
SCaR: Refining Skill Chaining for Long-Horizon Robotic Manipulation via Dual Regularization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Long-horizon robotic manipulation tasks typically involve a series of interrelated
2 sub-tasks spanning multiple execution stages. Skill chaining offers a feasible
3 solution for these tasks by pre-training the skills for each sub-task and linking
4 them sequentially. However, imperfections in skill learning or disturbances during
5 execution can lead to the accumulation of errors in skill chaining process, resulting
6 in execution failures. In this paper, we investigate how to achieve stable and
7 smooth skill chaining for long-horizon robotic manipulation tasks. Specifically,
8 we propose a novel skill chaining framework called **Skill Chaining via Dual**
9 **Regularization (SCaR)**. This framework applies dual regularization to sub-task skill
10 pre-training and fine-tuning, which not only enhances the *intra-skill dependencies*
11 within each sub-task skill but also reinforces the *inter-skill dependencies* between
12 sequential sub-task skills, thus ensuring smooth skill chaining and stable long-
13 horizon execution. We evaluate the SCaR framework on two representative long-
14 horizon robotic manipulation simulation benchmarks: IKEA furniture assembly
15 and kitchen organization. Additionally, we conduct real-world validation in desktop
16 robot pick-and-place tasks. The experimental results demonstrate that with the
17 support of SCaR, the robot performs long-horizon tasks with a higher success rate
18 than relevant baselines and is more robust to perturbations.

19 1 Introduction

20 Long-horizon robotic manipulation tasks are characterized by sequences of diverse and interdependent
21 sub-tasks, which makes it crucial to maintain the stability of multi-stage sequential execution. For
22 instance, in the robotic assembly of a stool (Fig. 1) involving two sub-tasks of leg installation, overall
23 success is evaluated based on both the sequential installation success and factors affecting the assembly
24 within environmental constraints. Although recent advances in deep reinforcement learning (RL) and
25 imitation learning (IL) show promise in training robots for such complex tasks [1, 2, 3, 4, 5, 6, 7],
26 managing long-horizon tasks with a scratch RL or IL policy remains challenging due to computational
27 demands, extensive exploration, and intricate step dependencies [8, 9]. Skill chaining, which involves
28 decomposing long-horizon tasks into smaller sub-tasks, pre-training skills for each, and executing
29 them sequentially, offers a practical solution [10, 11]. However, as shown in Fig. 1(a)(b), such
30 methods tend to fail when sub-task skills are insufficiently trained or unexpected states arise due
31 to disturbances, especially when applied to high-degree-of-freedom robots performing contact-rich,
32 long-horizon tasks. [12, 13, 14, 15, 16, 17].

33 In this paper, we argue that the coordination and enhancing of dependencies within and between sub-
34 task skills is necessary for stable and smooth skill chaining of long-horizon robotic manipulation [10].
35 For instance, as depicted in Fig. 1 (a)(b), the robot must consider following two points to ensure the
36 overall task is accomplished: 1) ensuring the gripper consistently grasps and installs the stool leg

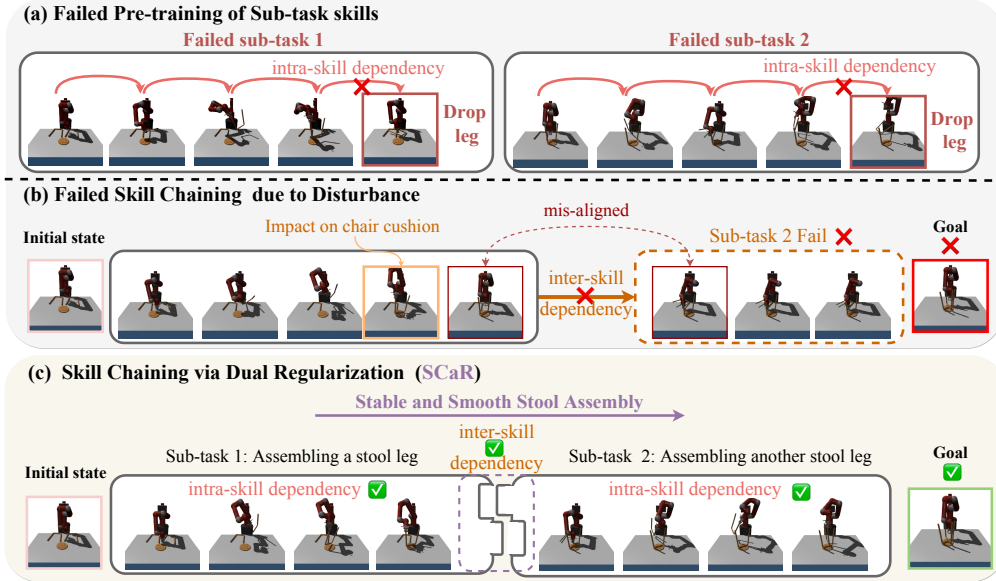


Figure 1: Illustration of the problem setting and the motivation of SCaR, using the example of a stool assembly task with two sub-tasks.

37 stably within each sub-task skill range. and 2) ensuring the terminal state of previous skill aligns with
 38 the initial state of next skill for smooth skill chaining. We define the above two points as *intra-skill*
 39 *dependencies* between sequential actions within each sub-task skill and *inter-skill dependencies*
 40 between sequential sub-task skills, respectively. In this context, we propose a novel robotic skill
 41 chaining framework, **Skill Chaining via Dual Regularization (SCaR)**. This framework enhances the
 42 aforementioned dependencies alternately through dual regularization during sub-task skill learning
 43 and chaining, aiming to provide stability for the execution of long-horizon robotic manipulation.

44 Specifically, in the pre-training phase of each sub-task skill, we propose the *adaptive sub-task skill*
 45 *learning* scheme, which employs a two-part policy learning objective that focuses on what sub-tasks
 46 the robot should perform (via RL) and how the robot should perform that task (via IL), and utilizes
 47 a novel adaptive equilibrium scheduling (AES) regularization to balance these two parts based on
 48 the robot’s learning progress. This process aims to reinforce the *intra-skill dependencies*, ensuring a
 49 coherent sequence of actions in each sub-task skill. Subsequently, *bi-directional adversarial learning*
 50 is introduced in the fine-tuning phase of SCaR for better chaining sequential sub-task skills. This
 51 mechanism uses bi-directional regularization to bring the terminal state of the current skill close to the
 52 initial state of its successor, and also to bring the initial state of the successor close to the terminal
 53 state of the current skill. This bi-directional alignment aims to reinforce robust *inter-skill dependencies*
 54 between sequential skills. Through the two innovative designs described, SCaR ensures coordination
 55 between the *intra-skill* and *inter-skill* dependencies, provides dual constraints for skill learning and
 56 skill chaining, as described in Fig. 1 (c), leading to a smooth skill chaining from the inside (**within the**
 57 **sub-task skills**) to the outside (**between sub-task skills**). Experimental results show that compared
 58 to scratch-training and skill chaining baselines, SCaR provides better task execution performance and
 59 stronger robustness to environmental perturbations in various long-horizon and contact-rich robotic
 60 manipulation simulation tasks. In addition, SCaR achieves higher task success rates in long-horizon
 61 real-robot pick-and-place tasks compared to previous skill chaining method.

62 The principal contributions of our work are delineated as follows: **1)** We propose a novel robotic skill
 63 chaining framework via dual regularization, SCaR, for smoothly executing long-horizon manipulation
 64 tasks. **2)** We introduce an adaptive sub-task skill learning scheme that acts as a regularization to
 65 enhance *intra-skill dependencies* between sequential actions within each sub-task skill. **3)** We develop
 66 a bi-directional adversarial learning mechanism that serves as a regularization for reinforcing *inter-*
 67 *skill dependencies* between sequential sub-task skills. **4)** In all eight simulated long-horizon robotic
 68 manipulation tasks, SCaR performs significantly better than scratch-training and skill chaining
 69 baselines. In addition, SCaR also shows better task completion performance compared to skill
 70 chaining baseline in real-robot long-horizon pick-and-place experiments. Video demonstrations are
 71 available at: <https://tinyurl.com/4333d6np>.

72 2 Related Work

73 2.1 Long-horizon Robotic Manipulation

74 Training robots from scratch for complex, long-horizon tasks using reinforcement learning (RL)
75 and imitation learning (IL) is challenging due to computational demands and distributional errors.
76 Solutions involve decomposing tasks into reusable sub-tasks [18]. Typically, such algorithms consist
77 of a set of sub-policies that can be obtained through various methods, such as unsupervised
78 exploration [19, 20, 21, 22, 23], learning from demonstrations [5, 6, 24, 25], and predefined mea-
79 sures [26, 27, 28, 29, 14]. Despite the merits of each of these approaches, they do not address well
80 the challenges of long-horizon robot manipulation in environments that are object-rich, contact-rich,
81 and characterized by multi-stage tasks [28, 29, 14]. Thus, even when pre-trained skills are provided,
82 ensuring a smooth connection between manipulation policies remains a formidable challenge.

83 2.2 Skill Chaining for Long-horizon Tasks

84 Previous skill chaining methods for long-horizon tasks mainly focus on updating each sub-task
85 policy to encompass the terminal state of the previous policy [11, 14, 30], implementing option
86 chains [11, 31, 32] to forge logical skill sequences, or utilizing modulated skills to facilitate smoother
87 transitions [33, 34, 35, 36, 14, 16]. However, these methods, while effective, often lead to a broad
88 range of skill start and end states, a challenge in complex robotic manipulation tasks. T-STAR [15] is
89 closely related to our work, addressing this by regularizing the learning process with a discriminator
90 to control the expansion of the terminal state space. However, it focuses only on uni-directional
91 dependencies between skills and ignores intra-skill dependencies within sub-task skills under long-
92 horizon goals. Sequential Dexterity [17] centers on dexterous hand manipulation, introducing an
93 optimization process to backpropagate long-term rewards across a policy chain. However, its scope
94 still primarily emphasizes strengthening the dependencies between sub-task skills. GSC [37] attempts
95 to solve skill chaining by employing diffusion models. It trains and chains primitive skills (pick,
96 place, push, pull) through a Transformer-based skill diffusion model. However, due to the use of
97 Transformer-based techniques, GSC requires high computational resources and cannot scale well to
98 task environments with object-rich and contact-rich conditions. Our method instead employs simple
99 and intuitive dual regularization constraints based on the lightweight policy network. By coordinating
100 the dependencies within and between skills, we achieve refinement within sub-task policies and
101 bi-directional alignment between them. This allows for stable skill chaining while also being scalable
102 to various long-horizon manipulation tasks.

103 3 Preliminaries

104 Among several related works on skill chaining, we consider a challenging yet practical problem
105 setting that *deals with long-horizon manipulation tasks through a combination of reinforcement*
106 *learning (RL) and imitation learning (IL)*. In each sub-task in the long-horizon task, we consider
107 robotic agents acting within a finite-horizon Markov Decision Process [38] $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma, d_{\mathcal{I}}, T)$,
108 where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{P}(s'|s, a)$ is the transition function, $r(s, a, s')$ is
109 the reward function, γ is the discount factor, $d_{\mathcal{I}}$ is the initial state distribution, and T is the episode
110 horizon of sub-task. We define a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maps states to actions and correspondingly
111 moves the robotic agent to a new state according to the transition probabilities. This sub-task policy is
112 trained to maximize the expected sum of discounted rewards $\mathbb{E}_{(s,a) \sim \pi} [\sum_{t=1}^T \gamma^t r(s_t, a_t, s_{t+1})]$. We
113 assume that each sub-task policy has an initial state set $\mathcal{I} \in \mathcal{S}$ and a terminal state set $\beta \in \mathcal{S}$, where
114 the initial set \mathcal{I} contains all the initial states that lead to the successful execution of the policy and
115 the terminal state set β contains all the final states of the successful execution. The environment
116 provides the environmental feedback for each step taken by the agent and success metrics for each
117 sub-task, derived from the terminal states of sub-task policy. For instance, as shown in Fig. 1(c),
118 the alignment of the back and legs of the stool triggers the connect action and the realization of the
119 sub-task goal, which indicates the successful completion of the sub-task. Additionally, we posit that
120 during each sub-task policy learning, the agent receives a set of pre-defined expert demonstrations,
121 $\mathbb{D}^E = \{\tau_1^E, \dots, \tau_N^E\}$, to facilitate the IL process. Here, N represents the number of episodes, and
122 each demonstration comprises a sequence of state-action pairs, $\tau^E = (s_1, a_1, \dots, s_{T-1}, a_{T-1}, s_T)$.

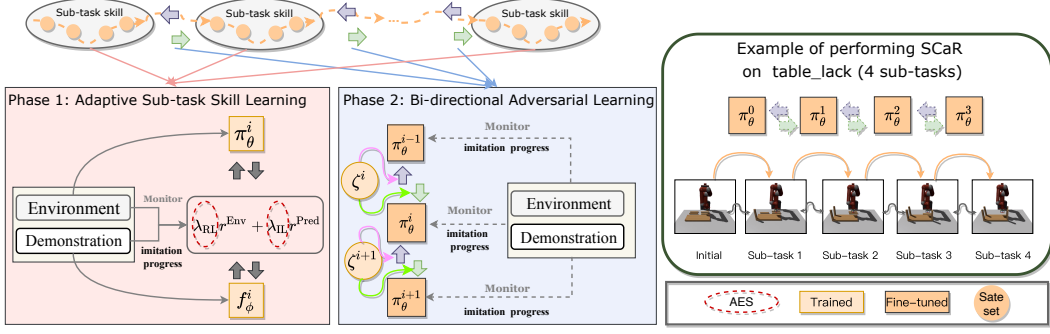


Figure 2: The Pipeline of Skill Chaining via Dual Regularization (SCaR). **(Left) Phase 1:** Sub-task skill pre-training () merges environmental feedback and expert guidance, using adaptive equilibrium scheduling (AES) regularization to balance learning, which enhances intra-skill dependencies within skills. **(Middle) Phase 2:** Bi-directional discriminators () coupled with AES to fine-tune pre-trained sub-task skills, as regularization for reinforcing inter-skill dependencies. **(Right) Evaluation:** Evaluation of SCaR on long-horizon manipulation.

123 4 Method

124 In Section 4.1, we present the pipeline of the **SCaR** framework. Sections 4.2 and 4.3 provide further
125 elaboration on the key design elements.

126 4.1 Overall Pipeline

127 As illustrated in Fig. 2, the **SCaR** framework has two phases: **(a) pre-training (adaptive sub-task
128 skill learning)** and **(b) fine-tuning (bi-directional adversarial learning)**. In the pre-training phase,
129 the agent co-learns sub-task skills by integrating environmental feedback and expert demonstrations.
130 In the fine-tuning phase, it refines these skills through bi-directional adversarial learning, enabling
131 sequential integration of sub-task skills. After fine-tuning, SCaR can smoothly chain sub-task skills
132 to complete long-horizon manipulation tasks. Specific modules and mechanisms for these phases are
133 detailed in Sections 4.2 and 4.3.

134 4.2 Adaptive Sub-task Skill Learning

135 **Weighted Reward Function** To learn sub-task skills better, we combine goal-conditional RL and
136 generative adversarial imitation learning (GAIL) [39], to pre-train skills that enable the agent to
137 perform challenging sub-tasks in a desired expert behavioral style [40, 15]. More specifically, we
138 consider the weighted reward function that is used to train each sub-task policy π_θ^i consists of two
139 components specifying: *what sub-task the agent should perform* - learning from environmental
140 feedback, and 2) *how the agent should perform that task* - learning from expert demonstrations:

$$r(s_t, a_t, s_{t+1}; \phi) = \lambda_{\text{RL}} r_i^{\text{Env}}(s_t, a_t, s_{t+1}, g) + \lambda_{\text{IL}} r_i^{\text{Pred}}(s_t, a_t; \phi). \quad (1)$$

141 As shown in Eq. 1, the first component is represented by a task-specific reward $r_i^{\text{Env}}(s_t, a_t, s_{t+1}, g)$,
142 which defines general objectives that the agent should satisfy to fulfill a given sub-task goal g for
143 current MDP \mathcal{M} (e.g. assembling a stool leg). The second component is represented through a learned
144 task-agnostic predict-reward $r_i^{\text{Pred}}(s_t, a_t; \phi)$, which specifies manipulation details of the behaviors
145 that the agent should adopt when performing the sub-task (e.g., the expert way to grab a stool leg
146 and attach it), and $r_i^{\text{Pred}}(s_t, a_t; \phi)$ is the predicted reward by a least-square GAIL discriminator
147 f_ϕ^i [41, 40, 15], which is more stable than the standard GAIL objective using the sigmoid cross-
148 entropy loss function. Therefore, the predicted reward is:

$$r_i^{\text{Pred}}(s_t, a_t; \phi) = \max[0, 1 - 0.25 \cdot [f_\phi^i(s_t, a_t) - 1]^2]. \quad (2)$$

149 We adopt the training objective of the least-squares GAIL discriminator with a gradient penalty
150 term [42, 43], This penalty term mitigates the instability of the training dynamics due to the interplay
151 between the discriminator and the policy [40], as follows:

$$\operatorname{argmin}_{f_\phi^i} \mathbb{E}_{(s) \sim \mathbb{D}^E} [(f_\phi^i(s) - 1)^2] + \mathbb{E}_{(s) \sim \pi_\theta^i} [(f_\phi^i(s) + 1)^2] + \frac{\eta^{\text{GP}}}{2} \mathbb{E}_{(s) \sim \mathbb{D}^E} [\|\nabla_s f_\phi^i(s)\|^2], \quad (3)$$

152 where η^{SP} is a manually-specified coefficient. The scales of r^{Env} and r^{Pred} in previous related
 153 works are set by fixed weights and linearly combined into the final reward function [40, 15]. This
 154 could lead to the agent rigidly imitating experts and curbing self-exploration, finding it difficult to
 155 adjust intra-skill dependencies and adapt to dynamic task perturbations. We propose a principle to
 156 counter this: ***If the agent fails to imitate the expert’s demonstration well, it should shift focus to
 157 self-learning from the environment. Conversely, effective imitation should continue, focusing on
 158 the expert to mitigate low sample efficiency in reinforcement learning.*** Accordingly, we extend
 159 the automatic discount scheduling (ADS) solution [9] to our problem setting, and propose adaptive
 160 equilibrium scheduling (AES) to regularize the scales of r^{Env} and r^{Pred} in sub-task skill learning for
 161 adaptive scheduling the focus of reinforcement and imitation learning, as shown in Fig. 3.

162 **Adaptive Equilibrium Scheduling (AES) Regularization** Specifically, AES balances the scales of
 163 r^{Env} and r^{Pred} during the learning process of each skill through adaptive scheduling of λ_{RL} and λ_{IL} ,
 164 according to how well the agent imitates the expert’s demonstration. To capture the agent’s imitation
 165 progress, AES refers to the solution in ADS [9] and uses the imitation identifier Φ to continuously
 166 monitor whether the agent is imitating the expert demonstration well enough.

167 At the beginning of training, the agent is assigned two initial
 168 balance factors $\lambda_{\text{RL}} = \alpha, \lambda_{\text{IL}} = 1 - \alpha$, where base exponent
 169 $\alpha \in [0, 1]$. We set $\alpha = 0.5$ in the experiments and the agent is
 170 assigned two identical balance factors $\lambda_{\text{RL}} = \lambda_{\text{IL}} = 0.5^1$, indicat-
 171 ing that at the beginning of learning, the agent imitates the expert’s
 172 behavior with the same weight as the behavior of environment
 173 exploration according to the task goal. As training progresses, the
 174 imitation progress recognizer Φ is queried periodically to monitor
 175 the progress of the agent’s imitation of the expert’s behavior. Φ
 176 receives the agent’s collected trajectories and infers the agent’s
 177 current imitation progress $p \in [0, T]$, where p in an integer and T
 178 is the step of the entire episode.

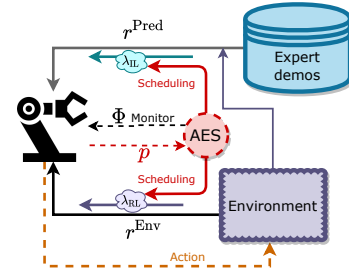


Figure 3: AES regularization for sub-task skill learning.

179 The construction of Φ , with reference to ADS, first requires
 180 the construction of a sequence $\mathbf{Q}(q_1, \dots, q_T)$, where $q_i =$
 181 $\text{argmin}_j c(s_i, s_j^E)$ is the index of the nearest neighbor of s_i in τ^E , c is the cosine similarity. The
 182 progress alignment between τ and τ_j^E is measured as the length of the longest increasing subsequence
 183 (LIS) in \mathbf{Q} , denoted as $LIS(\tau, \tau^E)$. Specifically, the agent’s imitation progress p is increased by 1 if
 184 the following inequality holds:

$$\max_{\hat{\tau}^E \in \mathbb{D}^E} LIS(\tau_{1:p+1}, \hat{\tau}_{1:p+1}^E) \geq \rho \times \min_{\hat{\tau}^E, \hat{\tau}^E \in \mathbb{D}^E} LIS(\hat{\tau}_{1:p+1}^E, \hat{\tau}_{1:p+1}^E), \quad (4)$$

185 where $\hat{\tau}^E \neq \tau^E$, the subscript $1 : p + 1$ denotes the first $p + 1$ steps of the trajectory, and $\rho \in [0, 1]$
 186 controls the strictness of the imitation progress monitoring. This suggests that the similarity of the
 187 agent trajectory to its best matching expert trajectory at time step $p + 1$ exceeds the minimal similarity
 188 criterion within the expert demonstration. See Appendix B for detailed explanation of AES.

189 After obtaining the current imitation progress p of the agent, AES then adopts a mapping function
 190 $\varphi_\lambda(p)$ to schedule the two new balance discount factors λ_{RL} and λ_{IL} . Straightforward idea of setting
 191 $\varphi_\lambda(p)$ is that ***If p is larger and reaches a certain threshold, i.e., the agent is able to imitate
 192 the expert behavior well, then the more the agent tends to imitate the expert’s behavior in
 193 subsequent training, and vice versa.*** Therefore, we set the threshold as $\frac{T}{2}$. If $p \in [0, \frac{T}{2})$, we

194 propose $\varphi_\lambda(p) = 1 - e^{-\frac{p}{k}}$; if $p \in [\frac{T}{2}, T)$, we propose $\varphi_\lambda(p) = e^{-\frac{p - \frac{T}{2}}{k}}$, where k is used to
 195 flatten the curve of the mapping function. Then λ_{RL} and λ_{IL} are scheduled to be :

$$\begin{cases} \lambda_{\text{RL}} = \alpha^{\varphi_\lambda(p)}, \lambda_{\text{IL}} = 1 - \alpha^{\varphi_\lambda(p)}. & \text{if } p \in [0, \frac{T}{2}) \\ \lambda_{\text{IL}} = \alpha^{\varphi_\lambda(p)}, \lambda_{\text{RL}} = 1 - \alpha^{\varphi_\lambda(p)}. & \text{if } p \in [\frac{T}{2}, T) \end{cases} \quad (5)$$

196 Consequently, the RL and IL components of sub-task skill learning can be adaptively scheduled and
 197 regularized through AES, effectively enhancing *intra-skill dependencies* between sequential actions.
 198 The pseudo-code of adaptive sub-task skill learning is outlined in Algorithm 1 in Appendix A.1.

¹We further explore what effect different α would have in the Ablation Experiments.

199 **4.3 Bi-directional Adversarial Learning for Skill Chaining**

200 Executing pre-trained sub-task skills sequentially without considering inter-skill dependencies may
 201 lead to failure. To address this, we propose bi-directional adversarial learning to further refine and
 202 better integrate sequential sub-task skills. The pseudo-code of bi-directional adversarial learning is
 203 outlined in Algorithm 2 in Appendix A.2.

204 **Bi-directional Regularization** In contrast to previous uni-directional regularization schemes that
 205 only augment the initial state set \mathcal{I}_i or regularize the terminal state set β_i [12, 15], we impose the
 206 *bi-directional constraints* ($\mathcal{C}_1, \mathcal{C}_2$) on inter-skill dependencies, facilitating smooth skill chaining, as
 207 shown in Fig 4. With the bi-directional constraint, we implement the bi-directional adversarial
 208 learning, centered on the joint training of a *bi-directional discriminator*, denoted by ζ_ω^i , which is
 209 adept at distinguishing between the terminal state set of the preceding policy and the initial state set
 210 of the subsequent policy. The bi-directional constraints $\mathcal{C}_1, \mathcal{C}_2$ are defined as Eq. 10:

$$\begin{aligned} \text{next initial} \rightarrow \text{previous terminal: } \mathcal{C}_1 &= \mathbb{E}_{s_{\mathcal{I}} \sim \mathcal{I}_i} [\zeta_\omega^i(s_{\mathcal{I}}) - 1]^2 + \mathbb{E}_{s_T \sim \beta_{i-1}} [\zeta_\omega^i(s_T)]^2 \\ \text{previous terminal} \rightarrow \text{next initial: } \mathcal{C}_2 &= \mathbb{E}_{s_T \sim \beta_i} [\zeta_\omega^i(s_T) - 1]^2 + \mathbb{E}_{s_{\mathcal{I}} \sim \mathcal{I}_{i+1}} [\zeta_\omega^i(s_{\mathcal{I}})]^2 \end{aligned} \quad (6)$$

211 ζ_ω^i is trained for each policy to minimize the objective func-
 212 tion²: $\mathcal{L}_i(\omega) = \frac{1}{2}\mathcal{C}_1 + \frac{1}{2}\mathcal{C}_2$. Guided by ζ_ω^i , the bi-directional
 213 adversarial learning not only steers the terminal state set of the
 214 current policy towards the initial state set of the subsequent policy,
 215 but also ensures alignment of the initial state set of the subse-
 216 quent policy with the terminal state set of current policy. This dual
 217 alignment establishes a balanced mapping between the initial and
 218 terminal states of sequential skills to reinforce inter-skill depen-
 219 dencies, ensure consistency and stability in multi-stage tasks, and
 220 guarantee smooth transitions between sequential skills. Accord-
 221 ingly, the *bi-directional regularization* can be added to the overall
 222 objective function of policy learning in the form of the following
 223 reward term: $r_i^{\text{Bi}}(s; \omega) = \mathbb{1}_{s \in \beta_i} \zeta^{i+1}(s) + \mathbb{1}_{s \in \mathcal{I}_i} \zeta^{i-1}(s)$.

224 **Overall Objective Function** So far, the objective function via
 225 dual regularization, i.e., AES regularization and bi-directional
 226 regularization, to pre-train, fine-tune and chain sub-task skills can be rewritten as a weighted sum of
 227 the individual reward terms:

$$r_i(s_t, a_t, s_{t+1}; \phi) = \underbrace{\lambda_{\text{RL}} r_i^{\text{Env}}(s_t, a_t, s_{t+1}, g)}_{\text{AES regularization}} + \underbrace{\lambda_{\text{IL}} r_i^{\text{Pred}}(s_t, a_t; \phi)}_{\text{bi-directional regularization}} + \underbrace{\lambda_{\text{Bi}} r_i^{\text{Bi}}(s_{t+1}; \omega)}_{\text{bi-directional regularization}}, \quad (7)$$

228 where λ_{Re} is the weighting factor of the bi-directional regularization. The objective function features
 229 AES regularization and bi-directional regularization to enhance intra- and inter-skill dependencies.
 230 It enables the agent to adaptively pre-train skills that can solve different sub-tasks well through
 231 environmental feedback and expert guidance, and further fine-tune them through the bi-directional
 232 discriminator to achieve dual alignment between sequential skills. At the same time, the fine-tuned
 233 sub-task skills help to collect terminal and initial states to refine the bi-directional discriminator. This
 234 iterative process ensures smooth long-horizon task skill chaining.

235 **5 Experiments**

236 **5.1 Experiment Setup**

237 We conduct simulation experiments on six IKEA furniture assembly tasks and two kitchen organi-
 238 zation tasks, and also perform long-horizon pick-and-place experiments on the real Sagittarius K1
 239 robot. Please refer to the Appendix for more detailed simulation experiment setup (Appendix G),
 240 network architecture (Appendix H), training details (Appendix I), more quantitative (Appendix D) and
 241 qualitative results (Appendix E) of the simulation tasks, and the real-robot experiments (Appendix F).

242 **Furniture Assembly** We conduct experiments in six IKEA furniture assembly tasks in [44]:
 243 *chair_agne*, *chair_bernhard*, *chair_ingolf*, *toy_table*, *table_lack* and *table_bjorkudden*.

²We explore the impact of different scales of \mathcal{C}_1 and \mathcal{C}_2 in Appendix D.3

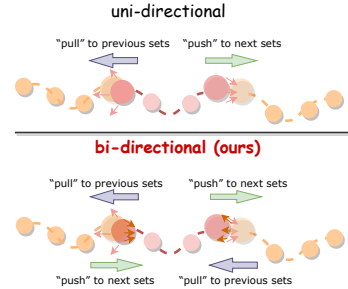


Figure 4: Bi-directional regularization for sub-task skill chaining.

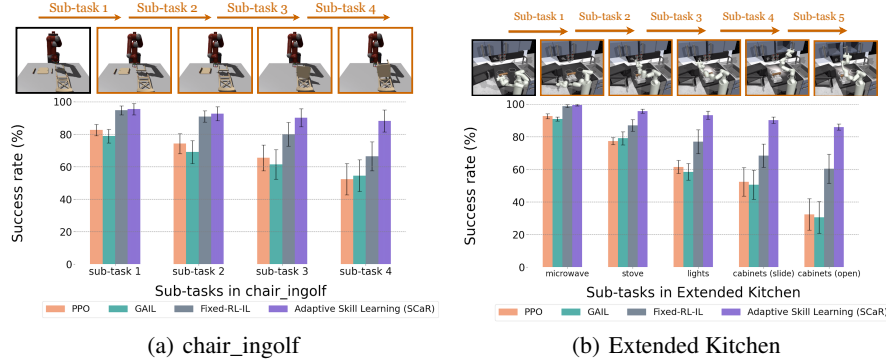


Figure 5: Evaluation Performance of Sub-task Skill Learning. Best viewed zoomed.

244 1) *chair_agne*: Two stool legs need to be picked up and aligned with the cross notches on the stool
 245 back. 2) *chair_bernhard*: The two chair supports need to be taken and aligned with the slots at the
 246 bottom of the chair surface. 3) *chair_ingolf*: Two chair supports and front legs need to be attached to
 247 the chair seat, which must then be secured to the chair back while avoiding collision with each other.
 248 4) *table_lack*: The four table legs need to be picked up and aligned with the corners of the tabletop. 5)
 249 *toy_table*: The four table legs need to be picked up and aimed and inserted with the four notches on
 250 the table back. 6) *table_dockstra*: After supporting the two bases with table leg, the table top needs to
 251 be mounted while preventing collision. For each assembly task, we define the assembly of individual
 252 parts as sub-tasks. We collect 200 demonstrations per sub-task using a procedural assembly policy
 253 for imitation learning. Each demonstration consists of 150 steps.

254 **Kitchen Organization** We use the Franka Kitchen tasks in D4RL [45] and collect 200 demonstra-
 255 tions per sub-task for imitation learning. Specifically, we refer to the kitchen task in [46] and further
 256 extend the task sequence: in the **Kitchen task**, the 7-DoF Franka Emika Panda arm needs to perform
 257 4 sequential sub-tasks, namely *Turn on the microwave - Move the kettle - Turn on the stove - Turn on*
 258 *the light*. In the **Extended Kitchen task**, the robot needs to perform 5 sequential sub-tasks: *Turn*
 259 *on the microwave - Turn on the stove - Turn on the light - Slide the cabinet to the right - Open the*
 260 *cabinet*, in which the sub-tasks have a lower probability of switching and is more challenging.

261 **Baselines** We compare SCaR with the following two types of baselines:

262 **Scratch Training:** 1) **PPO** is a model-free RL algorithm [47] that utilizes environmental rewards
 263 to learn tasks from scratch. 2) **GAIL** [39] is an adversarial imitation learning method to learn tasks
 264 from scratch, with a trained discriminator for distinguishing state-action distributions of experts and
 265 agents. 3) **Fixed-RL-IL** [40] uses fixed-weight environmental rewards and GAIL rewards to train
 266 policies from scratch. 4) **SkiMo** [46] is a model-based hierarchical RL approach that learns dynamic
 267 skill models for predicting outcomes in downstream tasks, which is used to test if modularly skill
 268 chaining method can surpass model-based scratch-training method on long-horizon tasks.

269 **Skill Chaining:** 1) **Policy Sequencing** [12] focuses on sequentially expanding the initial sets
 270 in skill chaining. 2) **T-STAR** [15] incorporates a discriminator to uni-directionally regularize the
 271 terminal states of sub-skills in a skill chaining. 3) **SCaR w/o Bi** reference to T-STAR during the
 272 fine-tuning phase, only uni-directional regularization of the terminal state set is performed to verify
 273 the validity of the proposed bi-directional regularization. 4) **SCaR w/o AES** fixes the scales of the
 274 two reward terms at 0.5 at all times to verify the effectiveness of the proposed AES regularization.

275 5.2 Quantitative Results

276 **Sub-task Skill Learning Performance** First, we evaluate the proposed adaptive sub-task skill
 277 learning scheme in the sub-tasks of furniture assembly and kitchen organization. Specifically, we
 278 treat each sub-task as a separate task for policy learning and take the success rate of the trained
 279 policy tested in the reset sub-task as the criterion. All methods are trained in each sub-task with 5
 280 random seeds, 150 million environment steps, and evaluated with the average success rate over 100
 281 testing episodes. As shown in the Fig. 5, in *chair_ingolf* and Extended Kitchen tasks, even with the
 282 increase of objects in the environment and the increase of unpredictable perturbations, our proposed
 283 adaptive skill learning learns good sub-task skills and consistently maintains a task success rate of
 284 more than 85% in all stages of the sub-task. In contrast, the PPO (only RL rewards), GAIL (only IL

Table 1: Long-horizon tasks execution performance (varies by sub-task completion progress): **tasks* with 2 sub-tasks progress by 0.5 per sub-task, **tasks* with 4 sub-tasks by 0.25, **tasks* with 5 sub-tasks by 0.2, and *table_dockstra* with 3 sub-tasks by 0.3, where 0.9 indicates completion of all tasks. Best viewed zoomed.

Method	Furniture Assembly						Kitchen Organization			
	<i>chair_agne</i>	<i>chair_bernhard</i>	<i>chair_ingolf</i>	<i>table_lack</i>	<i>toy_table</i>	<i>table_dockstra</i>	All	Kitchen	E-Kitchen	All
PPO (Scratch RL)	0.54±0.18	0.42±0.12	0.14±0.03	0.09±0.01	0.00±0.00	0.31±0.12	0.25±0.15	0.13±0.05	0.03±0.00	0.08±0.04
GAIL (Scratch IL)	0.31±0.05	0.23±0.02	0.00±0.00	0.00±0.00	0.00±0.00	0.21±0.04	0.12±0.09	0.00±0.00	0.00±0.00	0.00±0.00
Fixed-RL-IL	0.68±0.12	0.53±0.07	0.22±0.08	0.21±0.11	0.13±0.02	0.43±0.07	0.37±0.15	0.33±0.06	0.18±0.02	0.26±0.06
SkiMo	0.75±0.09	0.62±0.05	0.47±0.03	0.58±0.14	0.34±0.06	0.62±0.11	0.56±0.11	0.57±0.08	0.21±0.04	0.39±0.13
Policy Sequencing	0.89±0.08	0.82±0.09	0.77±0.12	0.63±0.28	0.45±0.18	0.61±0.14	0.70±0.16	0.53±0.11	0.36±0.09	0.44±0.09
T-STAR	0.92±0.02	0.90±0.04	0.89±0.04	0.90±0.07	0.71±0.21	0.77±0.09	0.85±0.09	0.68±0.13	0.48±0.08	0.58±0.10
SCaR w/o Bi	0.93±0.04	0.92±0.02	0.91±0.01	0.93±0.02	0.80±0.10	0.79±0.02	0.88±0.05	0.75±0.08	0.57±0.14	0.66±0.09
SCaR w/o AES	0.95±0.03	0.94 ±0.03	0.93±0.02	0.95±0.04	0.85±0.06	0.80±0.03	0.91±0.05	0.77±0.07	0.61±0.13	0.74±0.05
SCaR (Ours)	0.98 ±0.02	0.96 ±0.04	0.95 ±0.03	0.97 ±0.03	0.92 ±0.05	0.88 ±0.02	0.94 ±0.03 (12% ↑)	0.84 ±0.16	0.73 ±0.17	0.78 ±0.12 (18% ↑)

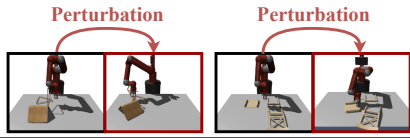
rewards), and Fixed-RL-IL (fixed RL and IL reward weights) baselines fail to maintain good sub-task success rates as the number of sub-task stages increases. This result well validates that our proposed adaptive weighted reward function based on AES regularization enhances *intra-skill dependencies* for multi-stage sub-task learning and brings effectiveness and stability.

Long-horizon Execution Performance We then demonstrate the performance of SCaR in performing 8 long-horizon tasks in IKEA furniture assembly and kitchen organization. Table 1 shows the mean and standard deviation for these 8 tasks across 200 testing episodes with 5 different seeds. The PPO and GAIL baselines show minimal success on tasks with 4 and 5 sub-tasks, indicating the difficulty of learning complex multi-stage tasks solely from reward signals or expert demonstrations. The fixed RL-IL baseline, although improved compared to PPO and GAIL, mostly completed only one sub-task, which highlights the limitations of using fixed RL and IL reward weights in long-horizon tasks. While SkiMo achieves better success rates than model-free methods by building dynamic skill models, its performance remains inconsistent on long-horizon tasks due to its scratch learning nature. The performance of these scratch baselines demonstrates the importance of effective staged sub-task learning for long-horizon tasks. The results in Table 1 further highlight the superiority of the SCaR framework. By reinforcing *intra- and inter-skill dependencies*, task success rates are considerably higher than previous skill chaining approaches such as Policy Sequencing and T-STAR, which primarily address uni-directional inter-skill dependencies. Compared to T-STAR, SCaR increases average success rates by more than 12% on six furniture assembly tasks and 18% on two kitchen organization tasks.³

5.3 Robustness to Perturbations

Perturbation tests are conducted to evaluate the robustness of skill chaining for two furniture assembly tasks. As shown in the top figure of Table 2, for the *chair_bernhard* task, the perturbation involves applying external joint torque to the robotic arm, moving the chair back before assembling the second support. For the *chair_ingolf* task, the perturbation is applied by exerting external torque on the robotic arms, causing them to move slightly before mounting the assembled chair seat to the chair back. The results in Table 2 highlight the detrimental impact of environmental perturbations on the success rates of baseline methods during the execution of multiple sub-task skills. Methods like Policy Sequencing and T-STAR, which focus solely on inter-skill dependencies through uni-directional regularization, struggle to complete tasks after perturbations. In contrast, SCaR, demonstrates more robust performance even under unseen perturbations. These results further support the advantages of our proposed *dual regularization* for stable skill chaining on long-horizon manipulation tasks.

Table 2: Comparison of the robustness of skill chaining in perturbed environments.



Method	<i>chair_bernhard</i>		<i>chair_ingolf</i>	
	No Perturb	Perturb	No Perturb	Perturb
Policy Sequencing	0.82±0.09	0.51±0.04	0.77±0.12	0.50±0.10
T-STAR	0.90±0.04	0.60±0.08	0.89±0.04	0.59±0.04
SCaR w/o Bi	0.92±0.02	0.65±0.11	0.91±0.01	0.63±0.05
SCaR w/o AES	0.94±0.03	0.74±0.09	0.93±0.02	0.71±0.07
SCaR (Ours)	0.96 ±0.04	0.85 ±0.11	0.95 ±0.03	0.80 ±0.13

³The overall increase is somewhat modest due to averaging the success rates of the 2, 3, and 4 sub-tasks and the 4 and 5 sub-tasks, respectively.

325 **5.4 Ablations and Analysis**

326 We perform ablation studies to explore the important factors that affect the performance of SCaR.

327 **Modular Ablation** We investigate how the adaptive sub-task skill learning and bi-directional
 328 adversarial learning impact skill chaining through SCaR w/o Bi and SCaR w/o AES. As shown in
 329 Table 1, without bi-directional regularization, SCaR w/o Bi experiences significant performance
 330 drops in tasks with more than two sub-tasks but still outperforms T-STAR. This is because SCaR
 331 w/o Bi maintains the adaptive scheduling of AES during sub-task skill learning, underscoring the
 332 importance of focusing on the *intra-skill dependencies* between successive actions. Similarly, the
 333 absence of AES regularization reduces SCaR w/o AES’s performance, though it still maintains stable
 334 outcomes. This underscores the importance of reinforcing *inter-skill dependencies* on long-horizon
 335 tasks and reaffirms the contribution of bi-directional regularization. As shown in Table 2, SCaR w/o
 336 Bi, though slightly more robust than T-STAR due to the presence of AES, still faces challenges in
 337 adapting to perturbations and maintaining stable skill chaining because of its uni-directional fine-
 338 tuning limitations. SCaR w/o AES manages to maintain a certain level of performance stability under
 339 perturbations, thanks to bi-directional regularization, which ensures the bi-directional alignment of
 340 initial and terminal states between skills. The results show that the pre-trained skills via AES exhibit
 341 enhanced *intra-skill dependencies* within sub-tasks, and bi-directional regularization ensures stable
 342 long-horizon execution, even in the presence of perturbations, by reinforcing *inter-skill dependencies*.
 343

344 **Parametric Ablation** We further in-
 345 vestigate the impact of different scales
 346 of RL and IL reward terms, as well
 347 as the size of expert demonstration
 348 datasets. The effect of varying the
 349 base exponent α on task success rates
 350 is tested across four tasks: *chair_agne*,
 351 *chair_ingolf*, *table_dockstra*, and *extend_kitchen*.
 352 As depicted in Fig. 6(a),
 353 SCaR achieves the highest success
 354 rates in all four tasks when $\alpha = 0.5$,
 355 indicating a balance between RL and IL at the beginning of learning. When α becomes smaller,
 356 emphasizing IL at the start, performance decreases more steeply. Conversely, as α becomes larger,
 357 giving more weight to RL, performance also declines but at a slower rate. We also evaluate the impact
 358 of different sizes of expert datasets on three skill chaining methods: Policy Sequencing, T-STAR,
 359 and SCaR, specifically in the *chair_ingolf* task. We vary the overall task expert data size from 80,
 360 120, 200, 400, 600, to 800 demos. As shown in Fig. 6(b), the results indicate significant performance
 361 improvement when increasing the dataset size from 400 to 800 demos, while the improvement is
 362 less pronounced when going from 80 to 120 demos. This demonstrates the importance of the demo
 363 dataset size in the effectiveness of data-driven approaches like skill chaining.

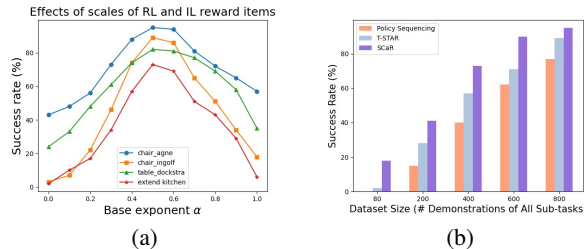


Figure 6: Ablation experiments. Best viewed zoomed.

364 **6 Discussion**

365 **Limitation and future directions** Limitations of our work are that the sub-task division of the
 366 long-horizon task is predefined and does not involve visual and semantic processing of objects.
 367 Scaling up our framework to address longer-horizon visual manipulation tasks is a direction we aim
 368 to investigate in future work. For instance, incorporating a more scalable architecture [48] along
 369 with large-scale pre-training on large datasets [49, 50] would be an interesting direction. Another
 370 compelling direction is applying our framework to actual robotic furniture assembly tasks, beyond
 371 just staged robotic pick-and-place tasks. Building a real-world deployment environment for furniture
 372 assembly and being able to guarantee full insertion of each furniture module are huge challenges.

373 **Conclusion** In this paper, we introduce SCaR, a novel skill chaining framework that ensures smooth
 374 and stable execution of long-horizon robotic manipulation tasks via dual regularization within and
 375 between sub-task skills. Extensive experiments demonstrate that the SCaR framework achieves better
 376 task success rates than the baseline methods in both simulated and real-robot manipulation tasks,
 377 while being robust against perturbations. We hope this work will inspire future research to further
 378 explore the potential of skill chaining for long-horizon robotic manipulation.

379 **References**

- 380 [1] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep
381 visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- 382 [2] Francisco Suárez-Ruiz and Quang-Cuong Pham. A framework for fine robotic assembly. In
383 *2016 IEEE international conference on robotics and automation (ICRA)*, pages 421–426. IEEE,
384 2016.
- 385 [3] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
386 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement
387 learning and demonstrations. *Robotics: Science and Systems XIV*, 2018.
- 388 [4] Divye Jain, Andrew Li, Shivam Singhal, Aravind Rajeswaran, Vikash Kumar, and Emanuel
389 Todorov. Learning deep visuomotor policies for dexterous hand manipulation. In *2019 interna-*
390 *tional conference on robotics and automation (ICRA)*, pages 3636–3643. IEEE, 2019.
- 391 [5] George Konidaris, Scott Kuindersma, Roderic Grupen, and Andrew Barto. Robot learning
392 from demonstration by constructing skill trees. *The International Journal of Robotics Research*,
393 31(3):360–375, 2012.
- 394 [6] Thomas Kipf, Yujia Li, Hanjun Dai, Vinicius Zambaldi, Alvaro Sanchez-Gonzalez, Edward
395 Grefenstette, Pushmeet Kohli, and Peter Battaglia. Compile: Compositional imitation learning
396 and execution. In *International Conference on Machine Learning*, pages 3418–3428. PMLR,
397 2019.
- 398 [7] Chen Wang, Linxi Fan, Jiankai Sun, Ruohan Zhang, Li Fei-Fei, Danfei Xu, Yuke Zhu, and
399 Anima Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play.
400 *arXiv preprint arXiv:2302.12422*, 2023.
- 401 [8] Marcin Andrychowicz, Anton Raichuk, Piotr Stańczyk, Manu Orsini, Sertan Girgin, Raphaël
402 Marinier, Leonard Hussenot, Matthieu Geist, Olivier Pietquin, Marcin Michalski, et al. What
403 matters for on-policy deep actor-critic methods? a large-scale study. In *International conference*
404 *on learning representations*, 2020.
- 405 [9] Yuyang Liu, Weijun Dong, Yingdong Hu, Chuan Wen, Zhao-Heng Yin, Chongjie Zhang, and
406 Yang Gao. Imitation learning from observation with automatic discount scheduling. *arXiv*
407 *preprint arXiv:2310.07433*, 2023.
- 408 [10] Caelan Reed Garrett, Rohan Chitnis, Rachel Holladay, Beomjoon Kim, Tom Silver, Leslie Pack
409 Kaelbling, and Tomás Lozano-Pérez. Integrated task and motion planning. *Annual review of*
410 *control, robotics, and autonomous systems*, 4:265–293, 2021.
- 411 [11] George Konidaris and Andrew Barto. Skill discovery in continuous reinforcement learning
412 domains using skill chaining. *Advances in neural information processing systems*, 22, 2009.
- 413 [12] Alexander Clegg, Wenhao Yu, Jie Tan, C Karen Liu, and Greg Turk. Learning to dress:
414 Synthesizing human dressing motion via deep reinforcement learning. *ACM Transactions on*
415 *Graphics (TOG)*, 37(6):1–10, 2018.
- 416 [13] Youngwoon Lee, Shao-Hua Sun, Sriram Somasundaram, Edward S Hu, and Joseph J Lim.
417 Composing complex skills by learning transition policies. In *International Conference on*
418 *Learning Representations*, 2018.
- 419 [14] Youngwoon Lee, Jingyun Yang, and Joseph J Lim. Learning to coordinate manipulation skills
420 via skill behavior diversification. In *International conference on learning representations*, 2019.
- 421 [15] Youngwoon Lee, Joseph J Lim, Anima Anandkumar, and Yuke Zhu. Adversarial skill chaining
422 for long-horizon robot manipulation via terminal state regularization. In *Conference on Robot*
423 *Learning (CoRL 2022)*, pages 406–416. PMLR, 2022.
- 424 [16] Jiayuan Gu, Devendra Singh Chaplot, Hao Su, and Jitendra Malik. Multi-skill mobile ma-
425 nipulation for object rearrangement. In *The Eleventh International Conference on Learning*
426 *Representations*, 2022.

- 427 [17] Yuanpei Chen, Chen Wang, Li Fei-Fei, and Karen Liu. Sequential dexterity: Chaining dexterous
428 policies for long-horizon manipulation. In *Conference on Robot Learning*, pages 3809–3829.
429 PMLR, 2023.
- 430 [18] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A
431 framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-
432 2):181–211, 1999.
- 433 [19] Jürgen Schmidhuber. *Towards compositional learning with dynamic neural networks*. Inst. für
434 Informatik, 1990.
- 435 [20] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings*
436 *of the AAAI conference on artificial intelligence*, volume 31, 2017.
- 437 [21] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical
438 reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- 439 [22] Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Learning multi-level hierar-
440 chies with hindsight. *arXiv preprint arXiv:1712.00948*, 2017.
- 441 [23] Visak CV Kumar, Sehoon Ha, and C Karen Liu. Expanding motor skills using relay networks.
442 In *Conference on Robot Learning*, pages 744–756. PMLR, 2018.
- 443 [24] Yuchen Lu, Yikang Shen, Siyuan Zhou, Aaron Courville, Joshua B Tenenbaum, and Chuang
444 Gan. Learning task decomposition with ordered memory policy network. In *International*
445 *Conference on Learning Representations*, 2020.
- 446 [25] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel
447 Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. Do as i can, not as i say: Grounding
448 language in robotic affordances. In *Conference on Robot Learning*, pages 287–318. PMLR,
449 2023.
- 450 [26] Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep
451 reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in*
452 *neural information processing systems*, 29, 2016.
- 453 [27] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization
454 with multi-task deep reinforcement learning. In *International Conference on Machine Learning*,
455 pages 2661–2670. PMLR, 2017.
- 456 [28] Josh Merel, Arun Ahuja, Vu Pham, Saran Tunyasuvunakool, Siqi Liu, Dhruva Tirumala, Nicolas
457 Heess, and Greg Wayne. Hierarchical visuomotor control of humanoids. In *International*
458 *Conference on Learning Representations*, 2018.
- 459 [29] Chen Wang, Danfei Xu, and Li Fei-Fei. Generalizable task planning through representation
460 pretraining. *IEEE Robotics and Automation Letters*, 7(3):8299–8306, 2022.
- 461 [30] Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. Mpc: Learning
462 composable hierarchical control with multiplicative compositional policies. *Advances in Neural*
463 *Information Processing Systems*, 32, 2019.
- 464 [31] Akhil Bagaria and George Konidaris. Option discovery using deep skill chaining. In *Internat-*
465 *ional Conference on Learning Representations*, 2019.
- 466 [32] Zixuan Chen, Ze Ji, Shuyang Liu, Jing Huo, Yiyu Chen, and Yang Gao. Cognizing and imitating
467 robotic skills via a dual cognition-action architecture. In *Proceedings of the 23rd International*
468 *Conference on Autonomous Agents and Multiagent Systems*, pages 2204–2206, 2024.
- 469 [33] Peter Pastor, Heiko Hoffmann, Tamim Asfour, and Stefan Schaal. Learning and generalization
470 of motor skills by learning from demonstration. In *2009 IEEE International Conference on*
471 *Robotics and Automation*, pages 763–768. IEEE, 2009.
- 472 [34] Jens Kober, Jan Peters, Jens Kober, and Jan Peters. Movement templates for learning of hitting
473 and batting. *Learning Motor Skills: From Algorithms to Robot Experiments*, pages 69–82, 2014.

- 474 [35] Katharina Mülling, Jens Kober, Oliver Kroemer, and Jan Peters. Learning to select and
475 generalize striking movements in robot table tennis. *The International Journal of Robotics*
476 *Research*, 32(3):263–279, 2013.
- 477 [36] Karol Hausman, Jost Tobias Springenberg, Ziyu Wang, Nicolas Heess, and Martin Riedmiller.
478 Learning an embedding space for transferable robot skills. In *International Conference on*
479 *Learning Representations*, 2018.
- 480 [37] Utkarsh Aashu Mishra, Shangjie Xue, Yongxin Chen, and Danfei Xu. Generative skill chaining:
481 Long-horizon skill planning with diffusion models. In *Conference on Robot Learning*, pages
482 2905–2925. PMLR, 2023.
- 483 [38] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A
484 survey. *Journal of artificial intelligence research*, 4:237–285, 1996.
- 485 [39] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural*
486 *information processing systems*, 29, 2016.
- 487 [40] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial
488 motion priors for stylized physics-based character control. *ACM Transactions on Graphics*
489 *(ToG)*, 40(4):1–20, 2021.
- 490 [41] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley.
491 Least squares generative adversarial networks. In *Proceedings of the IEEE international*
492 *conference on computer vision*, pages 2794–2802, 2017.
- 493 [42] Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira. How to train your dragan. *arXiv*
494 *preprint arXiv:1705.07215*, 2(4), 2017.
- 495 [43] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do
496 actually converge? In *International conference on machine learning*, pages 3481–3490. PMLR,
497 2018.
- 498 [44] Youngwoon Lee, Edward S Hu, and Joseph J Lim. Ikea furniture assembly environment for
499 long-horizon complex manipulation tasks. In *2021 IEEE International Conference on Robotics*
500 *and Automation (ICRA)*, pages 6343–6349. IEEE, 2021.
- 501 [45] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for
502 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 503 [46] Lucy Xiaoyang Shi, Joseph J Lim, and Youngwoon Lee. Skill-based model-based reinforcement
504 learning. In *Conference on Robot Learning*, pages 2262–2272. PMLR, 2023.
- 505 [47] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
506 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 507 [48] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao
508 Carreira. Perceiver: General perception with iterative attention. In *International Conference on*
509 *Machine Learning*, 2021.
- 510 [49] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper,
511 Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning.
512 In *Conference on Robot Learning*, 2019.
- 513 [50] Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan,
514 Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment:
515 Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.
- 516 [51] Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy
517 learning: Solving long-horizon tasks via imitation and reinforcement learning. In *Conference*
518 *on Robot Learning*, pages 1025–1037. PMLR, 2020.
- 519 [52] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
520 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in
521 pytorch. 2017.

- 522 [53] Dhruv Shah, Błażej Osipiński, Sergey Levine, et al. Lm-nav: Robotic navigation with large
523 pre-trained models of language, vision, and action. In *Conference on robot learning*, pages
524 492–504. PMLR, 2023.
- 525 [54] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn,
526 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics
527 transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- 528 [55] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,
529 Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied
530 multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.

531 Technical Appendix

532 A Pseudo-code

533 Pseudo-code for adaptive sub-task skill learning and bi-directional adversarial learning are shown in
534 Algorithm 1 and Algorithm 2 respectively. We highlight the key differences between our method and
535 the most relevant T-STAR with a gray background.

536 A.1 Adaptive Sub-task Skill Learning

537 As shown in Algorithm 1, the innovation of the sub-task skill learning scheme we propose, compared
538 to previous methods, consists of two parts: 1) We use a more stable weighted reward function for
539 policy learning of sub-task skills, as shown in Eq. 1 and Eq. 3 in the main paper. 2) We introduce
540 AES regularization constraints into this weighted reward function to periodically adaptively schedule
541 the scale of the two reward terms, as shown in line 11-14 of Algorithm 1, allowing the robot to fully
explore and learn from both the environment and the expert behaviors.

Algorithm 1 Adaptive Sub-task Skill Learning.
Key differences to T-STAR [15] in gray.

```
1: Require: expert demonstrations  $\mathbb{D}_1^E, \dots, \mathbb{D}_K^E$ , sub-task MDPs  $\mathcal{M}_1, \dots, \mathcal{M}_K$ 
2: Initialize sub-task policies  $\pi_\theta^1, \dots, \pi_\theta^K$ , least-squares GAIL discriminator  $f_\phi^1, \dots, f_\phi^K$ .
3: Initialize imitation progress recognizer  $\Phi$  with  $\mathbb{D}^E$ , balance discount factor  $\lambda_{\text{RL}} \leftarrow \alpha$ ,  $\lambda_{\text{IL}} \leftarrow 1 - \alpha$ .
4: for each sub-task  $i = 1, \dots, K$  do
5:   for episode = 1, 2,  $\dots$ ,  $N$  do
6:     Rollout trajectories  $\tau = (s_1, a_1, r_1^{\text{Env}}, \dots, s_T)$  with  $\pi_\theta^i$ 
7:     // WEIGHTED REWARD FUNCTION
8:     Compute balanced reward  $\{r_1, \dots, r_{T-1}\} \leftarrow \lambda_{\text{RL}} r^{\text{Env}} + \lambda_{\text{IL}} r^{\text{Pred}}$ 
9:     Update  $f_\phi^i$  with  $\tau$  and  $\tau^E \sim \mathbb{D}_i^E$  using Eq. 3
10:    Update  $\pi_\theta^i$  with the rewarded trajectories  $\{s_1, a_1, r_1, \dots, s_T\}$ 
11:    // ADAPTIVE EQUILIBRIUM SCHEDULING REGULARIZATION
12:    Update imitation progress recognizer  $\Phi$  with  $\tau$  and  $\tau^E \sim \mathbb{D}_i^E$ 
13:    Query  $\Phi$  about the current imitation progress  $p$ 
14:    Update balance discount factor  $\lambda_{\text{RL}}, \lambda_{\text{IL}} \leftarrow \varphi_\lambda(p)$ 
15:   end for
16: end for
```

542

543 A.2 Bi-directional Adversarial Learning

544 As shown in Algorithm 2, the innovation of the bi-directional adversarial learning mechanism consists
545 of two parts: 1) We propose a bi-directional regularization which is trained by two balanced bi-
546 directional constraints to better chain sequential skills, as shown in line 16-17 of Algorithm 2. 2)
547 We also employ the adaptive sub-skill learning scheme during the bi-directional adversarial learning
548 process in order to ensure inter-skill alignment while enabling the sub-task skills to be adaptively
549 adjusted to task changes during fine-tuning as well, as shown in line 10-12 of Algorithm 2.

550 B More Details on AES Regularization

551 Following the mechanism described in ADS [9], our AES also employs an imitation progress
552 recognizer Φ to monitor the extent to which the agent has assimilated the expert’s behaviors. The
553 main idea is to assess the closeness of the pair of trajectories by evaluating the agent-collected
554 trajectory $\tau = (s_0, \dots, s_T)$ and the expert trajectory $\tau^E = (s_0^E, \dots, s_T^E)$ through a monotonic
555 state-by-state alignment.

556 To be specific, Φ receives the agent’s collected trajectories τ (line 12 in Algorithm 1) and infers
557 the agent’s current imitation progress p , $p \in [0, T)$ (line 13 in Algorithm 1). The construction

Algorithm 2 Bi-directional Adversarial Learning
Key differences to T-STAR [15] in gray.

- 1: **Require:** expert demonstrations $\mathbb{D}_1^E, \dots, \mathbb{D}_K^E$, sub-task MDPs $\mathcal{M}_1, \dots, \mathcal{M}_K$, pre-trained sub-task policies $\pi_\theta^1, \dots, \pi_\theta^K$, pre-trained GAIL discriminator $f_\phi^1, \dots, f_\phi^K$.
 - 2: Initialize dual set discriminator $\zeta_\omega^1, \dots, \zeta_\omega^K$, imitation identifier Φ with \mathbb{D}^E , balance discount factor $\lambda_{\text{RL}} \leftarrow \alpha, \lambda_{\text{IL}} \leftarrow 1 - \alpha$.
 - 3: Initialize initial state buffers $\mathcal{B}_I^1, \dots, \mathcal{B}_I^K$, and terminal state buffers $\mathcal{B}_\beta^1, \dots, \mathcal{B}_\beta^K$.
 - 4: **for** iteration $m = 0, 1, \dots, M$ **do**
 - 5: **for** each sub-task $i = 1, \dots, K$ **do**
 - 6: Sample s_0 from environment or \mathcal{B}_β^{i-1}
 - 7: Rollout trajectories $\tau = (s_1, a_1, r_1, \dots, s_T)$ with pre-trained π_θ^i
 - 8: **if** τ is successful **then**
 - 9: $\mathcal{B}_I^i \leftarrow \mathcal{B}_I^i \cup s_1, \mathcal{B}_\beta^i \leftarrow \mathcal{B}_\beta^i \cup s_T$
 - 10: // ADAPTIVE EQUILIBRIUM SCHEDULING
 - 11: Update imitation identifier Φ with τ
 - 12: Query Φ about the current imitation progress p
 - 13: **end if**
 - 14: Update balance discount factor $\lambda_{\text{RL}}, \lambda_{\text{IL}} \leftarrow \varphi_\lambda(p)$
 - 15: Fine-tune f_ϕ^i with τ and $\tau^E \sim \mathbb{D}_i^E$
 - 16: // TRAIN BI-DIRECTIONAL DISCRIMINATOR
 - 17: Update ζ_ω^i with $s_\beta \sim \mathcal{B}_\beta^{i-1}$ and $s_I \sim \mathcal{B}_I^i$ with $\mathcal{L}_i(\omega) = \frac{1}{2}\mathcal{C}_1 + \frac{1}{2}\mathcal{C}_2$
 - 18: // FINE-TUNE WITH DUAL REGULARIZATION
 - 19: Update π_θ^i with $r_i(s_t, a_t, s_{t+1}; \phi, \omega)$ using Eq. 7
 - 20: **end for**
 - 21: **end for**
-

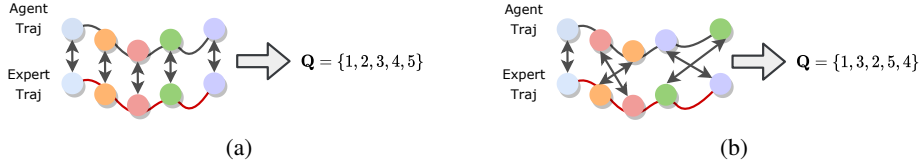


Figure 7: Visualization of the construction of the sequence \mathbf{Q} . To be more intuitive, we directly represent the minimum cosine similarity with double arrows.

558 of Φ , with reference to ADS, first requires the construction of a sequence $\mathbf{Q}(q_1, \dots, q_T)$, where
559 $q_i = \arg\min_j c(s_i, s_j^E)$ is the index of the nearest neighbor of s_i in τ^E , c is the cosine similarity. As
560 shown in Fig. 7, If τ and τ^E are exactly the same, then \mathbf{Q} becomes a strictly increasing sequence
561 (Fig 7(a)). On the contrary, if τ and τ^E characterize some different behaviors, there are some
562 unordered sequences in \mathbf{Q} (Fig 7(b)).

563 After constructing \mathbf{Q} , the progress alignment between τ and τ^E is measured as the length of
564 the longest increasing subsequence (LIS) in \mathbf{Q} , denoted as $LIS(\tau, \tau^E)$. For instance, if $\mathbf{Q} =$
565 $\{1, 3, 2, 5, 4\}$ as in Fig 7(b), then its LIS can be $\{1, 3, 5\}$, $\{1, 2, 5\}$, $\{1, 3, 4\}$ or $\{1, 2, 4\}$. The LIS
566 measurement concentrates on the consistency of the macroscopic trends in these trajectories, thereby
567 preventing overfitting to the microscopic features in the observation [9].

568 Further, if the following inequality Eq. 8 holds, this indicates that at this time step, the agent's
569 imitation of the expert's action is equivalent to the level of the expert's performance, then the agent's
570 imitation progress p will increase by 1:

$$\max_{\hat{\tau}^E \in \mathbb{D}^E} LIS(\tau_{1:p+1}, \hat{\tau}_{1:p+1}^E) \geq \rho \times \min_{\hat{\tau}^E, \hat{\tau}_{1:p+1}^E \in \mathbb{D}^E} LIS(\hat{\tau}_{1:p+1}^E, \hat{\tau}_{1:p+1}^E), \quad (8)$$

571 where $\hat{\tau}^E \neq \tau^E$, the subscript $1:p+1$ denotes the first $p+1$ steps of the extracted trajectory, and
572 $\rho \in [0, 1]$ controls the stringency of the imitation progress monitoring.

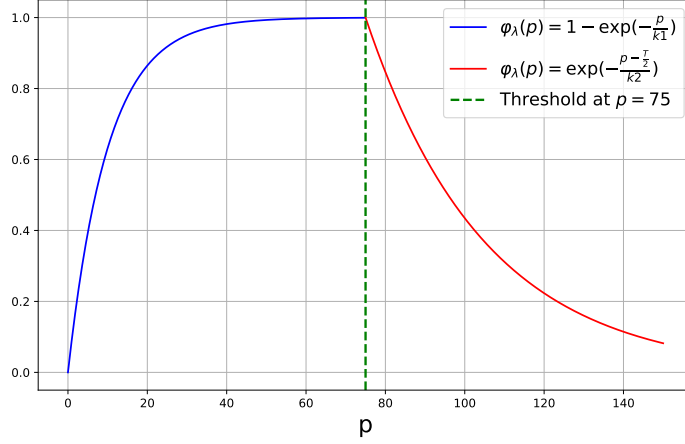


Figure 8: Visualization of the mapping function $\varphi_\lambda(p)$. In this example, we assume that $T = 150$.

573 After obtaining the current imitation progress p of the agent, AES then adopts a mapping function
574 $\varphi_\lambda(p)$ to schedule the two new balance discount factors λ_{RL} and λ_{IL} . Straightforward idea of setting
575 $\varphi_\lambda(p)$ is that **If p reaches a certain threshold, i.e., the agent is able to imitate the expert's behavior
576 well, then the more the agent tends to imitate the expert's behavior in subsequent training, and
577 vice versa.** Therefore, we set the threshold as $\frac{T}{2}$. If $p \in [0, \frac{T}{2})$, we propose $\varphi_\lambda(p) = 1 - e^{-\frac{p}{k_1}}$;
578 if $p \in [\frac{T}{2}, T)$, we propose $\varphi_\lambda(p) = e^{-\frac{p - \frac{T}{2}}{k_2}}$, where k is used to flatten the curve of the mapping
579 function. The mapping function shown in Fig. 8, where $T = 150$. In our experiments, we use
580 different flatten factors for the two stages, where $k_1 = 10$ and $k_2 = 30$.
581 Then λ_{RL} and λ_{IL} are scheduled to be :

$$\begin{cases} \lambda_{\text{RL}} = \alpha^{\tanh(\frac{p}{k_1})}, \lambda_{\text{IL}} = 1 - \alpha^{\tanh(\frac{p}{k_1})}. & \text{if } p \in [0, \frac{T}{2}) \\ \lambda_{\text{IL}} = \alpha^{\tanh(\frac{p - \frac{T}{2}}{k_2})}, \lambda_{\text{RL}} = 1 - \alpha^{\tanh(\frac{p - \frac{T}{2}}{k_2})}. & \text{if } p \in [\frac{T}{2}, T) \end{cases} \quad (9)$$

582 As can be seen from Eq. 9, Fig 8 and
583 Fig. 9, the scale of λ_{RL} is scheduled
584 to be larger than λ_{IL} when p does not
585 reach the imitation process threshold,
586 but this gap gets smaller and smaller
587 as p gets larger. When p reaches the
588 threshold $\frac{T}{2}$, the scale of λ_{IL} is sched-
589 uled to be larger than λ_{RL} , while the
590 scale of λ_{IL} increases as the agent imi-
591 tates better.

592 **Thus, if p is larger and reaches a
593 threshold step, i.e., the agent is able
594 to imitate the expert's behavior well,
595 then the more the agent tends to imi-
596 tate the expert's behavior in subse-
597 quent training, and vice versa.** The
598 entire process is adaptively scheduled
599 based on Φ periodic monitoring of
600 the agent's imitation process. Con-
601 sequently, the RL and IL components of sub-task skill learning can be adaptively scheduled and
602 regularized through AES, effectively enhancing *intra-skill dependencies* between sequential actions.

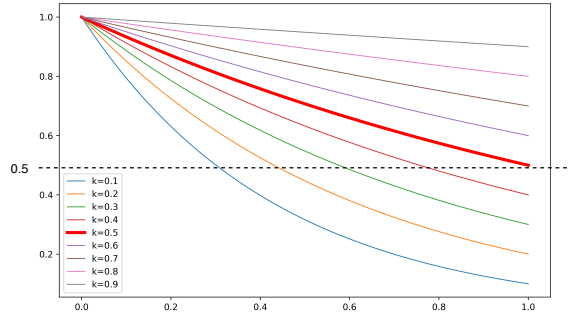


Figure 9: $\alpha^{\varphi_\lambda(p)}$ based on the variation of different α sizes in $\varphi_\lambda(p) \in [0, 1]$. We use $\alpha = 0.5$ as the base in our experiments.

603 C Sub-task Skills

604 In our simulation experiments, we use sequences of sub-tasks defined internally by the environ-
605 ment [44, 45] as task decomposition sub-tasks. Here we list these sequential skills to emphasize the
606 difficulty of long-horizon tasks. Each skill takes a 3D position as the input g_* .

607 **IKEA Furniture Assembly:**

608 **Chair_agne (2 sub-task skills):** Assemble stool leg 0 to target position $g_*^0 \rightarrow$ Assemble stool leg 1
609 to target position g_*^1

610 **Chair_bernhard (2 sub-task skills):** Assemble support leg 0 to target position $g_*^0 \rightarrow$ Assemble
611 support leg 1 to target position g_*^1

612 **Table_dockstra (3 sub-task skills):** Assemble table leg 0 to target position $g_*^0 \rightarrow$ Assemble table leg
613 1 to target position $g_*^1 \rightarrow$ Assemble table top to target position g_*^3

614 **Chair_ingolf (4 sub-task skills):** Assemble chair support 0 to target position $g_*^0 \rightarrow$ Assemble chair
615 support 1 to target position $g_*^1 \rightarrow$ Assemble front leg 0 to target position $g_*^3 \rightarrow$ Assemble front leg 1
616 to target position g_*^4

617 **Table_lack (4 sub-task skills):** Assemble table leg 0 to target position $g_*^0 \rightarrow$ Assemble table leg 1
618 to target position $g_*^1 \rightarrow$ Assemble table leg 2 to target position $g_*^3 \rightarrow$ Assemble table leg 3 to target
619 position g_*^4

620 **Toy_table (4 sub-task skills):** Assemble table leg 0 insert to target position $g_*^0 \rightarrow$ Assemble table leg
621 1 insert to target position $g_*^1 \rightarrow$ Assemble table leg 2 insert to target position $g_*^3 \rightarrow$ Assemble table
622 leg 3 insert to target position g_*^4

623 **Kitchen Organization:**

624 **Kitchen (4 sub-task skills):** Turn on the microwave to target position $g_*^0 \rightarrow$ Move the kettle to target
625 position $g_*^1 \rightarrow$ Turn on the stove (rotate the stove button to target position g_*^2) \rightarrow Turn on the light
626 (rotate the light button to target position g_*^3)

627 **Extended Kitchen (5 sub-task skills):** Turn on the microwave to target position $g_*^0 \rightarrow$ Turn on the
628 stove (rotate the stove button to target position g_*^1) \rightarrow Turn on the light (rotate the light button to
629 target position g_*^2) \rightarrow Slide the cabinet to the right target position $g_*^3 \rightarrow$ Open the cabinet to target
630 position g_*^4

631 D More Quantitative Results

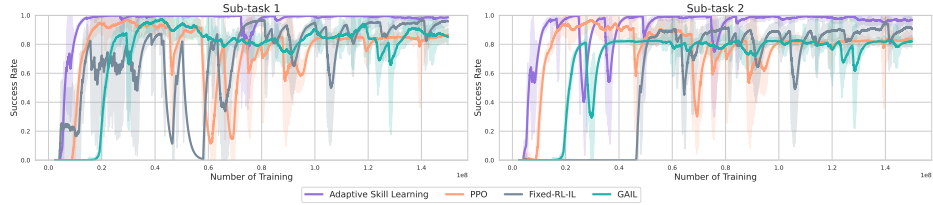
632 We present the training curves with different skill learning methods for sub-task skills in *chair_ingolf*
633 task, and we further present the evaluation performance of the pre-trained skills with different methods
634 across sub-tasks in the other 6 long-horizon simulation tasks. Also, we test the algorithms trained
635 from scratch in the presence of perturbations to further illustrate the importance of the execution of
636 sub-tasks on long-horizon tasks.

637 Additionally, the main paper does not delve into the loss function $\mathcal{L}_i(\omega)$ concerning the different
638 scales of the bi-directional constraints in bi-directional adversarial training. Therefore, we conduct
639 further ablation experiments to examine the impact of different scales of the two constraints in the
640 bi-directional discriminator.

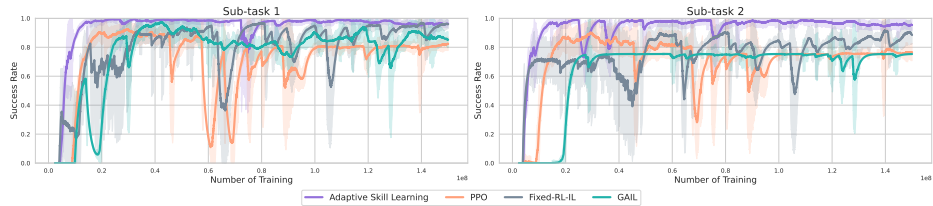
641 **D.1 Sub-task Skill Learning Performance**

642 **D.1.1 Training performance**

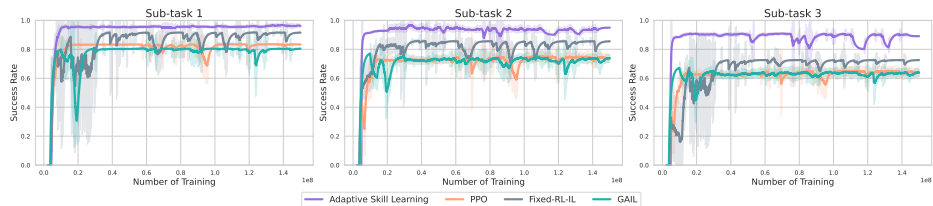
643 Fig. 10 shows the sub-task skill training curves in IKEA furniture assembly tasks. All methods are
644 trained in each sub-task with 5 random seeds, 15M environment steps. As can be seen, the sub-task
645 skill training based on PPO (learning only from environmental feedback), GAIL (learning only from
646 expert demonstrations) and Fixed-RL-IL learning from a fixed scale of environmental feedback
647 and expert demonstration) cannot maintain stability and exhibits significant training performance
648 degradation as the sub-task stage increases. In contrast, the sub-task skill training process using our
649 proposed adaptive sub-skill learning scheme has always been relatively stable and better performing.



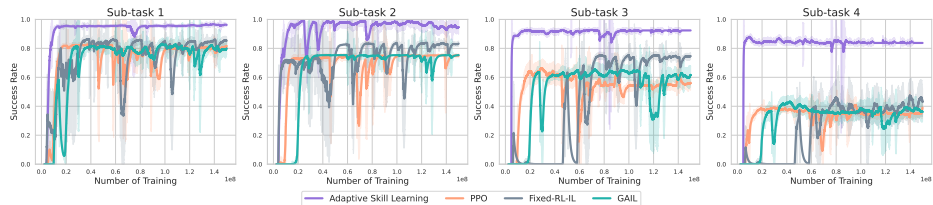
(a) chair_agne



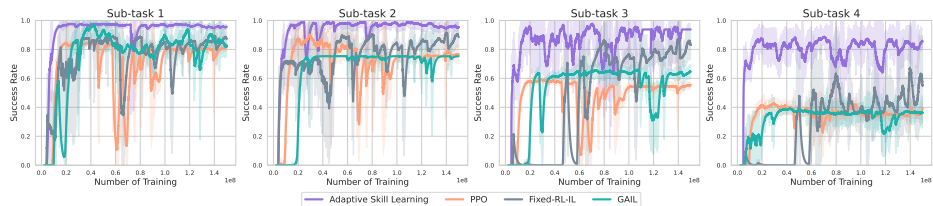
(b) chair_bernhard



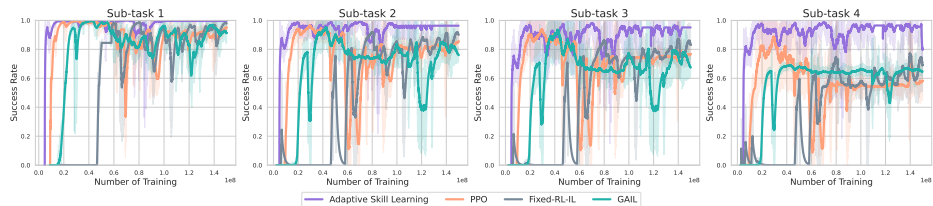
(c) table_dockstra



(d) table_lack



(e) toy_table



(f) chair_ingolf

Figure 10: Training curves for sub-task skills in IKEA furniture assembly tasks. The y-axis represents the success rate of the sub-task.

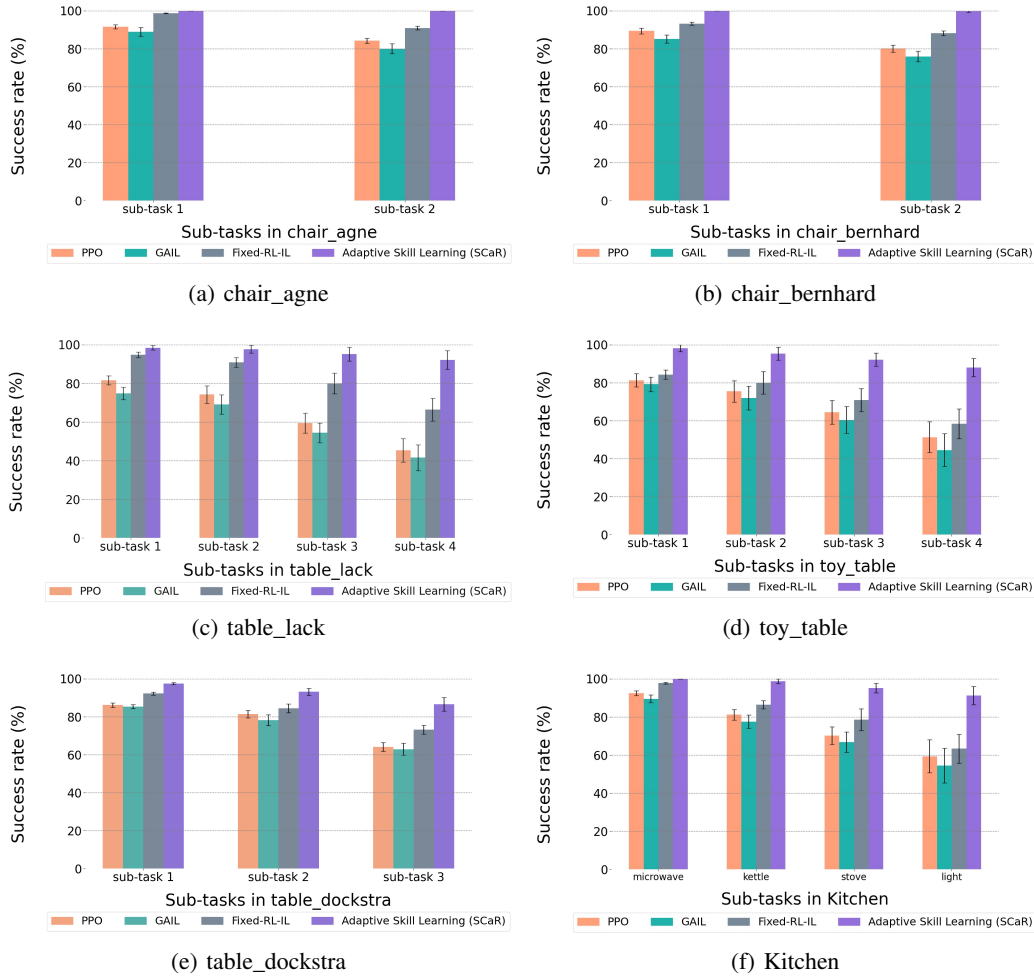


Figure 11: Evaluation Performance Comparison of Sub-task Skill Learning.

650 D.1.2 More evaluation performance

651 As shown in Fig. 11, in *chair_agne*, *chair_bernhard*, *table_lack*, *toy_table*, *table_dockstra*, and
 652 Kitchen tasks, even with the increase of objects in the environment - and the increase of unpredictable
 653 perturbations - our proposed adaptive skill learning learns better sub-task skills. In contrast, the PPO,
 654 GAIL, and Fixed-RL-IL baselines fail to maintain well-learning sub-task skills.

655 These results further corroborate that our proposed AES regularization can reinforce *inter-step*
 656 *dependencies* to the sequential actions within each sub-task skill, and thus pre-train better sub-task
 657 skills for long-horizon tasks.

658 D.2 Robustness to Perturbations

659 We test the algorithms trained from scratch in the presence of perturbations. As shown in Table 3,
 660 algorithms trained from scratch fail to successfully complete the task when environment perturbations
 661 occur during execution. This further illustrates the importance of dividing sub-tasks for multi-stage
 662 execution on long-horizon manipulation tasks that are contact-rich and subject to unanticipated
 663 perturbations. It also supports the significance of our work on long-horizon robotic manipulation
 664 tasks.

Table 3: Success rates of completing the two sub-tasks *chair_bernhard* and four sub-tasks *chair_ingolf* in stationary and perturbed environments.

Method	<i>chair_bernhard</i>		<i>chair_ingolf</i>	
	No Perturb	Perturb	No Perturb	Perturb
PPO (Scratch RL)	0.42 ± 0.12	0.01 ± 0.00	0.14 ± 0.03	0.00 ± 0.00
GAIL (Scratch IL)	0.23 ± 0.02	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
Fixed-RL-IL	0.53 ± 0.07	0.05 ± 0.00	0.22 ± 0.00	0.00 ± 0.00
SkiMo	0.62 ± 0.05	0.10 ± 0.00	0.47 ± 0.03	0.00 ± 0.00
Policy Sequencing	0.82 ± 0.09	0.51 ± 0.04	0.77 ± 0.12	0.50 ± 0.10
T-STAR	0.90 ± 0.04	0.60 ± 0.08	0.89 ± 0.04	0.59 ± 0.04
SCaR w/o Bi	0.92 ± 0.02	0.65 ± 0.11	0.91 ± 0.01	0.63 ± 0.05
SCaR w/o AES	0.94 ± 0.03	0.74 ± 0.09	0.93 ± 0.02	0.71 ± 0.07
SCaR (Ours)	0.96 ± 0.04	0.85 ± 0.11	0.95 ± 0.03	0.80 ± 0.13

665 D.3 Further Ablation

666 We set the loss function for the bi-directional discriminator in the main paper as $\mathcal{L}_i(\omega) = \frac{1}{2}\mathcal{C}_1 + \frac{1}{2}\mathcal{C}_2$,
 667 where the bi-directional constraints $\mathcal{C}_1, \mathcal{C}_2$ are defined as:

$$\begin{aligned}
 \text{next initial} \rightarrow \text{previous terminal: } \mathcal{C}_1 &= \mathbb{E}_{s_T \sim \mathcal{I}_i} [\zeta_\omega^i(s_T) - 1]^2 + \mathbb{E}_{s_T \sim \beta_{i-1}} [\zeta_\omega^i(s_T)]^2 \\
 \text{previous terminal} \rightarrow \text{next initial: } \mathcal{C}_2 &= \mathbb{E}_{s_T \sim \beta_i} [\zeta_\omega^i(s_T) - 1]^2 + \mathbb{E}_{s_T \sim \mathcal{I}_{i+1}} [\zeta_\omega^i(s_T)]^2
 \end{aligned} \tag{10}$$

668 The first constraint \mathcal{C}_1 trains the policy to have the initial states approach the terminal states of the
 669 previous policy, while the second constraint \mathcal{C}_2 trains the policy to have the terminal states close to
 670 the initial states of the next policy. In the experiments, these two constraints have the same scale in
 671 the training process of the bi-directional discriminator.

672 We wonder whether different scales of these two terms would lead to different performances, and
 673 for this reason, we conduct further parametric ablation experiments to explore this. Specifically,
 674 we define the scale parameter of the first term \mathcal{C}_1 as d_1 , and the second term \mathcal{C}_2 as $d_2 = 1 - d_1$,
 675 and set **0.1, 0.3, 0.5, 0.7, 0.9** for d_1 respectively for comparison experiments. We test the effect of
 676 different scales of bi-directional adversarial training items d_1 and d_2 on the success rate of SCaR in
 677 each of the four tasks: *chair_agne*, *chair_ingolf*, *table_dockstra*, and extend kitchen. As shown in
 678 Fig. 12, the experimental result is also in line with our intuition that when the ratio of the two terms
 679 **initial** \rightarrow **previous terminal** and **terminal** \rightarrow **next initial** is the same, the performance is the best
 680 among the four tasks, whereas when the more imbalanced the scale of the two terms is, the worse the
 681 performance is.

682 This ablation result further demonstrate our statement in Sec. 4.3 in the main paper: **The purpose**
 683 **of the bi-directional discriminator is to establish a balanced mapping relationship between**
 684 **the initial states and terminal states to ensure the coherence and stability of the policy.** If the
 685 constraint in one direction (e.g., from initial states to terminal states) is stronger than the constraint in
 686 the other direction (e.g., from terminal states to initial states), the information transmission becomes
 687 asymmetric. This asymmetry results in better training in one direction and insufficient training in the
 688 other, thereby affecting overall performance.

689 E More Qualitative Results

690 Fig 13 shows the qualitative comparison of skill chaining methods. Their animated versions can be
 691 found on our project website.

692 F Real-Robot Long-Horizon Manipulation via Sim-to-Real Transfer

693 **Real-robot Experiment Setup** We also evaluate the skill chaining performance of real-robot for
 694 solving simple yet intuitive real-world long-horizon manipulation. We set up two types of desktop-

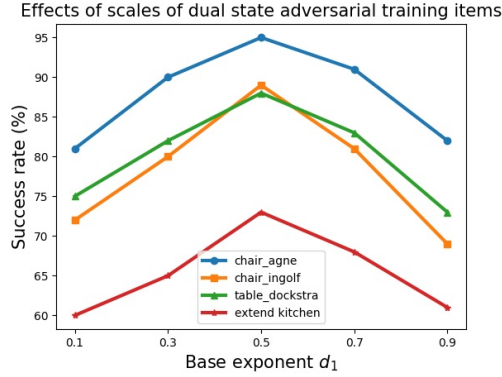


Figure 12: Impact on skill chaining performance of different scales of bi-directional constraints in SCaR.

Table 4: Skill chaining performance of real-world long-horizon robotic manipulation tasks.

Method	Success rate
T-STAR	70% (2 sub-tasks) / 50% (3 sub-tasks)
SCaR	90% (2 sub-tasks) / 70% (3 sub-tasks)

695 level long-horizon manipulation tasks. The robotic arm needs to pick-and-place 2 and 3 blue squares
 696 in sequence, as shown in the top figures in Table 4.

697 We built the corresponding task environment using the gazebo simulation that accompanies the
 698 K1 robot⁴, and collect 50 demonstrations of grasping skills for each square for training. With
 699 camera calibration, we deploy agents trained under simulation in a real robot desktop task to solve
 700 2-square as well as 3-square pick-and-place tasks without the need for adaptation processes. We
 701 conduct experiments with the Sagittarius K1 and use MoveIt2 library based on ROS 2 framework for
 702 controlling the arm. We use RGB observations from RealSense D435i camera on the wrist of the
 703 robotic arm.

704 **Results** For evaluation, we measure the success rate across 10 randomized square positions for
 705 each task. As shown in Table. 4, SCaR can solve the two long-horizon tasks and outperforms T-STAR
 706 baseline. Fig. 14 and Fig. 15 show the qualitative results of successful skill chaining in the 2 and
 707 3-blue-square pick-and-place tasks using SCaR. Video demonstrations are available at our webpage:
 708 <https://tinyurl.com/4333d6np>.

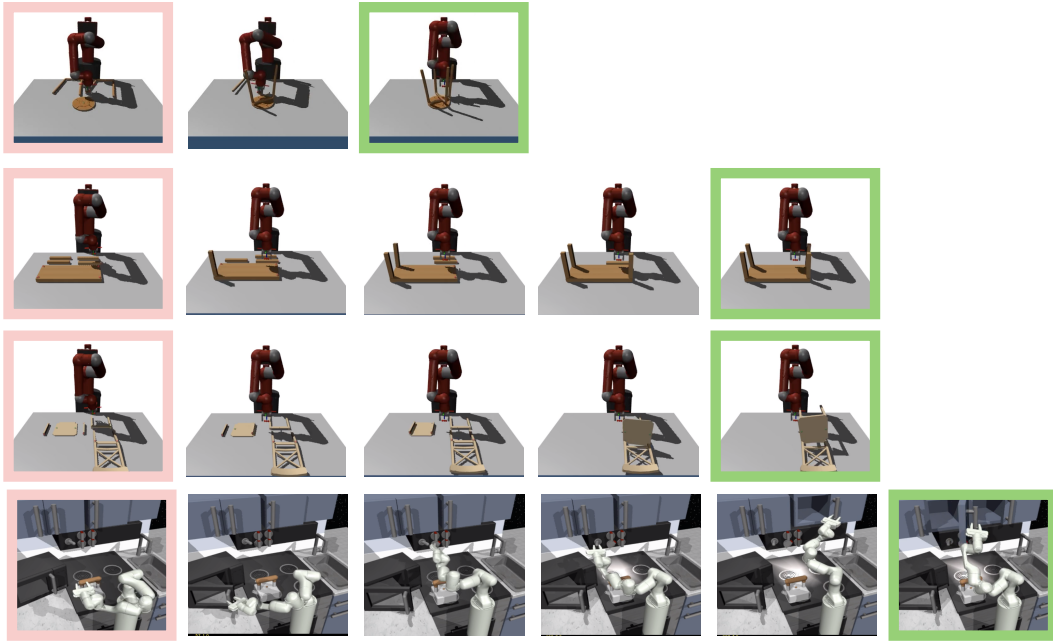
709 G Environment Details

710 G.1 IKEA Furniture Assembly

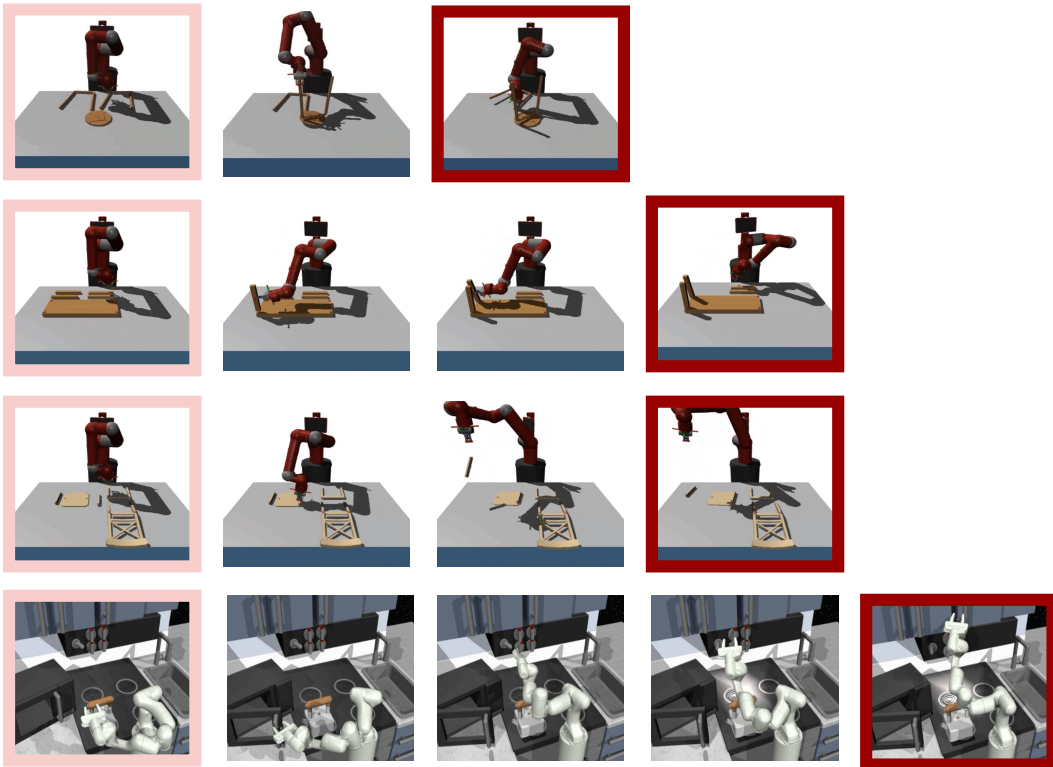
711 We choose six tasks, *chair_agne*, *chair_bernhard*, *chair_ingolf*, *toy_table*, *table_lack* and *ta-*
 712 *ble_bjorkudden* from the IKEA furniture assembly environment⁵ [44] as the focal points of our

⁴https://github.com/NXROBO/sagittarius_ws

⁵<https://github.com/clvrai/furniture>



(a) SCaR - Successful



(b) T-STAR - Failed

Figure 13: Qualitative results of successful skill chaining performance with SCaR and failed skill chaining performance with T-STAR. More qualitative results can be found on our project website <https://tinyurl.com/4333d6np>.

713 experiments, as shown in Fig. 17. Our chosen robotic platform is the 7-DoF Rethink Sawyer robot,
 714 and we control it using joint velocity commands.

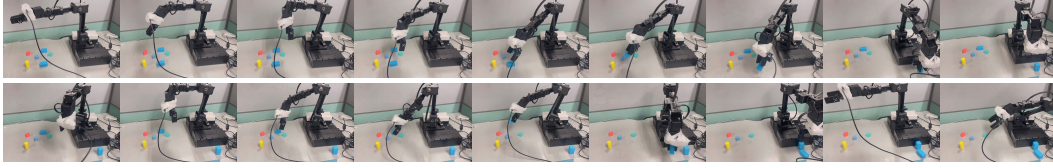


Figure 14: Visualization of the successful skill chaining in the 2-blue-square pick-and-place tasks using SCaR.

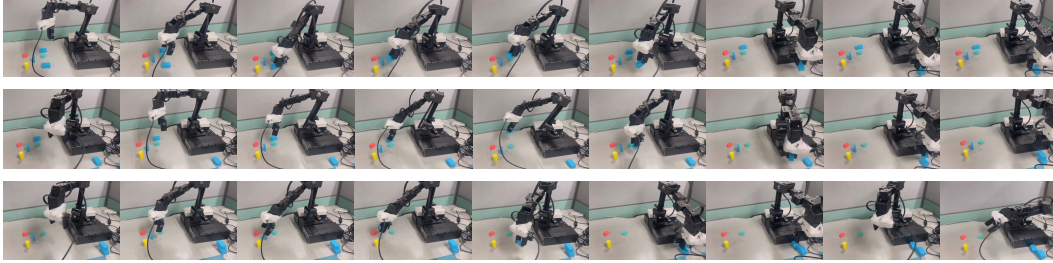


Figure 15: Visualization of the successful skill chaining in the 3-blue-square pick-and-place tasks using SCaR.

715 **Observation Space** The observation space comprises three key components: robot observations
 716 (29 dimensions), object observations (35 dimensions), and task phase information (8 dimensions).
 717 Robot observations encompass robot joint angles (7 dimensions), joint velocities (7 dimensions),
 718 gripper state (2 dimensions), gripper position (3 dimensions), gripper quaternion (4 dimensions),
 719 gripper velocity (3 dimensions), and gripper angular velocity (3 dimensions). Object observations
 720 include the positions (3 dimensions) and quaternions (4 dimensions) of all five furniture pieces in the
 721 scene. Task information, an 8-dimensional one-hot encoding, represents the current phase, including
 722 actions like reaching, grasping, lifting, moving, and aligning.

723 **Action space** The action space includes arm movement, gripper control, and the connect action,
 724 which can vary based on different control modes: 6D end-effector space control using inverse
 725 kinematics, joint velocity control, and joint torque control.

726 In the context of reinforcement learning (RL), we utilize a heavily shaped multi-phase dense reward
 727 obtained from the IKEA Furniture Assembly Environment [44].

728 **Environmental Reward Function** The IKEA furniture assembly environmental reward function is
 729 a multi-phase reward defined with respect to a pair of furniture parts to attach (e.g., a table leg and a
 730 table top) and the corresponding manually annotated way-points, such as a target gripping point g
 731 for each part. The reward function for a pair of furniture parts consists of eight different phases as
 732 follows:

- 733 • **Initial phase:** The robot has to reconfigure its arm pose to an appropriate pose \mathbf{p}_{init} for
 734 grasping a new furniture part. The reward is proportional to the negative distance between
 735 the end-effector \mathbf{p}_{eff} and \mathbf{p}_{init} .
- 736 • **Reach phase:** The robot reaches above a target furniture part. The reward is proportional to
 737 the negative distance between the end-effector \mathbf{p}_{eff} and a point $\mathbf{p}_{\text{reach}}$ 5 cm above the gripping
 738 point g .
- 739 • **Lower phase:** The gripper is lowered onto the target part. The phase reward is proportional
 740 to the negative distance between \mathbf{p}_{eff} and the target gripping points.
- 741 • **Grasp phase:** The robot learns to grasp the target part. The reward is given if the gripper
 742 contacts the part, and is proportional to the force exerted by the grippers.
- 743 • **Lift phase:** The robot lifts the gripped part up to \mathbf{p}_{lift} . The reward is proportional to the
 744 negative distance between the gripped part \mathbf{p}_{part} and the target point \mathbf{p}_{lift} .

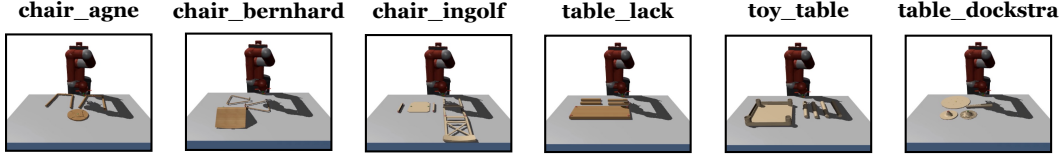


Figure 16: IKEA Furniture Assembly Environment for Long-Horizon Complex Manipulation Tasks.

- 745 • **Align phase:** The robot roughly rotates the gripped part before moving it. The reward is
746 proportional to the cosine similarity between up vectors \mathbf{u}_A , \mathbf{u}_B and forward vectors \mathbf{f}_A , \mathbf{f}_B
747 of the two connectors.
- 748 • **Move phase:** The robot moves and aligns the gripped part to another part. The reward is
749 proportional to the negative distance between the connector of the gripped part and a point
750 $\mathbf{p}_{\text{move_to}}$ 5 cm above the connector of another part, and the cosine similarity between two
751 connector up vectors, \mathbf{u}_A and \mathbf{u}_B , and forward vectors \mathbf{f}_A and \mathbf{f}_B . Note that all connectors
752 are labeled with aligned up vectors and forward vectors.
- 753 • **Fine-grained move phase:** The robot must finely align two connectors until attached. The
754 same reward is used as the **move phase** with a higher coefficient, making the reward more
755 sensitive to small changes. In addition, when the part is connectable, a reward is provided
756 based on the activation of the connect action $a[\text{connect}]$.

757 Upon completion of each phase, completion rewards are given to encourage the agent to move on to
758 the next phase. In addition to stage-based rewards, control penalties, stabilizing wrist pose rewards,
759 and grasping rewards (i.e., opening the grasping hand only during the initial, arrival, and lower stages)
760 are provided throughout the process. If the robot releases the grasped object, the phase ends early
761 and a negative reward is provided. Phase completion depends on the robot and part configurations
762 satisfying distance and angle constraints with respect to the goal configuration. After all stages are
763 completed, the stage resets to the initial stage. This process repeats until all parts are connected.

764 **Demonstration Collection** For imitation learning (IL), we gathered 200 demonstrations for each
765 furniture part assembly using a programmatic assembly policy. Each demonstration for single-part
766 assembly typically spans 150 steps, reflecting the overall task’s inherently long-horizon nature.

767 **Sub-tasks** In our experiments, we define a sub-task as the process of assembling one part to another.
768 Thus, the *chair_agne* and *chair_bernhard* tasks have two distinct sub-tasks, *table_dockstra* has
769 three distinct sub-tasks, and *chair_ingolf*, *table_lack*, and *toy_table* have four distinct sub-tasks.
770 These sub-tasks are trained independently, with their initial state sampled from the environment and
771 random noise introduced in the $[-2\text{cm}, 2\text{cm}]$ and $[-3^\circ, 3^\circ]$ ranges of the (x, y) plane. Importantly, the
772 decomposition of the sub-tasks is pre-determined, which means that the environment is initialized for
773 each sub-task, and the agent receives a notification when a sub-task is successfully completed. Once
774 the two components are firmly connected, the corresponding sub-task is considered completed and
775 the robotic arm is guided back to its initial pose, i.e., at the center of the workspace.

776 **Assembly Difficulty** The difficulty of modeling furniture depends largely on the shape of the
777 furniture. For example, the *toy_table* task with cylindrical legs is more difficult to grasp, whereas
778 the *table_lack* task with rectangular legs is easier to grasp. Chairs are generally more difficult to
779 assemble because of their irregular shape (e.g., seat and back). This is the reason why the success
780 rates of the *toy_table* and *chair_ingolf* tasks are lower than the success rates of *table_lack*.

781 G.2 Kitchen Organization

782 We use the Franka Kitchen tasks in D4RL [45] and refer to the experimental setup in SkiMo [46] for
783 the sub-task extensions. Including the following two tasks: **Kitchen task** and **Extended Kitchen**
784 **task**, as shown in Fig. 17.

785 **Kitchen** The 7-DoF Franka Emika Panda robot arm is tasked with performing four sequential
786 sub-tasks: *Turn on the microwave - Move the kettle - Turn on the stove - Turn on the lights*.

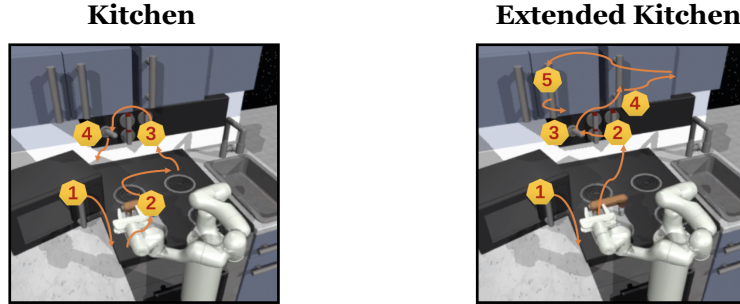


Figure 17: **Kitchen Organization Environment for Long-Horizon Complex Manipulation Tasks.**

787 **Extended Kitchen** The environment and task-agnostic data used in this experiment are consistent
 788 with those employed in the **Kitchen** scenario. However, we introduce a different set of sub-tasks for
 789 this experiment, namely: *Turn on the microwave - Turn on the stove - Turn on the lights - Slide the*
 790 *cabinets to the right - Open the cabinets*, as depicted in Fig. 17 (right). It’s worth noting that this
 791 sequence of tasks is not aligned with the sub-task transition probabilities observed in the task-agnostic
 792 data, posing a challenge for exploration based on prior data.

793 **Observation Space** The agent operates within a 30-dimensional observation space, which includes
 794 an 11-dimensional robot proprioceptive state and 19-dimensional object states. This modified
 795 observation space removes a constant 30-dimensional goal state found in the original environment.

796 **Action Space** The agent’s action space consists of 9 dimensions, encompassing 7-dimensional joint
 797 velocity control and 2-dimensional gripper velocity control.

798 **Environmental Reward Function** In terms of the environmental rewards, the agent receives a
 799 reward of +1 for each completed sub-task. The total episode length is set to 280 steps, and an episode
 800 concludes once all sub-tasks are successfully accomplished. The initial state is initialized with slight
 801 noise introduced in each state dimension.

802 **Demonstration Collection** For imitation learning, we collect 200 demonstrations per sub-task with
 803 reference to the dataset in [51] that obtained through teleoperation. This dataset covers interactions
 804 with all seven manipulatable objects within the environment.

805 H Network Architecture

806 For a fair comparison, our method and the benchmark methods use the same network structure.
 807 The policy network and the critic network consist of two layers of 128 and 256 hidden units fully
 808 connected with ReLU nonlinear properties, respectively. The output layer of the actor network
 809 outputs an action distribution, which consists of the mean and standard deviation of a Gaussian
 810 distribution. The critic network outputs only one critic value. The discriminator of GAIL [39] and
 811 the bi-directional discriminator of our proposed approach use a two-layer fully connected network
 812 with 256 hidden units. The outputs of these discriminators are clipped between [0, 1], following the
 813 least-square GAIL proposed by [40].

814 I Training Details

815 I.1 Computing Resources

816 Our method and all baselines were implemented using PyTorch [52]. All experiments were carried
 817 out on workstations equipped with Intel(R) Xeon(R) Gold 5218 CPUs and NVIDIA GeForce RTX
 818 3080 2 GPUs. Pre-training of each sub-task skill policy in SCaR (150M time steps) took about 10
 819 hours. Testing and evaluation of skill chaining for the entire long-horizon task, approximately 10 to
 820 15 hours, depending on the difficulty of the task. Training of the skill dynamics model in SkiMo [46]
 821 took approximately 24 hours (100M steps), and PPO [47], GAIL [39], and Fix-RL-IL were slower

822 (about 48 hours) because they all train the entire long-horizon task from scratch, with 450M time
 823 steps for the overall long-horizon task.

824 I.2 Algorithm Implementation Details

825 We report the hyperparameters used in our experiments in Table 5.

Table 5: Hyperparameters used in our experiments.

Hyperparameter	Value
Rollout Size	1024
Learning Rate	0.0003
Learning Rate Decay	Linear decay
Mini-batch Size	128
Discount Factor	0.99
Entropy Coefficient	0.003
Reward Scale	0.05
State Normalization	True
Discriminator learning rate	$1e^{-4}$
Sub-task training steps	150000000
# Workers	20
# Epochs per Update	10
Base exponent for balancing α	0.5
k_1 (used to flatten the mapping function during $p \in [0, \frac{T}{2})$)	10
k_2 (used to flatten the mapping function during $p \in [\frac{T}{2}, T)$)	30
Weighting factor λ_{Bi}	10000
ρ (for imitation progress recognizer Φ)	0.9
Penalty coefficient η^{EP}	10

826 For the baseline implementations, we use the official code for PPO [47], GAIL [39], Fixed-RL-IL [40],
 827 SkiMo [46], Policy Sequencing [12] and T-STAR [15]. The table below (Table 6) compares key
 828 components of **SCaR** with model-based, model-free and skill-based baselines and ablated methods,
 829 where *joint training* indicates whether or not reinforcement learning combined with imitation learning
 830 is used for training.

831 **PPO [47]** Any reinforcement learning algorithm can be used for policy optimization, in this paper
 832 we choose to use Proximal Policy Optimization (PPO) and use the default hyperparameters of
 833 PPO [47].

834 **GAIL [39]** In this paper we choose to use Generative Adversarial Imitation Learning (GAIL) [39]
 835 as the learning algorithm for imitation learning and use the default hyperparameters of GAIL [39]. We
 836 specifically use an agent states s to discriminate agent and expert trajectories, instead of state-action
 837 pairs (s, a) .

838 **Fixed-RL-IL [12]** We adopt the AMP [40] solution combining environmental rewards and least
 839 square GAIL with $\lambda_{RL} = \lambda_{IL} = 0.5$. For implementation details of least square GAIL training and
 840 GAIL rewards, see original paper [40].

841 **SkiMo [46]** We use the official implementation of the original paper and use the hyperparameters
 842 suggested in the official implementation.

843 **Policy Sequencing [12]** We employ the official implementation and the hyperparameters provided
 844 by [15].

845 **T-STAR [15]** We use the official implementation of the original paper and use the hyperparameters
 846 suggested in the official implementation [15].

847 **SCaR (ours)** We refer to T-STAR and use $\lambda_{Re} = 10000$ for bi-directional regularization. We take
 848 50% of the initial state samples from the start environment of each policy, 50% of the terminal state
 849 samples at the end, and 50% of the initial state buffer and 50% of the terminal state buffer from the
 850 previous skill, respectively.

Table 6: Comparison to prior work and ablated methods.

Method	Model-based	Skill-based	Scratch training	Joint training
PPO [47] and GAIL [39]	✗	✗	✓	✗
Fixed-RL-IL [40]	✗	✗	✓	✓
SkiMo [46]	✓	✓	✓	✓
Policy Sequencing [12]	✓	✓	✗	✓
T-STAR [15]	✗	✓	✗	✓
SCaR (Ours) and SCaR w/o Bi and SCaR w/o AES	✓	✓	✗	✓

851 J Potential negative impacts

852 Since our method is currently limited to applications in simulated environments and simple desktop-
853 level robot manipulation, it is not expected to have a significant negative impact on society. However,
854 privacy concerns may arise if our method is applied to real-world long time-series tasks with mobility,
855 as imitation learning agents used in applications such as autonomous driving [53] or real-time
856 control [54, 55] require large amounts of data that often contain controversial information. In
857 addition, the imitation learning policy is a challenge because it imitates a specified demonstration
858 that may include bad behavior. If the expert demonstration includes some nefarious behaviors
859 (e.g., training data for a mobile manipulation task includes behaviors that may be violent towards
860 pedestrians), then the policy may have a significant negative impact on the user. To address this issue,
861 future directions should focus on developing agents with safety adaptations in addition to improving
862 performance.

863 **NeurIPS Paper Checklist**

864 **1. Claims**

865 Question: Do the main claims made in the abstract and introduction accurately reflect the
866 paper's contributions and scope?

867 Answer: [\[Yes\]](#)

868 Justification: We propose a new skill chaining framework for long time-series robotic
869 manipulation tasks that improves overall task completion performance by providing dual
870 regularization for intra- and inter-skill dependencies. We hope this work will inspire
871 future research to further explore the potential of skill chaining for long-horizon robotic
872 manipulation.

873 Guidelines:

- 874 • The answer NA means that the abstract and introduction do not include the claims
875 made in the paper.
- 876 • The abstract and/or introduction should clearly state the claims made, including the
877 contributions made in the paper and important assumptions and limitations. A No or
878 NA answer to this question will not be perceived well by the reviewers.
- 879 • The claims made should match theoretical and experimental results, and reflect how
880 much the results can be expected to generalize to other settings.
- 881 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
882 are not attained by the paper.

883 **2. Limitations**

884 Question: Does the paper discuss the limitations of the work performed by the authors?

885 Answer: [\[Yes\]](#)

886 Justification: We discuss limitations in the last section of the main paper: limitations mainly
887 exist in that 1) the sub-tasks in our framework are predefined, 2) we did not test our method
888 on a more challenging real robot furniture assembly task due to limited hardware.

889 Guidelines:

- 890 • The answer NA means that the paper has no limitation while the answer No means that
891 the paper has limitations, but those are not discussed in the paper.
- 892 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 893 • The paper should point out any strong assumptions and how robust the results are to
894 violations of these assumptions (e.g., independence assumptions, noiseless settings,
895 model well-specification, asymptotic approximations only holding locally). The authors
896 should reflect on how these assumptions might be violated in practice and what the
897 implications would be.
- 898 • The authors should reflect on the scope of the claims made, e.g., if the approach was
899 only tested on a few datasets or with a few runs. In general, empirical results often
900 depend on implicit assumptions, which should be articulated.
- 901 • The authors should reflect on the factors that influence the performance of the approach.
902 For example, a facial recognition algorithm may perform poorly when image resolution
903 is low or images are taken in low lighting. Or a speech-to-text system might not be
904 used reliably to provide closed captions for online lectures because it fails to handle
905 technical jargon.
- 906 • The authors should discuss the computational efficiency of the proposed algorithms
907 and how they scale with dataset size.
- 908 • If applicable, the authors should discuss possible limitations of their approach to
909 address problems of privacy and fairness.
- 910 • While the authors might fear that complete honesty about limitations might be used by
911 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
912 limitations that aren't acknowledged in the paper. The authors should use their best
913 judgment and recognize that individual actions in favor of transparency play an impor-
914 tant role in developing norms that preserve the integrity of the community. Reviewers
915 will be specifically instructed to not penalize honesty concerning limitations.

916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: We provide further explanation in the Appendix to explain the assumptions presented in the main paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We further describe the network architecture, training details, dataset, and the open source codebase on which the method is based in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

969 In the case of closed-source models, it may be that access to the model is limited in
970 some way (e.g., to registered users), but it should be possible for other researchers
971 to have some path to reproducing or verifying the results.

972 5. Open access to data and code

973 Question: Does the paper provide open access to the data and code, with sufficient instruc-
974 tions to faithfully reproduce the main experimental results, as described in supplemental
975 material?

976 Answer: [NA]

977 Justification: All simulation environments, datasets, and open source code libraries that can
978 reproduce our method have been described in the Appendix.

979 Guidelines:

- 980 • The answer NA means that paper does not include experiments requiring code.
- 981 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
982 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 983 • While we encourage the release of code and data, we understand that this might not be
984 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
985 including code, unless this is central to the contribution (e.g., for a new open-source
986 benchmark).
- 987 • The instructions should contain the exact command and environment needed to run to
988 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
989 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 990 • The authors should provide instructions on data access and preparation, including how
991 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 992 • The authors should provide scripts to reproduce all experimental results for the new
993 proposed method and baselines. If only a subset of experiments are reproducible, they
994 should state which ones are omitted from the script and why.
- 995 • At submission time, to preserve anonymity, the authors should release anonymized
996 versions (if applicable).
- 997 • Providing as much information as possible in supplemental material (appended to the
998 paper) is recommended, but including URLs to data and code is permitted.

999 6. Experimental Setting/Details

1000 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
1001 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
1002 results?

1003 Answer: [Yes]

1004 Justification: We describe partial details in the main paper and provide further details in the
1005 Appendix.

1006 Guidelines:

- 1007 • The answer NA means that the paper does not include experiments.
- 1008 • The experimental setting should be presented in the core of the paper to a level of detail
1009 that is necessary to appreciate the results and make sense of them.
- 1010 • The full details can be provided either with the code, in appendix, or as supplemental
1011 material.

1012 7. Experiment Statistical Significance

1013 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1014 information about the statistical significance of the experiments?

1015 Answer: [Yes]

1016 Justification: Our experiments perform means and standard deviations for the five seed
1017 results.

1018 Guidelines:

- 1019 • The answer NA means that the paper does not include experiments.

- 1020 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
1021 dence intervals, or statistical significance tests, at least for the experiments that support
1022 the main claims of the paper.
- 1023 • The factors of variability that the error bars are capturing should be clearly stated (for
1024 example, train/test split, initialization, random drawing of some parameter, or overall
1025 run with given experimental conditions).
- 1026 • The method for calculating the error bars should be explained (closed form formula,
1027 call to a library function, bootstrap, etc.)
- 1028 • The assumptions made should be given (e.g., Normally distributed errors).
- 1029 • It should be clear whether the error bar is the standard deviation or the standard error
1030 of the mean.
- 1031 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1032 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1033 of Normality of errors is not verified.
- 1034 • For asymmetric distributions, the authors should be careful not to show in tables or
1035 figures symmetric error bars that would yield results that are out of range (e.g. negative
1036 error rates).
- 1037 • If error bars are reported in tables or plots, The authors should explain in the text how
1038 they were calculated and reference the corresponding figures or tables in the text.

1039 8. Experiments Compute Resources

1040 Question: For each experiment, does the paper provide sufficient information on the com-
1041 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1042 the experiments?

1043 Answer: [Yes]

1044 Justification: We illustrate the computational resources and the time required for the experi-
1045 ments in the Appendix.

1046 Guidelines:

- 1047 • The answer NA means that the paper does not include experiments.
- 1048 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1049 or cloud provider, including relevant memory and storage.
- 1050 • The paper should provide the amount of compute required for each of the individual
1051 experimental runs as well as estimate the total compute.
- 1052 • The paper should disclose whether the full research project required more compute
1053 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1054 didn't make it into the paper).

1055 9. Code Of Ethics

1056 Question: Does the research conducted in the paper conform, in every respect, with the
1057 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1058 Answer: [Yes]

1059 Justification: [TODO]

1060 Guidelines:

- 1061 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1062 • If the authors answer No, they should explain the special circumstances that require a
1063 deviation from the Code of Ethics.
- 1064 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1065 eration due to laws or regulations in their jurisdiction).

1066 10. Broader Impacts

1067 Question: Does the paper discuss both potential positive societal impacts and negative
1068 societal impacts of the work performed?

1069 Answer: [Yes]

1070 Justification: We elaborate on these in the final section of the Appendix.

1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our paper does not release data or models with high risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have cited or provided URLs to all the code, data, and models used in the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 1124 • For scraped data from a particular source (e.g., website), the copyright and terms of
1125 service of that source should be provided.
- 1126 • If assets are released, the license, copyright information, and terms of use in the
1127 package should be provided. For popular datasets, `paperswithcode.com/datasets`
1128 has curated licenses for some datasets. Their licensing guide can help determine the
1129 license of a dataset.
- 1130 • For existing datasets that are re-packaged, both the original license and the license of
1131 the derived asset (if it has changed) should be provided.
- 1132 • If this information is not available online, the authors are encouraged to reach out to
1133 the asset's creators.

1134 13. **New Assets**

1135 Question: Are new assets introduced in the paper well documented and is the documentation
1136 provided alongside the assets?

1137 Answer: [NA]

1138 Justification: Our paper does not release new assets.

1139 Guidelines:

- 1140 • The answer NA means that the paper does not release new assets.
- 1141 • Researchers should communicate the details of the dataset/code/model as part of their
1142 submissions via structured templates. This includes details about training, license,
1143 limitations, etc.
- 1144 • The paper should discuss whether and how consent was obtained from people whose
1145 asset is used.
- 1146 • At submission time, remember to anonymize your assets (if applicable). You can either
1147 create an anonymized URL or include an anonymized zip file.

1148 14. **Crowdsourcing and Research with Human Subjects**

1149 Question: For crowdsourcing experiments and research with human subjects, does the paper
1150 include the full text of instructions given to participants and screenshots, if applicable, as
1151 well as details about compensation (if any)?

1152 Answer: [NA]

1153 Justification: Our paper does not do crowdsourcing experiments and research on human
1154 subjects.

1155 Guidelines:

- 1156 • The answer NA means that the paper does not involve crowdsourcing nor research with
1157 human subjects.
- 1158 • Including this information in the supplemental material is fine, but if the main contribu-
1159 tion of the paper involves human subjects, then as much detail as possible should be
1160 included in the main paper.
- 1161 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1162 or other labor should be paid at least the minimum wage in the country of the data
1163 collector.

1164 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1165 Subjects**

1166 Question: Does the paper describe potential risks incurred by study participants, whether
1167 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1168 approvals (or an equivalent approval/review based on the requirements of your country or
1169 institution) were obtained?

1170 Answer: [NA]

1171 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

1172 Guidelines:

- 1173 • The answer NA means that the paper does not involve crowdsourcing nor research with
1174 human subjects.

- 1175 • Depending on the country in which research is conducted, IRB approval (or equivalent)
- 1176 may be required for any human subjects research. If you obtained IRB approval, you
- 1177 should clearly state this in the paper.
- 1178 • We recognize that the procedures for this may vary significantly between institutions
- 1179 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1180 guidelines for their institution.
- 1181 • For initial submissions, do not include any information that would break anonymity (if
- 1182 applicable), such as the institution conducting the review.