

Sim-to-Real for Robotic Tactile Sensing via Physics-Based Simulation and Learned Latent Projections

Yashraj Narang^{*1}, Balakumar Sundaralingam^{*1}, Miles Macklin¹, Arsalan Mousavian¹, Dieter Fox^{1,2}

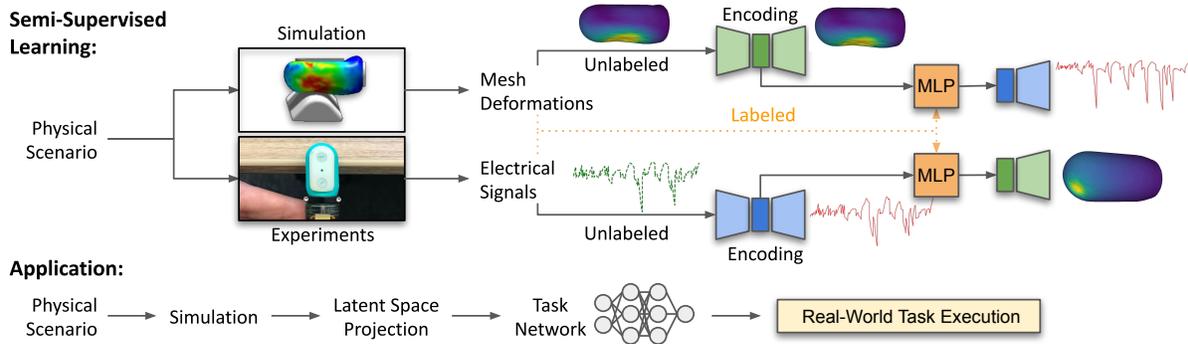


Fig. 1: Overview. We develop an efficient 3D FEM model of a SynTouch BioTac sensor to simulate contact interactions, and we conduct similar real-world experiments. In a learning phase, we train autoencoders to reconstruct unlabeled FEM deformations and real-world electrical signals. With a small amount of labeled data, we subsequently train MLPs to project between the FEM and electrical latent spaces. At test time, we use these learned latent projections to perform cross-modal transfer between FEM and electrical data for unseen contact interactions. During downstream application, we 1) accurately synthesize BioTac electrical signals, and 2) estimate the shape and location of contact patches, facilitating real-world task execution.

Abstract—Tactile sensing is critical for robotic grasping and manipulation of objects under visual occlusion. However, in contrast to simulations of robot arms and cameras, current simulations of tactile sensors have limited accuracy, speed, and utility. In this work, we develop an efficient 3D finite element method (FEM) model of the SynTouch BioTac sensor using an open-access, GPU-based robotics simulator. Our simulations closely reproduce results from an experimentally-validated model in an industry-standard, CPU-based simulator, but at 75x the speed. We then learn latent representations for simulated BioTac deformations and real-world electrical output through self-supervision, as well as projections between the latent spaces using a small supervised dataset. Using these learned latent projections, we accurately synthesize real-world BioTac electrical output and estimate contact patches, both for unseen contact interactions. This work contributes an efficient, freely-accessible FEM model of the BioTac and comprises one of the first efforts to combine self-supervision, cross-modal transfer, and sim-to-real transfer for tactile sensors.

I. INTRODUCTION

Tactile sensing is critical for grasping and manipulation under visual occlusion, as well as for handling delicate objects [1]. For example, humans leverage tactile sensing when retrieving keys, striking a match, holding a wine glass, and grasping fresh fruit without damage. In robotics, researchers are actively developing a wide variety of tactile sensors (e.g., [2]–[13]). These sensors have been used for tasks such as slip detection [14]–[16], object classification [17], [18], parameter estimation [19], force estimation [20]–[23], contour following [24], and reacting to humans [25].

For other aspects of robotics, such as robot kinematics, dynamics, and cameras, accurate and efficient simulators

have advanced the state-of-the-art in task performance. For example, simulators have enabled accurate testing of algorithms for perception, localization, planning, and control [26]; generation of synthetic datasets for learning such algorithms [27]–[29]; efficient training of control policies via reinforcement learning (RL) [30]–[32]; and execution of online algorithms, with the simulator as a model [33]. These capabilities have in turn reduced the need for costly, time-consuming, dangerous, or intractable experiments.

However, simulators for tactile sensing are still nascent. For the SynTouch BioTac sensor [2], [3], as well as vision-based tactile sensors, most simulation studies have focused on the *inverse* problem of interpreting sensor output in terms of physical quantities (e.g., [20], [22], [23], [34]). Far fewer efforts have addressed the *forward* problem of synthesizing sensor output, and perhaps none have accurately generalized to diverse contact scenarios. Forward simulation is invaluable for simulation-based training, which coupled with domain adaptation, can enable effective policy generation.

The dearth of tactile simulation capabilities is a result of its inherent challenges. Accurate tactile sensor simulators must model numerous contacts, complex geometries, and elastic deformation, which can be computationally prohibitive [35]. Simulators must also capture multiphysics behavior, as tactile sensors are cross-modal: for instance, the BioTac transduces skin deformations to fluidic impedances. Furthermore, whereas a small parameter set (e.g., camera intrinsics) can describe variation among visual sensors, no equivalent exists for tactile sensors. Due to manufacturing variability, even sensors of the same type can behave disparately [22], [36].

In this work, we address forward simulation and sim-to-real transfer for the BioTac (Fig. 1). We first develop 3D finite element method (FEM) simulations of the BioTac

^{*}These authors contributed equally.

¹NVIDIA Corporation, Seattle, USA. ²Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, USA.

using a GPU-based robotics simulator [37]; the FEM simulations predict contact forces and deformation fields for the sensor for arbitrary contact interactions. These simulations are designed to reproduce our previous results [22], [23], which utilized an industry-standard, commercial, CPU-based simulator and were carefully validated against real-world experiments. However, the new simulator is freely accessible, and the simulations execute 75x faster.

We then map FEM output to real-world BioTac electrical signals by leveraging recent methods in self-supervised representation learning. Specifically, we collect a large unlabeled dataset of sensor deformation fields from simulation, as well as a smaller dataset of electrical signals from real-world experiments; we then learn latent representations for each modality using variational autoencoders (VAE) [38], [39]. Next, we learn a cross-modal projection between the latent spaces using a small amount of supervised data. We demonstrate that this learned latent projection allows us to accurately predict BioTac electrical signals from simulated deformation fields for unseen contact interactions, including unseen objects. We can also execute the inverse mapping (from signals to deformations) with higher fidelity than in [22], [23], illustrated via a contact-patch estimator.

To summarize, our key contributions are the following:

- 1) An accurate, efficient, and freely-accessible 3D FEM-based simulation model for the BioTac
- 2) A novel application of self-supervision and learned latent-space projections for facilitating cross-modal transfer between FEM and electrical data
- 3) Demonstrations of sim-to-real transfer through accurate synthesis of BioTac electrical data and estimation of contact patches, both for unseen contact interactions

The simulation model, as well as additional implementation details and analyses, will be posted on our website.¹

II. RELATED WORKS

In this section, we review research efforts in sim-to-real transfer, self-supervision, and cross-modal transfer involving tactile sensors that are widely used in current robotics projects. For a recent comprehensive review, see [40].

A. Sim-to-Real for Vision-Based Tactile Sensors

In [41], visual output of the GelSight [4] was simulated using a depth camera in Gazebo [42], a calibrated reflectance model, and blurring to approximate gel contact. Quantitative evaluation was limited. In [43], a custom marker-based tactile sensor [9] was simulated using FEM, optics models, and synthetic noise. A U-net was trained on the synthetic data to regress to contact force fields with a resultant error of $0.14N$; almost all training and testing was conducted on 1 normally-oriented indenter. In [44], the pin locations of the TacTip [5] were simulated using an approximate deformation model in Unity [45]. Parameters were tuned for a plausible baseline, and domain randomization was performed. A fully connected network (FCN) was trained on synthetic data to

regress to contact location and angle with a minimum error of $0.5\text{-}0.7\text{mm}$ and 0.25rad . The sensor examined was large, and contact was made over limited orientations.

In comparison to these studies, we use a compact sensor, accurately simulate contact and deformation, do not perform domain adaptation beyond projection, and conduct simulations and experiments on 17 objects over diverse kinematics.

B. Sim-to-Real for Non-Vision-Based Tactile Sensors

In [46], the electrical output of the BarrettHand [47] capacitive sensor array was simulated using soft contact in PyBullet [48]. RL policies for stable grasps were trained on the simulator and transferred to the real world. Binary thresholding was applied to tactile signals, limiting precision. In [49], an electrical resistance tomography sensor was simulated using a simplified FEM deformation model, an empirical piezoresistive model, and an FEM conductivity model. FCNs were trained on synthetic data to predict discrete contact location and resultant force, with an 82% success rate for contact and a mean force error of $0.51N$. For unseen contact scenarios, the error increased to $5.0N$.

Finally, for the BioTac sensor, multiple studies have addressed the inverse problem of converting electrical output to physical quantities, such as contact location [22], [23], [50], force [14], [21]–[23], [50], and deformation [22], [23]. In particular, our previous work [22] presented a 3D FEM model of the BioTac, which was built with the industry-standard, CPU-based simulator ANSYS [51] and carefully validated against experimental data. Electrode signals were then mapped to simulator outputs via PointNet++ [52]. The simulations were slow (7min each on 6 CPUs). In addition, the forward problem of synthesizing electrical output was not addressed; some progress was made in our extension [23].

The forward problem has been further explored in [36], [53], [54]. In [53], the BioTac was simulated with an approximate contact model in Gazebo. A deep network was trained to regress from estimated force and real-world contact location to electrical outputs. An existing location estimator was tested on synthetic and real-world data with a mean difference of 0.83mm . Training and testing was conducted only on 1 spherical indenter. In [36], PointNet was used to regress from RGB-D images and grasp parameters to tactile readings, and in [54], semi-supervised learning was applied. Electrode prediction errors were relatively high for both the simple case of unseen *samples*, as well as the difficult case of unseen *objects*; numerical comparison is provided later.

In contrast to the preceding works, we develop an efficient FEM model, conduct simulations and experiments with 17 objects, regress to continuous electrical signals, and demonstrate accurate predictions for unseen objects.

C. Tactile Self-Supervision and Cross-Modal Transfer

Numerous studies in robotics have established the utility of multimodal data in task-specific learning. For example, in [55], [56], vision and tactile sensing were combined to predict grasp outcomes and select adjustments, and in [57], vision and force/torque (F/T) sensing were combined to

¹<https://sites.google.com/nvidia.com/tactiledata2>

perform manipulation tasks. In addition, in [58], vision, F/T sensing, and proprioception were used to learn a joint latent space via self-supervision with autoencoders; the output served as perception input to an RL agent for peg-in-hole.

Simultaneously, recent works both outside and within robotics have investigated cross-modal transfer. In [59], audio-video transfer was performed by learning a shared representation via a bimodal autoencoder; networks trained on one modality were then able to classify the other. In [60], audio-image transfer was achieved via generative adversarial networks (GAN). In [61], [62], cross-modal transfer was performed between data from cameras and vision-based tactile sensors, and in [36], [54], transfer was achieved between cameras and output from the BioTac SP.

Finally, previous efforts have applied distinct network architectures to encode BioTac-specific data. In our previous work [21]–[23], a 3D voxel-grid network, PointNet++, and an FCN were implemented to encode BioTac electrode data and regress to physical quantities such as forces and deformations. In [63], a 2D CNN was used to predict tactile directionality, and in [64], a graph convolutional network (GCN) was used to predict grasp stability.

We draw upon the preceding works with some distinctions. Analogous to [59], we learn latent representations of FEM and BioTac electrical data via self-supervision with autoencoders for cross-modal transfer. Like [58], we learn modality-specific representations, reducing training time and eliminating zero-inputs for non-present modalities. Unlike both, we never formulate a joint representation, but instead learn a projection between the latent spaces using a small amount of supervised data. To encode BioTac electrical data, we use VAEs, as for vision-based tactile sensors in [6].

III. METHODS

In this section, we first discuss our 3D FEM model, which predicts BioTac contact forces and deformation fields for arbitrary contact interactions. We then discuss our implementation of self-supervision and latent-space projection, which can synthesize BioTac electrical output from unlabeled FEM output and predict contact patches from electrical input. Finally, we describe the simulations and experiments used to collect the FEM and electrical data used in this paper.

A. Finite Element Modeling

FEM is a variational numerical formulation that 1) divides complex global geometries into simple subdomains, 2) solves the weak form of the governing PDEs in each subdomain, and 3) assembles the solutions into a global one. In 3D FEM, objects are represented as volumetric meshes, which consist of 3D elements (e.g., tetrahedrons or hexahedrons) and their associated nodes. With careful model design, high-quality meshes, and small timesteps, FEM predictions for deformable bodies can be exceptionally accurate [65], [66].

In this work, 3D FEM was performed using NVIDIA’s GPU-based Isaac Gym simulator [37]. Isaac Gym models the internal dynamics of deformable bodies using a co-rotational linear-elastic constitutive model; these bodies interact with

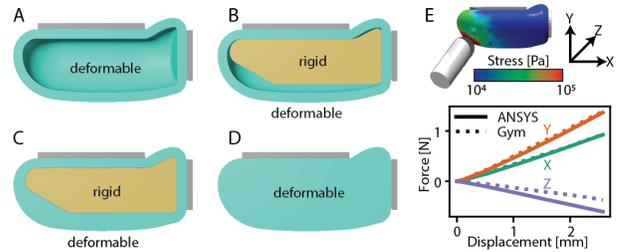


Fig. 2. A-D) 4 geometric configurations of the BioTac. Cross-sections are shown. A) Deformable shell, modeling the real-world rubber skin. B) Deformable shell with rigid core, modeling the skin and real-world plastic core. C) Deformable solid with rigid core. D) Deformable solid. Grey rigid bodies were used to apply fixed boundary conditions. E) Optimization results for FEM model. The material parameters of representation B were tuned in Isaac Gym to reproduce the force profile of a shear-rich indentation from ANSYS in [22], [23]. The von Mises stress distribution is visualized for the indentation midpoint. The mean ℓ^2 -norm of the force error vector was $0.125N$. Indenter displacement is relative to the point of initial contact.

external rigid objects via an isotropic Coulomb contact model [67]. The resulting nonlinear complementarity problem is solved via implicit time integration using a GPU-based Newton solver [68]. At each timestep, Isaac Gym returns nodal positions (i.e., deformation fields), nodal contact forces (used to compute resultant forces), and element stress tensors.

To create the FEM model for the BioTac, high-resolution triangular meshes for the external and internal surfaces of the BioTac skin were first extracted from the ANSYS model in [22], [23] and simplified via quadric edge collapse decimation in MeshLab [69]. A volumetric mesh was then generated with fTetWild [70]; similar to [22], [23], the mesh had ≈ 4000 nodes. Fixed boundary conditions (BCs) were applied to 2 sides of the skin to model the BioTac nail and clamp, respectively; these BCs were implemented by introducing thin rigid bodies at the corresponding locations (visible in Fig. 2A-D), which were attached to adjacent nodes on the skin. External rigid objects (e.g., indenters) were driven into the BioTac via a stiff proportional controller.

Relative to the experimentally-validated ANSYS model in [22], [23], the Isaac Gym model makes 3 critical approximations: 1) collisions are resolved via boundary-layer expanded meshes [71], rather than a normal-Lagrange method, 2) a compressible linear-elastic model is used for the skin, rather than a Neo-Hookean model [72], 3) the internal fluid is not modeled. The effects of the first approximation are mitigated by using small collision thicknesses and timesteps ($1e-4s$). However, the second and third approximations are mitigated by endowing the Isaac Gym model with sufficient expressivity and optimizing it to reproduce ANSYS results.

Specifically, 4 distinct geometric configurations of the BioTac were considered in Isaac Gym (Fig. 2A-D). The material properties (elastic modulus E , Poisson’s ratio ν , and friction μ of the BioTac skin) were designated as free parameters. In [22], [23], 1 shear-rich indentation of the BioTac was used to calibrate the ANSYS model against real-world data; in this work, the same indentation was resimulated in Isaac Gym and used to calibrate the Gym model against ANSYS data (Fig. 2E). For each configuration, the material properties were optimized via sequential least-squares programming to

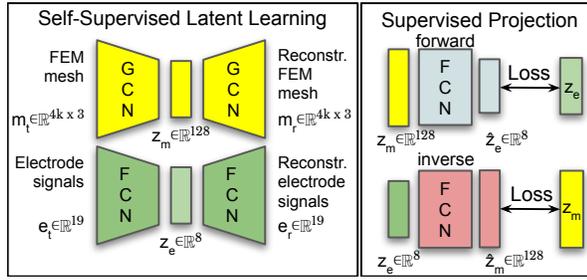


Fig. 3. Learning structure. To map between FEM deformations and BioTac electrode signals, modality-specific latent representations were learned via self-supervision. Specifically, graph convolutional networks (GCN) compressed deformed meshes with 4000 nodes to a 128-dim. latent space, and fully-connected networks (FCN) compressed the 19 electrode signals to an 8-dim. latent space. Next, FCNs were used on a small supervised dataset to learn forward and inverse projections between the latent spaces.

reproduce the force-deflection profile from ANSYS for the indentation. The cost was defined as the RMS ℓ^2 -norm of the force error vector over time. Subsequently, each optimized configuration was evaluated by resimulating 358 additional indentations from 8 indenters in [22], [23], and comparing the results to the force-deflection profiles from ANSYS.

Among the 4 configurations, the *deformable solid* (Fig. 2D) produced the lowest cost during optimization; however, the *deformable shell with rigid core* (Fig. 2B) produced the lowest cost during evaluation and was thus selected. For this representation, the optimal values of E , ν , and μ were $1.55e6Pa$, 0.316 , and 0.783 , respectively; the mean ℓ^2 -norm of the force error vector was $0.125N$ (Fig. 2E). In comparison, the optimal values for the ANSYS model in [22] were $2.80e5Pa$, 0.5 , and 0.186 , indicating that the Isaac Gym model compensated for its linearity and lack of fluid by increasing elastic modulus, compressibility, and friction.

B. Learning Latent Space Projections

Although FEM captures the effects of contact on BioTac deformations, the BioTac then transduces deformations to fluidic impedances measured at electrodes. Simulating the mapping between deformations and impedances is complex; thus, this mapping was learned. Specifically, modality-specific latent representations were learned using self-supervision, facilitating compression, mitigating noise, and reducing overfitting. Projections were then learned between the latent spaces, enabling cross-modal transfer (Fig. 3).

To learn a latent representation for the FEM deformations, convolutional mesh autoencoders from [39] were trained, which applied graph convolutional networks (GCN) and reduced the mesh data from 4000 nodal positions ($\mathbb{R}^{4k \times 3}$) to 128 units. To learn a latent representation for real-world BioTac electrode signals, an autoencoder composed of FCNs was trained, which reduced the electrode data from 19 impedances to 8 units. Latent dimensions were chosen via hyperparameter search. Both networks were trained as VAEs to generate smooth mappings to the latent space [38].

To learn the projections between the latent spaces, 2 FCNs were trained. The first network projected forward from the FEM mesh latent space z_m to the BioTac electrode latent space z_e , whereas the second network projected inversely

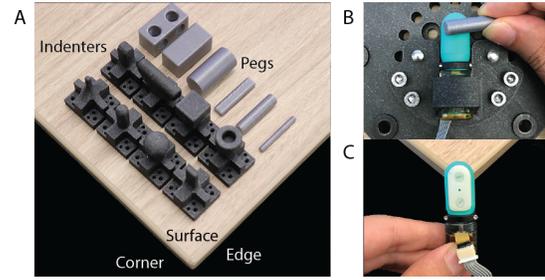


Fig. 4. Data collection. A) 17 objects and table features were used to generate data in simulation and the real world. The objects were designed to mimic diverse physical features (e.g., corners, edges, buttons, holes). B-C) The BioTac was constrained when applying kinematically-randomized indentations with the pegs, but unconstrained when contacting table features. from z_e to z_m . During training, the previously-described VAEs were frozen and provided with supervised data from [22], generating latent pairs of z_m and z_e . These pairs were used to train the projection networks with an RMS loss. Without freezing the VAEs, the networks overfit.

In the FEM deformation VAE, the encoder consisted of an initial convolution with filter size 128, 4 “convolve+downsample” layers with sizes [128, 128, 256, 64] and downsampling factors [4, 4, 4, 2], a convolution with size 64, and 2 fully-connected layers with dimensions [512, 128]. The decoder was symmetric with the encoder, using “upsample+convolve” instead of “convolve+downsample”. In the BioTac electrode VAE, the encoder consisted of 4 fully connected layers with [256, 128, 64, 8] neurons, respectively. The decoder was again symmetric with the encoder. The forward and inverse projection networks consisted of 3 fully connected layers with [256, 128, 128] and [128, 128, 256] neurons, respectively, and dropout of 0.3. Exponential linear unit (ELU) activations were applied, and the Adam optimizer was used with initial learning rate of $1e-3$ and decay of 0.95.

C. Dataset Collection

For the preceding learning steps, data was collected in both simulation and the real world. For learning the latent representations, unlabeled mesh data was collected by simulating kinematically-randomized interactions on the optimized BioTac model with 6 pegs and 3 table features (surfaces, edges, and corners) in Isaac Gym. Unlabeled experimental electrode data was collected by manually indenting these objects in the real world. For learning the latent projections, labeled data was collected by exactly resimulating 359 indentations from [22], [23] on the optimized BioTac model in Isaac Gym (as stated in Sec. III-A); these were directly matched to corresponding experimental electrode data in the dataset from [22], [23]. 72% of contact interactions were allocated for training, 18% for validation, and 10% for testing.

In total, 2.6k unique contact interactions were executed and 50k timesteps of FEM data were sampled. All objects used in simulation and experiments are shown in Fig. 4. Data collection visualizations are in the supplementary video.

IV. EXPERIMENTS & RESULTS

In this section, we present our results on FEM validation, synthesis of BioTac electrode signals from FEM deformations, and contact patch estimation from electrode signals.

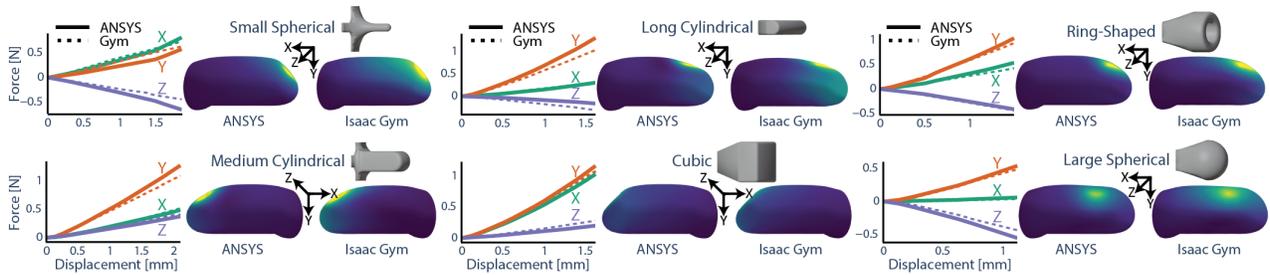


Fig. 5. Validation of FEM model. The optimized model from Isaac Gym was compared to an experimentally-validated model from ANSYS. Force-deflection profiles, as well as nodal deformation fields at maximum indentation depth, are illustrated for 6 randomly-selected indentations with 6 different indenters. Nodal deformation fields are colored according to corresponding displacements (i.e., change in nodal positions relative to the no-contact state). The 2 not-pictured indenters are a medium-sized spherical indenter and a small cylindrical indenter. Larger, higher-contrast graphics are on our website.

A. FEM Validation

As described earlier, the force-deflection profiles produced by the optimized model in Isaac Gym were compared to those produced by ANSYS over 358 indentations distributed across 8 indenters (Fig. 5). The mean ℓ^2 -norm of the force error vectors over all indentations ranged from $0.0876N$ for a medium-sized cylindrical indenter (*less* than the training error of $0.125N$) to $0.259N$ for a medium-sized spherical indenter, with a mean of $0.153N$ across all indenters. Thus, despite being optimized using force-deflection data from only a single indentation, the Isaac Gym model strongly generalized across a diverse range of objects and indentations.

The corresponding FEM deformation fields (i.e., the nodal positions of the deformed FEM meshes) were also compared between Isaac Gym and ANSYS (Fig. 5). For each dataset, the maximum and mean ℓ^2 -norms of the nodal displacement vectors were computed over all indentations. The mean error between the maximum norms across all indenters was $1.41e-4m$, and the mean error between the mean norms was $1.81e-5m$. Thus, again, the Isaac Gym model was shown to strongly generalize. These low errors were particularly important, as the nodal deformation fields from Isaac Gym were used as input for subsequent learning.

Finally, simulation speed was compared between the Isaac Gym and ANSYS models. The total simulation time for all 359 indentations was approximately 42 hours (7.08 minutes per sim) on 6 CPUs in ANSYS, but 33 minutes (5.57 seconds per sim) using 8 parallel environments (1 per indenter) on 1 GPU in Isaac Gym. For clarity, Isaac Gym can only currently simulate deformable solids with a linear-elastic model and

linear tetrahedral elements; such a model comprises only a small fraction of those that can be simulated within state-of-the-art FEM software such as ANSYS. However, for the current application, Isaac Gym is highly favorable.

B. Learning-Based Regression and Estimation

For regression from FEM deformations to BioTac electrode signals, 2 learning methods were evaluated: 1) the method proposed in this paper, denoted *Latent Projection*, which used unlabeled data for latent representation and labeled data for projection, and 2) a PointNet++ baseline [52], denoted *Fully Supervised*, which used only labeled data, with 128 nodes sampled from the FEM mesh as in [22], [23]. For reference, output is also shown for the *Electrode VAE*, which is used when training *Latent Projection*.

When evaluating generalization to novel objects, networks were trained on all objects *except* the ring (see Fig. 5) and tested on this indenter; the ring has the most unique (thus, challenging) geometry in the dataset. These experiments are denoted “Unseen Object.” (Results for other unseen objects are on our website.) When evaluating generalization to novel contact interactions with seen objects, the trained networks are tested on unseen interactions with the other indenters. These experiments are denoted “Unseen Trajectory.”

A qualitative comparison of regression results between the learning methods is depicted in Fig. 6A for 2 high-signal, high-variation electrodes over numerous interactions. Raw electrode values were between $[0, 4095]$ digital output units and were tared and normalized to $[-1, 1]$. For the challenging “Unseen Object” case, *Latent Projection* could predict several signal peaks over multiple electrodes that *Fully*

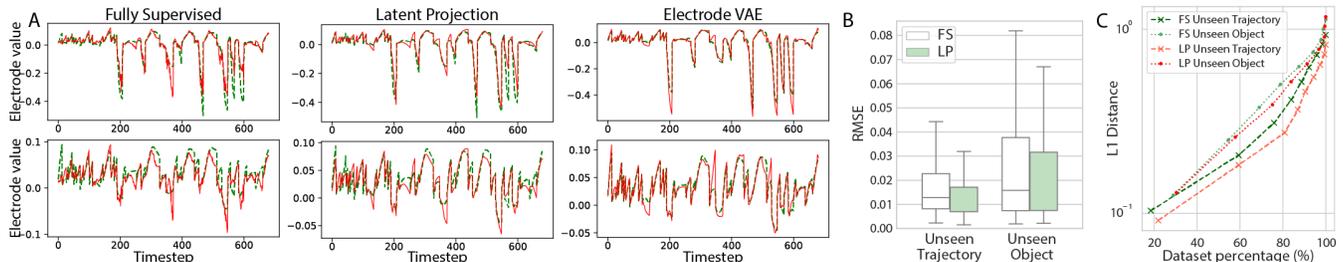


Fig. 6. Electrode prediction results. A) Visual comparisons for 2 high-magnitude, high-variation electrodes during contact interactions with an unseen object. Electrode range corresponds to a force range of $0.5\text{--}19.8N$. Each peak corresponds to a distinct interaction. Green and red lines indicate ground-truth and prediction, respectively. Our latent projection (LP) approach can predict peaks more accurately than a fully-supervised (FS) baseline, and often outperforms the VAE used in training. B) RMS error over all electrodes and interactions, for unseen trajectories and objects. Our approach has lower median errors and interquartile ranges. C) Coverage plot, with ℓ^1 distance to ground-truth. Our approach has lower errors over nearly the full data distribution.

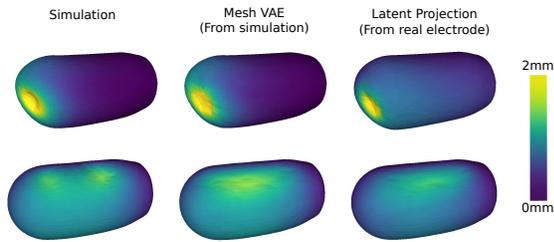


Fig. 7. Predicted FEM deformations for unseen trajectories and objects. The columns show the raw FEM output, the output of the mesh VAE, and the prediction from real-world electrode signals through our latent projection. The colorbar shows the Euclidean distance from the undeformed state. Top) “Unseen Trajectory” case, for a randomly-selected indentation. Predictions are consistently accurate. Bottom) “Unseen Object” case, for a randomly-selected indentation with the most distinct indenter, the ring. Predictions do not capture the bimodal deformation distribution due to limitations of the VAE, which has not seen any examples of such distributions in training.

Supervised could not capture. Additionally, *Latent Projection* predictions were consistently noise-free, whereas *Fully Supervised* ones exhibited low-magnitude, high-frequency noise. Finally, *Latent Projection* often outperformed *Electrode VAE*, showing the importance of mesh information in electrode signal synthesis. Predictions for the easier “Unseen Trajectory” case have higher fidelity and are thus not shown.

Quantitative comparisons between the learning methods are illustrated in Fig. 6B-C. RMS errors over all electrodes and interactions are compared in Fig. 6B. The *Latent Projection* approach performs better than *Fully Supervised* for both “Unseen Trajectory” and “Unseen Object,” in terms of both median error and interquartile range. Median errors for “Unseen Trajectory” were *Fully Supervised*: 0.012 (25 raw units) and *Latent Projection*: 0.010 (20 units), and those for “Unseen Object” were 0.016 (32 units) and 0.015 (31 units), respectively. Our errors are substantially lower than the errors from [54], which were 195 units for an easier case of unseen *samples*, and 305 units for unseen objects. A coverage plot is shown in Fig. 6C, where the ℓ^1 distance to the ground-truth electrode signals is depicted. For both “Unseen Trajectory” and “Unseen Object,” *Latent Projection* outperforms *Fully Supervised* for nearly the entire distribution of data.

For regression from BioTac electrode signals to FEM deformations, we again evaluate *Latent Projection*. For reference, we also show output of the FEM *Mesh VAE*, which is used when training *Latent Projection*. A visual comparison is shown against ground-truth FEM output in Fig. 7 for random indentations. For “Unseen Trajectory,” *Latent Projection* consistently predicted ground-truth, capturing deformation magnitudes and distributions. For “Unseen Object,” *Latent Projection* consistently captured magnitudes and distributions, but not bimodality. As seen from *Mesh VAE*, this limit is due to the mesh autoencoder (specifically, bottleneck size) rather than the projection. As before, we only show results for the most challenging unseen object, the ring; strong performance on other objects is shown in the video.

As a final demonstration, we also conducted free-form interactions of the BioTac with unseen objects and visualized the estimated contact patches (Fig. 8). For diverse pegs and table features, predicted deformations were visually accurate. For instance, contact patch locations were accurately

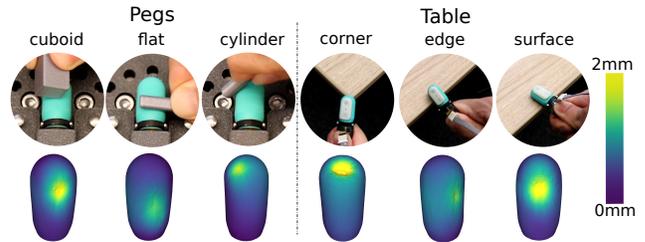


Fig. 8. Contact patch estimation for free-form interactions. The colorbar shows the Euclidean distance from the undeformed state. Top) Real-world contact with various unseen objects and table features. Bottom) Predicted contact patches using our learned latent projections. Contact patches match contacting features; for example, interactions with the corners of the flat peg and table produce highly-localized, high-magnitude deformation.

predicted across the full spatial limits of the BioTac, and interactions with the corners of a cuboid peg and table showed high-magnitude, highly-localized patch deformations.

V. DISCUSSION

In this paper, we present a framework for synthesizing BioTac electrical signals and estimating contact patches for novel contact interactions. The framework consists of 1) a 3D FEM model, which simulates contact between the BioTac and objects and outputs BioTac deformation fields, 2) VAEs that output compressed representations of the deformation fields and electrode signals, and 3) projection networks that perform cross-modal transfer between the representations to facilitate regression of electrode signals or contact patches.

This work has several key contributions. First, compared to our previously-presented, experimentally-validated FEM model [22], the current model is nearly equivalent, available in an open-access robotics simulator, and 75x faster. Second, our work presents one of the first applications of cross-modal self-supervision for tactile sensing; we show that this approach outperforms supervised-only methods for regressing to BioTac electrical signals. Third, for the first time, we accurately predict these signals for unseen interactions, including unseen objects. Finally, we can reconstruct BioTac deformations from real electrical signals with high fidelity.

The present study also has limitations. First, although our FEM model is substantially faster than previous efforts, it currently takes approximately 5.6s to simulate 6mm of indentation, which prohibits dynamic model-predictive control applications. Furthermore, although our networks accurately predicted electrode signals for unseen trajectories and objects, evaluation was performed for 1 BioTac; to compensate for manufacturing variation, unlabeled and labeled data from more BioTacs may be necessary to fine-tune the VAEs and projection networks. Future work will focus on applying our simulation and learning framework to non-BioTac sensors. In the process, we aim to develop powerful, generalized representations of tactile data that can serve as the foundation for transfer learning across sensors of entirely different modalities, such as the BioTac and GelSight.

ACKNOWLEDGMENT

We thank V. Makoviychuk, K. Guo, and A. Bakshi for their collaboration with Isaac Gym, as well as K. Van Wyk, A. Handa, and T. Hermans for their feedback.

REFERENCES

- [1] A. Billard and D. Kragic, "Trends and challenges in robot manipulation," *Science*, Jun. 2019.
- [2] SynTouch. [Online]. Available: <https://syntouchinc.com/>
- [3] N. Wettels, V. J. Santos, R. S. Johansson, and G. E. Loeb, "Biomimetic tactile sensor array," *Advanced Robotics*, 2008.
- [4] W. Yuan, S. Dong, and E. H. Adelson, "GelSight: High-resolution robot tactile sensors for estimating geometry and force," *Sensors*, Nov. 2017.
- [5] B. Ward-Cherrier, N. Pestell, L. Cramphorn, B. Winstone, M. E. Giannaccini, J. Rossiter, and N. F. Lepora, "The TacTip family: Soft optical tactile sensors with 3D-printed biomimetic morphologies," *Soft Robotics*, Apr. 2018.
- [6] 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*, Jul. 2020.
- [7] A. Alspach, K. Hashimoto, N. Kuppuswamy, and R. Tedrake, "Soft-bubble: A highly compliant dense geometry tactile sensor for robot manipulation," in *IEEE International Conference on Soft Robotics (RoboSoft)*, Apr. 2019.
- [8] A. Padmanabha, F. Ebert, S. Tian, R. Calandra, C. Finn, and S. Levine, "OmniTact: A multi-directional high resolution touch sensor," in *International Conference on Robotics and Automation (ICRA)*, May 2020.
- [9] C. Sferrazza and R. D'Andrea, "Design, motivation and evaluation of a full-resolution optical tactile sensor," *Sensors*, Feb. 2019.
- [10] A. Yamaguchi and C. G. Atkeson, "Combining finger vision and optical tactile sensing: Reducing and handling errors while cutting vegetables," in *IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, Nov. 2016.
- [11] I. Huang, J. Liu, and R. Bajcsy, "A depth camera-based soft fingertip device for contact region estimation and perception-action coupling," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2019.
- [12] B. W. McInroe, C. L. Chen, K. Y. Goldberg, R. Bajcsy, and R. S. Fearing, "Towards a soft fingertip with integrated sensing and actuation," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2018.
- [13] P. Piacenza, K. Behrman, B. Schifferer, I. Kymissis, and M. Ciocarlie, "A sensorized multicurved robot finger with data-driven touch sensing via overlapping light signals," *IEEE/ASME Transactions on Mechatronics*, Feb. 2020.
- [14] Z. Su, K. Hausman, Y. Chebotar, A. Molchanov, G. E. Loeb, G. S. Sukhatme, and S. Schaal, "Force estimation and slip detection/classification for grip control using a biomimetic tactile sensor," in *IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, Nov. 2015.
- [15] F. Veiga, J. Peters, and T. Hermans, "Grip stabilization of novel objects using slip prediction," *IEEE Transactions on Haptics*, vol. 11, no. 4, 2018.
- [16] J. W. James, N. Pestell, and N. F. Lepora, "Slip detection with a biomimetic tactile sensor," *IEEE Robotics and Automation Letters*, Jul. 2018.
- [17] J. Hoelscher, J. Peters, and T. Hermans, "Evaluation of tactile feature extraction for interactive object recognition," in *IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, 2015.
- [18] W. Yuan, Y. Mo, S. Wang, and E. H. Adelson, "Active clothing material perception using tactile sensing and deep learning," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2018.
- [19] B. Sundaralingam and T. Hermans, "In-hand object-dynamics inference using tactile fingertips," *arXiv preprint arXiv:2003.13165*, Mar. 2020.
- [20] 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)*, May 2019.
- [21] B. Sundaralingam, A. Lambert, A. Handa, B. Boots, T. Hermans, S. Birchfield, N. Ratliff, and D. Fox, "Robust learning of tactile force estimation through robot interaction," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2019.
- [22] Y. Narang, K. V. Wyk, A. Mousavian, and D. Fox, "Interpreting and predicting tactile signals via a physics-based and data-driven framework," in *Robotics: Science and Systems (RSS)*, Jul. 2020.
- [23] Y. Narang, B. Sundaralingam, K. V. Wyk, A. Mousavian, and D. Fox, "Interpreting and predicting tactile signals for the SynTouch Biotac," *arXiv preprint arXiv:2101.05452*, Jan. 2021.
- [24] N. F. Lepora, A. Church, C. de Kerckhove, R. Hadsell, and J. Lloyd, "From pixels to percepts: Highly robust edge perception and contour following using deep learning and an optical biomimetic tactile sensor," *IEEE Robotics and Automation Letters*, Apr. 2019.
- [25] I. Huang and R. Bajcsy, "High resolution soft tactile interface for physical human-robot interaction," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2020.
- [26] A. Afzal, D. S. Katz, C. L. Goues, and C. S. Timperley, "A study on the challenges of using robotics simulators for testing," *arXiv*, Apr. 2020.
- [27] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield, "Deep object pose estimation for semantic robotic grasping of household objects," in *Conference on Robot Learning (CoRL)*, Oct. 2018.
- [28] C. Matl, Y. Narang, R. Bajcsy, F. Ramos, and D. Fox, "Inferring the material properties of granular media for robotic tasks," in *International Conference on Robotics and Automation (ICRA)*, May 2020.
- [29] C. Matl, Y. Narang, D. Fox, R. Bajcsy, and F. Ramos, "STRSSD: Sim-to-real from sound for stochastic dynamics," in *Conference on Robot Learning (CoRL)*, Nov. 2020.
- [30] W. Yu, J. Tan, C. K. Liu, and G. Turk, "Preparing for the unknown: Learning a universal policy with online system identification," in *Robotics: Science and Systems (RSS)*, July 2017.
- [31] Y. Chebotar, A. Handa, V. Makovychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox, "Closing the sim-to-real loop: Adapting simulation randomization with real world experience," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2019.
- [32] I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang, "Solving Rubik's cube with a robot hand," *arXiv preprint arXiv:1910.07113*, Oct. 2019.
- [33] K. Lowrey, S. Kolev, J. Dao, A. Rajeswaran, and E. Todorov, "Reinforcement learning for non-prehensile manipulation: Transfer from simulation to physical system," in *IEEE International Conference on Simulation, Modeling, and Programming for Autonomous Robots (SIMPAP)*, 2018.
- [34] C. Sferrazza, A. Wahlsten, C. Trueeb, and R. D'Andrea, "Ground truth force distribution for learning-based tactile sensing: A finite element approach," *IEEE Access*, Nov. 2019.
- [35] M. M. Zhang, "Necessity for more realistic contact simulation," in *Robotics: Science and Systems (RSS) Workshop on Visuotactile Sensors for Robust Manipulation*, Jul. 2020.
- [36] B. S. Zapata-Impata, P. Gil, Y. Mezouar, and F. Torres, "Generation of tactile data from 3D vision and target robotic grasps," *IEEE Transactions on Haptics*, 2020.
- [37] NVIDIA, "Isaac Gym," <https://developer.nvidia.com/isaac-gym>, 2020.
- [38] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *International Conference on Learning Representations (ICLR)*, Apr. 2014.
- [39] A. Ranjan, T. Bolkart, S. Sanyal, and M. J. Black, "Generating 3D faces using convolutional mesh autoencoders," in *European Conference on Computer Vision (ECCV)*, 2018.
- [40] Q. Li, O. Kroemer, Z. Su, F. F. Veiga, M. Kaboli, and H. J. Ritter, "A review of tactile information: Perception and action through touch," *IEEE Transactions on Robotics*, 2020.
- [41] D. F. Gomes, A. Wilson, and S. Luo, "GelSight simulation for sim2real learning," in *IEEE International Conference on Robotics and Automation (ICRA) Workshop on Integrating Vision and Touch for Multimodal and Cross-modal Perception (ViTac)*, May 2019.
- [42] N. Koenig and A. Howard, "Design and use paradigms for Gazebo, an open-source multi-robot simulator," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sep. 2004.
- [43] 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*, Oct. 2020.
- [44] Z. Ding, N. F. Lepora, and E. Johns, "Sim-to-real transfer for optical tactile sensing," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2020.
- [45] A. Juliani, V.-P. Berges, E. Teng, A. Cohen, J. Harper, C. Elion, C. Goy, H. H. Yuan Gao, M. Mattar, and D. Lange, "Unity: A general

- platform for intelligent agents,” *arXiv preprint arXiv:1809.02627v2*, May 2020.
- [46] B. Wu, I. Akinola, J. Varley, and P. Allen, “MAT: Multi-fingered adaptive tactile grasping via deep reinforcement learning,” in *Conference on Robot Learning (CoRL)*, Oct. 2019.
- [47] Barrett Technology, “BarrettHand,” <https://advanced.barrett.com/barrethand>, 2020.
- [48] E. Coumans and Y. Bai, “PyBullet: A Python module for physics simulation for games, robotics and machine learning,” <http://pybullet.org>, 2016–2020.
- [49] H. Lee, H. Park, G. Serhat, H. Sun, and K. J. Kuchenbecker, “Calibrating a soft ERT-based tactile sensor with a multiphysics model and sim-to-real transfer learning,” in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2020.
- [50] C.-H. Lin, J. A. Fishel, and G. E. Loeb, “Estimating point of contact, force and torque in a biomimetic tactile sensor with deformable skin,” SynTouch LLC, Tech. Rep., 2013.
- [51] ANSYS, Inc., “Ansys Mechanical: Finite element analysis (FEA) software,” <https://www.ansys.com/products/structures/ansys-mechanical>, 2020.
- [52] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, “PointNet++: Deep hierarchical feature learning on point sets in a metric space,” *arXiv preprint arXiv:1706.02413*, Jun. 2017.
- [53] P. Ruppel, Y. Jonetzko, M. Gormer, N. Hendrich, and J. Zhang, “Simulation of the SynTouch BioTac sensor,” in *International Conference on Intelligent Autonomous Systems (IAS)*, Jun. 2018.
- [54] B. S. Zapata-Impata and P. Gil, “Prediction of tactile perception from vision on deformable objects,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) Workshop on Robotic Manipulation of Deformable Objects (ROMADO)*, 2020.
- [55] R. Calandra, A. Owens, M. Upadhyaya, W. Yuan, J. Lin, E. H. Adelson, and S. Levine, “More than a feeling: Learning to grasp and regrasp using vision and touch,” *Conference on Robot Learning*, Nov. 2017.
- [56] R. Calandra, D. J. Andrew Owens, J. Lin, W. Yuan, J. Malik, E. H. Adelson, and S. Levine, “More than a feeling: Learning to grasp and regrasp using vision and touch,” *IEEE Robotics and Automation Letters*, Oct. 2018.
- [57] N. Fazeli, M. Oller, J. Wu, Z. Wu, J. B. Tenenbaum, and A. Rodriguez, “See, feel, act: Hierarchical learning for complex manipulation skills with multisensory fusion,” *Science Robotics*, Jan. 2019.
- [58] M. A. Lee, Y. Zhu, K. Srinivasan, P. Shah, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg, “Making sense of vision and touch: Self-supervised learning of multimodal representations for contact-rich tasks,” in *International Conference on Robotics and Automation (ICRA)*, 2019.
- [59] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, “Multimodal deep learning,” in *International Conference on Machine Learning (ICML)*, 2011.
- [60] L. Chen, S. Srivastava, Z. Duan, and C. Xu, “Deep cross-modal audio-visual generation,” in *Thematic Workshops of ACM Multimedia*, Oct. 2019.
- [61] J.-T. Lee, D. Bollegala, and S. Luo, ““Touching to see” and “seeing to feel”: Robotic cross-modal sensory data generation for visual-tactile perception,” in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2019.
- [62] Y. Li, J.-Y. Zhu, R. Tedrake, and A. Torralba, “Connecting touch and vision via cross-modal prediction,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019.
- [63] K. Gutierrez and V. J. Santos, “Perception of tactile directionality via artificial fingerpad deformation and convolutional neural networks,” *IEEE Transactions on Haptics*, 2020.
- [64] A. Garcia-Garcia, B. S. Zapata-Impata, S. Orts-Escolano, P. Gil, and J. Garcia-Rodriguez, “TactileGCN: A graph convolutional network for predicting grasp stability with tactile sensors,” in *International Joint Conference on Neural Networks (IJCNN)*, 2019.
- [65] J. N. Reddy, *Introduction to the Finite Element Method*. McGraw Hill, 2019.
- [66] Y. S. Narang, J. J. Vlassak, and R. D. Howe, “Mechanically versatile soft machines through laminar jamming,” *Advanced Functional Materials*, 2018.
- [67] D. E. Stewart, “Rigid-body dynamics with friction and impact,” *SIAM Review*, 2000.
- [68] P. Cignoni, M. Callieri, M. Corsini, M. Dellepiane, F. Ganovelli, and G. Ranzuglia, “MeshLab: An open-source mesh processing tool,” in *Eurographics Italian Chapter Conference*, 2008.
- [70] Y. Hu, T. Schneider, B. Wang, D. Zorin, and D. Panozzo, “Fast tetrahedral meshing in the wild,” *ACM Transactions on Graphics*, Jul. 2020.
- [71] K. Hauser, “Robust contact generation for robot simulation with unstructured meshes,” *Springer Tracts in Advanced Robotics*, 2016.
- [72] R. W. Ogden, *Non-Linear Elastic Deformations*. Courier Corporation, 2013.