
Tactile MNIST: Benchmarking Active Tactile Perception

Tim Schneider^{1,2}, Guillaume Duret², Cristiana de Farias¹,
Roberto Calandra³, Liming Chen², and Jan Peters⁴

¹Department of Computer Science, TU Darmstadt, Germany.

²LIRIS, CNRS UMR5205, Ecole Centrale de Lyon, France.

³LASR Lab & CeTI, TU Dresden, Germany.

⁴DFKI, Hessian.AI, and Centre for Cognitive Science, TU Darmstadt, Germany.

Abstract

1 Tactile perception has the potential to significantly enhance dexterous robotic
2 manipulation by providing rich local information that can complement or substitute
3 for other sensory modalities such as vision. However, because tactile sensing is
4 inherently local, it is not well-suited for tasks that require broad spatial awareness
5 or global scene understanding on its own. A human-inspired strategy to address
6 this issue is to consider active perception techniques instead. That is, to actively
7 guide sensors toward regions with more informative or significant features and
8 integrate such information over time in order to understand a scene or complete
9 a task. Both active perception and different methods for tactile sensing have
10 received significant attention recently. Yet, despite advancements, both fields lack
11 standardized benchmarks. To bridge this gap, we introduce the *Tactile MNIST*
12 *Benchmark Suite*, an open-source, Gymnasium-compatible benchmark specifically
13 designed for active tactile perception tasks, including localization, classification,
14 and volume estimation. Our benchmark suite offers diverse simulation scenarios,
15 from simple toy environments all the way to complex tactile perception tasks using
16 vision-based tactile sensors. Furthermore, we also offer a comprehensive dataset
17 comprising 13,500 synthetic 3D MNIST digit models and 153,600 real-world
18 tactile samples collected from 600 3D printed digits. Using this dataset, we train a
19 CycleGAN for realistic tactile simulation rendering. By providing standardized
20 protocols and reproducible evaluation frameworks, our benchmark suite facilitates
21 systematic progress in the fields of tactile sensing and active perception.

22 **Project page:** <https://sites.google.com/robot-learning.de/tactile-mnist>

23 1 Introduction

24 Tactile perception is fundamental for enabling agents to interact effectively with their environments.
25 Studies of humans with impaired touch reveal that they face significant challenges in grasping and
26 performing routine manipulation tasks due to insufficient feedback about contact states between
27 fingers and objects [4]. Moreover, touch often complements—or even substitutes—other sensory
28 modalities such as vision: we feel the shape of a hard-to-see object on a high shelf, count cookies in
29 a jar without looking, or locate a key inside a bag purely by touch. Unlike vision, which typically
30 offers a broad field of view, touch provides highly localized yet information-rich feedback confined
31 to the point of contact [5, 6]. It is this inherently interactive nature that enables agents (such as
32 robots) to explore visually occluded areas, classify textures, infer material properties like stiffness

Corresponding Author: tim@robot-learning.de

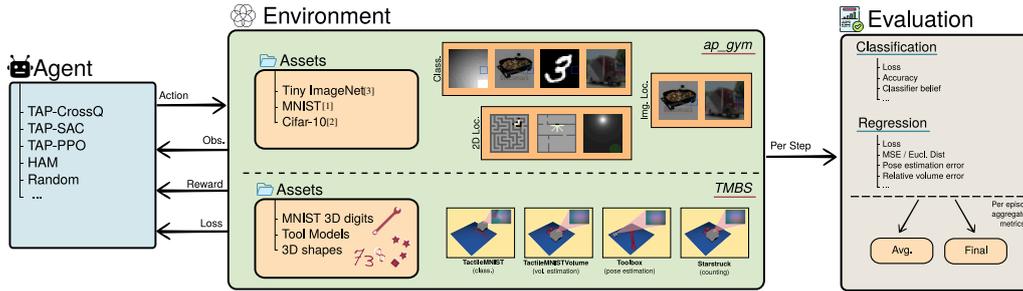


Figure 1: Overview of the *Active Perception Gym* (*ap_gym*), the *Tactile MNIST Benchmark Suite* (TMBS), and their associated assets. In the center we depict each environment from *ap_gym* and TMBS, along with the included asset sets (both custom and external). On the left, an agent (e.g., an active perception algorithm) interacts with the environment by receiving observations, rewards, and losses, and returning actions. On the right, we show task-specific evaluation metrics, available at each step, with support for both per-step outputs and aggregated performance scores.

33 and friction, and detect fine local features of objects [7, 8, 9, 10, 11, 12, 13, 14]. However, without
 34 standardized benchmarks and widely adopted datasets, it becomes difficult to rigorously evaluate new
 35 algorithms, reproduce results, or compare different approaches in a fair way. Yet, despite its growing
 36 importance, the field of tactile sensing still lacks such structured benchmarks and community-wide
 37 datasets tailored to touch-related tasks

38 In contrast to tactile sensing, computer vision has significantly benefited from such standardized
 39 datasets and clearly defined evaluation protocols. MNIST [1] established an early baseline for digit
 40 recognition that is still relevant today. Datasets such as ImageNet [15] and COCO [16] expanded
 41 both scale and complexity, driving major advances in object classification, detection, and scene
 42 understanding. Additionally, domain-specific benchmarks such as KITTI [17] and Omni3D [18] have
 43 enabled focused progress on challenges like fine-grained categorization and robustness in real-world
 44 applications. By comparison, the few tactile perception benchmarks in the literature still rely on
 45 custom hardware or narrowly scoped datasets, which limits their generalizability and slows broader
 46 adoption [14, 19].

47 In this work, we focus explicitly on benchmarking *active tactile perception* tasks where agents
 48 deliberately interact with their environment to gather task-relevant information. Active perception
 49 involves strategic decisions about where and when to sense, efficiently choosing actions that maximize
 50 information gain [10, 9, 20, 8]. Here, as each contact takes time, resources, and may cause wear
 51 on sensors or environments, efficiency becomes a central concern: agents must extract as much
 52 information as possible using as few interactions as necessary. A well-designed benchmark for active
 53 tactile perception can help answer key questions such as: How should an agent select contact points to
 54 maximize information gain? What policies enable accurate inference with minimal touch? How does
 55 uncertainty (from sensor noise, ambiguous contact, or environmental variability) affect the trade-off
 56 between exploration and confidence?

57 We introduce *Active Perception Gym* (*ap_gym*)¹, a framework compatible with Gymnasium [21],
 58 designed to benchmark active perception algorithms. *ap_gym* includes nine toy scenarios where an
 59 agent must learn to efficiently extract information to solve perception tasks. Building on *ap_gym*, we
 60 present the *Tactile MNIST Benchmark Suite* (TMBS)², which extends the framework to simulated
 61 active tactile perception problems. In TMBS, agents control a simulated GelSight Mini sensor [22]
 62 without access to visual inputs. Tactile perception tasks include MNIST-style classification of 3D
 63 digit models, pose estimation of tools on platforms, and object counting. Across all tasks, the agent
 64 must actively determine *what* to sense and strategically select *where and when* to explore through
 65 touch. Thus, solving these tasks requires solving a dual problem: making accurate predictions from
 66 past observations while optimizing an exploration policy to maximize information gain. An overall
 67 visualization of our framework and tasks in *ap_gym* and TMBS is shown in Fig. 1.

¹ <https://github.com/TimSchneider42/active-perception-gym>

² <https://github.com/TimSchneider42/tactile-mnist>

Table 1: Comparison of Active Tactile Perception Benchmarks. Tactile modalities are **bolded**.

Method	Dataset Available	Active Benchmark	Sensor Modality
Active Vision Grasp [24]	✗	✓	Vis
R3ED [25]	✓	✓	Vis
Robotic Vision Challenge [26]	✗	✓	Vis
NBV_Bench [27]	✓	✓	Vis
AVD [28, 29]	✓	✓	Vis
Active Object Search [30]	✓	✓	Vis
ActiView [31]	✓	✓	Vis+Text
TIP Bench. [32]	✗	✗	Tac
FoTa [33]	✓	✗	Tac
YCB-Slide [34]	✓	✗	Tac
TacBench [14]	✓	✗	Tac
ActiveCloth [35]	✓	✗	Tac
Touch and Go [36]	✓	✗	Tac+Vis
FeelSight [19]	✓	✗	Tac+Vis
ViTac [37]	✗	✗	Tac+Vis
SSVTP [13]	✓	✗	Tac+Vis
GelFabric [38]	✓	✗	Tac+Vis
VisGel [39]	✓	✗	Tac+Vis
PHYSICLEAR [40]	✓	✗	Tac+Text
TVL [41]	✓	✗	Tac+Vis+Text
Touch100k [42]	✓	✗	Tac+Vis+Text
ObjectFolder [43, 44, 45]	✓	✗	Tac+Vis+Audio
Ours	✓	✓	Tac

68 With this benchmark, our goal is to provide a reproducible and accessible evaluation framework for
 69 the tactile perception community. To that end, TMBS is a fully *simulated* environment that is easy to
 70 set up and enables rapid, consistent comparisons, without the need to replicate a complex real-world
 71 setup. However, we acknowledge the inherent challenges and noise present in real-world tactile data
 72 that are often absent in simulation. To help bridge this sim-to-real gap, we complement our simulated
 73 environment with a large-scale dataset of 13,580 high-fidelity 3D object models. From this collection,
 74 we 3D-printed 600 objects and constructed a curated real-world dataset comprising 153,600 tactile
 75 contacts, each annotated with detailed temporal and spatial metadata. We further leverage this dataset
 76 to train a CycleGAN [23], enabling the rendering of realistic tactile signals within the simulation.

77 In summary, our main contributions are:

- 78 • To the best of our knowledge, we introduce the first benchmark suite specifically for active
 79 *tactile* perception. It offers a range of tasks from simple toy problems to challenging,
 80 high-dimensional scenarios.
- 81 • We further introduce `ap_gym`, an extensible framework for generic active perception al-
 82 gorithms. `ap_gym` is Gymnasium compatible and is, thus, easy to integrate with existing
 83 reinforcement-learning pipelines, enabling fair evaluation.
- 84 • We provide an open dataset of 13,580 high-resolution 3D models of handwritten digits,
 85 designed for both simulated tactile-image generation and physical 3D printing.
- 86 • From our library of 3D models, we 3D-printed 600 objects and captured 153,600 tactile
 87 contacts using a GelSight Mini sensor, each annotated with spatial location and class labels.

88 2 Related Work

89 **(Active) Tactile Perception:** Tactile sensing enables robots to infer object geometry, texture and
 90 material properties through physical contact, complementing or, in some cases, substituting vision.
 91 Vision-based tactile sensors such as GelSight [22] and DIGIT [46] have become widely available,
 92 producing high-resolution “tactile images” that capture local surface features and force distributions.
 93 These rich signals have been exploited for shape reconstruction and material recognition [7, 32, 41,
 94 40], as well as advanced dexterous manipulation [47, 48, 19].

95 Much of the work in tactile sensing, however, focuses on passive touch: the robot either registers a
 96 single contact or follows a predefined exploration policy. Inspired by the successes of active vision
 97 (as well as by early research on active touch in robotics [49, 50]), recent efforts have revisited the idea
 98 of tactile exploration as an active, decision-driven process. In this framing, the robot dynamically
 99 selects where and how to touch next, rather than relying on a fixed sequence of actions. Gaussian

100 process and Bayesian optimization have been employed to drive this active exploration, yielding
 101 significant improvements in tasks such as shape reconstruction, texture classification and grasp
 102 planning [10, 7, 51]. Reinforcement-learning based approaches [9] tackle active exploration in high-
 103 dimensional tactile state spaces. Furthermore, [8] introduced HAM a selective-attention mechanism
 104 to optimize scene exploration, and in [52] task-agnostic strategies generalize active touch across
 105 different objectives. Additionally, [35] demonstrated how Kinect-based vision can guide active touch
 106 for material classification, and [13] proposed a self-supervised visuo-tactile pretraining scheme that
 107 benefits both passive and active perception tasks.

108 **Benchmarking Methods for Tactile Sensing & Active Perception:** Over the past decade, the
 109 active vision community has produced a number of datasets and challenges to evaluate a range of
 110 tasks. Early work such as the Active Vision Dataset (AVD) provided large-scale Kinect captures for
 111 navigation and class-incremental learning tasks [28, 29]. Subsequent efforts explored next-best-view
 112 planning for classification (NBV_Bench [27]), heuristic and data-driven view selection for grasp
 113 synthesis on YCB objects [24], and embodied 3D exploration in real indoor scenes (R3ED [25]).
 114 Simulation-based approaches such as the Robotic Vision Scene Understanding Challenge [26] and
 115 Active Object Search [30] have further expanded evaluation protocols for semantic SLAM and object
 116 detection tasks. More recently, multi-modal active perception has been addressed by ActiView, which
 117 tests an agent’s ability to zoom and pan to answer vision-language queries [31].

118 In parallel, tactile sensing research has released a diverse set of characterization benchmarks (TIP
 119 Bench. [32]), texture and material recognition datasets (ActiveCloth [35], ViTac [37], SSVTP [13],
 120 GelFabric [38]), and cross-modal vision–touch benchmarks and datasets (VisGel [39], Touch and
 121 Go [36], PHYSICLEAR [40], FeelSight [19], TacBench [14]). Large-scale multimodal datasets
 122 such as Touch100k [42], TVL [41], and FoTa [33] now exceed tens of thousands of samples across
 123 vision, touch, and language modalities. Multisensory datasets like ObjectFolder [43, 44, 53] further
 124 integrate tactile, visual, and audio data. In a similar manner, efforts to standardize vision-based tactile
 125 simulation have led to platforms like TACTO [54] and Taxim [55], which enable high-resolution
 126 visuo-tactile data generation. Despite these efforts, only the datasets provided by ActiveCloth [35]
 127 and SSVTP [13] support downstream active tactile perception tasks. However, these works do not
 128 provide standardized evaluation protocols or benchmarks to systematically evaluate active perception
 129 approaches. A comparison of existing datasets and benchmarks is summarized in Table 1. To the best
 130 of our knowledge, our proposed Tactile MNIST Benchmark Suite is the first to introduce a dedicated
 131 and reproducible benchmark for *active tactile perception*, where tactile exploration is an integral
 132 component of the perceptual process.

Table 2: Overview of the environments, including task types, descriptions, and assets.

Suite	Environment	Task Type	Description	Assets
ap_gym	TinyImageNet	Classification	Classify natural images into 200 categories by moving a limited field-of-view glimpse.	Tiny ImageNet [3]
	CIFAR10	Classification	Classify natural images into 10 categories by moving a limited field-of-view glimpse.	CIFAR-10 [2]
	CircleSquare	Classification	Determine whether a given image contains a circle or a square using limited agent visibility.	Geometric shapes
	MNIST	Classification	Digit recognition task using standard MNIST digits.	MNIST [1]
	LightDark	Regression (2D localization)	Position estimation from brightness-dependent noisy observations, requiring movement to light.	None
	LIDARLocRooms	Regression (2D localization)	Navigate procedurally generated maps with ambiguous LIDAR readings to localize.	None
	LIDARLocMaze	Regression (2D localization)	Navigate procedurally generated mazes with ambiguous LIDAR readings to localize.	None
	TinyImageNetLoc	Regression (2D patch localization)	Localize a glimpse within a natural image by moving a limited field-of-view.	Tiny ImageNet [3]
CIFAR10Loc	Regression (2D patch localization)	Localize a glimpse within a natural image by moving a limited field-of-view.	CIFAR-10 [2]	
TMBS	TactileMNIST	Classification	Touch-based digit classification using a vision-based tactile sensor.	MNIST 3D digits
	TactileMNISTVolume	Regression (volume estimation)	Estimate volume of digits using a vision-based tactile sensor.	MNIST 3D digits
	Toolbox	Regression (object pose estimation)	Estimate 6D pose of tools (e.g., a wrench) using a vision-based tactile sensor.	3D tools
	Starstruck	Classification (counting)	Count stars among other objects using a vision-based tactile sensor.	3D shapes

133 3 Framework: Benchmarking Active Perception

134 In this section, we introduce both the *Active Perception Gym* (ap_gym), a framework for benchmarking
 135 active perception algorithms, and the *Tactile MNIST Benchmark Suite* (TMBS), which introduces
 136 environments for four active tactile perception tasks. Fig. 1 depicts an overview of our framework.

137 3.1 Active Perception

138 In active perception tasks, an agent’s main objective is to gather information and make predictions
 139 about a desired property of the environment, e.g., the class label or pose of an object. Examples of
 140 such properties could be the location of an object in case of a search task or the class of an object
 141 for the agent in case of a classification task. To gather information, the agent must interact with the
 142 environment, e.g., by moving a sensor around a platform in case of TMBS.

143 We model active perception as an episodic process, where the agent can take a number of time-
 144 discrete sequential actions until the episode terminates and is reset. At every step, the agent obtains
 145 an observation (e.g., a tactile glimpse) from the environment that reveals some information but never
 146 the full state at once. Formally, that makes active perception problems a special case of Partially
 147 Observable Markov Decision Processes (POMDPs).

148 POMDPs are defined by the tuple $(S, A, T, R, \Omega, O, \gamma)$, consisting hidden states S , actions A , a
 149 transition function $T : S \times A \times S \rightarrow [0, 1]$, a reward function $R : S \times A \rightarrow \mathbb{R}$, a set of observations
 150 Ω , an observation function $O : S \times A \times \Omega \rightarrow [0, 1]$, and a discount factor $\gamma \in [0, 1]$. The objective
 151 of the agent in a POMDP is to maximize the expected cumulative reward over time by selecting
 152 actions based on its belief about the underlying state. Since the agent does not have direct access
 153 to the true state, it maintains a belief distribution over states, updating it using observations and the
 154 observation function. The environment evolves according to the transition function, where taking an
 155 action leads to a probabilistic transition to a new state, which in turn generates an observation based
 156 on the observation function.

157 In case of active perception problems, we assume that the hidden state S , the action A , the reward
 158 function R , and the transition function T have specific structures. First, we assume that the target
 159 property the agent is tasked to predict is part of the hidden state. Hence, S is defined as $S = S_{\text{base}} \times Y^*$,
 160 where S_{base} is the set of base (hidden) states of the environment and Y^* is the set of prediction
 161 targets. E.g., Y^* could be the set of classes in a classification task or the set of possible locations
 162 in a localization task, while S_{base} contains all the other hidden state information. To allow the
 163 agent to make predictions, the action space A is defined as $A_{\text{base}} \times Y$, where A_{base} is the base
 164 action space and Y is the prediction space. The base action space A_{base} contains all the actions
 165 the agent can take to interact with the environment, while Y is the set of possible predictions
 166 the agent can make. Crucially, environments are defined in a way that the agent’s prediction
 167 never influences the hidden state of the environment. Thus, the transition function T is defined as
 168 $T(s, a, s') = T(s, (a_{\text{base}}, y), s') = T_{\text{base}}(s, a_{\text{base}}, s')$. An example of a base action could be a desired
 169 movement of the tactile sensor, while the prediction could be the logits of the agent’s current class
 170 prediction.

171 Finally, the reward function is defined as $R(s, a) = R((s_{\text{base}}, y^*), (a_{\text{base}}, y)) = R_{\text{base}}(s_{\text{base}}, a_{\text{base}}) -$
 172 $\ell(y^*, y)$, where R_{base} is the base reward function and ℓ is a differentiable loss function. An example

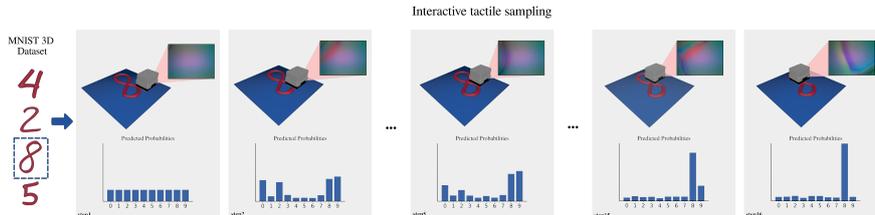


Figure 2: Illustration of the TactileMNIST classification task. In each episode of the TactileMNIST classification task, one random digit from the MNIST 3D dataset is selected and presented to the agent. It can then move the sensor around and touch the object 16 times before the episode terminates. Notably, it does not receive any visual input but has to rely solely on the readings from its tactile sensor. After every touch, it has to make a prediction about the class label, the digit’s numeric value, and its performance is measured by the average prediction accuracy throughout the episode.

173 for a base reward could be an action regularization term, while the loss function ℓ could be a
 174 cross-entropy loss in a classification task. Hence, the agent has to make a prediction in every
 175 step, encouraging it to gather information quickly to maximize its prediction reward early on. A
 176 visualization of this process on the TactileMNIST digit classification task is shown in Fig. 2.

177 3.2 Active Perception Gym

178 `ap_gym` models active perception tasks as episodic processes in a way that is fully compatible with
 179 Gymnasium [21]. Each task is defined as a Gymnasium environment, bundled with the differentiable
 180 loss function $\ell(y^*, y)$ and the prediction target y^* . Since the loss functions need to convey gradient
 181 information to the learning algorithm, we currently provide them either JAX [56] or PyTorch [57]
 182 functions, but more autograd frameworks might be supported in the future. During roll-outs, `ap_gym`
 183 automatically computes task-dependent metrics, such as accuracy for classification, or Euclidean
 184 distance for regression.

185 `ap_gym` provides a family of lightweight environments designed to isolate core exploration and
 186 decision-making behaviors in active perception. As summarized in Table 2, `ap_gym` includes 11
 187 environments spanning both classification and regression tasks. Four progressively harder image-
 188 based classification benchmarks (CircleSquare, MNIST, CIFAR-10, TinyImageNet) evaluate an
 189 agent’s ability to select informative glimpses from natural or synthetic visuals. Two image-localization
 190 tasks (TinyImageNetLoc, CIFAR10Loc) require the agent to infer the position of a limited-field-
 191 of-view patch within a larger image. Finally, five non-visual regression tasks — LightDark and
 192 four LIDAR-based 2D localization environments (Rooms and Maze) — challenge agents to reduce
 193 state uncertainty by navigating procedurally generated maps or lighting fields. Crucially, in the
 194 self-localization environments, the agent influences the property it is trying to infer — its position —
 195 through its actions. Hence, the prediction target changes over time in these environments, which is an
 196 explicitly supported aspect of `ap_gym` environments.

197 For the image-based classification and localization tasks, `ap_gym` relies on a mix of third-party assets:
 198 Tiny ImageNet [3], CIFAR-10 [2], and MNIST [1] datasets. The LIDARLoc (rooms and maze)
 199 environments generate map layouts procedurally and require no external data. Whenever applicable,
 200 `ap_gym` defines two versions of each environment, one for training with the training split of the
 201 respective dataset, and one for evaluation with the test split.

202 By abstracting away complex contact models and dynamics, `ap_gym` environments enable rapid
 203 prototyping of active perception strategies and provide a controlled baseline for more complex
 204 scenarios in the TMBS suite. However, despite their simplistic appearance, all `ap_gym` environments
 205 impose significant challenges due to partial observability and non-immediate action payoffs. The tasks
 206 differ substantially from each other, testing the algorithm’s capability to handle diverse scenarios.

207 3.3 The Tactile MNIST Benchmark Suite

208 Tactile sensing presents unique challenges for
 209 perception: as an interactive modality, touch can
 210 unintentionally shift objects during exploration,
 211 and modern vision-based tactile sensors, such
 212 as GelSight [22] and DIGIT [46], produce in-
 213 herently high-dimensional observations. At the
 214 same time, tactile sensing provides highly local-
 215 ized information confined to points of contact,
 216 necessitating active perception.

217 The Tactile MNIST Benchmark Suite (TMBS)
 218 extends `ap_gym` to tactile tasks using vision-
 219 based tactile sensors. It contains four environ-
 220 ments: *TactileMNIST*, where the agent must classify a 3D model of a hand-written digit (see
 221 Section 3.4), *TactileMNISTVolume*, which tasks the agent to infer the volume of a given digit, *Toolbox*,
 222 where the agent must estimate the pose of a tool, and *Starstruck*, in which the agent must count the
 223 number of stars among other shapes (see Table 2 for an overview). In each environment, the agent
 224 controls a simulated GelSight Mini tactile sensor [22] and is presented with one or more objects on a
 225 platform. By interacting with the object, the agent must infer task-dependent properties, such as the

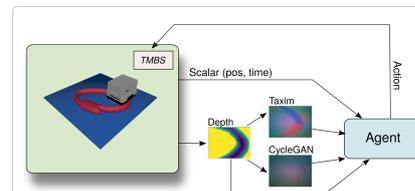


Figure 3: The agent receives the sensory information from the environment. This observation consists of scalar values and a tactile image. This can have three rendering modes, depth, Taxim, or CycleGAN.

226 class of the object, its pose, or its volume. Particularly, aside from tactile data and proprioception, the
 227 agent does not receive any additional sensory data, so it has to infer the required property from touch
 228 alone. Although, for simplicity and performance reasons, we do not simulate the physical interaction
 229 between the sensor and the object, we shift the objects around randomly to simulate unintended
 230 object movements.

231 To simulate the tactile sensor, we support three rendering modes:

232 **Taxim:** Taxim [55] computes an approximation of the gel deformation and afterwards applies
 233 a data-driven rendering algorithm.

234 **CycleGAN:** With data collected on 3D printed MNIST 3D objects (see Section 3.5), we train a
 235 CycleGAN [23] for a style transfer between a depth image and tactile image [58].
 236 The resulting images are visually much more realistic than the Taxim renderings, and
 237 thus might be beneficial for sim-to-real transfer. However, this mode is currently only
 238 available for the *TactileMNIST* environment. For more details, refer to Appendix C.

239 **Depth:** Here, the agent receives a depth image clipped to the GelSight gel thickness (4.25mm).

240 A comparison of all rendering modes is shown in Fig. 3, a visualization of the TactileMNIST digit
 241 classification task is provided in Fig. 2, and an overview of all tasks in TMBS is given in Table 2.

242 3.4 The MNIST 3D Dataset

243 *MNIST 3D*³ is a collection of 13,580 auto-
 244 generated 3d-printable meshes derived
 245 from a 500×500 - pixel high-resolution
 246 MNIST variant [59] and scaled to fit in a
 247 10×10 cm square. The MNIST 3D dataset
 248 poses an exciting tactile classification chal-
 249 lenge, as it has significant variability in
 250 shape and size within the classes, while
 251 also being large enough to facilitate learn-
 252 ing from data. A single touch is rarely
 253 enough to classify objects from this dataset,
 254 as segments of hand-written digits are usu-
 255 ally ambiguous. Hence, even after finding
 256 the object, the agent has to apply some strategy (e.g., contour following) to gather enough information
 257 for a successful classification. In addition to tactile sensing, this dataset could also be used as a
 258 benchmark for 3D mesh classification methods. More details on the generation of the MNIST 3D
 259 dataset can be found in Appendix D.

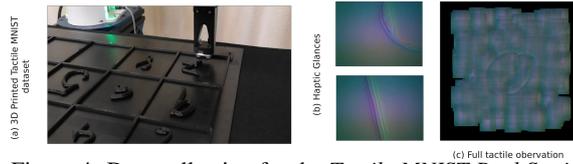


Figure 4: Data collection for the *Tactile MNIST Real Static* dataset. We mounted a GelSight Mini sensor on a Franka Research 3 (a) and collected 153,600 touches across 600 3D-printed MNIST digits. Examples of individual touches are visible in (b), and a collection of 256 touches overlaid based on their positions in (c).

260 3.5 A Large Dataset of Real Tactile Interactions

261 To complement simulated renders, TMBS includes a real-world, static tactile dataset captured with
 262 a GelSight sensor on 3D-printed MNIST 3D digits: the *Real Tactile MNIST Dataset*⁴. The dataset
 263 contains video sequences of 153,600 touches across 600 digits, which amounts to 256 touches per
 264 object collected in sequence. For data acquisition, we laid each 3D-printed MNIST digit in a 12×12 cm
 265 grid on a rubber mat and used a Franka Research 3 robot arm [60], with a GelSight tactile sensor to
 266 press the sensor down at random locations in the cell. Once we measured a normal force exceeding
 267 5N, we stopped pressing and registered the time stamp. To prevent degradation of the elastomer gel,
 268 we replaced the GelSight sensor’s gel pad after every 76,800 touches (i.e., halfway through each
 269 dataset). Finally, we partitioned each dataset into training (90 %) and test (10 %) splits, ensuring
 270 uniform class distributions across each split. Note that we also provide two processed versions of this
 271 dataset, where we replaced the videos with still images at the time of contact, one in full resolution at
 272 320×240 px⁵ and one scaled to 64×64 px⁶ for faster loading and training.

³ <https://huggingface.co/datasets/TimSchneider42/tactile-mnist-mnist3d>

⁴ <https://huggingface.co/datasets/TimSchneider42/tactile-mnist-touch-real-seq-t256-320x240>

⁵ <https://huggingface.co/datasets/TimSchneider42/tactile-mnist-touch-real-single-t256-320x240>

⁶ <https://huggingface.co/datasets/TimSchneider42/tactile-mnist-touch-real-single-t256-64x64>

273 We note that an additional challenge introduced during data collection is the possibility of the digit
274 shifting slightly due to contact with the sensor. This variability makes the dataset more representative
275 of real-world scenarios and provides an opportunity for methods to learn robustness to object
276 movement and misalignment. Thus, the Real Tactile MNIST Dataset serves several key roles in the
277 TMBS benchmark. First, it enables training of a CycleGAN for realistic simulation of tactile images
278 and sim-to-real transfer. Second, it can be used for both pretraining and fine-tuning of learning-based
279 perception models, enabling models to acquire basic tactile features before being deployed on a robot
280 for active learning tasks and again facilitating sim-to-real transfer. Finally, the dataset provides a
281 reproducible, offline benchmark for validating and comparing active perception algorithms under
282 realistic sensor noise and material artifacts. For additional details on the data collection procedure,
283 refer to Appendix E.

284 3.6 Evaluation Protocols

285 In `ap_gym` environments, there are two levels of exploration: (1) during an episode, the agent must
286 explore to gather information, and (2), over the course of the training, the agent must explore the
287 effects of its actions to optimize its model and policy. To disambiguate the measures of performance
288 in these two levels, we will call the first one *exploration efficiency* and the second *sample efficiency*.
289 Importantly, the former is a quality measure of the policy, while the latter is a quality measure of the
290 learning algorithm. Here we borrow the term *sample efficiency* from the RL literature, where it refers
291 to the number of environment interactions the agent needs in order to learn to solve the given task.
292 By *exploration efficiency*, on the other hand, we refer to the efficiency with which the agent collects
293 information within an episode.

294 Regarding *exploration efficiency*, we consider two measures: the *average* prediction accuracy and
295 the *final* prediction accuracy. Here, *average* prediction accuracy means the prediction accuracy the
296 agent exhibited throughout an episode on average, while *final* prediction accuracy means the accuracy
297 the agent exhibited at the final step of the episode. Normally, the agent starts each episode with
298 little to no information and then keeps gathering information as the episode progresses. Hence, for a
299 rational agent, we expect the prediction accuracy to increase over the course of the episode and to be
300 highest at the end of the episode. Thus, the *average* prediction accuracy could be seen as a measure
301 of how quickly the agent explored, while the *final* prediction accuracy could be seen as a measure
302 of how thoroughly the agent explored throughout the episode. In `ap_gym` environments, both of
303 these measures are tracked for a number of environment-specific prediction accuracy metrics, such as
304 classification accuracy, mean-squared-error, pose error, and others. For a detailed list of metrics for
305 each environment, refer to Appendix F.

306 Approaches evaluating on `ap_gym` or the Tactile MNIST benchmark suite should report both *average*
307 and *final* metric values over the course of the training. If applicable, the metrics should be computed
308 on the test variants of the environments, which use the test split instead of the training split. The
309 objective is to maximize both *sample efficiency* and *exploration efficiency*. Section 4 serves as an
310 example for an evaluation report on `ap_gym` and Tactile MNIST environments.

311 4 Experiments

312 In this section, we highlight experiments across selected environments from `ap_gym` and TMBS for
313 various baseline methods, including TAP [52] and HAM [8]. Both TAP and HAM are RL-based active
314 perception methods and, thus, well suited for evaluation on TMBS. The main difference between
315 them is that TAP employs an actor-critic RL approach in combination with a transformer architecture,
316 while HAM relies on a REINFORCE gradient in combination with an LSTM model. TAP provides two
317 variants: TAP-SAC and TAP-CrossQ, based on SAC [61] and CrossQ [62]. We additionally evaluate
318 a baseline that uses TAP’s transformer model with PPO [63], which we call TAP-PP0.

319 We highlight experiments on four environments in total. In the CircleSquare environment, the agent
320 can move a glimpse around an image and has to find and classify an object that can be either a circle
321 or a square. The TactileMNIST environment tasks the agent to classify MNIST 3D models by touch
322 alone, and in the Starstruck environment, the agent must count the number of stars (1-3) among other
323 objects on the platform. In the Toolbox environment, the agent must find a tool and determine its 2D
324 pose and orientation on the platform. These environments represent a diverse set of challenges, and
325 solving them requires the agent to adopt efficient exploration strategies, which is made evident by the

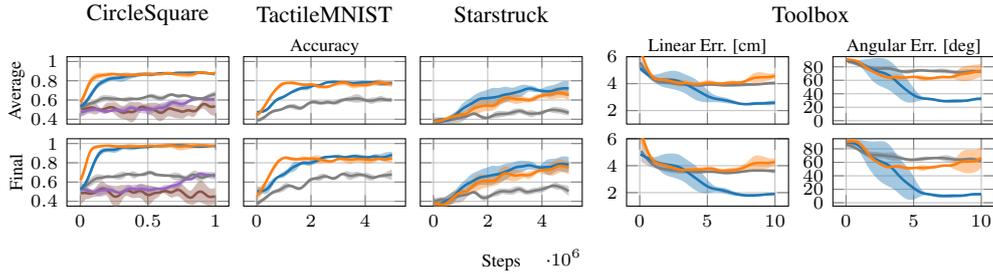


Figure 5: Average and final prediction accuracies for the baseline methods TAP-SAC, TAP-CrossQ, TAP-PP0, HAM [8], and a random baseline TAP-RND for the CircleSquare (ap_gym), TactileMNIST (TMBS), Starstruck (TMBS), and Toolbox (TMBS) environments. All methods were trained on 5 seeds for up to 10M. Shaded areas represent one standard deviation. Metrics are computed on evaluation tasks with unseen objects, except for Circle-Square and Toolbox, which have only two and one, respectively. For Starstruck, a correct classification requires predicting the exact number of stars. For Toolbox, we compute the linear and angular displacement between the prediction and the actual object pose as a metric. As HAM does not have a vision encoder, we evaluate it on non-tactile environments only. For TAP-PP0, we found it to be unstable in combination with a vision encoder, so we resort to only evaluating it on non-tactile environments as well.

326 consistently poor performance of TAP’s random baseline TAP-RND in Fig. 5. Further experiments and
 327 details for the training can be found in Appendix G.

328 As visible in Fig. 5 and Appendix G, TAP’s off-policy methods perform consistently better than
 329 the on-policy baselines HAM and TAP-PP0. This gap is likely due to on-policy methods generally
 330 being more sample-efficient than off-policy methods, as on-policy methods cannot reuse previously
 331 collected samples. However, despite the better performance of TAP, neither the ap_gym tasks nor the
 332 TMBS can be considered solved. TAP requires millions of environment interactions to learn viable
 333 exploration policies and falls short of perfect accuracy. More research in the area of sample-efficient
 334 RL and active perception is needed to improve sample efficiency to allow for the deployment of such
 335 methods in the real world.

336 5 Limitations and Conclusion

337 In this paper, we have introduced *Active Perception Gym* (ap_gym), a Gymnasium-compatible
 338 benchmark suite tailored for evaluating active perception tasks, and the *Tactile MNIST Benchmark*
 339 *Suite* (TMBS), which contains four tactile-specific tasks designed for robust exploration. To support
 340 these benchmarks, we have released a dataset comprising 13,580 high-resolution 3D digit models and
 341 an extensive real-world dataset of 153,600 tactile samples collected from 600 3D-printed digits using
 342 a GelSight Mini sensor. Provided as an open-source framework, ap_gym and TMBS offer a structured,
 343 standardized, and reproducible benchmark intended to facilitate advancements in active perception
 344 research, including efficient exploration strategies, sim-to-real adaptation through CycleGAN training,
 345 and the pretraining and fine-tuning of tactile models.

346 However, our benchmark has limitations, most notably the absence of online, real-world evalua-
 347 tion scenarios and metrics beyond the static dataset. While there are established datasets such as
 348 YCB [64] enable benchmarking of contact-rich manipulation tasks using real-world objects, our
 349 suite deliberately focuses on controlled, simulation-based exploration to ensure reproducibility and
 350 interoperability. In future work, we aim to address this limitation by designing carefully structured
 351 real-world tactile exploration experiments that extend the benchmark’s relevance to physical robotic
 352 systems and support the study of sim-to-real transfer in active perception. Another key limitation is
 353 the absence of a physics engine in our simulation environment, which currently prevents modeling of
 354 more complex, contact-rich interactions. As a result, tasks such as grasping, 3D object reconstruction
 355 (where objects may tumble upon contact), object retrieval in cluttered scenes, and in-hand pose
 356 estimation remain out of scope. Extending our framework to incorporate physics engines would open
 357 the door to these richer interaction scenarios and significantly broaden the benchmark’s applicability.
 358 Ultimately, we view this benchmark as a foundational resource for advancing active tactile perception
 359 and enabling the development of more robust, efficient algorithms for robotic tactile manipulation.

360 References

- 361 [1] Yann LeCun. The mnist database of handwritten digits. *http://yann.lecun.com/exdb/mnist/*,
362 1998.
- 363 [2] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
364 2009.
- 365 [3] Yann Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 366 [4] Roland S Johansson and J Randall Flanagan. Coding and use of tactile signals from the fingertips
367 in object manipulation tasks. *Nature Reviews Neuroscience*, 10(5):345–359, 2009.
- 368 [5] Susan J Lederman and Roberta L Klatzky. Hand movements: A window into haptic object
369 recognition. *Cognitive psychology*, 19(3):342–368, 1987.
- 370 [6] Tony J. Prescott, Mathew E. Diamond, and Alan M. Wing. Active touch sensing. *Phil. Trans. R.*
371 *Soc. B*, 366(1581):2989–2995, Nov 2011.
- 372 [7] A. Boehm, T. Schneider, B. Belousov, A. Kshirsagar, L. Lin, K. Doerschner, K. Drawing, C.A.
373 Rothkopf, and J. Peters. What matters for active texture recognition with vision-based tactile
374 sensors. In *Proceedings of the IEEE International Conference on Robotics and Automation*
375 *(ICRA)*, 2024.
- 376 [8] Sascha Fleer, Alexandra Moringen, Roberta L Klatzky, and Helge Ritter. Learning efficient
377 haptic shape exploration with a rigid tactile sensor array. *PloS one*, 15(1):e0226880, 2020.
- 378 [9] Jingxi Xu, Shuran Song, and Matei Ciocarlie. Tandem: Learning joint exploration and decision
379 making with tactile sensors. *IEEE Robotics and Automation Letters*, 7(4):10391–10398, 2022.
- 380 [10] Zhengkun Yi, Roberto Calandra, Filipe Veiga, Herke van Hoof, Tucker Hermans, Yilei Zhang,
381 and Jan Peters. Active tactile object exploration with gaussian processes. In *2016 IEEE/RSJ*
382 *International Conference on Intelligent Robots and Systems (IROS)*, pages 4925–4930. IEEE,
383 2016.
- 384 [11] Marten Björkman, Yasemin Bekiroglu, Virgile Högman, and Danica Kragic. Enhancing visual
385 perception of shape through tactile glances. In *2013 IEEE/RSJ International Conference on*
386 *Intelligent Robots and Systems*, pages 3180–3186. IEEE, 2013.
- 387 [12] Niklas Funk, Erik Helmut, Georgia Chalvatzaki, Roberto Calandra, and Jan Peters. Evetac: An
388 Event-Based Optical Tactile Sensor for Robotic Manipulation. *IEEE Trans. Rob.*, 40:3812–3832,
389 July 2024.
- 390 [13] Justin Kerr, Huang Huang, Albert Wilcox, Ryan Hoque, Jeffrey Ichnowski, Roberto Calandra,
391 and Ken Goldberg. Self-supervised visuo-tactile pretraining to locate and follow garment
392 features. *arXiv preprint arXiv:2209.13042*, 2022.
- 393 [14] Carolina Higuera, Akash Sharma, Chaithanya Krishna Bodduluri, Taosha Fan, Patrick Lancaster,
394 Mrinal Kalakrishnan, Michael Kaess, Byron Boots, Mike Lambeta, Tingfan Wu, and Mustafa
395 Mukadam. Sparsh: Self-supervised touch representations for vision-based tactile sensing. In
396 *8th Annual Conference on Robot Learning*, 2024.
- 397 [15] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale
398 Hierarchical Image Database. In *CVPR09*, 2009.
- 399 [16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
400 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*
401 *Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014,*
402 *Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- 403 [17] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving?
404 the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition*
405 *(CVPR)*, 2012.

- 406 [18] Garrick Brazil, Abhinav Kumar, Julian Straub, Nikhila Ravi, Justin Johnson, and Georgia
407 Gkioxari. Omni3d: A large benchmark and model for 3d object detection in the wild. In
408 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
409 13154–13164, 2023.
- 410 [19] Sudharshan Suresh, Haozhi Qi, Tingfan Wu, Taosha Fan, Luis Pineda, Mike Lambeta, Jitendra
411 Malik, Mrinal Kalakrishnan, Roberto Calandra, Michael Kaess, et al. Neuralfeels with neural
412 fields: Visuotactile perception for in-hand manipulation. *Science Robotics*, 9(96):eadl0628,
413 2024.
- 414 [20] Amir-Hossein Shahidzadeh, Seong Jong Yoo, Pavan Mantripragada, Chahat Deep Singh, Cor-
415 nelia Fermüller, and Yiannis Aloimonos. Actexplore: Active tactile exploration on unknown
416 objects. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages
417 3411–3418. IEEE, 2024.
- 418 [21] Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu,
419 Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A
420 standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*,
421 2024.
- 422 [22] Wenzhen Yuan, Siyuan Dong, and Edward H Adelson. Gelsight: High-resolution robot tactile
423 sensors for estimating geometry and force. *Sensors*, 17(12):2762, 2017.
- 424 [23] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image
425 translation using cycle-consistent adversarial networks, 2020.
- 426 [24] Sabhari Natarajan, Galen Brown, and Berk Calli. Aiding grasp synthesis for novel objects using
427 heuristic-based and data-driven active vision methods. *Frontiers in Robotics and AI*, 8:696587,
428 2021.
- 429 [25] Qianfan Zhao, Lu Zhang, Lingxi Wu, Hong Qiao, and Zhiyong Liu. A real 3d embodied dataset
430 for robotic active visual learning. *IEEE Robotics and Automation Letters*, 7(3):6646–6652,
431 2022.
- 432 [26] David Hall, Ben Talbot, Suman Raj Bista, Haoyang Zhang, Rohan Smith, Feras Dayoub,
433 and Niko Sünderhauf. The robotic vision scene understanding challenge. *arXiv preprint*
434 *arXiv:2009.05246*, 2020.
- 435 [27] Pourya Hoseini, Shuvo Kumar Paul, Mircea Nicolescu, and Monica Nicolescu. Next best view
436 planning in a single glance: An approach to improve object recognition. *SN Computer Science*,
437 4(1):51, 2022.
- 438 [28] Phil Ammirato, Patrick Poirson, Eunbyung Park, Jana Kosecka, and Alexander C. Berg. A
439 dataset for developing and benchmarking active vision. In *IEEE International Conference on*
440 *Robotics and Automation (ICRA)*, 2017.
- 441 [29] Phil Ammirato, Alexander C Berg, and Jana Kosecka. Active vision dataset benchmark. In
442 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*,
443 pages 2046–2049, 2018.
- 444 [30] Jie Wu, Tianshui Chen, Lishan Huang, Hefeng Wu, Guanbin Li, Ling Tian, and Liang Lin.
445 Active object search. In *Proceedings of the 28th ACM International Conference on Multimedia*,
446 pages 973–981, 2020.
- 447 [31] Ziyue Wang, Chi Chen, Fuwen Luo, Yurui Dong, Yuanchi Zhang, Yuzhuang Xu, Xiaolong
448 Wang, Peng Li, and Yang Liu. Actiview: Evaluating active perception ability for multimodal
449 large language models, 2024.
- 450 [32] Tianyi Liu and Benjamin Ward-Cherrier. The tip benchmark: A tactile image-based
451 psychophysics-inspired benchmark for artificial tactile sensors. In *International Conference on*
452 *Human Haptic Sensing and Touch Enabled Computer Applications*, pages 94–106. Springer,
453 2024.

- 454 [33] Jialiang Zhao, Yuxiang Ma, Lirui Wang, and Edward H. Adelson. Transferable tactile trans-
455 formers for representation learning across diverse sensors and tasks, 2024.
- 456 [34] Sudharshan Suresh, Zilin Si, Stuart Anderson, Michael Kaess, and Mustafa Mukadam. Midas-
457 Touch: Monte-Carlo inference over distributions across sliding touch. In *Proc. Conf. on Robot*
458 *Learning, CoRL*, Auckland, NZ, December 2022.
- 459 [35] Wenzhen Yuan, Yuchen Mo, Shaoxiong Wang, and Edward H Adelson. Active clothing material
460 perception using tactile sensing and deep learning. In *2018 IEEE International Conference on*
461 *Robotics and Automation (ICRA)*, pages 4842–4849. IEEE, 2018.
- 462 [36] Fengyu Yang, Chenyang Ma, Jiacheng Zhang, Jing Zhu, Wenzhen Yuan, and Andrew
463 Owens. Touch and go: Learning from human-collected vision and touch. *arXiv preprint*
464 *arXiv:2211.12498*, 2022.
- 465 [37] Shan Luo, Wenzhen Yuan, Edward Adelson, Anthony G Cohn, and Raul Fuentes. Vitac:
466 Feature sharing between vision and tactile sensing for cloth texture recognition. In *2018 IEEE*
467 *International Conference on Robotics and Automation (ICRA)*, pages 2722–2727. IEEE, 2018.
- 468 [38] Wenzhen Yuan, Shaoxiong Wang, Siyuan Dong, and Edward Adelson. Connecting look and
469 feel: Associating the visual and tactile properties of physical materials. In *Proceedings of the*
470 *IEEE conference on computer vision and pattern recognition*, pages 5580–5588, 2017.
- 471 [39] Yunzhu Li, Jun-Yan Zhu, Russ Tedrake, and Antonio Torralba. Connecting touch and vision via
472 cross-modal prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
473 *Pattern Recognition*, pages 10609–10618, 2019.
- 474 [40] Samson Yu, Kelvin Lin, Anxing Xiao, Jiafei Duan, and Harold Soh. Octopi: Object property
475 reasoning with large tactile-language models. *arXiv preprint arXiv:2405.02794*, 2024.
- 476 [41] Letian Fu, Gaurav Datta, Huang Huang, William Chung-Ho Panitch, Jaimyn Drake, Joseph
477 Ortiz, Mustafa Mukadam, Mike Lambeta, Roberto Calandra, and Ken Goldberg. A touch,
478 vision, and language dataset for multimodal alignment. *arXiv preprint arXiv:2402.13232*, 2024.
- 479 [42] Ning Cheng, Changhao Guan, Jing Gao, Weihao Wang, You Li, Fandong Meng, Jie Zhou, Bin
480 Fang, Jinan Xu, and Wenjuan Han. Touch100k: A large-scale touch-language-vision dataset for
481 touch-centric multimodal representation. *arXiv preprint arXiv:2406.03813*, 2024.
- 482 [43] Ruohan Gao, Yiming Dou, Hao Li, Tanmay Agarwal, Jeannette Bohg, Yunzhu Li, Li Fei-Fei,
483 and Jiajun Wu. The objectfolder benchmark: Multisensory learning with neural and real objects.
484 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
485 pages 17276–17286, 2023.
- 486 [44] Ruohan Gao, Zilin Si, Yen-Yu Chang, Samuel Clarke, Jeannette Bohg, Li Fei-Fei, Wenzhen
487 Yuan, and Jiajun Wu. Objectfolder 2.0: A multisensory object dataset for sim2real transfer. In
488 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
489 pages 10598–10608, 2022.
- 490 [45] Ruohan Gao, Yen-Yu Chang, Shivani Mall, Li Fei-Fei, and Jiajun Wu. Objectfolder: A
491 dataset of objects with implicit visual, auditory, and tactile representations. *arXiv preprint*
492 *arXiv:2109.07991*, 2021.
- 493 [46] Mike Lambeta, Po-Wei Chou, Stephen Tian, Brian Yang, Benjamin Maloon, Victoria Rose
494 Most, Dave Stroud, Raymond Santos, Ahmad Byagowi, Gregg Kammerer, et al. Digit: A
495 novel design for a low-cost compact high-resolution tactile sensor with application to in-hand
496 manipulation. *IEEE Robotics and Automation Letters*, 5(3):3838–3845, 2020.
- 497 [47] Zhao-Heng Yin, Binghao Huang, Yuzhe Qin, Qifeng Chen, and Xiaolong Wang. Rotating
498 without seeing: Towards in-hand dexterity through touch. *arXiv preprint arXiv:2303.10880*,
499 2023.
- 500 [48] Haozhi Qi, Brent Yi, Sudharshan Suresh, Mike Lambeta, Yi Ma, Roberto Calandra, and Jitendra
501 Malik. General in-hand object rotation with vision and touch. In *Conference on Robot Learning*,
502 pages 2549–2564. PMLR, 2023.

- 503 [49] Kenneth Y Goldberg and Ruzena Bajcsy. Active touch and robot perception. *Cognition and*
504 *Brain Theory*, 7(2):199–214, 1984.
- 505 [50] Ruzena Bajcsy, Yiannis Aloimonos, and John K Tsotsos. Revisiting active perception. *Au-*
506 *tonomous Robots*, 42:177–196, 2018.
- 507 [51] Cristiana De Farias, Naresh Marturi, Rustam Stolkin, and Yasemin Bekiroglu. Simultaneous
508 tactile exploration and grasp refinement for unknown objects. *IEEE Robotics and Automation*
509 *Letters*, 6(2):3349–3356, 2021.
- 510 [52] Tim Schneider, Cristiana de Farias, Roberto Calandra, Liming Chen, and Jan Peters. Active
511 perception for tactile sensing: A task-agnostic attention-based approach, 2025.
- 512 [53] Ruohan Gao, Zilin Si, Yen-Yu Chang, Samuel Clarke, Jeannette Bohg, Li Fei-Fei, Wenzhen
513 Yuan, and Jiajun Wu. Objectfolder 2.0: A multisensory object dataset for sim2real transfer. In
514 *CVPR*, 2022.
- 515 [54] Shaoxiong Wang, Mike Lambeta, Po-Wei Chou, and Roberto Calandra. Tacto: A fast, flexible,
516 and open-source simulator for high-resolution vision-based tactile sensors. *IEEE Robotics and*
517 *Automation Letters*, 7(2):3930–3937, 2022.
- 518 [55] Zilin Si and Wenzhen Yuan. Taxim: An example-based simulation model for gelsight tactile
519 sensors. *IEEE Robotics and Automation Letters*, 7(2):2361–2368, 2022.
- 520 [56] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
521 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
522 Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- 523 [57] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
524 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
525 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
526 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,
527 high-performance deep learning library. In *Advances in Neural Information Processing Systems*
528 32, pages 8024–8035. Curran Associates, Inc., 2019.
- 529 [58] Weihang Chen, Yuan Xu, Zhenyang Chen, Peiyu Zeng, Renjun Dang, Rui Chen, and Jing Xu.
530 Bidirectional sim-to-real transfer for gelsight tactile sensors with cyclegan. *IEEE Robotics and*
531 *Automation Letters*, 7(3):6187–6194, 2022.
- 532 [59] Cédric Beaulac and Jeffrey S. Rosenthal. Introducing a new high-resolution handwritten digits
533 data set with writer characteristics. *SN Computer Science*, 4(1), November 2022.
- 534 [60] Franka Inc. Homepage, January 2025. [Online; accessed 28. Jan. 2025].
- 535 [61] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan,
536 Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms
537 and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- 538 [62] Aditya Bhatt, Daniel Palenicek, Boris Belousov, Max Argus, Artemij Amiranashvili, Thomas
539 Brox, and Jan Peters. Crossq: Batch normalization in deep reinforcement learning for greater
540 sample efficiency and simplicity. *arXiv preprint arXiv:1902.05605*, 2019.
- 541 [63] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
542 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 543 [64] Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M.
544 Dollar. The ycb object and model set: Towards common benchmarks for manipulation research.
545 In *2015 International Conference on Advanced Robotics (ICAR)*, pages 510–517, 2015.

546 **NeurIPS Paper Checklist**

547 **1. Claims**

548 Question: Do the main claims made in the abstract and introduction accurately reflect the
549 paper’s contributions and scope?

550 Answer: [Yes]

551 Justification: The main claim of the paper is the introduction of the first benchmarking
552 protocol for active tactile perception. Section 2 supports the claim we are the first with a
553 comprehensive literature review. Furthermore, our two benchmark suites and the datasets
554 are introduced in Section 3, with additional details in the appendix and are provided online.

555 Guidelines:

- 556 • The answer NA means that the abstract and introduction do not include the claims
557 made in the paper.
- 558 • The abstract and/or introduction should clearly state the claims made, including the
559 contributions made in the paper and important assumptions and limitations. A No or
560 NA answer to this question will not be perceived well by the reviewers.
- 561 • The claims made should match theoretical and experimental results, and reflect how
562 much the results can be expected to generalize to other settings.
- 563 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
564 are not attained by the paper.

565 **2. Limitations**

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

567 Answer: [Yes]

568 Justification: In Section 5 we discussed limitations and future directions of the work.

569 Guidelines:

- 570 • The answer NA means that the paper has no limitation while the answer No means that
571 the paper has limitations, but those are not discussed in the paper.
- 572 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 573 • The paper should point out any strong assumptions and how robust the results are to
574 violations of these assumptions (e.g., independence assumptions, noiseless settings,
575 model well-specification, asymptotic approximations only holding locally). The authors
576 should reflect on how these assumptions might be violated in practice and what the
577 implications would be.
- 578 • The authors should reflect on the scope of the claims made, e.g., if the approach was
579 only tested on a few datasets or with a few runs. In general, empirical results often
580 depend on implicit assumptions, which should be articulated.
- 581 • The authors should reflect on the factors that influence the performance of the approach.
582 For example, a facial recognition algorithm may perform poorly when image resolution
583 is low or images are taken in low lighting. Or a speech-to-text system might not be
584 used reliably to provide closed captions for online lectures because it fails to handle
585 technical jargon.
- 586 • The authors should discuss the computational efficiency of the proposed algorithms
587 and how they scale with dataset size.
- 588 • If applicable, the authors should discuss possible limitations of their approach to
589 address problems of privacy and fairness.
- 590 • While the authors might fear that complete honesty about limitations might be used by
591 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
592 limitations that aren’t acknowledged in the paper. The authors should use their best
593 judgment and recognize that individual actions in favor of transparency play an impor-
594 tant role in developing norms that preserve the integrity of the community. Reviewers
595 will be specifically instructed to not penalize honesty concerning limitations.

596 **3. Theory assumptions and proofs**

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

599 Answer: [NA]

600 Justification: [NA] .

601 Guidelines:

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

612 4. Experimental result reproducibility

613 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

614 perimental results of the paper to the extent that it affects the main claims and/or conclusions

615 of the paper (regardless of whether the code and data are provided or not)?

616 Answer: [Yes]

617 Justification: Throughout Section 3 we present our benchmark suites. Throughout the

618 appendix, and particularly in Appendix F and Appendix H we present detailed description on

619 how each of our environments work as well as implementation details. We add that our code

620 is also available online with detailed documentation. Furthermore, for the dataset collection,

621 we fully describe the steps taken to generate the MNIST 3D Meshes (in Section 3.4 and

622 Appendix D) and the real tactile data collection (in Section 3.5 and Appendix E).

623 Guidelines:

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

653 some way (e.g., to registered users), but it should be possible for other researchers
654 to have some path to reproducing or verifying the results.

655 5. Open access to data and code

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

659 Answer: [Yes]

660 Justification: We provide all the code for our benchmark environments with extensive
661 documentation.

662 Guidelines:

- 663 • The answer NA means that paper does not include experiments requiring code.
- 664 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
665 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 666 • While we encourage the release of code and data, we understand that this might not be
667 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
668 including code, unless this is central to the contribution (e.g., for a new open-source
669 benchmark).
- 670 • The instructions should contain the exact command and environment needed to run to
671 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
672 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 673 • The authors should provide instructions on data access and preparation, including how
674 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 675 • The authors should provide scripts to reproduce all experimental results for the new
676 proposed method and baselines. If only a subset of experiments are reproducible, they
677 should state which ones are omitted from the script and why.
- 678 • At submission time, to preserve anonymity, the authors should release anonymized
679 versions (if applicable).
- 680 • Providing as much information as possible in supplemental material (appended to the
681 paper) is recommended, but including URLs to data and code is permitted.

682 6. Experimental setting/details

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

686 Answer: [Yes]

687 Justification: We provide details both in Section 4, Appendix G and Appendix H

688 Guidelines:

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

694 7. Experiment statistical significance

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

697 Answer: [Yes]

698 Justification: Our results in Section 4 and Appendix G show one standard deviation for each
699 learning curve.

700 Guidelines:

- 701 • The answer NA means that the paper does not include experiments.
- 702 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
703 dence intervals, or statistical significance tests, at least for the experiments that support
704 the main claims of the paper.

- 705 • The factors of variability that the error bars are capturing should be clearly stated (for
706 example, train/test split, initialization, random drawing of some parameter, or overall
707 run with given experimental conditions).
- 708 • The method for calculating the error bars should be explained (closed form formula,
709 call to a library function, bootstrap, etc.)
- 710 • The assumptions made should be given (e.g., Normally distributed errors).
- 711 • It should be clear whether the error bar is the standard deviation or the standard error
712 of the mean.
- 713 • It is OK to report 1-sigma error bars, but one should state it. The authors should
714 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
715 of Normality of errors is not verified.
- 716 • For asymmetric distributions, the authors should be careful not to show in tables or
717 figures symmetric error bars that would yield results that are out of range (e.g. negative
718 error rates).
- 719 • If error bars are reported in tables or plots, The authors should explain in the text how
720 they were calculated and reference the corresponding figures or tables in the text.

721 8. Experiments compute resources

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

725 Answer: [Yes]

726 Justification: All implementation and performance details for the experiments are present in
727 Section 4, Appendix G and Appendix H

728 Guidelines:

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

737 9. Code of ethics

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

740 Answer: [Yes]

741 Justification: Our research fully conforms to the NeurIPS Code of Ethics. No human
742 subjects or participants were involved in the creation of the datasets used, so there are no
743 related privacy or consent concerns. All data and code are released publicly on GitHub
744 and Hugging Face under permissive licenses (MIT for code, CC-BY for data), ensuring
745 long-term accessibility, transparency, and reproducibility. We have verified that none of the
746 datasets used are deprecated and that they comply with licensing terms. Our work poses no
747 foreseeable risks related to safety, security, discrimination, or environmental harm.

748 Guidelines:

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

754 10. Broader impacts

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

757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810

Answer: [NA] .

Justification: This work introduces a benchmark for active tactile perception in robotic systems. It is foundational in nature and not tied to any direct applications or deployments. While it may contribute to progress in robotic manipulation over the long term, we do not foresee any immediate societal impact from this work.

Guidelines:

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

11. Safeguards

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

Answer: [NA] .

Justification: None of our data or benchmark suites has any significant risk of misuse.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: Yes, all external assets used in this work are properly cited in both the main paper and the appendix. We explicitly list the main assets used in our environment in Table 2, along with their associated licenses. All listed assets are released under a CC-BY license, and we have ensured that their terms of use are fully respected.

- 811 Guidelines:
- 812 • The answer NA means that the paper does not use existing assets.
 - 813 • The authors should cite the original paper that produced the code package or dataset.
 - 814 • The authors should state which version of the asset is used and, if possible, include a
 - 815 URL.
 - 816 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - 817 • For scraped data from a particular source (e.g., website), the copyright and terms of
 - 818 service of that source should be provided.
 - 819 • If assets are released, the license, copyright information, and terms of use in the
 - 820 package should be provided. For popular datasets, `paperswithcode.com/datasets`
 - 821 has curated licenses for some datasets. Their licensing guide can help determine the
 - 822 license of a dataset.
 - 823 • For existing datasets that are re-packaged, both the original license and the license of
 - 824 the derived asset (if it has changed) should be provided.
 - 825 • If this information is not available online, the authors are encouraged to reach out to
 - 826 the asset’s creators.

827 **13. New assets**

828 Question: Are new assets introduced in the paper well documented and is the documentation

829 provided alongside the assets?

830 Answer: [Yes]

831 Justification: All the code, models, and datasets we have released come with detailed docu-

832 mentation. Additionally, Section 3, Appendix F, Appendix D, Appendix E and Appendix H

833 provide details of the released assets.

834 Guidelines:

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

843 **14. Crowdsourcing and research with human subjects**

844 Question: For crowdsourcing experiments and research with human subjects, does the paper

845 include the full text of instructions given to participants and screenshots, if applicable, as

846 well as details about compensation (if any)?

847 Answer: [NA] .

848 Justification: -

849 Guidelines:

- 850 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 851 human subjects.
- 852 • Including this information in the supplemental material is fine, but if the main contribu-
- 853 tion of the paper involves human subjects, then as much detail as possible should be
- 854 included in the main paper.
- 855 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 856 or other labor should be paid at least the minimum wage in the country of the data
- 857 collector.

858 **15. Institutional review board (IRB) approvals or equivalent for research with human**

859 **subjects**

860 Question: Does the paper describe potential risks incurred by study participants, whether

861 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

862 approvals (or an equivalent approval/review based on the requirements of your country or

863 institution) were obtained?

864 Answer: [NA]
865 Justification: The paper does not involve crowdsourcing or research with human subjects
866 Guidelines:
867 • The answer NA means that the paper does not involve crowdsourcing nor research with
868 human subjects.
869 • Depending on the country in which research is conducted, IRB approval (or equivalent)
870 may be required for any human subjects research. If you obtained IRB approval, you
871 should clearly state this in the paper.
872 • We recognize that the procedures for this may vary significantly between institutions
873 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
874 guidelines for their institution.
875 • For initial submissions, do not include any information that would break anonymity (if
876 applicable), such as the institution conducting the review.

877 **16. Declaration of LLM usage**

878 Question: Does the paper describe the usage of LLMs if it is an important, original, or
879 non-standard component of the core methods in this research? Note that if the LLM is used
880 only for writing, editing, or formatting purposes and does not impact the core methodology,
881 scientific rigorousness, or originality of the research, declaration is not required.

882 Answer: [NA].

883 Justification: This research does not involve LLMs as any important, original, or non-
884 standard components.

885 Guidelines:
886 • The answer NA means that the core method development in this research does not
887 involve LLMs as any important, original, or non-standard components.
888 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
889 for what should or should not be described.