

Dynamic-Resolution Model Learning for Object Pile Manipulation

Anonymous CVPR submission

Paper ID 17

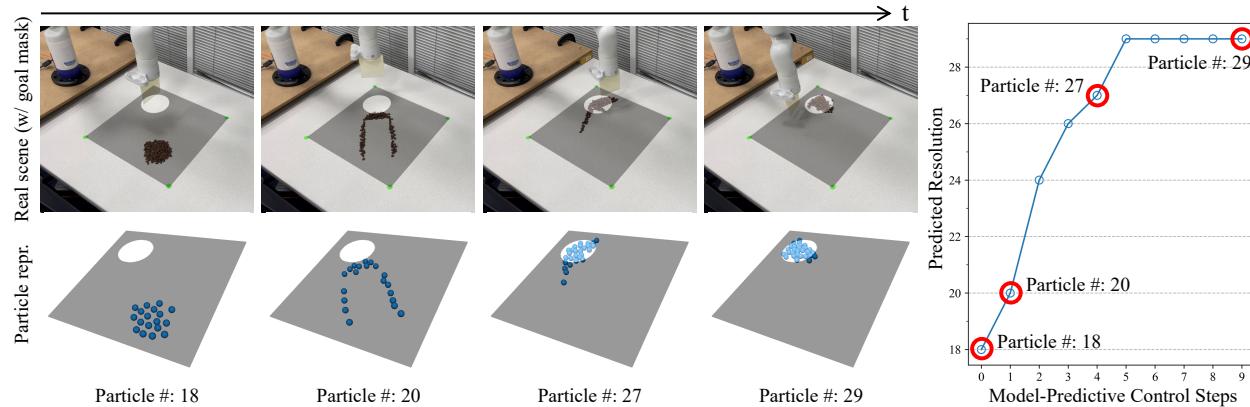


Figure 1. **Dynamic-Resolution Model Learning for Object Pile Manipulation in the Real World.** Depending on the progression of a task, representations at different granularity levels may be needed at each model-predictive control (MPC) step to make the most effective progress on the overall task. In this work, we construct dynamic-resolution particle representations of the environment and learn a *unified* dynamics model using graph neural networks (GNNs) that allows adaptive selection of the abstraction level. In this figure, we demonstrate a real-world task of gathering the object pile into a target region. Figures on the left show the task execution process and the corresponding particle representation. The plot on the right shows the predicted optimal resolution at each MPC step, where the red circles correspond to the frames on the left.

Abstract

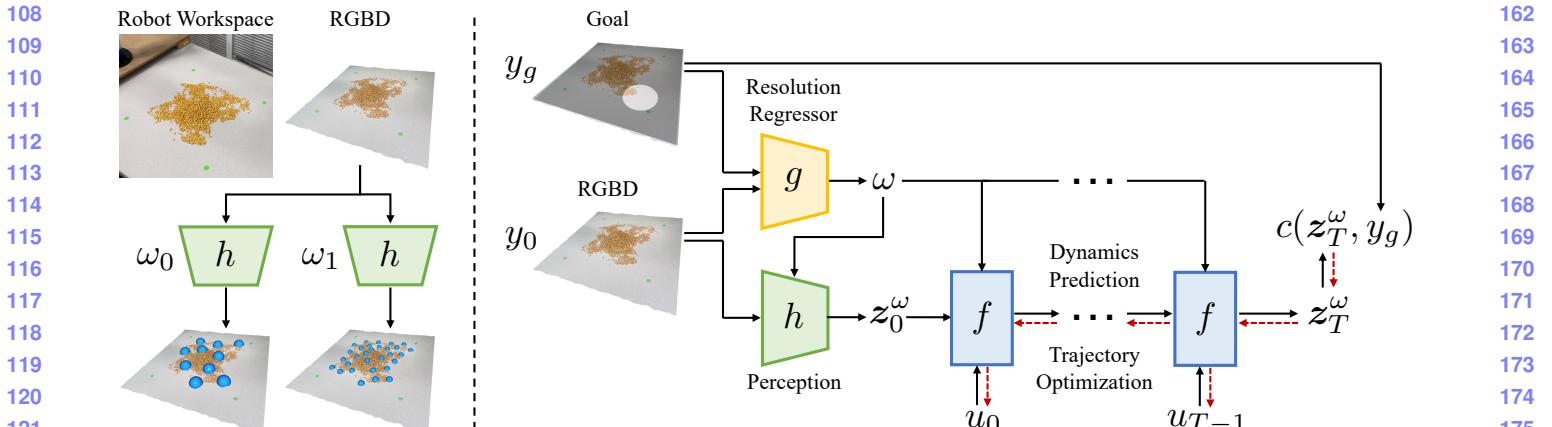
Dynamics models learned from visual observations have shown to be effective in various robotic manipulation tasks. One of the key questions for learning such dynamics models is what scene representation to use. Prior works typically assume representation at a fixed dimension or resolution, which may be inefficient for simple tasks and ineffective for more complicated tasks. In this work, we investigate how to learn dynamic and adaptive representations at different levels of abstraction to achieve the optimal trade-off between efficiency and effectiveness. Specifically, we construct dynamic-resolution particle representations of the environment and learn a unified dynamics model using graph neural networks (GNNs) that allows continuous selection of the abstraction level. During test time, the agent can adaptively determine the optimal resolution at each model-predictive control (MPC) step. We evaluate our method in object pile manipulation, a task we commonly encounter in cooking, agriculture, manufacturing, and pharmaceutical applications. Through comprehensive evaluations both in

the simulation and the real world, we show that our method achieves significantly better performance than state-of-the-art fixed-resolution baselines at the gathering, sorting, and redistribution of granular object piles made with various instances like coffee beans, almonds, corn, etc.

1. Introduction

Predictive models are core to robotic systems for navigation [25], locomotion [28], and manipulation [22, 71]. In robotic manipulation, learned dynamics models have demonstrated impressive results. A learning-based dynamics model includes an encoder and a predictive model. Scene representation choices (e.g., latent vectors [19, 20, 33], object-centric [15, 69] or keypoint representations [37, 40, 66]) affect expressiveness and generalization capabilities, which makes it crucial for a given task.

Prior work uses a fixed representation for the entire task, but the optimal representation may differ depending on the object, task, or stage. An ideal representation balances efficiency and effectiveness [5, 61]. For instance, in object pile



(a) Representations at different resolutions

Figure 2. Overview of the proposed framework. (a) Our perception module h processes the input RGBD image and generates particle representations at different levels of abstraction depending on the resolution ω . (b) The resolution regressor g takes the current observation y_0 and the goal y_g as input. It then predicts the resolution ω we intend to represent the environment. The dynamics model f , conditioned on the dynamically-selected resolution ω and the input action u_t , predicts the temporal evolution of the scene representation z_t^ω . During planning time, we calculate the task objective $c(z_T^\omega, y_g)$ and backpropagate the gradients to optimize the action sequence $\{u_t\}$.

manipulation, a more complex target configuration needs a finer model to capture all the details. While for the same targets, we might want representations at different abstraction levels for the most effective actions at different stages, as shown in Figure 1.

We focus on manipulating object piles, a crucial task in cooking, agriculture, manufacturing, and pharmaceutical scenarios. This task is highly challenging due to the environment’s extremely high degrees of freedom [52], making it an ideal scenario to demonstrate how we can learn dynamics models at different levels of abstraction to achieve the optimal trade-off between efficiency and effectiveness.

Our aim is to learn a dynamics model that can adaptively express the world at different granularity levels based on the task objective and observation. To achieve this, we introduce a resolution regressor that predicts the optimal resolution using self-supervised learning with labels from Bayesian optimization [18]. Besides the resolution regressor, our model also includes perception, dynamics, and planning modules (Figure 2).

During task execution, we follow a model-predictive control (MPC) framework. At each MPC step, the resolution regressor predicts the resolution most effective for control optimization. The perception module then samples particles from the RGBD visual observation based on the predicted resolution. The derived particle-based scene representation, together with the robot action, will be the input to the dynamics model to predict the environment’s evolution. The dynamics model can then be used for trajectory optimization to derive the action sequence. Specifically, the dynamics model is instantiated as a graph neural network consisting of node and edge encoders. Such compositional structures naturally generalize to particle sets of different

sizes and densities—a unified graph-based dynamics model can support model-predictive control at various abstraction levels, selected continuously by the resolution regressor.

Our contributions are threefold: (1) a framework that dynamically determines the scene representation at different abstraction levels, (2) comprehensive evaluations showing the superiority of our dynamic scene representation over fixed resolution, and (3) a unified robotic manipulation system for various object pile manipulation tasks.

2. Method

In this section, we first present the problem formulation in Section A.1. We then discuss the structure of our dynamic-resolution dynamics models, how we learn a resolution regressor to automatically select the scene representation, and how we use the model for the downstream planning tasks in Section A.2, A.3 and A.4 respectively.

3. Experiments

In this section, we evaluate the proposed framework in various object pile manipulation tasks. In particular, we aim to answer the following three questions through the experiments. (1) Does a trade-off exist between efficiency and effectiveness as we navigate through representations at different abstraction levels? (2) Is a fixed-resolution dynamics model sufficient, or do we need to dynamically select the resolution at each MPC step? (3) Can our dynamic-resolution model accomplish three challenging object pile manipulation tasks: **Gather**, **Redistribute**, and **Sort**? Experiments details can be found in Section B

216 **References** 270
217

- 218 [1] Edward H Adelson, Charles H Anderson, James R Bergen, Peter J Burt, and Joan M Ogden. Pyramid 271 methods in image processing. *RCA engineer*, 29(6): 272 33–41, 1984. 12 273
- 219 [2] Andrew G Barto and Sridhar Mahadevan. Recent 274 advances in hierarchical reinforcement learning. *Discrete 275 event dynamic systems*, 13(1-2):41–77, 2003. 12 276
- 220 [3] Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo 277 Jimenez Rezende, et al. Interaction networks for learning 278 about objects, relations and physics. *Advances in neural 279 information processing systems*, 29, 2016. 12 280
- 221 [4] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, 281 Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz 282 Malinowski, Andrea Tacchetti, David Raposo, Adam 283 Santoro, Ryan Faulkner, et al. Relational inductive 284 biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018. 12 285
- 222 [5] Yoshua Bengio, Aaron Courville, and Pascal Vincent. 286 Representation learning: A review and new perspectives. 287 *IEEE transactions on pattern analysis and machine 288 intelligence*, 35(8):1798–1828, 2013. 1 289
- 223 [6] Matthew Michael Botvinick. Hierarchical reinforcement 290 learning and decision making. *Current opinion in neurobiology*, 22(6):956–962, 2012. 12 291
- 224 [7] Peter J Burt and Edward H Adelson. The laplacian 292 pyramid as a compact image code. In *Readings in computer 293 vision*, pages 671–679. Elsevier, 1987. 12 294
- 225 [8] Eduardo F Camacho and Carlos Bordons Alba. *Model 295 predictive control*. Springer science & business media, 296 2013. 8 297
- 226 [9] Michael B Chang, Tomer Ullman, Antonio Torralba, 298 and Joshua B Tenenbaum. A compositional object-based 299 approach to learning physical dynamics. *arXiv preprint arXiv:1612.00341*, 2016. 12 300
- 227 [10] Peng Chang and Taşkin Padır. Model-based manipulation 301 of linear flexible objects with visual curvature feedback. 302 In *2020 IEEE/ASME International Conference on Advanced Intelligent Mechatronics (AIM)*, 303 pages 1406–1412. IEEE, 2020. 12 304
- 228 [11] Andrea Cherubini, Valerio Ortenzi, Akansel Cosgun, 305 Robert Lee, and Peter Corke. Model-free vision-based 306 shaping of deformable plastic materials. *The International Journal of Robotics Research*, 39(14):1739– 307 1759, 2020. 12 308
- 229 [12] Samuel Clarke, Travers Rhodes, Christopher G Atkeson, 309 and Oliver Kroemer. Learning audio feedback for 310 estimating amount and flow of granular material. *Proceedings of Machine Learning Research*, 87, 2018. 13 311
- 230 [13] Ingrid Daubechies. *Ten lectures on wavelets*. SIAM, 312 1992. 12 313
- 231 [14] Lokenath Debnath and Firdous Ahmad Shah. *Wavelet 314 transforms and their applications*. Springer, 2002. 12 315
- 232 [15] Danny Driess, Zhiao Huang, Yunzhu Li, Russ 316 Tedrake, and Marc Toussaint. Learning multi-object 317 dynamics with compositional neural radiance fields. 318 *arXiv preprint arXiv:2202.11855*, 2022. 1 319
- 233 [16] Nima Fazeli, Miquel Oller, Jiajun Wu, Zheng Wu, 320 Joshua B Tenenbaum, and Alberto Rodriguez. See, 321 feel, act: Hierarchical learning for complex manipulation 322 skills with multisensory fusion. *Science Robotics*, 4(26):eaav3123, 2019. 12 323
- 234 [17] Niklas Funk, Georgia Chalvatzaki, Boris Belousov, 324 and Jan Peters. Learn2assemble with structured 325 representations and search for robotic architectural 326 construction. In *Conference on Robot Learning*, pages 327 1401–1411. PMLR, 2022. 12 328
- 235 [18] Roman Garnett. *Bayesian optimization*. Cambridge 329 University Press, 2023. 2 330
- 236 [19] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and 331 Mohammad Norouzi. Dream to control: Learning 332 behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019. 1 333
- 237 [20] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben 334 Villegas, David Ha, Honglak Lee, and James Davidson. 335 Learning latent dynamics for planning from pixels. 336 In *International conference on machine learning*, 337 pages 2555–2565. PMLR, 2019. 1 338
- 238 [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian 339 Sun. Spatial pyramid pooling in deep convolutional 340 networks for visual recognition. *IEEE transactions on pattern 341 analysis and machine intelligence*, 37(9): 342 1904–1916, 2015. 12 343
- 239 [22] François Robert Hogan and Alberto Rodriguez. Feed- 344 back control of the pusher-slider system: A story of 345 hybrid and underactuated contact dynamics. *arXiv preprint arXiv:1611.08268*, 2016. 1, 12 346
- 240 [23] Jun-Ting Hsieh, Bingbin Liu, De-An Huang, Li F Fei- 347 Fei, and Juan Carlos Niebles. Learning to decompose 348 and disentangle representations for video prediction. 349 *Advances in neural information processing systems*, 350 31, 2018. 12 351

- 324 [24] Zixuan Huang, Xingyu Lin, and David Held. Mesh- 378
325 based dynamics with occlusion reasoning for cloth 379
326 manipulation. *arXiv preprint arXiv:2206.02881*, 380
327 2022. 12 381
328
- 329 [25] Boris Ivanovic, Amine Elhafsi, Guy Rosman, Adrien 382
330 Gaidon, and Marco Pavone. Mats: An interpretable 383
331 trajectory forecasting representation for planning and 384
332 control. *arXiv preprint arXiv:2009.07517*, 2020. 1 385
333
- 334 [26] Ioannis G Kevrekidis, C William Gear, James M 386
335 Hyman, Panagiotis G Kevrekidis, Olof Runborg, 387
336 Constantinos Theodoropoulos, et al. Equation-free, 388
337 coarse-grained multiscale computation: enabling 389
338 microscopic simulators to perform system-level analysis. 390
339 *Commun. Math. Sci.*, 1(4):715–762, 2003. 12
340
- 341 [27] Ioannis G Kevrekidis, C William Gear, and Gerhard 391
342 Hummer. Equation-free: The computer-aided analysis 392
343 of complex multiscale systems. *AIChE Journal*, 50(7): 393
344 1346–1355, 2004. 12 394
- 345 [28] Scott Kuindersma, Robin Deits, Maurice Fallon, 395
346 Andrés Valenzuela, Hongkai Dai, Frank Permenter, 396
347 Twan Koolen, Pat Marion, and Russ Tedrake. 397
348 Optimization-based locomotion planning, estimation, 398
349 and control design for the atlas humanoid robot. *Autonom. 399
350 robots*, 40:429–455, 2016. 1
351
- 352 [29] J Nathan Kutz. *Data-driven modeling & scientific 400
353 computation: methods for complex systems & big 401
354 data*. Oxford University Press, 2013. 12
355
- 356 [30] J Nathan Kutz, Xing Fu, and Steven L Brunton. 402
357 Multiresolution dynamic mode decomposition. *SIAM 403
358 Journal on Applied Dynamical Systems*, 15(2):713– 404
359 735, 2016. 12
360
- 361 [31] Yunzhu Li, Jiajun Wu, Russ Tedrake, Joshua B Tenen- 405
362 baum, and Antonio Torralba. Learning particle 406
363 dynamics for manipulating rigid bodies, deformable 407
364 objects, and fluids. *arXiv preprint arXiv:1810.01566*, 408
365 2018. 10, 12
366
- 367 [32] Yunzhu Li, Jiajun Wu, Jun-Yan Zhu, Joshua B Tenen- 409
368 baum, Antonio Torralba, and Russ Tedrake. Propaga- 410
369 tion networks for model-based control under partial 411
370 observation. In *2019 International Conference on 412
371 Robotics and Automation (ICRA)*, pages 1205–1211. 413
372 IEEE, 2019. 7, 12
373
- 374 [33] Yunzhu Li, Shuang Li, Vincent Sitzmann, Pulkit 414
375 Agrawal, and Antonio Torralba. 3d neural scene 415
376 representations for visuomotor control. In *Conference 416
377 on Robot Learning*, pages 112–123. PMLR, 2022. 1, 12
378
- 379 [34] Xingyu Lin, Yufei Wang, Zixuan Huang, and David 380
381 Held. Learning visible connectivity dynamics for 382
383 cloth smoothing. In *Conference on Robot Learning*, 384
385 pages 256–266. PMLR, 2022. 12
386
- 387 [35] Qingkai Lu and Liangjun Zhang. Excavation learning 388
388 for rigid objects in clutter. *IEEE Robotics and Au- 389
389 tomaton Letters*, 6(4):7373–7380, 2021. 12
390
- 391 [36] Miles Macklin, Matthias Müller, Nuttapon Chentanez, 392
392 and Tae-Yong Kim. Unified particle physics for 393
393 real-time applications. *ACM Transactions on Graph- 394
394 ics (TOG)*, 33(4):1–12, 2014. 10
395
- 396 [37] Lucas Manuelli, Yunzhu Li, Pete Florence, and Russ 397
397 Tedrake. Keypoints into the future: Self-supervised 398
398 correspondence in model-based reinforcement learn- 399
399 ing. *arXiv preprint arXiv:2009.05085*, 2020. 1
399
- 400 [38] David Marr. *Vision: A computational investigation 401
401 into the human representation and processing of visual 402
402 information*. MIT press, 2010. 12
403
- 404 [39] Carolyn Matl, Yashraj Narang, Ruzena Bajcsy, Fabio 405
405 Ramos, and Dieter Fox. Inferring the material prop- 406
406 erties of granular media for robotic tasks. In *2020 ieee 407
407 international conference on robotics and automation 408
408 (icra)*, pages 2770–2777. IEEE, 2020. 13
409
- 410 [40] Matthias Minderer, Chen Sun, Ruben Villegas, For- 411
411 rester Cole, Kevin P Murphy, and Honglak Lee. Un- 412
412 supervised learning of object structure and dynamics 413
413 from videos. *Advances in Neural Information Pro- 414
414 cessing Systems*, 32, 2019. 1
415
- 416 [41] Peter Mitrano, Dale McConachie, and Dmitry Beren- 417
417 son. Learning where to trust unreliable models in an 418
418 unstructured world for deformable object manipula- 419
419 tion. *Science Robotics*, 6(54):eabd8170, 2021. 12
420
- 421 [42] Carsten Moenning and Neil A Dodgson. Fast march- 422
422 ing farthest point sampling. Technical report, Univer- 423
423 sity of Cambridge, Computer Laboratory, 2003. 7
424
- 425 [43] Damian Mrowca, Chengxu Zhuang, Elias Wang, Nick 426
426 Haber, Li F Fei-Fei, Josh Tenenbaum, and Daniel L 427
427 Yamins. Flexible neural representation for physics 428
428 prediction. *Advances in neural information process- 429
429 ing systems*, 31, 2018. 12
430
- 431 [44] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and 432
432 Sergey Levine. Data-efficient hierarchical reinforce- 433
433 ment learning. *Advances in neural information pro- 434
434 cessing systems*, 31, 2018. 12
435
- 436 [45] Monica N Nicolescu and Maja J Matarić. A hierar- 437
437 chical architecture for behavior-based robots. In *Pro- 438
438 ceedings of the first international joint conference on 439
439*

- 432 Autonomous agents and multiagent systems: part 1,
433 pages 227–233, 2002. 12 486
- 434
- 435 [46] Tao Pang, HJ Suh, Lujie Yang, and Russ Tedrake.
436 Global planning for contact-rich manipulation via lo-
437 cal smoothing of quasi-dynamic contact models. *arXiv*
438 preprint arXiv:2206.10787, 2022. 12 487
- 439
- 440 [47] Shubham Pateria, Budhitama Subagdja, Ah-hwee Tan,
441 and Chai Quek. Hierarchical reinforcement learning:
442 A comprehensive survey. *ACM Computing Surveys
(CSUR)*, 54(5):1–35, 2021. 12 488
- 443
- 444 [48] Tobias Pfaff, Meire Fortunato, Alvaro Sanchez-
445 Gonzalez, and Peter W Battaglia. Learning mesh-
446 based simulation with graph networks. *arXiv preprint
447 arXiv:2010.03409*, 2020. 12 489
- 448
- 449 [49] Haozhi Qi, Xiaolong Wang, Deepak Pathak, Yi Ma,
450 and Jitendra Malik. Learning long-term visual dynam-
451 ics with region proposal interaction networks. *arXiv
452 preprint arXiv:2008.02265*, 2020. 12 490
- 453
- 454 [50] Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias
455 Springenberg, Josh Merel, Martin Riedmiller, Raia
456 Hadsell, and Peter Battaglia. Graph networks as learn-
457 able physics engines for inference and control. In *Inter-
458 national Conference on Machine Learning*, pages
459 4470–4479. PMLR, 2018. 12 491
- 460
- 461 [51] Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias
462 Pfaff, Rex Ying, Jure Leskovec, and Peter Battaglia.
463 Learning to simulate complex physics with graph net-
464 works. In *International conference on machine learn-
465 ing*, pages 8459–8468. PMLR, 2020. 12 492
- 466
- 467 [52] Connor Schenck, Jonathan Tompson, Sergey Levine,
468 and Dieter Fox. Learning robotic manipulation of
469 granular media. In *Conference on Robot Learning*,
470 pages 239–248. PMLR, 2017. 2, 12 493
- 471
- 472 [53] Bokui Shen, Zhenyu Jiang, Christopher Choy,
473 Leonidas J Guibas, Silvio Savarese, Anima Anandku-
474 mar, and Yuke Zhu. Acid: Action-conditional implicit
475 visual dynamics for deformable object manipulation.
476 *arXiv preprint arXiv:2203.06856*, 2022. 12 494
- 477
- 478 [54] Haochen Shi, Huazhe Xu, Zhiao Huang, Yunzhu Li,
479 and Jiajun Wu. Robocraft: Learning to see, simulate,
480 and shape elasto-plastic objects with graph networks.
481 *arXiv preprint arXiv:2205.02909*, 2022. 12 495
- 482
- 483 [55] Tom Silver, Rohan Chitnis, Aidan Curtis, Joshua B
484 Tenenbaum, Tomás Lozano-Pérez, and Leslie Pack
485 Kaelbling. Planning with learned object importance in
large problem instances using graph neural networks.
In *Proceedings of the AAAI conference on artificial
intelligence*, volume 35, pages 11962–11971, 2021. 12 496
- 486 [56] Jasper Snoek, Hugo Larochelle, and Ryan P Adams.
487 Practical bayesian optimization of machine learning
488 algorithms. *Advances in neural information process-
489 ing systems*, 25, 2012. 8 497
- 490
- 491 [57] HJ Suh and Russ Tedrake. The surprising effective-
492 ness of linear models for visual foresight in object
493 pile manipulation. *arXiv preprint arXiv:2002.09093*,
494 2020. 12 498
- 495
- 496 [58] Hyung Ju Terry Suh, Tao Pang, and Russ Tedrake.
497 Bundled gradients through contact via randomized
498 smoothing. *IEEE Robotics and Automation Letters*,
499 7(2):4000–4007, 2022. 12 500
- 500
- 501 [59] Kuniyuki Takahashi, Wilson Ko, Avinash Umma-
502 disingu, and Shin-ichi Maeda. Uncertainty-aware self-
503 supervised target-mass grasping of granular foods. In
504 *2021 IEEE International Conference on Robotics and
505 Automation (ICRA)*, pages 2620–2626. IEEE, 2021.
506 12 507
- 507
- 508 [60] Russ Tedrake. Underactuated robotics: Learning,
509 planning, and control for efficient and agile machines
510 course notes for mit 6.832. *Working draft edition*, 3:4,
511 2009. 8 512
- 512
- 513 [61] Naftali Tishby, Fernando C Pereira, and William
514 Bialek. The information bottleneck method. *arXiv
515 preprint physics/0004057*, 2000. 1 516
- 516
- 517 [62] Hsiao-Yu Fish Tung, Zhou Xian, Mihir Prabhude-
518 sai, Shamit Lal, and Katerina Fragkiadaki. 3d-
519 odes: Viewpoint-invariant object-factorized environ-
520 ment simulators. *arXiv preprint arXiv:2011.06464*,
521 2020. 12 522
- 522
- 523 [63] Neea Tuomainen, David Blanco-Mulero, and Ville
524 Kyrki. Manipulation of granular materials by learning
525 particle interactions. *IEEE Robotics and Automation
526 Letters*, 7(2):5663–5670, 2022. 12 527
- 527
- 528 [64] Benjamin Ummenhofer, Lukas Prantl, Nils Thuerey,
529 and Vladlen Koltun. Lagrangian fluid simulation with
530 continuous convolutions. In *International Conference
531 on Learning Representations*, 2020. 12 532
- 532
- 533 [65] Alexander Sasha Vezhnevets, Simon Osindero, Tom
534 Schaul, Nicolas Heess, Max Jaderberg, David Sil-
535 ver, and Koray Kavukcuoglu. Feudal networks for
536 hierarchical reinforcement learning. In *International
537 Conference on Machine Learning*, pages 3540–3549.
538 PMLR, 2017. 12 539
- 539

- 540 diance field. In *2022 IEEE 18th International Conference on Automation Science and Engineering (CASE)*,
541 pages 1138–1143. IEEE, 2022. 1 594
542
543
- 544 [67] Nicholas Watters, Daniel Zoran, Theophane Weber,
545 Peter Battaglia, Razvan Pascanu, and Andrea Tac-
546 chetti. Visual interaction networks: Learning a
547 physics simulator from video. *Advances in neural in-*
548 *formation processing systems*, 30, 2017. 12 598
549
550 [68] Yufei Ye, Maneesh Singh, Abhinav Gupta, and Shub-
551 ham Tulsiani. Compositional video prediction. In *Pro-*
552 *ceedings of the IEEE/CVF International Conference on*
553 *Computer Vision*, pages 10353–10362, 2019. 12 600
554
555 [69] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli,
556 Jiajun Wu, Antonio Torralba, and Joshua B Tenen-
557 baum. Clevrer: Collision events for video represen-
558 tation and reasoning. *arXiv preprint arXiv:1910.01442*,
559 2019. 1, 12 601
560
561 [70] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiao-
562 gang Wang, and Jiaya Jia. Pyramid scene parsing
563 network. In *Proceedings of the IEEE conference on*
564 *computer vision and pattern recognition*, pages 2881–
565 2890, 2017. 12 602
566
567 [71] Jiaji Zhou, Yifan Hou, and Matthew T Mason. Push-
568 ing revisited: Differential flatness, trajectory plan-
569 ning, and stabilization. *The International Journal of*
570 *Robotics Research*, 38(12-13):1477–1489, 2019. 1,
571 12 603
572
573 [72] Guangxiang Zhu, Zhiao Huang, and Chongjie Zhang.
574 Object-oriented dynamics predictor. *Advances in Neu-*
575 *ral Information Processing Systems*, 31, 2018. 12 604
576
577 [73] Yifan Zhu, Laith Abdulmajeid, and Kris Hauser. A
578 data-driven approach for fast simulation of robot lo-
579 comotion on granular media. In *2019 international*
580 *conference on robotics and automation (ICRA)*, pages
581 7653–7659. IEEE, 2019. 13 605
582
583
584
585
586
587
588
589
590
591
592
593

648 Appendix

649 650 A. Method

651 652 A.1. Problem Formulation

653 Our goal is to derive the resolution ω to represent the
 654 environment to achieve the best trade-off between efficiency
 655 and effectiveness for control optimization. We define the
 656 following trajectory optimization problem over a horizon T :

$$\begin{aligned} \min_{\{u_t\}} \quad & c(\mathbf{z}_T^\omega, y_g), \\ \text{s.t.} \quad & \omega = g(y_0, y_g), \\ & \mathbf{z}_0^\omega = h(y_0, \omega), \\ & \mathbf{z}_{t+1}^\omega = f(\mathbf{z}_t^\omega, u_t, \omega), \end{aligned} \quad (1)$$

664 where the resolution regressor $g(\cdot, \cdot)$ takes the current observation
 665 y_0 and the goal configuration y_g as input and predicts
 666 the model resolution. $h(\cdot, \cdot)$, the perception module, takes
 667 in the current observation y_0 and the predicted resolution
 668 ω , then derives the scene representation \mathbf{z}_0^ω for the current
 669 time step. The dynamics module $f(\cdot, \cdot, \cdot)$ takes the current
 670 scene representation \mathbf{z}_t^ω , the input action u_t , and the
 671 resolution ω as inputs, and then predicts the representation's
 672 evolution at the next time step \mathbf{z}_{t+1}^ω . The optimization aims
 673 to find the action sequence $\{u_t\}$ to minimize the task objective
 674 $c(\mathbf{z}_T^\omega, y_g)$.

675 In the following sections, we describe (1) the details
 676 of the perception module $h(\cdot, \cdot)$ and the dynamics mod-
 677 ule $f(\cdot, \cdot, \cdot)$ in Section A.2, (2) how we obtain the self-
 678 supervision for the resolution regressor $g(\cdot, \cdot)$ in Sec-
 679 tion A.3, and (3) how we solve Equation 1 in a closed plan-
 680 ning loop in Section A.4.

681 682 A.2. Dynamic-Resolution Model Learning

683 To instantiate the optimization problem defined in Equa-
 684 tion 1, we use graphs of different sizes as the representation
 685 $\mathbf{z}_t^\omega = (\mathcal{O}_t, \mathcal{E}_t)$, where ω indicates the number of vertices
 686 in the graph. The vertices $\mathcal{O}_t = \{o_t^i\}_{i=1,\dots,|\mathcal{O}_t|}$ denote the
 687 particle set and o_t^i represents the 3D position of the i^{th}
 688 particle. The edge set $\mathcal{E}_t = \{e_t^j\}_{j=1,\dots,|\mathcal{E}_t|}$ denotes the relations
 689 between the particles, where $e_t^j = (u_t^j, v_t^j)$ denotes an edge
 690 pointing from particle of index v_t^j to u_t^j .

691 To obtain the particle set \mathcal{O}_t from the RGBD visual
 692 observation y_t , we first transform the RGBD image into
 693 a point cloud and then segment the point cloud to obtain
 694 the foreground according to color and depth information
 695 $\bar{y}_t \in \mathbb{R}^{N \times 3}$. We then deploy the farthest point sampling
 696 technique [42] to subsample the foreground but ensure suf-
 697 ficient coverage of \bar{y}_t . Specifically, given already sampled
 698 particles $o_t^{1,\dots,i-1}$, we apply

$$o_t^i = \arg \max_{y^k \in \bar{y}_t} \min_{o_t^j \in o_t^{1,\dots,i-1}} \|y^k - o_t^j\|_2^2 \quad (2)$$

702 to find the i^{th} particle o_t^i . We iteratively apply this process
 703 until we reach ω particles. Different choices of ω indicate
 704 scene representations at different abstraction levels, as illus-
 705 trated in Figure 2a. The edge set is constructed dynamically
 706 over time and connects particles within a predefined dis-
 707 tance while limiting the maximum number of edges a node
 708 can have.

709 We instantiate the dynamics model $f(\cdot, \cdot, \cdot)$ as graph neu-
 710 tral networks (GNNs) that predict the evolution of the graph
 711 representation \mathbf{z}_t^ω under external actions u_t and the selected
 712 resolution ω . $f(\cdot, \cdot, \cdot)$ consists of node and edge encoders
 713 $f_{\mathcal{O}}^{\text{enc}}(\cdot, \cdot, \cdot)$, $f_{\mathcal{E}}^{\text{enc}}(\cdot, \cdot, \cdot)$ to obtain node and edge representa-
 714 tions:

$$\begin{aligned} p_t^i &= f_{\mathcal{O}}^{\text{enc}}(o_t^i, u_t, \omega), \quad i = 1, \dots, |\mathcal{O}_t|, \\ q_t^j &= f_{\mathcal{E}}^{\text{enc}}(o_t^{u_t^j}, o_t^{v_t^j}, \omega), \quad j = 1, \dots, |\mathcal{E}_t|. \end{aligned} \quad (3)$$

715 We then have node and edge decoders $f_{\mathcal{O}}^{\text{dec}}(\cdot, \cdot)$, $f_{\mathcal{E}}^{\text{dec}}(\cdot, \cdot)$
 716 to obtain the corresponding representations and predict the
 717 representation at the next time step:

$$\begin{aligned} r_t^j &= f_{\mathcal{E}}^{\text{dec}}(q_t^j, \omega), \quad j = 1, \dots, |\mathcal{E}_t|, \\ \hat{o}_{t+1}^i &= f_{\mathcal{O}}^{\text{dec}}(p_t^i, \sum_{j \in \mathcal{N}_i} r_t^j), \quad i = 1, \dots, |\mathcal{O}_t|, \end{aligned} \quad (4)$$

722 where \mathcal{N}_i is the index set of the edges, in which particle i
 723 is the receiver. In practice, we follow Li et al. [32] and use
 724 multi-step message passing over the graph to approximate
 725 the instantaneous propagation of forces.

726 To train the dynamics model, we iteratively predict fu-
 727 ture particle states over a time horizon of T and then op-
 728 timize the neural network's parameters by minimizing the
 729 mean squared error (MSE) between the predictions and the
 730 ground truth future states:

$$\mathcal{L} = \frac{1}{T \cdot |\mathcal{O}_t|} \sum_{t'=1}^T \sum_{i=1}^{|\mathcal{O}_t|} \|\hat{o}_{t+t'}^i - o_{t+t'}^i\|_2^2. \quad (5)$$

731 732 A.3. Adaptive Resolution Selection

733 The previous sections discussed how to obtain the parti-
 734 cle set and how we predict its evolution given a resolution
 735 ω . In this section, we present how we learn the resolution
 736 regressor $g(\cdot, \cdot)$ in Equation 1 that can automatically deter-
 737 mine the resolution in a self-supervised manner. Specifi-
 738 cally, we intend to find the resolution ω that is the most
 739 effective for minimizing the task objective given the cur-
 740 rent observation y_0 and the goal y_g . We reformulate the
 741 optimization problem in Equation 1 by considering ω as a
 742 variable of the objective function as the following:

$$\begin{aligned} c^*(y_0, y_g, \omega) &= \min_{\{u_t\}} c(\mathbf{z}_T^\omega, y_g), \\ \text{s.t.} \quad & \mathbf{z}_0^\omega = h(y_0, \omega), \\ & \mathbf{z}_{t+1}^\omega = f(\mathbf{z}_t^\omega, u_t, \omega). \end{aligned} \quad (6)$$

756 For a given ω , we solve the above optimization problem
 757 via a combination of sampling and gradient descent using
 758 shooting methods [60] under a given time budget—the
 759 higher resolution representation will go through fewer optimiza-
 760 tion iterations. For simplicity, we denote the objective in
 761 Equation 6 as $c^*(\omega)$ in the following part of this section.
 762

763 Given the formulation, we are then interested in finding
 764 the parameter ω that can minimize the following objective:
 765

$$\begin{aligned} \min_{\omega} \quad & c^+(\omega) = c^*(\omega) + R(\omega), \\ \text{s.t.} \quad & \omega \in (\omega_{\min}, \omega_{\max}), \end{aligned} \quad (7)$$

766 where $R(\omega)$ is a regularizer penalizing the choice of an
 767 excessively large ω to encourage efficiency. Regularizer
 768 details can be found in supplementary materials. We use
 769 Bayesian optimization [56] to find the optimal ω by
 770 iteratively sampling ω and approximating $c^+(\omega)$ using the
 771 Gaussian process. At each sampling stage, we sample
 772 one or more data points ω_i according to the expected im-
 773 provement of the objective function and evaluate their value
 774 $c^+(\omega_i)$. Then, at the approximation stage, we assume
 775 the distribution of $c^+(\omega)$ follows the Gaussian distribu-
 776 tion $\mathcal{N}(\mu(\omega), \sigma^2)$; thus, the joint distribution of the eval-
 777 uated points $\Omega_{\text{train}} = [\omega_1, \dots, \omega_n]$ and the testing points
 778 $\Omega_{\text{test}} = [\omega'_1, \dots, \omega'_m]$ can be expressed as the following:
 779

$$\begin{aligned} \mathbf{C}_{\text{train}} &= [c^+(\omega_1), \dots, c^+(\omega_n)], \\ \mathbf{M}_{\text{train}} &= [\mu(\omega_1), \dots, \mu(\omega_n)], \\ \mathbf{C}_{\text{test}} &= [c^+(\omega'_1), \dots, c^+(\omega'_m)], \\ \mathbf{M}_{\text{test}} &= [\mu(\omega'_1), \dots, \mu(\omega'_m)], \\ \begin{bmatrix} \mathbf{C}_{\text{train}} \\ \mathbf{C}_{\text{test}} \end{bmatrix} &\sim \mathcal{N}\left(\begin{bmatrix} \mathbf{M}_{\text{train}} \\ \mathbf{M}_{\text{test}} \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^\top & \mathbf{K}_{**} \end{bmatrix}\right), \end{aligned} \quad (8)$$

780 where \mathbf{K} is a kernel function matrix derived via $\mathbf{K} = K(\Omega_{\text{train}}, \Omega_{\text{train}})$. $K(\cdot, \cdot)$ is the kernel function used to
 781 compute the covariance. Similarly, $\mathbf{K}_* = K(\Omega_{\text{train}}, \Omega_{\text{test}})$ and
 782 $\mathbf{K}_{**} = K(\Omega_{\text{test}}, \Omega_{\text{test}})$.

783 Equation 8 shows the joint probability of $\mathbf{C}_{\text{train}}$ and \mathbf{C}_{test}
 784 conditioned on Ω_{train} and Ω_{test} . Through marginalization,
 785 we could fit $c^+(\omega)$ using the following conditional distri-
 786 bution:

$$\mathbf{C}_{\text{test}} | \mathbf{C}_{\text{train}}, \Omega_{\text{train}}, \Omega_{\text{test}} \sim \mathcal{N}(\mathbf{K}_*^\top \mathbf{K} \mathbf{C}_{\text{train}}, \mathbf{K}_{**} - \mathbf{K}_*^\top \mathbf{K}^{-1} \mathbf{K}_*). \quad (9)$$

801 We can then use the mean value of the Gaussian distribution
 802 in Equation 9 as the metric to minimize $c^+(\omega)$. Therefore,
 803 the solution to Equation 7 is approximated as the following:
 804

$$\begin{aligned} \omega^* &= \arg \min_{\omega} \mathbf{K}_*^\top \mathbf{K} \mathbf{C}_{\text{train}} \\ \text{s.t.} \quad & \omega \in (\omega_{\min}, \omega_{\max}). \end{aligned} \quad (10)$$

805 To train the resolution regressor, we randomly generate
 806 a dataset containing the observation and goal pairs (y_0, y_g) .
 807

810 For each pair, we follow the above optimization process to
 811 generate the optimal resolution label ω^* . We then train the
 812 resolution regressor $\omega = g(y_0, y_g)$ to predict the resolution
 813 based on the observation and the goal via supervised learn-
 814 ing. Training the ω regressor is a self-supervised learning
 815 process, as the labels are automatically generated via an op-
 816 timization process without any human labeling.
 817

A.4. Closed-Loop Planning on Adaptive Repr.

818 Now that we have obtained the resolution regressor g ,
 819 the perception module h , and the dynamics module f . We
 820 can wire things together to solve Equation 1 and use the
 821 optimized action sequence in a closed loop within a model-
 822 predictive control (MPC) framework [8]. Specifically, for
 823 each MPC step, we follow Algorithm 1, which first deter-
 824 mines the resolution to represent the environment, then uses
 825 a combination of sampling and gradient descent to derive
 826 the action sequence through trajectory optimization using
 827 the shooting method. We then execute the first action from
 828 the action sequence in the real world, obtain new observa-
 829 tions, and apply Algorithm 1 again. Such a process allows
 830 us to take feedback from the environment and adaptively se-
 831 lect the most appropriate resolution at each step as the task
 832 progresses. Figure 2b also shows an overview of the future
 833 prediction and inverse planning process. Details including
 834 task objective definition and MPC hyperparameter are in-
 835 cluded in supplementary materials.

Algorithm 1 Trajectory optimization at each MPC step

836 **Input:** Current observation y_0 , goal y_g , time horizon T ,
 837 the resolution regressor g , the perception module h ,
 838 the dynamics module f , and gradient descent iteration
 839 count N
Output: Actions $u_{0:T-1}$

840 Predict the resolution $\omega \leftarrow g(y_0, y_g)$
 841 Obtain the current representation $z_0^\omega \leftarrow h(y_0, \omega)$
 842 Sample M action sequences $\hat{u}_{0:T-1}^{1:M}$
 843 **for** $m = 1, \dots, M$ **do**
 844 **for** $i = 1, \dots, N$ **do**
 845 **for** $t = 0, \dots, T-1$ **do**
 846 Predict the next step $z_{t+1}^\omega \leftarrow f(z_t^\omega, \hat{u}_t^m, \omega)$
 847 **end for**
 848 Calculate the task loss $c^m \leftarrow c(z_T^\omega, y_g)$
 849 **if** $i < N$ **then**
 850 Update $\hat{u}_{0:T-1}^m$ using gradients $\nabla_{\hat{u}_{0:T-1}^m} c^m$
 851 **end if**
 852 **end for**
 853 **end for**
 854 $m^* \leftarrow \arg \min_m c^m$
 855 Return $\hat{u}_{0:T-1}^{m^*}$

864

A.5. Perception Module Details

Building graph $\mathcal{Z}_t^\omega = (\mathcal{O}_t, \mathcal{E}_t)$ from point cloud $\bar{y}_t \in \mathbb{R}^{N \times 3}$ contains two phases. First, we sample $\mathcal{O}_t^{\text{fps}}$ from the point cloud using farthest-point sampling. Specifically, given already sampled particles $o_t^{1, \dots, i-1}$, we apply

$$o_t^i = \arg \max_{y^k \in \bar{y}_t} \min_{o_t^j \in o_t^{1, \dots, i-1}} \|y^k - o_t^j\|_2^2 \quad (11)$$

to find the i^{th} particle o_t^i . We iteratively apply this process until we reach ω particles. We found that sampled particles from the point cloud are likely to be at the edge of underneath objects. To bias sampled particles towards object centers, we define \mathcal{O}_t as mass centers of $\mathcal{O}_t^{\text{fps}}$ neighbor points. For example, for $o_t^{i, \text{fps}} \in \mathcal{O}_t^{\text{fps}}$, we find the corresponding $o_t^i \in \mathcal{O}_t$ by applying

$$\begin{aligned} \bar{y}_t' &= \{y | y \in \bar{y}_t, \|y - o_t^{i, \text{fps}}\|_2 \leq r_{\text{center}}\} \\ o_t^i &= \frac{1}{|\bar{y}_t'|} \sum_{y \in \bar{y}_t'} y. \end{aligned} \quad (12)$$

r_{center} is the hyperparameter to determine the neighbor points.

To find edges \mathcal{E}_t , we first approximate particle displacements $\Delta \mathcal{O}_t$ using the action u_t . We compute the sweeping region given the action u_t . If o_t^i is within the sweeping region, Δo_t^i is the vector from o_t^i to the pusher end. We define $\hat{\mathcal{O}}_t = \Delta \mathcal{O}_t + \mathcal{O}_t$. For $\hat{o}_t^i, (\hat{o}_t^i, \hat{o}_t^j) \in \mathcal{E}_t$ if it satisfies the following criteria:

$$\begin{aligned} \|\hat{o}_t^i - \hat{o}_t^j\|_2 &< r_{\text{edge}} \\ \hat{o}_t^j &\in \text{kNN}(\hat{o}_t^i). \end{aligned} \quad (13)$$

r_{edge} is the hyperparameter to determine the distance needed for two nodes to interact. $\text{kNN}(\hat{o}_t^i)$ is the set of k-nearest-neighbor of the node \hat{o}_t^i and k is a hyperparameter.

A.6. Distribution Distance

We compute distribution distance to evaluate the performance of different methods on gathering tasks quantitatively. The distribution distance is computed similarly to the task objective. We use the RGBD observation to segment out the foreground object and obtain foreground pixels $\mathcal{F}_t \in \mathbb{R}^{F \times 2}$. Similar to Equation 14, the distribution distance d is defined using the following equation:

$$d = \sum_{f_i \in \mathcal{F}_t} \min_{q_j \in \mathcal{Q}_t} \|f_i - q_j\|_2 + \sum_{q_j \in \mathcal{Q}_t} \min_{f_i \in \mathcal{F}_t} \|f_i - q_j\|_2 \quad (14)$$

A.7. High-Level Planner for Sort Task

In the **sort** task, we use a high-level planner to avoid colliding between two object piles. We use the A* search algorithm to find the high-level path. We represent every blob

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

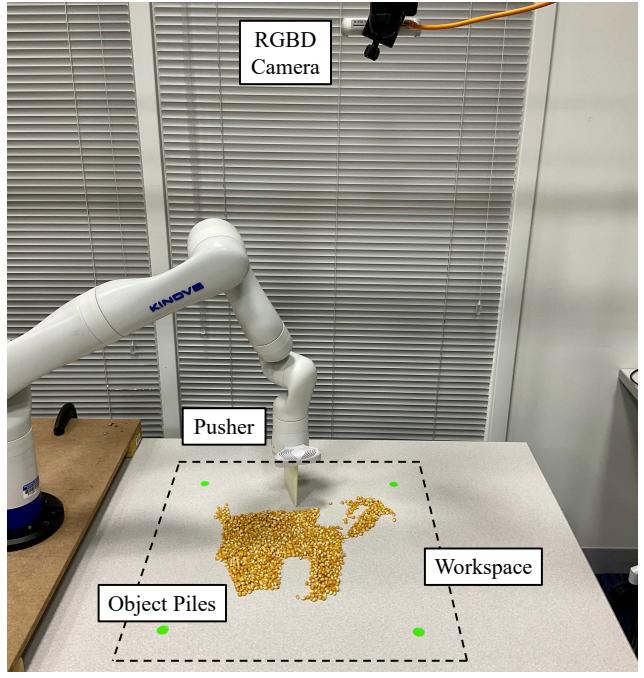
967

968

969

970

971



(a) Robot setup

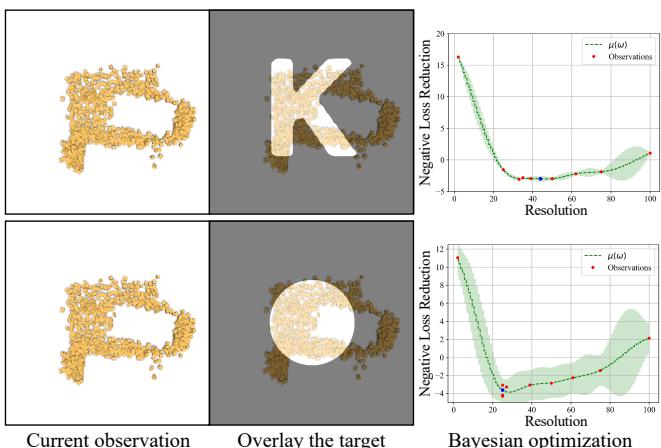


(b) Object piles considered in this work

Figure 3. **Robot setup and the testing object piles.** (a) The dashed black square shows the robot’s workspace. The robotic manipulator, equipped with a pusher at the end effector, pushes the object piles within the workspace. A calibrated RGBD camera mounted at the top provides visual observations of the environment. (b) We show the object piles considered in this work, including M&M, almond, granola, candy, carrot, rice, corn, and coffee beans.

of object piles as a circle with a fixed radius in image space. Therefore, the state for one blob can be represented as its 2D blob center in image space. For the search algorithm, if there are k blobs in the scene, the node is the concatenation of k blob centers. One node is connected to nodes that can be reached by changing one blob center in a single step with no collision.

To accelerate motion planning, we divide the image space into a sparse grid. The blob center will only be on the grid. In addition, to encourage a path with fewer steps, we add a constant cost for each path so that a path with fewer steps is preferred.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987

(a) Same initial but different goal configurations

Figure 4. **Optimal resolution differs depending on the initial and goal configurations.** (a) We show two examples with the same initial but different goal configurations. We apply Bayesian optimization to solve the problem discussed in Section A.3 to find the optimal resolution for both cases. The example with a more complicated target shape requires a higher-resolution representation to be the most effective at making task progress. (b) When the goal is to gather the pieces in the center of the workspace, a coarse representation is sufficient for examples with spread-out pieces. The task progresses as long as the agent pushes any outlying pieces toward the goal region. In contrast, a higher-resolution representation is needed to reveal the subtle difference between the initial and goal configurations when they are close.

994

B. Experiment

B.1. Setup

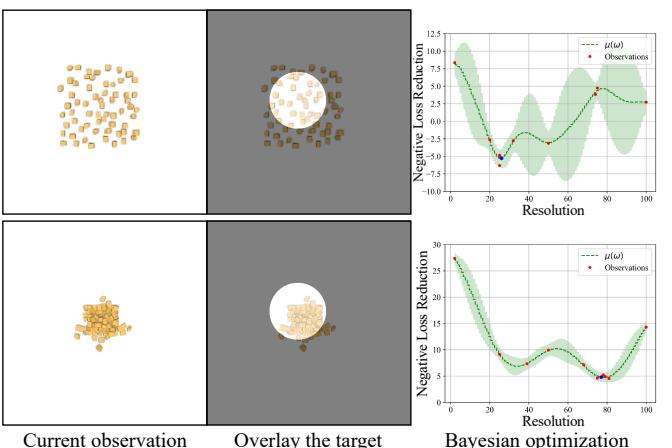
We conduct experiments in both the simulation environment and the real world. The simulation environment is built using NVIDIA FleX [31, 36], a position-based simulator capable of simulating the interactions between a large number of object pieces. In the real world, we conducted experiments using the setup shown in Figure 3a. We use RealSense D455 as the top-down camera to capture the RGBD visual observations of the workspace. We attach a flat pusher to the robotic manipulator’s end effector to manipulate the object piles.

B.2. Tasks

We evaluate our methods on three object pile manipulation tasks that are common in daily life.

- **Gather:** The robot needs to push the object pile into a target blob with different locations and radii.
- **Redistribute:** The robot is tasked to manipulate the object piles into many complex target shapes, such as letters.
- **Sort:** The robot has to move two different object piles to target locations without mixing each other.

We use a unified dynamics model for all three tasks, which involve objects pieces of different granularities, appearances, and physical properties (Figure 3b).



(b) Same goal but different initial configurations

B.3. Trade-Off Between Efficiency and Effectiveness

The trade-off between efficiency and effectiveness can vary depending on the tasks, the current, and the goal configurations. As we have discussed in Section A.3, given the resolution ω , we set a fixed time budget to solve the optimization problem defined in Equation 6. Intuitively, if the resolution is too low, the representation will not contain sufficiently detailed information about the environment to accomplish the task, the optimization of which is efficient but not effective enough to finish the task. On the contrary, if we choose an excessively high resolution, the representation will carry redundant information not necessary for the task and can be inefficient in optimization. We thus conduct experiments evaluating whether the trade-off exists (i.e., whether the optimal resolution ω calculated from Equation 10 is different for different initial and goal configurations).

We use Bayesian optimization and follow the algorithm described in Section A.3 to find the optimal trade-off on **Gather** and **Redistribute** tasks in the simulation. As shown in Figure 4a, higher-resolution dynamics models do not necessarily lead to better performance due to their optimization inefficiency. Compared between goal configurations, even if the current observation is the same, a more complicated goal typically requires a higher resolution representation to make the most effective task progression. More specifically, when the target region is a plain circle, the coarse representation captures the rough shape of the object pile, sufficient for the task objective, allowing more efficient optimization

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

than the higher-resolution counterparts. However, when the target region has a more complicated shape, low-resolution representation fails to inform downstream MPC of detailed object pile shapes. Therefore, high-resolution representation is necessary for effective trajectory optimization.

The desired representation does not only depend on goal configurations. Even if the goal configurations are the same, different initial configurations can also lead to different optimal resolutions, as illustrated in Figure 4b. When the initial configuration is more spread out, the most effective way of decreasing the loss is by pushing the outlying pieces to the goal region. Our farthest sampling strategy, even with just a few particles, could capture outlying pieces and helps the agent to make good progress. Therefore, when pieces are sufficiently spread out, higher particle resolution does not necessarily contain more useful information for the task but makes the optimization process inefficient. On the other hand, when the initial configuration concentrates on the goal region, to effectively decrease the task objective, MPC needs more detailed information about the object pile’s geometry to pinpoint the mismatching area. For example, the agent needs to know more precise contours of the goal region and the outlying part of object piles to decide how to improve the planning results further. Low-resolution representations will be less effective in revealing the difference between the current observation and the goal, thus less helpful in guiding the agent to make action decisions.

B.4. Is a Single Resolution Dynamics Model Sufficient?

Although there is a trade-off between representation resolution and task progression, can we benefit from this trade-off in trajectory optimization? We compared our dynamic-resolution dynamics model with fixed-resolution dynamics models on **Gather** and **Redistribute** tasks. Figure 1 shows how our model changes its resolution prediction as MPC proceeds in the real world. Trained on the generated dataset of optimal ω , our regressor learned that fixing a resolution throughout the MPC process is not optimal. Instead, our regressor learns to adapt the resolution according to the current observation feedback. In addition, for the example shown in Figure 1, we can see that the resolution increases as object piles approach the goal. This matches our expectation as explained in Section B.3.

We quantitatively evaluate different fixed-resolution baselines and our adaptive representation learning algorithms in simulation. We record the final step distribution distance between object piles and the goal. Specifically, given a distance threshold τ_p , the number of tasks with a distance lower than τ_p is N_p , and the total number of tasks is N . The task score is then defined as N_p/N (i.e., y-axis in Figure 5). Our adaptive resolution model almost always achieves the highest task score, regardless of the threshold

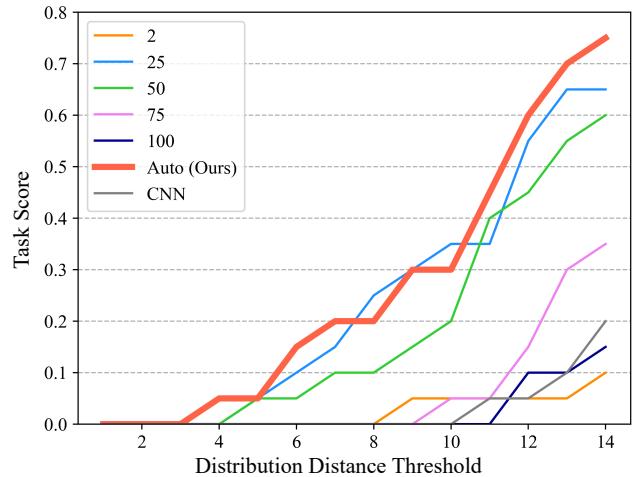


Figure 5. **Model-predictive control (MPC) results.** We evaluated the MPC performance on different representation choices. We use the task score as the evaluation metric. Task execution trial results in a distribution distance lower than the threshold is considered a success. The task score is the number of successful trials divided by the total number of task trials for both the **Gather** and **Redistribute** tasks. Our method automatically and adaptively selects the scene representation, which achieves the best overall performance compared with the scores of fixed-resolution baselines and a method that uses convolutional neural networks (CNN) as the dynamics model class.

used.

Figure 6a shows a qualitative comparison between the fixed-resolution baselines and our dynamic-resolution selection method on the **Gather** task in the real world. All methods start from a near-identical configuration. We can see from the qualitative results that our method manipulates the object pile to a configuration closest to the goal region, whereas the best-performing fixed-resolution baseline still has some outlying pieces far from the goal region. In addition, we could see from the quantitative evaluation curve in Figure 6b that our model is always the best throughout the whole MPC process. Representations with an excessively high resolution are unlikely to converge to a decent solution within the time budget, as demonstrated by resolutions 75 and 100. Conversely, if the representation is too low resolution, it will converge to a loss much higher than our model. A resolution of 25 reached a comparable final loss to our method. However, because the same resolution was ineffective for initial timesteps, its loss does not reduce as rapidly as our adaptive approach. Because our model could adapt to different resolutions in different scenes, making it more effective at control optimization.

That is why our model could reach the goal region faster than all other fixed-resolution models and consistently performs better at all timestamps, highlighting the benefits of adaptive resolution selection.

| | | |
|------|---|------|
| 1188 | B.5. Can a Unified Dynamics Model Achieve All | 1242 |
| 1189 | Three Tasks? | 1243 |
| 1190 | | 1244 |
| 1191 | We further demonstrate that our method could work on | 1245 |
| 1192 | all three tasks and diverse object piles. For the Gather task, | 1246 |
| 1193 | we test our method on different objects with different initial | 1247 |
| 1194 | and goal configurations. From left to right in Figure 6c, our | 1248 |
| 1195 | agent gathers different object piles made with almond, granola, | 1249 |
| 1196 | , or M&M™. Different appearances and physical prop- | 1250 |
| 1197 | erties challenge our method’s generalization capability. For | 1251 |
| 1198 | example, while almonds and granola are almost quasi-static | 1252 |
| 1199 | during the manipulation, M&M™ will roll around and have | 1253 |
| 1200 | high uncertainties in its dynamics. In addition, unlike al- | 1254 |
| 1201 | mmonds and M&M™, granola pieces are non-uniform. Our | 1255 |
| 1202 | method has a good performance for all these objects and | 1256 |
| 1203 | configurations. | 1257 |
| 1204 | For the Redistribute task, we redistribute carrots and al- | 1258 |
| 1205 | monds into target letters ‘J’, ‘T’, and ‘U’ with spread-out | 1259 |
| 1206 | initial configurations. The final results match the desired | 1260 |
| 1207 | letter shape. Please check our supplementary materials for | 1261 |
| 1208 | video illustrations of the manipulation process. | 1262 |
| 1209 | For the Sort task, we use a high-level motion planner to | 1263 |
| 1210 | find the intermediate waypoints in the image space. Then | 1264 |
| 1211 | we use a similar method as Gather task to push the ob- | 1265 |
| 1212 | ject pile into the target location. For the three examples | 1266 |
| 1213 | shown in Figure 6e, we require object piles to go to their | 1267 |
| 1214 | own target locations while not mixing with each other. Ob- | 1268 |
| 1215 | jects with different scales and shapes are present here. For | 1269 |
| 1216 | example, coffee beans have smaller granularity and round | 1270 |
| 1217 | shapes, while candies are relatively large and square. Here | 1271 |
| 1218 | we demonstrate success trials of manipulating the object | 1272 |
| 1219 | piles to accomplish the Sort task for different objects and | 1273 |
| 1220 | goal configurations. Please check our video for the mani- | 1274 |
| 1221 | pulation process. | 1275 |
| 1222 | | 1276 |
| 1223 | C. Related Work | 1277 |
| 1224 | | 1278 |
| 1225 | C.1. Scene Representation at Different Abstraction | 1279 |
| 1226 | Levels | 1280 |
| 1227 | To build multi-scale models of the dynamical sys- | 1281 |
| 1228 | tems, prior works have adopted wavelet-based methods and | 1282 |
| 1229 | windowed Fourier Transforms to perform multi-resolution | 1283 |
| 1230 | analysis [13, 14, 29]. Kevrekidis et al. [26, 27] investi- | 1284 |
| 1231 | gated equation-free, multi-scale modeling methods via | 1285 |
| 1232 | computer-aided analysis. Kutz et al. [30] also combined | 1286 |
| 1233 | multi-resolution analysis with dynamic mode decomposi- | 1287 |
| 1234 | tion for the decomposition of multi-scale dynamical data. | 1288 |
| 1235 | Our method is different in that we directly learn from vision | 1289 |
| 1236 | data for the modeling and planning of real-world manipu- | 1290 |
| 1237 | lation systems. | 1291 |
| 1238 | In computer vision, Marr [38] laid the foundation by | 1292 |
| 1239 | proposing a multi-level representational framework back | 1293 |
| 1240 | in 1982. Since then, people have investigated pyramid | 1294 |
| 1241 | methods in image processing [1, 7] using Gaussian, Lapla- | 1295 |
| | cian, and Steerable filters. Combined with deep neural | |
| | networks, the multi-resolution visual representation also | |
| | showed stunning performance in various visual recogni- | |
| | tion tasks [21, 70]. In the field of robotics, reinforcement | |
| | learning researchers have also studied task- or behavior-level | |
| | abstractions and come up with various hierarchical reinfor- | |
| | cegment learning algorithms [2, 6, 16, 44, 45, 47, 65]. Our | |
| | method instead focuses on spatial abstractions from vision, | |
| | where we learned structured representations based on par- | |
| | ticles to model the object interactions within the environment | |
| | at different levels. | |
| | | |
| | C.2. Compositional Model Learning for Robotic | |
| | Manipulation | |
| | | |
| | Physics-based models have demonstrated their effective- | |
| | ness in many robotic manipulation tasks (e.g., [22, 46, 58, | |
| | 71]). However, they typically rely on complete information | |
| | about the environment, limiting their use in scenarios where | |
| | full-state estimation is hard or impossible to acquire (e.g., | |
| | precise shape and pose estimation of each one of the object | |
| | pieces in Figure 1). Learning-based approaches provide a | |
| | way of building dynamics models directly from visual ob- | |
| | servations. Prior methods have investigated various scene | |
| | representations for dynamics modeling and manipulation | |
| | of objects with complicated physical properties, including | |
| | clothes [24, 34], ropes [10, 41], fluids [33], softbodies [53], | |
| | and plasticine [54]. Among the methods, graph-structured | |
| | neural networks (GNNs) have shown great promise by in- | |
| | roducing explicit relational inductive biases [4]. Prior | |
| | works have shown GNNs’ effectiveness in modeling com- | |
| | positional dynamical systems involving the interaction be- | |
| | tween multiple objects [3, 9, 17, 32, 50, 55], systems repre- | |
| | sented using particles or meshes [31, 43, 48, 51, 64], or for | |
| | compositional video prediction [23, 49, 62, 67, 68, 69, 72]. | |
| | However, these works typically assume scene representa- | |
| | tion at a fixed resolution, whereas our method learns a uni- | |
| | fied graph dynamics model that can generalize to scene rep- | |
| | resentations at different levels of abstraction. | |
| | | |
| | C.3. Object Pile Manipulation | |
| | | |
| | Robotic manipulation of object piles and granular pieces | |
| | has been one of the core capabilities if we want to deploy | |
| | robot systems for complicated tasks like cooking and manu- | |
| | facturing. Suh and Tedrake [57] proposed to learn visual dy- | |
| | namics based on linear models for redistributing the object | |
| | pieces. Along the lines of learning the dynamics of granu- | |
| | lar pieces, Tuomainen et al. [63] and Schenck et al. [52] | |
| | also proposed the use of GNNs or convolutional neural dy- | |
| | namics models for scooping and dumpling granular pieces. | |
| | Other works introduced success predictors for excavation | |
| | tasks [35], a self-supervised mass predictor for grasping | |
| | granular foods [59], visual serving for shaping deformable | |
| | plastic materials [11], or data-driven methods to calibrate | |

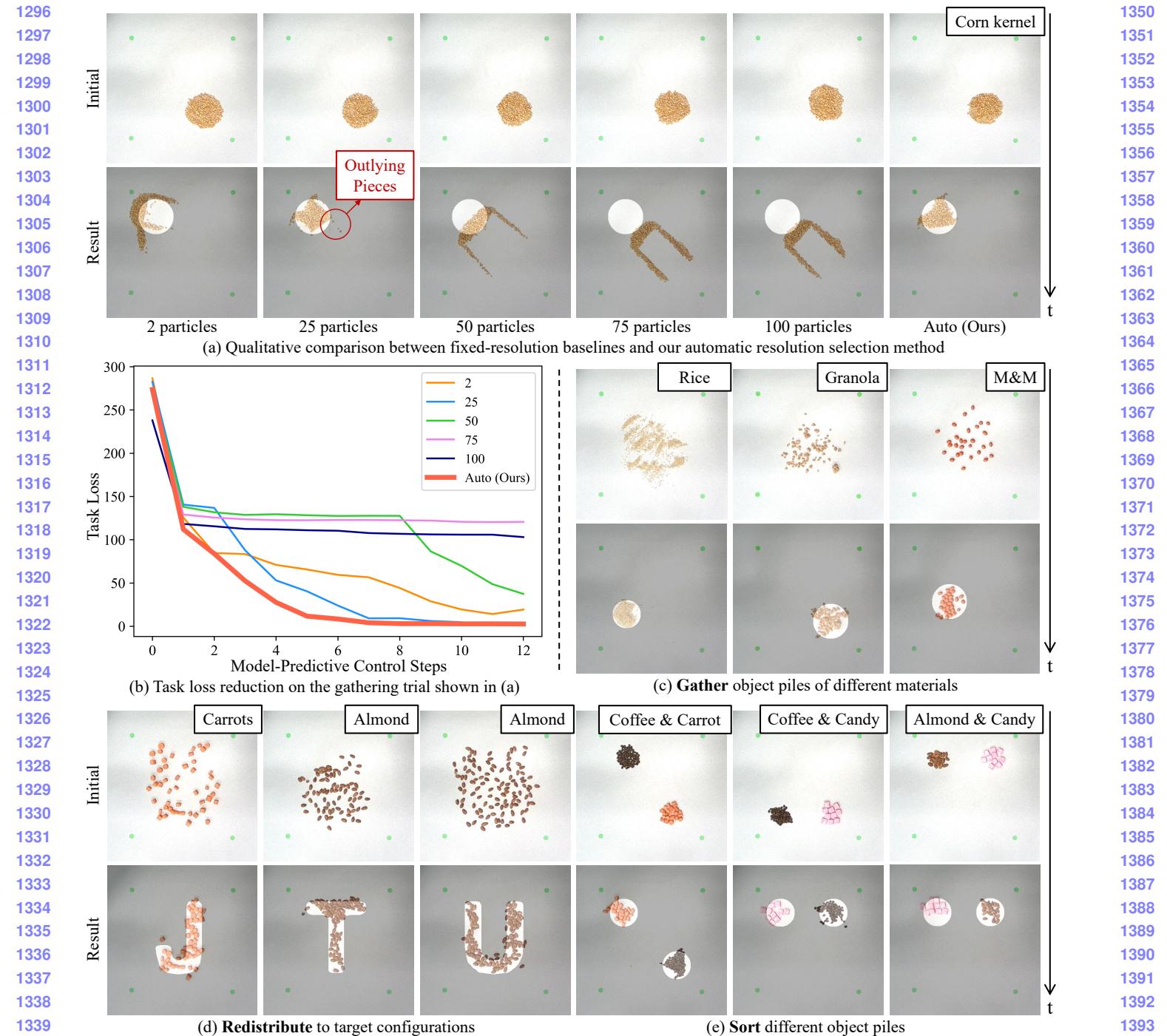


Figure 6. Qualitative results in the real world. (a) Our method outperforms in gathering corn pieces into the target region qualitatively. (b) Quantitative comparisons for the qualitative results in (a). Starting from similar initial configurations, our automatic resolution selection method performs the best throughout the MPC steps. (c) Our method is evaluated on the **Gather** task with various objects varying in scales and physical properties. (d) Our method can **redistribute** object pieces into complicated target configurations, such as letter shapes. (e) Our method can be combined with a high-level planner for more complex tasks, such as sorting different object piles into target regions.

the physics-based simulators for both manipulation and locomotion tasks [39, 73]. Audio feedback has also shown to be effective at estimating the amount and flow of granular

materials [12]. Our work instead focuses on three tasks using a unified dynamic-resolution graph dynamics to balance efficiency and effectiveness for real-world deployment.