

FMB: a Functional Manipulation Benchmark for Generalizable Robotic Learning

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** In this paper, we propose a benchmark for studying robotic learning
2 for functional manipulation. We identify handling complex contact dynamics and
3 generalization as two central challenges in practical robotic manipulation. While
4 many prior works have addressed one challenge or the other, few have studied both
5 in combination. We hypothesize that making progress on the combination of these
6 challenges requires a set of real-world benchmark tasks that balance complexity
7 with accessibility, providing a set of tasks that are sufficiently narrowly scoped
8 that models and datasets of reasonable scale can be used to make progress, but
9 sufficiently varied that they present a meaningful generalization challenge not just
10 in terms of basic and imprecise skills such as grasping, but also more complex
11 and precise behaviors that require functional manipulation, such as repositioning
12 and reorienting an object for a precise assembly task. Our functional manipulation
13 benchmark consists of a variety of 3D printed objects that can be reproduced pre-
14 cisely by other researchers, each one requiring a sequence of grasping, reorientation,
15 and assembly behaviors. Generalization can be evaluated on test objects and varied
16 positions, as well as more complex multi-stage assembly tasks. We also provide an
17 imitation learning system that provides a basic set of policies for each skill, allow-
18 ing researchers to use our tasks as a toolkit for studying any portion of the pipeline
19 – for example by proposing a better design for a grasping controller and evaluating
20 it in combination with our baseline reorientation and assembly controllers. Our
21 dataset, object CAD files and evaluation videos can be found on our project website:
22 <https://sites.google.com/view/manipulationbenchmark>

23 **Keywords:** manipulation, imitation learning, benchmarking

24 1 Introduction

25 Manipulation is one of the foundational problems in robotics research, but enabling robots to perform
26 dexterous manipulation skills that reflect the capabilities of humans is still out of reach. In fact,
27 even matching the performance of human *teleoperation* remains a major challenge, particularly
28 in environments that require generalization and are not constrained to a specific fixed set of well-
29 characterized objects. As Cui and Trinkle [1] point out, two primary sources of difficulty in robotic
30 manipulation lie in handling complex contact mechanics and intelligently handling the variability
31 in the environment and objects. While robotic learning techniques hold potential to address these
32 challenges, effective progress will require tasks that are accessible enough for current methods while
33 still exposing the key challenges of complex contact mechanics and object generalization.

34 While significant recent research in robotic learning has made progress on various aspects of the
35 manipulation problem [2, 3], much of the emphasis on recent works has either been on broad gener-
36 alization with relatively simple tasks, which often do not capture the many physical challenges of
37 manipulation (e.g., focusing on picking or imprecise pick-and-place tasks) [4], or else training policies
38 for narrow tasks that are physically more complex but not do demand extensive generalization [5].

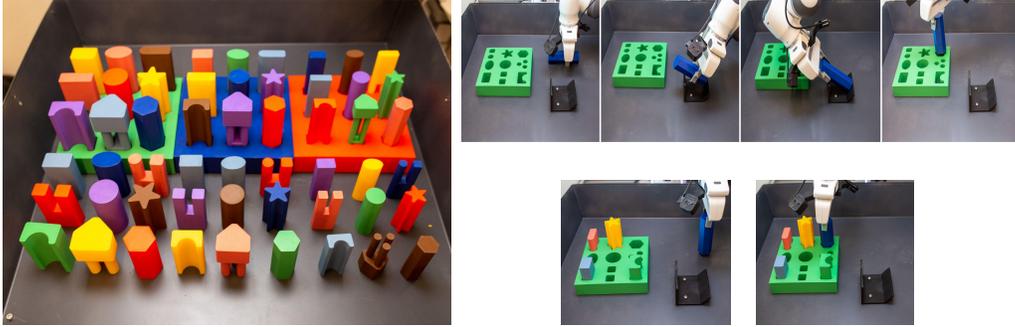


Figure 1: **Left:** The 3D-printed parts for the simpler insertion tasks. **Right:** An illustration of the steps for inserting a single part, which requires grasping the part, reorienting it (potentially using an environment fixture), and then inserting it into the appropriate slot. Note that the full task requires grasping, reorientation, and insertion to be performed in concert.

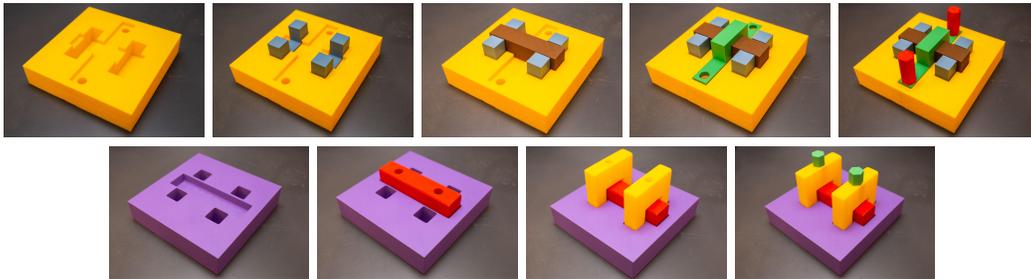


Figure 2: Two instantiations of the complex assembly task. These tasks require similar functional manipulation behaviors as the simpler set of tasks, but with multiple interlocking objects and a more complex higher-level structure that requires assembling the parts in the right order.

39 This is not unreasonable: it is very difficult to simultaneously make progress on broad generalization
 40 (which often requires huge datasets), and tackle the full physical complexity of dexterous manipula-
 41 tion. So how can we take a step toward facilitating robotic learning research that emphasizes both
 42 generalization and physically intricate skills, while still keeping the problem constrained enough so
 43 as to enable meaningful progress?

44 In this paper, we propose a family of benchmarks that aims to cover the important dimensions of
 45 physical complexity and object generalization, while still providing a degree of accessibility by
 46 carefully restricting the scope to a domain where we can make progress with reasonable sized datasets
 47 and models. We approach the design of this benchmark by defining functional object manipulation as
 48 the problem of picking up an object in a functionality relevant way, positioning it in an appropriate
 49 pose, and then using it for a physical interaction. While this definition is more restrictive, we believe
 50 it captures a broad range of practical manipulation tasks, and includes both the challenges of complex
 51 contact dynamics and object generalization.

52 The specific tasks we instantiate to capture functional manipulation are themed around assembly
 53 problems, including simpler pick-and-place tasks and more complex multi-part assemblies. These
 54 tasks, illustrated in Fig. 1, require picking up the individual pieces, reorienting them (potentially using
 55 environment affordances and regrasping), and then slotting them into their required location. Each
 56 phase requires addressing both the challenge of complex contacts and the challenge of generalization.
 57 The objects may vary between training and test-time, and their locations are randomized. The grasping
 58 phase requires selecting a grasp that is suitable for reorienting the object, the reorientation phase
 59 requires positioning the object so that contact with the environment changes its pose in the desired
 60 way, and the assembly phase requires compliant insertion and proper accounting for the contact
 61 forces on the object. Each phase requires handling different objects (including new test-time objects)
 62 and different poses. The robotic assembly task has been long seen as a representative manipulation

63 benchmark [6, 7, 8, 9]; however, the generalization effect across such tasks has been less studied
64 comprehensively. Thus, performing such tasks among the pool of diverse shapes would be an ideal
65 candidate for benchmarking generalizable dexterity.

66 To ensure reproducibility and portability of our benchmark task, we use 54 3D-printed objects with
67 diverse shapes and sizes that can be reproduced by other researchers, and a widely used Franka
68 robotic arm. We collected a dataset of 9000 human demonstrations of grasping, repositioning, and
69 inserting these objects, and trained a baseline imitation learning system to perform each stage of the
70 task. Our dataset also contains a variety sensory modalities as presented in Fig. 3: we record RGB
71 and depth images from eye-in-hand and eye-to-hand views. These make our environment modular so
72 that other researchers can repurpose it for a variety of methods that they may wish to develop, and
73 can focus on any stage or aspect of the task. For example, some researchers might choose to focus on
74 better functional grasping methods, while the other stages are handled by our baseline system, while
75 others might focus on compliant insertion, utilizing our baseline system for the grasping stage. Our
76 tasks are also designed to accommodate pretraining with finetuning to other downstream behaviors.
77 To this end, we also provide a set of more complex assembly objects, as shown in Fig. 2, which can be
78 handled by policies adapted from pretraining on the main dataset. We describe our benchmark tasks,
79 and conduct a comprehensive evaluation studying the performance of imitation learning methods
80 trained on our data, evaluating both training object and test set performance. Our hope is that our
81 functional manipulation benchmark (FMB) will provide a toolkit for robotic learning researchers to
82 study manipulation both in terms of complex contact dynamics and generalization.

83 2 Related Work

84 Considerable recent progress on robotic manipulation has studied generalization, though often in
85 the context of simpler tasks such as grasping [10, 2], pushing [10], and imprecise repositioning [10].
86 A number of other works have studied tasks that are dynamic [11], precise (e.g., insertion) [12], or
87 otherwise physically challenging [5]. However, few works have studied these factors in combination.
88 We believe many of the central challenges in robotic manipulation lie at the confluence of these
89 two challenges: tasks that require handling complex contact dynamics, not by memorizing the
90 particular pattern needed for a single narrow task, but by learning general behaviors for handling
91 object interaction that can generalize to new objects. Our aim is to propose a benchmark that can study
92 this combination of challenges, while keeping the scope narrow enough that it remains accessible to
93 many researchers.

94 Our tasks combine aspects of grasping, repositioning, and peg insertion or assembly. A number
95 of works have studied these individual stages [2, 13]. Our goal is not to attain the best possible
96 performance in narrow settings for any of these stages (e.g., ultra-high-precision industrial insertion),
97 but to use these tasks as a lens through which to gauge general manipulation capabilities learned via
98 general-purpose robotic learning methods.

99 A number of prior works have proposed datasets for robotic learning, including datasets consisting of
100 demonstrations [4] and autonomously collected data [2, 14], as well as annotated datasets of grasp
101 points [15], object geometries [16, 17], and simulated environments [18]. However, there has been
102 comparatively little work on standard and accessible object sets that are combined with multi-stage
103 tasks for studying generalization. The YCB object set [19] comes with a number of evaluation
104 protocols [19], but these protocols generally focus on object repositioning tasks that do not evaluate
105 the complex contacts challenges that we discuss in the previous section. A number of existing
106 demonstration datasets cover many different behaviors [4, 20], but also focus on behaviors that
107 emphasize basic pick-and-place skills rather than precise or contact-rich manipulation. Some works
108 have focused on insertion skills in particular (e.g., connector insertion) [21]. While our benchmark
109 is related, we aim specifically to cover a range of skills, including grasping and repositioning, that
110 we believe cover a basis of basic manipulation capabilities. We also emphasize generalization as a
111 primary challenge for our benchmark.

112 We use 3D printed objects to facilitate reproducibility. Other prior works have also proposed standard
113 meshes and 3D printed parts for benchmarking and reproducibility [19], typically focusing on object

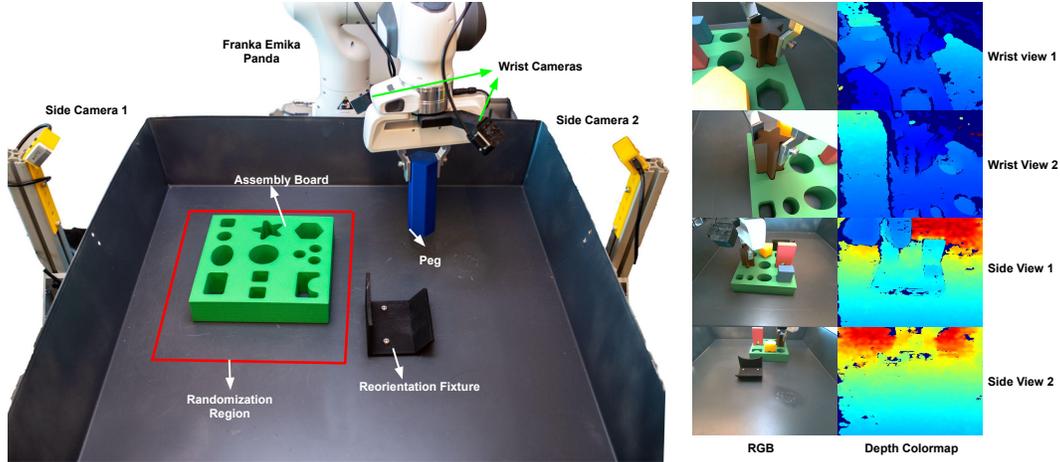


Figure 3: Illustration of the robot setup, with a standard Franka arm equipped with four cameras (two on the wrist and two attached to the environment), each with RGB and depth, positioned in front of a workspace containing an object, reorientation fixture, and assembly board. The board is placed into a random pose within the randomization region, and the object is located in a randomized pose on the table, from where it must be picked up, reoriented, and inserted.

114 grasping. These efforts are related, but our aim is to provide parts that are specifically well suited for
 115 evaluating all of the stages: grasping, reorientation, and assembly, rather than only grasping.

116 3 Functional Manipulation Benchmark

117 In this section, we introduce the basic principles behind FMB and the protocols to evaluate different
 118 methods on this benchmark. We are mostly concerned with studying the generalization of each
 119 individual functional manipulation task as well as the combinatorial ways of composing them to
 120 achieve novel behaviors. Therefore we collect a diverse dataset of robotic behaviors with different
 121 objects, viewpoints, and robot initial poses. We also additionally provide novel objects for the purpose
 122 of benchmarking the generalization capability of individual skills, as well as the ability for a method
 123 to compose these skills to solve unseen long-horizon tasks.

124 3.1 Object Set

125 We designed 54 3D-printed objects of different sizes, shapes, and colors, with examples shown in
 126 Figure 1.

127 In total, we have 9 different basic shapes, and for each shape there are 6 different sizes. The parts
 128 are assigned 8 different random colors. There are three boards with matching holes for the objects.
 129 We additionally designed two more complex boards to facilitate multi-stage assembly tasks, shown
 130 in Figure 2, where multiple parts must be fitted together. The tolerance for mating these objects
 131 is consistently 1mm to 1.5mm. All of our CAD files including those for environment fixtures and
 132 camera mount are publicly available on our project website.

133 3.2 Functional Manipulation Tasks

134 In this section, we describe the individual tasks that we propose to evaluate with our benchmark.
 135 For each type of tasks, we provide demonstration trajectories collected with a Franka robot (see
 136 Figure 3), and an evaluation protocol. The modular design of our benchmark facilitates extension
 137 to add new tasks with the provided objects, but the tasks we describe here are suitable both for
 138 evaluating generalization and for testing a range of manipulation capabilities.

139 **Grasping.** The grasping task in our benchmark is a *functional* grasping task, in the sense that the
 140 robot must grasp the object in a way that facilitates downstream reorientation, rather than simply
 141 picking the object in any pose. We illustrate this task in Fig. 4. A top-down grasp is reasonable if the
 142 object is placed in a vertical pose, as shown on the right side of Fig. 4. However, a horizontal grasp is



Figure 4: Objects may need to be grasped from a variety of poses, particularly when using the reorientation fixture, where they might lie at an angle.

143 much more desirable if the object is positioned as on the left side of Fig. 4. In case such a grasp is
 144 infeasible due to the robot’s kinematic constraints, the robot needs to perform additional repositioning
 145 steps to adjust the feasible grasp pose. The robot must learn grasping skills that deploy the appropriate
 146 grasp for the object’s current configuration, and also generalize across different object shapes, colors,
 147 and sizes. Our demonstration dataset for the grasping task consists of 50 trajectories per object, with
 148 varying object rest poses in the bin, for a total of 2700 trajectories performing functional grasping
 149 over the 54-object set.

150 **Repositioning.** A repositioning step is sometimes necessary to adjust the grasping pose so that the
 151 object is held in a way that is suitable for downstream assembly. Manipulating and reorienting objects
 152 by leveraging environment affordances (e.g., tilting the object in the gripper by levering it against a
 153 table or wall) may often be necessary for fluent and complex manipulation, and this reorientation task
 154 exercises this capability. We provide a simple fixture that can serve as an environment affordance
 155 to rest the object at angle, as shown in Fig. ???. To reorient the objects into the right pose, the robot
 156 may need to use this fixture, resting the object on it and then regrasping it in a more appropriate
 157 pose for reorientation. We collected 3000 demonstrations for placing and regrasping, which can be
 158 used to learn strategies for using environmental affordances for regrasping and reorientation. Since
 159 objects land in the fixture in a relatively deterministic fashion, we partially script our demonstration
 160 collection process while maintaining a certain degree of randomness for the purpose of data diversity.

161 **Assembly.** Our assembly tasks consider assembling objects of
 162 diverse shapes into their matching slots, which requires performing
 163 fine-grained precise manipulation. An illustrative example is shown
 164 in Fig. 5. Here, having completed the preceding two steps, the robot
 165 is holding an object, and needs to insert it into the matching slot in
 166 the blue board. For each object, we collect 50 human demonstrations
 167 that include various robot initial poses and board positions, for a
 168 total of 3000 demonstrations performing the assembly task from
 169 various initial conditions. Note that the board is located in different
 170 places on the table for different episodes, requiring a reactive strategy that localizes the board and the
 171 appropriate opening, and guides the object into the correct location.

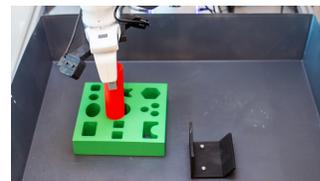


Figure 5: An illustration of the assembly task.

172 **Long-horizon manipulation.** Aside from performing individual steps, such as grasping, reorienta-
 173 tion, and assembly, our benchmark and demonstrations can be used to learn the entire long-horizon
 174 sequence, performing the steps in turn to insert one or multiple objects into the board. The difficulty
 175 of this task mainly comes from the compounding errors accumulated over each individual step which
 176 gets even more magnified when switching between tasks.

177 **Multi-step interlocking assembly.** We also present two sets of novel objects for benchmarking
 178 much broader generalization capability. The pieces in Fig. 2 are largely different from our original
 179 set of objects, and would require adaptation to perform grasping or insertion which can be achieved
 180 by pretraining on the collected dataset and finetuning on a few new demonstrations. The major
 181 challenge with these tasks is that these objects need to be put together in a specific order, such as in
 182 an interlocking fashion. While it may not be too hard to perform individual steps alone, the difficulty
 183 increases rapidly when a policy needs to simultaneously reason the manipulation sequence as well as
 184 accounting for compounding manipulation errors introduced by individual steps.

185 3.3 Robotic system and data collection

186 We now describe the robotic system and the process we used to collect the training data.

187 **Robotic system overview** Our system can be seen in Fig. 3. We use a Franka Panda robot to collect
188 our dataset, since it is widely adopted for research and offers a torque control interface which is
189 very desirable in contact-rich manipulation tasks. To record demonstrations, we use a SpaceMouse
190 to control the robot at 10 HZ. In total, we have four Intel RealSense D405 cameras, two of which
191 are mounted on the robot end-effector, and the rest are placed on each side of the bin to provide a
192 complementary view of objects in the bin. We concurrently capture RGB and depth images from
193 these cameras.

194 **Data collection protocol.** Our dataset consists of 2700 demonstrations for the full long-horizon
195 task of grasping, reorientation, and assembly for these 54 objects. For each object, we collect around
196 50 demonstrations per task. Each such demonstration trajectory is around 20 to 30 seconds long, and
197 thus it’s more practical to break them into individual “primitives” of shorter horizons. In fact, we
198 automatically add indicators at the end of a manipulation skill such as grasping so that we can segment
199 these long-horizon trajectories. In our dataset, these primitives include grasping, reorientation,
200 move, insertion; so in that sense, we have 8100 demonstrations of each primitive with horizons
201 around 5 seconds. For the grasping task, the object of interest is randomized around a 20cm x 30cm
202 rectangular area in the bin; whereas for the insertion task, the board is randomized around a 40cm
203 by 60cm area. We also include distractors (i.e. objects not needed for a task) when performing the
204 insertion task, half of the insertion demonstrations were carried out when there are distractors present
205 to gain robustness.

206 4 Using the FMB in imitation learning

207 To illustrate the utility of our benchmark in imitation learning. We describe a few example usages
208 of our dataset and the corresponding evaluation protocol. The detailed evaluation protocol and met-
209 ric can be found on our website <https://sites.google.com/view/manipulationbenchmark>.
210 Although in principle our data can also be easily altered to study other approaches such as offline
211 reinforcement learning.

212 4.1 Training and Evaluation of Individual Skills

213 Generally speaking, we expect to see the emergence of generalization by training on a large, diverse
214 dataset. To verify this hypothesis, we refer to two ways of testing generalization. For grasping and
215 insertion, we can hold out a specific object in the training set, train a policy without seeing any data
216 associated with that object, and then test on the held-out object. Alternatively, we also provide five
217 novel objects that are not contained in the dataset for which we can directly evaluate trained policies.

218 4.2 Pre-training and Finetuning

219 Pretraining and finetuning visuomotor skills is an open and important research question. By having a
220 large-scale diverse dataset of robot manipulation behaviors, it’s possible for us to study this problem.
221 We can pretrain on a set of robot behaviors associated with some objects, and then finetune on data
222 from objects that are not present in the pretraining dataset. If that object is entirely novel, such as the
223 more complex assembly objects in Figure 2, we can collect some additional demonstrations using our
224 setup for finetuning.

225 4.3 Composing Skills to Solve Long-Horizon tasks

226 FMB also supports studying long-horizon tasks in various ways: one can train “flat” style imitation
227 learning methods on all the data or hierarchical style methods that trigger individual primitives in
228 some intelligent ways. In addition to the original “grasp-reorient-assembly” task, it’s also possible
229 to study more complex novel tasks such as the one shown in Fig. 2, by finetuning on the new objects.



Figure 6: Unseen test objects used for evaluating generalization in our protocol. The robot must generalize to new combinations of shapes, colors, and sizes using the diverse training set.

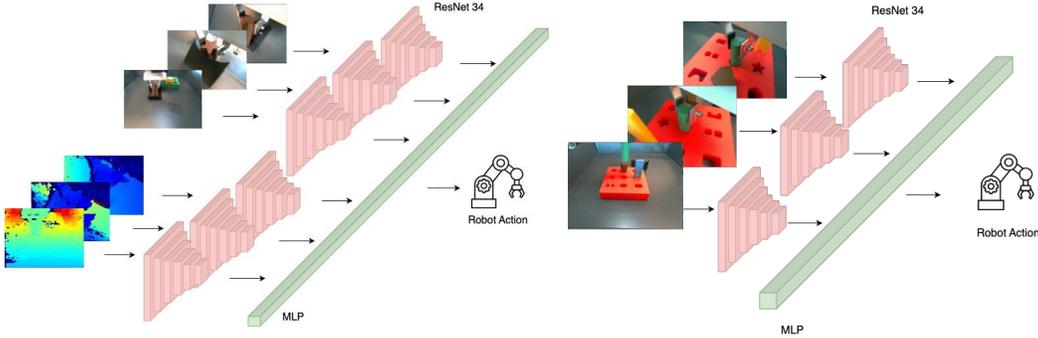


Figure 7: Architecture diagrams for the grasping and reorientation tasks (left), and the assembly task (right). The models encode each observation with a ResNet34 encoder, and then fuse the modalities with fully connected layers.

230 5 Experiments

231 In this section, we conduct experiments first to verify our proposed tasks are actually feasible. We
 232 then seek to answer the following questions: (1) For each individual task, does training on a diverse
 233 manipulation dataset generalize across object properties? (2) When do multi-modal inputs help for
 234 which manipulation skill? (3) What are the necessary ingredients for solving long-horizon complex
 235 manipulation tasks?

236 5.1 Grasping Task

237 To show that it is feasible to perform the grasping tasks using
 238 our dataset, we first train a grasping BC policy specifically for the oval object. We obtain 12 successful grasps out of 30 trials,
 239 which amounts to 40% success rate. During the evaluation, we
 240 test on all six oval objects, performing five trials per each object
 241 so that generalization can be fairly tested.
 242

243 Then we train two grasping BC policies on all the data and
 244 test the trained policy on both in-distribution objects as well
 245 as novel objects. We present results in Fig. 8. One grasping
 246 policy is trained with RGB images, the other one we provide
 247 additional depth information; their neural network architecture
 248 can be seen in Fig. 7. We find that depth information is crucial
 249 in helping achieve better grasping performance.

250 5.2 Repositioning Task

251 For the repositioning task, we train BC policies to first place
 252 the object on the fixture and then try to re-grasp the object from
 253 the other end. The policy’s success rate is 0% if trained to solve
 254 place and reorient at once with all the data. If we train only
 255 on placing data, the policy can achieve 33.3% success rate out
 256 of 30 trials; however, the re-grasping policy is 5% success rate
 257 trained on corresponding re-grasping data. The failure mode
 258 includes missing the object, flipping over the object, and the

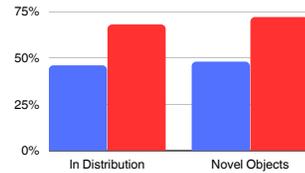


Figure 8: Comparison of grasping success rates on in distribution and novel objects when trained with (Red) and without (Blue) depth information.

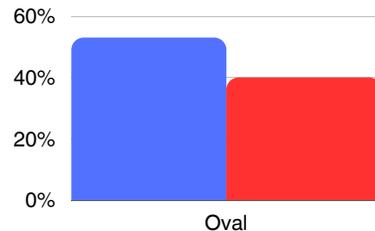


Figure 9: Success rates of grasping oval when training without oval data (Blue) and only oval data (Red).

259 arm wasn't able to turn over. This is reasonable since we only train policies on visual data without
 260 robot state, and the re-grasping part is particularly challenging due to its multi-modal nature.

261 5.3 Assembly Task

262 For the assembly task, we first train three separate insertion
 263 policies on the shapes of round, hexagon, and two-square. We
 264 train these policies on a small portion of our dataset that only
 265 contains their corresponding shapes, e.g., the policy to insert
 266 round objects is trained on only round data. We also vary each
 267 policy's input modality differentiating by depth information.
 268 The results can be seen in Fig. 10. We can see the success
 269 rate does decrease as the shape becomes more complex, this
 270 implies the chosen assembly task is indeed a challenging robotic
 271 manipulation task thus worthwhile benchmarking. We also
 272 find that naively adding depth information doesn't help for the
 273 insertion tasks.

274 To carry out an initial study on the generalization of training
 275 on diverse shapes, we train an insertion BC policy with all the
 276 data and test it on the round and star shapes; which we get 4/30
 277 and 0/30 success rates respectively. However, when we train
 278 individual policies just with data from that particular shape,
 279 we are able to get 9/30 and 0/30. This is reasonable because
 280 this task is very precise, naively mixing the data will cause the
 281 uni-modal BC model confused so that performance gets hurt.
 282 We present this result in Fig. 11.

283 5.4 Long-Horizon Task

284 In addition to training BC policies on only primitives, we train
 285 an end-to-end long-horizon Behavioral Cloning (BC) policy
 286 with all the long horizon demos and transitions. We provide all
 287 the RGB camera views for this policy. The goal of this policy
 288 is to successfully grasp, reorient, and perform assembly. We
 289 evaluate the end-to-end Behavioral Cloning policy on 5 different
 290 pegs for a total of 10 trials; this policy achieves a success rate
 291 of 0/10. This is reasonable due to the accumulation of errors in
 292 long-horizon end-to-end behavioral cloning.

291 We explore alternatives to naive Behavioral Cloning and modify
 292 our approach to include our previously trained grasping and
 293 insertion policies. By manually triggering the grasp and insertion
 294 policies, as well as using a scripted reorientation motion, we
 295 achieve a success rate of 2/10. We present this result in Fig. 12.

296 6 Discussion and Limitation

297 In this paper, we present a benchmark for functional manipu-
 298 lation. We open-source all the data as well as the object CAD
 299 files to facilitate reproducibility. We evaluate imitation learning
 300 methods with different input modalities and their abilities to
 301 generalize across objects. We hope our benchmark FMB would
 302 encourage and contribute to in robotic manipulation research.

303 **Limitations and future work.** Although our dataset has a
 304 variety of diverse objects, we still only have one scene; it would
 305 be helpful if we can include more background scenes. Addi-
 306 tionally, 3D-printed objects are easy to reproduce, however, it
 307 would be more useful if we include real standardized objects in
 the future so we can study much broader generalization.

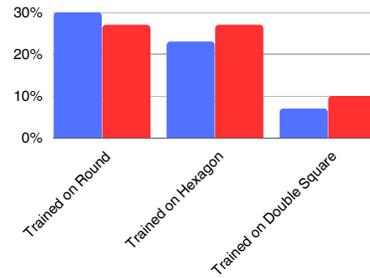


Figure 10: Comparing insertion success rates when training only on one peg data with (Red) and without (Blue) depth then evaluating only on trained peg.

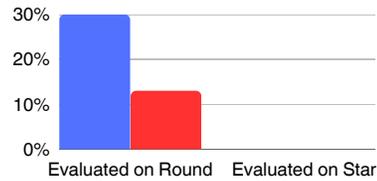


Figure 11: Comparing insertion on specific pegs when training on round with RGB only (Blue) vs training on all data with RGB only (Red).

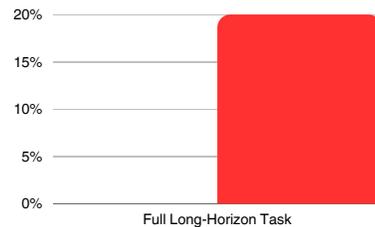


Figure 12: Comparison between end-to-end behavioral cloning (Blue) for grasping, reorientation, and assembly and manually triggered learned skills (Red) and scripted skills

References

- 308
- 309 [1] J. Cui and J. Trinkle. Toward next-generation learned robot manipulation. *Sci-*
310 *ence Robotics*, 6(54), 2021. URL [https://www.science.org/doi/abs/10.1126/](https://www.science.org/doi/abs/10.1126/scirobotics.abd9461)
311 [scirobotics.abd9461](https://www.science.org/doi/abs/10.1126/scirobotics.abd9461).
- 312 [2] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen. Learning hand-eye coordination for robotic
313 grasping with deep learning and large-scale data collection, 2016.
- 314 [3] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakr-
315 ishnan, V. Vanhoucke, and S. Levine. Qt-opt: Scalable deep reinforcement learning for vision-
316 based robotic manipulation, 2018.
- 317 [4] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and
318 S. Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets.
319 *arXiv preprint arXiv:2109.13396*, 2021.
- 320 [5] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino,
321 M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng,
322 Q. Yuan, W. Zaremba, and L. Zhang. Solving rubik’s cube with a robot hand, 2019.
- 323 [6] I.-W. Kim, D.-J. Lim, and K.-I. Kim. Active peg-in-hole of chamferless parts using
324 force/moment sensor. *Proceedings 1999 IEEE/RSJ International Conference on Intelligent*
325 *Robots and Systems. Human and Environment Friendly Robots with High Intelligence and*
326 *Emotional Quotients (Cat. No.99CH36289)*. doi:10.1109/iros.1999.812802.
- 327 [7] W. Newman, Y. Zhao, and Y.-H. Pao. Interpretation of force and moment signals for compliant
328 peg-in-hole assembly. *Proceedings 2001 ICRA. IEEE International Conference on Robotics*
329 *and Automation (Cat. No.01CH37164)*. doi:10.1109/robot.2001.932611.
- 330 [8] V. Gullapalli, R. Grupen, and A. Barto. Learning reactive admittance control. *Proceedings 1992*
331 *IEEE International Conference on Robotics and Automation*. doi:10.1109/robot.1992.220143.
- 332 [9] M. Majors and R. Richards. A neural-network-based flexible assembly controller. In *1995*
333 *Fourth International Conference on Artificial Neural Networks*, pages 268–273, 1995. doi:
334 [10.1049/cp:19950566](https://doi.org/10.1049/cp:19950566).
- 335 [10] S. Dasari, F. Ebert, S. Tian, S. Nair, B. Bucher, K. Schmeckpeper, S. Singh, S. Levine, and
336 C. Finn. Robonet: Large-scale multi-robot learning, 2020.
- 337 [11] D. Seita, N. Jamali, M. Laskey, A. K. Tanwani, R. Berenstein, P. Baskaran, S. Iba, J. Canny, and
338 K. Goldberg. Deep transfer learning of pick points on fabric for robot bed-making, 2019.
- 339 [12] K. Zakka, A. Zeng, J. Lee, and S. Song. Form2fit: Learning shape priors for generalizable
340 assembly from disassembly, 2020.
- 341 [13] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. A. Ojea, and K. Goldberg. Dex-net
342 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics,
343 2017.
- 344 [14] L. Pinto and A. Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700
345 robot hours, 2015.
- 346 [15] H.-S. Fang, C. Wang, M. Gou, and C. Lu. Graspnet: A large-scale clustered and densely
347 annotated dataset for object grasping, 2020.
- 348 [16] S. Tyree, J. Tremblay, T. To, J. Cheng, T. Mosier, J. Smith, and S. Birchfield. 6-dof pose
349 estimation of household objects for robotic manipulation: An accessible dataset and benchmark,
350 2022.

- 351 [17] J. J. P. Y.-W. Chao, and Y. Xiang. Fewsol: A dataset for few-shot object learning in robotic
352 environments, 2023.
- 353 [18] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell,
354 and K. Bousmalis. Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-
355 to-canonical adaptation networks, 2019.
- 356 [19] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar. Benchmarking in
357 manipulation research: Using the yale-CMU-berkeley object and model set. *IEEE Robotics
358 & Automation Magazine*, 22(3):36–52, sep 2015. doi:10.1109/mra.2015.2448951. URL
359 [https://doi.org/10.1109%2Fmra.2015.2448951](https://doi.org/10.1109/2Fmra.2015.2448951).
- 360 [20] A. Mandlekar, J. Booher, M. Spero, A. Tung, A. Gupta, Y. Zhu, A. Garg, S. Savarese, and
361 L. Fei-Fei. Scaling robot supervision to hundreds of hours with roboturk: Robotic manipulation
362 dataset through human reasoning and dexterity, 2019.
- 363 [21] G. D. Magistris, A. Munawar, T.-H. Pham, T. Inoue, P. Vinayavekhin, and R. Tachibana.
364 Experimental force-torque dataset for robot learning of multi-shape insertion, 2018.