

TACTO: A Fast, Flexible, and Open-source Simulator for High-Resolution Vision-based Tactile Sensors

Shaoxiong Wang^{1,*}, Mike Lambeta², Po-Wei Chou², and Roberto Calandra²

Abstract—Simulators perform an important role in prototyping, debugging, and benchmarking new advances in robotics and learning for control. Although many physics engines exist, some aspects of the real world are harder than others to simulate. One of the aspects that have so far eluded accurate simulation is touch sensing. To address this gap, we present TACTO – a fast, flexible, and open-source simulator for vision-based tactile sensors. This simulator allows to render realistic high-resolution touch readings at hundreds of frames per second, and can be easily configured to simulate different vision-based tactile sensors, including DIGIT and OmniTact. In this paper, we detail the principles that drove the implementation of TACTO and how they are reflected in its architecture. We demonstrate TACTO on a perceptual task, by learning to predict grasp stability using touch from 1 million grasps, and on a marble manipulation control task. Moreover, we provide a proof-of-concept that TACTO can be successfully used for Sim2Real applications. We believe that TACTO is a step towards the widespread adoption of touch sensing in robotic applications, and to enable machine learning practitioners interested in multi-modal learning and control. TACTO is open-source at <https://github.com/facebookresearch/tacto>.

Index Terms—Simulation and Animation; Perception for Grasping and Manipulation; Force and Tactile Sensing; Learning and Adaptive Systems; Deep Learning in Robotics and Automation

I. INTRODUCTION

SIMULATORS play an important role in prototyping, debugging and benchmarking new advances in robotics. With an appropriate simulator, expensive and time-consuming experiments in the real world can be approximated inexpensively on our computers. This allows to perform orders of magnitude more experiments at a fraction of the effort, and in many cases of the time. The robotics community has traditionally made extensive use of simulators for control, and many different physics engines are available to researchers and practitioners [1]. One aspect that has proven so far to be difficult to simulate is tactile sensing, and in particular vision-based tactile sensors [2], [3], [4] which provide rich high-resolution measurements. This is because to accurately model this family of tactile sensors it is necessary not only to model

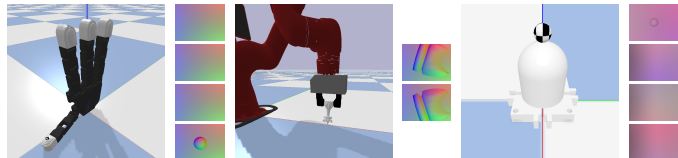


Fig. 1: We open-source TACTO – a simulator of vision-based tactile sensors. TACTO produces high-resolution and high-fidelity reading from tactile sensors at high-frequency (>100 Hz). Its modular structure allows to model different vision-based tactile sensors and to be integrated with different physics engines. We believe that such a tool can benefit the touch sensing and robotic community, as well as researchers in machine learning that can now access a new sensor modality.

the dynamics of the contact, but also to model the optical properties of the sensors and the corresponding illumination to obtain realistic perceptual outputs. All of this while keeping the simulator flexible enough to implement various sensors with different form factors, and fast enough to be of practical use.

To fill this lack of touch sensing simulators, we introduce TACTO – a simulator of vision-based tactile sensors explicitly designed to be fast and flexible. Our main contribution is to develop and open-source this simulator of vision-based tactile sensors. In this paper, we describe the design choices adopted, the resulting software architecture of the simulator, and discuss some of its most important features. Following, we present simulated experiments to demonstrate the capabilities of TACTO on perception and control tasks. Fig. 1 shows examples of TACTO in different scenarios.

TACTO can natively be used in conjunction with PyBullet [5], but can also be interfaced with other physics engines. Particular care was dedicated to having a simulator that could be both fast and flexible. Through our design, TACTO can render hundreds of frames per second, thus making the simulator practical for many control and learning to control applications. In addition, thanks to its modularity it is easy to implement vision-based tactile sensors with different form factors and lighting properties. Currently, TACTO implements two recent vision-based tactile sensors: OmniTact [3] and DIGIT [4]. While the ideal touch simulator would provide both realistic perceptual outputs and accurate contact dynamics, TACTO is aimed at tackling the rendering of realistic perceptual outputs and the creation of accurate contact dynamics is currently left to the underlying physical engine being used by the user. We demonstrate our simulator by learning grasp stability models from touch, and by learning in-hand marble manipulation. For learning grasp stability from touch, we collected a simulated

Manuscript received: September, 9, 2021; Revised December, 16, 2021; Accepted January, 9, 2022.

This paper was recommended for publication by Editor Dan Popa upon evaluation of the Associate Editor and Reviewers' comments.

* Work done during an internship at Meta AI.

¹ Massachusetts Institute of Technology

² Meta AI, Menlo Park, CA, USA

Corresponding author: Roberto Calandra (rcalandra@fb.com)

Digital Object Identifier (DOI): 10.1109/LRA.2022.3146945

dataset of 1 million grasps, which is orders of magnitude more data than the largest dataset previously used for this task [6]. This highlights the advantages of having a fast simulator of vision-based tactile sensors and allows us to better understand the performance of machine learning grasp stability models trained from large datasets that have, so far, been infeasible to collect in the real world. Moreover, we apply Bayesian optimization to learn in simulation control to manipulate marbles between two fingers from touch.

Finally, we present a proof-of-concept of Sim2Real on a pose-estimation task to show the comparison between simulated and real signals, and various ways to improve the performance. With the Sim2Real gap in mind, our major goal is to provide a playground, similar to OpenAI Gym [7], to test different design or learning algorithms for the touch modality. The model can be later trained with real data. But some fundamental understanding and experience learned in simulation can be helpful for real robots, e.g. exploring the possible representation/policy to combine vision and touch for robotic tasks.

We believe that TACTO can be of practical value for different communities: 1) To hardware designers, it provides a preliminary method to simulate and evaluate their design choices for future sensors; 2) To the robotic community, it provides a way to simulate and study the integration of touch sensing into control scheme; 3) To the machine learning community, it provides an easy-to-use tool to generate multi-modal inputs which would otherwise require real-world hardware. To enable and stimulate researchers and practitioners in these fields, we open-source TACTO at <https://github.com/facebookresearch/tacto>.

II. RELATED WORK

While there is a large number of physics engines available to robotic practitioners [1], the choice when simulating tactile sensors is more limited. This is mostly due to the difficulty of accurately and efficiently simulating touch.

Several traditional low-dimensional sensors have been simulated in the literature, including BioTac [8] and fabric-based tactile sensors [9]. [10] simulated the iCub’s iSkin in Gazebo through the use of an array of tactels appropriately distributed on the kinematic structure of the robot. The speed of the simulation was however severely impacted by modeling such numbers of sensors independently. Similarly, [11] made use of Gazebo to model general-purpose robotic skin, but modeled each element of the array of tactels as a spring-mass-damper system to provide an improved characterization of the mechanical properties of the skin. In contrast to these works that simulate low-dimensional tactile sensors, we aim to simulate high-resolution vision-based tactile sensors, which can possess millions of tactels. This creates novel challenges and requires different modeling approaches.

Exploiting the nature of vision-based tactile sensors, it is possible to use ray-tracing models from computer graphics to render sensor output. There are emergent simulators that shows promise for vision-based tactile sensors based on Phong’s model [12], [13], Mitsuba2[14], Unity [15], optical flow [16], and finite elements models [17]. Among them, the closest

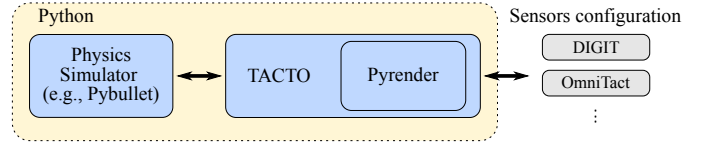


Fig. 2: Software Architecture. TACTO bridges between physics simulator and back-end rendering engine, and can be configured to model different sensor designs through configuration files.

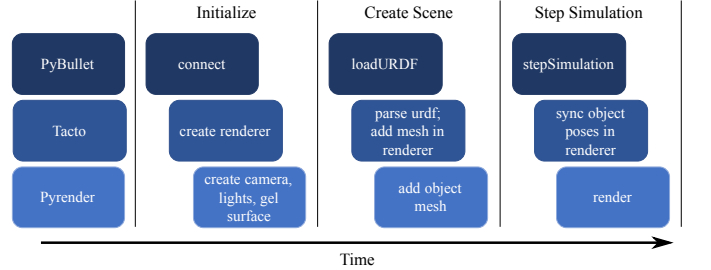


Fig. 3: Workflow showing the functionality of TACTO at three major phases. (1) Initialize: create the sensor structure in the renderer; (2) Create scene: parse the objects URDF and add them into the renderer; (3) Step simulation: synchronize the object poses from physics engine to the renderer.

concurrent works are [14], [13]. For comparison, the authors in [14] focused on simulating realistic tactile imprints using Mitsuba2 renderer[18], but have not worked on efficiently integrating it with physical simulators so far. The authors in [13] combined the Phong’s model with the Gazebo simulator, but the experiments were mostly performed with a tactile sensor pressing on fixed objects. In contrast, we focus on exploring the scenarios where the sensors actively interact with the objects during the grasping or manipulation. Overall, we aim to provide an open-source simulator that is fast, flexible, and can be efficiently integrated with the physics engine for experimenting with different learning and control algorithms with touch modality.

III. A FAST AND FLEXIBLE SIMULATOR OF VISION-BASED TACTILE SENSORS

We now detail the desired design considered when designing our simulators, and the corresponding architectural choices made to achieve these desiderata. Following, we discuss some of the salient features that TACTO offers.

A. Design Desiderata

High-throughput: simulators must be as fast as possible to reduce the real-world time of running simulations. Reproducing touch sensing, even for low-dimensional sensors, has traditionally been a very computationally intensive operation often leading to simulations barely faster than real-time [10]. Obtaining a simulator that could perform hundreds of frames per second was, from the beginning, one of our most important design desiderata.

Flexible: since there are different sensor designs of vision-based tactile sensors, where some of them have complex geometry [3], mirrors [19], and transparent case for light piping [20], it is desirable for any simulator to be flexible and powerful enough to support a large choice of optical components, and mechanical designs.

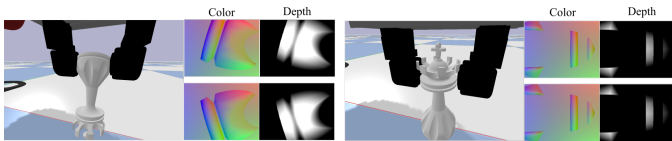


Fig. 4: Example images of simulated DIGIT imprints. TACTO is able to generate color and depth images at the same time with details of the local geometry at high speed.

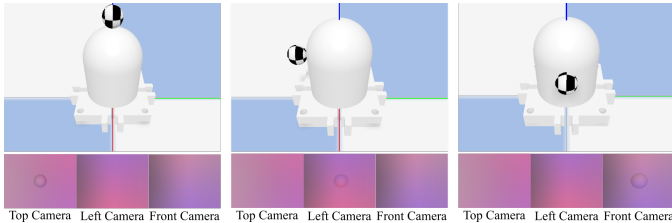


Fig. 5: Example images of a simulated OmniTact [3] touching a sphere. We show only 3 of the 5 cameras mounted on the sensor. It is visible how the specific light/camera placement in the OmniTact results in more pinkish images also in our simulator.

Realistic: it is desirable that the simulator produces outputs as close as possible to real measurements, including illumination of the contact region, global lighting distribution, and details like shadows and deformation on the contact boundary.

Easy to Use: finally, it is from a practical point of view very important that the simulator is easy to install, use, modify and set up for different perception and control tasks.

B. Architectural Choices

Here we compare several architectural choices. The simulator needs to calculate the local contact geometry (Depth), and the corresponding rendering (RGB) that best meets the desiderata.

1. Phong’s model for RGB rendering from Depth (simple but less powerful): PyBullet built-in camera can provide a depth map of the contact area. To render the RGB image from the depth map, researchers from [13] implemented their own renderer based on Phong’s reflection model. It generated promising results and should be easy to use. However, it assumed that light only bounces once, directly from the gel surface to the camera. Hence, it is difficult to adapt to existing and future sensor designs that require advanced functionalities, like reflection, refraction, and shadows with fast speed.

Although it can be extended to support more ray-tracing functionalities, this may require non-trivial engineering time to re-implement methods with GPU acceleration which are already provided and tested by open-source graphics libraries, like OpenGL[21].

2. OpenGL for RGB rendering from Depth (powerful but slow): Alternatively, we can leverage the power of OpenGL [21]. Besides the features in [13], OpenGL also supports mirrors, transparent objects, shadow, GPU acceleration, etc, which opens up the possibility for rendering advanced sensor design at high speed, with minimal effort. To use the power of OpenGL, one can get the depth image from PyBullet first, and pass the depth map to OpenGL for rendering.

The rendering alone is fast and can be sped up in GPU, however, I/O speed becomes the bottleneck. In preliminary

experiments, a significant amount of time was spent on loading the mesh generated by the depth map into OpenGL, and this limited the overall speed of this method to only 20 frames per second, even on GPU.

3. OpenGL for RGB rendering from synchronized scenes (proposed, powerful and fast): Our proposed system design builds a synchronized scene from the physical simulator, and directly renders both depth and RGB images in OpenGL. It can achieve high speed with powerful rendering functionalities.

Specifically, to avoid the I/O bottleneck of loading the mesh from depthmaps repeatedly, TACTO preloads the gel surface and object meshes into the OpenGL scene, then keeps synchronizing their poses from the physical simulator. At each step, TACTO overlaps the original gel geometry with the contact object, with the respective poses, in OpenGL to extract depth and RGB images. Although loading meshes is slow, it is very fast to change their poses and re-render afterwards. In this way, the system can render at very high speed (up to 200 frames per second in our experiments).

Note that this is designed to speed up the computation for interacting the sensor with rigid objects, where the deformation of the object itself is negligible. Since the silicone gel is much softer than many everyday objects, the method can still be applied to these objects. To manipulate very deformable objects, TACTO provided the functionality to render RGB from depth, as described in the second option, with slower speed though.

One limitation of the method is that it is difficult to model the deformation of the gel on the contact boundary, because the RGB and depth images are calculated at the same time. However, this mostly affects sharp edges, and can be approximated by smoothing objects’ meshes during pre-processing. We can also add data augmentation [22] or refine the rendering afterward with generative models to make it more realistic [23], [24]. Besides, we expect that contact boundary contains relatively little useful information compared to other aspects like contact location, contact mask, object contact pose, normal forces, etc. It can still be a good playground to study different algorithms in simulation for touch. Overall, we think it is worthwhile to trade this for speed.

For implementation, we use Pyrender [25] as the renderer in TACTO since it provides a lightweight python interface for deploying OpenGL with GPU support, making TACTO not only powerful, but also easy to use.

Based on extensive experiments and analysis of each method, we decided to use this third option. This results in TACTO being closely aligned to our ideal simulator, which is fast, flexible, and powerful. In Section IV, we will demonstrate that TACTO is also easy to set up for perception and control tasks.

C. Overview of the Software Architecture

Fig. 2 shows the overall software architecture of TACTO. TACTO takes the responsibility of bridging between the physics simulator and the back-end rendering engine. Note that the architecture we propose leverages the advanced ray-tracing tools, which enables simulating high-quality tactile

	Res	1 obj (1 in contact)				100 objs (1 in contact)				100 objs (10 in contact)			
		Step	Sync	Render	FPS	Step	Sync	Render	FPS	Step	Sync	Render	FPS
GPU	160×120	0.2	2.2	2.5	220	5.8	2.4	2.8	200	11.6	6.3	6.9	80
	320×240	0.2	2.2	5.0	140	5.8	2.5	7.6	100	11.5	6.3	9.5	60
	640×480	0.2	2.3	9.4	90	6.1	2.6	8.6	90	11.7	6.3	15.0	50
CPU	160×120	0.4	4.2	11.8	60	9.8	4.5	11.3	60	19.5	10.7	16.7	40
	320×240	0.4	4.2	25.8	30	9.6	4.6	26.1	30	18.0	9.8	26.2	30
	640×480	0.4	4.5	78.3	10	9.0	4.8	78.5	10	17.2	10.4	80.1	10

Table I: Breakdown time of each stage (in millisecond) and overall speed (in frames per second, excluding physics simulation) for TACTO with PyBullet when simulating a DIGIT sensor with multiple objects (each mesh with 12K faces) in the scene. Step: PyBullet simulates physics; Sync: TACTO synchronizes the scene; Render: Pyrender renders the scene. Note that physics simulation can run asynchronously to speed up. CPU machine: Intel Core i7-6820HQ. GPU machine: Nvidia RTX 2080 Super GPU with Intel Core i9-9900K CPU.

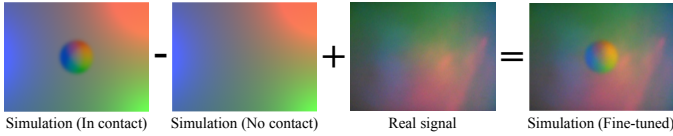


Fig. 6: If readings from a real-world sensor are available, TACTO allows to fine-tune the simulator using the real-world data. This is achieved by calculating the pixel-wise difference of the simulated images with and without touch, and then adding the reference real-world image.

signals efficiently, with minimal efforts. It can render different sensors by loading corresponding sensor configurations, which include parameters such as the camera, lights, and gel surface.

Fig. 3 shows a summary of the overall workflow. TACTO includes three major phases (for simplicity, we here assume PyBullet as the underlying physics engine). 1) Initialize: TACTO loads the sensor configuration, and setup up the sensor (camera, lights, gel mesh) in the rendering engine. 2) Create scene: PyBullet loads object URDF into the scene, and TACTO parses the URDF (by urdfpy package in our case), and adds the analyzed mesh into the rendering engine. 3) Step simulation: PyBullet first calculates physics simulation. Then TACTO loads the poses of each link from PyBullet, synchronizes the poses of objects and sensors in the rendering engine, and fetches the rendered tactile imprints.

D. Salient Features

Fast: TACTO is very fast for being a tactile sensor simulator. Table I shows the speed of TACTO. When interacting with an object mesh with 12K faces, TACTO is able to render a single DIGIT sensor at 200 frames per second on GPU with an output resolution of 160x120 pixels. It can also render multiple sensor outputs, e.g., rendering four DIGIT sensors mounted on an Allegro hand at 50 frames per second on GPU with an output resolution of 160x120 pixels. We optimized the TACTO so that only the number of objects in contact influences Pyrender’s speed, making it maintain a high speed even in a cluttered environment. The rendering time grows linearly with higher output resolution on GPU, and the speed stays the same when the mesh size scales from 2K to 12K faces.

Flexible: TACTO can adapt to different sensor designs by changing configuration files easily. Besides rendering DIGIT [4] signals, as shown in Fig. 4, we also demonstrate the possibility to render OmniTact [3] signals, as shown in Fig. 5, which has a substantially more complicated sensor structure including round surface, 5 cameras, and 11 light sources. For

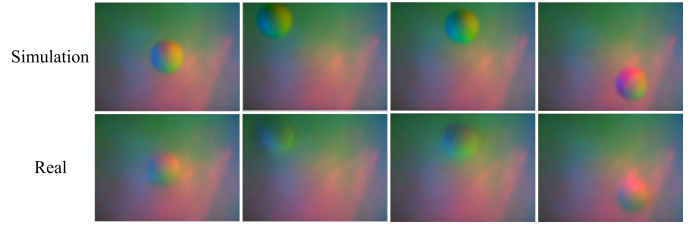


Fig. 7: Comparison of simulation and real signals with contacts across the sensor. TACTO captures the non-uniform light distribution similar to the real signals. The real-world readings are collected from a DIGIT sensor [4] touching a ball of 5.3 mm diameter. Simulated images are fine-tuned with a background real-sensor image.



Fig. 8: TACTO supports rendering shadows to obtain more realistic simulations. The real-world measurement is collected from a DIGIT sensor touching a ball of 3.7 mm diameter.

the configuration file details, it requires parameters of the gel (pose, mesh), a list of camera(s) (pose, field of view, clipping plane), a list of lights (pose, color, intensity), and user-specific details like noise level, and the mapping from force to deformation. The example configuration files in the repository include the explanations for each parameter. The configuration file and renderer can be further extended to support more functionality of OpenGL/Pyrender and sensors.

Force dependent: In TACTO, the contact forces are generated from the physics engine of choice (currently by default PyBullet with a rigid body contact model). However, within TACTO we also apply a deformation function to simulate the dynamics range of the deformation of the gel by mapping the force measured in the physics engine to different deformations of the mesh used to simulate the gel. This means that light forces will yield fewer deformations, and higher forces to more deformation of the gel. The current deformation function used is a piece-wise linear mapping from normal forces to deformation depth, which approximates the linear elasticity of real sensors within a reasonable range [26] and consider both a lower (below which the sensor is not capable of sensing deformation) and an upper threshold (above which the deformation is saturated) to the force sensed. However, it is straightforward for the end-user to provide their own deformation function, for example, by characterizing the gel

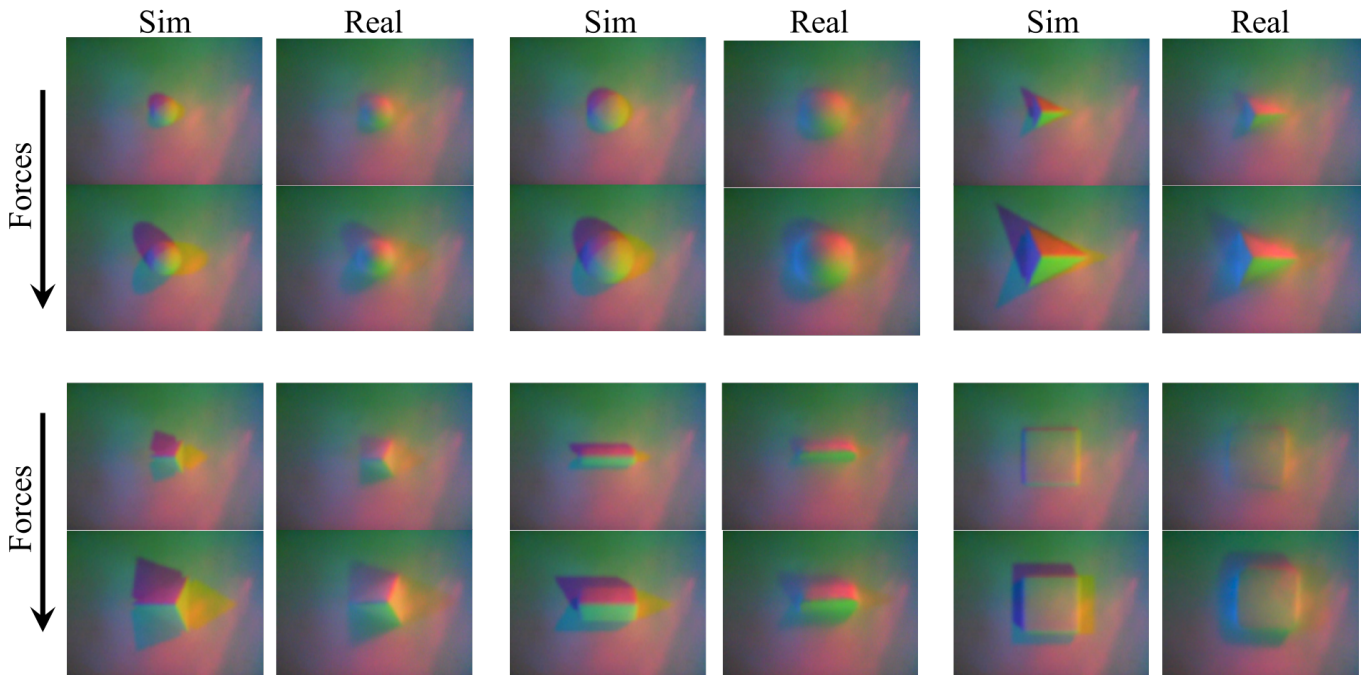


Fig. 9: Comparison of simulated and real tactile readings. Real readings are collected from a DIGIT sensor touching objects with different shapes (spheres, corners, edges, planes), with different forces. TACTO generates realistic tactile signals for both contact surfaces and shadows. The simulated images are overlaid on a background real sensor image. We smooth the cube’s mesh to better model the deformations on the contact boundaries, during pre-processing. A 15×15 Gaussian filter is applied on the rendered RGB image to match the imperfection from camera/light source in the real sensor. Note that it takes only 2 extra milliseconds for rendering shadows (with meshes of 19k faces, output resolution of 640x480, GPU).

with a non-linear function. In the future, we plan to provide more realistic transfer functions for the sensors to which we have access.

Rendering from depth: TACTO also supports rendering tactile imprints from depth images as an option. Given the depth image, it generates corresponding meshes to replace the original gel surface in the rendering engine. As mentioned in Section III-B, the speed is limited due to I/O bottleneck, however, it can be helpful in some cases where it requires modification on depth images before rendering for more realistic tactile imprints.

Calibration from real sensors: To make the rendering more realistic the simulator also supports a procedure, sketched in Fig. 6, that allows to fine-tune the rendering using readings collected from real-world sensors. This results in highly realistic renderings that can be easily customized to each sensor. In addition, TACTO captures the illumination changes across the sensor similar to real measurements. In real sensors, the light becomes dimmer when traveling longer, which generates non-uniform light distributions. Fig. 7 shows a comparison across different contact regions. Finally, TACTO supports rendering shadows to match real signals, as shown in Fig. 8. This option is easy to enable with our framework, and takes only ≈ 2 extra milliseconds to render on GPU. In contrast, it can take non-trivial time to re-implement shadow rendering from scratch with Phong’s reflection model [13], especially to optimize for GPU acceleration. Fig. 9 shows a comparison between simulated and real imprints. TACTO can simulate realistic rendering for both contact surfaces and shadows, with different contact geometry and forces.

Compatibility to different physics engines: We use an open-source physics engine, PyBullet, to demonstrate the framework. But TACTO can also support other physics engines. Specifically, there are two ways: 1) is to synchronize the scene with the functions provided by the selected physics engine. The required functions include getting object/link poses to synchronize the scene, and getting contact forces to simulate deformation. 2) is to render from depth as described in Section III-D. Provided the depth information, TACTO can generate the mesh in the scene and render corresponding images.

IV. SIMULATED EXPERIMENTS

We now demonstrate TACTO on a perception task to learn grasp stability from touch, and on a control task to manipulate a marble between two fingers using touch. We choose these two tasks based on previous works with real robots[6], [27] for better comparison between simulated and real environments. The experiments in simulation achieve similar results to the ones with real robots, which demonstrates the effectiveness and potentials of the simulated environment. For these experiments, we use TACTO in conjunction with the PyBullet physics engine [5].

A. Learning Grasp Stability in Simulation

In this task, we learn in simulation a classifier of grasp stability from vision and touch readings, following the previous work on real robots [6]. The goal is to predict whether a grasped object will be successfully lifted, based on the touch

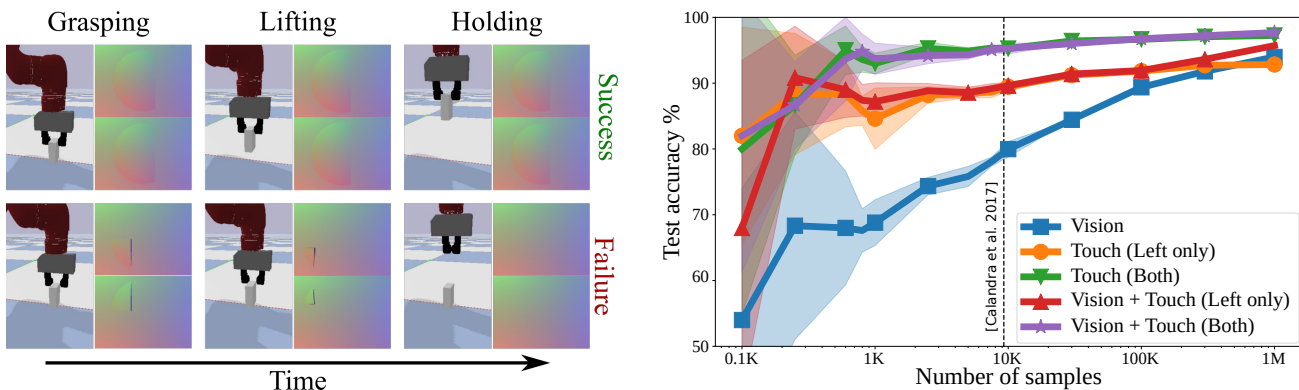


Fig. 10: Learning Grasp Stability. (Left) Examples of a successful grasp and a failure grasp. In the failure grasp, the object is only grasped by the corner and begins to slip after being lifted. (Right) Median and 68% percentile of the learned models when varying the number of data used. We compare using only vision, only touch and both vision and touch as inputs of the models. Results show that learning grasp stability from touch needs significantly less amount of data to achieve relative high accuracy compared to vision, and that increasing the amount of data helps to improve performance. The vertical dashed line shows the largest dataset collected on real robot [6]. In the simulation, we can experiment with a dataset more than two orders of magnitude larger.

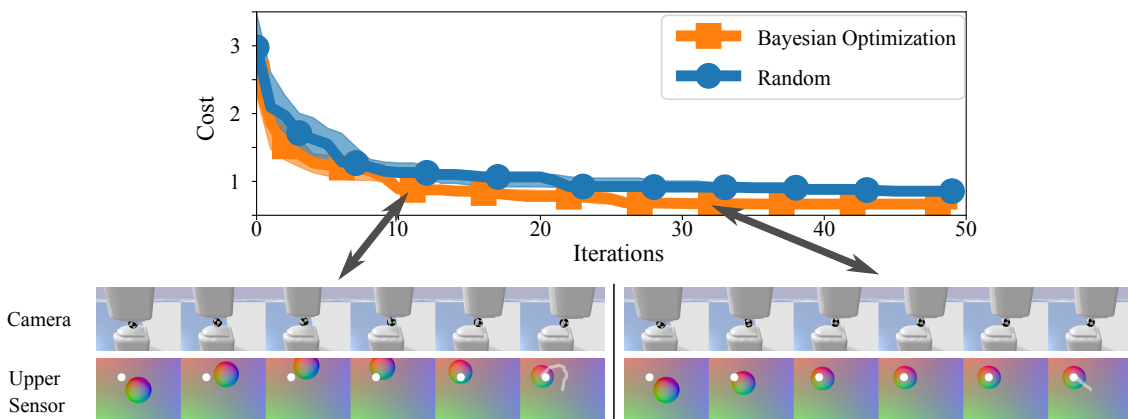


Fig. 11: Learning in-hand marble manipulation. (Upper) The learning curve (median and 68th percentile over 30 experiments) for Bayesian optimization and random search. (Lower) Examples of a marble rolling into different target locations within the fingers at the early and later learning stage. At the later stages of learning, the controller can roll the marble between fingers with faster speed and higher accuracy, while avoid dropping the marble.

and vision readings before the lifting. We aim to investigate whether the results in simulation match the real ones in [6].

Setting Our setup consists of two DIGIT sensors [4] mounted on a WSG-50 parallel-jaw gripper, and an external camera. We collected 1 million grasps in a self-supervised manner by randomizing the position, orientation and force applied by the gripper. The ground truth for each grasp was labeled depending on whether the object was still between the fingers after being lifted. Some examples are shown in Fig. 10. It took only one day to collect 1 million grasps with 5 threads. In our experiment, we used a single box object for demonstration purposes, but it would be straightforward to extend to various object datasets, such as YCB dataset [28], or Dex-Net [29] dataset. The resulting dataset collected with TACTO is several orders of magnitude larger than any publicly available dataset of grasps using tactile sensing [6]. As such, it allows us to evaluate the performance of learned grasp stability models in data regimes that are still unexplored. To learn the grasp stability, we trained ResNet-18 [30] neural network models that, given the raw vision and touch signals, would predict whether the grasp is stable or not after lift-off. The training procedure followed previous work [6]. For details, to fuse vision and touch signals, the feature vectors produced by

ResNet-18 were concatenated and fed to two fully connected layers, with 512 and 256 hidden units, to predict final results. To speed up training, we used ResNet-18 model pre-trained on ImageNet [31]. And we used vision and touch images with 160×120 pixels. The images were resized to 256×256 and randomly cropped to 224×224 for data augmentation.

To evaluate the performance of different dataset sizes, we used K-fold cross-validation and computed the median and 68% percentile of the classification accuracy. Due to computational limits, we used $K = 10$ up to 1000 datapoints, $k = 5$ up to 10k datapoints, and above 10k datapoints we evaluated a single train/test split 80/20%. We trained 10 epochs for each dataset size using Adam optimizer [32] with a learning rate of $5e^{-4}$ and batch size of 32.

Results The results shown in Fig. 10 suggests: (1) the model learned fast from touch with only a little amount of data, while vision requires 3 or 4 orders of magnitude more data to catch up; (2) single tactile sensor worked significantly worse than two tactile sensors, because the object may look stable from one tactile sensor while it unstably contacts the other side; (3) on the low-data regime, the result agrees with previous real-world experiments [6]: combining vision and touch worked best in most of the cases. Although touch-only

achieved comparable results to the combined model, we think the difference can increase with a larger object set; on the high-data regime, we can evaluate with 2 orders of magnitude more data in simulation compared to [6], and observe the trends that all the models keep improving still, and vision’s potentials with more data.

B. Learning In-hand Marble Manipulation

In this task, we learn in simulation to roll a marble to target locations in the sensor coordinate, following the previous work on real robots [27]. The setup includes two DIGIT sensors as shown in Fig. 11. The lower sensor is fixed, while the upper sensor is controlled to roll the marble. Since there are rich contacts with friction happening during the rolling, we aim to investigate how stable the TACTO and PyBullet are, and explore whether we could achieve similar performance to the real one in [27].

Setting We use position control with a maximal force for controlling the upper sensor. Specifically, we control the horizontal position of the upper sensor, and push the sensor vertically with maximal force to keep the marble in hand while rolling. We parameterize the controller as $\mathbf{u} = \mathbf{K}\bar{\mathbf{x}}$, where $\mathbf{u} \in \mathbb{R}^2$ is the desired velocity of the upper sensor in horizontal plane, $\bar{\mathbf{x}} \in \mathbb{R}^2$ is the error state between the current marble center location \mathbf{x} and the goal location \mathbf{x}^* in tactile space, and $\mathbf{K} \in \mathbb{R}^{2 \times 2}$ is the parameter to learn. The cost is defined as cumulative error distance in tactile space $\sum_t \|\bar{\mathbf{x}}_t\|$, and we set eight different target locations and take the average cost for robustness. We apply Bayesian optimization [33] with upper confidence bound to optimize the parameter \mathbf{K} automatically. Our main purpose here is to validate the simulation system, and provide benchmark experiments, however, the controller can be replaced by model predictive control and/or reinforcement learning to manipulate more complex objects [27], [34] and for dexterous hands [4].

Results Fig. 11 shows the quantitative and qualitative results of rolling a marble into desired locations. The system can learn to roll the marble into different target locations with few iterations and roll faster with more iterations. During the experiments, we validate that both PyBullet and TACTO run as expected without abnormal situations. Because the simulation is fast, it only takes 8 minutes for Bayesian optimization to learn marble manipulation with 50 iterations. It includes 6 minutes for optimizing the acquisition function, and 2 minutes for simulation, where there are 50 iterations, and each iteration includes 50 steps for rolling into each of 8 directions, rendering 20,000 tactile imprints of 160×120 resolution in total.

V. SIM2REAL EXPERIMENTS

The major goal with TACTO is to provide a platform to study representations and algorithms for robot learning using touch as a sensor modality, which the researchers can use to test algorithms before experimenting on real sensors/robots. However, another interesting venue for experimentation with tactile sensors is Sim2Real, where data from the simulator are used to directly train models that are either applied in

	position error (mm)	rotation error (degrees)
Random	11.75 ± 1.15	46.56 ± 5.91
Sim2Sim	0.41 ± 0.01	3.48 ± 0.34
Real2Real (16 datapoints)	4.45 ± 0.86	33.85 ± 1.07
Real2Real (32 datapoints)	3.48 ± 0.56	25.45 ± 2.07
Real2Real (64 datapoints)	2.01 ± 0.14	10.96 ± 0.24
Real2Real (128 datapoints)	0.76 ± 0.07	4.96 ± 0.70
Sim2Real (without augmentation)	4.56 ± 0.40	17.64 ± 2.34
Sim2Real (with augmentation)	1.66 ± 0.16	11.60 ± 4.65
Sim+Real (16 real datapoints)	1.55 ± 0.12	9.08 ± 1.65
Sim+Real (32 real datapoints)	1.36 ± 0.05	7.95 ± 1.60
Sim+Real (64 real datapoints)	1.24 ± 0.03	8.25 ± 0.77
Sim+Real (128 real datapoints)	0.52 ± 0.03	4.14 ± 0.57

Table II: Sim2Real experiment results for pose estimation of a cylinder, evaluated in position error of the contact center, and rotation error. We repeat the experiments 5 times, and report the mean error \pm standard deviation. We compare the results of Sim2Sim, Real2Real (with different training samples), and Sim2Real (with/without data augmentation and with different amount of mixed real data).

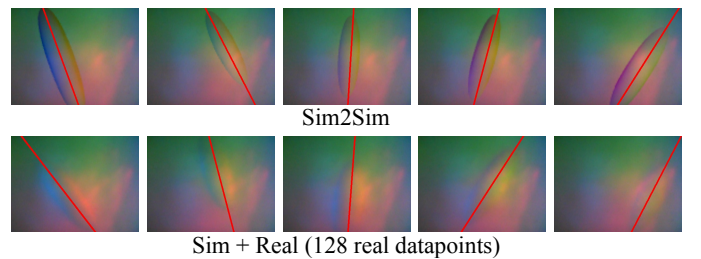


Fig. 12: Examples of tactile reading from simulated (*Top row*) and real data (*Bottom row*), and the corresponding estimated pose (red line).

the real-world or used in conjunction with a small number of real-world data. Here, we present a proof-of-concept of Sim2Real on a pose-estimation task, where the goal is to estimate the contact center and angle of a pen w.r.t. the sensor being touched.

Setting We generated in TACTO 10,000 simulated tactile imprints with corresponding poses, and collected 200 real tactile imprints with manually annotated poses. Following, we trained and validated the same convolutional neural networks on several combinations of datasets: 1) Sim2Sim trained and evaluated on simulated data. 2) Real2Real trained and evaluated on real data 3) Sim2Real trained on simulated data and evaluated on real data. 4) Sim+Real trained on simulated data mixed with a small amount of real data and evaluated on real data.

Results We show quantitative and qualitative results in Table II and Fig. 12 respectively. From the results in Table II, we can observe the sim2real gap (Sim2Real without augmentation). To bridge the gap, we add data augmentation. We find color jittering very helpful, where we randomly change the brightness and contrast for each RGB channel. It makes the model more robust to a variety of illumination conditions (Sim2Real with augmentation vs. Sim2Real without augmentation). We also compare the results between Sim2Real and Real2Real. Without any real data, Sim2Real with augmentation can achieve comparable results with Real2Real (64). When mixed with real data, Sim2Real consistently outperforms Real2Real with the same amount of real data. These show the potentials for increasing data efficiency by simulated data.

VI. CONCLUSION

In this paper, we introduce TACTO, a simulator for vision-based tactile sensors. TACTO is designed to provide an easy-to-use, fast and flexible simulator capable of generating realistic high-resolution readings. We demonstrate and validate TACTO on a perceptual task for learning grasp stability, and a control task for marble manipulation. Moreover, we provide a proof-of-concept that TACTO can be successfully used for Sim2Real applications. To foster the tactile sensing community and to enable robotics and machine learning researchers to make use of touch in simulation, we open-source TACTO at <https://github.com/facebookresearch/tacto>.

Future work will focus on improving the modeling of the effects of forces through the deformation of the elastomers, and ultimately generating more realistic readings. Besides, different experiment variations can be evaluated, such as the use of different deformation models, the effects of image filtering, the comparison of simulated tactile sensors with real-world sensors on varying geometries, and learning and transferring grasp policies on different objects.

ACKNOWLEDGMENT

We thank Edward Smith and Huazhe Xu for feedback on early versions of the simulator; Frederik Ebert, Stephen Tian and Akhil Padmanabha for providing the models of the OmniTact;

REFERENCES

- [1] T. Erez, Y. Tassa, and E. Todorov, "Simulation tools for model-based robotics: Comparison of bullet, havok, mujoco, ode and physx," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2015, pp. 4397–4404.
- [2] W. Yuan, S. Dong, and E. H. Adelson, "Gelsight: High-resolution robot tactile sensors for estimating geometry and force," *Sensors*, 2017.
- [3] A. Padmanabha, F. Ebert, S. Tian, R. Calandra, C. Finn, and S. Levine, "OmniTact: A multi-directional high-resolution touch sensor," *IEEE International Conference on Robotics and Automation (ICRA)*, 2020.
- [4] M. Lambeta, P.-W. Chou, S. Tian, B. Yang, B. Maloon, V. R. Most, D. Stroud, R. Santos, A. Byagowi, G. Kammerer, D. Jayaraman, and R. Calandra, "DIGIT: A novel design for a low-cost compact high-resolution tactile sensor with application to in-hand manipulation," *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, no. 3, pp. 3838–3845, 2020.
- [5] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," <http://pybullet.org>, 2016–2019.
- [6] R. Calandra, A. Owens, M. Upadhyaya, W. Yuan, J. Lin, E. H. Adelson, and S. Levine, "The feeling of success: Does touch sensing help predict grasp outcomes?" *Conference on Robot Learning (CORL)*, pp. 314–323, 2017.
- [7] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.
- [8] P. Ruppel, Y. Jonetzko, M. Görner, N. Hendrich, and J. Zhang, "Simulation of the syntouch biotac sensor," in *International Conference on Intelligent Autonomous Systems*. Springer, 2018, pp. 374–387.
- [9] A. Melnik, L. Lach, M. Plappert, T. Korthals, R. Haschke, and H. Ritter, "Using tactile sensing to improve the sample efficiency and performance of deep deterministic policy gradients for simulated in-hand manipulation tasks," *Frontiers in Robotics and AI*, vol. 8, 2021.
- [10] J. Geukes, M. Nakatenus, and R. Calandra, "Gazebo plugin for the icub skin," <https://github.com/robertocalandra/icub-gazebo-skin>, 2017.
- [11] A. Habib, I. Ranatunga, K. Shook, and D. O. Popa, "Skinsim: A simulation environment for multimodal robot skin," in *IEEE International Conference on Automation Science and Engineering (CASE)*, 2014, pp. 1226–1231.
- [12] F. R. Hogan, M. Jenkin, S. Rezaei-Shoshtari, Y. Girdhar, D. Meger, and G. Dudek, "Seeing through your skin: Recognizing objects with a novel visuotactile sensor," in *IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 1218–1227.
- [13] D. F. Gomes, A. Wilson, and S. Luo, "Gelsight simulation for sim2real learning," *ICRA Workshop on ViTac: Integrating Vision and Touch for Multimodal and Cross-modal Perception*, 2019.
- [14] P. Sodhi, M. Kaess, M. Mukadam, and S. Anderson, "Learning tactile models for factor graph-based estimation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 13 686–13 692.
- [15] Z. Ding, N. F. Lepora, and E. Johns, "Sim-to-real transfer for optical tactile sensing," *IEEE International Conference on Robotics and Automation (ICRA)*, 2020.
- [16] C. Sferrazza, T. Bi, and R. D'Andrea, "Learning the sense of touch in simulation: a sim-to-real strategy for vision-based tactile sensing," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 4389–4396.
- [17] N. Kuppaswamy, A. Castro, C. Phillips-Grafflin, A. Alspach, and R. Tedrake, "Fast model-based contact patch and pose estimation for highly deformable dense-geometry tactile sensors," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1811–1818, 2019.
- [18] M. Nimier-David, D. Vicini, T. Zeltner, and W. Jakob, "Mitsuba 2: A retargetable forward and inverse renderer," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–17, 2019.
- [19] E. Donlon, S. Dong, M. Liu, J. Li, E. Adelson, and A. Rodriguez, "Gelslim: A high-resolution, compact, robust, and calibrated tactile-sensing finger," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 1927–1934.
- [20] B. Romero, F. Veiga, and E. Adelson, "Soft, round, high resolution tactile fingertip sensors for dexterous robotic manipulation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 4796–4802.
- [21] D. Shreiner, G. Sellers, J. Kessenich, and B. Licea-Kane, *OpenGL programming guide: The Official guide to learning OpenGL, version 4.3*. Addison-Wesley, 2013.
- [22] L. Perez and J. Wang, "The effectiveness of data augmentation in image classification using deep learning," *arXiv preprint arXiv:1712.04621*, 2017.
- [23] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb, "Learning from simulated and unsupervised images through adversarial training," in *IEEE conference on computer vision and pattern recognition (CVPR)*, 2017, pp. 2107–2116.
- [24] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, "Cycada: Cycle-consistent adversarial domain adaptation," in *International conference on machine learning*, 2018, pp. 1989–1998.
- [25] M. Matl, "Pyrender," <https://github.com/mmatl/pyrender>.
- [26] D. Ma, E. Donlon, S. Dong, and A. Rodriguez, "Dense tactile force estimation using gelslim and inverse fem," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2019, pp. 5418–5424.
- [27] S. Tian, F. Ebert, D. Jayaraman, M. Mudigonda, C. Finn, R. Calandra, and S. Levine, "Manipulation by feel: Touch-based control with deep predictive models," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2019, pp. 818–824.
- [28] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, "The ycb object and model set: Towards common benchmarks for manipulation research," in *IEEE International conference on advanced robotics (ICAR)*, 2015, pp. 510–517.
- [29] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. Aparicio, and K. Goldberg, "Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics," in *Robotics: Science and Systems (RSS)*, July 2017.
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [31] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE conference on computer vision and pattern recognition*, 2009, pp. 248–255.
- [32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2015.
- [33] J. Snoek, H. Larochelle, and R. P. Adams, "Practical bayesian optimization of machine learning algorithms," in *Advances in neural information processing systems*, 2012, pp. 2951–2959.
- [34] M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray et al., "Learning dexterous in-hand manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 3–20, 2020.