
TOUCH and GO: A Real-World Multisensory Dataset

Chenyang Ma*
University of Michigan
dannymcy@umich.edu

Fengyu Yang*
University of Michigan
fredyang@umich.edu

Andrew Owens
University of Michigan
ahowens@umich.edu

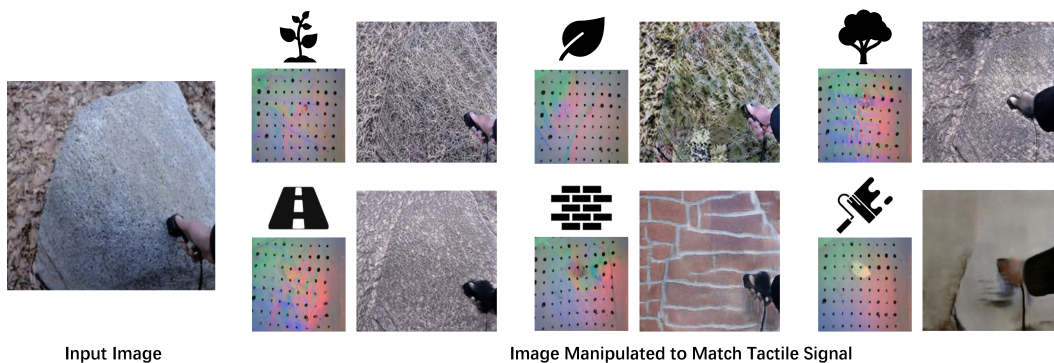


Figure 1: *Tactile-Guided Image Stylization* (TGIS). We present a method for manipulating the appearance of an object to match its material property indicated by the tactile signal.

Abstract

1 Unlike how humans perceive the world from associations between senses and
2 through a series of inanimate objects, contemporary research on robot perception
3 problem mainly rely on vision units or visual inputs to teach the robots interact
4 with the world. We identify that this is due to the lack of real-world multisensory
5 rich object dataset. To tackle this challenge, we present **TOUCH and GO**, a
6 multisensory dataset containing real-world synchronized high-quality video and
7 tactile data containing 12600 object instances over 37800 touches and 30 hours of
8 video captured from egocentric viewpoint, greatly exceeding the size of existing
9 real-world multisensory datasets. All objects in our dataset are originated from
10 real environments with fine-grained textures retained. We propose and apply our
11 dataset on two novel tasks, tactile-guided image stylization and multi-modal video
12 prediction on tactile images.

13 1 Introduction

14 Humans perceive the world not using a single modality. Instead, we have access to many sensory
15 streams and learn from associations between senses. When a child eats an apple, for instance, she'll
16 not only taste it—she'll also hear it crunch, see its shiny skin, and feel its smooth surface [48]. In
17 addition, humans perceive the world not as a single giant entity but often through a series of inanimate
18 objects, which exist as bounded wholes and move on connected paths. We interact with these objects

*The first two authors contributed equally to this paper.

19 through an array of different sensory systems—vision, touch, audition, smell, taste, and proprioception.
20 These multisensory inputs shape our daily experiences.

21 It is so nature for humans to learn knowledge through interactions with different objects with multiple
22 senses. Cognitive science studies [49, 48] show that both object representation and multisensory
23 perception play a crucial role in early human cognitive development. However, for robots, it may
24 not be the case. Contemporary research on robot perception problem mainly rely on vision units or
25 visual inputs to teach the robots perceive and interact with the world. This focus on learning from
26 vision alone makes the perception problem harder because some of the most important spectrum of
27 physical object properties and sensory modes — such as touching — are lost.

28 We identify that this is due to the lack of real-world multisensory rich object dataset. Several works
29 have been done regarding simulated multisensory dataset [16, 17, 14]. However, we argue there
30 are two fundamental differences between the quality and utility of simulated and real-world dataset.
31 First, simulated data fail to perfectly represent reality. Models trained purely on synthetic data do not
32 generalize to the real world due to the discrepancy between simulated and real environments, in terms
33 of both visual and physical properties. In fact, the more we increase the fidelity of our simulations,
34 the more effort we have to expend in order to build them, both in terms of implementing complex
35 physical phenomena and in terms of creating the content (e.g., objects, backgrounds) to populate these
36 simulations. This difficulty is compounded by the fact that powerful optimization methods based
37 on deep learning are exceptionally proficient at exploiting simulator flaws: the more powerful the
38 machine learning algorithm, the more likely it is to discover how to "cheat" the simulator to succeed
39 in ways that are infeasible in the real world [4]. Second, "reality gap" exists by transferring simulated
40 experience into the real world. While simulated data continue to improve in fidelity, the peculiar and
41 pathological regularities of synthetic data, and the wide, unpredictable diversity of real-world objects,
42 makes bridging the reality gap particularly difficult when the robot use its sensors to perceive the
43 world, as is the case for example in many manipulation tasks [4, 29].

44 Therefore, our goal is to establish a real-world multisensory dataset containing rich objects that are
45 1) easily accessible to the community as a standard benchmark, 2) high-quality in terms of visual
46 textures, and 3) augmented with real data from the perspective of human beings. To this end, we
47 introduce TOUCH and GO — an egocentric multisensory dataset of synchronized video and tactile
48 sensing. We take inspiration from the way infants explore the physical properties of a scene by poking
49 and prodding at the objects in front of them [3, 46], a process that may help them learn an intuitive
50 theory of physics. The egocentric viewpoint enables our dataset to contain enough details to observe
51 the fine-grained texture of objects, and mimics the perception of a real human.

52 More specifically, we collect over 30 hours of real-world synchronized high-quality video and
53 tactile data containing 12600 object instances over 37800 touches. Our dataset contains rich objects
54 categories from both indoor and outdoor scenes (none of the existing real-world multi-modal dataset
55 contains data from outdoor scenes). TOUCH and GO enables many applications. We present a
56 method for manipulating the appearance of an object to match its material property indicated by the
57 tactile signal, a problem we term *Tactile-Guided Image Stylization* (TGIS), as shown in Figure 1. We
58 also propose a novel multi-modal video prediction problem on tactile image deformation. For the
59 first task, We design a deep neural network based on CUT [44], which fuses data from video and
60 tactile streams. For both tasks, experimental results suggest better results are achieved by leveraging
61 our TOUCH and GO dataset.

62 Our main contributions can be concluded as the followings: 1) We introduce TOUCH and GO, a
63 real-world dataset that makes multisensory learning with vision and touch easily accessible to the
64 research community. 2) All objects in our dataset are originated from real environments and will
65 be made publicly available as a standard testbed for robotic multisensory learning. 3) We propose
66 and apply our dataset on two novel tasks including tactile-guided image stylization and multi-modal
67 video prediction on tactile images.

68 2 Related Work

69 **Multisensory Datasets** There is a mixture of real and simulated data across different single-modal
70 datasets. ImageNet [12], MS COCO [35], ObjectNet [1], and OpenImages [30] focus on the collection
71 of large-scale real 2D images. ModelNet [56] and ShapeNet [9] contain synthetic 3D CAD models,
72 emphasizing on geometry of 3D objects but pay less attention to fine-grained visual textures. BigBIRD
73 [47], YCB [8], and ABO [10] model real-world 3D objects with limited object instances. The majority
74 of multi-modal datasets incorporate simulated data. Pix3D [51], IKEA Objects [34], and Object3D
75 [57] match synthetic 3D CAD models to objects in real images. OBJECT-FOLDER 1.0 [16] contains
76 multisensory simulated data as implicit neural representations. Built upon it, OBJECT-FOLDER
77 2.0 [17] is ten times larger than the previous version with encoding of more realistic data. A few
78 real-world multi-modal datasets exist. VisGel [33] comprises real-world data of videos and touches
79 collected by robotics arm, thus has very restricted scenes and bias introduced by the arm. Greatest
80 Hits [43] contains high-quality egocentric videos of humans probing environments with a drumstick,
81 but its goal is not on scale expansion and generalization. Our TOUCH and GO dataset contains
82 high-quality synchronized RGB video and tactile data, with over 30 hours of videos, 37800 touches,
83 and 12600 object instances, which greatly exceeds the size of existing real-world egocentric datasets.

84 **Touch and Vision** Researches are conducted on the types of haptic, force, and tactile sensors to
85 give robots tactile sensing ability [11, 25, 32, 31]. GelSight [23, 22, 59, 7] is widely adapted as a
86 high-resolution tactile sensor for computer vision and robotics applications, which includes improving
87 grasp stability with rotation measurement [28], the study of the physical and material properties of
88 fabrics [60], predicting the grasping success through both vision and tactile sensing [7], and cloth
89 texture recognition [38]. Here we introduce the novel application of tactile-guided image stylization.

90 **Image-to-Image Stylization** Image stylization (translation) translates an input image from one
91 domain to a photo realistic output in the target domain [20, 36, 61]. The key to the success of this
92 task is due to the emergence of generative adversarial networks (GAN) [18, 41], which have been
93 vigorously researched in the last several years with many applications including generating photos
94 from sketches [20, 45], changing time of a day [20, 63], and translating semantic meanings into
95 scenes [20, 55]. While most of the image stylization tasks have paired image-to-image translation,
96 in certain cases, the corresponding examples from domains are unavailable, resulting in unpaired
97 image-to-image stylization. Cycle consistency [26, 58, 62], as one of the approaches, enforce the
98 correspondence between the input and output image domain by adopting the underlying bijective
99 assumption, which may be too restrictive in cases when images from one domain contain additional
100 information compared to the other domain. CUT [44] adapts contrastive learning to make each
101 patch in the output reflect the content of the corresponding patch in the input by maximizing mutual
102 information between the two. We propose a new model based on CUT, which receives multi-modal
103 data as inputs and learns to build the tactile-visual style associations without any human supervision,
104 for our proposed novel task of tactile-guided image stylization.

105 **Video Prediction** Approaches for video prediction are diverse, evolving from the modeling of long-
106 range dependencies recurrent networks [24, 40, 42, 50, 52, 5] to photorealistic video prediction using
107 large convolutional neural networks [37, 54, 39]. Time-agnostic prediction [21], which enables model
108 to predict any future frames in a video, is also proposed. In addition, methods based variational
109 autoencoders (VAEs) [27, 2, 15, 13, 53] are introduced to tackle the challenges of uncertainty in video
110 prediction. Our approach uses VAE-based video prediction model [13, 53] to combine multi-modal
111 data as inputs and predict the next frame tactile images.

112 3 TOUCH and GO Dataset

113 We collect a real-world vision-tactile dataset that contains egocentric videos of human (the authors)
114 pressing environments using a tactile sensor, Gelsight, and the tactile information from the Gelsight
115 that is simultaneously recorded with the RGB video. The touch of the environment contains useful
116 information about an object associated with the visual information, including hardness, shape,



Figure 2: TOUCH and GO Dataset. What do these objects feel like when they are touched? Here, we show some images from a selection of videos from our dataset for a subset of the object instances.

Table 1: Comparison of touch datasets.

	Hours	Touches	Object Inst.	Real-World	Indoor	Outdoor
More Than a Feeling [6]	-	6450	65	✓	✓	✗
VisGel [33]	20-30	12000	195	✓	✓	✗
The Feeling of Success [7]	-	9269	106	✓	✓	✗
Object Folder [16]	-	-	100	✗	-	-
Object Folder 2.0 [17]	-	-	1000	✗	-	-
TOUCH and GO (Ours)	>30	>37800*	>12600*	✓	✓	✓

117 material etc, which can be useful in various downstream tasks. Unlike traditional scene-centric
 118 datasets focusing on the full scene, our dataset is taken from an egocentric viewpoint which contains
 119 enough details to observe the fine-grained texture of an object.

120 3.1 Dataset Description

121 We collect over 30 hours of videos consisting of over 37800 touches from over 12600 objects under
 122 both indoor (58%) and outdoor (42%) scenes. In total, there are 1.89M frames containing touches
 123 of an object and each touch is composed of 50 frames on average. Our dataset contains daily seen
 124 objects, both hard and deformable, from indoor and outdoor scenes including rock, grass, road, brick,
 125 carpet, chair, table and so on. All the touch frames are annotated with the name of the object.

126 3.2 Comparison with other datasets

127 We compare TOUCH and GO with existing multisensory datasets in Table 1. Compared to the largest
 128 real-world dataset collected by robot, VisGel [33], our dataset comprises of longer hours of video,
 129 more touches, and most importantly much more diverse of object instances, where the total object
 130 instances is 65 times larger. It is worth noting that VisGel [33] and other robot collected datasets only
 131 contain indoor scene and the background is mostly fixed to the robotic operating station, which is far

132 from the place where the object actually exists in the real world. Our dataset is collected by human
133 from exact the real world where each object is recorded under the natural environment. With respect
134 to the largest object-centric multisensory dataset Object Folder 2.0 [17], our dataset contains object
135 instances 10 times larger and all tactile inputs are completely recorded by touching the actual object,
136 which provides more realistic tactile data compared to the synthetic or simulated images.

137 3.3 Data Collection Setup

138 As shown in Figure 2, we utilize a webcam to record the RGB video and a GelSight sensor to capture
139 the tactile signals, which are both connected to our PC. We record the timestamp of each frame to
140 synchronize visual and tactile images. GelSight sensor [23, 22, 59, 7] is an optical tactile sensor that
141 enables high spatial resolution measurement of the texture and geometry of a contact surface. The
142 sensor consists of a 1.5cm \times 1.5cm surface of a soft elastomer painted with a reflective membrane,
143 which deforms to the shape of the object upon contact. There exists an ordinary camera beneath the
144 elastomer so that we can view the deformed gel. The gels are illuminated by colored LEDs from
145 different directions, producing a three-channel surface normal image. Thus, we can observe the
146 texture of a surface undergoing the deformation process via consecutive 2D images. We can then
147 treat the tactile images as normal 2D images, and pass them to visual backbone network to extract
148 tactile information.

149 3.4 Detecting Touch Onset

150 According to our dataset, we have approximately 1/3 of the video frames that the GelSight sensor
151 is not touching the object. This is because our dataset is collected by human moving around and
152 touching objects seen during the movement. Thus the GelSight sensor has no deformation during the
153 interval when human is moving from one object to another. However, at the mean time, the RGB
154 camera still records the scene during the video. Under this circumstance, the scenes captured by RGB
155 camera will be incorrectly linked to the tactile signal of no deformation, which will negatively impact
156 the downstream tasks. Thus, to alleviate this issue, we train a binary classifier to classify whether the
157 frame is at touch onset. We hand label 10,000 frames from the dataset and finetune ResNet-18 [19]
158 initialized by weights pretrained on the ImageNet [12] on our dataset as our classifier. We report a
159 97% accuracy on our test set that is 20% of our labeled frames.

160 4 Applications

161 4.1 Tactile-Guided Image Stylization

162 The sense of touch conveys useful information about an object, including hardness, shape, material
163 etc., which creates an inherent association with the visual input in the video. This connection between
164 visual and tactile signals is embedded in our dataset and the neural network is able to build the
165 tactile-visual style associations without human supervision. Moreover, tactile signals may provide us
166 subtle distinction about objects that visual input can not capture. As shown in Figure 3, even when
167 objects of different materials share similar visual appearance, the tactile signals are able to reveal
168 their subtle difference. Given the unique properties of tactile signals and its association with visual
169 input, we propose *Tactile-Guided Image Stylization* (TGIS) application on our dataset, which, to the
170 best of our knowledge, is firstly considered in the current literature.

171 4.1.1 Proposed Method

172 Given an source domain $\mathcal{X} \in \mathbf{R}^{\mathbf{H} \times \mathbf{W} \times \mathbf{C}}$, our goal in TGIS is to learn the translation from \mathcal{X} to look
173 visually similar to an image from the target domain \mathcal{Y} that is corresponded to the tactile domain \mathcal{T} .
174 During the training time, we randomly sample two visual images from \mathcal{X} , \mathcal{Y} and a tactile image
175 from \mathcal{T} corresponding to the target image \mathcal{Y} . It is worth noting that our training requires no human
176 annotation and it can be done under self-supervision.

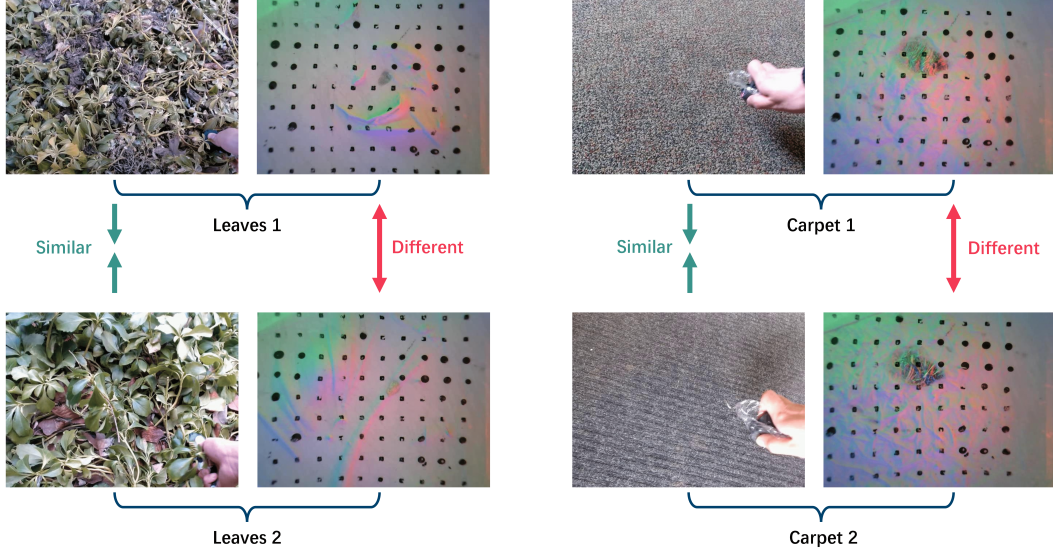


Figure 3: Subtle distinction from gelsight. Although it is hard to distinguish some object instances from visual appearance, tactile signals may convey enough information.

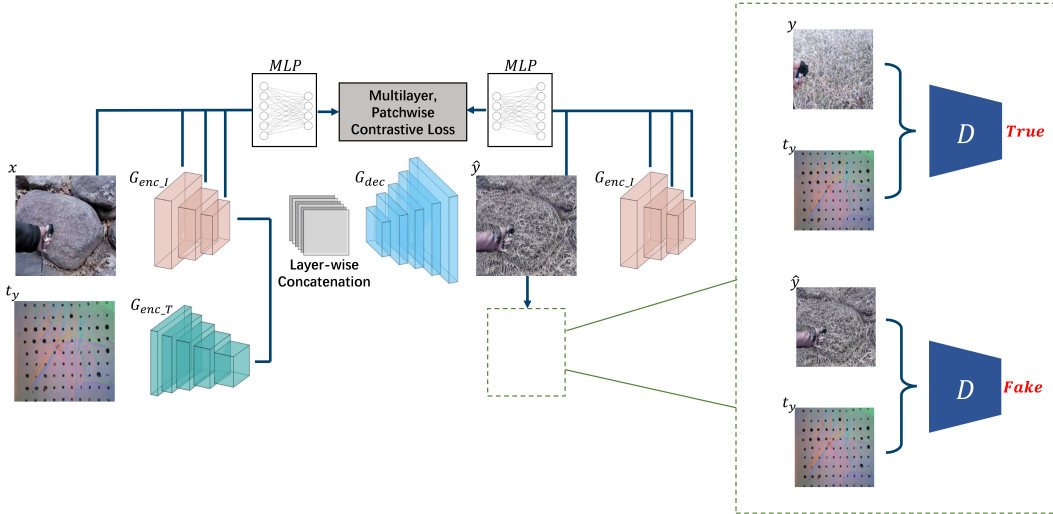


Figure 4: Pipeline of tactile-guided image stylization.

177 As shown in the Figure 4 about our pipeline, our model consists of a multi-modal generator, a tactile-
 178 visual texture discriminator, and a patch-wise structure discriminator. We can further break up our
 179 multi-modal generator into three components, an image encoder G_{enc_I} , a tactile encoder G_{enc_T} , and
 180 a decoder G_{dec} . Given our dataset that contains unpaired instances $X = \{\mathbf{x} \in \mathcal{X}\}$, $Y = \{\mathbf{y} \in \mathcal{Y}\}$
 181 and tactile input $T_y = \{\mathbf{t}_y \in \mathcal{T}\}$, the output image $\hat{\mathbf{y}}$ can be expressed as $\hat{\mathbf{y}} = G(\mathbf{x}, \mathbf{t}_y) =$
 182 $G_{dec}(\text{concat}(G_{enc_I}(\mathbf{x}), G_{enc_T}(\mathbf{t}_y)))$.

183 **Tactile-Visual Adversarial Loss** To leverage the association between visual input and tactile input,
 184 we propose a tactile-visual adversarial loss between $\hat{\mathbf{y}}$ and \mathbf{t}_y . In formal terms:

$$\mathcal{L}_{GAN} D(G_{X \rightarrow Y}, D_Y) = \mathbb{E}_{\mathbf{y} \sim Y} \log D(\mathbf{y}, \mathbf{t}_y) + \mathbb{E}_{\mathbf{x} \sim X} \log(1 - D(G(\mathbf{x}, \mathbf{t}_y), \mathbf{t}_y)) \quad (1)$$

185 where D is the discriminator. For the discriminator D , we adopt the early fusion where we first directly
 186 concatenate the generated image $\hat{\mathbf{y}}$ with the tactile input \mathbf{t}_y and then feed into the discriminator D .

187 **Structure Preservation via Contrastive Learning** Our goal in this tactile-guided image stylization
188 is to restyle the source image with the textures that are associated with the target tactile input while
189 preserving the source structure. However, structure and texture of an image are often entangled
190 with each other. With only tactile-visual adversarial loss, it becomes a trivial solution to completely
191 transfer the source image to the target domain without preserving the original structure. Thus, we
192 introduce an additional constraint called noise contrastive estimation (NCE) [44] to preserve the
193 structural information between the visual input x and the generated image \hat{y} .

194 4.2 Multi-modal Video Prediction

195 This section is still in progress. We are in the process of conducting more experiments and ablation
196 studies. Explanations and results will be completed very soon.

197 5 Experiments

198 5.1 Tactile-Guided Image Stylization

199 5.1.1 Experimental Setup

200 **Implementation Details** Our image encoder and decoder of the generator are fully convolutional
201 neural network consisting of 9 blocks of ResNet-based CNN bottlenecks. The first convolution layer is
202 set to 7×7 and the rest are set to 3×3 . For the tactile encoder, we adopt a ResNet-18 [19] backbone
203 pretrained on the ImageNet [12]. For the discriminator we adopt the PatchGAN architecture [20].
204 To compute the NCE loss, we extract features from five different layers: the input image layer, the
205 first and second downsampling convolution layer and the first and fifth residual blocks. We set the
206 hyperparameter λ and μ equal to 0.5. We train our model on 4 Nvidia 2080-Ti GPUs for 100 epochs
207 with the batch size of 8 and the learning rate of 0.0002. For input visual images, we employ a random
208 crop and an horizontal flip.

209 5.1.2 Results

210 We show the qualitative results in the Figure 5. All of the results are generated from the single model
211 (i.e., by one-to-many relation). With input of tactile signals, our model is capable to distinguish and
212 capture the subtle distinction between the input category and the output category without any label.

213 5.2 Multi-modal Video Prediction

214 This section is still in progress. We are in the process of conducting more experiments and ablation
215 studies. Explanations and results will be completed very soon.

216 6 Conclusion

217 We introduce TOUCH and GO, a multisensory dataset containing real-world synchronized high-
218 quality video and tactile data captured from egocentric viewpoint. Compared to existing real-world
219 multisensory datasets, our work contains much greater hours of videos, object instances, and touches.
220 We propose two novel applications including tactile-guided image stylization and multi-modal video
221 prediction. Leveraging TOUCH and GO dataset, experimental results indicate our designed models
222 outperform label-based counterparts in both quantitative and qualitative evaluations. We hope our
223 dataset, which is easily accessible to the community, will drive more multisensory applications and
224 serve as a standard benchmark.

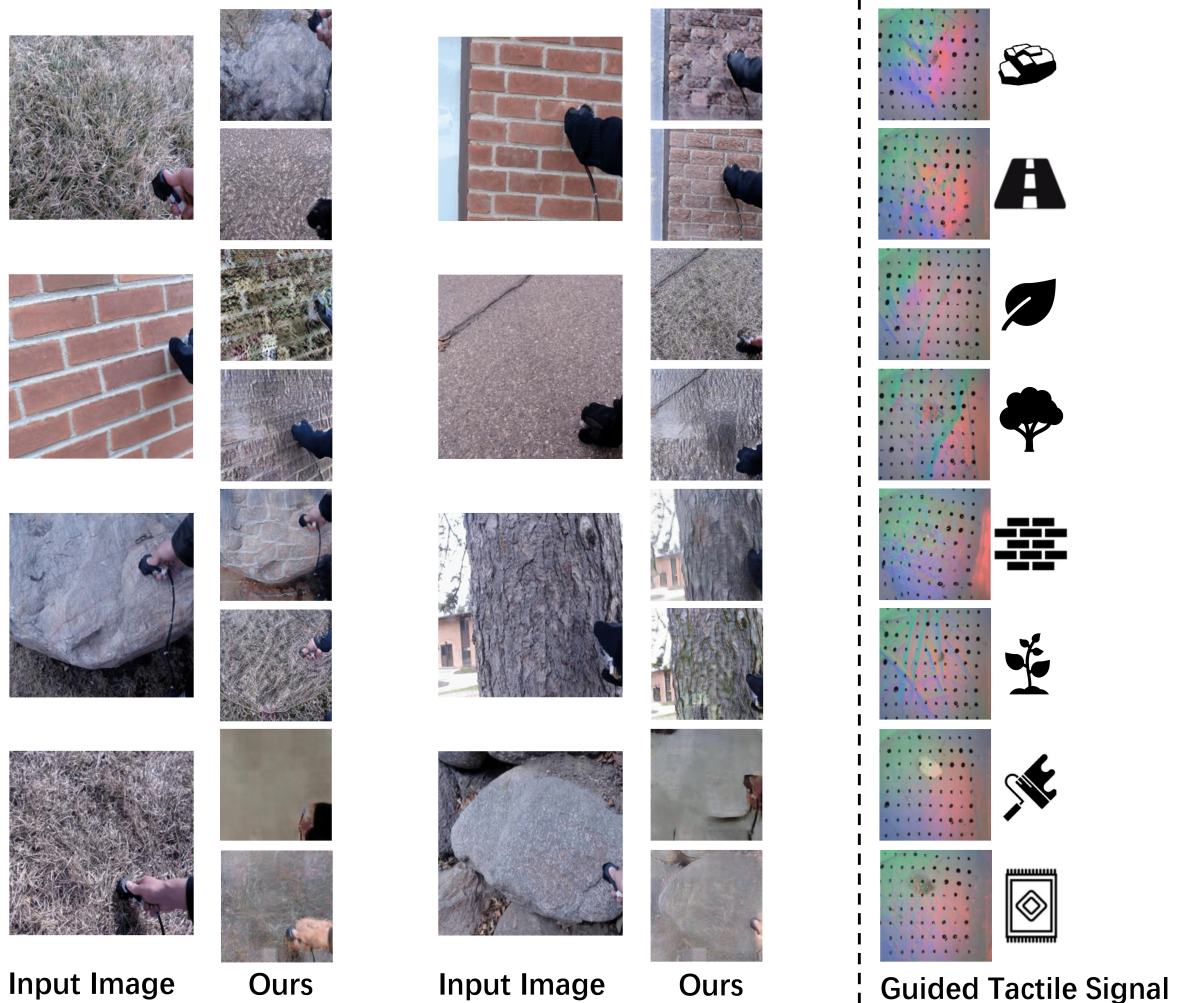


Figure 5: Qualitative results our model on tactile-guided image stylization. For reference, we show guided tactile signals as well as their corresponding images in the last column.

References

- 225
- 226 [1] *ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models*,
 227 Vancouver, Canada, 11/2019 2019.
- 228 [2] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. Stochastic
 229 variational video prediction. In *International Conference on Learning Representations*, 2018.
- 230 [3] Renée Baillargeon. *The Acquisition of Physical Knowledge in Infancy: A Summary in Eight Lessons*, pages
 231 46–83. 11 2002.
- 232 [4] Konstantinos Bousmalis and Sergey Levine. Closing the simulation-to-reality gap for deep robotic learning,
 233 2017.
- 234 [5] Wonmin Byeon, Qin Wang, Rupesh Kumar Srivastava, and Petros Koumoutsakos. Contextvp: Fully
 235 context-aware video prediction. In *ECCV*, 2018.
- 236 [6] Roberto Calandra, Andrew Owens, Dinesh Jayaraman, Justin Lin, Wenzhen Yuan, Jitendra Malik, Ed-
 237 ward H. Adelson, and Sergey Levine. More than a feeling: Learning to grasp and regrasp using vision and
 238 touch. *IEEE Robotics and Automation Letters*, 3:3300–3307, 2018.
- 239 [7] Roberto Calandra, Andrew Owens, Manu Upadhyaya, Wenzhen Yuan, Justin Lin, Edward Adelson, and
 240 Sergey Levine. The feeling of success: Does touch sensing help predict grasp outcomes? 10 2017.

- 241 [8] Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M. Dollar. The
242 ycb object and model set: Towards common benchmarks for manipulation research. In *2015 International
243 Conference on Advanced Robotics (ICAR)*, pages 510–517, 2015.
- 244 [9] Angel Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio
245 Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An
246 information-rich 3d model repository. 12 2015.
- 247 [10] Jasmine Collins, Shubham Goel, Achleshwar Luthra, Leon Xu, Kenan Deng, Xi Zhang, Tomas F Yago Vi-
248 cente, Himanshu Arora, Thomas Dideriksen, Matthieu Guillaumin, and Jitendra Malik. Abo: Dataset and
249 benchmarks for real-world 3d object understanding. *arXiv preprint arXiv:2110.06199*, 2021.
- 250 [11] Mark R. Cutkosky, Robert D. Howe, and William R. Provancher. *Force and Tactile Sensors*, pages 455–476.
251 Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.
- 252 [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
253 image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255,
254 2009.
- 255 [13] Emily L. Denton and Rob Fergus. Stochastic video generation with a learned prior. In *ICML*, 2018.
- 256 [14] Zihan Ding, Nathan Lepora, and Edward Johns. Sim-to-real transfer for optical tactile sensing. pages
257 1639–1645, 05 2020.
- 258 [15] Piotr Dollar, Christian Wojek, Bernt Schiele, and Pietro Perona. Pedestrian detection: A benchmark. In
259 *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 304–311, 2009.
- 260 [16] Ruohan Gao, Yen-Yu Chang, Shivani Mall, Li Fei-Fei, and Jiajun Wu. Objectfolder: A dataset of objects
261 with implicit visual, auditory, and tactile representations. In *CoRL*, 2021.
- 262 [17] Ruohan Gao, Zilin Si, Yen-Yu Chang, Samuel Clarke, Jeannette Bohg, Li Fei-Fei, Wenzhen Yuan, and
263 Jiajun Wu. Objectfolder 2.0: A multisensory object dataset for sim2real transfer. In *CVPR*, 2022.
- 264 [18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron
265 Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes,
266 N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, vol-
267 ume 27. Curran Associates, Inc., 2014.
- 268 [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
269 *arXiv preprint arXiv:1512.03385*, 2015.
- 270 [20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional
271 adversarial networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
272 pages 5967–5976, 2017.
- 273 [21] Dinesh Jayaraman, Frederik Ebert, Alexei A. Efros, and Sergey Levine. Time-agnostic prediction:
274 Predicting predictable video frames. *ArXiv*, abs/1808.07784, 2019.
- 275 [22] Micah Johnson, Forrester Cole, Alvin Raj, and Edward Adelson. Microgeometry capture using an
276 elastomeric sensor. *ACM Trans. Graph.*, 30:46, 07 2011.
- 277 [23] Micah K. Johnson and Edward H. Adelson. Retrographic sensing for the measurement of surface texture
278 and shape. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1070–1077,
279 2009.
- 280 [24] Nal Kalchbrenner, Aäron van den Oord, Karen Simonyan, Ivo Danihelka, Oriol Vinyals, Alex Graves, and
281 Koray Kavukcuoglu. Video pixel networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of
282 the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning
283 Research*, pages 1771–1779. PMLR, 06–11 Aug 2017.
- 284 [25] Zhanat Kappasov, Juan-Antonio Corrales, and Véronique Perdereau. Tactile sensing in dexterous robot
285 hands — review. *Robotics and Autonomous Systems*, 74:195–220, 2015.
- 286 [26] Taeksoo Kim, Moon-su Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. Learning to discover
287 cross-domain relations with generative adversarial networks. In *ICML*, 2017.
- 288 [27] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2014.

- 289 [28] Raj Kolamuri, Zilin Si, Yufan Zhang, Arpit Agarwal, and Wenzhen Yuan. Improving grasp stability with
290 rotation measurement from tactile sensing. In *Proceedings of (IROS) IEEE/RSJ International Conference*
291 *on Intelligent Robots and Systems*, October 2021.
- 292 [29] Oliver Kroemer, Scott Niekum, and George Dimitri Konidaris. A review of robot learning for manipulation:
293 Challenges, representations, and algorithms. *J. Mach. Learn. Res.*, 22:30:1–30:82, 2021.
- 294 [30] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab
295 Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open
296 images dataset v4: Unified image classification, object detection, and visual relationship detection at scale.
297 *IJCV*, 2020.
- 298 [31] S Lederman and Roberta Klatzky. Haptic perception: A tutorial. volume 71, pages 1439–59, 10 2009.
- 299 [32] Susan J Lederman and Roberta L Klatzky. Hand movements: A window into haptic object recognition.
300 *Cognitive Psychology*, 19(3):342–368, 1987.
- 301 [33] Yunzhu Li, Jun-Yan Zhu, Russ Tedrake, and Antonio Torralba. Connecting touch and vision via cross-
302 modal prediction. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June
303 2019.
- 304 [34] Joseph J. Lim, Hamed Pirsiavash, and Antonio Torralba. Parsing ikea objects: Fine pose estimation. In
305 *2013 IEEE International Conference on Computer Vision*, pages 2992–2999, 2013.
- 306 [35] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
307 and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla,
308 Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014.
309 Springer International Publishing.
- 310 [36] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In
311 *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*,
312 page 700–708, Red Hook, NY, USA, 2017. Curran Associates Inc.
- 313 [37] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction
314 and unsupervised learning. 05 2016.
- 315 [38] Shan Luo, Wenzhen Yuan, Edward Adelson, Anthony Cohn, and Raul Fuentes. Cloth texture recognition
316 using vision and tactile sensing. 05 2018.
- 317 [39] Michaël Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean
318 square error. *CoRR*, abs/1511.05440, 2016.
- 319 [40] Vincent Michalski, Roland Memisevic, and Kishore Konda. Modeling deep temporal dependencies with
320 recurrent grammar cells". In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger,
321 editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- 322 [41] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets, 2014.
- 323 [42] Roni Mittelman, Benjamin Kuipers, Silvio Savarese, and Honglak Lee. Structured recurrent temporal
324 restricted boltzmann machines. In *ICML*, pages 1647–1655, 2014.
- 325 [43] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H. Adelson, and William T.
326 Freeman. Visually indicated sounds. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*
327 *(CVPR)*, pages 2405–2413, 2016.
- 328 [44] Taesung Park, Alexei A. Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for conditional
329 image synthesis. In *ECCV*, 2020.
- 330 [45] Patsorn Sangkloy, Jingwan Lu, Chen Fang, Fisher Yu, and James Hays. Scribbler: Controlling deep image
331 synthesis with sketch and color. *2017 IEEE Conference on Computer Vision and Pattern Recognition*
332 *(CVPR)*, pages 6836–6845, 2017.
- 333 [46] Laura Schulz. The origins of inquiry: Inductive inference and exploration in early childhood. *Trends in*
334 *cognitive sciences*, 16:382–9, 06 2012.
- 335 [47] Arjun Singh, James Sha, Karthik S. Narayan, Tudor Achim, and P. Abbeel. Bigbird: A large-scale 3d
336 database of object instances. *2014 IEEE International Conference on Robotics and Automation (ICRA)*,
337 pages 509–516, 2014.

- 338 [48] Linda B. Smith and Michael Gasser. The development of embodied cognition: Six lessons from babies.
339 *Artificial Life*, 11:13–29, 2005.
- 340 [49] Elizabeth S. Spelke and Katherine D. Kinzler. Core knowledge. *Developmental science*, 10 1:89–96, 2007.
- 341 [50] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. Unsupervised learning of video represen-
342 tations using LSTMs. In *ICML*, 2015.
- 343 [51] Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B
344 Tenenbaum, and William T Freeman. Pix3d: Dataset and methods for single-image 3d shape modeling. In
345 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- 346 [52] Ilya Sutskever, Geoffrey E. Hinton, and Graham W. Taylor. The recurrent temporal restricted boltzmann
347 machine. In *NIPS*, 2008.
- 348 [53] Ruben Villegas, Arkanath Pathak, Harini Kannan, D. Erhan, Quoc V. Le, and Honglak Lee. High fidelity
349 video prediction with large stochastic recurrent neural networks. In *NeurIPS*, 2019.
- 350 [54] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro.
351 Video-to-video synthesis. *NIPS’18*, page 1152–1164, Red Hook, NY, USA, 2018. Curran Associates Inc.
- 352 [55] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-
353 resolution image synthesis and semantic manipulation with conditional gans. pages 8798–8807, 06
354 2018.
- 355 [56] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao.
356 3d shapenets: A deep representation for volumetric shapes. pages 1912–1920, 06 2015.
- 357 [57] Yu Xiang, Wonhui Kim, Wei Chen, Jingwei Ji, Chris Choy, Hao Su, Roozbeh Mottaghi, Leonidas Guibas,
358 and Silvio Savarese. Objectnet3d: A large scale database for 3d object recognition. volume 9912, pages
359 160–176, 10 2016.
- 360 [58] Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dualgan: Unsupervised dual learning for image-to-image
361 translation. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2868–2876, 2017.
- 362 [59] Wenzhen Yuan, Siyuan Dong, and Edward Adelson. Gelsight: High-resolution robot tactile sensors for
363 estimating geometry and force. *Sensors*, 17:2762, 11 2017.
- 364 [60] Wenzhen Yuan, Shaoxiong Wang, Siyuan Dong, and Edward Adelson. Connecting look and feel: Associat-
365 ing the visual and tactile properties of physical materials. 04 2017.
- 366 [61] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei Efros. Generative visual manipulation on
367 the natural image manifold. volume 9909, 08 2016.
- 368 [62] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei Efros. Unpaired image-to-image translation using
369 cycle-consistent adversarial networks. pages 2242–2251, 10 2017.
- 370 [63] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shecht-
371 man. Toward multimodal image-to-image translation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach,
372 R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*
373 30, pages 465–476. Curran Associates, Inc., 2017.