

Interactive Visuo-Tactile Learning to Estimate Properties of Articulated Objects

Anonymous Author(s)

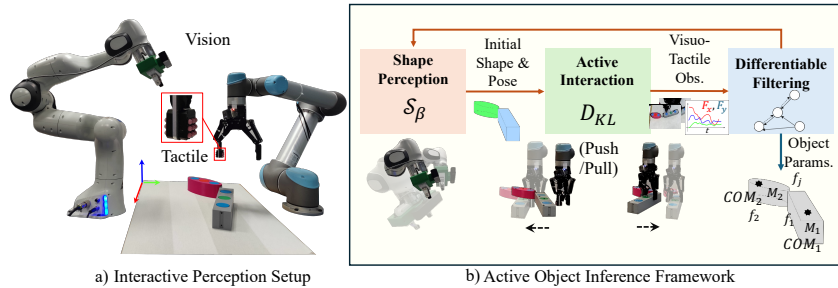
Affiliation

Address

email

1 **Abstract:** Robotic systems operating in unstructured environments must infer
2 key physical properties of objects, such as stiffness, mass, center of mass, friction,
3 and shape, to ensure stable manipulation. Accurate estimation of these properties
4 is crucial for predicting and effective planning manipulation outcomes. In this
5 work, we present a novel framework for identifying the properties of challenging
6 objects which are articulated through versatile, non-prehensile push-pull actions
7 and using visuo-tactile observation. Our approach introduces a differentiable
8 filtering method that incorporates embedding interaction physics into graph neural
9 networks, enabling the system to actively learn object-robot interactions and con-
10 sistently infer both directly observable pose information and indirectly observable
11 physical parameters. Experimental results on real robotic systems show that our
12 method outperforms existing baselines in efficiency and accuracy.

13 **Keywords:** Perception for Grasp & Manipulation, Visuo-Tactile Sensing, Active
14 Learning, Interaction Dynamics



15 1 Introduction

16 Robotic systems engaged in contact-rich object manipulation tasks need to perceive the physical
17 properties of the object, such as mass, center of mass, and surface friction, to perform effectively.
18 However, these properties are difficult to estimate, as they are not directly observable in static en-
19 vironments and become salient only during specific object-robot interactions [1]. Current visual or
20 tactile perception frameworks struggle to handle previously unseen objects [2, 3], necessitating the
21 use of simple and robust interaction strategies to infer these physical properties prior to manipulation
22 tasks [4, 5]. In this study, we propose a novel interactive learning and perception framework for
23 inferring the properties of articulated objects using both vision and tactile sensing seamlessly using
24 versatile push-pull interactions.

25 2 Related Work

26 Estimating inertial and surface properties of rigid objects is a long-standing problem in control the-
27 ory, particularly for rigid body identification [6, 7]. The early methods relied on rigidly attached

28 objects to manipulators [8, 9], limiting their applicability in unstructured environments due to spe-
 29 cialized mechanisms and known object geometry. Interactive techniques like grasping or pushing
 30 [10, 11, 12] tried to overcome these issues but relied on simplified assumptions. Recent research ex-
 31 plores data-driven [13, 14, 15] and physics-based approaches [16], with studies [17, 18] showing
 32 the potential of graph networks to capture object-robot interactions. However, current GNNs fo-
 33 cus on spatial relationships and kinematics but fail to capture contact forces influenced by physical
 34 properties and robot actions. This highlights the need for a graph-based model incorporating tactile
 35 information with a stronger inductive bias. Moreover, existing data-driven methods require exten-
 36 sive training and often lack strategic interaction, limiting their use to simulation environments. This
 37 motivates us to investigate possible active/informative interaction techniques [19, 20, 21], which is
 38 addressed in this work. Furthermore, prior research has relied mainly on visual [1, 22] or tactile
 39 [12, 11] methods to estimate physical properties, each with limitations. Tactile sensing can infer
 40 multiple properties of objects, but requires precise information and prior knowledge, while vision
 41 offers a limited range of observable properties, but provides a global view of the shape and move-
 42 ment of an object. Recent works [23, 24] combining vision and tactile approaches have shown
 43 improvements in pose estimation and contact-rich manipulation tasks. Building on these advances,
 44 our framework integrates sensing modalities with active exploratory actions: *non-prehensile push-*
 45 *ing*, and *prehensile pulling* to enhance object perception. By encoding object-robot interactions into
 46 Probabilistic Markov Models and using a learned interaction model (differentiable filter [25, 22]),
 47 our system predicts visuo-tactile observations and estimates key physical parameters of articulated
 48 objects in a Bayesian Inference setting. The learned models capture not only the complex interaction
 49 dynamics but also modality-specific noise, improving the efficiency of inference.

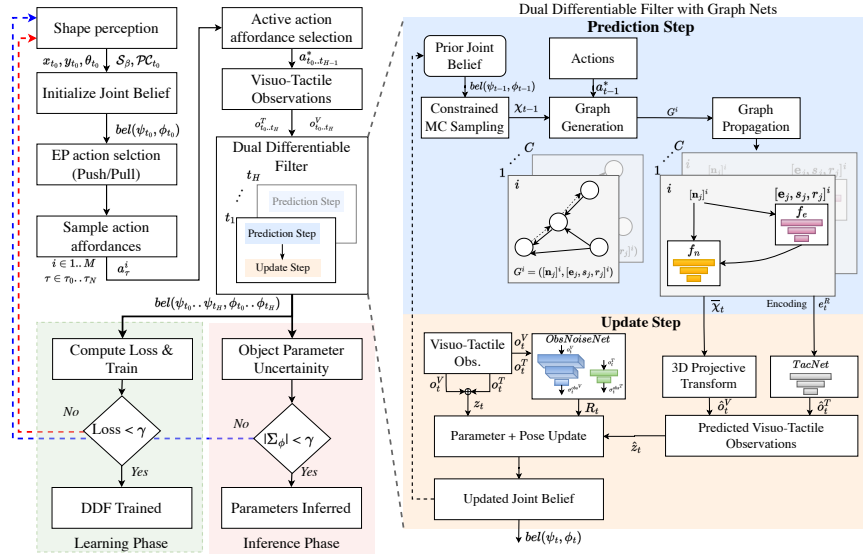


Figure 1: Proposed framework for interactively learning & inferring the properties of articulated objects using visuo-tactile sensing.

50 3 Methods

51 3.1 Problem Formulation

52 We tackle estimating the state s of an unknown rigid object on a support surface using visual (o^V)
 53 and tactile (o^T) inputs along with actions (a). The objects are articulated, with multiple links con-
 54 nected through rotational joints. At time t , the state of the object s_t is composed of $l \in 1, \dots, L$
 55 links, expressed as $s_t = \{s_t^1 \dots s_t^L\}$. The state of each link l , $s_t^l = \{\psi_t^l, \phi^l\}$, includes time-varying
 56 elements: the 2D *pose* and *twist*, $\psi_t^l = \{x_t, y_t, \theta_t, v_{x_t}, v_{y_t}, \omega_t\}$, and time-invariant elements ϕ^l ,
 57 involving *inertial parameters* like $\{m, CoM_x, CoM_y\}$ mass and center of mass vector, and *interac-*
 58 *tion parameters* $\{f, f_r, f_j\}$ for friction with the table, robot, and adjacent link. The visual data o_t^V

59 includes RGB-D images of the robot-object interface, while the tactile data o_t^T is 2D contact forces
60 from the robotic gripper’s interaction (fingertip forces). The push/pull action is defined by the tuple
61 *contact point* (cp), *direction* (pd) and *velocity* (u). In addition, for autonomous and seamless explo-
62 ration of the object, the shape of each link \mathcal{S}^l is estimated via superquadrics [26]. The belief about
63 the current state of the object s_t is represented by a distribution conditioned on previous actions $a_{1:t}$
64 and observations $o_{1:t}$ and employs recursive Bayesian filtering.

$$bel(s_t) = p(s_t|o_{1:t}, a_{1:t}) = \eta p(o_t|s_t, a_t) \int p(s_t|s_{t-1}, a_{t-1}) bel(s_{t-1}) ds_{t-1} \quad (1)$$

65 where η is a normalizing factor. We employ a data-driven strategy to learn the process, observation,
66 and noise models. Since object pose intricately relies on inertial and interaction parameters, joint
67 filtering for pose and parameters [27] is found to be ineffective and we adopt a dual filter design to
68 maintain consistent filtering and infer object parameters.

69 3.2 Dual Differentiable Filter

70 For the dual filter formulation, we explicitly represent the state of the object (joint distribution of
71 pose and twist) via Multivariate Gaussian distribution:

$$bel(\psi_t, \phi_t) \doteq \mathcal{N}(\psi_t, \phi_t | \mu_t, \Sigma_t), \quad \mu_t = \begin{pmatrix} \mu_{\psi_t} \\ \mu_{\phi_t} \end{pmatrix}, \quad \Sigma_t = \begin{pmatrix} \Sigma_{\psi_t} & \Sigma_{\psi_t \phi_t} \\ \Sigma_{\phi_t \psi_t} & \Sigma_{\phi_t} \end{pmatrix} \quad (2)$$

72 with dimensions $\mu_t \in \mathbb{R}^{11L-1}$ and $\Sigma_t \in \mathbb{R}^{(11L-1) \times (11L-1)}$. The dual filter as shown in Fig.1
73 follows the structure of a Kalman filter with a *prediction step* and an *update step*, with the proposed
74 novelty explained in this section.

75 3.2.1 Prediction Step

76 In the prediction step, the next joint belief is predicted based on the prior belief and actions. Since
77 the object’s inertial and interaction parameters have physical constraints (e.g., $m, f, f_j > 0$, CoM_x ,
78 CoM_y must lie within the object boundary), constrained Monte Carlo sigma point sampling is per-
79 formed to maintain these constraints and the Gaussian variance. A differentiable sampling method
80 [28] is used to sample C sigma points $\chi_{t-1}^i, i = \{1..C\}$ from the joint distribution $bel(\psi_{t-1}, \phi_{t-1})$,
81 with an associated weight $w_{t-1}^i = 1/C$.

82 We employ Graph Neural Networks (GNNs) to model the interaction between the object, the support
83 surface, and the robot. Using the sigma points χ_{t-1}^i and the robot action a_{t-1} , a directed graph
84 $G_t^i = (\{\mathbf{n}_l\}, \{e_j, s_j, r_j\})$ is constructed, where \mathbf{n}_l represents the nodes for each link of the object,
85 the robot and the support surface, and e_j represents the directed edges. Each node $\mathbf{n}_l \in \mathbb{R}^{L+2}$
86 contains features including dynamic (pose, twist) and static (inertial) parameters, populated from
87 the sigma points for the object links, with default values for the robot and surface. The edges
88 $e_j \in \mathbb{R}^{3L}$ capture the interaction between the links between the objects, the robot and the support
89 surface, with features such as friction coefficients. To update node and edge features from time $t-1$
90 to t , we use a novel graph propagation algorithm (see the Appendix) with two functions: f_n for node
91 updates and f_e for edge updates.

92 3.2.2 Update Step

93 The dual filter employs a separate update of the parameter belief similar to the parameter update
94 presented in [29] and the conditional pose belief update based on the UKF update [30]. To reduce
95 the complexity of predicting raw RGB-D images, we use the initial segmented point cloud \mathcal{PC}_{t_0}
96 from the shape perception method to transform it using the predicted pose and generate expected
97 RGB-D images using the standard 3D to 2D projective transformation approach [31] involving
98 the intrinsic and extrinsic values of the camera, also avoiding generalization issues. For the tactile
99 counterpart, a three-layer feedforward network is utilized to predict the contact force information
100 from the edge encoding directed towards the robot. The filtering step is used end-to-end for both
101 learning and inference.

102 3.2.3 Active Interaction: N -step Information Gain

103 To make the framework more sample efficient for real robot scenarios, we employ active action
 104 selection by formulating an N -step information gain criteria [32] under the filtering setting. We
 105 recursively use the prediction step of the dual differentiable filter without the update step to compute
 106 the expected Information Gain for both model learning and object parameter inference for each
 107 sampled non-prehensile pushing or prehensile pulling action $\pi^{[i]} = a_{\tau_0:\tau_N}^i$ over N -step in future
 108 $\tau = \tau_0 \dots \tau_N$

$$IG_N(\pi^{[i]}) \approx -\mathbb{E}_{p(\psi_{\tau_N}, \phi_{\tau_N} | \pi^{[i]})} [\ln(\overline{bel}^{[i]}(\psi_{\tau_N}, \phi_{\tau_N})) - \ln(\overline{bel}^{[i]}(\psi_{\tau_0}, \phi_{\tau_0}))] \quad (3)$$

109 where, $\overline{bel}^{[i]}(\psi_{\tau_N}, \phi_{\tau_N})$ is the hypothetical predictive joint distribution after N -step by taking action
 110 $\pi^{[i]}$ without taking account the actual observation. At every step $\pi^* = \arg \max_{\pi^i} IG_N(\pi^{[i]})$ is
 111 selected for interaction.

112 4 Results & Conclusion

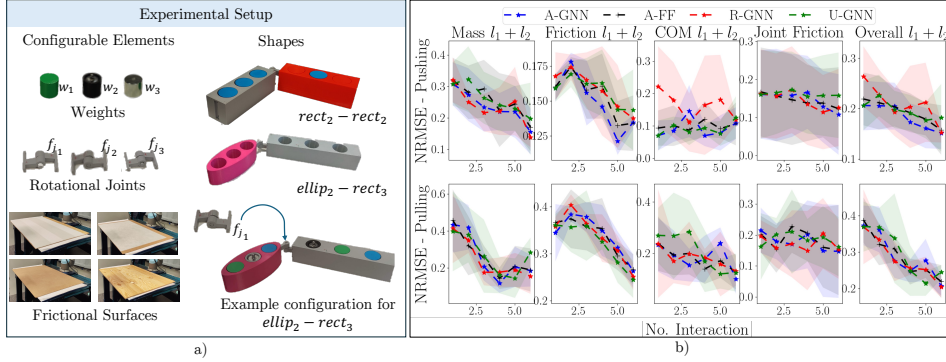


Figure 2: A) Experimental Setup with configurable objects B) Parameter estimation error across multiple interactions for articulated objects comparing proposed $A-GNN$ with $A-FF$, $R-GNN$ and $U-GNN$

113 We compare our proposed method, $A-GNN$, with the baseline $A-FF$ from
 114 [14] and conducted ablation studies to evaluate active action selection against
 115 uniform ($U-GNN$) and random ($R-GNN$) selection for model learning and
 116 inference. We designed 60 3D-printed articulated objects by varying weights,
 117 frictional surfaces, and joint friction
 118 (Fig.2a)). The networks were trained using negative log-likelihood (\mathcal{L}_{NLL}), mean squared error
 119 (\mathcal{L}_{MSE}), and observed noise log-likelihood (\mathcal{L}_{NLL}^{obs}). \mathcal{L}_{NLL} and \mathcal{L}_{MSE} compared ground truth and
 120 predicted poses, parameters, and forces, while \mathcal{L}_{NLL}^{obs} was calculated using predicted and real obser-
 121 vations. To account for different inertial and interaction parameters, we used normalized root mean
 122 squared error (NRMSE) to report estimation errors. Table 1.A shows that $A-GNN$ with active
 123 action selection improved data efficiency by 25% over uniform selection and 9% over $A-FF$,
 124 particularly for complex articulated objects and push interactions. Moreover, inference accuracy
 125 remains consistent with low SD and surpasses baseline methods (see Fig.2.B and Table 1.B).

	Pulling		Pushing	
	A	B	A	B
A-GNN	2830	0.21 ± 0.02	4390	0.15 ± 0.03
A-FF	3410	0.22 ± 0.04	4810	0.16 ± 0.05
R-GNN	2922	0.21 ± 0.03	5295	0.15 ± 0.02
U-GNN	3405	0.25 ± 0.02	6000	0.18 ± 0.07

Table 1: Col. A) presents the no. of interactions required for training convergence, and B) presents $NRMSE$ of the overall parameter inference

130 Although this study considers a single object in isolation, future work will address more complex
 131 clutter scenarios and include interactive perception for prismatic-rotational joint identification. We
 132 also assumed the objects were planar and each articulated link was at most connected by two links.
 133 In conclusion, the proposed novel framework enables the robotic system to estimate the properties
 134 of intricate articulated objects autonomously using simple and efficient (active) interactive actions:
 135 non-prehensile push and prehensile pull.

References

- [1] Z. Xu, J. Wu, A. Zeng, J. B. Tenenbaum, and S. Song. Densephysnet: Learning dense physical object representations via multi-step dynamic interactions. *arXiv preprint arXiv:1906.03853*, 2019.
- [2] N. Navarro-Guerrero et al. Visuo-haptic object perception for robots: an overview. *AuRo*, 47(4):377–403, 2023.
- [3] Q. Li et al. A review of tactile information: Perception and action through touch. *IEEE Transactions on Rob.*, 36(6):1619–1634, 2020. doi:10.1109/TRO.2020.3003230.
- [4] P. Agrawal, A. V. Nair, P. Abbeel, J. Malik, and S. Levine. Learning to poke by poking: Experiential learning of intuitive physics. *Advances in neural information processing systems*, 29, 2016.
- [5] J. Bohg, K. Hausman, B. Sankaran, O. Brock, D. Kragic, S. Schaal, and G. S. Sukhatme. Interactive perception: Leveraging action in perception and perception in action. *IEEE Transactions on Robotics*, 33(6):1273–1291, 2017. doi:10.1109/TRO.2017.2721939.
- [6] L. Ljung. *System Identification*, pages 163–173. Birkhäuser Boston, Boston, MA, 1998. ISBN 978-1-4612-1768-8. doi:10.1007/978-1-4612-1768-8_11. URL https://doi.org/10.1007/978-1-4612-1768-8_11.
- [7] M. Niebergall and H. Hahn. Identification of the ten inertia parameters of a rigid body. *Non-linear Dynamics*, 13:361–372, 1997.
- [8] C. G. Atkeson, C. H. An, and J. M. Hollerbach. Estimation of inertial parameters of manipulator loads and links. *The International Journal of Robotics Research*, 5(3):101–119, 1986.
- [9] C. Wang, X. Zang, X. Zhang, Y. Liu, and J. Zhao. Parameter estimation and object gripping based on fingertip force/torque sensors. *Measurement*, 179:109479, 2021.
- [10] S. Tanaka, T. Tanigawa, Y. Abe, M. Uejo, and H. T. Tanaka. Active mass estimation with haptic vision. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, volume 3, pages 256–261. IEEE, 2004.
- [11] K. Yao, M. Kaboli, and G. Cheng. Tactile-based object center of mass exploration and discrimination. In *2017 IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)*, pages 876–881, 2017. doi:10.1109/HUMANOIDS.2017.8246975.
- [12] B. Sundaralingam and T. Hermans. In-hand object-dynamics inference using tactile fingertips. *IEEE Transactions on Robotics*, 37(4):1115–1126, 2021.
- [13] C. Song and A. Boularias. A probabilistic model for planar sliding of objects with unknown material properties: Identification and robust planning. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5311–5318. IEEE, 2020.
- [14] A. Dutta, E. Burdet, and M. Kaboli. Push to know!-visuo-tactile based active object parameter inference with dual differentiable filtering. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3137–3144. IEEE, 2023.
- [15] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia. Learning to simulate complex physics with graph networks. In *International conference on machine learning*, pages 8459–8468. PMLR, 2020.
- [16] J. Wu, I. Yildirim, J. J. Lim, B. Freeman, and J. Tenenbaum. Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. *Advances in neural information processing systems*, 28, 2015.

- 179 [17] F. Paus, T. Huang, and T. Asfour. Predicting pushing action effects on spatial object relations
180 by learning internal prediction models. In *2020 IEEE International Conference on Robotics
181 and Automation (ICRA)*, pages 10584–10590, 2020. doi:10.1109/ICRA40945.2020.9197295.
- 182 [18] A. E. Tekden, A. Erdem, E. Erdem, T. Asfour, and E. Ugur. Object and relation centric repre-
183 sentations for push effect prediction. *Robotics and Autonomous Systems*, page 104632, 2024.
- 184 [19] J. Xu, H. Lin, S. Song, and M. Ciocarlie. Tandem3d: Active tactile exploration for 3d object
185 recognition. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*,
186 pages 10401–10407. IEEE, 2023.
- 187 [20] T. Parr, G. Pezzulo, and K. J. Friston. *Active inference: the free energy principle in mind,
188 brain, and behavior*. MIT Press, 2022.
- 189 [21] M. Kaboli, K. Yao, D. Feng, and G. Cheng. Tactile-based active object discrimination and
190 target object search in an unknown workspace. *Autonomous Robots*, 43:123–152, 2019.
- 191 [22] A. Kloss, M. Bauza, J. Wu, J. B. Tenenbaum, A. Rodriguez, and J. Bohg. Accurate vision-
192 based manipulation through contact reasoning. In *2020 IEEE International Conference on
193 Robotics and Automation (ICRA)*, pages 6738–6744. IEEE, 2020.
- 194 [23] P. K. Murali, A. Dutta, M. Gentner, E. Burdet, R. Dahiya, and M. Kaboli. Active visuo-
195 tactile interactive robotic perception for accurate object pose estimation in dense clutter. *IEEE
196 Robotics and Automation Letters*, 7(2):4686–4693, 2022. doi:10.1109/LRA.2022.3150045.
- 197 [24] M. A. Lee, Y. Zhu, K. Srinivasan, P. Shah, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg.
198 Making sense of vision and touch: Self-supervised learning of multimodal representations for
199 contact-rich tasks. In *2019 International Conference on Robotics and Automation (ICRA)*,
200 pages 8943–8950. IEEE, 2019.
- 201 [25] R. Jonschkowski, D. Rastogi, and O. Brock. Differentiable particle filters: End-to-end learning
202 with algorithmic priors. *arXiv preprint arXiv:1805.11122*, 2018.
- 203 [26] W. Liu, Y. Wu, S. Ruan, and G. S. Chirikjian. Robust and accurate superquadric recovery: A
204 probabilistic approach. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
205 Pattern Recognition*, pages 2676–2685, 2022.
- 206 [27] A. Kloss, G. Martius, and J. Bohg. How to train your differentiable filter. *Autonomous Robots*,
207 45(4):561–578, 2021.
- 208 [28] M. Wüthrich, C. G. Cifuentes, S. Trimpe, F. Meier, J. Bohg, J. Issac, and S. Schaal. Robust
209 gaussian filtering using a pseudo measurement. In *2016 American Control Conference (ACC)*,
210 pages 3606–3613. IEEE, 2016.
- 211 [29] J. Liu and M. West. Combined parameter and state estimation in simulation-based filtering. In
212 *Sequential Monte Carlo methods in practice*, pages 197–223. Springer, 2001.
- 213 [30] S. Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002.
- 214 [31] I. Nematollahi, E. Rosete-Beas, S. M. B. Azad, R. Rajan, F. Hutter, and W. Burgard. T3vip:
215 Transformation-based 3D video prediction. In *2022 IEEE/RSJ International Conference on
216 Intelligent Robots and Systems (IROS)*, pages 4174–4181, 2022. doi:10.1109/IROS47612.
217 2022.9981187.
- 218 [32] K. Friston. The free-energy principle: a unified brain theory? *Nature reviews neuroscience*, 11
219 (2):127–138, 2010.

220 **5 Appendix**

221 **5.1 Graph Propagation Algorithm**

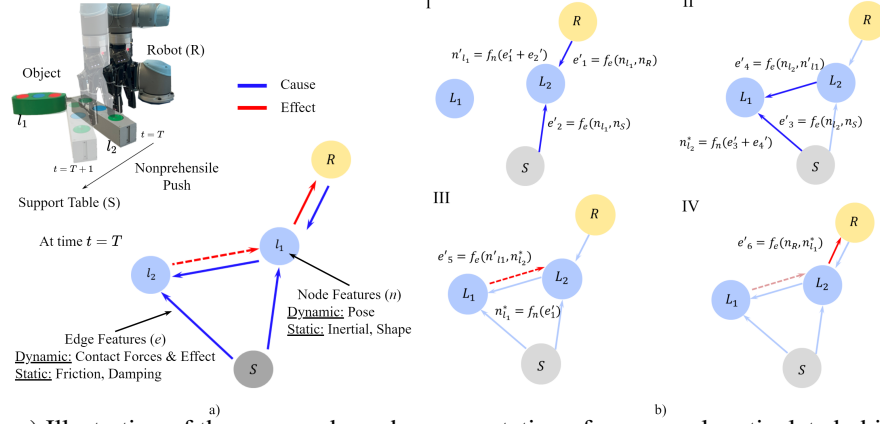


Figure 3: a) Illustration of the proposed graph representation of an example articulated object with two links b) Novel graph propagation for updating the graphical model from time $t - 1$ to t for the example object. The support edges e_1, e_2 , the edge e_6 contains contact force or tactile information

Algorithm 1 Graph Propagation Algorithm (*GP*)

Input: Graph $G_{t-1} = (\{\mathbf{n}_i\}, \{\mathbf{e}_j, s_j, r_j\})$

Initialize Stacks (LIFO)

$NTV \leftarrow \mathbf{n}_R$

$NV \leftarrow \emptyset$

$EN \leftarrow \emptyset$

▷ Nodes to visit

▷ Nodes visited

▷ End nodes

Propagate cause

while do $NTV \neq \emptyset$

$\mathbf{n}_i = \text{Pop } NTV$

$\mathbf{n}_{r_j} = \text{Gather receiver nodes of } \mathbf{n}_i$

$\mathbf{n}_{r_j} = \mathbf{n}_{r_j} \setminus NV$

▷ Remove nodes already visited

if $\mathbf{n}_{r_j} \neq \emptyset$ **then**

 Push $\mathbf{n}_i \rightarrow NV$

 Push $\mathbf{n}_{r_j} \rightarrow NTV$

for each node \mathbf{n}_{r_j} **do**

 Compute causal edges, $\mathbf{e}_j^* = f_e(\mathbf{n}_i, \mathbf{n}_{r_j}, \mathbf{e}_{s_j})$

 ▷ \mathbf{e}_{s_j} is static edge feature (friction values)

 Compute support edges, $\mathbf{e}_k^* = f_e(\mathbf{n}_S, \mathbf{n}_{r_j}, \mathbf{e}_{s_k})$

 Compute node features, $\mathbf{n}_i^* = f_n(\mathbf{n}_i, \mathbf{e}_j^* + \mathbf{e}_k^*)$

end for

else

 Push $\mathbf{n}_i \rightarrow EN$

end if

end while

Propagate effect

while do $NV \neq \emptyset$

$\mathbf{n}_i = \text{Pop } NV$

$\mathbf{n}_{s_j}^* = \text{Gather sender nodes of } \mathbf{n}_i$

$\mathbf{n}_{s_j}^* = \mathbf{n}_{s_j}^* \setminus NV$

 Aggregate effect edges, $\mathbf{e}_j^* = f_e(\mathbf{n}_{s_j}^*, \mathbf{n}_i, \mathbf{e}_{s_j})$

 Update node features, $\mathbf{n}_i^* = f_n(\mathbf{n}_i, \sum_{j/s_j} \mathbf{e}_j^*)$

end while

Output: Graph $G_t = (\{\mathbf{n}_i^*\}, \{\mathbf{e}_j^*, s_j, r_j\})$
