

# Flipbot: Learning Continuous Paper Flipping via Coarse-to-Fine Exteroceptive-Proprioceptive Exploration

Chao Zhao<sup>\*1</sup>, Chunli Jiang<sup>\*1</sup>, Junhao Cai<sup>1</sup>, Michael Yu Wang<sup>1,2</sup>, Hongyu Yu<sup>1,2</sup>, and Qifeng Chen<sup>1</sup>

**Abstract**—This paper tackles the task of singulating and grasping paper-like deformable objects. We refer to such tasks as paper-flipping. In contrast to manipulating deformable objects that lack compression strength (such as shirts and ropes), minor variations in the physical properties of the paper-like deformable objects significantly impact the results, making manipulation highly challenging. Here, we present Flipbot, a novel solution for flipping paper-like deformable objects. Flipbot allows the robot to capture object physical properties by integrating exteroceptive and proprioceptive perceptions that are indispensable for manipulating deformable objects. Furthermore, by incorporating a proposed coarse-to-fine exploration process, the system is capable of learning the optimal control parameters for effective paper-flipping through proprioceptive and exteroceptive inputs. We deploy our method on a real-world robot with a soft gripper and learn in a self-supervised manner. The resulting policy demonstrates the effectiveness of Flipbot on paper-flipping tasks with various settings beyond the reach of prior studies, including but not limited to flipping pages throughout a book and emptying paper sheets in a box. The code is available here : <https://robot11.github.io/Flipbot/>

## I. INTRODUCTION

Deformable object manipulation has achieved notable progress in robotics. However, until now, robots could not match the generalization and robustness of humans in manipulating thin and flexible objects. One of these tasks is flipping book pages, as shown in Fig. 1, which requires singulating and grasping paper page by page. Humans can briskly turn pages of a book by watching the target and using the tactile sensations on their fingertips to adjust their actions. In this process, human instinctively combines exteroceptive and proprioceptive perception to accommodate the irregular paper thickness and physical properties, such as slipperiness, stiffness, and friction. Endowing robots to have such capability is a grand challenge in the field of robotics.

One of the foremost challenges in manipulating thin and flexible objects is incomplete and noisy perception [1]. For example, a stack of paper is unstable, and the contact between each layer is not observable. Therefore, the robot may have to perceive physical properties between paper, such as friction, and elasticity, to successfully singulate and grasp a sheet from a stack. Exteroceptive perception obtained from camera sensors is incomplete for such tasks and unreliable in real-world conditions. The depth sensors, which most existing works rely on, cannot distinguish the different layers of

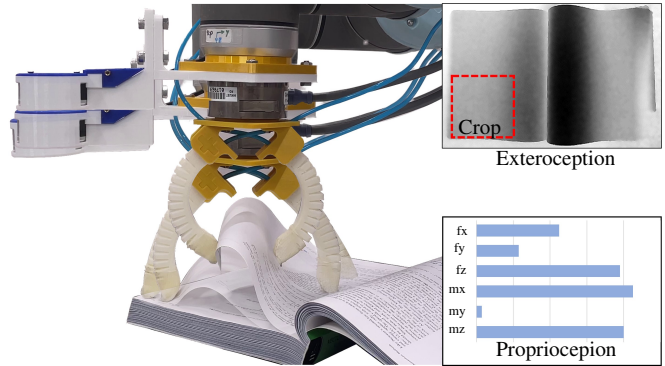


Fig. 1. A soft gripper with the learned policy flips a book. The time-lapse image depicts the operation of the gripper as it interacts with the book to singulate and grasp a piece of paper. The cropped depth image in the red line box located at the upper right corner presents the exteroceptive observation from the depth camera. The readings on the bottom right show the proprioceptive observation from the force-torque sensor.

stacked paper due to the paper thickness. Depth sensors are also inherently incapable of capturing the surface’s physical properties, such as hardness and flexibility [2]. Some works use tactile sensors as proprioception to estimate deformable objects’ physical properties. For example, [3] uses a high-precision tactile sensor to measure the geometry of the contact surface and the object’s hardness. [4] manipulates cables with a pair of robotic grippers using real-time tactile feedback. Nevertheless, high-precision tactile sensors are often expensive and require specific finger shapes to fit. In addition to the challenge in environment perception, manipulating thin and flexible objects may desire the gripper with the dexterity and compliance of human fingers, which further adds to the difficulty [5].

To address the above challenges, we present Flipbot, a self-supervised method for singulating and grasping paper-like deformable objects at unprecedented robustness, enabling continuous paper flipping. At its core, Flipbot is based on a principled solution integrating exteroceptive and proprioceptive perceptions into policy learning. We obtain proprioception from the Force/Torque (F/T) sensor readings and exteroception from a depth camera. We use a procedural motion, referred to as “Swipe” to actively interact with the environment. When a “Swipe” motion is applied to a piece of paper, the deformation brought about by the interaction between the finger and object reveals imperceptible physical characteristics like mass, flexural rigidity, and friction. Meanwhile, visual observation provides global information on the environment. We design a cross-sensory encoder to integrate exteroceptive and proprioceptive perceptions into an

<sup>\*</sup> Authors with equal contribution.

<sup>1</sup>The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong.

<sup>2</sup>HKUST Shenzhen-Hong Kong Collaborative Innovation Research Institute, Futian, Shenzhen.

implicit state representation. The encoder is trained end-to-end in a self-supervised manner as a part of policy learning. By incorporating exteroceptive-proprioceptive information into policy learning, the robot is able to discover the optimal policy for paper-flipping through continuous exploration. Furthermore, the reward signal for policy learning is derived from visual observation; Flipbot is fully trained by self-exploration without human demonstration or annotation.

The primary contribution of the presented work is the proposed new approach, Flipbot, for singulating and grasping paper-like objects. It achieves substantial improvements over the prior studies while maintaining exceptional robustness. Our extensive experiments show that Flipbot is able to perform page-flipping from the beginning to the end of a book accurately and consistently, and exhibits remarkable zero-shot generalization under conditions never encountered during training: novel paper materials such as coated and plastic paper and tasks such as emptying a box full filled with paper.

## II. RELATED WORK

Deformable object manipulation presents a persistent and enduring challenge within the field of robotics. Conventional analytic approaches rely on modeling object dynamics and then using model predictive control [6], or trajectory optimization [7] for manipulation. However, analytic approaches require substantial prior knowledge of geometry, and the physical properties of the object [1]. For example, [8] presents an approach for manipulating a piece of thin deformable object by analyzing the object’s internal energy exchange concerning object poses. And [9] proposes a close-loop shape control method utilizing visual markers, which limits the generality. Moreover, the high-dimensional state representation and complex dynamics of the deformable object provide additional challenges to generalizing novel objects and environments.

Recently, learning-based methods have become increasingly popular alternatives to perform deformable object manipulation. Most work [10], [11] learns the object dynamic from visual features rather than explicit modeling physical processes. For example, [12] encodes visual observation into latent space with self-supervision, followed by model-based planning. Another line of approach defines a set of primitives for deformable object manipulation and learns a mapping from image to predefined primitives [13]. Such image-to-primitive formulation has been applied across various tasks including manipulating rope [14], smoothing fabric [15], and blowing bags [16]. However, the physical information of the environment, which necessitates deformable object manipulation, is challenging to be obtained from visual perception. In this regard, [17] estimates the physical properties of fabric materials through a high-resolution tactile sensor, GelSight [18]. Further, [4] proposes an approach to manipulate a cable based on tactile feedback without vision sensory. [19] employs tactile sensors to manually collect data for training a classifier that can differentiate between towels with thicknesses of 1-3 layers. Then a heuristic approach is

used to consistently attempt to grasp specific layers of towels based on the classifier’s prediction outcomes. Nevertheless, tactile sensors alone are hard to provide global information about the environment, which inevitably restricts the range of manipulation or requires prior knowledge of objects.

More recently, a small number of papers have explored the use of soft grippers in deformable object manipulation, which is known for its ease of grasping objects without high precision control [8], [20]. The authors of [21] demonstrated a soft gripper system that is capable of handling a wide range of food products by reconfiguring fingers into different poses. In addition, [5] quantitatively indicates that the compliance of the soft gripper can facilitate the manipulation of thin deformable objects.

Compared with the above studies, our presented approach, Flipbot, incorporates exteroceptive and proprioceptive feedback in deformable object manipulation rather than relying on a single perception source. Flipbot thus combines the best of both worlds: the global information about the environment afforded by exteroception and the local information about physic property afforded by proprioception. Flipbot also leverages the compliance from a soft pneumatic gripper for performing dexterous behavior. The resulting policy has taken the real robot to various tasks surpassing prior published work in the field of deformable object manipulation.

## III. METHOD

The goal of Flipbot aims to empower robots to effectively singulate and grasp thin and flexible objects through exteroceptive and proprioceptive perception. Our key insight is that global information about positions and shapes on a large scale provided by vision and local information about contact and force provided by proprioceptive perception are indispensable parts of manipulating deformable objects like paper. Also, proprioception and exteroception fusion reveals physical information that helps robots better explore and make decisions. The overview of Flipbot is shown in Fig. 2.

First, we utilize a simple soft gripper for manipulation (see Fig. 4(c)). The natural compliance of the soft gripper provides unique benefits for manipulating thin and flexible objects while avoiding damage to the object. Another advantage is that the soft gripper has a more straightforward actuation strategy in movements such as bending the fingers, compared with fully actuated rigid grippers.

Then, we use a coarse-to-fine exploration process to obtain unobservable physical information about deformable objects. In this process, first, the depth camera provides a rough observation of the object. We then use a procedural motion “Swipe” and an F/T sensor to monitor the object’s state. One advantage of using the F/T sensor instead of a tactile sensor is that the force sensor can be assembled seamlessly with soft hands without a specific finger design.

Last, we use a cross-sensory encoder to fuse the proprioception and exteroception and use model-free reinforcement learning (RL) to learn the policy that avoids explicit modeling of diverse and frequent transitions in the contact state between the object and the soft hand.

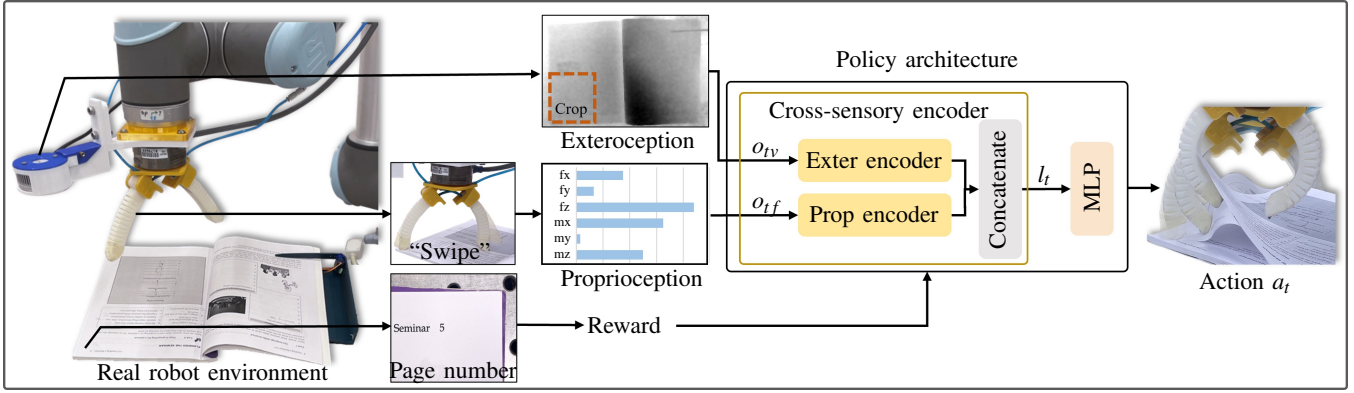


Fig. 2. **System Overview.** We train the policy using SAC in the real world. We follow a coarse-to-fine exploration process to obtain exteroception and proprioception. First, the camera captures the depth image, and the cropped area is used as extrinsic perception. Next, the soft finger “Swipe” on paper captures force  $(f_x, f_y, f_z)$  and torque  $(m_x, m_y, m_z)$  values from force sensors as proprioception. The RL agent receives the observations and predicts the actions to be performed by the robot, and receives the reward based on changes in page numbers.

### A. Problem Formulation

We formulate the problem of the paper-flipping as a Markov Decision Process (MDP). An MDP consists of four components: a state space  $S$ , an action space  $A$ , a reward function  $R(s_t, s_{t+1})$ , and a transition probability  $P(s_{t+1}|s_t, a_t)$ . In our framework, an agent uses a policy  $\pi(a_t|s_t)$  to select an action  $a_t$  for controlling the robot and receives rewards  $r_t$ . The goal of the reinforcement learning framework is to obtain the optimal policy  $\pi^*$ , which maximizes the expected discounted sum of rewards over a finite time horizon. To achieve this objective, we utilize the Soft Actor-Critic [22] (SAC) algorithm for training. SAC requires the learning of an actor network that maps observations to actions and a critic network that estimates the expected future rewards based on the input.

### B. Observations via Coarse-to-Fine Exploration

The state is defined as  $s_t = (o_{tv}, o_{tf})$ , where  $o_{tv}$  refers to the exteroceptive observation,  $o_{tf}$  refers to the proprioceptive observation, shown in Fig. 2. We deploy a coarse-to-fine exploration procedure with two steps for obtaining observations  $o_{tv}$  and  $o_{tf}$ . First, a wrist-mounted camera takes the environment’s point cloud  $p_t$  from a height and converts the point cloud to a depth image. We then use a  $60 \times 60$  resolution window to crop the depth image, as the exteroceptive observation  $o_{tv}$ . Next, we perform an exploratory “Swipe” motion, to obtain physical information about the contact surface between the paper and the finger. The robot first descends a certain distance that the finger of a soft hand approaches the surface of the top right corner of the paper diagonally, where the distance is calculated according to the point cloud  $p_t$ . Then, we give the soft gripper a positive air pressure so that fingers touch and interact with the paper. After this process, we record readings from the F/T sensor, including forces  $(f_x, f_y, f_z)$  in  $x, y, z$  axes and three simultaneous torques  $(m_x, m_y, m_z)$  about the same axes. Thus, the proprioceptive observation  $o_{tf}$  is defined as a tuple of  $(f_x, f_y, f_z, m_x, m_y, m_z)$ . Fig. 3 shows forces and torques after “Swipe” on different pages in the book. By incorporating an F/T sensor and exploratory action,  $o_{tf}$  latently contains rich information related to contact states

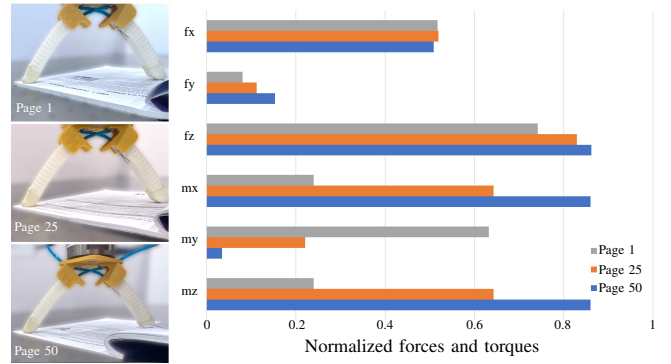


Fig. 3. Visualization of forces and torques after “Swipe” on the different page numbers.

between the fingers and the object, such as gripper-object friction.

### C. Action and Reward

After the coarse-to-fine exploration procedure, the robot predicts the action based on observations to singulate and grasp the paper. The action includes a gripper displacement, denote as  $(x_t, z_t, \theta_t)$ , as shown in Fig. 4(a). The gripper displacement refers to the relative difference between the current pose after the “Swipe” exploration procedure and the desired one. Specifically,  $x_t \in [-6mm, 6mm]$  is the relative displacement on the line  $\alpha$  connecting the two fingertips, where  $\alpha$  belongs to the longitudinal plane  $A$  formed by two fingers.  $\theta_t \in [0^\circ, 3^\circ]$  is the orientation of the gripper about the normal  $\beta \perp A$ .  $z_t \in [-6mm, 6mm]$  is the the relative displacement on the line  $\gamma$ , where  $\gamma \perp (\alpha \times \beta)$ . Furthermore, an additional action component  $\Lambda$  is utilized to govern the closing or opening of the gripper. Operationally, we control the gripper aperture by commanding the pressure change. Thus, the action is formally defined as  $a_t = (x_t, z_t, \theta_t, \Lambda)$ , where each coordinate of the action is discretized based on the characteristics of the workspace.

At the end of an episode, the reward is given, 1 for successfully flipping a single layer of paper and 0 for otherwise. In other words, flipping two or more layers of paper simultaneously is treated as a failure. The reward is

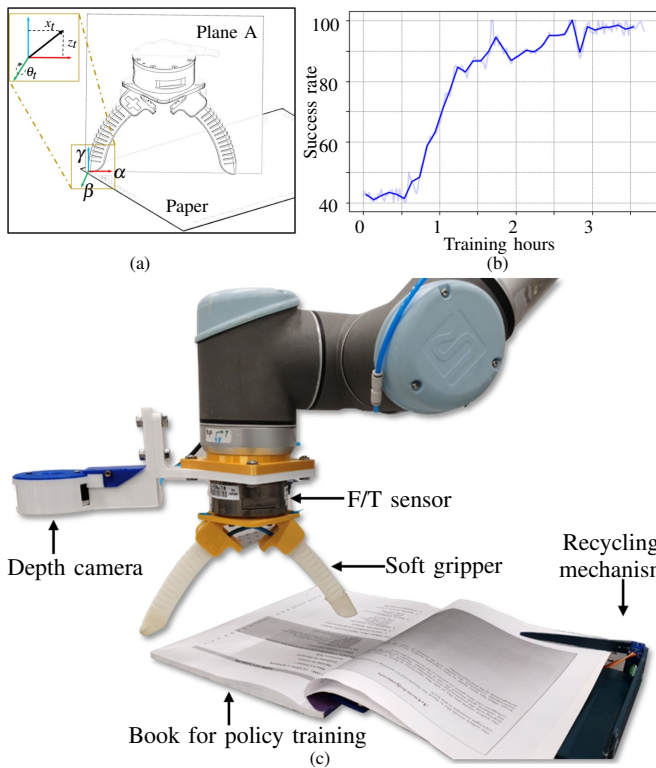


Fig. 4. (a): Visualization of our action coordinate system. (b): Success rate curve of our policy training. (c): Our hardware setting for policy training.

automatically determined by identifying page numbers on the book, which we describe further in Sec. III-E

#### D. Policy architecture

The policy  $\pi(a_t|s_t)$  is modeled with a cross-sensory encoder and a multilayer perceptron (MLP) block, as shown in Fig. 2. The cross-sensory encoder takes the exteroceptive observation  $o_{tV}$  and proprioceptive observation  $o_{tF}$  as inputs and embeds them into a latent vector, which represents the abstraction of proprioception and exteroception. More specifically,  $o_{tV}$  is processed by a global pooling layer and concatenated with  $o_{tF}$  to be a vector of size  $1 \times 7$ . Then, the concatenated vector is fed into subsequent an MLP block to compress inputs to a more compact representation  $l_t$ . At last, the  $l_t$  is fed through the subsequent MLP layer to predict actions.

#### E. Training via self-supervision

We train the policy in a real robot platform. Fig. 4(c) shows our hardware setting for training, including the following major components: a Universal Robot 10 robot arm equipped with a 3D printed thermoplastic polyurethane soft gripper, an ATI gamma F/T sensor, and an Intel Realsense L515 depth camera, as well as a recycling mechanism. During the whole training, we only use a book assembled with printer paper, as shown in Fig. 4(c). We train our model through trial-and-error with the following procedure:

At each training step, the robot starts to execute coarse-to-fine exploration from an initial pose. In the process of “Swipe”, the wrist-mounted camera captures an RGB-D image. The depth channel is used to construct exteroceptive observation  $o_{tV}$ , and the page number  $n_t$  is recognized from



Fig. 5. A subset (9 of 27) of our test scene settings. Columns from left to right show different paper materials: printing paper, coated paper, and plastic paper. Rows from top to bottom show different test scenarios and workspace tilt angles

RGB channels for reward calculation. The robot then descends a certain distance that the finger approaches the paper’s surface to perform an exploratory action to obtain the physical observation  $o_{tF}$  from the readings of F/T sensors. After this, the robot downloads the latest policy parameters from the optimizer to predict action  $a_t$  and executes. We automatically calculate rewards according to the change of page numbers without human intervention, the reward  $r_t$  is 1 if  $n_{t+1} = n_t + 2$ , otherwise 0. The page number identification benefits from Tesseract [23]. At last, the generated episode is added to a replay buffer, and the optimizer sampling from this replay buffer to update the policy. We use the Adam optimizer with a learning rate of  $3 \times 10^{-3}$ . The robot then continuously collects episodes until it reaches the last page of the book, at which point the book is reset to the first page again using the recycling mechanism. In this way, human intervention is kept at a minimum during the training process.

The final model training took four hours, with the learning curves for the training presented in Fig. 4(b).

## IV. EXPERIMENTS

We design a set of experiments in real-world settings to evaluate the system’s generalization ability to novel object physical parameters and the advantage of using exteroceptive and proprioceptive exploration. For all following experiments, we use the same robot hardware setting and the same model trained with the book assembled from printer paper, described in Sec. III-E. The system’s performance is evaluated on its generalization to unseen paper types (i.e., flipping different types of paper when only trained on printer paper) and unseen scenarios (e.g., emptying paper in a box) and its efficiency (i.e., the speed and accuracy of paper-flipping).

**Scene setup:** We investigate the performance of our system across various object settings and scene configurations. In total, we have 27 different test scenes with the combination of test scenarios, paper types and tilt angles. We test with the following three scenarios:

- Full Book page flipping. It is a similar scenario as in policy training, where the robot needs to flip book pages one by one throughout the book.

TABLE I  
RESULTS OF EXPERIMENTS IN THE REAL WORLD.

Method	Tilt angle	Full Book page flipping						Paper-box emptying						Single paper grasping					
		Printer SR	Paper PPH	Coated SR	Paper PPH	Plastic SR	Paper PPH	Printer SR	Paper PPH	Coated SR	Paper PPH	Plastic SR	Paper PPH	Printer SR	Paper PPH	Coated SR	Paper PPH	Plastic SR	Paper PPH
Flex&Flip [8]		72%	223	77%	239	52%	161	69%	214	82%	254	49%	152	83%	260	91%	282	74%	229
Flipbot-w/o prop	0°	85%	264	93%	288	66%	205	81%	251	91%	282	60%	186	95%	295	<b>98%</b>	<b>304</b>	85%	264
<b>Flipbot</b>		<b>94%</b>	<b>291</b>	<b>96%</b>	<b>298</b>	<b>82%</b>	<b>254</b>	<b>90%</b>	<b>279</b>	<b>94%</b>	<b>291</b>	<b>68%</b>	<b>211</b>	<b>99%</b>	<b>307</b>	<b>98%</b>	<b>304</b>	<b>92%</b>	<b>285</b>
Flex&Flip [8]		76%	236	74%	229	44%	136	62%	192	72%	223	42%	130	80%	248	87%	270	76%	236
Flipbot-w/o prop	30°	88%	273	87%	270	63%	195	84%	260	88%	273	55%	171	85%	264	92%	295	86%	267
<b>Flipbot</b>		<b>93%</b>	<b>288</b>	<b>91%</b>	<b>282</b>	<b>72%</b>	<b>223</b>	<b>88%</b>	<b>273</b>	<b>91%</b>	<b>282</b>	<b>62%</b>	<b>192</b>	<b>92%</b>	<b>285</b>	<b>95%</b>	<b>295</b>	<b>90%</b>	<b>279</b>
Flex&Flip [8]		64%	198	56%	174	47%	192	56%	174	58%	180	38%	118	84%	260	82%	254	83%	257
Flipbot-w/o prop	60°	76%	236	72%	223	62%	192	77%	239	70%	217	58%	179	86%	267	85%	264	91%	282
<b>Flipbot</b>		<b>84%</b>	<b>260</b>	<b>82%</b>	<b>253</b>	<b>70%</b>	<b>217</b>	<b>82%</b>	<b>254</b>	<b>80%</b>	<b>248</b>	<b>66%</b>	<b>205</b>	<b>96%</b>	<b>298</b>	<b>92%</b>	<b>285</b>	<b>94%</b>	<b>291</b>

\* SR stands for success rate.

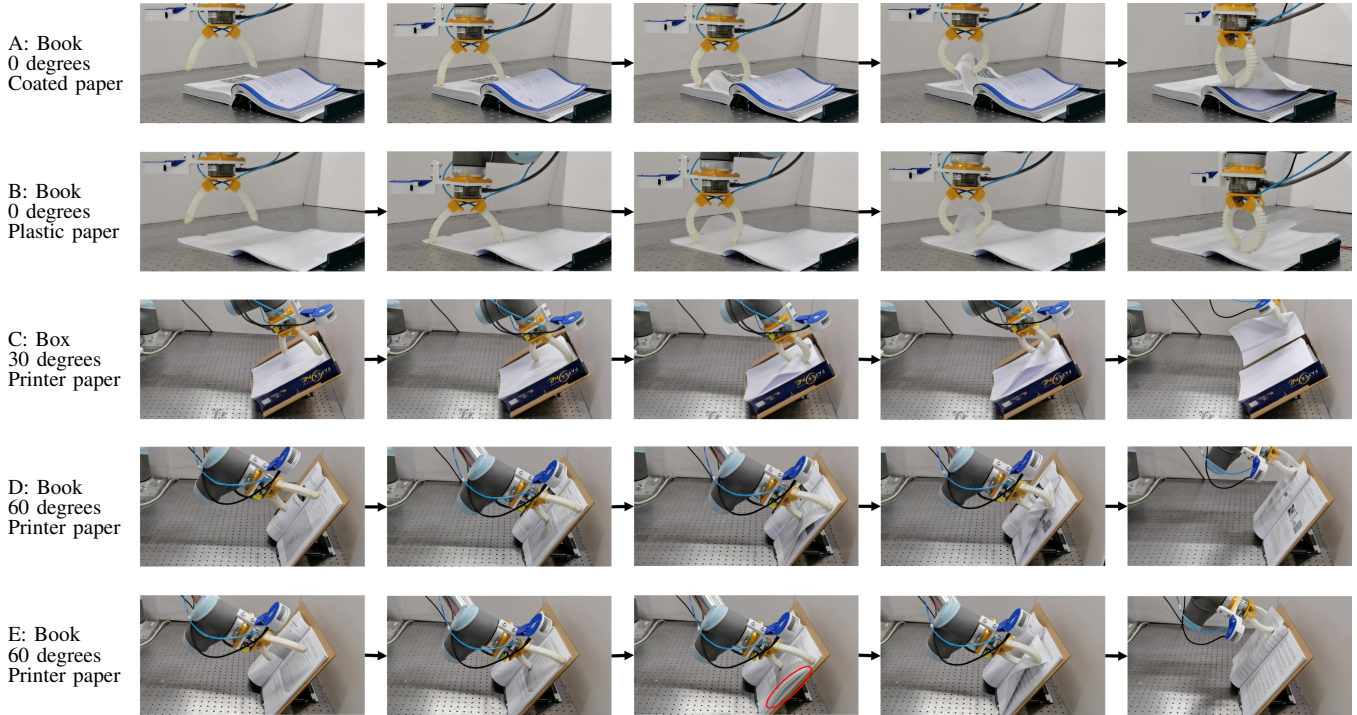


Fig. 6. Flipbot performs paper-flipping in different scenes. A-D: Flipbot successfully singulates and grasps a piece of paper in various settings; E: Flipbot fails to singulate and grasp a piece of printer paper with a 60-degree tilt angle. The circled area in red denotes that two layers of paper were flipped.

- Paper-box emptying. The robot grasps each sheet one by one from a pile of paper dumped into a box until emptying it. This is more challenging than the book setup because the physical interaction between the paper is more complex without the constraints of the spine.
- Single paper grasping. The robot grasps a single piece of paper lying on a flat surface.

In each scenario, we use three types of paper that have different physical properties, including the printer paper, coated paper, and plastic paper. The physical property of printer paper is the same as we have used during training, which has the highest friction coefficient among the three types. The coated paper and plastic paper are unseen paper types. The coated paper has the lowest friction coefficient and the plastic paper has medium friction coefficient. The detailed physical properties of these three paper types are

Physical properties	Printer paper (seen type)	Coated paper (unseen type)	Plastic paper (unseen type)
Static Coefficient of Friction	<b>0.462±0.0087</b>	0.283±0.0104	0.334±0.0066
Kinetic Coefficient of Friction	<b>0.417±0.0542</b>	0.174±0.0229	0.259±0.0263
Young's Modulus in Machine Direction(GPa)	<b>2.84±0.17</b>	2.62±0.14	1.54±0.23
Density (g/m <sup>2</sup> )	102.5±2.32	59.8±0.93	<b>385.4±1.74</b>
Thickness (mm)	0.096±0.006	0.057±0.012	<b>0.151±0.017</b>

shown in Tab. II. Meanwhile, we also vary tilt angles (0, 30, 60 degrees) for the workspace to test the effect of gravity on paper flipping.

**Metric:** We utilize two evaluation metrics for validating algorithm performance: success rates (successful paper

flips/total attempts) and PPH (successful paper flips per hour). The success of paper flipping for each attempt is measured by whether the gripper detaches and flips strictly one piece of paper. For example, in the book page flipping task, the robot detaches and flips two pieces of paper simultaneously is considered a failure. PPH is the product of the speed of flipping in an hour and the success rate, which includes the time of perception, network inference, and robot execution in enabling paper-flipping manipulation. It is important to note that our Flipbot implementation is not optimized for high-speed execution; thus, the reported PPH is solely used to compare relative performance.

**Algorithm comparisons:** We compare with the following methods:

- **Flex&Flip [8]:** it simplifies a piece of paper as a linear object and uses a physical model to analyze the motion. Its original version could only grasp a single piece of paper lying on a flat surface. We adapt and extend the physical model provided by the authors and hardcode the thickness of different paper types to allow for multi-layered paper flipping.
- **Flipbot-w/o prop:** policy learns from only exteroceptive sensory (i.e., depth camera), which directly maps the visual observation to action.
- **Flipbot:** policy learns with coarse-to-fine exteroceptive-proprioceptive exploration, which is the full non-ablated method we propose in this article.

#### A. Experimental Results

**Comparison to prior work.** We first compare the performance of our approach with Flex&Flip [8] with different paper types and scenarios (row 1 vs. row 3 in Tab. I). Note that Flex&Flip [8] is the state-of-the-art method for single-layer paper grasping, and we extend it to multi-layered paper scenarios (i.e., paper-box emptying and full book page flipping). In the single paper grasping case, Flipbot performs better (+16%) than Flex&Flip [8] on printer paper. The advantage is much more pronounced in multi-layered paper cases, with Flipbot outperforming Flex&Flip [8] around 20%. In all three test scenarios, quantitative results in Tab. I suggest that our method (Flipbot) maintains comparable success rates on unseen paper types (i.e., coated and plastic paper) with respect to the seen paper type (i.e., printer paper). In contrast, the performance of Flex&Flip [8] on the plastic paper type degrades significantly (up to -20%) on unseen paper types.

**Effectiveness of exteroceptive-proprioceptive exploration.** We conduct controlled experiments to evaluate the contribution of exteroceptive-proprioceptive exploration quantitatively. The proprioceptive perception provides information on the unobservable physical features, facilitating policy learning effectiveness. As a result, compared with Flipbot-w/o prop that does not use proprioceptive, Flipbot achieved a higher success rate. Quantitative results in Tab. I indicate that compared to Flipbot-w/o prop, the success rate of Flipbot increases at most 24% and at least 4% across test cases.

**Generalization to novel tilt angles of workspace.** In this experiment, we investigate the generalization ability of these

methods to gravity changes by varying tilt angles (0, 30, 60 degrees) of the workspace (see Fig. 6C-D). In different tilt angle setups, detaching a single sheet of paper becomes more challenging as the physical properties between the different layers of the paper change with the direction of gravity. Quantitative results in Tab. I show that the performance of our learned policy degrades slightly as the tilt angle increases. We hypothesize this happened since the physics in these test scenes differ from the training, increasing the difficulty of generalization. Nevertheless, Flipbot still outperforms other methods in terms of success rate and PPH in all test cases.

Overall, our experimental evaluation demonstrates that Flipbot is an efficient approach for paper-flipping tasks. We find the exteroceptive and proprioceptive perceptions are essential for paper-flipping, particularly for singulating and detaching a sheet from a pile of paper. The learned policy has been demonstrated to outperform state-of-the-art methods and is also applicable to tasks beyond the reach of prior studies, such as turning pages throughout a book. Our work is not without limitations. First, when the working area is at a larger inclination angle, the friction between the paper tends to be smaller. Hence, multiple layers of paper are easy to be grasped simultaneously (see Fig. 6E). Also, two layers of paper sometimes stick together. We assume it happens because of Van der Waals forces. A dual-arm system may be essential to address this issue, suggesting exciting opportunities for future study.

## V. CONCLUSION

We have presented a novel solution for singulating and grasping thin and flexible deformable objects that utilize the cross-sensory encoding of exteroceptive and proprioceptive perceptions, which we term Flipbot. Meanwhile, the system takes advantage of the under actuation and compliance of the soft pneumatic actuator to control contact forces precisely for the singulation of a thin layer of deformable objects. We deploy the algorithm on a real-robot system and show that integrating exteroceptive and proprioceptive inputs can effectively facilitate deformable object manipulation. Extensive controlled experiments demonstrated the robustness and effectiveness of Flipbot. Beyond the experiment results, our work extends frontiers in deformable object manipulation, and the methodology presented in this work can have broad applications. A future direction is to extend the proposed approach to long-horizon deformable object manipulation tasks, such as origami folding, cleaning messy desktops, collecting mail and letters, etc.

## REFERENCES

- [1] J. Zhu, A. Cherubini, C. Dune, D. Navarro-Alarcon, F. Alambeigi, D. Berenson, F. Ficuciello, K. Harada, J. Kober, X. LI, J. Pan, W. Yuan, and M. Gienger, "Challenges and outlook in robotic manipulation of deformable objects," *IEEE Robotics Automation Magazine*, pp. 2–12, 2022.
- [2] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning robust perceptive locomotion for quadrupedal robots in the wild," *Science Robotics*, vol. 7, no. 62, p. eabk2822, 2022.

- [3] W. Yuan, M. A. Srinivasan, and E. H. Adelson, "Estimating object hardness with a gelsight touch sensor," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 208–215.
- [4] Y. She, S. Wang, S. Dong, N. Sunil, A. Rodriguez, and E. Adelson, "Cable manipulation with a tactile-reactive gripper," *The International Journal of Robotics Research*, vol. 40, no. 12-14, pp. 1385–1401, 2021.
- [5] C. B. Teeple, J. Werfel, and R. J. Wood, "Multi-dimensional compliance of soft grippers enables gentle interaction with thin, flexible objects," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 728–734.
- [6] F. Allgöwer and A. Zheng, *Nonlinear model predictive control*. Birkhäuser, 2012, vol. 26.
- [7] S. Zimmermann, R. Poranne, and S. Coros, "Dynamic manipulation of deformable objects with implicit integration," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 4209–4216, 2021.
- [8] C. Jiang, A. Nazir, G. Abbasnejad, and J. Seo, "Dynamic flex-and-flip manipulation of deformable linear objects," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 3158–3163.
- [9] Y. Guo, X. Jiang, and Y. Liu, "Deformation control of a deformable object based on visual and tactile feedback," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 675–681.
- [10] H. Shi, H. Xu, Z. Huang, Y. Li, and J. Wu, "Robocraft: Learning to see, simulate, and shape elasto-plastic objects with graph networks," *Robotics: Science and Systems (RSS)*, 2022.
- [11] B. Shen, Z. Jiang, C. Choy, L. J. Guibas, S. Savarese, A. Anandkumar, and Y. Zhu, "Acid: Action-conditional implicit visual dynamics for deformable object manipulation," *Robotics: Science and Systems (RSS)*, 2022.
- [12] W. Yan, A. Vangipuram, P. Abbeel, and L. Pinto, "Learning predictive representations for deformable objects using contrastive estimation," in *Conference on Robot Learning*. PMLR, 2022.
- [13] H. Ha and S. Song, "Flingbot: The unreasonable effectiveness of dynamic manipulation for cloth unfolding," in *Conference on Robot Learning*. PMLR, 2022, pp. 24–33.
- [14] P. Sundaresan, J. Grannen, B. Thananjeyan, A. Balakrishna, M. Laskey, K. Stone, J. E. Gonzalez, and K. Goldberg, "Learning rope manipulation policies using dense object descriptors trained on synthetic depth data," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 9411–9418.
- [15] Y. Avigal, L. Berscheid, T. Asfour, T. Kröger, and K. Goldberg, "Speedfolding: Learning efficient bimanual folding of garments," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 1–8.
- [16] Z. Xu, C. Chi, B. Burchfiel, E. Cousineau, S. Feng, and S. Song, "Dexterity: Deformable manipulation can be a breeze," in *Proceedings of Robotics: Science and Systems (RSS)*, 2022.
- [17] W. Yuan, Y. Mo, S. Wang, and E. H. Adelson, "Active clothing material perception using tactile sensing and deep learning," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 4842–4849.
- [18] W. Yuan, S. Dong, and E. H. Adelson, "Gelsight: High-resolution robot tactile sensors for estimating geometry and force," *Sensors*, vol. 17, no. 12, p. 2762, 2017.
- [19] S. Tirumala, T. Weng, D. Seita, O. Kroemer, Z. Temel, and D. Held, "Learning to singulate layers of cloth using tactile feedback," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 7773–7780.
- [20] J. Hughes, U. Culha, F. Giardina, F. Guenther, A. Rosendo, and F. Iida, "Soft manipulators and grippers: a review," *Frontiers in Robotics and AI*, vol. 3, p. 69, 2016.
- [21] J. H. Low, P. M. Khin, Q. Q. Han, H. Yao, Y. S. Teoh, Y. Zeng, S. Li, J. Liu, Z. Liu, P. V. y Alvarado, *et al.*, "Sensorized reconfigurable soft robotic gripper system for automated food handling," *IEEE/ASME Transactions On Mechatronics*, 2021.
- [22] T. Haamoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.
- [23] R. Smith, "An overview of the tesseract ocr engine," in *Ninth international conference on document analysis and recognition (ICDAR 2007)*, vol. 2. IEEE, 2007, pp. 629–633.