 ± 22.08	33.42 ± 22.08	33.42 ± 22.08	33.42 ± 22.08	33.42 ± 22.08
Hiking-5	249.21 ± 88.59	307.84 ± 82.31	242.99 ± 89.89	245.57 ± 82.69	129.46 ± 46.49	51.85 ± 24.45	32.50 ± 16.10	32.50 ± 16.10	32.50 ± 16.10	32.50 ± 16.10	32.50 ± 16.10	32.50 ± 16.10
Hiking-6	205.35 ± 94.93	304.81 ± 88.85	286.80 ± 75.52	242.85 ± 86.73	190.41 ± 82.29	122.60 ± 57.10	49.54 ± 29.70	49.54 ± 29.70	49.54 ± 29.70	49.54 ± 29.70	49.54 ± 29.70	49.54 ± 29.70
Hiking-7	328.49 ± 86.19	300.20 ± 91.31	351.31 ± 79.39	88.03 ± 83.03	183.38 ± 66.13	85.75 ± 44.21	86.11 ± 44.21	86.11 ± 44.21	86.11 ± 44.21	86.11 ± 44.21	86.11 ± 44.21	86.11 ± 44.21
Hiking-8	307.88 ± 93.16	334.19 ± 75.55	286.90 ± 92.37	310.95 ± 92.37	310.95 ± 92.37	310.40 ± 86.14	195.27 ± 70.20	203.27 ± 90.25	197.81 ± 76.28	369.17 ± 83.43	320.51 ± 96.15	375.05 ± 91.59
Hiking-9	406.20 ± 68.53	382.85 ± 74.61	362.57 ± 90.48	445.42 ± 55.05	278.41 ± 93.47	214.42 ± 82.90	74.61 ± 49.33	74.61 ± 49.33	74.61 ± 49.33	74.61 ± 49.33	74.61 ± 49.33	74.61 ± 49.33
Hiking-10	391.34 ± 83.86	380.98 ± 75.06	435.03 ± 54.62	408.69 ± 69.59	325.13 ± 65.63	281.65 ± 51.50	67.86 ± 31.07	67.86 ± 31.07	67.86 ± 31.07	67.86 ± 31.07	67.86 ± 31.07	67.86 ± 31.07
train envs: Stacking-1												
Stacking-1	90.41 ± 41.76	87.51 ± 44.45	19.51 ± 11.43	166.28 ± 73.29	71.22 ± 32.48	74.28 ± 28.20	27.21 ± 18.82	26.19 ± 14.26	25.77 ± 17.88	32.08 ± 17.80	14.20 ± 10.46	19.37 ± 12.85
Stacking-2	240.60 ± 72.85	169.45 ± 68.58	47.80 ± 25.53	266.71 ± 94.60	126.27 ± 45.95	101.49 ± 44.17	48.74 ± 28.97	79.51 ± 41.77	46.87 ± 25.50	48.18 ± 26.84	32.19 ± 18.09	104.42 ± 40.26
Stacking-3	289.80 ± 84.97	248.66 ± 72.74	135.28 ± 53.73	425.83 ± 57.79	162.87 ± 52.82	144.99 ± 50.55	52.89 ± 35.30	126.83 ± 62.01	55.97 ± 26.02	80.80 ± 38.70	29.29 ± 15.58	114.96 ± 37.90
Stacking-4	345.80 ± 76.79	264.50 ± 93.31	231.37 ± 80.40	425.63 ± 68.33	136.72 ± 72.30	157.20 ± 66.69	86.39 ± 46.09	136.09 ± 62.68	87.34 ± 47.91	99.84 ± 47.91	64.03 ± 40.23	151.54 ± 71.11
Stacking-5	268.01 ± 95.69	279.98 ± 95.12	289.93 ± 91.97	326.34 ± 93.23	91.97 ± 22.60	226.09 ± 83.10	262.06 ± 82.81	130.11 ± 82.81	333.28 ± 86.11	147.54 ± 49.16	168.47 ± 73.32	136.06 ± 45.22
Stacking-6	326.05 ± 88.54	363.28 ± 78.66	337.54 ± 83.76	372.62 ± 90.63	267.98 ± 86.31	292.27 ± 90.25	197.81 ± 76.28	369.17 ± 83.43	201.48 ± 70.20	192.93 ± 74.46	170.47 ± 84.27	226.06 ± 64.15
Stacking-7	338.64 ± 88.71	307.31 ± 93.23	311.89 ± 98.22	354.78 ± 94.50	282.60 ± 91.27	293.55 ± 88.88	172.09 ± 77.38	333.48 ± 88.99	218.81 ± 83.58	226.31 ± 90.85	294.14 ± 85.79	290.35 ± 91.14
Stacking-8	318.77 ± 97.15	353.58	308.02 ± 73.55	409.03 ± 76.63	329.78 ± 79.07	300.53 ± 96.64	334.58 ± 86.22	434.79 ± 50.15	267.47 ± 66.48	322.58 ± 88.01	317.98 ± 71.15	341.84 ± 71.79
train envs: Stacking-2												
Stacking-1	13.41 ± 6.25	12.48 ± 4.97	72.08 ± 37.10	26.41 ± 14.20	25.16 ± 16.05	15.79 ± 8.81	33.42 ± 22.08	train envs: Stacking-3				
Stacking-2	27.35 ± 12.62	32.31 ± 11.39	71.22 ± 41.71	29.12 ± 17.36	39.25 ± 23.56	29.87 ± 18.09	20.64 ± 16.32	train envs: RepilePRL				

Table 2: Performance comparison for evaluating policies in PandaReach, PandaPush, PandaSlide and PandaStack, averaged over 5 random seeds.

test envs	PPO	SAC	HIRO	PRL	PPO	SAC	HIRO	PRL	PPO	SAC	HIRO	PRL	
train envs: PandaReach-1													
PandaReach-1	27.65 ± 15.34	26.02 ± 14.83	23.49 ± 13.82	30.9 ± 15.55	16.87 ± 8.45	19.11 ± 10.50	15.89 ± 10.70	24.97 ± 7.67	48.72 ± 24.93	86.49 ± 51.67	80.27 ± 51.55	24.74 ± 21.73	
PandaReach-2	215.13 ± 76.61	243.50 ± 83.84	229.22 ± 95.67	69.33 ± 77.11	59.21 ± 50.16	36.11 ± 18.80	33.70 ± 19.72	292.30 ± 82.68	48.87 ± 35.71	173.13 ± 97.27	139.91 ± 101.63	52.40 ± 34.68	
PandaReach-3	327.78 ± 85.91	318.04 ± 95.15	309.12 ± 94.56	67.97 ± 77.13	53.21 ± 50.16	126.36 ± 50.36	132.29 ± 48.02	379.95 ± 74.64	116.90 ± 61.86	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16	
PandaReach-4	385.07 ± 90.15	350.92 ± 90.15	344.04 ± 94.50	67.97 ± 77.13	53.21 ± 50.16	126.50 ± 76.91	126.62 ± 85.99	383.72 ± 86.40	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16	
PandaReach-5	413.22 ± 77.28	381.32 ± 86.85	428.32 ± 86.85	57.70 ± 77.39	53.35 ± 50.16	131.11 ± 77.23	131.82 ± 90.63	333.75 ± 86.73	230.23 ± 86.41	233.06 ± 74.24	222.15 ± 90.92	220.81 ± 61.86	
PandaReach-6	388.58 ± 89.73	405.11 ± 81.32	427.58 ± 86.85	57.80 ± 77.69	53.35 ± 50.16	131.20 ± 77.23	131.49 ± 89.82	302.89 ± 84.12	236.75 ± 86.43	280.60 ± 93.04	295.19 ± 87.74	268.13 ± 84.74	
PandaReach-7	383.77 ± 90.66	405.46 ± 55.39	336.20 ± 93.74	440.01 ± 68.43	351.70 ± 92.55	307.00 ± 94.53	302.29 ± 89.82	375.56 ± 84.34	270.14 ± 85.77	310.22 ± 78.27	346.92 ± 98.42	362.72 ± 82.23	
PandaReach-8	437.66 ± 55.39	336.20 ± 93.74	440.01 ± 68.43	351.70 ± 92.55	321.13 ± 96.69	346.57 ± 88.34	354.70 ± 94.26	223.86 ± 81.56	338.47 ± 65.41	313.26 ± 95.68	359.77 ± 62.23	440.76 ± 50.25	
train envs: PandaReach-4													
PandaReach-1	24.80 ± 15.34	26.02 ± 14.83	23.49 ± 13.82	31.50 ± 15.55	16.87 ± 8.45	19.11 ± 10.50	15.89 ± 10.70	24.97 ± 7.67	48.72 ± 24.93	86.49 ± 51.67	80.27 ± 51.55	24.74 ± 21.73	
PandaReach-2	42.46 ± 23.49	53.09 ± 33.55	40.21 ± 20.01	196.56 ± 91.42	30.92 ± 14.35	42.25 ± 47.77	42.12 ± 12.48	70.36 ± 30.16	40.35 ± 16.18	36.66 ± 21.26	36.66 ± 21.26	36.66 ± 21.26	
PandaReach-3	128.58 ± 46.00	127.28 ± 46.00	182.33 ± 80.03	97.12 ± 44.09	234.85 ± 93.90	80.81 \text{vtablashkash-1}	pm25.04	136.00 ± 58.69	57.72 ± 35.93	116.90 ± 61.86	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaReach-4	211.02 ± 76.84	168.89 ± 278.76	147.22 ± 58.62	332.58 ± 88.90	96.40 ± 27.72	174.34 ± 87.23	174.34 ± 87.23	220.98 ± 84.79	84.58 ± 47.08	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaReach-5	288.48 ± 89.40	313.39 ± 87.00	222.29 ± 78.94	332.22 ± 85.15	175.77 ± 62.65	233.20 ± 61.01	233.20 ± 61.01	288.01 ± 86.91	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaReach-6	331.08 ± 89.40	315.66 ± 89.36	220.20 ± 77.66	343.70 ± 82.83	233.20 ± 61.01	288.01 ± 86.91	233.20 ± 61.01	346.57 ± 88.34	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaReach-7	331.08 ± 89.40	315.66 ± 89.36	220.20 ± 77.66	343.70 ± 82.83	233.20 ± 61.01	288.01 ± 86.91	233.20 ± 61.01	347.97 ± 92.67	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaReach-8	406.25 ± 69.44	376.90 ± 75.28	224.63 ± 89.15	336.01 ± 93.85	373.00 ± 52.69	346.57 ± 88.34	346.57 ± 88.34	354.89 ± 94.26	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
train envs: PandaPush-4													
PandaPush-1	215.13 ± 76.61	243.50 ± 83.84	229.22 ± 95.67	69.33 ± 77.11	59.21 ± 50.16	36.11 ± 18.80	33.70 ± 19.72	292.30 ± 82.68	48.87 ± 35.71	116.90 ± 61.86	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaPush-2	327.78 ± 85.91	318.04 ± 95.15	309.12 ± 94.56	67.97 ± 77.13	53.21 ± 50.16	126.36 ± 50.36	132.29 ± 48.02	379.95 ± 74.64	116.90 ± 61.86	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16	124.49 ± 84.71
PandaPush-3	385.07 ± 90.15	350.92 ± 90.15	344.04 ± 94.50	67.97 ± 77.13	53.21 ± 50.16	126.50 ± 76.91	126.62 ± 85.99	383.72 ± 86.40	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16	124.49 ± 84.71
PandaPush-4	413.22 ± 77.28	381.32 ± 86.85	428.32 ± 86.85	57.70 ± 77.39	53.35 ± 50.16	131.11 ± 77.23	131.49 ± 89.82	320.89 ± 84.43	236.75 ± 86.43	280.60 ± 93.04	295.19 ± 87.74	268.13 ± 84.74	220.81 ± 61.86
PandaPush-5	388.58 ± 89.73	405.11 ± 81.32	427.58 ± 86.85	57.80 ± 77.69	53.35 ± 50.16	131.20 ± 77.23	131.49 ± 89.82	320.89 ± 84.43	236.75 ± 86.43	280.60 ± 93.04	295.19 ± 87.74	268.13 ± 84.74	220.81 ± 61.86
PandaPush-6	437.66 ± 55.39	336.20 ± 93.74	440.01 ± 68.43	351.70 ± 92.55	321.13 ± 96.69	346.57 ± 88.34	346.57 ± 88.34	354.89 ± 94.26	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaPush-7	24.80 ± 15.34	26.02 ± 14.83	23.49 ± 13.82	30.9 ± 15.55	16.87 ± 8.45	19.11 ± 10.50	15.89 ± 10.70	24.97 ± 7.67	48.72 ± 24.93	86.49 ± 51.67	80.27 ± 51.55	24.74 ± 21.73	
PandaPush-8	42.46 ± 23.49	53.09 ± 33.55	40.21 ± 20.01	196.56 ± 91.42	30.92 ± 14.35	42.25 ± 47.77	42.12 ± 12.48	70.36 ± 30.16	40.35 ± 16.18	36.66 ± 21.26	36.66 ± 21.26	36.66 ± 21.26	
train envs: PandaSlide-4													
PandaSlide-1	22.20 ± 18.91	23.49 ± 18.91	21.80 ± 18.91	32.01 ± 20.01	19.66 ± 91.42	30.92 ± 14.35	29.70 ± 14.35	42.31 ± 24.31	36.66 ± 21.26	39.47 ± 20.23	39.47 ± 20.23	39.47 ± 20.23	
PandaSlide-2	45.46 ± 33.49	53.09 ± 33.55	53.09 ± 33.55	53.09 ± 33.55	53.09 ± 33.55	36.30 ± 21.04	37.50 ± 19.75	70.36 ± 30.16	40.35 ± 16.18	36.66 ± 21.26	36.66 ± 21.26	36.66 ± 21.26	
PandaSlide-3	128.58 ± 46.00	127.28 ± 46.00	182.33 ± 80.03	97.12 ± 44.09	234.85 ± 93.90	80.81 \text{vtablashkash-1}	pm25.04	136.00 ± 58.69	57.72 ± 35.93	116.90 ± 61.86	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaSlide-4	211.02 ± 76.84	168.89 ± 278.76	147.22 ± 58.62	332.58 ± 88.90	96.40 ± 27.72	174.34 ± 50.19	174.34 ± 50.19	220.98 ± 78.94	53.10 ± 25.06	121.71 ± 76.46	161.35 ± 39.21	121.71 ± 76.46	161.35 ± 39.21
PandaSlide-5	241.12 ± 77.28	381.32 ± 86.85	428.32 ± 86.85	57.70 ± 77.39	53.35 ± 50.16	131.11 ± 77.23	131.49 ± 89.82	320.89 ± 84.43	236.75 ± 86.43	280.60 ± 93.04	295.19 ± 87.74	268.13 ± 84.74	220.81 ± 61.86
PandaSlide-6	348.48 ± 89.40	313.39 ± 87.00	222.29 ± 78.94	332.22 ± 85.15	175.77 ± 62.65	188.21 ± 70.04	188.21 ± 70.04	233.20 ± 61.01	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaSlide-7	331.08 ± 89.40	315.66 ± 89.36	220.20 ± 77.66	343.70 ± 82.83	233.20 ± 61.01	288.01 ± 86.91	233.20 ± 61.01	347.97 ± 92.67	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaSlide-8	406.25 ± 69.44	376.90 ± 75.28	224.63 ± 89.15	336.01 ± 93.85	373.00 ± 52.69	346.57 ± 88.34	346.57 ± 88.34	354.89 ± 94.26	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
train envs: PandaSlide-4													
PandaSlide-1	28.20 ± 13.28	32.01 ± 20.01	30.9 ± 15.55	21.80 ± 15.55	16.87 ± 8.45	19.11 ± 10.50	15.89 ± 10.70	24.97 ± 7.67	48.72 ± 24.93	86.49 ± 51.67	80.27 ± 51.55	24.74 ± 21.73	
PandaSlide-2	45.46 ± 33.49	53.09 ± 33.55	53.09 ± 33.55	53.09 ± 33.55	53.09 ± 33.55	37.50 ± 19.75	37.50 ± 19.75	70.36 ± 30.16	40.35 ± 16.18	36.66 ± 21.26	36.66 ± 21.26	36.66 ± 21.26	
PandaSlide-3	128.58 ± 46.00	127.28 ± 46.00	182.33 ± 80.03	97.12 ± 44.09	234.85 ± 93.90	80.81 \text{vtablashkash-1}	pm25.04	136.00 ± 58.69	57.72 ± 35.93	116.90 ± 61.86	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaSlide-4	211.02 ± 76.84	168.89 ± 278.76	147.22 ± 58.62	332.58 ± 88.90	96.40 ± 27.72	174.34 ± 50.19	174.34 ± 50.19	220.98 ± 78.94	53.10 ± 25.06	121.71 ± 76.46	161.35 ± 39.21	121.71 ± 76.46	161.35 ± 39.21
PandaSlide-5	241.12 ± 77.28	381.32 ± 86.85	428.32 ± 86.85	57.70 ± 77.39	53.35 ± 50.16	131.11 ± 77.23	131.49 ± 89.82	320.89 ± 84.43	236.75 ± 86.43	280.60 ± 93.04	295.19 ± 87.74	268.13 ± 84.74	220.81 ± 61.86
PandaSlide-6	348.48 ± 89.40	313.39 ± 87.00	222.29 ± 78.94	332.22 ± 85.15	175.77 ± 62.65	188.21 ± 70.04	188.21 ± 70.04	233.20 ± 61.01	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaSlide-7	331.08 ± 89.40	315.66 ± 89.36	220.20 ± 77.66	343.70 ± 82.83	233.20 ± 61.01	288.01 ± 86.91	233.20 ± 61.01	347.97 ± 92.67	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaSlide-8	406.25 ± 69.44	376.90 ± 75.28	224.63 ± 89.15	336.01 ± 93.85	373.00 ± 52.69	346.57 ± 88.34	346.57 ± 88.34	354.89 ± 94.26	98.50 ± 51.62	155.11 ± 78.13	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
train envs: PandaSlide-4													
PandaSlide-1	18.51 ± 11.72	38.19 ± 89.66	24.97 ± 12.12	38.84 ± 89.42	16.87 ± 8.45	19.11 ± 10.50	15.89 ± 10.70	24.97 ± 7.67	48.72 ± 24.93	86.49 ± 51.67	80.27 ± 51.55	24.74 ± 21.73	
PandaSlide-2	45.46 ± 33.49	53.09 ± 33.55	53.09 ± 33.55	53.09 ± 33.55	53.09 ± 33.55	38.57 ± 11.73	38.57 ± 11.73	70.36 ± 30.16	40.35 ± 16.18	36.66 ± 21.26	36.66 ± 21.26	36.66 ± 21.26	
PandaSlide-3	128.58 ± 46.00	127.28 ± 46.00	182.33 ± 80.03	97.12 ± 44.09	234.85 ± 93.90	80.81 \text{vtablashkash-1}	pm25.04	136.00 ± 58.69	57.72 ± 35.93	116.90 ± 61.86	192.15 ± 66.30	124.49 ± 84.71	119.99 ± 88.16
PandaSlide-4	211.02 ± 76.84	168.89 ± 278.76	147.22 ± 58.62	332.58 ± 88.90	96.40 ± 27.72								