

Overcoming State and Action Space Disparities in Multi-Domain, Multi-Task Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Current multi-task reinforcement learning (MTRL) methods have the
2 ability to perform a large number of tasks with a single policy. However when
3 attempting to interact with a new domain, the MTRL agent would need to be re-
4 trained due to differences in domain dynamics and structure. Because of these
5 limitations, we are forced to train multiple policies even though tasks may have
6 shared dynamics, leading to needing more samples and is thus sample inefficient.
7 In this work, we explore the ability of MTRL agents to learn in various domains
8 with various dynamics by simultaneously learning in multiple domains, without
9 the need to fine-tune extra policies. In doing so we find that a MTRL agent trained
10 in multiple domains induces an increase in sample efficiency of up to 70% while
11 maintaining the overall success rate of the MTRL agent.

12 **Keywords:** Multi-Task Reinforcement Learning, Morphology-Agnostic Policies

13 1 Introduction

14 Reinforcement learning (RL) has shown remarkable success in its application to many different
15 scenarios including game playing [1], controlling stratospheric balloons [2], and controlling fusion
16 reactors [3]. In order to scale RL to accomplish multiple tasks with a single policy, multi-task
17 RL learns across multiple tasks which may accelerate learning in additional tasks which improves
18 sample efficiency. However, current online multi-task RL algorithms are unable to learn across tasks
19 that have different state and action spaces, and different semantics.

20 To overcome this gap in the reuse of RL policies in tasks and domains, this paper aims to study the
21 design decisions and training protocols required for training a single policy across multiple domains,
22 and evaluates their performance in robotic manipulation domains in different domains. For a RL
23 agent to have generalized skills, the RL agent must be able to leverage it’s previous experiences
24 in new tasks. An example of this would be a RL agent that can pick up a frying pan in a kitchen
25 domain and can also pick up a remote control in a living room domain. Various approaches have
26 been proposed to address the challenge of transferring skills across different domains and robot
27 morphologies. Recent works have leveraged offline reinforcement learning where a large dataset of
28 trajectories are used to enable effective learning on different robots [4][5][6][7][8]. However, these
29 policies are trained using datasets of successful trajectories limiting their ability to acquire new
30 skills. Alternatively, online RL approaches generally either receive joint specific observations [9] or
31 attempt to infer robot morphology from the reinforcement learning objective [10]. We differ from
32 these lines of work as we perform all learning in an online fashion while utilizing an architecture that
33 can handle varying state and action spaces without requiring joint-specific observations or explicit
34 morphology inference. Our approach, which we call SAL (Shared, domain-Agnostic, Latent space),
35 enables effective skill transfer across different domains and robot morphologies in an online setting.

36 To address these challenges, SAL employs a unique architecture designed to overcome the limita-
37 tions of varying state and action spaces. A central challenge impeding training a single policy across

38 domains is the difference in state and action space sizes that prevents any direct reuse of a trained
 39 policy, as the policy can have different input and output sizes for each specific domain. This work
 40 aims to enable the transfer of skills between domains with different robotic morphologies, where
 41 both the dimensions and semantic meanings of state and action spaces differ. This domain-agnostic
 42 transfer is enabled by learning the appropriate latent representations through state and action transla-
 43 tion layers, which map the states to a latent state space and then extract an action from latent space to
 44 an domain-specific action space, respectively. This architecture is an extension of the multi-headed
 45 architecture from [11] and [12]. We demonstrate the effectiveness of these layers in continuous
 46 control tasks that transfer high-level manipulation concepts between the domains. Our experiments
 47 show that learning a policy in a shared, domain-agnostic latent space by translating states to this
 48 SAL space and decoding latent skills into actions from the SAL space yields sample efficiency gains
 49 anywhere from 7% to 72% compared to training on a single domain’s tasks alone.

50 The proposed method, and thus the efficient skill transfer across domains, can be achieved through
 51 minimal modifications to existing RL algorithms. We demonstrate skill transfer between the Meta-
 52 World [11] and Franka Kitchen [13] domains, and show how SAL can be applied to state of the art
 53 algorithms, such as Soft-Actor Critic (SAC) [14].

54 Our contributions are as follows:

- 55 • We propose the shared, domain-agnostic, latent policy architecture, which enables multi-
 56 task reinforcement learning across domains with different state and action spaces by using
 57 a single policy in Section 3.
- 58 • We show that learning with our proposed architecture can accelerate learning by up to 72%
 59 in Section 4.4.
- 60 • We provide insights into how the latent space of the SAL architecture is organized to better
 61 understand the limitations and future work associated with learning across differing state
 62 and action spaces in Section 4.6.
- 63 • We provide an open-source implementation to apply SAL on existing RL algorithms.

64 2 Problem Statement

65 Reinforcement learning is formulated using a Markov decision process (MDP) [15], where an MDP
 66 \mathcal{M} is a tuple of (S, A, P, r, γ, p) , where S is the state space, A is the action space, the probability
 67 transition function $P : S \times A \rightarrow [0, 1]^{|S|}$, $r : S \times A \rightarrow \mathbb{R}$ the reward function, γ in $[0, 1)$ is the
 68 discount factor, and p is the initial state distribution.

69 At each time step, the agent observes the state at time t , chooses an action sampled according to
 70 some policy $\pi : S \rightarrow A$ based on s_t , receives a reward for landing in-state s_{t+1} , and observes state
 71 s_{t+1} . The goal of the agent is to find a policy that maximizes expected returns for the current task
 72 $E[R(\tau)]$ where $R(\tau)$ is the sum of discounted rewards along the trajectory induced by following the
 73 policy π .

74 In the multitask reinforcement learning problem, a task distribution must be chosen where N tasks
 75 are sampled from a task set T according to $t \sim n(t)$. Each task can be viewed as having its own
 76 MDP: $(S, A, P_i, r_i, \gamma, p_i)$ where the state space, action space, and discount factors are held constant
 77 across tasks while the probability transition function, reward function, and initial state distribution
 78 are specific to each task t_i . The goal of the multitask reinforcement learning agent is to maximize the
 79 expected sum of discounted rewards across each task $E_{t \sim p(t)}[E_{\tau \sim \pi}[R(\tau)]]$ using policy $\pi_\eta(a|s, z)$,
 80 where z is the task identifier and η are learnable parameters of the policy.

81 In this paper, the problem we consider is extending multitask RL to multiple domains that don’t
 82 share the same state and action spaces. Let E be a set of M domains, where each domain $m \in E$
 83 has a set of T_m distinct tasks. Each task $t_m \in T_m$ in each domain m has an MDP $\mathcal{M}_t^{(m)}$ for each
 84 task t . Therefore, the MDP $\mathcal{M}_j^{(i)}$ is now comprised of the tuple $(S^{(i)}, A^{(i)}, P_j^{(i)}, r_j^{(i)}, \gamma, p_j^{(i)})$, where
 85 the elements in the tuple can change with respect to the domain.

86 The domain and the task can be jointly sampled from a joint distribution $m, t \sim n(m, t)$. Since the
 87 state space $S^{(m)}$ and action space $A^{(m)}$ are domain-dependent, traditional multitask reinforcement
 88 learning algorithms cannot be applied to share policy parameters π_η between tasks because the input
 89 state and output actions can vary in dimensions. In order to overcome this problem, we propose using
 90 translation layers to the input and action heads for the output of the policy. The input translation
 91 layer maps varied state representations to a common latent state space, while the output action heads
 92 decodes actions from this latent space back to domain-specific action spaces. This approach allows
 93 us to maintain a consistent internal representation for the policy, regardless of the specific domain
 94 it’s operating in.

95 To aid the reader in understanding the difference between domains and tasks, we use the follow-
 96 ing definition of domains and tasks. In Meta-World, we leverage MT10 which contains 10 tasks
 97 available in Meta-World. Each task t in MT10 are distinct tasks for the MTRL agent to accomplish
 98 such as: opening a window, closing a drawer, grasping an object and moving it to a goal. Thus the
 99 domain is Meta-World, where the tasks are the individual tasks available in MT10. Similarly, for
 100 the Franka Kitchen domain, there are several tasks for the MTRL agent to accomplish such as slide
 101 a cabinet, open the microwave, grasp a kettle by the handle and move it to a goal location. Thus
 102 in the Franka Kitchen domain, there are several tasks available for the MTRL agent to accomplish.
 103 Additional information on the Meta-World and Franka Kitchen domains can be found in Appendix
 104 B.2 and B.3.

105 With the use of translation layers and action heads, our goals are two-fold. The first goal is to
 106 determine if we can demonstrate skills transfer across domains. The second goal is to find the best
 107 training approach to learn a single policy π that can act in all domains, in terms of mean success rate
 108 and sample complexity.

109 3 Enabling Skill Transfer

110 In this section, we introduce our method for enabling the sharing of parameters across differing state
 111 and action spaces. In order to train using this multi-task algorithm we must first overcome the issue
 112 of differing state and action spaces.

113 3.1 Translation Layers

114 In order to overcome the difference in state spaces between domains, we propose the use of trans-
 115 lation layers. A translation layer f_{θ_i} is parameterized by parameters θ_i where i is the index of an
 116 domain. The translation layer f_{θ_i} receives the state-vector $S_t^{(i)}$ from domain i and encodes the state-
 117 vector into the shared latent space of the SAL architecture. This allows for the translation layer f_{θ_i}
 118 to process the state-vector for domain specific features before embedding the features into the shared
 119 latent space. By using these translation layers, we have overcome the issue of differing state spaces
 120 by using domain specific modules.

121 3.2 Action Heads

122 To overcome the difference in action spaces, we propose the use of a shared space to generate
 123 actions inspired by [8], and then decode these shared action representations via domain specific
 124 action heads. Once the translation layer f_{θ_i} processes the state-vector $S_t^{(i)}$ from domain i , the
 125 shared, domain-agnostic latent space $a_{sh} = \pi_\eta(f_{\theta_i}(S_t^{(i)}))$ encodes the features into a shared latent
 126 action space where a_{sh} is the output of the final parameters of the shared domain agnostic latent
 127 space π_η . Once the features are embedded into the shared action representation, the domain specific
 128 action heads decode this action into domain specific actions $g_{\phi_i}(a_{sh})$ where the action head i is
 129 parameterized by a set of parameters ϕ_i to generate action $a_t^{(i)}$.

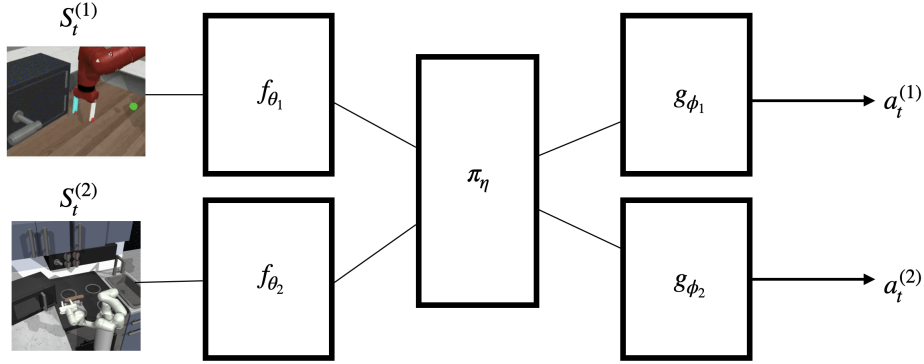


Figure 1: Control architecture for multi-task learning. The shared, domain-agnostic, latent (SAL) policy.

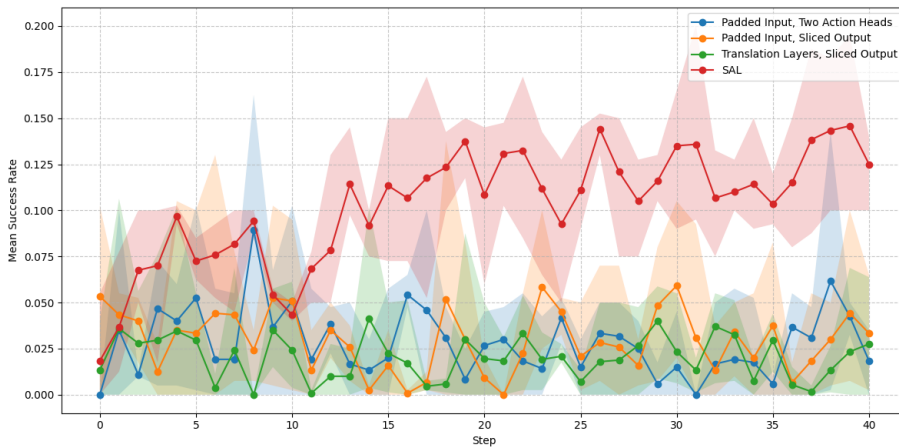


Figure 2: Mean success rate plots across all tasks from MT10 and Franka Kitchen. Blue plot is a network with padded inputs, and action heads. Orange plot is padded inputs, and outputting actions in the larger action space & only using needed action dimensions. Green plot uses translation layers and actions in the larger action space, like Orange. Red is the SAL architecture. Shaded regions indicate the minimum and maximum success rates.

130 3.3 Shared, domain-agnostic, latent (SAL) policy

131 To have an domain-agnostic policy that can learn in new domains while sharing parameters across
 132 diverse domains, we propose the SAL policy found in Figure 1. For some domain i to produce
 133 an action $a_t^{(i)}$ at time step t , a forward pass must be completed as $g_{\phi_i} \left(\pi_\eta \left(f_{\theta_i} (s_t^{(i)}) \right) \right)$. These
 134 parameters are then optimized using the loss of the RL algorithm. In this paper we use Soft Actor
 135 Critic [14] for each domain but any RL algorithm can be used. The optimizer for each domain i is
 136 operating on the translation layer f_{θ_i} , the action head g_{ϕ_i} , and the shared policy parameters π_η . The
 137 SAL policy is designed and optimized in a way that shares parameters across diverse domains in
 138 the policy layers π_η allowing for relevant similarities across domains to be encoded in the Shared,
 139 domain-Agnostic, Latent policy space.

140 4 How to train a SAL policy

141 In order to learn the dynamics of multiple domains within a shared parameter space, we propose to
 142 use the SAL policy architecture. We use the domains of Meta-World [11] and Franka Kitchen [13].
 143 Further information about these domains can be found in Appendix B.2 and B.3. The following sub-

144 sections outline the process of searching for the best performance when using the SAL architecture
145 and some of the current limitations of the method.

146 **4.1 Issues when training on multiple domains**

147 Due to the static nature of feed-forward neural network input and output dimensions, they cannot
148 handle inputs or outputs of various sizes. One approach to overcome this issue would be to pad the
149 state and action spaces with zeros to ensure that input and output dimensions are the same across
150 domains similar to the approach in [4]. In these experiments we semantically align common state
151 features, such as goal and task IDs, and pad missing dimensions with zeros. In order to overcome
152 the issue of differing action spaces, we have our policy output actions the size of the largest action
153 space, and then only use the appropriate action size for any domain that has a smaller action space.
154 In Figure 2 we find that performance of a multi-task RL algorithm using this method is extremely
155 limited. In order to learn multiple sets of dynamics in a single policy network, we propose the use of
156 SAL which aims to separate domain specific parameters into the translation layers and action heads,
157 while sharing common features of domain dynamics in the shared parameters.

158 **4.2 Training a SAL policy**

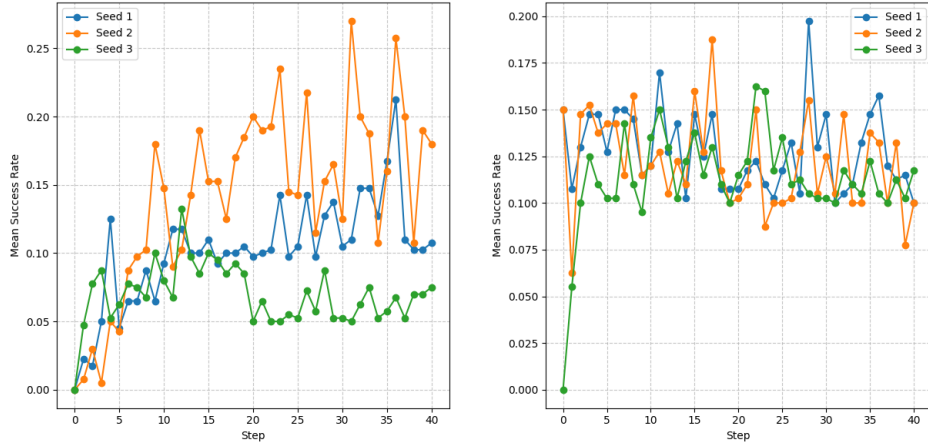
159 In Figure 2 we show the performance of SAL in reference to the padded state and action space
160 baseline approaches. The difference in performance is about 10%, showing the benefit of extracting
161 domain relevant features in the translation layers & action heads while also sharing some features in
162 the shared parameters. One of the limitations of this current iteration of SAL is that it’s difficult for
163 the SAL shared parameters π_η to align task relevant features across domains. For example if task X
164 in domain i is similar to task Y in domain j , there isn’t a mechanism to enforce similarity between
165 learned features for those tasks X and Y .

166 **4.3 Task alignment**

167 One approach to relate one task to another could be to do what we call task alignment on the task
168 IDs that are input to the translation layers. This alignment process would enable re-use of policy
169 parameters for similar tasks, leading to more sample efficient learning of similar tasks. To the best
170 of the authors knowledge, there is no method for determining this alignment a priori to training a
171 policy. Thus we explore random alignments of task IDs in order to see if there are any similar tasks
172 that can be learned across domains. In Figure 3 we show the success rates of different alignments of
173 tasks across multiple random seeds. We can see that there does seem to be some grouping in Figure
174 3a or Figure 3b, however it is not present for all tasks limiting the overall success rate of the agent
175 when attempting to align all of the tasks present.

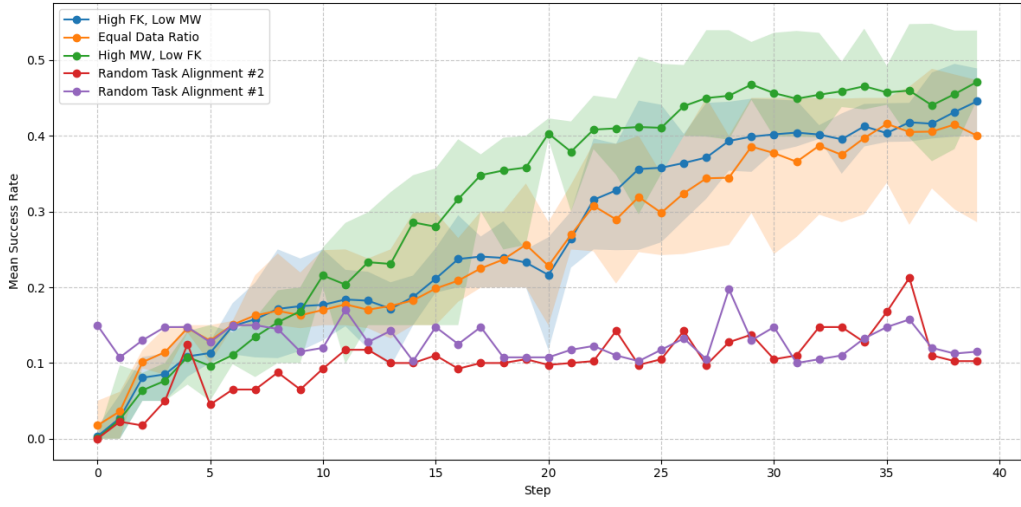
176 **4.4 Manual task alignment**

177 In order to better leverage the information gained from the random task alignments, we implement
178 what we refer to as the manual task alignment. In this alignment technique we leverage the grouping
179 information of the random task alignment experiments in Section 4.3, while also introducing some
180 human bias towards similar tasks by manually choosing tasks that are similar across domains. For
181 example, the pick place task in Meta-World is similar to the kettle task in Franka Kitchen as both
182 tasks require the agent to grasp and object and move it to a goal location. Therefore we align those
183 task IDs to have the same value. For any tasks that do not have a similar task in another domain,
184 we give these tasks unique task IDs. In Figure 3c we can see that by leveraging this manual task
185 alignment method we have a substantial increase in success rate compared to the results found in
186 Section 4.2 or 4.3. However, we also find that the sample complexity for learning tasks has wors-
187 ened using this manual task alignment compared to learning the tasks individually. In Table 1 we
188 find that by using the SAL architecture we increase the amount of samples the policy network needs
189 to be trained on to learn tasks. While the manual task alignment shows promise in improving overall
190 success rates, the increased sample complexity highlights a new challenge: balancing the learning



(a)

(b)



(c)

Figure 3: Figure 3a & 3b: success rates across 2 random task alignments, over 3 random seeds each. Additional seeds are included in Section 6. Figure 3c: mean success rate from the manual task alignment. Results indicate that SAL can learn an effective policy. Shaded regions indicate the minimum and maximum success rates.

191 across multiple domains efficiently. This trade-off between improved task generalization and in-
 192 creased sample requirements suggests that further refinement of our learning approach is necessary.
 193 To address this, we turn our attention to a more nuanced strategy for managing the learning process
 194 across multiple domains.

195 **4.5 Data manipulation**

196 While the task alignment approach in Section 4.4 showed some promise, its limitations in consis-
 197 tently grouping similar tasks across domains highlight the need for alternative strategies to enhance
 198 SAL’s learning capabilities. One key challenge in multi-domain learning is balancing the exposure
 199 to different tasks to achieve efficient and robust learning. To address this, we propose a dynamic
 200 data sampling strategy that gradually shifts the policy’s focus across domains. Initially, we create
 201 a limited data scenario where the policy uses 90% of the data from domain i and 10% of the data
 202 from domain j for its update. This imbalanced data ratio allows the policy to first build a strong
 203 foundation in one domain before gradually incorporating knowledge from the second. As training

204 progresses, we slowly adjust this ratio towards a 50/50 split, enabling the policy to transfer and inte-
205 grate knowledge across domains more effectively. This approach aims to build on the improvement
206 of the manual task alignment in Section 4.4 by providing a more controlled and gradual exposure
207 to multiple tasks. We hypothesize that SAL can develop more robust shared parameters while still
208 maintaining domain-specific adaptations through its translation layers and action heads with this
209 method.

210 Table 1 shows the results of two configurations of this data manipulation method, one where we
211 use 90% Franka Kitchen data and 10% Meta-World data in the policy update and one where we
212 use 90% Meta-World data and 10% Franka Kitchen data. These configurations are labeled High
213 FK and High MW respectively. We find that the High FK induces more sample efficient learning
214 across the Meta-World and Franka Kitchen tasks, while the High MW configuration induces some
215 more sample efficient learning. These results validate the use of the SAL architecture and the data
216 manipulation training configuration, as the goal of training simultaneously in these domains is to
217 learn with a relatively high success rate while also learning different tasks in a more sample efficient
218 manner. One of the limitations of these methods is the decrease in sample efficiency for non-aligned
219 tasks. Table 4 and Table 5 show the number of samples needed to learn all tasks in either Meta-World
220 or Franka Kitchen. For the un-aligned tasks, there is a decrease in sample efficiency.

221 4.6 SAL Latent Space Analysis

222 Finally, to examine the groupings of tasks made by training using the SAL policy architecture we
223 plot the outputs of both the translation layers f_{θ_i} and the output of shared policy layers π_{η} . We use
224 T-SNE embeddings [16] to visualize the latent vectors. In Figure 7 we show the T-SNE embeddings
225 after the first epoch and after the last epoch. Figure 7c and Figure 7d show the outputs of the
226 translation layers, while Figure 7a and Figure 7b show the outputs of the SAL policy layers. The
227 SAL architecture shows the ability to group related tasks together, however it is still uncertain how
228 to handle tasks that don't have any similar tasks in another domain.

229 5 Limitations and Future Work

230 We would like to highlight the limitations of the current approaches that were successful as they are
231 important for both the future of this work, and important for the field of multi-task reinforcement
232 learning.

233 5.1 Task Alignment

234 In Section 4.3 and Section 4.4 we outlined the methods that we used to align the task IDs of individ-
235 ual tasks across the various tasks from Meta-World and Franka Kitchen. However, these alignment
236 methods have several limitations. As the number of tasks increases, manually aligning the tasks
237 becomes more time-consuming and challenging. By allowing for a human to align the tasks, it's
238 possible that the human will introduce biases that will limit the capabilities of the MTRL algorithm.
239 The alignment method may not generalize well to new domains or the introduction of additional
240 domains during training. As noted in Section 4.4, the manual alignment method without the data
241 manipulation method increased the sample complexity, suggesting that there is a trade-off between
242 generalization and learning efficiency. Thus these limitations highlight the need for further research
243 into being able to determine how the actions from task i may transfer to task j .

244 5.2 Data Manipulation

245 In Section 4.5, we outline the method we use to learn in a more sample efficient manner where we
246 slowly introduce data from one domain during a policy update. This method may have additional
247 issues when transitioning to new domains. The strategy of using 90% data from one domain initially
248 may lead to temporarily poor performance in the under-sampled domain. The performance of the
249 system may be sensitive to how quickly or slowly the data ratio is adjusted, and finding the opti-

mal schedule could be challenging. The optimal starting ratio and adjustment schedule might vary significantly between different sets of domains, potentially requiring extensive tuning. Our current implementation focuses on balancing between two domains. Extending this to multiple domains may introduce additional complexities. In some cases, the data manipulation strategy might force the model to learn from less relevant data, potentially leading to negative transfer between tasks.

Table 1: Millions of network updates for each training procedure using the SAL architecture. Bolded results indicate a decrease in number of network updates for that task to reach 90% success rate. Single env is the results from training a single multi-task policy on Meta-World or Franka Kitchen. MW or FK beside the domain names indicate where that task is from. SAL (MTA) indicates the results of using SAL with the manual task alignment. High FK is the policy update data manipulation experiment with a ratio of 90% Franka Kitchen data and 10% Meta-World data at the start of training. High MW is the opposite.

Task	Single Env	SAL (MTA)	High FK	High MW
Drawer Close (MW)	6.4	10.7	2.4 (-62.76%)	11.5
Drawer Open (MW)	42.7	211.2	24.8 (-42.0%)	53.9
Reach (MW)	17.1	19.2	12.7 (-25.88%)	27.9
Slide Cabinet (FK)	14.9	19.2	11.5 (-23.0%)	4.3 (-71.0%)
Microwave (FK)	76.8	147.2	38.3 (-50.13%)	71.5 (-6.94%)

254

255 5.3 Future Work

256 In this work we have highlighted several successful methods for training a single policy across
 257 reinforcement learning domains with differing state and action spaces. However, we have also
 258 highlighted several limitations of our proposed method, namely the task alignment process and the
 259 data manipulation experiments.

260 In order to align the tasks for training using our SAL policy, we randomly sampled alignments
 261 and then trained policies using those alignments. As the number of tasks increases this would be
 262 an unusable method. Instead of this method, it may be possible to learn a set of options [17] and
 263 compare the trajectories induced by these options across all tasks. If an option induces similar
 264 trajectories on the task it was trained on and some new task, they may be similar tasks.

265 One of the methods that we found increased performance of the SAL policy from both a success
 266 rate and sample efficiency perspective was the manipulation of data ratios in the SAL policy update.
 267 However, it’s possible that this method would only work when the included domains are robotic
 268 manipulation tasks. A useful future work could be to explore the usefulness of training on multiple
 269 domain tasks types, such as robotic manipulation, navigation, and other continuous control domains
 270 similar to [8].

271 6 Conclusion

272 In this paper we explored the problem of applying multi-task reinforcement learning algorithms
 273 to domains with differing state and action space sizes, in addition to the state and action spaces
 274 having differing semantic meaning. We propose the Shared, Domain-Agnostic, Latent (SAL) policy
 275 to overcome these issues which leverage separate input translation layers and output action heads
 276 for each domain. We explored the usefulness of SAL and the methods that are needed to train an
 277 effective policy. We found that by applying these methods, we can increase the sample efficiency of
 278 learning in Meta-World and Franka Kitchen domains by up to 70%. We also highlighted some of
 279 the limitations and future directions for this work, including a method of determining if two tasks
 280 are similar and an investigation into training on multiple tasks with potentially different input types.
 281 This work serves as a foundation for further exploration into flexible learning architectures that can
 282 bridge the gap between diverse task domains and robot configurations.

References

- 283
- 284 [1] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
285 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou,
286 H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through
287 deep reinforcement learning. *Nature*, 518(7540):529–533, Feb. 2015. ISSN 1476-4687. doi:
288 10.1038/nature14236. URL <https://doi.org/10.1038/nature14236>.
- 289 [2] M. G. Bellemare, S. Candido, P. S. Castro, J. Gong, M. C. Machado, S. Moitra, S. S. Ponda,
290 and Z. Wang. Autonomous navigation of stratospheric balloons using reinforcement learning.
291 *Nature*, 588(7836):77–82, Dec. 2020. ISSN 1476-4687. doi:10.1038/s41586-020-2939-8.
292 URL <https://doi.org/10.1038/s41586-020-2939-8>.
- 293 [3] J. Degraeve, F. Felici, J. Buchli, M. Neunert, B. Tracey, F. Carpanese, T. Ewalds, R. Hafner,
294 A. Abdolmaleki, D. de las Casas, C. Donner, L. Fritz, C. Galperti, A. Huber, J. Keeling,
295 M. Tsimpoukelli, J. Kay, A. Merle, J.-M. Moret, S. Noury, F. Pesamosca, D. Pfau, O. Sauter,
296 C. Sommariva, S. Coda, B. Duval, A. Fasoli, P. Kohli, K. Kavukcuoglu, D. Hassabis, and
297 M. Riedmiller. Magnetic control of tokamak plasmas through deep reinforcement learning.
298 *Nature*, 602(7897):414–419, Feb. 2022. ISSN 1476-4687. doi:10.1038/s41586-021-04301-9.
299 URL <https://www.nature.com/articles/s41586-021-04301-9>. Number: 7897 Pub-
300 lisher: Nature Publishing Group.
- 301 [4] T. Chen, A. Murali, and A. Gupta. Hardware Conditioned Policies for Multi-Robot Transfer
302 Learning, Jan. 2019. URL <http://arxiv.org/abs/1811.09864>. arXiv:1811.09864 [cs].
- 303 [5] Q. Zhang, T. Xiao, A. A. Efros, L. Pinto, and X. Wang. Learning Cross-Domain Correspondence
304 for Control with Dynamics Cycle-Consistency, Dec. 2020. URL <http://arxiv.org/abs/2012.09811>.
305 arXiv:2012.09811 [cs].
- 306 [6] A. Singh, H. Liu, G. Zhou, A. Yu, N. Rhinehart, and S. Levine. PARROT: DATA-DRIVEN
307 BEHAVIORAL PRIORS FOR REINFORCEMENT LEARNING. 2021.
- 308 [7] D. Tanneberg, K. Ploeger, E. Rueckert, and J. Peters. SKID RAW: Skill Discovery from Raw
309 Trajectories. *IEEE Robotics and Automation Letters*, 6(3):4696–4703, July 2021. ISSN 2377-
310 3766, 2377-3774. doi:10.1109/LRA.2021.3068891. URL [http://arxiv.org/abs/2103.](http://arxiv.org/abs/2103.14610)
311 [14610](http://arxiv.org/abs/2103.14610). arXiv:2103.14610 [cs].
- 312 [8] D. Shah, A. Sridhar, A. Bhorkar, N. Hirose, and S. Levine. GNM: A General Navigation
313 Model to Drive Any Robot, May 2023. URL <http://arxiv.org/abs/2210.03370>.
314 arXiv:2210.03370 [cs].
- 315 [9] N. Bohlinger, G. Czechmanowski, M. Krupka, P. Kicki, K. Walas, J. Peters, and D. Tateo. One
316 Policy to Run Them All: an End-to-end Learning Approach to Multi-Embodiment Locomotion,
317 Sept. 2024. URL <http://arxiv.org/abs/2409.06366>. arXiv:2409.06366 [cs].
- 318 [10] B. Trabucco, M. Phielipp, and G. Berseth. AnyMorph: Learning Transferable Policies
319 By Inferring Agent Morphology, June 2022. URL <http://arxiv.org/abs/2206.12279>.
320 arXiv:2206.12279 [cs].
- 321 [11] T. Yu, D. Quillen, Z. He, R. Julian, A. Narayan, H. Shively, A. Bellathur, K. Hausman, C. Finn,
322 and S. Levine. Meta-World: A Benchmark and Evaluation for Multi-Task and Meta Reinforce-
323 ment Learning. *arXiv:1910.10897 [cs, stat]*, June 2021. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1910.10897)
324 [1910.10897](http://arxiv.org/abs/1910.10897). arXiv: 1910.10897.
- 325 [12] C. D’Eramo, D. Tateo, A. Bonarini, M. Restelli, and J. Peters. SHARING KNOWLEDGE IN
326 MULTI-TASK DEEP REINFORCEMENT LEARNING. 2020.
- 327 [13] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay Policy Learning: Solving
328 Long-Horizon Tasks via Imitation and Reinforcement Learning, Oct. 2019. URL <https://arxiv.org/abs/1910.11956v1>.
329

- 330 [14] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft Actor-Critic: Off-Policy Maximum
331 Entropy Deep Reinforcement Learning with a Stochastic Actor, Aug. 2018. URL [http://](http://arxiv.org/abs/1801.01290)
332 arxiv.org/abs/1801.01290. arXiv:1801.01290 [cs, stat].
- 333 [15] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 334 [16] G. E. Hinton and S. Roweis. Stochastic Neighbor Embedding. In S. Becker, S. Thrun, and
335 K. Obermayer, editors, *Advances in Neural Information Processing Systems*, volume 15. MIT
336 Press, 2002. URL [https://proceedings.neurips.cc/paper_files/paper/2002/](https://proceedings.neurips.cc/paper_files/paper/2002/file/6150ccc6069bea6b5716254057a194ef-Paper.pdf)
337 [file/6150ccc6069bea6b5716254057a194ef-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2002/file/6150ccc6069bea6b5716254057a194ef-Paper.pdf).
- 338 [17] R. S. Sutton, D. Precup, and S. Singh. Between MDPs and semi-MDPs: A framework for
339 temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–211, 1999.
340 ISSN 0004-3702. doi:[https://doi.org/10.1016/S0004-3702\(99\)00052-1](https://doi.org/10.1016/S0004-3702(99)00052-1). URL [https://www.](https://www.sciencedirect.com/science/article/pii/S0004370299000521)
341 [sciencedirect.com/science/article/pii/S0004370299000521](https://www.sciencedirect.com/science/article/pii/S0004370299000521).
- 342 [18] R. Yang, H. Xu, Y. WU, and X. Wang. Multi-Task Reinforcement Learning with Soft Mod-
343 ularization. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors,
344 *Advances in Neural Information Processing Systems*, volume 33, pages 4767–4777. Curran
345 Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2020/file/32cfdce9631d8c7906e8e9d6e68b514b-Paper.pdf)
346 [2020/file/32cfdce9631d8c7906e8e9d6e68b514b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/32cfdce9631d8c7906e8e9d6e68b514b-Paper.pdf).
- 347 [19] T. Yu, S. Kumar, A. Gupta, S. Levine, K. Hausman, and C. Finn. Gradient surgery for multi-
348 task learning. *Advances in Neural Information Processing Systems*, 33:5824–5836, 2020.
- 349 [20] M. Cho, W. Jung, and Y. Sung. Multi-task reinforcement learning with task representation
350 method. In *ICLR 2022 Workshop on Generalizable Policy Learning in Physical World*, 2022.
351 URL <https://openreview.net/forum?id=rV2zaEpNybc>.
- 352 [21] H. He, C. Bai, K. Xu, Z. Yang, W. Zhang, D. Wang, B. Zhao, and X. Li. Diffusion
353 model is an effective planner and data synthesizer for multi-task reinforcement learning.
354 *arXiv:2305.18459*, 2023.
- 355 [22] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven
356 reinforcement learning, 2020.
- 357 [23] M. Xu, M. Veloso, and S. Song. ASPiRe: Adaptive Skill Priors for Reinforcement Learning.
358 2022.
- 359 [24] M. Xu, Z. Xu, C. Chi, M. Veloso, and S. Song. XSkill: Cross Embodiment Skill Discovery.
360 2023.
- 361 [25] J. Yang, C. Glossop, A. Bhorkar, D. Shah, Q. Vuong, C. Finn, D. Sadigh, and S. Levine.
362 Pushing the Limits of Cross-Embodiment Learning for Manipulation and Navigation, Feb.
363 2024. URL <http://arxiv.org/abs/2402.19432>. arXiv:2402.19432 [cs].
- 364 [26] K. Hausman, J. T. Springenberg, Z. Wang, N. Heess, and M. Riedmiller. LEARNING AN
365 EMBEDDING SPACE FOR TRANSFERABLE ROBOT SKILLS. 2018.
- 366 [27] W. Huang, I. Mordatch, and D. Pathak. One Policy to Control Them All: Shared Modular Poli-
367 cies for Agent-Agnostic Control, July 2020. URL <http://arxiv.org/abs/2007.04976>.
368 arXiv:2007.04976 [cs, stat].
- 369 [28] A. Gupta, L. Fan, S. Ganguli, and L. Fei-Fei. MetaMorph: Learning Universal Controllers with
370 Transformers, Mar. 2022. URL <http://arxiv.org/abs/2203.11931>. arXiv:2203.11931
371 [cs].
- 372 [29] A. Gupta, C. Devin, Y. Liu, P. Abbeel, and S. Levine. Learning invariant feature spaces to
373 transfer skills with reinforcement learning. In *International Conference on Learning Repre-*
374 *sentations*, 2017. URL <https://openreview.net/forum?id=Hyq4yhile>.

375 A Related Work

376 One of the first comprehensive benchmarks for testing in the multitask RL domain was Meta-World
377 [11]. The benefit of using the Meta-World benchmark is that it has a high degree of shared do-
378 main and control structure which allows for efficient learning of distinct but related tasks [11].
379 Another benefit of Meta-World is the dense reward function available for each individual task. [11]
380 propose several benchmark algorithms for multitask reinforcement learning including Multi-Task
381 Multi-Head Soft Actor Critic (MTMHSAC). The MTMHSAC algorithm modifies the base Soft-
382 Actor Critic algorithm by adding an entropy head for each task, allowing for different levels of
383 exploration per task [11]. In Yu et al. [11] the states are augmented with a one-hot vector that in-
384 dicates which domain the state belongs to. Other recent approaches include Soft-Modularization
385 which uses the one-hot vector as input to a routing network that outputs probabilities for how data
386 is routed through the policy network [18], Yu et al. [19] developed a method of projecting con-
387 flicting gradients onto the same plane thus making network optimization more efficient, Cho et al.
388 [20] developed a variational based method as well as a measure of negative transfer. Recently, He
389 et al. [21] have shown diffusion models to be effective planners and data synthesizer in multitask
390 reinforcement learning settings.

391 Gupta et al. [13] originally designed Franka Kitchen as a benchmark for algorithms that can solve
392 long-horizon, multitask problems. Franka Kitchen was then modified by [22]. The robotic arm in
393 Franka Kitchen is a 9-DOF Franka robot that is placed in a kitchen domain with many different
394 household kitchen objects. The goal of the domain is to achieve some desired configuration of the
395 objects through manipulation [13].

396 We formulate our problem with the context of agent-agnostic reinforcement learning. There has been
397 numerous works that take different perspectives on learning policies that can be applied, or quickly
398 adapted, for running on robotics hardware with a different morphology than the policy was trained
399 on. There are two main approaches to these works. The first line of work is interested in leveraging
400 large datasets, with or without actions, and learning how to extract useful skill information from the
401 dataset. [6] uses a dataset of successful trajectories to learn a prior over robotic tasks that can then
402 be applied to unseen new tasks. [7] learns a latent variable model that is able to segment unlabelled
403 trajectories into subtasks used across all tasks in the dataset, with these learned subtasks being
404 transferrable to physical robotics manipulations tasks. Similarly, [23] learns a dictionary of skill
405 priors from expert demonstrations and then trains a policy conditioned on these skill priors enabling
406 effective long-horizon model-free reinforcement learning. XSkill [24] has been proposed for cross
407 embodiment skill discovery, where a dataset of videos across multiple embodiments are used to
408 discover skills, and then the learned skills are transferred to the current embodiment and aligned
409 with the current tasks for downstream skill-conditioned visuomotor policy learning. Lastly [25]
410 trained a single, goal-conditioned, policy across manipulation, navigation, and driving datasets with
411 various embodiments. [25] found that there was an increase in goal completion when combining
412 the manipulation and navigation data, while postulating that the policy must understand where the
413 current state is with respect to it’s goal, as well as understanding how to navigate cluttered spaces.
414 We differ from each of these previous works as we aim to do reinforcement learning without any
415 priors, and completely in an online fashion.

416 In addition to the work that learns on large datasets, there is also a line of work that aims to con-
417 dition the policy learning on a vector representation that captures task and/or robotic morphology
418 information. [26] was one of the first works to learn a representation of tasks, which allowed for
419 interpolation between learned embeddings for zero-shot transfer of skills. However, when applied to
420 Meta-World, the performance of this method degrades significantly [11]. In order to facilitate cross
421 embodiment learning, [27] proposed to learn a modular policy for each component of an agent’s
422 morphology that passes messages from module to module. [28] leverages a transformer architecture
423 to combat the inefficiencies in graph neural networks for incompatible MTRL by injecting explicit
424 morphological information into the model. [10] proposes to infer robot morphology using the re-
425 inforcement learning objective directly, instead of a representation learning objective. Recently,
426 [9] proposed a multi-task reinforcement learning method to learn across various quadreped robots.

427 Concurrently to our work, [9] make use of observation and joint specific information to train a single policy across multiple legged robots. Our work differs from [9] in the types of embodiments, 428 the tasks, and the types of observations. We don’t supply the SAL policy with information about 429 specifics of the joint it is controlling. We differ from this line of work, in general, as we don’t inject 430 or infer any embodiment information. We aim to share all overlapping task information within the 431 SAL policy network, with the translation layers and action heads handling the cardinalities of the 432 state and action spaces. 433

434 B Background information

435 B.1 Reinforcement Learning

436 In this work, the Reinforcement Learning (RL) algorithm that we use is Soft-Actor Critic (SAC). 437 SAC is an off-policy RL algorithm that is based on the maximum entropy learning framework [14]. 438 This framework leads the agent to maximize both expected rewards and entropy, which leads the 439 agent to succeed in the task while acting as randomly as possible [14]. In this work there are three 440 sets of parameters to optimize: the soft Q-function $Q_\theta(s_t, a_t)$ with Q parameterized by θ , the policy 441 parameters $\pi_\phi(a_t|s_t)$ where π is parameterized by ϕ , and the entropy penalty coefficient α [14]. The 442 objective for policy optimization is:

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim D}[\alpha \log \pi_\phi(a_t|s_t) - Q_\theta(s_t, a_t)] \quad (1)$$

443 with α controlling the entropy penalty coefficient. The coefficient α is learned using the objective:

$$J(\alpha) = \mathbb{E}_{a_t \sim \pi_\phi}[-\alpha \log \pi_\phi(a_t|s_t) - \alpha \bar{\mathcal{H}}] \quad (2)$$

444 where $\bar{\mathcal{H}}$ is the minimum target entropy.

445 Previous work has modified the SAC algorithm to include an entropy term for each of the N tasks to 446 guide exploration in each task individually, as well as to introduce replay buffers per task (Yu et al. 447 [11], Yang et al. [18]).

448 B.2 Meta-World

449 The first set of domains used in this work are from Meta-World [11]. Meta-World is a suite of 450 multitask RL and meta-RL domains that consist of a number of robotic manipulation tasks. These 451 domains are subdivided into different sets with any number of domains and different goals. This 452 work focuses on the Multi-Task 10 (MT10) set of domains. In the MT10 set, there are 10 tasks that 453 the RL agent can interact with. Some domains are closely related to each other, such as window 454 open and window close, while other domains are not as closely related to each other, such as pick 455 and place, and window close. This difference in tasks allows for a wide variety of skills to be learned 456 across these domains with robust learning happening due to the number of goals available for each 457 of the individual tasks[11]. Meta-World also provides a dense and smooth reward function to help 458 RL agents learn[11]. The different tasks of the MT10 set of domains can be found in Figure 4.

459 B.3 Franka Kitchen

460 The Franka Kitchen domain is the other domain to be used in this work. The Franka Kitchen 461 domain has typically been used in hierarchical RL where the goal is to complete a number of tasks 462 sequentially [13]. This work uses the Franka Kitchen domain in a slightly different manner where 463 each of the available tasks in Franka Kitchen are used individually to create a similar set of tasks to 464 MT10 from Meta-World. This can be verified in Figure 5. Thus, the agent only needs to solve one 465 of the available tasks in a single domain. The default reward function in Franka Kitchen is a sparse 466 reward function where the agent only receives a reward for solving the task [13].

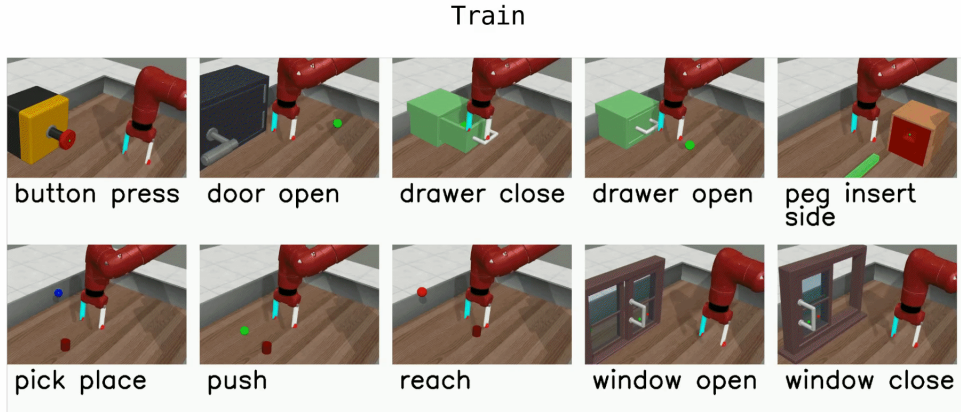


Figure 4: Meta-World MT10 tasks, image from [11].

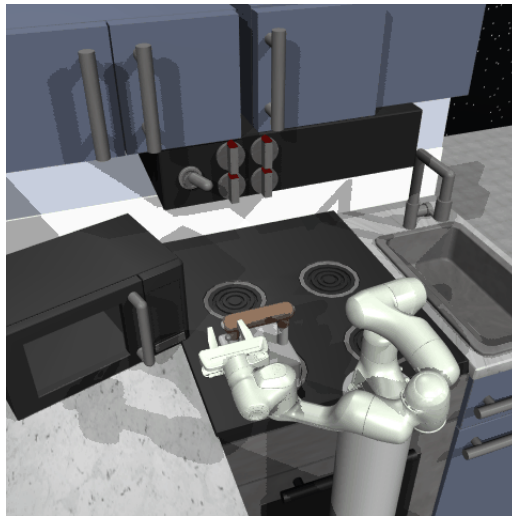


Figure 5: Franka Kitchen domain, image from [29].

467 C Evaluation procedure

468 We evaluate the performance and transfer capabilities of the RL agent in two domains: Meta World
 469 [11] and Franka Kitchen [13](see Section B.2 and B.3). In all experiments, we report the success
 470 rate across 100 episodes per task with 50,000 gradient steps each to ensure that our results are
 471 statistically significant. These results are reported across 2 different random seeds. For our multitask
 472 policy $\pi_{\theta}(a|s, t)$, the task identifier t is a one-hot encoding of the task to inform the agent what task
 473 it is solving.

474 In addition to the success rate, we also report the number of gradient update steps to our policy that
 475 it takes to learn specific tasks. In order to calculate the number of updates it takes to learn a task, we
 476 choose a success threshold of 90%, once a policy learns a task past this threshold for the first time
 477 we can then calculate the number of network gradient updates it took to learn this task. This value is
 478 calculated by the following formula: $C * GS * BS$, where C is the current epoch, GS is the number
 479 of gradient steps per epoch, and BS is the number of samples of data for this task.

Table 2: Dense reward functions used for Franka Kitchen from Meta-World.

Franka Kitchen task	Meta-World Reward Function
Microwave	Door Open
Right Hinge Door	Door Open
Left Hinge Door	Door Open
Light Switch	Sweep
Top Right Burner	Dial Turn
Top Left Burner	Dial Turn
Bottom Right Burner	Dial Turn
Bottom Left Burner	Dial Turn
Slide Cabinet	Push
Kettle	Pick Place

480 D Reward Function for Franka Kitchen

481 The Franka Kitchen domain uses a sparse reward function that gives the agent a reward of 0.3
 482 for completing the desired task and 0 otherwise. This limits the ability to do online RL as it is
 483 extremely difficult for the agent to learn the correct sequence of actions to complete any of the tasks.
 484 To overcome this challenge in creating our baseline approach, we adopt a modified version of the
 485 dense reward function that is used in Meta-World for Franka Kitchen. Table 2 shows the mapping
 486 between the Franka Kitchen task to be solved and the reward function used from Meta-World to
 487 provide dense rewards.

488 Some slight modifications are made to the reward functions to use them in Franka Kitchen. Both
 489 Meta-World and Franka Kitchen were designed using Mujoco. However, the objects that a RL
 490 agent interacts with in Meta-World are Mujoco bodies, while most of the Franka Kitchen objects
 491 are attached to bodies. The exception to this in Franka Kitchen is the kettle as it is a body itself. To
 492 overcome this issue, the reward functions for each task are modified slightly to use different Mujoco
 493 sites of the object of interest. For example, the microwave reward function uses the handle site as a
 494 method of determining how open or closed the door is. We refer the interested reader to our public
 495 implementation for further details on how the reward function was modified to use sites. To denote
 496 that a task has been completed, we used a threshold of 5 cm.

Table 3: domains that were aligned across MW and FK are shown. The Button Press Topdown, Window Close, Peg Insert Side, Reach, Drawer Open, Top Right Burner, Top Left Burner, Bottom Right Burner, and Bottom Left Burner tasks all received unique one-hot IDs.

Franka Kitchen task(s)	Meta-World task(s)
Slide Cabinet	Drawer Close
Kettle	Pick Place
Left Hinge Door, Right Hinge Door, Microwave	Door Open
Light Switch	Push

497 E Implementation and Compute Details

498 Code will be open-sourced upon acceptance.

499 F Task Alignment

500 Figure 6 shows the performance of 3 additional seeds of random task alignment, while Figure 7
 501 shows the latent vector alignment after 1 epoch of training, and then after a full experiment run.

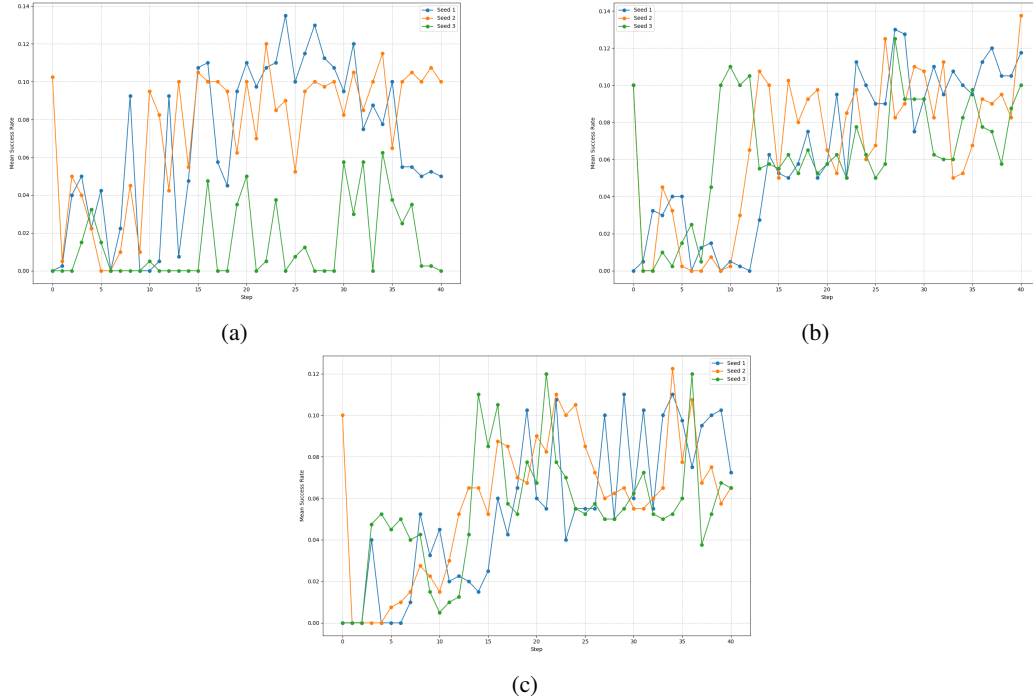


Figure 6: Results from 3 additional random seeds of alignment experiments

502 **G Network Updates Data**

Table 4: Millions of network updates for each training procedure on Meta-World tasks. This number is based on how many network updates it took to get to a 90% success rate. 0's indicate that the task was not learned in that training configuration.

Task	Single domain	FK&MW(50/50)	High FK	High MW
Door Open	32.0	168.5	46.9	69.1
Drawer Close	6.4	10.7	2.4	11.5
Drawer Open	42.7	211.2	24.8	53.9
Window Open	25.6	185.6	46.9	72.7
Window Close	14.9	155.7	43.7	64.8
Peg Insert Side	0.0	0.0	0.0	0.0
Pick Place	0.0	0.0	0.0	0.0
Push	0.0	0.0	0.0	140.8
Button Press Topdown	49.1	215.5	91.7	110.9
Reach	17.1	19.2	12.7	27.9
Total Network Updates	187.7	966.4	269.1	551.7

Figure 7: T-SNE embeddings from the translation layers (a + b) and the SAL policy (c + d). We find that the policy is able to maintain connections for tasks where there is an increase in sample efficiency.

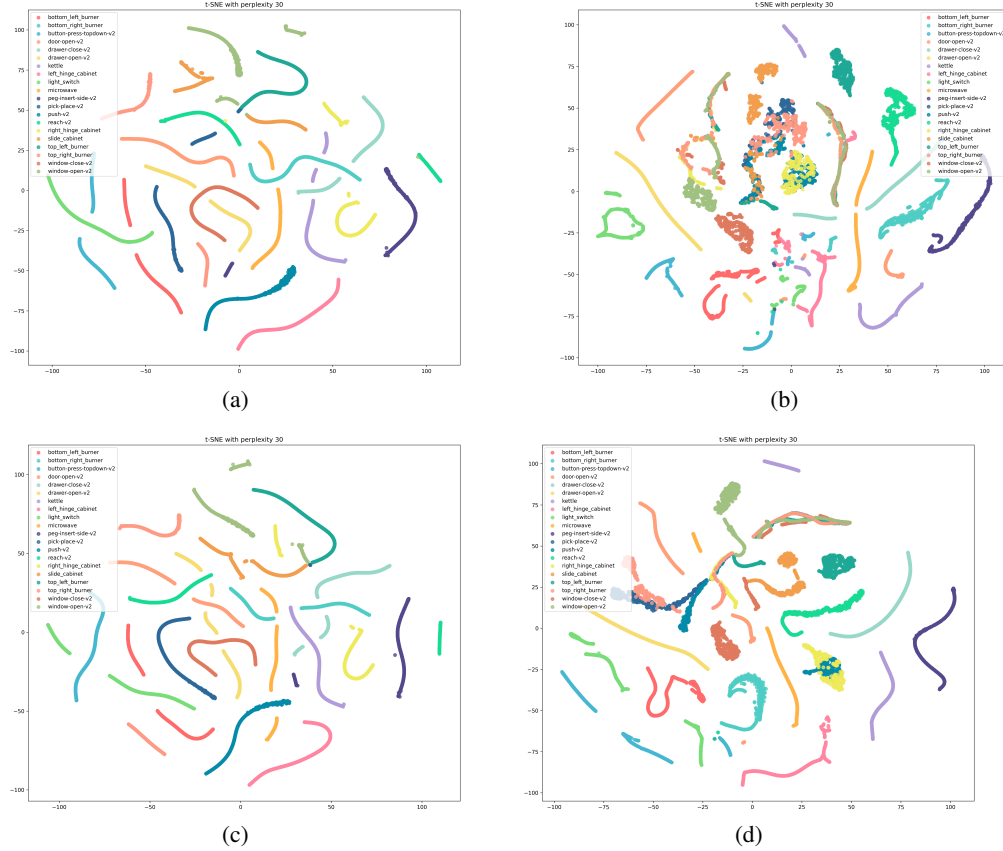


Table 5: Millions of network updates for each training procedure on Franka Kitchen tasks. This number is based on how many network updates it took to get to a 90% success rate. 0's indicate that the task was not learned in that training configuration.

Task	Single domain	FK&MW(50/50)	High FK	High MW
Slide Cabinet	14.9	19.2	11.5	4.3
Microwave	76.8	147.2	38.3	71.5
Top Right Hinge Cabinet	0.0	0.0	0.0	0.0
Top Left Hinge Cabinet	125.9	243.2	0.0	0.0
Top Right Burner	19.2	0.0	0.0	0.0
Top Left Burner	185.6	0.0	0.0	0.0
Bottom Right Burner	6.4	0.0	0.0	68.4
Bottom Left Burner	138.7	0.0	0.0	0.0
Kettle	0.0	0.0	0.0	0.0
Light Switch	14.9	125.9	110.1	41.2
Total Network Updates	582.4	535.5	159.9	185.4