

Motion Policy Networks

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Collision-free motion generation in unknown environments is a core
2 building block for robot manipulation. Generating such motions is challenging
3 due to multiple objectives; not only should the solutions be optimal, the mo-
4 tion generator itself must be fast enough for real-time performance and reliable
5 enough for practical deployment. A wide variety of methods have been proposed
6 ranging from local controllers to global planners, often being combined to offset
7 their shortcomings. We present an end-to-end neural model called Motion Policy
8 Networks ($M\pi$ Nets) to generate collision-free, smooth motion from just a single
9 depth camera observation. $M\pi$ Nets are trained on over 3 million motion plan-
10 ning problems in 500,000 environments. Our experiments show that $M\pi$ Nets are
11 significantly faster than global planners while exhibiting the reactivity needed to
12 deal with dynamic scenes. They are 46% better than prior neural planners and
13 more robust than local control policies. Despite being only trained in simulation,
14 $M\pi$ Nets transfer well to the real robot with noisy partial point clouds.

15 **Keywords:** Motion Control, Imitation Learning, End-to-end Learning

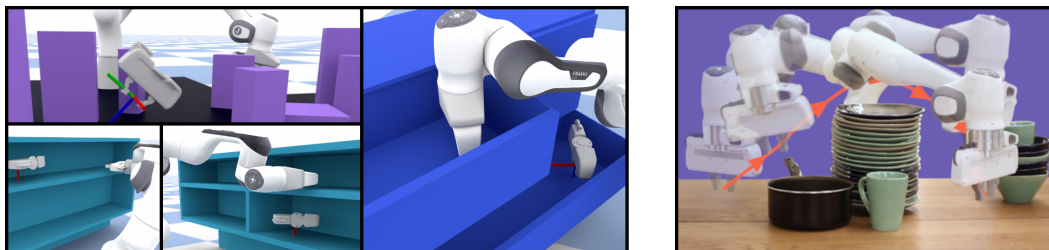


Figure 1: $M\pi$ Nets are trained on a large dataset of synthetic demonstrations (*left*) and can solve complex motion planning problems using raw point cloud observations (*right*).

16 1 Introduction

17 Generating fast and legible motions for a robotic manipulator in unknown environments is still an
18 open problem. Decades of research have established many well-studied algorithms, but there are
19 two practical issues that prevent motion planning methods from being widely adopted in industrial
20 applications and home environments that require real-time control. First, it is challenging for any
21 single approach to satisfy multiple planning considerations: speed, completeness, optimality, ease-
22 of-use, legibility (from the perspective of a human operator), determinism, and smoothness. Second,
23 existing approaches enforce strong assumptions about the visual obstacle representations and hence
24 are not flexible enough to operate in novel scenes directly using raw sensor observations.

25 Global planners such as RRT [1] optimize for speed and completeness but fall short on optimality.
26 Other sampling-based approaches choose to forgo speed for the sake of optimality [2, 3, 4, 5].
27 Similarly, optimization-based approaches [6, 7, 8] embrace locally optimal behavior in exchange for

28 completeness. Recent methods such as Geometric Fabrics [9] and STORM [10] deploy reactive local
29 policies and assume that local decisions will lead to globally acceptable paths. Unfortunately, as we
30 show in our experiments, the performance of these local approaches degrades in more geometrically
31 complex environments as they get stuck in local minima. Motivated by the success of deep learning,
32 neural motion planning approaches such as Motion Planning Networks [11] have been proposed to
33 greatly improve the sampling of an RRT planner with imitation learning. However, they still require
34 a planner and a collision checker with known models at test time.

35 Planners have traditionally been evaluated with known environment models and perfect state esti-
36 mation. When deploying them in practice, however, one would have to create one of several scene
37 representations: a static or dynamic mesh, occupancy grids [12, 13], signed distance fields, etc. Re-
38 construction systems such as SLAM and KinectFusion [14] have a large system start-up time, require
39 a moving camera to aggregate many viewpoints, and ultimately require costly updates in the pres-
40 ence of dynamic objects. Recent implicit deep learning methods like DeepSDF [15] and NERF [16]
41 are slow or do not yet generalize to novel scenes. Methods such as SceneCollisionNet [17] provide
42 fast collision checks but require expensive MPC rollouts at test time. It also draws samples from a
43 straight line path in configuration space which may not generalize to challenging environments be-
44 yond a tabletop. Other RL-based methods learn a latent representation from observations but have
45 only been applied to simple 2D [18, 19] or 3D [20] environments in simulation.

46 We present *Motion Policy Networks (M π Nets)*, a novel method for learning an end-to-end policy for
47 motion planning. Our approach circumvents the challenges of traditional motion planning and is
48 flexible enough to be applied in unknown environments. Our contributions are as follows:

- 49 • We present a large-scale effort in neural motion planning for manipulation. Specifically, we
50 learn from over 3 million motion planning problems across over 500,000 instances of three
51 types of environments, nearly 300x larger than prior work [11].
- 52 • We train a reactive, end-to-end neural policy that operates on point clouds of the environment
53 and moves to task space targets while avoiding obstacles. Our policy is significantly faster than
54 other baseline configuration space planners and succeeds more than local task space controllers.
- 55 • On our challenging dataset benchmarks, we show that M π Nets is nearly 46% better than prior
56 work [11] without even needing the full scene collision model.
- 57 • Finally, we demonstrate *sim2real* transfer to real robot partial point cloud observations.

58 2 Related Work

59 **Global Planning:** Robotic motion planning typically splits into three camps: search, sampling, and
60 optimization-based planning. Search-based planning algorithms, such as A* [21, 22, 23], discretize
61 the state space and perform a graph search to find an optimal path. These algorithms are slow but
62 find least-cost paths. Sampling-based planners [1] function in a continuous state space by drawing
63 samples and building a tree. When the tree has sufficient coverage of the planning problem, the
64 algorithm traverses the tree to produce the final plan. Sampling based planners are typically much
65 faster than search-based planners and some are even *asymptotically optimal* [2, 3, 4], but due to their
66 random nature can produce erratic—though valid—paths.

67 Both of the aforementioned planner types are designed to optimize for path length in the given state
68 space (*e.g.* configuration space) while avoiding collisions. However, optimal paths in configuration
69 space can lead to unintuitive paths in a task space (*e.g.* Cartesian position of the robot’s end effector).
70 Motion Optimization [6, 7, 24] on the other hand, generates paths with non-linear optimization and
71 considers multiple objectives such as smoothness of the motion, obstacle avoidance and convergence
72 to an end effector pose. These algorithms require careful tuning of the respective cost functions to
73 ensure convergence to a desirable path and are prone to local minima. Furthermore, non-linear
74 optimization is computationally complex and can be slow for difficult planning problems.

75 **Local Control:** In contrast to global planners, local controllers have long been applied to create
76 collision-free motions [25, 26, 9, 10]. While they prioritize speed and smoothness, they are highly

77 local and may fail to find a valid path in complex environments. We demonstrate in our experiments
 78 that $M\pi$ Nets are more effective at producing convergent motions in these types of environments,
 79 including in dynamic and in partially observed settings.

80 **Neural Motion Planning:** Many deep planning methods [12, 27, 28, 29] seek to learn efficient sam-
 81 plers to speed up traditional planners. Motion Planning Networks (MPNets) [11] learn to directly
 82 plan through imitation of a standard sampling based RRT* planner [2] and is used in conjunction
 83 with a traditional planner for stronger guarantees. While these works greatly improve the speed of
 84 the planning search, they have the same requirements as a standard planning system: targets in con-
 85 figuration space and an explicit collision checker to connect the path. Our work operates based on
 86 task-space targets and perceptual observations from a depth sensor without explicit state estimation.

87 Novel architectures have been proposed, such as differentiable planning modules in Value Iteration
 88 Networks [19], transformers by Chaplot et al. [30] and goal-conditioned RL policies [31]. These
 89 methods are challenging to generalize to unknown environments or have only been shown in simple
 90 2D [18] or 3D settings [20]. In contrast, we illustrate our approach in the challenging domain of
 91 controlling a 7 degrees of freedom (DOF) manipulator in unknown, dynamic environments.

92 3 Learning from Motion Planning

93 3.1 Problem Formulation

94 $M\pi$ Nets expect two inputs, a robot configuration q_t and a segmented, calibrated point cloud z_t .
 95 Before passing q_t through the network, we normalize each element to be within $[-1, 1]$ according
 96 to the limits for the corresponding joint. We call this $q_t^{\|\cdot\|}$. The point cloud is always assumed to be
 97 calibrated in the robot’s base frame, and it encodes three segmentation classes: the robot’s current
 98 geometry, the scene geometry, and the target pose. Targets are inserted into the point cloud via
 99 points sampled from the mesh of a floating end effector placed at the target pose.

100 The network produces a displacement within normalized configuration space $\hat{q}_t^{\|\cdot\|}$. To get the next
 101 predicted state \hat{q}_{t+1} , we take $q_t^{\|\cdot\|} + \hat{q}_t^{\|\cdot\|}$, clamp between $[-1, 1]$, and unnormalize. During training,
 102 we use \hat{q}_{t+1} to compute the loss, and when executing, we use \hat{q}_{t+1} as the next position target for the
 103 robot’s low-level controller.

104 3.2 Model Architecture

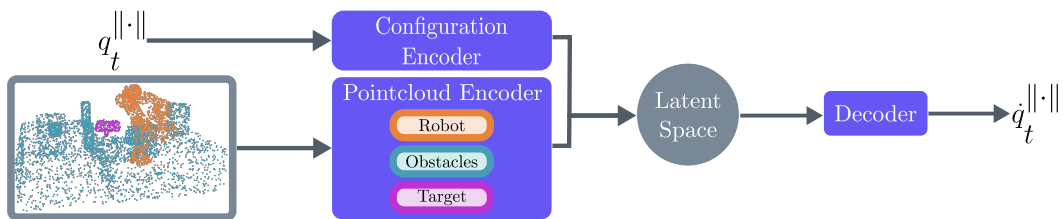


Figure 2: $M\pi$ Nets encodes state as a normalized robot configuration and segmented point cloud with three classes for the robot, the obstacles, and the target. The policy outputs a displacement in normalized joint space, which can then be applied to the input before unnormalizing to get q_{t+1} .

105 The network consists of two separate encoders, one for the point cloud and one for the robot’s
 106 current configuration, as well as a decoder, totaling 19M parameters. Our neural policy architecture
 107 is visualized in Fig. 2. We use PointNet++ [32] for our point cloud encoder. PointNet++ learns
 108 a hierarchical point cloud representation and can encode a point cloud’s 3D geometry, even with
 109 high variation in sampling density. PointNet++ architectures have been shown to be effective for
 110 a variety of point cloud processing tasks, such as segmentation [32], collision checking [17], and
 111 robotic grasping [33, 34]. The robot configuration encoder and the displacement decoder are both
 112 fully connected multilayer perceptrons. We discuss the architecture in detail in the Appendix.

113 3.3 Loss Function

114 The network is trained with a compound loss function with two constituent parts: a behavior cloning
115 loss to enforce accurate predictions and a collision loss to safeguard against catastrophic behavior.

116 **Geometric Loss for Behavior Cloning** To encourage alignment between the prediction and the
117 expert, we compute a geometric loss across many (1024) task spaces on the surface of the robot.

$$L_{\text{BC}}(\hat{\Delta}q_t) = \sum_i \|\hat{x}_{t+1}^i - x_{t+1}^i\|_2 + \|\hat{x}_{t+1}^i - x_{t+1}^i\|_1, \text{ where } \begin{matrix} \hat{x}_{t+1}^i = \phi^i(q_t + \hat{\Delta}q_t) \\ x_{t+1}^i = \phi^i(q_{t+1}) \end{matrix} \quad (1)$$

118 $\phi^i(\cdot)$ represents a forward kinematics mapping from the joint angles of the robot to point i defined
119 on the robot’s surface. The loss is computed as the sum of the $L1$ and $L2$ distances between cor-
120 responding points on the expert and the prediction after applying the predicted displacement. By
121 using both $L1$ and $L2$, we are able to penalize both large and small deviations.

122 We use a geometric, task-space loss because our goal is to ensure task-space consistency of our
123 policy. Configuration space loss appears in prior work [11], but does capture the accumulated error
124 of the kinematic chain as effectively (see Appendix).

125 **Collision Loss** In order to avoid collisions—a catastrophic failure—we apply an additional hinge-
126 loss inspired by motion optimization [35].

$$L_{\text{collision}} = \sum_i \sum_j \|h_j(\hat{x}_{t+1}^i)\|_2, \text{ where } h_j(\hat{x}_{t+1}^i) = \begin{cases} -D_j(\hat{x}_{t+1}^i), & \text{if } D_j(\hat{x}_{t+1}^i) \leq 0 \\ 0, & \text{if } D_j(\hat{x}_{t+1}^i) > 0 \end{cases} \quad (2)$$

127 The synthetic environments are fully-observable during training, giving us access to the signed-
128 distance functions (SDF), $\{D_j(\cdot)\}_j$, of the obstacles in each scene. For a given closed surface, its
129 SDF maps a point in Euclidean space to the minimum distance from the point to the surface. If the
130 point is inside the surface, the function returns negative.

131 3.4 Training Implementation Details

132 $M\pi$ Nets is only trained for single-step prediction, but we use it for closed-loop rollouts. The com-
133 pounded noise in subsequent inputs equates covariate shift [36, 37]. To promote robustness, we
134 augment our training data with random perturbations in two ways. We apply Gaussian noise to the
135 joint angles of each input configuration, which in turn affects the corresponding points in the point
136 cloud, passed as input to the network [38, 39]. We also sample novel scene point clouds on-the-fly
137 during training. We train our network with a single set of weights across our entire dataset.

138 4 Procedural Data Generation

139 4.1 Large-scale Motion Planning Problems

140 Each planning problem is defined by three components: the scene geometry, the start configuration,
141 and the goal pose. Our dataset consists of randomly generated problems across all three components,
142 totaling 3.27 million problems in over 575,000 environments. We have three classes of problems of
143 increasing difficulty: a cluttered tabletop with randomly placed objects, cubbies and dressers. Rep-
144 resentative examples of these environments are shown in Fig. 1. Once we build these environments,
145 we generate a set of potential end-effector targets and corresponding inverse kinematics solutions.
146 We then randomly choose pairs of these configurations and verify if a plan exists between them
147 using our expert pipeline, as detailed further in Sec. 4.2 and in the Appendix.

148 4.2 Expert Pipeline

149 Our expert pipeline is designed to produce high quality demonstrations we want to mimic, *i.e.* tra-
150 jectories with smooth, consistent motion and short path lengths. Here, *consistency* describes how

151 the motion would change in response to a perturbation in any dimension of the planning problem.
152 We considered two candidates for the expert - the *Global Planner* which is a typical state-of-the-art
153 configuration space planning pipeline [5] and a *Hybrid Planner* that we engineered specifically to
154 generate consistent motion in task space. For both planners, we reject any trajectories that produce
155 collisions, exceed the joint limits, exhibit erratic behavior (*i.e.* high jerk), or that have divergent
156 motion (*i.e.* final task space pose is more than 5cm from the target).

157 **Global Planner** consists of off-the-shelf components of a standard motion planning pipeline—inverse
158 kinematics [40], configuration-space AIT* [5], and spline-based, collision-aware trajectory smooth-
159 ing [41]. The *Global Planner* is theoretically complete and optimal, but in practice, we observed
160 several common failures types. The planner can time out without finding the optimal path or worse,
161 a solution at all. We utilize discrete collision checking for speed, which can cause the smoother to
162 produce collisions when very close to obstacles. Despite these types of failures, we generated 6.54
163 million trajectories for training across 773,068 environments.

164 **Hybrid Planner** is designed to produce consistent motion in task space. It consists of task-space
165 AIT* [5] and Geometric Fabrics [9]. AIT* ensures a minimal end-effector path and Geometric
166 Fabrics produce geometrically consistent motion. Finally, we fit a spline to the path and resample
167 to create a consistent configuration space velocity profile. As we discuss in Sec. 5.1, Geometric
168 Fabrics often fail to converge to a target, so we redefine the planning problem to have the same
169 target as the final position of the trajectory produced by the expert. Inspired by [42], we call this
170 technique *Hindsight Goal Revision (HGR)* and demonstrate its importance in Sec. 5.4. Using the
171 *Hybrid Planner*, we generated 3.27 million trajectories across 576,532 environments.

172 5 Experimental Evaluation

173 We evaluate our method and others using three test datasets: a set of problems solvable by the *Global*
174 *Planner*, problems solvable by the *Hybrid Planner*, and problems solvable by both. Each dataset has
175 1,800 problems, with 600 in each of the three types of environments. Each problem has a unique,
176 randomly generated environment, as well as a unique target and starting configuration.

177 **Quantitative Metrics:** To understand the performance of a policy, we roll it out until it matches
178 one of two termination conditions: 1) the Euclidean distance to the target is within 1cm or 2) the
179 trajectory has been executed for 20s (based on consultations with the authors of [9] and [10]). We
180 consider the following metrics (see Appendix for details):

- 181 • *Success Rate* - A trajectory is considered a success if its final position and orientation target
182 errors are below 1 cm and 15° respectively and there are no physical violations.
- 183 • *Time* - We measure the wall time for each *successful* trajectory. We also measure *Cold Start*
184 (*CS*) *Time*, the average time to react to a new planning problem.
- 185 • *Rollout Target Error* - The L2 position and orientation error (taken from [43]) between the
186 target and final end-effector pose in the trajectory.
- 187 • *Collision Rate* - The rate of fatal collisions, both self and scene collisions
- 188 • *Smoothness* - We use Spectral Arc Length (SPARC) [44] and consider a path to be smooth if
189 its SPARC values in joint and end-effector space are below -1.6 .

190 5.1 Comparison to Methods With Complete State

191 Most methods to generate motion in the literature assume access to complete state information in
192 order to perform collision checks. In each of the following experiments, we provide each baseline
193 method with an oracle collision checker. When running M π Nets, we use a point cloud sampled
194 uniformly from the surface of the entire scene. Results are shown in Table 1.

195 **Global Configuration Space Planner** The *Global Planner* is unmatched in its ability to reach a
196 target, but this comes at the cost of average computation time (16.46s) compared to M π Nets (0.33s).

	Soln. Time (s)	CS Time (s)	Success Rate (%)			
			Global	Hybrid	Both	Smooth (%)
Global Planner [5]	16.46 ± 0.90	16.46 ± 0.90	100	78.44	100	51.00
Hybrid Planner	7.37 ± 2.23	7.37 ± 2.23	50.22	100	100	99.26
G. Fabrics [9]	0.15 ± 0.09	2.4e-4 ± 3e-5	38.44	59.33	60.06	85.39
STORM [10]	4.03 ± 1.89	13.4e-3 ± 2.2e-3	50.22	74.50	76.00	62.26
MPNets [11]						
<i>Hybrid Expert</i>	4.95 ± 23.51	4.95 ± 23.51	41.33	65.28	67.67	99.97
<i>Random</i>	0.31 ± 3.55	0.31 ± 3.55	32.89	55.33	58.17	99.96
M π Nets (Ours)						
<i>Global Expert</i>	0.33 ± 0.08	6.8e-3 ± 7e-5	75.06	80.39	82.78	89.67
<i>Hybrid Expert</i>	0.33 ± 0.08	6.8e-3 ± 7e-5	75.78	95.33	95.06	93.81

Table 1: Algorithm performance on problems sets solvable by planner types. All prior methods use state-information and a oracle collision checker while M π Nets only needs a point cloud

	Training Set	Evaluation Set	
		MPNets-Style	Hybrid Expert Solvable (Ours)
MPNets [11]	MPNets-Style	78.70	49.89
M π Nets (Ours)	MPNets-Style	33.70	5.50
MPNets [11]	Hybrid Expert	88.90	65.28
M π Nets (Ours)	Hybrid Expert	89.50	95.33

Table 2: Success rates (%) of our method compared to Motion Planning Networks (MPNets) [11] trained and evaluated on different datasets

197 With a global planner, there is no option to partially solve a problem, meaning the Cold Start Time is
198 equal to the planning time. In a real system, optimizers [6, 7, 8] could be used to quickly replan once
199 an initial plan has been discovered. As discussed in Sec. 4.2, the *Global Planner* is theoretically
200 complete, but fails in practice on some of the *Hybrid Planner*-solvable problems due to system
201 timeouts and discrete collision checking during smoothing.

202 **Hybrid End-Effector Space Planner** Our *Hybrid Planner* struggles with a large proportion of
203 problems solvable by the *Global Planner*. Yet, its solutions are both faster and smoother than the
204 *Global Planner*. Surprisingly, M π Nets trained with data from the expert *outperformed* the expert
205 on the *Global Planner*-solvable test set. We attribute this to two features: 1) we use strict rejection
206 sampling to reduce erratic and divergent behavior in our expert dataset and train only on the filtered
207 data and 2) our use of Hindsight Goal Revision to turn an imperfect expert into a perfect one.

208 **Neural Motion Planning** Motion Planning Networks (MPNets) [11] proposed a similar method
209 for neural motion planning, but there are a few key differences in both problem setup and system
210 architecture. MPNets requires a ground-truth collision checker to connect sparse waypoints, plans
211 in configuration space, and is not reactive to changing conditions. In the architecture, MPNets uses
212 a trained neural sampler within a hierarchical bidirectional planner. The neural sampler is a fully-
213 connected network that accepts the start, goal, and a flattened representation of the obstacle points
214 as inputs and outputs a sample. MPNets guarantees completeness by using a traditional planner as a
215 fallback if the neural sampler fails to produce a valid plan.

216 In addition to our data, we generated a set of tabletop problems, which we call *MPNets-Style*, akin
217 to the Baxter experiments in [11], in order to fairly compare the two methods. The results of this
218 experiment can be seen in Table 2. M π Nets requires a large dataset for compelling performance,
219 while MPNets’ utilization of a traditional planning system is much more effective with a small
220 dataset. However, the MPNets architecture does not scale to more complex scenes, even with more
221 data, as we show in Fig. 3. When trained and evaluated on the Hybrid Planner-solvable dataset,

	% Env. Coll.	% Self Coll.	% Jnt Viol.	% Within			
				1cm	5cm	15°	30°
G. Fabrics [9]	8.61	0.11	0.44	69.89	75.17	83.44	85.11
STORM [10]	0.93	0.11	0.25	79.81	83.54	81.57	85.41
M π Nets (Ours)							
<i>Hybrid Expert</i>	0.94	0.00	0.00	98.94	99.72	98.22	99.00
<i>Global Expert</i>	13.78	0.06	0.00	98.67	99.89	97.56	99.11

Table 3: Failure Modes on problems solvable by both the global and hybrid planners

222 MPNets succeeds in 65.28% of the test set, whereas M π Nets succeeds in 95.33%, thus decreasing
 223 the failure rate by 7X. Furthermore, as we show in Table 1, using the MPNets neural sampler trained
 224 with the *Hybrid Planner* performs similarly to a uniform random sampler when both are embedded
 225 within the bidirectional MPNets planner.

226 **Local Task Space Controllers** Unlike planners, which succeed or fail in binary fashion, local
 227 policies will produce individual actions that, when rolled out, may fail for various reasons. We
 228 break down the various failure modes across the set of problems solvable by both experts in Table 3.

229 STORM [10] and Geometric Fabrics [9] make local decisions that can lead them to diverge from
 230 the target in complex scenarios, such as cluttered environments or those with pockets. For example,
 231 both STORM and Geometric Fabrics struggle to retract from a drawer and then reach into another
 232 drawer in a single motion without intermediate waypoints. Since our network compresses long-term
 233 planning information into a policy, it is much less prone to local minima and can reach between
 234 drawers with no additional scene information.

235 On problems solvable by the *Hybrid Planner*, M π Nets ties or outperforms these other methods
 236 across nearly all metrics. On the set of problems solvable by the *Global Planner*, M π Nets target
 237 convergence rate is consistently higher, while its collision rate (11%) is worse than either STORM
 238 (2.06%) or Geometric Fabrics (7.94%). Deteriorating performance on out-of-distribution problems
 239 is a typical downside of a supervised learning approach such as M π Nets. However, this could be
 240 improved with a more robust expert, *e.g.* one with the consistency of our *Hybrid Planner* but the
 241 success rate of the *Global Planner*, with finetuning, or with DAGger [45].

242 5.2 Importance of the Expert Pipeline

243 We observed that the choice of the expert pipeline affects the performance of M π Nets. We trained
 244 two policies: M π Nets-G with 6.54M demonstrations from the *Global Planner* and M π Nets-H with
 245 3.27M demonstrations from the *Hybrid Planner*. When evaluated on a test set of problems solvable
 246 by the *Global Planner*, M π Nets-G shows far better target convergence (97.94% vs. 87.72%) com-
 247 pared to M π Nets-H but worse obstacle avoidance (21.94% collision rate vs. 11%). Nonetheless,
 248 M π Nets-H is significantly better across all metrics when evaluated on problems solved by both ex-
 249 perts as shown in Table 3. We hypothesize that an expert combining the properties of these two—the
 250 consistency of the *Hybrid Planner* and the generality of the *Global Planner*, would further improve
 251 M π Nets’s performance. We refer to M π Nets-H as M π Nets throughout the rest of the paper.

252 5.3 Comparison to Methods With Partial Observations

253 In addition to demonstrating M π Nets’ performance on a real robot system, we also compared
 254 M π Nets to the *Global Planner* (AIT* [5]) in a single-view depth camera setting in simulation.
 255 We evaluated on the test set of problems solvable by both the *Global* and *Hybrid Planners*. M π Nets
 256 only has a minor drop in success rate when using a partial point cloud vs. a full point cloud— from
 257 95.06% to 93.22% though the collision rate increases from 0.94% to 3.06% due to occlusions. To
 258 evaluate AIT*, we omitted the smoothing step, which can cause erroneous collisions as discussed in
 259 Sec. 4.2. We used a voxel-based reconstruction akin to the standard perception pipeline packaged

260 with MoveIt [46]. On the same test set, AIT* produces plans with collisions on 16.41% of problems.
 261 In this setting, $M\pi$ Nets’s collision rate is over 5X smaller than that of the *Global Planner*.

262 5.4 Ablations

263 We perform several ablations to justify our design
 264 decisions. All ablations were trained using the *Hy-*
 265 *brid Planner* dataset and evaluated on the *Hybrid*
 266 *Planner*-solvable test set. More ablations and details
 267 can be found in the Appendix.

268 **$M\pi$ Nets Performance Scales with More Data** As
 269 shown in Fig. 3, the performance of $M\pi$ Nets contin-
 270 ues to improve with more data, although it saturates
 271 at 1.1M. Meanwhile, MPNets [11] has constant per-
 272 formance, demonstrating that our architecture is bet-
 273 ter able to scale with the data.

274 **Robot Point Representation Improves Perform-**
 275 **ance** Instead of representing the robot by its configura-
 276 tion vector, we insert the robot point cloud
 277 at the specific configuration. Without this representa-
 278 tion, the success rate decreases from 95.33% to
 279 65.06%.

280 **Hindsight Goal Revision Improves Convergence** When trained without *HGR*, *i.e.* with the plan-
 281 ner’s original target given to the network, we see 58.11% success rate vs. 95.33% when trained with
 282 *HGR*. In particular, only 60.28% of trajectories get within 1cm of the target during evaluation.

283 **Noise Injection Improves Robustness** When we train $M\pi$ Nets without injecting noise into the
 284 input q_t , the policy performance decreases by 10.72%.

283 5.5 Real Robot Evaluation

284 We ran multiple qualitative experiments on a 7-DOF Franka Emika Panda robot with an extrinsically
 285 calibrated Intel Realsense L515 RGB-D camera mounted next to it. Depth measurements belonging
 286 to the robot are removed before inference with $M\pi$ Nets. Rollouts are between 2 and 80 time steps
 287 long depending on the control loop frequency. Results can be viewed at <https://mpinets.github.io>
 288 and the attached video. As can be seen, $M\pi$ Nets achieve *sim2real* transfer on noisy real-world point
 289 clouds in unknown and changing scenes.

290 6 Limitations

291 While $M\pi$ Nets can handle a large class of problems, they are ultimately limited by the quality of
 292 the expert supervisor. It will also struggle to generalize to out-of-distribution settings typical of
 293 any supervised learning approach. In future work we aim to improve $M\pi$ Nets with DAgger [45] or
 294 domain adaptation. In order to further enable safe operation in real robot systems, $M\pi$ Nets could
 295 also be combined with a ground-truth or learnt collision checker such as SceneCollisionNet [17].

296 7 Conclusion

297 $M\pi$ Nets is a class of end-to-end neural policy policies that learn to navigate to pose targets in task
 298 space while avoiding obstacles. $M\pi$ Nets show robust, reactive performance on a real robot system
 299 using data from a single, static depth camera. We train $M\pi$ Nets with what is, as far as we are aware,
 300 the largest existing dataset of end-to-end motion for a robotic manipulator. Our experiments show
 301 that when applied to appropriate problems, $M\pi$ Nets are significantly faster than a global motion
 302 planner and more capable than prior neural planners and manually designed local control policies.
 303 We will release our code and data upon publication.

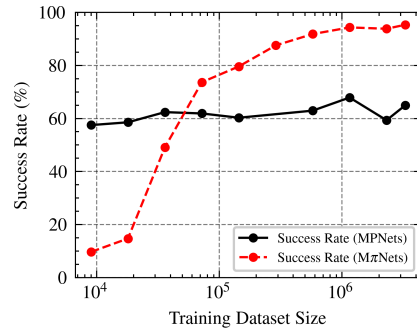


Figure 3: $M\pi$ Nets performance continues to increase with more training data, while MPNets performance stays relatively constant

References

- [1] S. M. LaValle. Rapidly-exploring random trees : a new tool for path planning. *The annual research report*, 1998.
- [2] S. Karaman and E. Frazzoli. Sampling-based algorithms for optimal motion planning. *The International Journal of Robotics Research*, 30:846 – 894, 2011.
- [3] J. D. Gammell, S. S. Srinivasa, and T. D. Barfoot. Batch informed trees (bit*): Sampling-based optimal planning via the heuristically guided search of implicit random geometric graphs. *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3067–3074, 2015.
- [4] M. P. Strub and J. D. Gammell. Advanced bit* (abit*): Sampling-based planning with advanced graph-search techniques. *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 130–136, 2020.
- [5] M. P. Strub and J. D. Gammell. Adaptively informed trees (ait*): Fast asymptotically optimal path planning through adaptive heuristics. *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3191–3198, 2020.
- [6] N. D. Ratliff, M. Zucker, J. A. Bagnell, and S. S. Srinivasa. Chomp: Gradient optimization techniques for efficient motion planning. *2009 IEEE International Conference on Robotics and Automation*, pages 489–494, 2009.
- [7] J. Schulman, Y. Duan, J. Ho, A. X. Lee, I. Awwal, H. Bradlow, J. Pan, S. Patil, K. Goldberg, and P. Abbeel. Motion planning with sequential convex optimization and convex collision checking. *The International Journal of Robotics Research*, 33:1251 – 1270, 2014.
- [8] M. Mukadam, J. Dong, X. Yan, F. Dellaert, and B. Boots. Continuous-time Gaussian process motion planning via probabilistic inference. volume 37, pages 1319–1340, 2018.
- [9] K. V. Wyk, M. Xie, A. Li, M. A. Rana, B. Babich, B. N. Peele, Q. Wan, I. Akinola, B. Sundaralingam, D. Fox, B. Boots, and N. D. Ratliff. Geometric fabrics: Generalizing classical mechanics to capture the physics of behavior. *IEEE Robotics and Automation Letters*, 7:3202–3209, 2022.
- [10] M. Bhardwaj, B. Sundaralingam, A. Mousavian, N. D. Ratliff, D. Fox, F. Ramos, and B. Boots. STORM: An integrated framework for fast joint-space model-predictive control for reactive manipulation. 2021.
- [11] A. H. Qureshi, M. J. Bency, and M. C. Yip. Motion planning networks. *2019 International Conference on Robotics and Automation (ICRA)*, pages 2118–2124, 2019.
- [12] B. Ichter, J. Harrison, and M. Pavone. Learning sampling distributions for robot motion planning. *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7087–7094, 2018.
- [13] J. C. Kew, B. Ichter, M. Bandari, T.-W. E. Lee, and A. Faust. Neural collision clearance estimator for batched motion planning. *arXiv preprint arXiv:1910.05917*, 2019.
- [14] R. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. *2011 10th IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, 2011.
- [15] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [16] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.

- 348 [17] M. Danielczuk, A. Mousavian, C. Eppner, and D. Fox. Object rearrangement using learned
349 implicit collision functions. *2021 IEEE International Conference on Robotics and Automation*
350 (*ICRA*), pages 6010–6017, 2021.
- 351 [18] T. Jurgenson and A. Tamar. Harnessing Reinforcement Learning for Neural Motion Planning.
352 In *Robotics Science and Systems*, 2019.
- 353 [19] A. Tamar, Y. Wu, G. Thomas, S. Levine, and P. Abbeel. Value iteration networks. In *Neural*
354 *Information Processing Systems*, 2016.
- 355 [20] R. A. M. Strudel, R. G. Pinel, J. Carpentier, J.-P. Laumond, I. Laptev, and C. Schmid. Learning
356 obstacle representations for neural motion planning. In *CoRL*, 2020.
- 357 [21] P. E. Hart, N. J. Nilsson, and B. Raphael. A formal basis for the heuristic determination of
358 minimum cost paths. *IEEE Trans. Syst. Sci. Cybern.*, 4:100–107, 1968.
- 359 [22] M. Likhachev, G. J. Gordon, and S. Thrun. Ara*: Anytime a* with provable bounds on sub-
360 optimality. In *NIPS*, 2003.
- 361 [23] M. Likhachev, D. Ferguson, G. J. Gordon, A. Stentz, and S. Thrun. Anytime dynamic a*: An
362 anytime, replanning algorithm. In *ICAPS*, 2005.
- 363 [24] N. D. Ratliff, M. Toussaint, and S. Schaal. Understanding the geometry of workspace obstacles
364 in motion optimization. *2015 IEEE International Conference on Robotics and Automation*
365 (*ICRA*), pages 4202–4209, 2015.
- 366 [25] O. Khatib. Real-time obstacle avoidance for manipulators and mobile robots. In *Autonomous*
367 *robot vehicles*, pages 396–404. Springer, 1986.
- 368 [26] N. D. Ratliff, J. Issac, and D. Kappler. Riemannian motion policies. *ArXiv*, abs/1801.02854,
369 2018.
- 370 [27] R. Kumar, A. Mandalika, S. Choudhury, and S. S. Srinivasa. Lego: Leveraging experience in
371 roadmap generation for sampling-based planning. *2019 IEEE/RSJ International Conference*
372 *on Intelligent Robots and Systems (IROS)*, pages 1488–1495, 2019.
- 373 [28] C. Zhang, J. Huh, and D. D. Lee. Learning implicit sampling distributions for motion planning.
374 *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages
375 3654–3661, 2018.
- 376 [29] C. Chamzas, Z. K. Kingston, C. Quintero-Peña, A. Shrivastava, and L. E. Kavraki. Learning
377 sampling distributions using local 3d workspace decompositions for motion planning in high
378 dimensions. *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages
379 1283–1289, 2021.
- 380 [30] D. S. Chaplot, D. Pathak, and J. Malik. Differentiable spatial planning using transformers. In
381 *International Conference on Machine Learning*, 2021.
- 382 [31] B. Eysenbach, R. Salakhutdinov, and S. Levine. C-learning: Learning to achieve goals via
383 recursive. In *International Conference on Learning Representations*, 2021.
- 384 [32] C. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point
385 sets in a metric space. In *NIPS*, 2017.
- 386 [33] A. Mousavian, C. Eppner, and D. Fox. 6-dof graspnet: Variational grasp generation for object
387 manipulation. In *International Conference on Computer Vision (ICCV)*, 2019.
- 388 [34] A. Murali, A. Mousavian, C. Eppner, C. Paxton, and D. Fox. 6-dof grasping for target-driven
389 object manipulation in clutter. In *IEEE International Conference on Robotics and Automation*
390 (*ICRA*), 2020.

- 391 [35] A. Fishman, C. Paxton, W. Yang, D. Fox, B. Boots, and N. D. Ratliff. Collaborative interaction
392 models for optimized human-robot teamwork. *2020 IEEE/RSJ International Conference on*
393 *Intelligent Robots and Systems (IROS)*, pages 11221–11228, 2020.
- 394 [36] S. Ross and A. Bagnell. Efficient reductions for imitation learning. In *International Conference*
395 *on Artificial Intelligence and Statistics*, page 661–668, 2010.
- 396 [37] S. Ross, G. Gordon, and A. Bagnell. A reduction of imitation learning and structured prediction
397 to no-regret online learning. In *arXiv preprint arXiv:1011.0686*, 2010.
- 398 [38] L. Ke, J. Wang, T. Bhattacharjee, B. Boots, and S. S. Srinivasa. Grasping with chopsticks:
399 Combating covariate shift in model-free imitation learning for fine manipulation. *2021 IEEE*
400 *International Conference on Robotics and Automation (ICRA)*, pages 6185–6191, 2021.
- 401 [39] M. Laskey, J. Lee, R. Fox, A. Dragan, and K. Goldberg. Dart: Noise injection for robust
402 imitation learning. *Conference on Robot Learning*, pages 143–156, 2017.
- 403 [40] R. Diankov. *Automated Construction of Robotic Manipulation Programs*. PhD the-
404 sis, Carnegie Mellon University, Robotics Institute, August 2010. URL [http://www.](http://www.programmivision.com/rosen_diankov_thesis.pdf)
405 [programmivision.com/rosen_diankov_thesis.pdf](http://www.programmivision.com/rosen_diankov_thesis.pdf).
- 406 [41] K. K. Hauser and V. Ng-Thow-Hing. Fast smoothing of manipulator trajectories using op-
407 timal bounded-acceleration shortcuts. *2010 IEEE International Conference on Robotics and*
408 *Automation*, pages 2493–2498, 2010.
- 409 [42] M. Andrychowicz, D. Crow, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,
410 P. Abbeel, and W. Zaremba. Hindsight experience replay. In *NIPS*, 2017.
- 411 [43] P. Wunsch, S. Winkler, and G. Hirzinger. Real-time pose estimation of 3d objects from cam-
412 era images using neural networks. *Proceedings of International Conference on Robotics and*
413 *Automation*, 4:3232–3237 vol.4, 1997.
- 414 [44] S. Balasubramanian, A. Melendez-Calderon, A. Roby-Brami, and E. Burdet. On the analysis
415 of movement smoothness. *Journal of NeuroEngineering and Rehabilitation*, 12, 2015.
- 416 [45] S. Ross, G. J. Gordon, and J. A. Bagnell. A reduction of imitation learning and structured
417 prediction to no-regret online learning. In *AISTATS*, 2011.
- 418 [46] S. Chitta, I. Sucas, and S. Cousins. Moveit![ros topics]. *IEEE Robotics & Automation Maga-*
419 *zine*, 19(1):18–19, 2012.