

LIBERO: Benchmarking Knowledge Transfer for Lifelong Robot Learning

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Lifelong learning offers a promising paradigm of building a generalist
2 agent that learns and adapts over its lifespan. Unlike traditional lifelong learning
3 problems in image and text domains, which primarily involve the transfer of declarative
4 knowledge of entities and concepts, lifelong learning in decision-making
5 (LLDM) also necessitates the transfer of procedural knowledge, such as actions
6 and behaviors. To advance research in LLDM, we introduce LIBERO, a novel
7 benchmark of lifelong learning for robot manipulation. Specifically, LIBERO
8 highlights five key research topics in LLDM: **1)** how to efficiently transfer declarative
9 knowledge, procedural knowledge, or the mixture of both; **2)** how to design
10 effective policy architectures and **3)** effective algorithms for LLDM; **4)** the robustness of a lifelong learner with respect to task ordering; and **5)** the effect of
11 model pretraining for LLDM. We develop an extendible *procedural generation*
12 pipeline that can in principle generate infinitely many tasks. For benchmarking
13 purpose, we create four task suites (130 tasks in total) that we use to investigate the
14 above-mentioned research topics. To support sample-efficient learning, we provide
15 high-quality human-teleoperated demonstration data for all tasks. Our extensive
16 experiments present several insightful or even *unexpected* discoveries: sequential
17 finetuning outperforms existing lifelong learning methods in forward transfer, no
18 single visual encoder architecture excels at all types of knowledge transfer, and
19 naive supervised pretraining can hinder agents’ performance in the subsequent
20 LLDM.
21

22 1 Introduction

23 A longstanding goal in machine learning is to develop a generalist agent that can perform a wide
24 range of tasks. While multitask learning [1] is one approach, it is computationally demanding and
25 not adaptable to ongoing changes. Lifelong learning [2], however, offers a practical solution by
26 amortizing the learning process over the agent’s lifespan. Its goal is to leverage prior knowledge to
27 facilitate learning new tasks (forward transfer) and use the newly acquired knowledge to enhance
28 performance on prior tasks (backward transfer).

29 The main body of the lifelong learning literature has focused on how agents transfer *declarative*
30 knowledge in visual or language tasks, which pertains to *declarative knowledge* about entities and
31 concepts [3, 4]. Yet it is understudied how agents transfer knowledge in decision-making tasks,
32 which involves a mixture of both *declarative* and *procedural* knowledge (knowledge about how to *do*
33 something). Consider a scenario where a robot, initially trained to retrieve juice from a fridge, fails
34 after learning new tasks. This could be due to forgetting the juice or fridge’s location (declarative
35 knowledge) or how to open the fridge or grasp the juice (procedural knowledge). So far, we lack
36 methods to systematically and quantitatively analyze this complex knowledge transfer.

37 To bridge this research gap, this paper introduces a new simulation benchmark, Lifelong learning
38 BEchmark on RObot manipulation tasks, LIBERO, to facilitate the systematic study of lifelong
39 learning in decision making (LLDM). An ideal LLDM testbed should enable continuous learning
40 across an expanding set of diverse tasks that share concepts and actions. LIBERO supports this

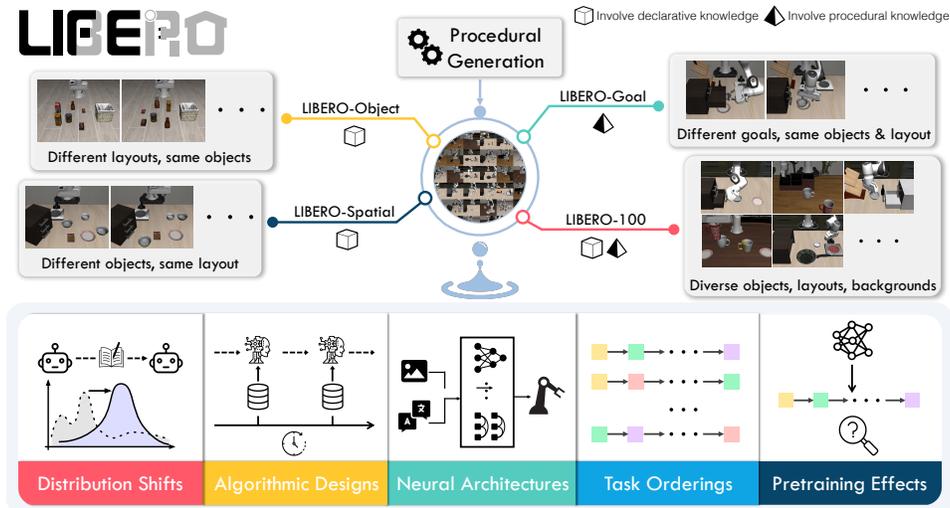


Figure 1: **Top**: LIBERO has four procedurally-generated task suites: LIBERO-SPATIAL, LIBERO-OBJECT, and LIBERO-GOAL have 10 tasks each and require transferring knowledge about spatial relationships, objects, and task goals; LIBERO-100 has 100 tasks and requires the transfer of entangled knowledge. **Bottom**: we investigate five key research topics in LLDM on LIBERO.

41 through a procedural generation pipeline for endless task creation, based on robot manipulation tasks
 42 with shared visual concepts (declarative knowledge) and interactions (procedural knowledge).

43 For benchmarking purpose, LIBERO generates 130 language-conditioned robot manipulation tasks
 44 inspired by human activities [5] and, grouped into four suites. The four task suites are designed to
 45 examine distribution shifts in the object types, the spatial arrangement of objects, the task goals,
 46 or the mixture of the previous three (top row of Figure 1). LIBERO is scalable, extendable, and
 47 designed explicitly for studying lifelong learning in robot manipulation. To support efficient learning,
 48 we provide high-quality, human-teleoperated demonstration data for all 130 tasks.

49 We present an initial study using LIBERO to investigate five major research topics in LLDM
 50 (Figure 1): **1**) knowledge transfer with different types of distribution shift; **2**) neural architecture
 51 design; **3**) lifelong learning algorithm design; **4**) robustness of the learner to task ordering; and **5**)
 52 how to leverage pre-trained models in LLDM (bottom row of Figure 1). We perform extensive
 53 experiments across different policy architectures and different lifelong learning algorithms. Based on
 54 our experiments, we make several insightful or even **unexpected** observations: (1) Policy architecture
 55 design is as crucial as lifelong learning algorithms. The transformer architecture is better at abstracting
 56 temporal information than a recurrent neural network. Vision transformers work well on tasks with
 57 rich visual information (e.g., a variety of objects). Convolution networks work well when tasks
 58 primarily need procedural knowledge. (2) While the lifelong learning algorithms we evaluated
 59 are effective at preventing forgetting, they generally perform *worse* than sequential finetuning in
 60 terms of forward transfer. (3) Our experiment shows that using pretrained language embeddings of
 61 semantically-rich task descriptions yields performance *no better* than using those of the task IDs.
 62 (4) Basic supervised pretraining on a large-scale offline dataset can have a *negative* impact on the
 63 learner’s downstream performance in LLDM.

64 2 Research Topics in LLDM

65 We outline five major research topics in LLDM that motivate the design of LIBERO and our study.

66 **(T1) Transfer of Different Types of Knowledge** In order to accomplish a task such as *put the*
 67 *ketchup next to the plate in the basket*, a robot must understand the concept *ketchup*, the location
 68 of the *plate/basket*, and how to *put* the ketchup in the basket. Indeed, robot manipulation tasks in
 69 general necessitate different types of knowledge, making it hard to determine the cause of failure.
 70 We present four task suites in Section 3.2: three task suites for studying the transfer of knowledge

71 about spatial relationships, object concepts, and task goals in a disentangled manner, and one suite
72 for studying the transfer of mixed types of knowledge.

73 **(T2) Neural Architecture Design** An important research question in LLDM is how to design
74 effective neural architectures to abstract the multi-modal observations (images, language descriptions,
75 and robot states) and transfer only relevant knowledge when learning new tasks.

76 **(T3) Lifelong Learning Algorithm Design** Given a policy architecture, it is crucial to determine
77 what learning algorithms to apply for LLDM. Specifically, the sequential nature of LLDM suggests
78 that even minor forgetting over successive steps can potentially lead to a total failure in execution. As
79 such, we consider the design of lifelong learning algorithms to be an open area of research in LLDM.

80 **(T4) Robustness to Task Ordering** It is well-known that task curriculum influences policy
81 learning [6, 7]. A robot in the real world, however, often cannot choose which task to encounter first.
82 Therefore, a good lifelong learning algorithm should be robust to different task orderings.

83 **(T5) Usage of Pretrained Models** In practice, robots will be most likely pretrained on large
84 datasets in factories before deployment [8]. However, it is not well-understood whether or how
85 pretraining could benefit subsequent LLDM.

86 **3 LIBERO**

87 **3.1 Procedural Generation of Tasks**

88 Research in LLDM requires a systematic way to create new tasks while maintaining task diversity
89 and relevance to existing tasks. LIBERO procedurally generates new tasks in three steps: **1)** extract
90 behavioral templates from language annotations of human activities and generate sampled tasks
91 described in natural language based on such templates; **2)** specify an initial object distribution given a
92 task description; and **3)** specify task goals using a propositional formula that aligns with the language
93 instructions. Our generation pipeline is built on top of Robosuite [9], a modular manipulation
94 simulator that offers seamless integration. Figure 2 illustrates an example of task creation using this
95 pipeline, and each component is expanded upon below.

96 **Behavioral Templates and Instruction Generation** Human activities serve as a fertile source of
97 tasks that can inspire and generate a vast number of manipulation tasks. We choose a large-scale
98 activity dataset, Ego4D [5], which includes a large variety of everyday activities with language
99 annotations. We pre-process the dataset by extracting the language descriptions and then summarize
100 them into a large set of commonly used language templates. After this pre-processing step, we use
101 the templates and select objects available in the simulator to generate a set of task descriptions in the
102 form of language instructions. For example, we can generate an instruction “Open the drawer of the
103 cabinet” from the template “Open ...”.

104 **Initial State Distribution (μ_0)** To specify μ_0 , we first sample a scene layout that matches the
105 objects/behaviors in a provided instruction. For instance, a kitchen scene is selected for an instruction
106 *Open the top drawer of the cabinet and put the bowl in it*. Then, the details about μ_0 are generated in
107 the PDDL language [10, 11]. Concretely, μ_0 contains information about object categories and their
108 placement (Figure 2-(A)), and their initial status (Figure 2-(B)).

109 **Goal Specifications (g)** Based on μ_0 and the language instruction, we specify the task goal using
110 a conjunction of predicates. Predicates include *unary predicates* that describe the properties of an
111 object, such as `Open(X)` or `TurnOff(X)`, and *binary predicates* that describe spatial relations between
112 objects, such as `On(A, B)` or `In(A, B)`. An example of the goal specification using PDDL language
113 can be found in Figure 2-(C). The simulation terminates when all predicates are verified true.

114 **3.2 Task Suites**

115 While the pipeline in Section 3.1 supports the generation of an unlimited number of tasks, we offer
116 fixed sets of tasks for benchmarking purposes. LIBERO has four task suites: LIBERO-SPATIAL,
117 LIBERO-OBJECT, LIBERO-GOAL, and LIBERO-100. The first three task suites are curated to
118 disentangle the transfer of *declarative* and *procedural* knowledge (as mentioned in (T1)), while
119 LIBERO-100 is a suite of 100 tasks with entangled knowledge transfer.

120 **LIBERO-X** LIBERO-SPATIAL, LIBERO-OBJECT, and LIBERO-GOAL all have 10 tasks¹
121 and are designed to investigate the controlled transfer of knowledge about spatial information
122 (declarative), objects (declarative), and task goals (procedural). Specifically, all tasks in LIBERO-
123 SPATIAL request the robot to place a bowl, among the same set of objects, on a plate. But there are
124 two identical bowls that differ only in their location or spatial relationship to other objects. Hence, to
125 successfully complete LIBERO-SPATIAL, the robot needs to continually learn and memorize new
126 spatial relationships. All tasks in LIBERO-OBJECT request the robot to pick-place a unique object.
127 Hence, to accomplish LIBERO-OBJECT, the robot needs to continually learn and memorize new
128 object types. All tasks in LIBERO-GOAL share the same objects with fixed spatial relationships but
129 differ only in the task goal. Hence, to accomplish LIBERO-GOAL, the robot needs to continually
130 learn new knowledge about motions and behaviors. More details are in Appendix G.

131 **LIBERO-100** LIBERO-100 contains 100 tasks that entail diverse object interactions and versatile
132 motor skills. In this paper, we split LIBERO-100 into 90 short-horizon tasks (LIBERO-90) and 10
133 long-horizon tasks (LIBERO-LONG). LIBERO-90 serves as the data source for pretraining (**T5**)
134 and LIBERO-LONG for downstream evaluation of lifelong learning algorithms.

135 3.3 Lifelong Learning Algorithms

136 We implement three representative lifelong learning algorithms to facilitate research in algorithmic
137 design for LLDM. Specifically, we implement Experience Replay (ER) [12], Elastic Weight
138 Consolidation (EWC) [13], and PACKNET [14]. We pick ER, EWC, and PACKNET because they
139 correspond to the memory-based, regularization-based, and dynamic-architecture-based methods for
140 lifelong learning. In addition, prior research [15] has discovered that they are state-of-the-art methods.
141 Besides these three methods, we also implement sequential finetuning (SQL) and multitask learning
142 (MTL), which serve as a lower bound and upper bound for lifelong learning algorithms, respectively.
143 More details about the algorithms are in Appendix F.1.

144 4 Experiments

145 Experiments are conducted as an initial study for the five research topics mentioned in Section 2.
146 Specifically, we focus on addressing the following research questions: **Q1**: How do different
147 architectures/LL algorithms perform under specific distribution shifts? **Q2**: To what extent does
148 neural architecture impact knowledge transfer in LLDM, and are there any discernible patterns in
149 the specialized capabilities of each architecture? **Q3**: How do existing algorithms from lifelong
150 supervised learning perform on LLDM tasks? **Q4**: To what extent does language embedding affect
151 knowledge transfer in LLDM? **Q5**: How robust are different LL algorithms to task ordering in
152 LLDM? **Q6**: Can supervised pretraining improve downstream lifelong learning performance in
153 LLDM? The detailed results/findings are in Appendix A.

154 5 Conclusion and Limitations

155 This paper introduces LIBERO, a new benchmark in the robot manipulation domain for supporting
156 research in LLDM. LIBERO includes a procedural generation pipeline that can create an infinite
157 number of manipulation tasks in the simulator. We use this pipeline to create 130 standardized tasks
158 and conduct a comprehensive set of experiments on policy and algorithm designs. The empirical
159 results suggest several future research directions: 1) how to design a better neural architecture to better
160 process spatial information or temporal information; 2) how to design a better algorithm to improve
161 forward transfer ability; and 3) how to use pretraining to help improve lifelong learning performance.
162 In the short term, we do not envision any negative societal impacts triggered by LIBERO. But
163 as the lifelong learner mainly learns from humans, studying how to preserve user privacy within
164 LLDM [16] is crucial in the long run.

¹A suite of 10 tasks is enough to observe catastrophic forgetting while maintaining computation efficiency.

References

- 165 [1] R. Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- 166 [2] S. Thrun and T. M. Mitchell. Lifelong robot learning. *Robotics and autonomous systems*, 15
167 (1-2):25–46, 1995.
- 168 [3] M. Biesialska, K. Biesialska, and M. R. Costa-Jussa. Continual lifelong learning in natural
169 language processing: A survey. *arXiv preprint arXiv:2012.09823*, 2020.
- 170 [4] Z. Mai, R. Li, J. Jeong, D. Quispe, H. Kim, and S. Sanner. Online continual learning in image
171 classification: An empirical survey. *Neurocomputing*, 469:28–51, 2022.
- 172 [5] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang,
173 M. Liu, X. Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In
174 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
175 18995–19012, 2022.
- 176 [6] Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum learning. In *Proceedings of
177 the 26th annual international conference on machine learning*, pages 41–48, 2009.
- 178 [7] S. Narvekar, B. Peng, M. Leonetti, J. Sinapov, M. E. Taylor, and P. Stone. Curriculum learning
179 for reinforcement learning domains: A framework and survey. *arXiv preprint arXiv:2003.04960*,
180 2020.
- 181 [8] L. P. Kaelbling. The foundation of efficient robot learning. *Science*, 369(6506):915–916, 2020.
- 182 [9] Y. Zhu, J. Wong, A. Mandlekar, and R. Martín-Martín. robosuite: A modular simulation
183 framework and benchmark for robot learning. *arXiv preprint arXiv:2009.12293*, 2020.
- 184 [10] D. McDermott, M. Ghallab, A. Howe, C. Knoblock, A. Ram, M. Veloso, D. Weld, and
185 D. Wilkins. Pddl-the planning domain definition language. 1998.
- 186 [11] S. Srivastava, C. Li, M. Lingelbach, R. Martín-Martín, F. Xia, K. E. Vainio, Z. Lian, C. Gokmen,
187 S. Buch, K. Liu, et al. Behavior: Benchmark for everyday household activities in virtual,
188 interactive, and ecological environments. In *Conference on Robot Learning*, pages 477–490.
189 PMLR, 2022.
- 190 [12] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. Torr, and M. Ran-
191 zato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*,
192 2019.
- 193 [13] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan,
194 J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in
195 neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- 196 [14] A. Mallya and S. Lazebnik. Packnet: Adding multiple tasks to a single network by iterative
197 pruning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*,
198 pages 7765–7773, 2018.
- 199 [15] M. Wołczyk, M. Zajkac, R. Pascanu, L. Kuciński, and P. Miloš. Continual world: A robotic
200 benchmark for continual reinforcement learning. In *Neural Information Processing Systems*,
201 2021.
- 202 [16] B. Liu, Q. Liu, and P. Stone. Continual learning and private unlearning. In *CoLLAs*, 2022.
- 203 [17] N. Díaz-Rodríguez, V. Lomonaco, D. Filliat, and D. Maltoni. Don’t forget, there is more than
204 forgetting: new metrics for continual learning. *arXiv preprint arXiv:1810.13166*, 2018.
- 205

- 206 [18] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese,
207 Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations for
208 robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.
- 209 [19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
210 transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 211 [20] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
212 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision.
213 In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 214 [21] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are
215 unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 216 [22] L. Deng. The mnist database of handwritten digit images for machine learning research. *IEEE*
217 *Signal Processing Magazine*, 29(6):141–142, 2012.
- 218 [23] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 219 [24] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical
220 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages
221 248–255. Ieee, 2009.
- 222 [25] V. Lomonaco and D. Maltoni. Core50: a new dataset and benchmark for continuous object
223 recognition. In *Conference on Robot Learning*, pages 17–26. PMLR, 2017.
- 224 [26] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman. Glue: A multi-
225 task benchmark and analysis platform for natural language understanding. *arXiv preprint*
226 *arXiv:1804.07461*, 2018.
- 227 [27] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich. Superglue: Learning feature
228 matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer*
229 *vision and pattern recognition*, pages 4938–4947, 2020.
- 230 [28] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller.
231 Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- 232 [29] O. E. L. Team, A. Stooke, A. Mahajan, C. Barros, C. Deck, J. Bauer, J. Sygnowski, M. Trebacz,
233 M. Jaderberg, M. Mathieu, et al. Open-ended learning leads to generally capable agents. *arXiv*
234 *preprint arXiv:2107.12808*, 2021.
- 235 [30] M. Kempka, M. Wydmuch, G. Runc, J. Toczek, and W. Jaśkowski. Vizdoom: A doom-based ai
236 research platform for visual reinforcement learning. In *2016 IEEE conference on computational*
237 *intelligence and games (CIG)*, pages 1–8. IEEE, 2016.
- 238 [31] S. Powers, E. Xing, E. Kolve, R. Mottaghi, and A. Gupta. Cora: Benchmarks, baselines, and met-
239 rics as a platform for continual reinforcement learning agents. *arXiv preprint arXiv:2110.10067*,
240 2021.
- 241 [32] K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging procedural generation to benchmark
242 reinforcement learning. In *International conference on machine learning*, pages 2048–2056.
243 PMLR, 2020.
- 244 [33] M. Samvelyan, R. Kirk, V. Kurin, J. Parker-Holder, M. Jiang, E. Hambro, F. Petroni, H. Küt-
245 tler, E. Grefenstette, and T. Rocktäschel. Minihack the planet: A sandbox for open-ended
246 reinforcement learning research. *arXiv preprint arXiv:2109.13202*, 2021.
- 247 [34] M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and
248 D. Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In
249 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
250 10740–10749, 2020.

- 251 [35] A. Ayub and A. R. Wagner. F-siol-310: A robotic dataset and benchmark for few-shot incre-
 252 mental object learning. In *2021 IEEE International Conference on Robotics and Automation*
 253 (*ICRA*), pages 13496–13502. IEEE, 2021.
- 254 [36] Q. She, F. Feng, X. Hao, Q. Yang, C. Lan, V. Lomonaco, X. Shi, Z. Wang, Y. Guo, Y. Zhang,
 255 et al. Openloris-object: A robotic vision dataset and benchmark for lifelong deep learning. In
 256 *2020 IEEE international conference on robotics and automation (ICRA)*, pages 4767–4773.
 257 IEEE, 2020.
- 258 [37] S. I. Mirzadeh, A. Chaudhry, D. Yin, T. Nguyen, R. Pascanu, D. Gorur, and M. Farajtabar.
 259 Architecture matters in continual learning. *arXiv preprint arXiv:2202.00275*, 2022.
- 260 [38] M. Wołczyk, M. Zajkac, R. Pascanu, L. Kuciński, and P. Miloš. Disentangling transfer in
 261 continual reinforcement learning. *ArXiv*, abs/2209.13900, 2022.
- 262 [39] B. Ermiš, G. Zappella, M. Wistuba, and C. Archambeau. Memory efficient continual learning
 263 with transformers. 2022.
- 264 [40] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A
 265 benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on*
 266 *robot learning*, pages 1094–1100. PMLR, 2020.
- 267 [41] O. Ahmed, F. Träuble, A. Goyal, A. Neitz, Y. Bengio, B. Schölkopf, M. Wüthrich, and S. Bauer.
 268 Causalworld: A robotic manipulation benchmark for causal structure and transfer learning.
 269 *arXiv preprint arXiv:2010.04296*, 2020.
- 270 [42] S. James, Z. Ma, D. R. Arrojo, and A. J. Davison. Rlbench: The robot learning benchmark &
 271 learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020.
- 272 [43] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine,
 273 M. Lingelbach, J. Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday
 274 activities and realistic simulation. In *Conference on Robot Learning*, pages 80–93. PMLR,
 275 2023.
- 276 [44] T. Mu, Z. Ling, F. Xiang, D. Yang, X. Li, S. Tao, Z. Huang, Z. Jia, and H. Su. Maniskill:
 277 Generalizable manipulation skill benchmark with large-scale demonstrations. *arXiv preprint*
 278 *arXiv:2107.14483*, 2021.
- 279 [45] J. Gu, F. Xiang, X. Li, Z. Ling, X. Liu, T. Mu, Y. Tang, S. Tao, X. Wei, Y. Yao, et al. Maniskill2:
 280 A unified benchmark for generalizable manipulation skills. *arXiv preprint arXiv:2302.04659*,
 281 2023.
- 282 [46] J. A. Mendez, M. Hussing, M. Gummadi, and E. Eaton. Composuite: A compositional
 283 reinforcement learning benchmark. *arXiv preprint arXiv:2207.04136*, 2022.
- 284 [47] M. Hausknecht and P. Stone. Deep recurrent q-learning for partially observable mdps. In *2015*
 285 *aaai fall symposium series*, 2015.
- 286 [48] M. Bain and C. Sammut. A framework for behavioural cloning. In *Machine Intelligence 15*,
 287 pages 103–129, 1995.
- 288 [49] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with a
 289 general conditioning layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 290 volume 32, 2018.
- 291 [50] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Imitation learning for vision-based manipulation
 292 with object proposal priors. *arXiv preprint arXiv:2210.11339*, 2022.

- 293 [51] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
294 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
295 30, 2017.
- 296 [52] W. Kim, B. Son, and I. Kim. Vilt: Vision-and-language transformer without convolution
297 or region supervision. In *International Conference on Machine Learning*, pages 5583–5594.
298 PMLR, 2021.
- 299 [53] C. M. Bishop. Mixture density networks. 1994.
- 300 [54] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay:
301 Long-horizon imitation learning by watching human play. *arXiv preprint arXiv:2302.12422*,
302 2023.
- 303 [55] A. Ayub and C. Fendley. Few-shot continual active learning by a robot. *arXiv preprint*
304 *arXiv:2210.04137*, 2022.
- 305 [56] M. Kang, J. Park, and B. Han. Class-incremental learning by knowledge distillation with
306 adaptive feature consolidation. In *Proceedings of the IEEE/CVF conference on computer vision*
307 *and pattern recognition*, pages 16071–16080, 2022.
- 308 [57] A. Rios and L. Itti. Lifelong learning without a task oracle. In *2020 IEEE 32nd International*
309 *Conference on Tools with Artificial Intelligence (ICTAI)*, pages 255–263. IEEE, 2020.
- 310 [58] D.-W. Zhou, F.-Y. Wang, H.-J. Ye, L. Ma, S. Pu, and D.-C. Zhan. Forward compatible few-shot
311 class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
312 *and Pattern Recognition*, pages 9046–9056, 2022.
- 313 [59] G. Saha, I. Garg, A. Ankit, and K. Roy. Space: Structured compression and sharing of
314 representational space for continual learning. *IEEE Access*, 9:150480–150494, 2021.
- 315 [60] B. Cheung, A. Terekhov, Y. Chen, P. Agrawal, and B. Olshausen. Superposition of many models
316 into one. *Advances in neural information processing systems*, 32, 2019.
- 317 [61] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu,
318 R. Pascanu, and R. Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*,
319 2016.
- 320 [62] J. Yoon, E. Yang, J. Lee, and S. J. Hwang. Lifelong learning with dynamically expandable
321 networks. *arXiv preprint arXiv:1708.01547*, 2017.
- 322 [63] C.-Y. Hung, C.-H. Tu, C.-E. Wu, C.-H. Chen, Y.-M. Chan, and C.-S. Chen. Compacting, picking
323 and growing for unforgetting continual learning. *Advances in Neural Information Processing*
324 *Systems*, 32, 2019.
- 325 [64] L. Wu, B. Liu, P. Stone, and Q. Liu. Firefly neural architecture descent: a general approach
326 for growing neural networks. *Advances in Neural Information Processing Systems*, 33:22373–
327 22383, 2020.
- 328 [65] E. Ben-Iwhiwhu, S. Nath, P. K. Pilly, S. Kolouri, and A. Soltoggio. Lifelong reinforcement
329 learning with modulating masks. *arXiv preprint arXiv:2212.11110*, 2022.
- 330 [66] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. Torr. Riemannian walk for incremental learn-
331 ing: Understanding forgetting and intransigence. In *Proceedings of the European Conference*
332 *on Computer Vision (ECCV)*, pages 532–547, 2018.
- 333 [67] J. Schwarz, W. Czarnecki, J. Luketina, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and
334 R. Hadsell. Progress & compress: A scalable framework for continual learning. In *International*
335 *Conference on Machine Learning*, pages 4528–4537. PMLR, 2018.

- 336 [68] H. Liu and H. Liu. Continual learning with recursive gradient optimization. *arXiv preprint*
337 *arXiv:2201.12522*, 2022.
- 338 [69] D. Lopez-Paz and M. Ranzato. Gradient episodic memory for continual learning. *Advances in*
339 *neural information processing systems*, 30, 2017.
- 340 [70] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny. Efficient lifelong learning with
341 a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- 342 [71] P. Buzzega, M. Boschini, A. Porrello, D. Abati, and S. Calderara. Dark experience for general
343 continual learning: a strong, simple baseline. *Advances in neural information processing*
344 *systems*, 33:15920–15930, 2020.
- 345 [72] M. De Lange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and
346 T. Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE*
347 *transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021.
- 348 [73] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter. Continual lifelong learning with
349 neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- 350 [74] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint*
351 *arXiv:1412.6980*, 2014.
- 352 [75] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction
353 to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial*
354 *intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings,
355 2011.
- 356 [76] S. Greydanus, A. Koul, J. Dodge, and A. Fern. Visualizing and understanding atari agents.
357 *ArXiv*, abs/1711.00138, 2017.

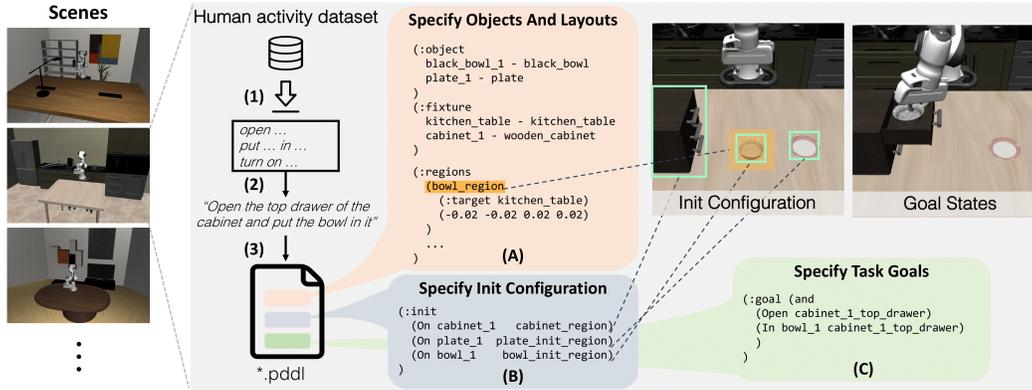


Figure 2: LIBERO’s procedural generation pipeline: Extracting behavioral templates from a large-scale human activity dataset (1), Ego4D, for generating task instructions (2); Based on the task description, selecting the scene and generating the PDDL description file (3) that specifies the objects and layouts (A), the initial object configurations (B), and the task goal (C).

358 A Experiments

359 Experiments are conducted as an initial study for the five research topics mentioned in Section 2. We
 360 first introduce the evaluation metric used in experiments, and present analysis of empirical results in
 361 LIBERO. The detailed experimental setup is in Appendix H and the study on Q5 is in Appendix I.2.
 362 Our experiments focus on addressing the following research questions:

- 363 **Q1:** How do different architectures/LL algorithms perform under specific distribution shifts?
 364 **Q2:** To what extent does neural architecture impact knowledge transfer in LLDM, and are there any
 365 discernible patterns in the specialized capabilities of each architecture?
 366 **Q3:** How do existing algorithms from lifelong supervised learning perform on LLDM tasks?
 367 **Q4:** To what extent does language embedding affect knowledge transfer in LLDM?
 368 **Q5:** How robust are different LL algorithms to task ordering in LLDM?
 369 **Q6:** Can supervised pretraining improve downstream lifelong learning performance in LLDM?

370 A.1 Evaluation Metrics

371 We report three metrics: FWT (forward transfer) [17], NBT (negative backward transfer), and
 372 AUC (area under the success rate curve). All metrics are computed in terms of success rate, as
 373 previous literature has shown that the success rate is a more reliable metric than training loss for
 374 manipulation policies [18] (Detailed explanation in Appendix I.3). Lower NBT means a policy
 375 has better performance in the previously seen tasks, higher FWT means a policy learns faster on a
 376 new task, and higher AUC means an overall better performance considering both NBT and FWT.
 377 Specifically, denote $c_{i,j,e}$ as the agent’s success rate on task j when it learned over $i-1$ previous tasks
 378 and has just learned e epochs ($e \in \{0, 5, \dots, 50\}$) on task i . Let $c_{i,i}$ be the best success rate over all
 379 evaluated epochs e for the current task i (i.e., $c_{i,i} = \max_e c_{i,i,e}$). Then, we find the earliest epoch e_i^*
 380 in which the agent achieves the best performance on task i (i.e., $e_i^* = \arg \min_e c_{i,i,e} = c_{i,i}$), and
 381 assume for all $e \geq e_i^*$, $c_{i,i,e} = c_{i,i}$.² Given a different task $j \neq i$, we define $c_{i,j} = c_{i,j,e_i^*}$. Then
 382 the three metrics are defined: $\text{FWT} = \sum_{k \in [K]} \frac{\text{FWT}_k}{K}$, $\text{FWT}_k = \frac{1}{11} \sum_{e \in \{0 \dots 50\}} c_{k,k,e}$, $\text{NBT} =$
 383 $\sum_{k \in [K]} \frac{\text{NBT}_k}{K}$, $\text{NBT}_k = \frac{1}{K-k} \sum_{\tau=k+1}^K (c_{k,k} - c_{\tau,k})$, and $\text{AUC} = \sum_{k \in [K]} \frac{\text{AUC}_k}{K}$, $\text{AUC}_k =$
 384 $\frac{1}{K-k+1} (\text{FWT}_k + \sum_{\tau=k+1}^K c_{\tau,k})$. A visualization of these metrics is provided in Figure 4.

385 A.2 Experimental Results

386 We present empirical results to address the research questions. Please refer to Appendix I.1 for the
 387 full results across all algorithms, policy architectures, and task suites.

²In practice, it’s possible that the agent’s performance on task i is not monotonically increasing due to the variance of learning. But we keep the best checkpoint among those saved at epochs $\{e\}$ as if the agent stops learning after e_i^* .

388 **Study on the Policy’s Neural Architectures (Q1, Q2)** Table 1 reports the agent’s lifelong learning
 389 performance using the three different neural architectures on the four task suites. Results are reported
 390 when ER and PACKNET are used as they demonstrate the best lifelong learning performance across
 all task suites.

| Policy Arch. | ER | | | PACKNET | | |
|----------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| | FWT(↑) | NBT(↓) | AUC(↑) | FWT(↑) | NBT(↓) | AUC(↑) |
| LIBERO-LONG | | | | | | |
| RESNET-RNN | 0.16 ± 0.02 | 0.16 ± 0.02 | 0.08 ± 0.01 | 0.13 ± 0.00 | 0.21 ± 0.01 | 0.03 ± 0.00 |
| RESNET-T | 0.48 ± 0.02 | 0.32 ± 0.04 | 0.32 ± 0.01 | 0.22 ± 0.01 | 0.08 ± 0.01 | 0.25 ± 0.00 |
| ViT-T | 0.38 ± 0.05 | 0.29 ± 0.06 | 0.25 ± 0.02 | 0.36 ± 0.01 | 0.14 ± 0.01 | 0.34 ± 0.01 |
| LIBERO-SPATIAL | | | | | | |
| RESNET-RNN | 0.40 ± 0.02 | 0.29 ± 0.02 | 0.29 ± 0.01 | 0.27 ± 0.03 | 0.38 ± 0.03 | 0.06 ± 0.01 |
| RESNET-T | 0.65 ± 0.03 | 0.27 ± 0.03 | 0.56 ± 0.01 | 0.55 ± 0.01 | 0.07 ± 0.02 | 0.63 ± 0.00 |
| ViT-T | 0.63 ± 0.01 | 0.29 ± 0.02 | 0.50 ± 0.02 | 0.57 ± 0.04 | 0.15 ± 0.00 | 0.59 ± 0.03 |
| LIBERO-OBJECT | | | | | | |
| RESNET-RNN | 0.30 ± 0.01 | 0.27 ± 0.05 | 0.17 ± 0.05 | 0.29 ± 0.02 | 0.35 ± 0.02 | 0.13 ± 0.01 |
| RESNET-T | 0.67 ± 0.07 | 0.43 ± 0.04 | 0.44 ± 0.06 | 0.60 ± 0.07 | 0.17 ± 0.05 | 0.60 ± 0.05 |
| ViT-T | 0.70 ± 0.02 | 0.28 ± 0.01 | 0.57 ± 0.01 | 0.58 ± 0.03 | 0.18 ± 0.02 | 0.56 ± 0.04 |
| LIBERO-GOAL | | | | | | |
| RESNET-RNN | 0.41 ± 0.00 | 0.35 ± 0.01 | 0.26 ± 0.01 | 0.32 ± 0.03 | 0.37 ± 0.04 | 0.11 ± 0.01 |
| RESNET-T | 0.64 ± 0.01 | 0.34 ± 0.02 | 0.49 ± 0.02 | 0.63 ± 0.02 | 0.06 ± 0.01 | 0.75 ± 0.01 |
| ViT-T | 0.57 ± 0.00 | 0.40 ± 0.02 | 0.38 ± 0.01 | 0.69 ± 0.02 | 0.08 ± 0.01 | 0.76 ± 0.02 |

Table 1: Performance of the three neural architectures using ER and PACKNET on the four task suites. Results are averaged over three seeds and we report the mean and standard error. The best performance is **bolded**, and colored in **purple** if the improvement is statistically significant over other neural architectures, when a two-tailed, Student’s t-test under equal sample sizes and unequal variance is applied with a p -value of 0.05.

392 *Findings:* First, we observe that RESNET-T and ViT-T work much better than RESNET-RNN on
 393 average, indicating that using a transformer on the “temporal” level could be a better option than
 394 using an RNN model. Second, the performance difference among different architectures depends
 395 on the underlying lifelong learning algorithm. If PACKNET (a dynamic architecture approach) is
 396 used, we observe no significant performance difference between RESNET-T and ViT-T except on
 397 the LIBERO-LONG task suite where ViT-T performs much better than RESNET-T. In contrast,
 398 if ER is used, we observe that RESNET-T performs better than ViT-T on all task suites except
 399 LIBERO-OBJECT. This potentially indicates that the ViT architecture is better at processing visual
 400 information with more object varieties than the ResNet architecture when the network capacity is
 401 sufficiently large (See the MTL results in Table 8 on LIBERO-OBJECT as the supporting evidence).
 402 The above findings shed light on how one can improve architecture design for better processing of
 403 spatial and temporal information in LLDM.

404 **Study on Lifelong Learning Algorithms (Q1, Q3)** Table 2 reports the lifelong learning per-
 405 formance of the three lifelong learning algorithms, together with the SEQL and MTL baselines.
 406 All experiments use the same RESNET-T architecture as it performs the best across all policy
 407 architectures.

408 *Findings:* We observed a series of interesting findings that could potentially benefit future research
 409 on algorithm design for LLDM: 1) SEQL shows the best FWT over all task suites. This is surprising
 410 since it indicates all lifelong learning algorithms we consider actually hurt forward transfer; 2)
 411 PACKNET outperforms other lifelong learning algorithms on LIBERO-X but is outperformed by
 412 ER significantly on LIBERO-LONG, mainly because of low forward transfer. This confirms that
 413 the dynamic architecture approach is good at preventing forgetting. But since PACKNET splits the
 414 network into different sub-networks, the essential capacity of the network for learning any individual
 415 task is smaller. Therefore, we conjecture that PACKNET is not rich enough to learn on LIBERO-

| Lifelong Algo. | FWT(\uparrow) | NBT(\downarrow) | AUC(\uparrow) | FWT(\uparrow) | NBT(\downarrow) | AUC(\uparrow) |
|----------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| | LIBERO-LONG | | | LIBERO-SPATIAL | | |
| SEQL | 0.54 \pm 0.01 | 0.63 \pm 0.01 | 0.15 \pm 0.00 | 0.72 \pm 0.01 | 0.81 \pm 0.01 | 0.20 \pm 0.01 |
| ER | 0.48 \pm 0.02 | 0.32 \pm 0.04 | 0.32 \pm 0.01 | 0.65 \pm 0.03 | 0.27 \pm 0.03 | 0.56 \pm 0.01 |
| EWC | 0.13 \pm 0.02 | 0.22 \pm 0.03 | 0.02 \pm 0.00 | 0.23 \pm 0.01 | 0.33 \pm 0.01 | 0.06 \pm 0.01 |
| PACKNET | 0.22 \pm 0.01 | 0.08 \pm 0.01 | 0.25 \pm 0.00 | 0.55 \pm 0.01 | 0.07 \pm 0.02 | 0.63 \pm 0.00 |
| MTL | | | 0.48 \pm 0.01 | | | 0.83 \pm 0.00 |
| | LIBERO-OBJECT | | | LIBERO-GOAL | | |
| SEQL | 0.78 \pm 0.04 | 0.76 \pm 0.04 | 0.26 \pm 0.02 | 0.77 \pm 0.01 | 0.82 \pm 0.01 | 0.22 \pm 0.00 |
| ER | 0.67 \pm 0.07 | 0.43 \pm 0.04 | 0.44 \pm 0.06 | 0.64 \pm 0.01 | 0.34 \pm 0.02 | 0.49 \pm 0.02 |
| EWC | 0.56 \pm 0.03 | 0.69 \pm 0.02 | 0.16 \pm 0.02 | 0.32 \pm 0.02 | 0.48 \pm 0.03 | 0.06 \pm 0.00 |
| PACKNET | 0.60 \pm 0.07 | 0.17 \pm 0.05 | 0.60 \pm 0.05 | 0.63 \pm 0.02 | 0.06 \pm 0.01 | 0.75 \pm 0.01 |
| MTL | | | 0.54 \pm 0.02 | | | 0.80 \pm 0.01 |

Table 2: Performance of three lifelong algorithms and the SEQL and MTL baselines on the four task suites, where the policy is fixed to be RESNET-T. Results are averaged over three seeds and we report the mean and standard error. The best performance is **bolded**, and colored in **purple** if the improvement is statistically significant over other algorithms, when a two-tailed, Student’s t-test under equal sample sizes and unequal variance is applied with a p -value of 0.05.

416 LONG; **3**) EWC works worse than SEQL, showing that the regularization on the loss term can actually
 417 impede the agent’s performance on LLDM problems (See Appendix I.3); and **4**) ER, the rehearsal
 418 method, is robust across all task suites.

419 **Study on Language Embeddings as the Task Identifier (Q4)** To investigate to what extent
 420 language embedding play a role in LLDM, we compare the performance of the same lifelong learner
 421 using four different pretrained language embeddings. Namely, we choose BERT [19], CLIP [20],
 422 GPT-2 [21] and the Task-ID embedding. Task-ID embeddings are produced by feeding a string such
 423 as “Task 5” into a pretrained BERT model.

| Embedding Type | Dimension | FWT(\uparrow) | NBT(\downarrow) | AUC(\uparrow) |
|----------------|-----------|------------------------|------------------------|------------------------|
| BERT | 768 | 0.48 \pm 0.02 | 0.32 \pm 0.04 | 0.32 \pm 0.01 |
| CLIP | 512 | 0.52 \pm 0.00 | 0.34 \pm 0.01 | 0.35 \pm 0.01 |
| GPT-2 | 768 | 0.46 \pm 0.01 | 0.34 \pm 0.02 | 0.30 \pm 0.01 |
| Task-ID | 768 | 0.50 \pm 0.01 | 0.37 \pm 0.01 | 0.33 \pm 0.01 |

Table 3: Performance of a lifelong learner using four different language embeddings on LIBERO-LONG, where we fix the policy architecture to RESNET-T and the lifelong learning algorithm to ER. The Task-ID embeddings are retrieved by feeding “Task + ID” into a pretrained BERT model. Results are averaged over three seeds and we report the mean and standard error. The best performance is **bolded**. No statistically significant difference is observed among the different language embeddings.

424 **Findings:** From Table 3, we observe *no* statistically significant difference among various language
 425 embeddings, including the Task-ID embedding. This, we believe, is due to sentence embeddings
 426 functioning as bag-of-words that differentiates different tasks. This insight calls for better language
 427 encoding to harness the semantic information in task descriptions. Despite the similar performance,
 428 we opt for BERT embeddings as our default task embedding.

429 **Study on How Pretraining Affects Downstream LLDM (Q6)** Fig 3 reports the results on
 430 LIBERO-LONG of five combinations of algorithms and policy architectures, when the underlying
 431 model is pretrained on the 90 short-horizon tasks in LIBERO-100 or learned from scratch. For
 432 pretraining, we apply behavioral cloning on the 90 tasks using the three policy architectures for 50
 433 epochs. We save a checkpoint every 5 epochs of training and then pick the checkpoint for each
 434 architecture that has the best performance as the pretrained model for downstream LLDM.

435 **Findings:** We observe that the basic supervised pretraining can *hurt* the model’s downstream lifelong
 436 learning performance. This, together with the results seen in Table 2 (e.g., naive sequential fine-tuning

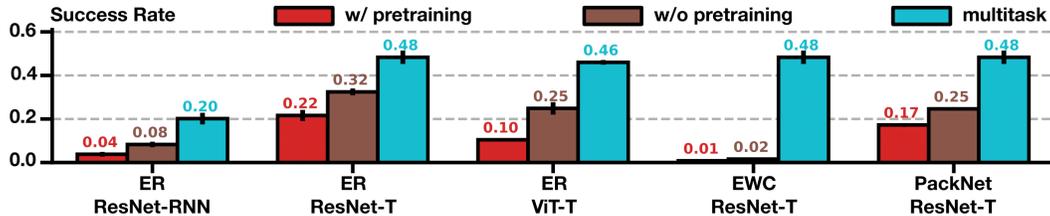


Figure 3: Performance of different combinations of algorithms and architectures without pretraining or with pretraining. The multi-task learning performance is also included for reference.

437 has better forward transfer than when lifelong learning algorithms are applied), indicates that better
 438 pretraining techniques are needed.

439 **Attention Visualization:** To better understand what type of knowledge the agent forgets during the
 440 lifelong learning process, we visualize the agent’s attention map on each observed image input. The
 441 visualized saliency maps and the discussion can be found in Appendix I.4.

442 B Related Work

443 This section provides an overview of existing benchmarks for lifelong learning and robot learning.
 444 We refer the reader to Appendix F.1 for a detailed review of lifelong learning algorithms.

445 **Lifelong Learning Benchmarks** Pioneering work has adapted standard vision or language
 446 datasets for studying LL. This line of work includes image classification datasets like MNIST [22],
 447 CIFAR [23], and ImageNet [24]; segmentation datasets like Core50 [25]; and natural language
 448 understanding datasets like GLUE [26] and SuperGLUE [27]. Besides supervised learning datasets,
 449 video game benchmarks (e.g., Atari [28], XLand [29], and VisDoom [30]) in reinforcement learning
 450 (RL) have also been used for studying LL. However, LL in standard supervised learning does not
 451 involve procedural knowledge transfer, while RL problems in games do not represent human activities.
 452 ContinualWorld [15] modifies the 50 manipulation tasks in MetaWorld for LL. CORA [31] builds
 453 four lifelong RL benchmarks based on Atari, Procgen [32], MiniHack [33], and ALFRED [34].
 454 F-SIOL-310 [35] and OpenLORIS [36] are challenging real-world lifelong object learning datasets
 455 that are captured from robotic vision systems. Prior works have also analyzed different components
 456 in a LL agent [37–39], but they do not focus on robot manipulation problems.

457 **Robot Learning Benchmarks** A variety of robot learning benchmarks have been proposed
 458 to address challenges in meta learning (MetaWorld [40]), causality learning (CausalWorld [41]),
 459 multi-task learning [42, 43], policy generalization to unseen objects [44, 45], and compositional
 460 learning [46]. Compared to existing benchmarks in lifelong learning and robot learning, the task
 461 suites in LIBERO are curated to address the research topics of LLDM. The benchmark includes a
 462 large number of tasks based on everyday human activities that feature rich interactive behaviors with
 463 a diverse range of objects. Additionally, the tasks in LIBERO are procedurally generated, making
 464 the benchmark scalable and adaptable. Moreover, the provided high-quality human demonstration
 465 dataset in LIBERO supports and encourages learning efficiency.

466 C Background

467 This section introduces the problem formulation and defines key terms used throughout the paper.

468 C.1 Markov Decision Process for Robot Learning

469 A robot learning problem can be formulated as a finite-horizon Markov Decision Process: $\mathcal{M} =$
 470 $(\mathcal{S}, \mathcal{A}, \mathcal{T}, H, \mu_0, R)$. Here, \mathcal{S} and \mathcal{A} are the state and action spaces of the robot. μ_0 is the initial
 471 state distribution, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is the transition
 472 function. In this work, we assume a sparse-reward setting and replace R with a goal predicate

473 $g : \mathcal{S} \rightarrow \{0, 1\}$. The robot’s objective is to learn a policy π that maximizes the expected return:
 474 $\max_{\pi} J(\pi) = \mathbb{E}_{s_t, a_t \sim \pi, \mu_0} [\sum_{t=1}^H g(s_t)]$.

475 C.2 Lifelong Robot Learning Problem

476 In a *lifelong robot learning problem*, a robot sequentially learns over K tasks $\{T^1, \dots, T^K\}$ with a
 477 single policy π . We assume π is conditioned on the task, i.e., $\pi(\cdot | s; T)$. For each task, $T^k \equiv (\mu_0^k, g^k)$
 478 is defined by the initial state distribution μ_0^k and the goal predicate g^k .³ We assume $\mathcal{S}, \mathcal{A}, \mathcal{T}, H$ are
 479 the same for all tasks. Up to the k -th task T^k , the robot aims to optimize

$$\max_{\pi} J_{\text{LRL}}(\pi) = \frac{1}{k} \sum_{p=1}^k \left[\mathbb{E}_{s_t^p, a_t^p \sim \pi(\cdot; T^p), \mu_0^p} \left[\sum_{t=1}^L g^p(s_t^p) \right] \right]. \quad (1)$$

480 An important feature of the lifelong setting is that the agent loses access to the previous $k - 1$ tasks
 481 when it learns on task T^k .

482 **Lifelong Imitation Learning** Due to the challenge of sparse-reward reinforcement learning, we
 483 consider a practical alternative setting where a user would provide a small demonstration dataset
 484 for each task in the sequence. Denote $D^k = \{\tau_i^k\}_{i=1}^N$ as N demonstrations for task T^k . Each
 485 $\tau_i^k = (o_0, a_0, o_1, a_1, \dots, o_{l^k})$ where $l^k \leq H$. Here, o_t is the robot’s sensory input, including the
 486 perceptual observation and the information about the robot’s joints and gripper. In practice, the
 487 observation o_t is often non-Markovian. Therefore, following works in partially observable MDPs [47],
 488 we represent s_t by the aggregated history of observations, i.e. $s_t \equiv o_{\leq t} \triangleq (o_0, o_1, \dots, o_t)$. This
 489 results in the *lifelong imitation learning problem* with the same objective as in Eq. (1). But during
 490 training, we perform behavioral cloning [48] with the following surrogate objective function:

$$\min_{\pi} J_{\text{BC}}(\pi) = \frac{1}{k} \sum_{p=1}^k \mathbb{E}_{o_t, a_t \sim D^p} \left[\sum_{t=0}^{l^p} \mathcal{L}(\pi(o_{\leq t}; T^p), a_t^p) \right], \quad (2)$$

491 where \mathcal{L} is a supervised learning loss, e.g., the negative log-likelihood loss, and π is a Gaussian
 492 mixture model. Similarly, we assume $\{D^p : p < k\}$ are not fully available when learning T^k .

493 D Metrics Visualization

We provide a visualization of the three metrics we compute in Figure 4. For completeness, we also

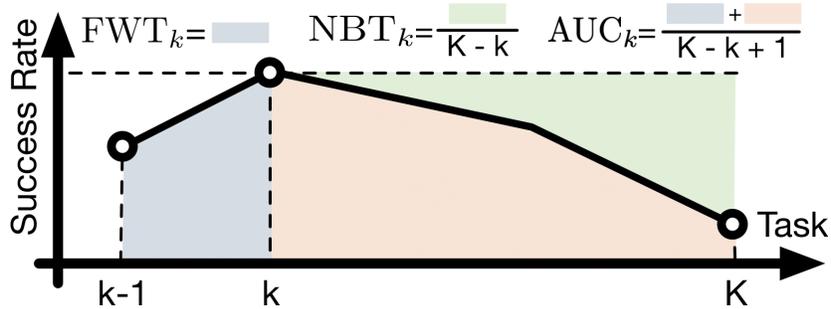


Figure 4: Metrics for LLDM

494 provide the formulas for the metrics here:
 495

³Throughout the paper, a superscript/subscript is used to index the task/time step.

$$\begin{aligned} \text{FWT} &= \sum_{k \in [K]} \frac{\text{FWT}_k}{K}, & \text{FWT}_k &= \frac{1}{11} \sum_{e \in \{0 \dots 50\}} c_{k,k,e} \\ \text{NBT} &= \sum_{k \in [K]} \frac{\text{NBT}_k}{K}, & \text{NBT}_k &= \frac{1}{K-k} \sum_{\tau=k+1}^K (c_{k,k} - c_{\tau,k}) \\ \text{AUC} &= \sum_{k \in [K]} \frac{\text{AUC}_k}{K}, & \text{AUC}_k &= \frac{1}{K-k+1} (\text{FWT}_k + \sum_{\tau=k+1}^K c_{\tau,k}). \end{aligned}$$

496 E Implemented Neural Architectures and Lifelong Learning Algorithms

| | |
|-------------------------|--------------|
| Neural Policy Arch. | RESNET-RNN |
| | RESNET-T |
| | VIT-T |
| Lifelong Learning Algo. | SEQL |
| | EWC [13] |
| | ER [12] |
| | PACKNET [14] |
| | MTL |

Table 4: The implemented neural policy architectures and the lifelong learning algorithms in LIBERO.

497 E.1 Neural Network Architectures

498 We implement three vision-language policy networks, RESNET-RNN, RESNET-T, and VIT-T, that
 499 integrate visual, temporal, and linguistic information for LLD. Language instructions of tasks
 500 are encoded using pretrained BERT embeddings [19]. The RESNET-RNN [18] uses a ResNet as the
 501 visual backbone that encodes per-step visual observations and an LSTM as the temporal backbone to
 502 process a sequence of encoded visual information. The language instruction is incorporated into the
 503 ResNet features using the FiLM method [49] and added to the LSTM inputs, respectively. RESNET-T
 504 architecture [50] uses a similar ResNet-based visual backbone, but a transformer decoder [51] as
 505 the temporal backbone to process outputs from ResNet, which are a temporal sequence of visual
 506 tokens. The language embedding is treated as a separate token in inputs to the transformer alongside
 507 the visual tokens. The VIT-T architecture [52], which is widely used in visual-language tasks, uses a
 508 Vision Transformer (ViT) as the visual backbone and a transformer decoder as the temporal backbone.
 509 The language embedding is treated as a separate token in inputs of both ViT and the transformer
 510 decoder. All the temporal backbones output a latent vector for every decision-making step. We
 511 compute the multi-modal distribution over manipulation actions using a Gaussian-Mixture-Model
 512 (GMM) based output head [53, 18, 54]. In the end, a robot executes a policy by sampling a continuous
 513 value for end-effector action from the output distribution. Figure 5 visualizes the three architectures.

514 For all the lifelong learning algorithms and neural architectures, we use behavioral cloning (BC) [48]
 515 to train policies for individual tasks (See (2)). BC allows for efficient policy learning such that we
 516 can study lifelong learning algorithms with limited computational resources. To train BC, we provide
 517 50 trajectories of high-quality demonstrations for every single task in the generated task suites. The
 518 demonstrations are collected by human experts through teleoperation with 3Dconnexion Spacemouse.

519 In Section E.1, we outlined the neural network architectures utilized in our experiments, namely
 520 RESNET-RNN, RESNET-T, and VIT-T. The specifics of each architecture are illustrated in Figure 5.
 521 Furthermore, Table 5, 6, and 7 display the hyperparameters for the architectures used throughout all
 522 of our experiments.

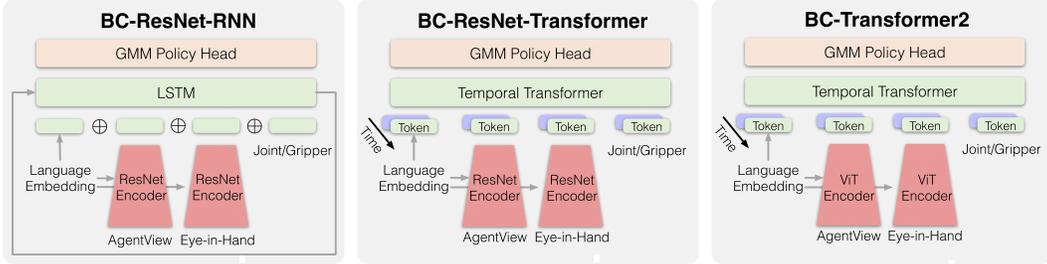


Figure 5: We provide visualizations of the architectures for RESNET-RNN, RESNET-T, and ViT-T, respectively. It is worth noting that each model architecture incorporates language embedding in distinct ways.

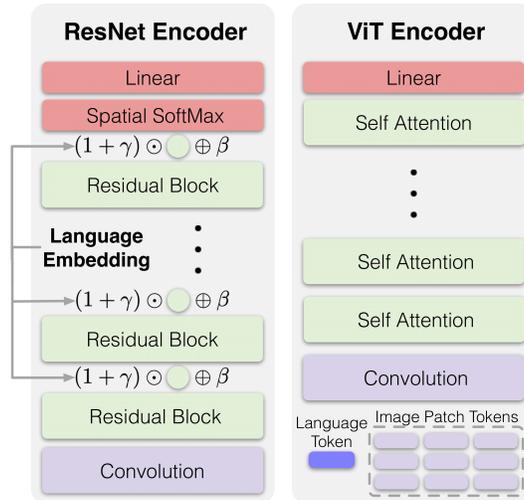


Figure 6: The image encoders: ResNet-based encoder and the vision transformer-based encoder.

523 **F Computation**

524 For all experiments, we use a single Nvidia A100 GPU or a single Nvidia A40 GPU (CUDA 11.7)
 525 with 8 16 CPUs for training and evaluation.

| Variable | Value |
|-------------------------|-------|
| resnet_image_embed_size | 64 |
| text_embed_size | 32 |
| rnn_hidden_size | 1024 |
| rnn_layer_num | 2 |
| rnn_dropout | 0.0 |

526

Table 5: Hyper parameters of RESNET-RNN.

| Variable | Value |
|------------------------------|-------|
| extra_info_hidden_size | 128 |
| img_embed_size | 64 |
| transformer_num_layers | 4 |
| transformer_num_heads | 6 |
| transformer_head_output_size | 64 |
| transformer_mlp_hidden_size | 256 |
| transformer_dropout | 0.1 |
| transformer_max_seq_len | 10 |

Table 6: Hyper parameters of RESNET-T.

| Variable | Value |
|---------------------------------------|-------|
| extra_info_hidden_size | 128 |
| img_embed_size | 128 |
| spatial_transformer_num_layers | 7 |
| spatial_transformer_num_heads | 8 |
| spatial_transformer_head_output_size | 120 |
| spatial_transformer_mlp_hidden_size | 256 |
| spatial_transformer_dropout | 0.1 |
| spatial_down_sample_embed_size | 64 |
| temporal_transformer_input_size | null |
| temporal_transformer_num_layers | 4 |
| temporal_transformer_num_heads | 6 |
| temporal_transformer_head_output_size | 64 |
| temporal_transformer_mlp_hidden_size | 256 |
| temporal_transformer_dropout | 0.1 |
| temporal_transformer_max_seq_len | 10 |

Table 7: Hyper parameters of ViT-T.

527 F.1 Lifelong Learning Algorithms

528 Lifelong learning (LL) is a field of study that aims to understand how an agent can continually
529 acquire and retain knowledge over an infinite sequence of tasks without catastrophically forgetting
530 previous knowledge. Recent literature proposes three main approaches to address the problem of
531 catastrophic forgetting in deep learning: Dynamic Architecture approaches, Regularization-Based
532 approaches, and Rehearsal approaches. Although some recent works explore the combination of
533 different approaches [55–57] or new strategies [58–60], our benchmark aims to provide an in-depth
534 analysis of these three basic lifelong learning directions to reveal their pros and cons on robot learning
535 tasks.

536 The dynamic architecture approach gradually expands the learning model to incorporate new knowl-
537 edge [61, 62, 14, 63–65]. Regularization-based methods, on the other hand, regularize the learner to
538 a previous checkpoint when it learns a new task [13, 66–68]. Rehearsal methods save exemplar data
539 from prior tasks and replay them with new data to consolidate the agent’s memory [12, 69–71]. For a
540 comprehensive review of LL methods, we refer readers to surveys [72, 73].

541 The following paragraphs provide details on the three lifelong learning algorithms that we have
542 implemented.

543 **ER** Experience Replay (ER) [12] is a **rehearsal-based** approach that maintains a memory buffer
544 of samples from previous tasks and leverages it to learn new tasks. After the completion of policy
545 learning for a task, ER stores a portion of the data into a storage memory. When training a new
546 task, ER samples data from the memory and combines it with the training data from the current task
547 so that the training data approximately represents the empirical distribution of all-task data. In our

548 implementation, we use a replay buffer to store a portion of the training data (up to 1000 trajectories)
 549 after training each task. For every training iteration during the training of a new task, we uniformly
 550 sample a fixed number of replay data from the memory (32 trajectories) along with each batch of
 551 training data from the new task.

EWC Elastic Weight Consolidation(EWC) [13] is a **regularization-based** approach that add a regularization term that constraints neural network update to the original single-task learning objective. Specifically, EWC uses the Fisher information matrix that quantify the importance of every neural network parameter. The loss function for task k is:

$$\mathcal{L}_k^{EWC}(\theta) = \mathcal{L}_K^{BC}(\theta) + \sum_i \frac{\lambda}{2} F_i (\theta_i - \theta_{k-1,i}^*)^2,$$

552 where λ is a penalty hyperparameter, and the coefficient F_i is the diagonal of the Fisher information
 553 matrix: $F_k = \mathbb{E}_{s \sim \mathcal{D}_k} \mathbb{E}_{a \sim p_{\theta}(\cdot|s)} (\nabla_{\theta_k} \log p_{\theta_k}(a|s))^2$. In this work, we use the online update version
 554 of EWC that updates the Fisher information matrix using exponential moving average along the
 555 lifelong learning process, and use the empirical estimation of above Fisher information matrix to
 556 stabilize the estimation. Formally, the actually used estimation of Fisher Information Matrix is
 557 $\tilde{F}_k = \gamma F_{k-1} + (1 - \gamma) F_k$, where $F_k = \mathbb{E}_{(s,a) \sim \mathcal{D}_k} (\nabla_{\theta_k} \log p_{\theta_k}(a|s))^2$ and k is the task number. We
 558 set $\gamma = 0.9$ and $\lambda = 5 \cdot 10^4$.

559 **PACKNET** PACKNET [14] is a **dynamic architecture-based** approach that aims to prevent changes
 560 to parameters that are important for previous tasks in lifelong learning. To achieve this, PACKNET
 561 iteratively trains, prunes, fine-tunes, and freezes parts of the network. The method theoretically
 562 completely avoids catastrophic forgetting, but for each new task, the number of available parameters
 563 shrinks. The pruning process in PACKNET involves two stages. First, the network is trained, and at the
 564 end of the training, a fixed proportion of the most important parameters (25% in our implementation)
 565 are chosen, and the rest are pruned. Second, the selected part of the network is fine-tuned and then
 566 frozen. In our implementation, we follow the original paper [14] and do not train all biases and
 567 normalization layers. We perform the same number of fine-tuning epochs as for training (50 epochs
 568 in our implementation). Note that all evaluation metrics are calculated *before* the fine-tuning stage.

569 **G LIBERO Task Suite Designs**

570 **G.1 Task Suites**

571 We visualize all the tasks from the four task suites in Figure 7- 10. Figure 7 visualizes the initial
572 states since the task goals are always the same. All the figures visualize the goal states of tasks except
573 for Figure 7, which visualizes the initial states since the task goals are always the same.

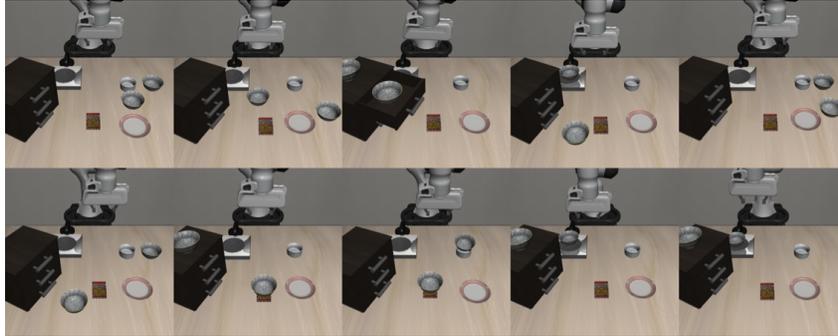


Figure 7: LIBERO-SPATIAL

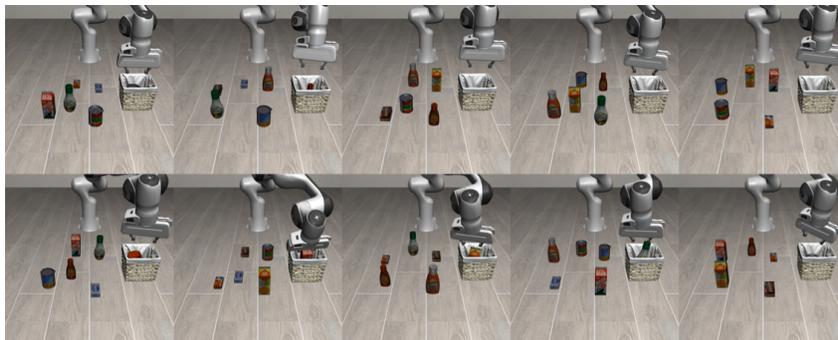


Figure 8: LIBERO-OBJECT

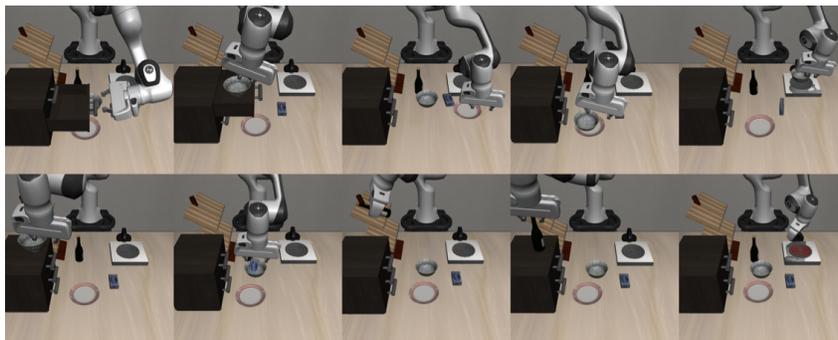


Figure 9: LIBERO-GOAL

574
575
576
577
578
579

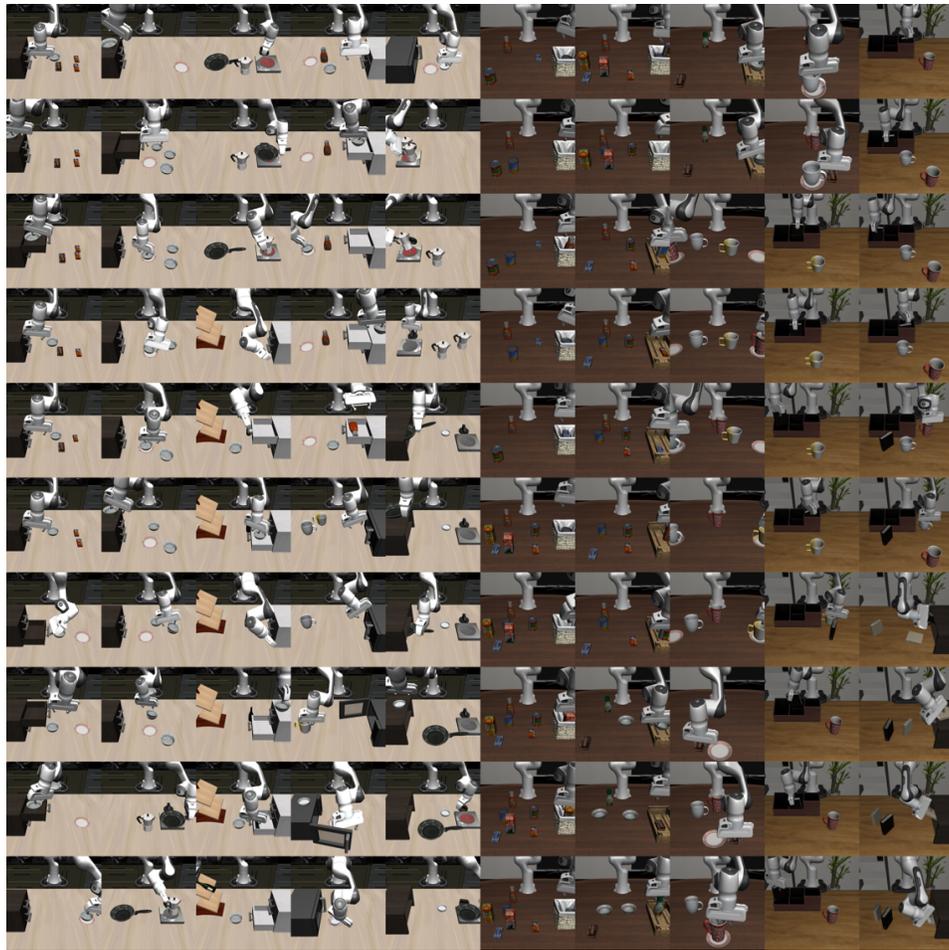


Figure 10: LIBERO-100

580

581

582 G.2 PDDL-based Scene Description File

583 Here we visualize the whole content of an example scene description file based on PDDL. This file
584 corresponds to the task shown in Figure 2.

585 **Example task:** *Open the top drawer of the cabinet and put the bowl in it.*

```
586 (define (problem LIBERO_Kitchen_Tabletop_Manipulation)
587   (:domain robosuite)
588   (:language open the top drawer of the cabinet and put the bowl in it)
589   (:regions
590     (wooden_cabinet_init_region
591       (:target kitchen_table)
592       (:ranges (
593         (-0.01 -0.31 0.01 -0.29)
594       ))
595     )
596     (:yaw_rotation (
597       (3.141592653589793 3.141592653589793)
598     ))
599   )
600 )
601 (akita_black_bowl_init_region
602   (:target kitchen_table)
603   (:ranges (
604     (-0.025 -0.025 0.025 0.025)
605   ))
606 )
607   (:yaw_rotation (
608     (0.0 0.0)
609   ))
610 )
611 )
612 (plate_init_region
613   (:target kitchen_table)
614   (:ranges (
615     (-0.025 0.225 0.025 0.275)
616   ))
617 )
618   (:yaw_rotation (
619     (0.0 0.0)
620   ))
621 )
622 )
623 (top_side
624   (:target wooden_cabinet_1)
625 )
626 (top_region
627   (:target wooden_cabinet_1)
628 )
629 (middle_region
630   (:target wooden_cabinet_1)
631 )
632 (bottom_region
633   (:target wooden_cabinet_1)
634 )
635 )
636 (:fixtures
637
```

```

638     kitchen_table - kitchen_table
639     wooden_cabinet_1 - wooden_cabinet
640 )
641
642 (: objects
643     akita_black_bowl_1 - akita_black_bowl
644     plate_1 - plate
645 )
646
647 (: obj_of_interest
648     wooden_cabinet_1
649     akita_black_bowl_1
650 )
651
652 (: init
653     (On akita_black_bowl_1 kitchen_table_akita_black_bowl_init_region)
654     (On plate_1 kitchen_table_plate_init_region)
655     (On wooden_cabinet_1 kitchen_table_wooden_cabinet_init_region)
656 )
657
658 (: goal
659     (And (Open wooden_cabinet_1_top_region)
660          (In akita_black_bowl_1 wooden_cabinet_1_top_region)
661         )
662 )
663
664 )

```

665 H Experimental Setup

666 We consider five lifelong learning algorithms: SEQL the sequential learning baseline where the
667 agent learns each task in the sequence directly without any further consideration, MTL the multitask
668 learning baseline where the agent learns all tasks in the sequence simultaneously, the regularization-
669 based method EWC [13], the replay-based method ER [12], and the dynamic architecture-based
670 method PACKNET [14]. SEQL and MTL can be seen as approximations of the lower and upper
671 bounds respectively for any lifelong learning algorithm. The other three methods represent the
672 three primary categories of lifelong learning algorithms. For the neural architectures, we consider
673 three vision-language policy architectures: RESNET-RNN, RESNET-T, ViT-T, which differ in how
674 spatial or temporal information is aggregated (See Appendix E.1 for more details). For each task,
675 the agent is trained over 50 epochs on the 50 demonstration trajectories. We evaluate the agent’s
676 average success rate over 20 test rollout trajectories of a maximum length of 600 every 5 epochs.
677 We use Adam optimizer [74] with a batch size of 32, and a cosine scheduled learning rate from
678 0.0001 to 0.00001 for each task. Following the convention of Robomimic [18], we pick the model
679 checkpoint that achieves the best success rate as the final policy for a given task. After 50 epochs
680 of training, the agent with the best checkpoint is then evaluated on all previously learned tasks,
681 with 20 test rollout trajectories for each task. All policy networks are matched in Floating Point
682 Operations Per Second (FLOPS): all policy architectures have $\sim 13.5\text{G}$ FLOPS. For each combination
683 of algorithm, policy architecture, and task suite, we run the lifelong learning method 3 times with
684 random seeds $\{100, 200, 300\}$ (180 experiments in total). See Table 4 for the implemented algorithms
685 and architectures.

686 I Additional Experiment Results

687 I.1 Full Results

688 We provide the full results across three different lifelong learning algorithms (e.g., EWC, ER,
689 PACKNET) and three different policy architectures (e.g., RESNET-RNN, RESNET-T, ViT-T) on the
690 four task suites in Table 8.

691 To better illustrate the performance of each lifelong learning agent throughout the learning process,
692 we present plots that show how the agent’s performance evolves over the stream of tasks. Firstly, we
693 provide plots that compare the performance of the agent using different lifelong learning algorithms
694 while fixing the policy architecture (refer to Figure 11, 12, and 13). Next, we provide plots that
695 compare the performance of the agent using different policy architectures while fixing the lifelong
696 learning algorithm (refer to Figure 14, 15, and 16)

| Algo. | Policy Arch. | FWT(\uparrow) | NBT(\downarrow) | AUC(\uparrow) | FWT(\uparrow) | NBT(\downarrow) | AUC(\uparrow) |
|---------|--------------|-------------------|---------------------|-------------------|-------------------|---------------------|-------------------|
| | | LIBERO-LONG | | | LIBERO-SPATIAL | | |
| SEQL | RESNET-RNN | 0.24 \pm 0.02 | 0.28 \pm 0.01 | 0.07 \pm 0.01 | 0.50 \pm 0.01 | 0.61 \pm 0.01 | 0.14 \pm 0.01 |
| | RESNET-T | 0.54 \pm 0.01 | 0.63 \pm 0.01 | 0.15 \pm 0.00 | 0.72 \pm 0.01 | 0.81 \pm 0.01 | 0.20 \pm 0.01 |
| | ViT-T | 0.44 \pm 0.04 | 0.50 \pm 0.05 | 0.13 \pm 0.01 | 0.63 \pm 0.02 | 0.76 \pm 0.01 | 0.16 \pm 0.01 |
| ER | RESNET-RNN | 0.16 \pm 0.02 | 0.16 \pm 0.02 | 0.08 \pm 0.01 | 0.40 \pm 0.02 | 0.29 \pm 0.02 | 0.29 \pm 0.01 |
| | RESNET-T | 0.48 \pm 0.02 | 0.32 \pm 0.04 | 0.32 \pm 0.01 | 0.65 \pm 0.03 | 0.27 \pm 0.03 | 0.56 \pm 0.01 |
| | ViT-T | 0.38 \pm 0.05 | 0.29 \pm 0.06 | 0.25 \pm 0.02 | 0.63 \pm 0.01 | 0.29 \pm 0.02 | 0.50 \pm 0.02 |
| EWC | RESNET-RNN | 0.02 \pm 0.00 | 0.04 \pm 0.01 | 0.00 \pm 0.00 | 0.14 \pm 0.02 | 0.23 \pm 0.02 | 0.03 \pm 0.00 |
| | RESNET-T | 0.13 \pm 0.02 | 0.22 \pm 0.03 | 0.02 \pm 0.00 | 0.23 \pm 0.01 | 0.33 \pm 0.01 | 0.06 \pm 0.01 |
| | ViT-T | 0.05 \pm 0.02 | 0.09 \pm 0.03 | 0.01 \pm 0.00 | 0.32 \pm 0.03 | 0.48 \pm 0.03 | 0.06 \pm 0.01 |
| PACKNET | RESNET-RNN | 0.13 \pm 0.00 | 0.21 \pm 0.01 | 0.03 \pm 0.00 | 0.27 \pm 0.03 | 0.38 \pm 0.03 | 0.06 \pm 0.01 |
| | RESNET-T | 0.22 \pm 0.01 | 0.08 \pm 0.01 | 0.25 \pm 0.00 | 0.55 \pm 0.01 | 0.07 \pm 0.02 | 0.63 \pm 0.00 |
| | ViT-T | 0.36 \pm 0.01 | 0.14 \pm 0.01 | 0.34 \pm 0.01 | 0.57 \pm 0.04 | 0.15 \pm 0.00 | 0.59 \pm 0.03 |
| MTL | RESNET-RNN | | | 0.20 \pm 0.01 | | | 0.61 \pm 0.00 |
| | RESNET-T | | | 0.48 \pm 0.01 | | | 0.83 \pm 0.00 |
| | ViT-T | | | 0.46 \pm 0.00 | | | 0.79 \pm 0.01 |
| | | LIBERO-OBJECT | | | LIBERO-GOAL | | |
| SEQL | RESNET-RNN | 0.48 \pm 0.03 | 0.53 \pm 0.04 | 0.15 \pm 0.01 | 0.61 \pm 0.01 | 0.73 \pm 0.01 | 0.16 \pm 0.00 |
| | RESNET-T | 0.78 \pm 0.04 | 0.76 \pm 0.04 | 0.26 \pm 0.02 | 0.77 \pm 0.01 | 0.82 \pm 0.01 | 0.22 \pm 0.00 |
| | ViT-T | 0.76 \pm 0.03 | 0.73 \pm 0.03 | 0.27 \pm 0.02 | 0.75 \pm 0.01 | 0.85 \pm 0.01 | 0.20 \pm 0.01 |
| ER | RESNET-RNN | 0.30 \pm 0.01 | 0.27 \pm 0.05 | 0.17 \pm 0.05 | 0.41 \pm 0.00 | 0.35 \pm 0.01 | 0.26 \pm 0.01 |
| | RESNET-T | 0.67 \pm 0.07 | 0.43 \pm 0.04 | 0.44 \pm 0.06 | 0.64 \pm 0.01 | 0.34 \pm 0.02 | 0.49 \pm 0.02 |
| | ViT-T | 0.70 \pm 0.02 | 0.28 \pm 0.01 | 0.57 \pm 0.01 | 0.57 \pm 0.00 | 0.40 \pm 0.02 | 0.38 \pm 0.01 |
| EWC | RESNET-RNN | 0.17 \pm 0.04 | 0.23 \pm 0.04 | 0.06 \pm 0.01 | 0.16 \pm 0.01 | 0.22 \pm 0.01 | 0.06 \pm 0.01 |
| | RESNET-T | 0.56 \pm 0.03 | 0.69 \pm 0.02 | 0.16 \pm 0.02 | 0.32 \pm 0.02 | 0.48 \pm 0.03 | 0.06 \pm 0.00 |
| | ViT-T | 0.57 \pm 0.03 | 0.64 \pm 0.03 | 0.23 \pm 0.00 | 0.32 \pm 0.04 | 0.45 \pm 0.04 | 0.07 \pm 0.01 |
| PACKNET | RESNET-RNN | 0.29 \pm 0.02 | 0.35 \pm 0.02 | 0.13 \pm 0.01 | 0.32 \pm 0.03 | 0.37 \pm 0.04 | 0.11 \pm 0.01 |
| | RESNET-T | 0.60 \pm 0.07 | 0.17 \pm 0.05 | 0.60 \pm 0.05 | 0.63 \pm 0.02 | 0.06 \pm 0.01 | 0.75 \pm 0.01 |
| | ViT-T | 0.58 \pm 0.03 | 0.18 \pm 0.02 | 0.56 \pm 0.04 | 0.69 \pm 0.02 | 0.08 \pm 0.01 | 0.76 \pm 0.02 |
| MTL | RESNET-RNN | | | 0.10 \pm 0.03 | | | 0.59 \pm 0.00 |
| | RESNET-T | | | 0.54 \pm 0.02 | | | 0.80 \pm 0.01 |
| | ViT-T | | | 0.78 \pm 0.02 | | | 0.82 \pm 0.01 |

Table 8: We present the full results of all networks and algorithms on all four task suites. For each task suite, we highlight the top three AUC scores among the combinations of the three lifelong learning algorithms and the three neural architectures. The best three results are highlighted in **magenta** (the best), **light magenta** (the second best), and **super light magenta** (the third best), respectively.

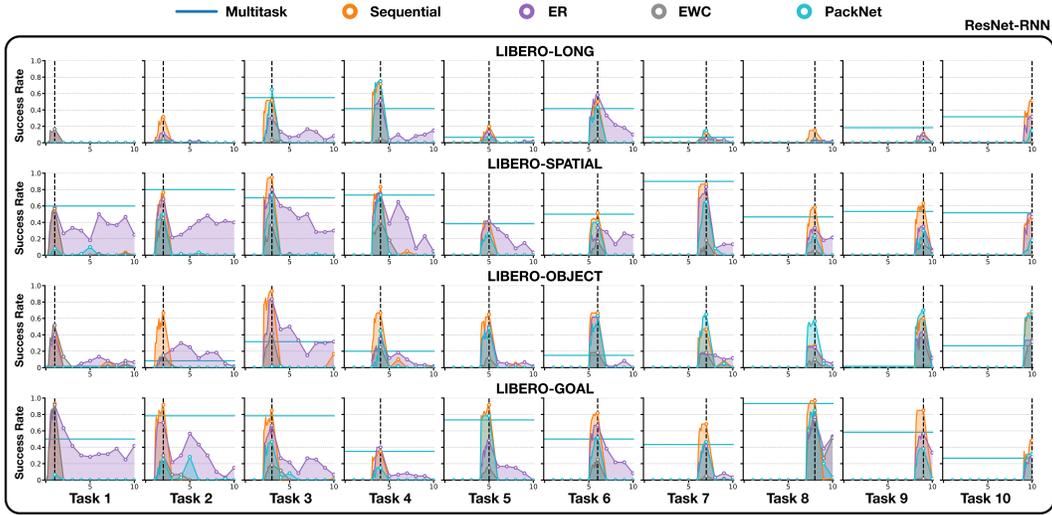


Figure 11: We compare the performance of different algorithms using the RESNET-RNN policy architecture in Figure 11. The y -axis represents the success rate, and the x -axis shows the agent's performance on each of the 10 tasks in a specific task suite over the course of learning. For example, the upper-left plot in the figure displays the agent's performance on the first task as it learns the 10 tasks sequentially.

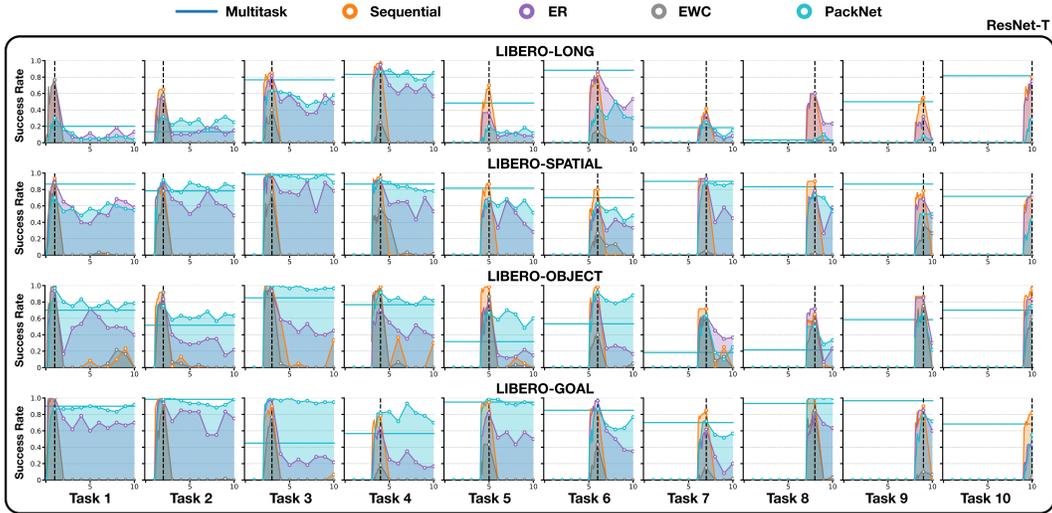


Figure 12: Comparison of different algorithms using the RESNET-T policy architecture. The y -axis represents the success rate, while the x -axis shows the agent's performance on each of the 10 tasks in a given task suite during the course of learning. For example, the plot in the upper-left corner depicts the agent's performance on the first task as it learns the 10 tasks sequentially.

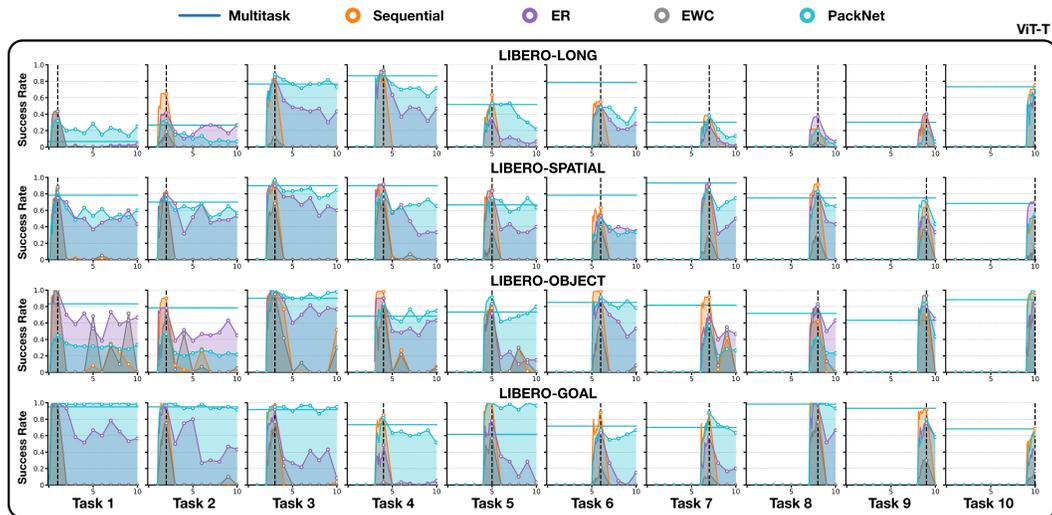


Figure 13: Comparison of different algorithms using the ViT-T policy architecture. The success rate is represented on the y -axis, while the x -axis shows the agent’s performance on the 10 tasks in a given task suite over the course of learning. For instance, the plot in the upper-left corner illustrates the agent’s performance on the first task when learning the 10 tasks sequentially.

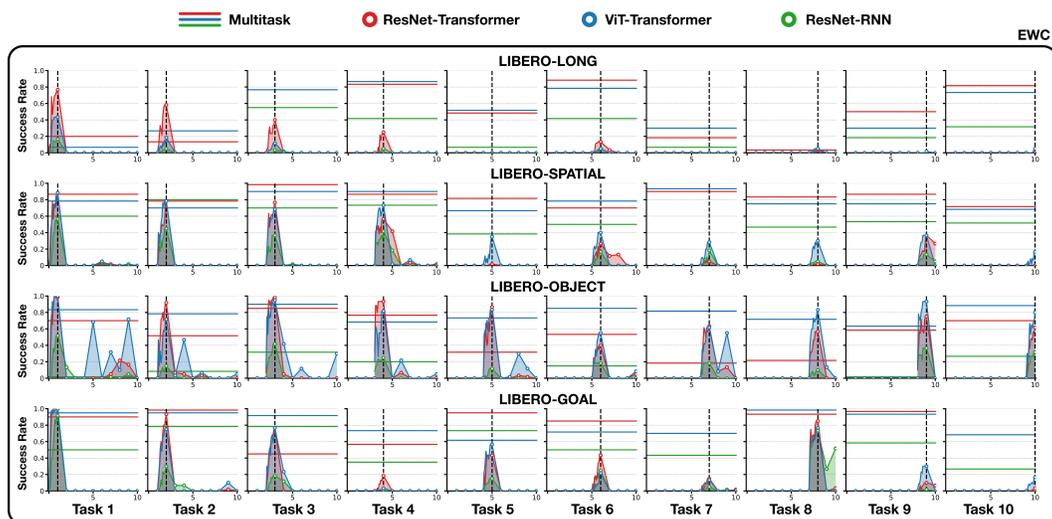


Figure 14: Comparison of different architectures with the EWC algorithm. The y -axis is the success rate, while the x -axis shows the agent’s performance on the 10 tasks in a given task suite over the course of learning. For instance, the upper-left plot shows the agent’s performance on the first task when learning the 10 tasks sequentially.

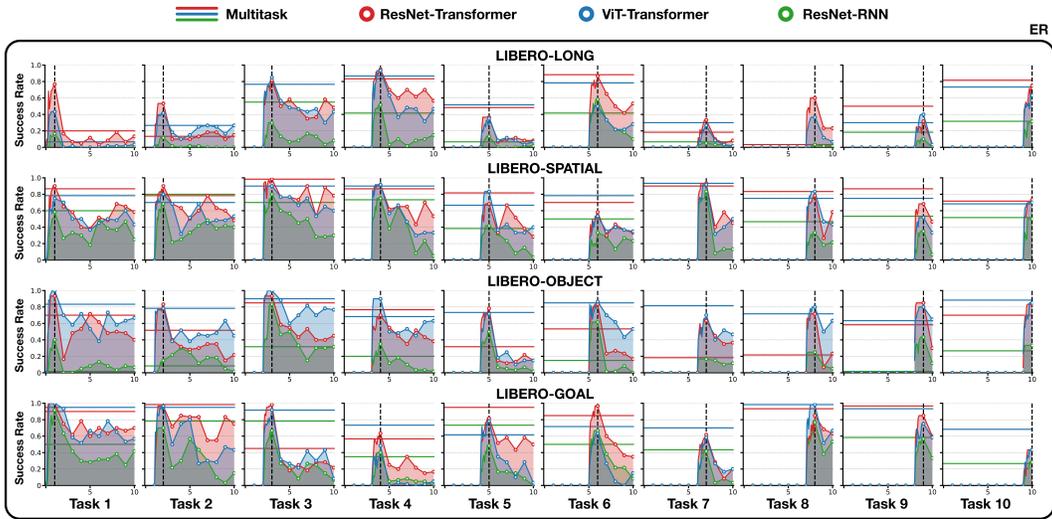


Figure 15: Comparison of different architectures with the ER algorithm. The y -axis is the success rate, while the x -axis shows the agent’s performance on the 10 tasks in a given task suite ver the course of learning. For instance, the upper-left plot shows the agent’s performance on the first task when learning the 10 tasks sequentially.

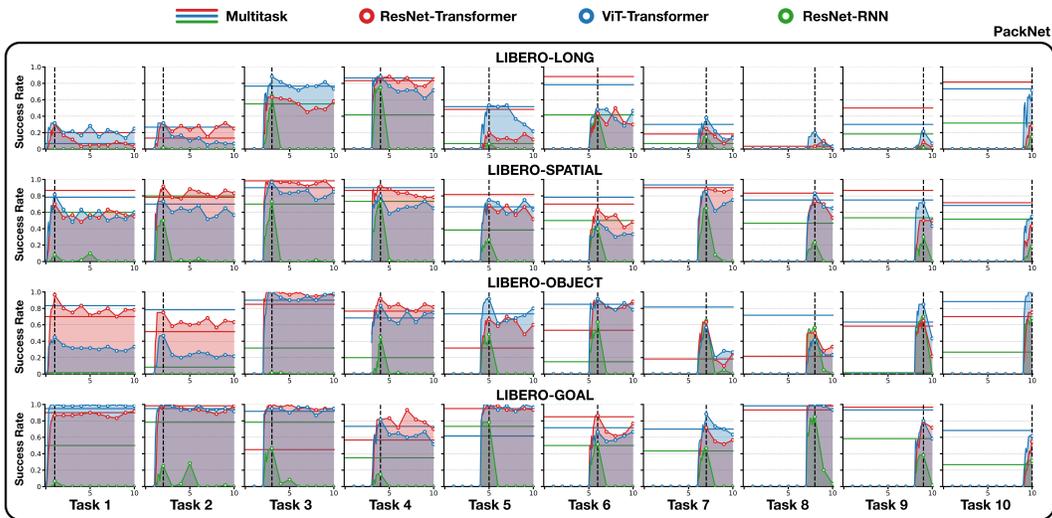


Figure 16: Comparison of different architectures with the PACKNET algorithm. The y -axis is the success rate, while the x -axis shows the agent’s performance on the 10 tasks in a given task suite over the course of learning. For instance, the upper-left plot shows the agent’s performance on the first task when learning the 10 tasks sequentially.

697 **I.2 Study on task ordering (Q4)**

698 Figure 17 shows the result of the study on Q4. For all experiments in this study, we used RESNET-
699 T as the neural architecture and evaluated both ER and PACKNET. As the figure illustrates, the
700 performance of both algorithms varies across different task orderings. This finding highlights an
701 important direction for future research: developing algorithms or architectures that are robust to
702 varying task orderings.

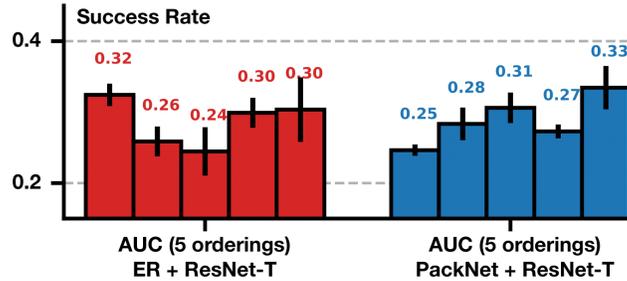


Figure 17: Performance of ER and PACKNET using RESNET-T on five different task orderings. An error bar shows the performance standard deviation for a fixed ordering.

703 *Findings:* From Figure 17, we observe that indeed different task ordering could result in very different
704 performances for the same algorithm. Specifically, such difference is statistically significant for
705 PACKNET.

706 **I.3 Loss v.s. Success Rates**

707 We demonstrate that behavioral cloning loss can be a misleading indicator of task success rate
708 in this section. In supervised learning tasks like image classifications, lower loss often indicates
709 better prediction accuracy. However, this is not, in general, true for decision-making tasks. This is
710 because errors can compound until failures during executing a robot [75]. Figure 18, 12 and 13 plots
711 the training loss and success rates of three lifelong learning methods (ER, EWC, and PACKNET)
712 for comparison. We evaluate the three algorithms on four task suites using three different neural
713 architectures.

714 *Findings:* We observe that though sometimes EWC has the **lowest** loss, it did not achieve good
715 success rate. ER, on the other hand, can have the highest loss but perform better than EWC. In
716 conclusion, success rates, instead of behavioral cloning loss, should be the right metric to evaluate
717 whether a model checkpoint is good or not.

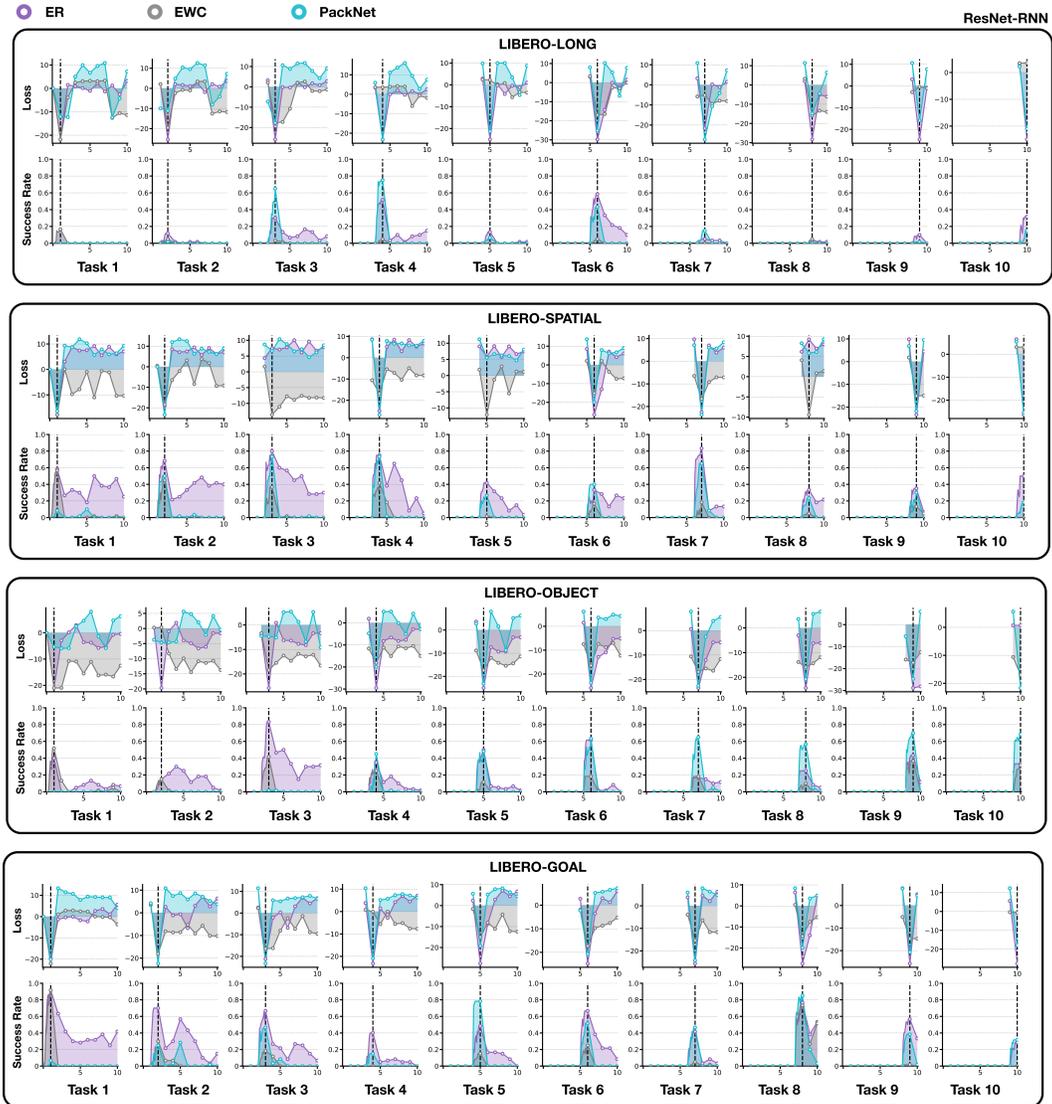


Figure 18: Losses and success rates of ER (violet), EWC (grey), and PACKNET (blue) on four task suites with RESNET-RNN policy. The first (second) row shows the loss (success rate) of the agent on task i throughout the LLDM procedure.

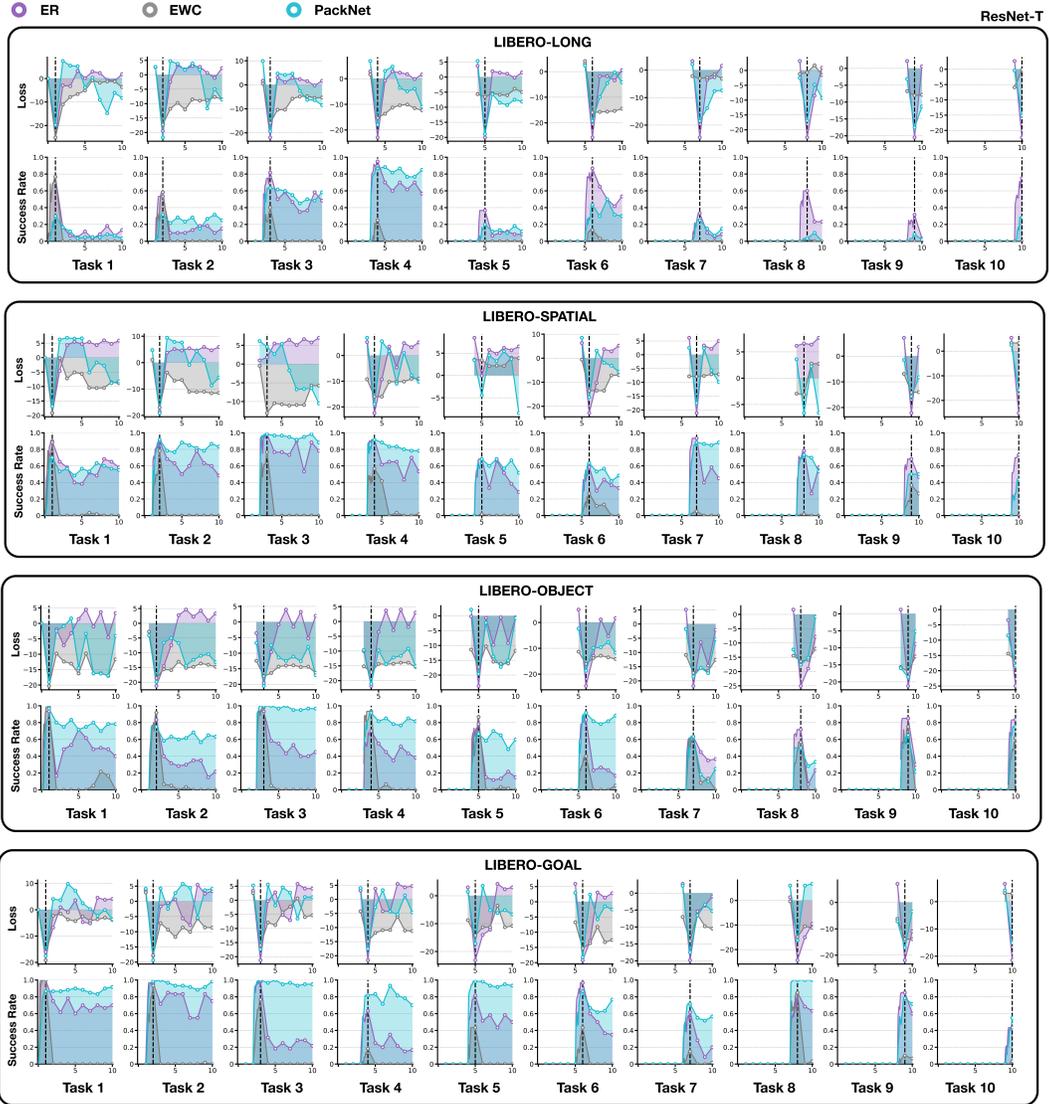


Figure 19: Losses and success rates of ER (violet), EWC (grey), and PACKNET (blue) on four task suites with RESNET-T policy. The first (second) row shows the loss (success rate) of the agent on task i throughout the LLDM procedure.

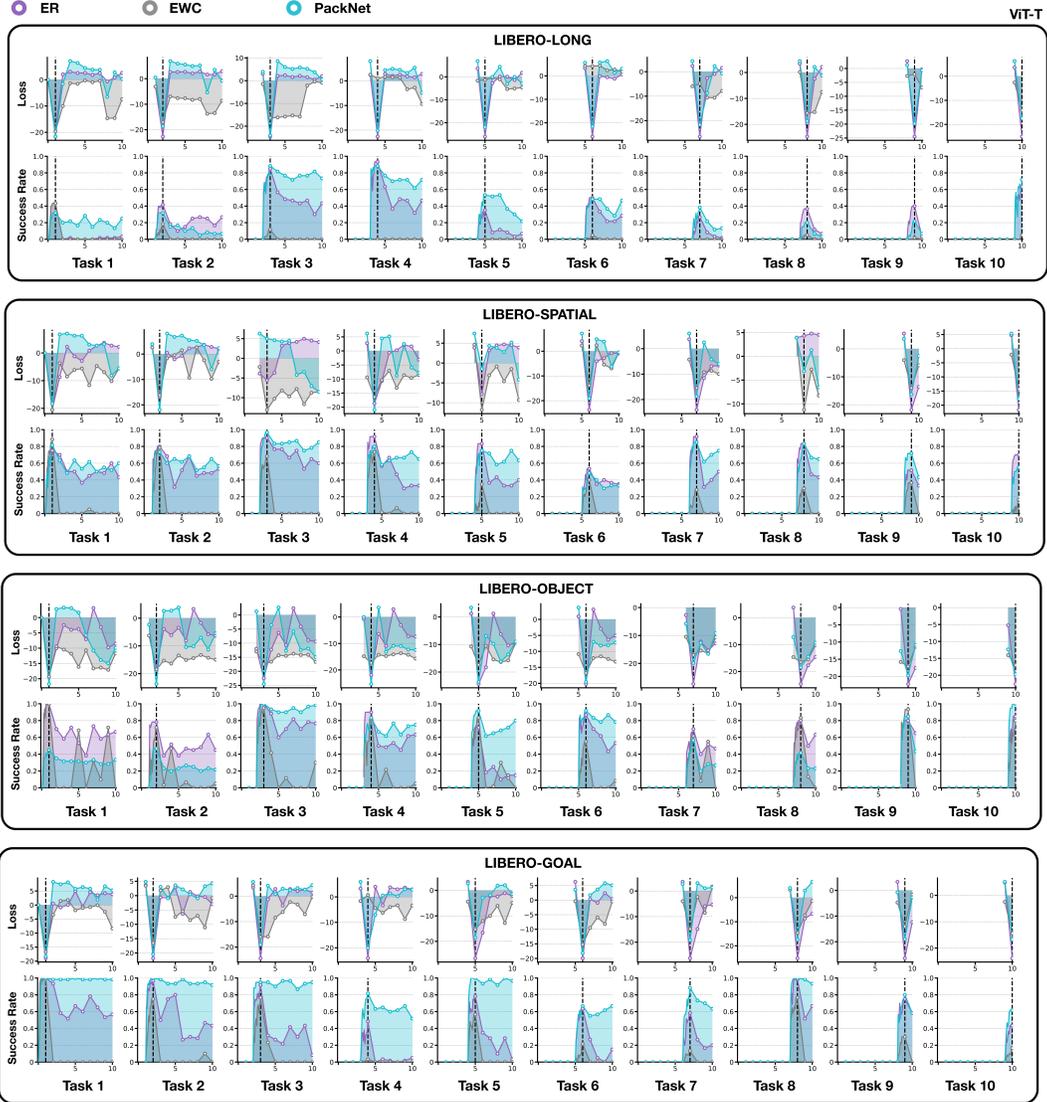


Figure 20: Losses and success rates of ER (violet), EWC (grey), and PACKNET (blue) on four task suites with ViT-T policy. The first (second) row shows the loss (success rate) of the agent on task i throughout the LLDM procedure.

718 **I.4 Attention Visualization**

719 It is also important to visualize the behavior of the robot and its attention maps during the completion
 720 of tasks in the lifelong learning process to give us intuition and qualitative feedback on the perfor-
 721 mance of different algorithms and architectures. We visualize the attention maps of learned policies
 722 with Greydanus et al. [76] and compare them in different studies as in A.2 to see if the robot correctly
 723 pays attention to the right regions of interest in each task.

724 **Perturbation-based attention visualization:** We use a perturbation-based method [76] to extract
 725 attention maps from agents. Given an input image I , the method applies a Gaussian filter to a pixel
 726 location (i, j) to blur the image partially, and produces the perturbed image $\Phi(I, i, j)$. Denote the
 727 learned policy as π and the inputs to the spatial module (e.g., the last latent representation of resnet
 728 or ViT encoder) $\pi_u(I)$ for image I . Then we define the saliency score as the Euclidean distance
 729 between the latent representations of the original and the blurred images:

$$S_\pi(i, j) = \frac{1}{2} \left\| \pi_u(I) - \pi_u(\Phi(I, i, j)) \right\|^2. \quad (3)$$

730 Intuitively, $S_\pi(i, j)$ describes *how much removing information from the region around location (i, j)*
 731 *changes the policy*. In other words, a large $S_\pi(i, j)$ indicates that the information around pixel (i, j)
 732 is important for the learning agent’s decision-making. Instead of calculating the score for every
 733 pixel, [76] found that computing a saliency score for pixel $i \bmod 5$ and $j \bmod 5$ produced good
 734 saliency maps at lower computational costs for Atari games. The final saliency map P is normalized
 735 as $P(i, j) = \frac{S_\pi(i, j)}{\sum_{i, j} S_\pi(i, j)}$.

736 We provide the visualization and our analysis on the following pages.

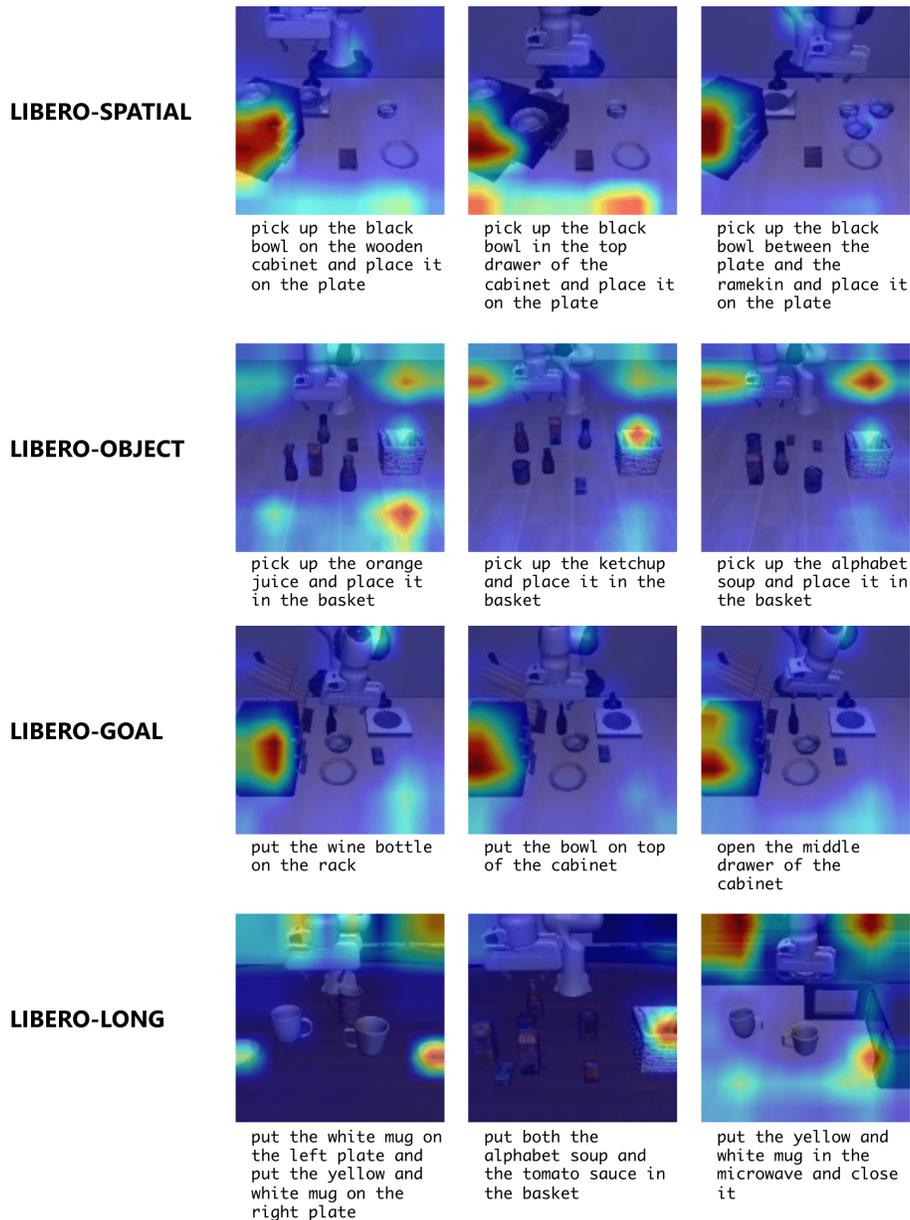


Figure 21: Attention map comparison among different task suites with ER and RESNET-T. Each row corresponds to a task suite.

739 **Findings:** Figure 21 shows attention visualization for 12 tasks across 4 task suites (e.g., 3 tasks per
740 suite). We observe that:

- 741 1. policies pay more attention to the robot arm and the target placement area than the target
742 object.
- 743 2. sometimes the policy pays attention to task-irrelevant areas, such as the blank area on the
744 table.

745 These observations demonstrate that the learned policy use perceptual data for decision-making
746 in a very different way from how humans do. The robot policies tends to spuriously correlate
747 task-irrelevant features with actions, a major reason why the policies overfit to the tasks and do not
748 generalize well across tasks.

749
750

The Same Task over the Course of Lifelong Learning

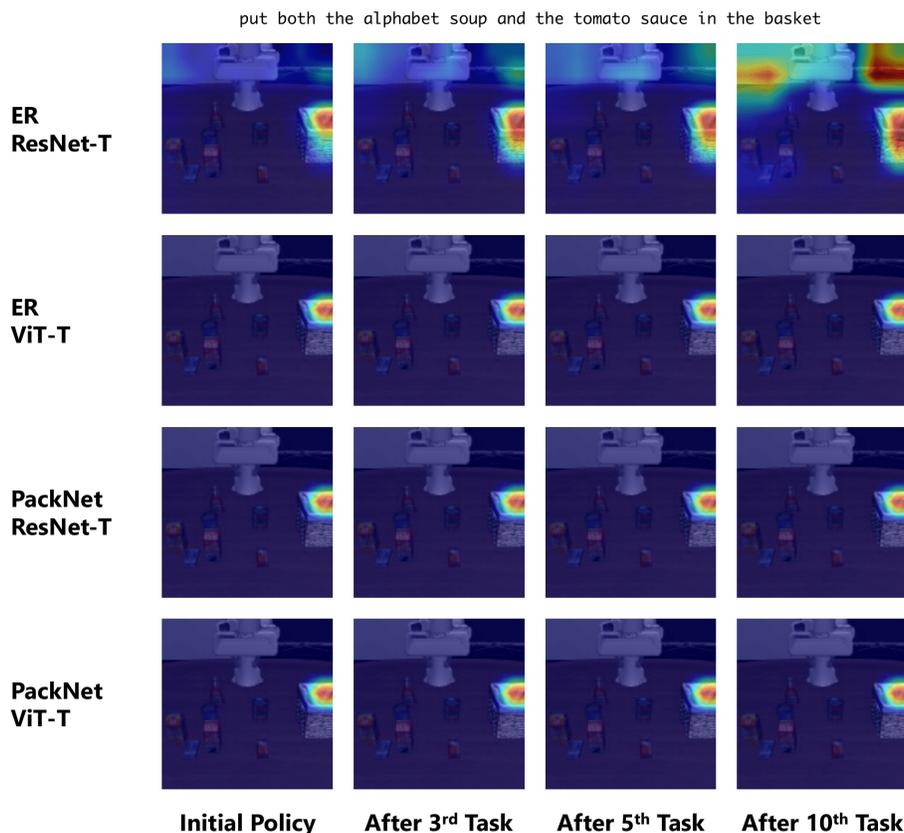


Figure 22: Attention map of the same state of the task *put both the alphabet soup and the tomato sauce in the basket* from LIBERO-LONG during lifelong learning. Each row visualizes how the attention maps change on the first task with one of the LL algorithms (ER and PACKNET) and one of the neural architectures (RESNET-T and ViT-T). Initial policy is the policy that is trained on the first task. And all the following attention maps correspond to policies after training on the third, fifth, and the tenth tasks.

751 **Findings:** Figure 22 shows attention visualizations from policies trained with ER and PACKNET
752 using the architectures RESNET-T and ViT-T respectively. We observe that:

- 753 1. The ViT visual encoder’s attention is more consistent over time, while the ResNet encoder’s
754 attention map gradually dilutes.
- 755 2. PackNet, as it splits the model capacity for different tasks, shows a more consistent attention
756 map over the course of learning.

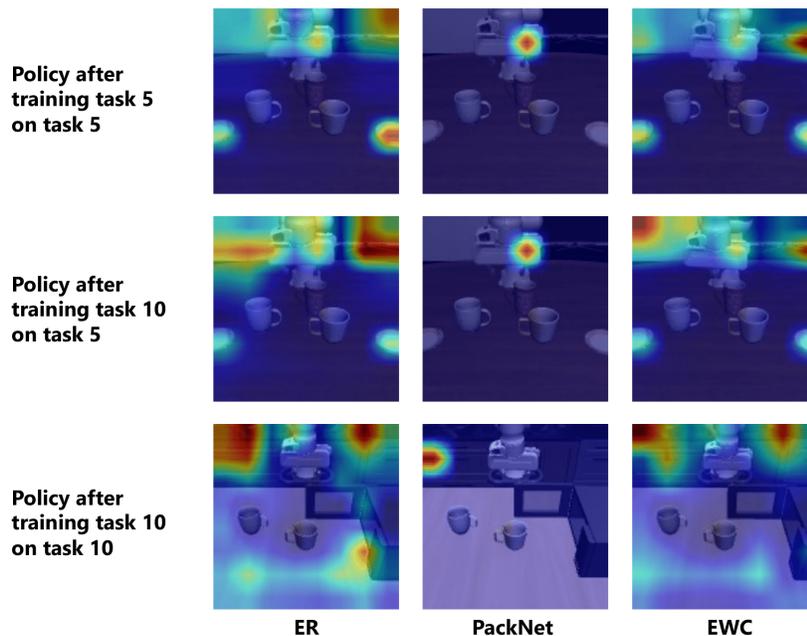


Figure 23: Comparison of attention maps of different lifelong learning algorithms with RESNET-T on LIBERO-LONG. Each row shows the same state of a task with different neural architectures. “Task 5” refers to the task *put the white mug on the left plate and put the yellow and white mug on the right plate*. “Task 10” refers to the task *put the yellow and white mug in the microwave and close it*. The second row shows the policy that is trained on task 10 and gets evaluated on task 5, showing the attention map differences in backward transfer.

759 **Findings:** Figure 23 shows the attention visualization of three lifelong learning algorithms on
 760 LIBERO-LONG with RESNET-T on two tasks (task 5 and task 10). The first and third rows show the
 761 attention of the policy on the same task it has just learned. While the second row shows the attention
 762 of the policy on the task it learned in the past. We observe that:

- 763 1. PACKNET shows more concentrated attention compared against ER and EWC (usually just
 764 a single mode).
- 765 2. ER shares similar attention map with EWC, but EWC performs much worse than ER.
 766 Therefore, attention can only assist the analysis but cannot be treated as a criterion for
 767 performance prediction.

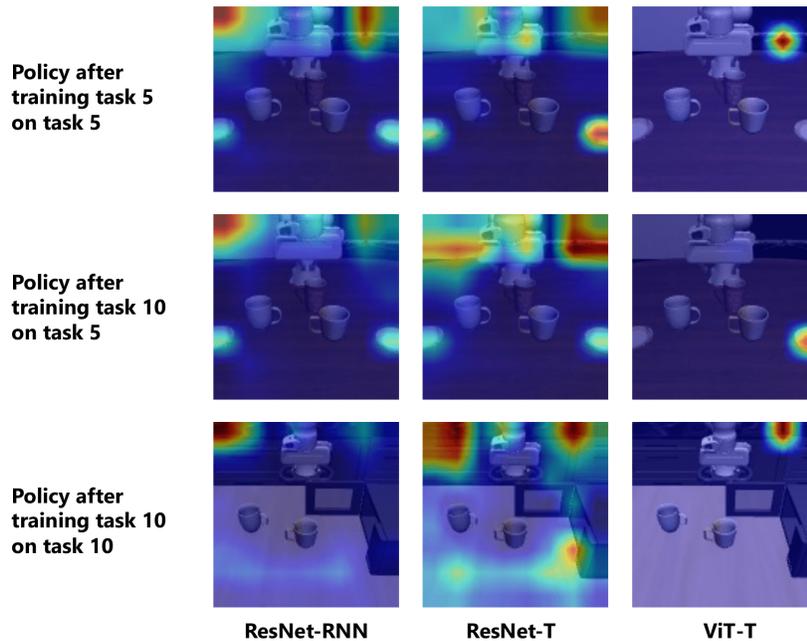


Figure 24: Comparison of attention maps of different neural architectures with ER on LIBERO-LONG. Each row shows the same state of a task with different neural architectures. “Task 5” refers to the task *put the white mug on the left plate and put the yellow and white mug on the right plate*. “Task 10” refers to the task *put the yellow and white mug in the microwave and close it*. The second row shows the policy that is trained on task 10 and gets evaluated on task 5, showing the attention map differences in backward transfer.

770 *Findings:* Figure 24 shows attention map comparisons of the three neural architectures on LIBERO-
771 LONG with ER on two tasks (task 5 and task 10). We observe that:

- 772 1. ViT has more concentrated attention than policies using ResNet.
- 773 2. When ResNet forgets, the attention is changing smoothly (more diluted). But for ViT, when
774 it forgets, the attention can completely shift to a different location.
- 775 3. When ResNet is combined with LSTM or a temporal transformer, the attention hints at the
776 "course of future trajectory". But we do not observe that when ViT is used as the encoder.

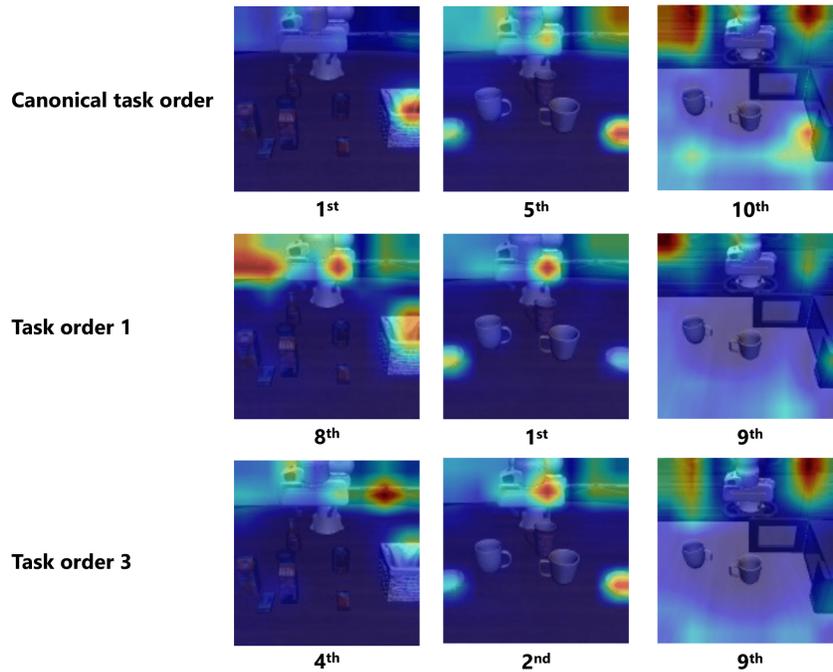


Figure 25: Attention map comparison among different orderings with ER and RESNET-T on three selected tasks from LIBERO-LONG: *put both the alphabet soup and the tomato sauce in the basket*, *put the white mug on the left plate and put the yellow and white mug on the right plate*, and *put the yellow and white mug in the microwave and close it*. Each row corresponds to a specific sequence of task ordering, and the caption of each attention map indicates the order of the task in that sequence.

779 **Findings:** Figure 25 shows attention map comparisons of three different task orderings. We show two
780 immediately learned tasks from LIBERO-LONG trained with ER and RESNET-T. We observe that:

- 781 1. As expected, learning the same task at different positions in the task stream results in
782 different attention visualization.
- 783 2. There seems to be a trend that the policy has a more spread-out attention when it learns on
784 tasks that are later in the sequence.

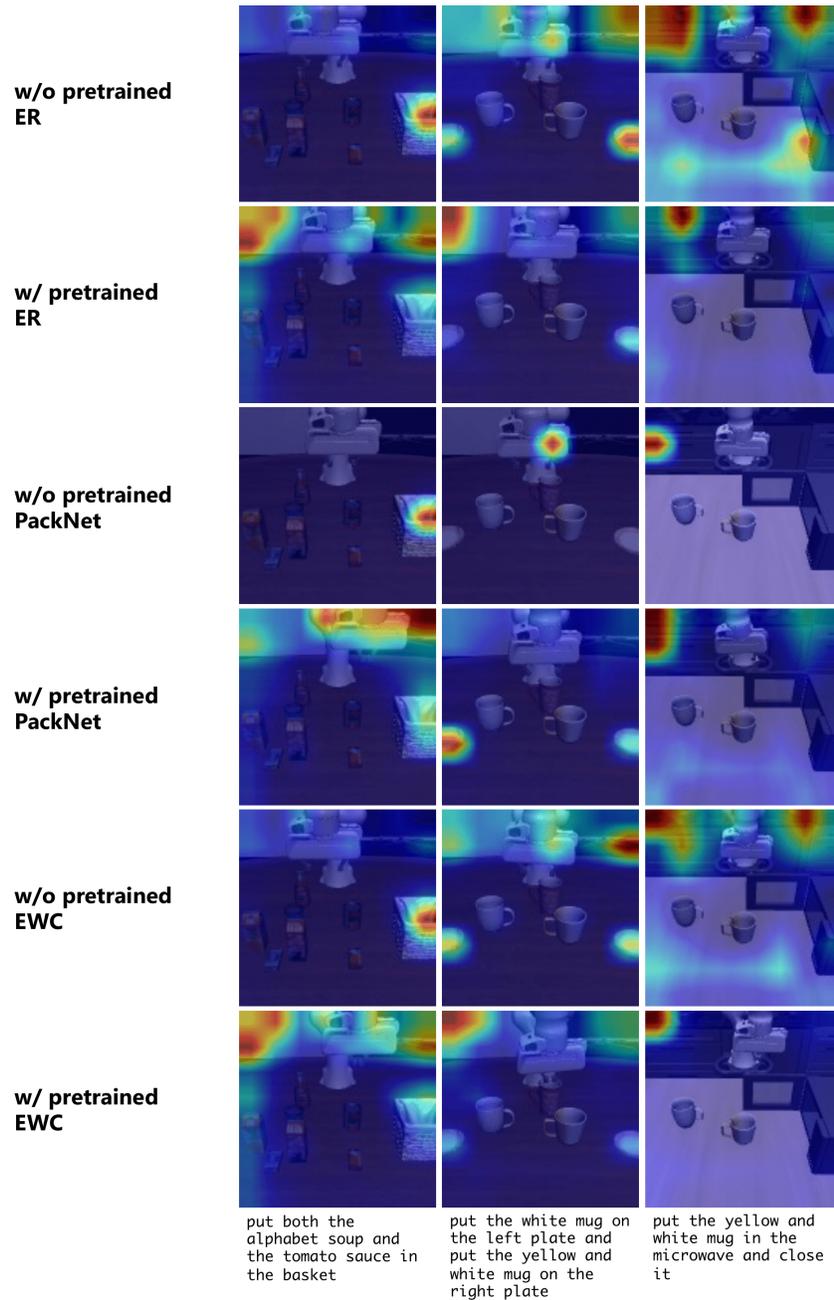


Figure 26: Attention map comparison between models without/with pretrained models using RESNET-T and different lifelong learning algorithms on three selected tasks from LIBERO-LONG.

787 *Findings:* Figure 26 shows attention map comparisons between models with/without pretrained
788 models on LIBERO-LONG with RESNET-T and all three LL algorithms. We observe that:

- 789 1. With pretraining, the policies attend to task-irrelevant regions more easily than those without
790 pretraining.
- 791 2. Some of the policies with pretraining have better attention to the task-relevant features than
792 their counterparts without pretraining, but their performance remains lower (the last in the
793 second row and the second in the fourth row). This observation, again, shows that there is
794 no positive correlation between semantically meaningful attention maps and the policy's
795 performance.