

More Than a Feeling: Learning to Grasp and Regrasp using Vision and Touch

Presenter: 王思源 苏兆文 林沛君 赵瑞涵 刘子羽

2023.3.2



AncoraSIR.com



SUSTech
Southern University
of Science and Technology

Motivation and Main Problem

High-level description of problem being solved

The problem being solved is how a **robot** can learn to use **tactile information** to **iteratively** and **efficiently** adjust its **grasp**, as **humans** heavily rely on **tactile feedback** during the process of grasping an object, but most recent robotic grasping work has been based **only on visual input**, which cannot easily benefit from feedback after **initiating contact**.

Motivation and Main Problem

Why is the problem important?

- *significance towards general-purpose robot autonomy*

Grasping is a fundamental aspect of many tasks, and the ability to adjust grasp based on tactile feedback can greatly improve the **success rate** and **efficiency** of the task.

- *potential application and societal impact of the problem*

By developing a method for **integrating touch sensing** into grasping, robots can become more **versatile and adaptable**, and can potentially perform a **wider range of tasks** in a variety of environments.

Motivation and Main Problem

Technical challenges arising from the problem

- *the role of the AI and machine learning in tackling this problem*

The authors propose a method based on learning an **action-conditioned** grasping model, which is trained **end-to-end** in a **self-supervised** manner using a robot to autonomously collect **grasp attempts**.

This approach allows the system to learn from its **own experiences** and adjust its behavior accordingly, leading to **improved** grasping performance **over time**.

In addition, **integrating tactile sensing** into the model needs advanced machine learning techniques to deal with the **intermittent** and **changing tactile input**. The authors use **GelSight** sensors to capture rich touch information, and their model **incorporates** this information robustly to the changing scene.

Motivation and Main Problem

why prior approaches didn't already solve? & Key insights

- *High-level idea of why prior approaches didn't already solve*

The incorporation of touch sensing into robotic grasping has been challenging due to **hardware limitations** such as sensor sensitivity and cost, as well as the difficulty of **integrating tactile inputs** into standard control schemes, so that the predominant input modalities currently used in the robotic grasping are **vision and depth**, which are **not always precise and robust**.

- *Key insights of the proposed work*

An end-to-end **action-conditional** model that learns regrasping policies from **raw visuo-tactile data**

Problem Setting

Problem formulation, key definitions and notations

- *Problem formulation*
 - I. Incorporating **vision and touch sensing** into **action-conditional** models.
 - II. Develop a **learning-based** approach to further improve the **grasping** performance.
- *Key definitions*
 - I. Regrasping: Adjusting its grip if necessary to achieve a better grasp
 - II. End-to-End: A system designed to work seamlessly from start to finish.
 - III. Action-Conditioned: Takes into account the current conditions or environment and adjusts its actions accordingly to achieve a desired outcome.

Related Work & Limitations of Prior Work

Learning to Grasp

- Analytic grasping models

- "Robot grasp synthesis algorithms: A survey" (Shimoga. 1996)
- "Data-driven grasping" (Goldfeder and Allen. 2011)
- "From caging to grasping" (Rodriguez, Mason and Ferry. 2012)

- Visual Data-driven approaches

- "Learning to grasp" (Saxena. 2015)
- "Deep learning for grasping" (Saxena. 2015)
- "Leveraging big data for grasping" (Schaal. 2015)
- "Deep learning a grasp function for grasping under gripper pose uncertainty" (Johns, Leutenegger and Davison. 2016)

- Model misspecification and unmodled factors reduce the actual performance

- A limited ability to reason about contact forces, pressures, and compliance

Most of these methods rely on selecting grasp configurations in advance, before coming into contact with the target

object

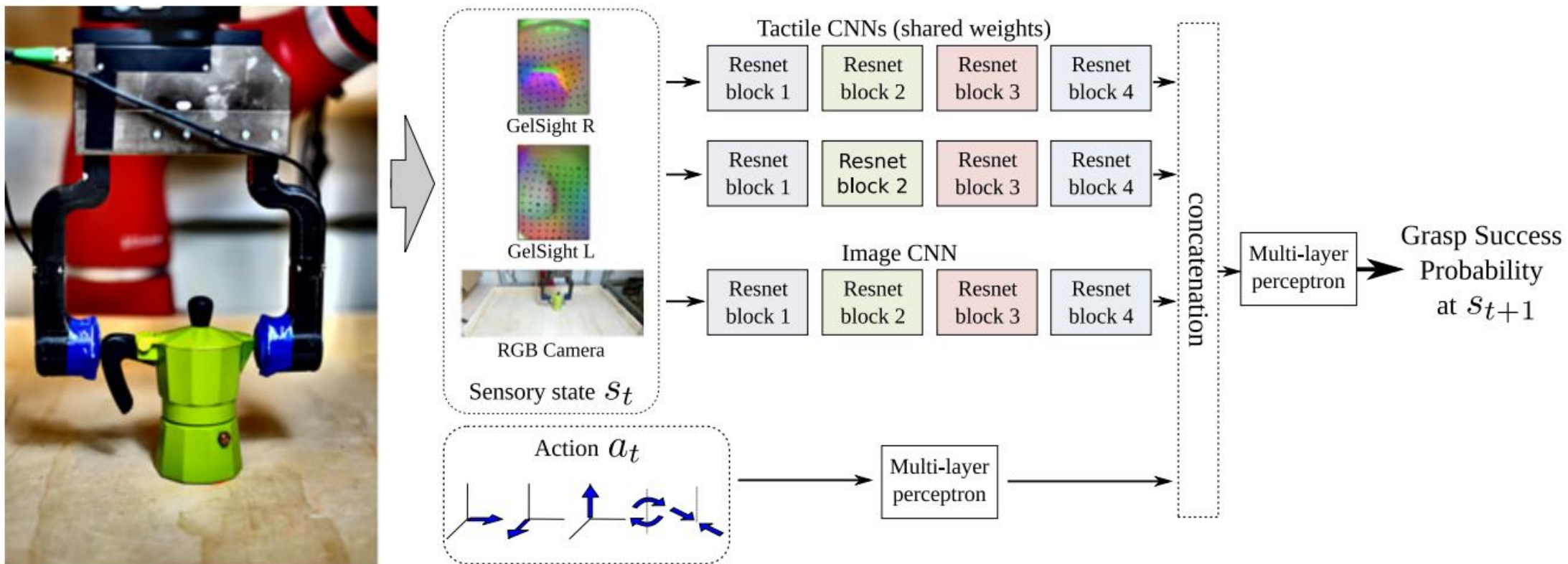
Related Work & Limitations of Prior Work

Tactile Sensors in Grasping

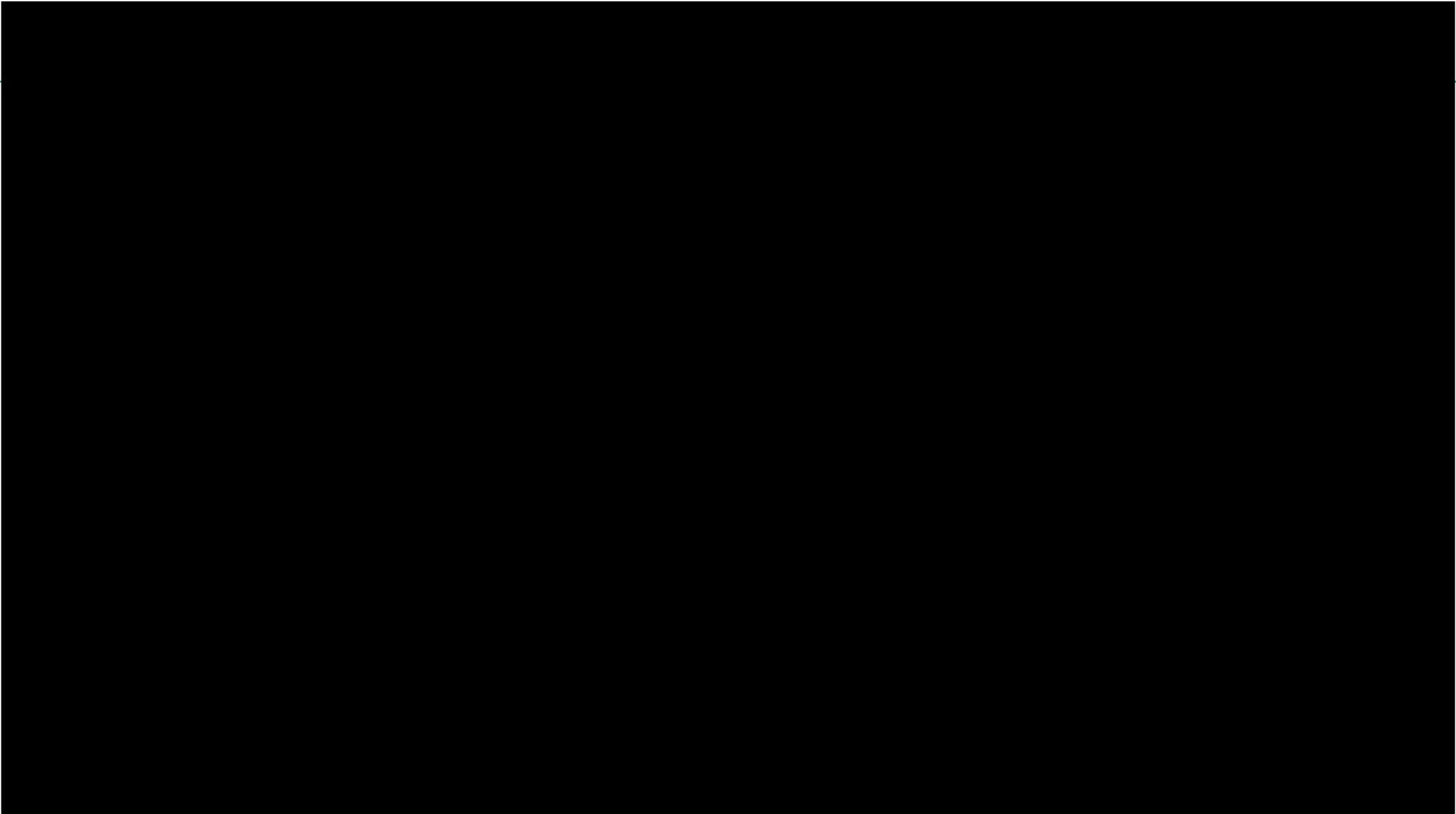
- The use of tactile sensors on grasp
 - "Assessing grasp stability based on learning and haptic data" (Bekiroglu, Laaksonen, Jorgensen, Kyrki and Kragic. 2011)
 - "Stable grasping under pose uncertainty using tactile feedback" (Dang and Allen. 2016)
- Model-based methods for integrating visual and tactile information
 - "Learning to assess grasp stability from vision touch and proprioception" (Bekiroglu. 2012)
 - "Robotic grasping using visual and tactile sensing" (Guo, Sun, Yang and Xi. 2017)
- regrasping policy
 - "Self-supervised regrasping using spatio-temporal tactile features and reinforcement learning" (Chebotar, Hausman, Su, Sukhatme and Schaal. 2016)
- Only to estimate the stability of an ongoing grasp, not produce a stable new one
- depend on the accurate models of the robots and objects to grasp
- require the data collection to be on-policy and to be intertwined with the policy update

Proposed Approach

Generation of the prediction model



Action-conditioned visuo-tactile model network architecture



Proposed Approach

Innovation of the framework and Algorithm

- Exploit **vision** and **tactile feedback** to produce a stable new grasp
- Learn entirely **end-to-end** from raw inputs, so it does not require any prior model or transition function.
- Utilize an **action-conditioned model** not a policy, which can match any data collected and adapt to a different objective at evaluation time.

Theory

Markov decision process (MDP)

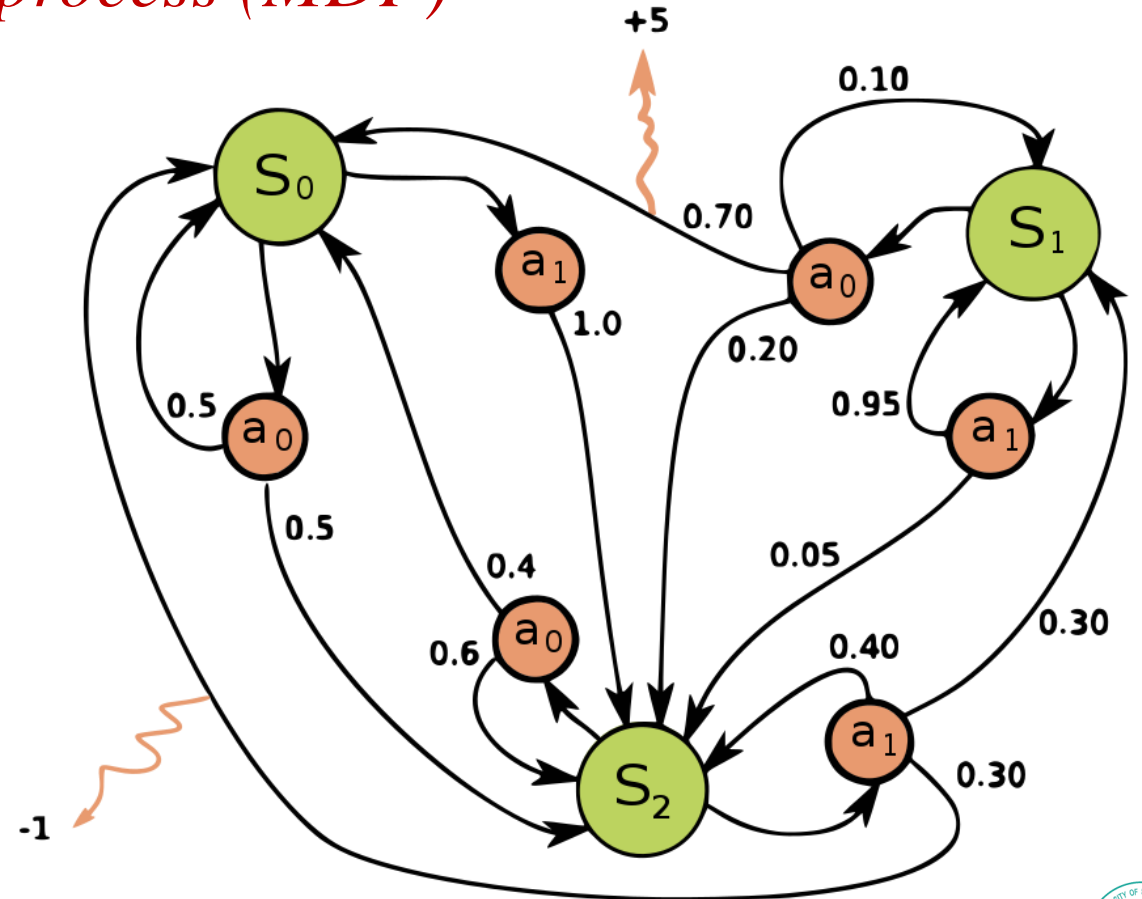
- The **Markov decision process** (MDP) is a mathematical model of sequential decisions and a dynamic optimization method.
- A Markov decision process is a 4-tuple (S, A, P_a, R_a) , where:
 - S is a set of states called the state space,
 - A is a set of actions called the action space (alternatively, A_s is the set of actions available from state S),
 - $P_a(s, s') = \Pr(s_{t+1} = s' \mid s_t = s, a_t = a)$ is the *probability* that action a in state s at time t will lead to state s' at time $t + 1$,
 - $R_a(s, s')$ is the immediate reward (or expected immediate reward) received after transitioning from state s to state s' , due to action a .

Theory

Markov decision process (MDP)

- Given the robot's current visuo-tactile observations s_t at time t , and an action a , we predict the probability that, after applying the action, the gripper will be in a configuration that leads to a successful grasp at time $t + 1$.
- the action that maximize the expected probability of success of the grasp

$$a_t = \arg \max_a f(s_t, a)$$



Theory

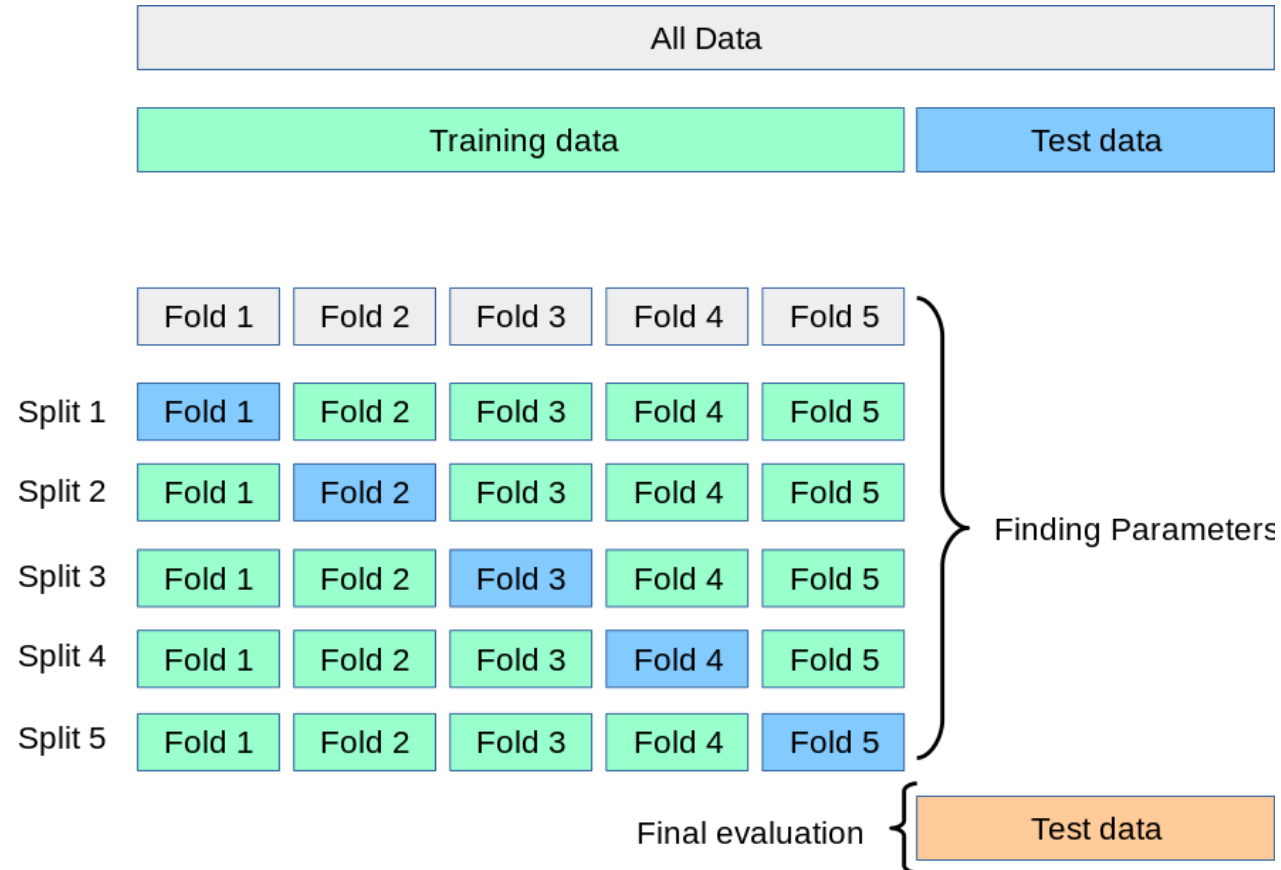
K-fold cross-validation

- *Why K-fold cross-validation?*
- To evaluate the performance of some model on a dataset, we need to measure how well the predictions made by the model match the observed data. The most common way to measure this is by using the mean squared error (MSE).
- The test MSE gives us an idea of how well a model will perform on data it hasn't previously seen.
- The drawback of using only one testing set is that the test MSE can vary greatly depending on which observations were used in the training and testing sets.

Theory

K-fold cross-validation

- How to do K-fold cross-validation*



Experimental Setup

Hardware setups

7-DoF Sawyer arm



Weiss WSG-50 parallel gripper



GelSight sensors



Experimental Setup

Data Collection

- In each trial, depth data from the front Kinect was used to approximately identify the starting position of the object and enclose it within a **cylinder**.
- Set the end-effector (x, y) coordinates to the position of the center of the cylinder plus a **small random perturbation**, and set its **height** to be a **random value** between the floor and the height of the cylinder.
- Moreover, the orientation ϕ and the gripping force F are *randomized* .
- After moving to the chosen position and orientation, the gripper close and attempt to lift the object and wait in the air for 4s. If the object was still in the gripper at the end of this time, the robot would place the object back at a randomized position, and a new trial would start.

Experimental Setup

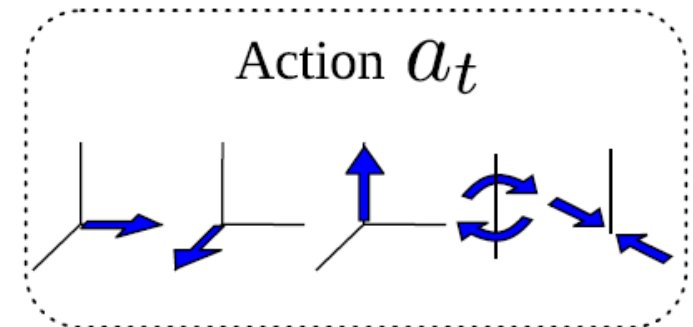
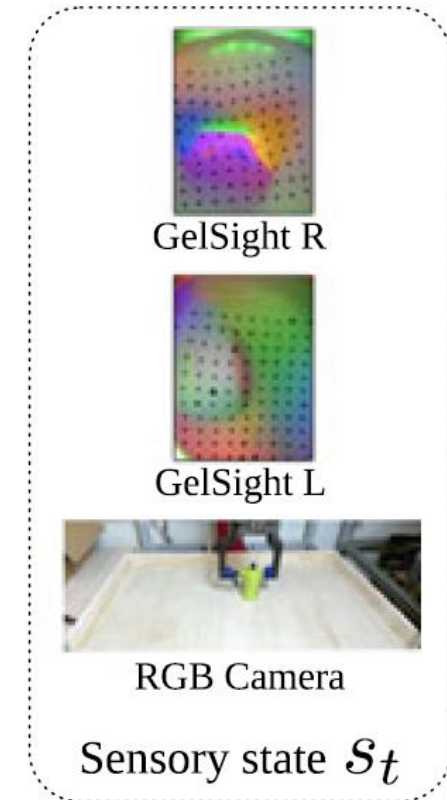
Data Collection

- The labels for this were also automatically generated using deep neural network classifiers trained to detect contacts using the raw GelSight images observed.
- Additional manual labeling on a small set of samples were performed, for which the automatic classification was borderline ambiguous, or in the rare cases when a visual inspection would indicate a wrong label.
- As the gripper moves from one position to another, the locations that it moves to along the way can provide additional data points for training.
- 6450 grasping trials from over 65 training objects are collected, and the datasets contains 18070 examples.

Experimental Setup

Datasets

- Raw visuo-tactile observations s are acquired from tactile sensors and the RGB camera. Each action a directs the gripper to a new pose relative to its current pose.
- let $o(s_t, a) \in \{0, 1\}$ be the binary grasp outcome at time $t + 1$ resulting from executing action a from grasp state s_t : if $o(s, a)$ is 1, the grasp is successful.
- At training time, the robot performs random trials to collect state-action-outcome tuples $(s_i, a_i, o_i) \in X$.



Experimental Evaluation

Model Evaluation

- Question1: Can the model successfully learn to predict future grasp success for novel objects?
- Question2: Whether the model learns to use actions to predict future outcomes, conditional on a relative adjustment from the current grasp?
- Evaluate by comparing the performance of a number of variations of the current model, though K-fold($K = 3$) cross-validation, using the dataset of grasps.

Experimental Evaluation

Model Evaluation

- Models:
 - This model
 - This model without action
 - Vision-only
 - Tactile-only

TABLE I
K-FOLD (K = 3) CROSS-VALIDATION ACCURACY OF THE
DIFFERENT MODELS TRAINED WITH 18,070 DATA POINTS

Model	Accuracy (mean \pm std. err.)
Chance	62.80% \pm 0.85%
Vision (+ action)	73.03% \pm 0.24%
Tactile (+ action)	79.34% \pm 0.66%
Tactile + Vision (+ action)	80.28% \pm 0.68%
Tactile + Vision (no action)	76.43% \pm 0.42%

Experimental Evaluation

Robot Grasp Evaluation

- Evaluated the learned models on the robot, testing the robot grasp a given object after executing a series of regrasp actions, Each grasp begins by **randomly sampling** an end-effector position and angle, used in data collection.
- **Compare this model with the vision-only variant of the model**, then use the learned models to select the next grasp, by solving the optimization equation.
- The optimization is performed by randomly sampling 4900 actions, plus 100 additional actions sweeping over the grasping force interval.
- Moreover, if the predicted grasp success probability is above the desired threshold(0.9), the re-grasp also includes lifting the object.

Experimental Evaluation

Robot Grasp Evaluation

- **Baseline:** also evaluated against an approach that fits a cylinder around the object using depth data and subsequently attempt to grasp the centroid of the object using a constant grasping force of 10 N.

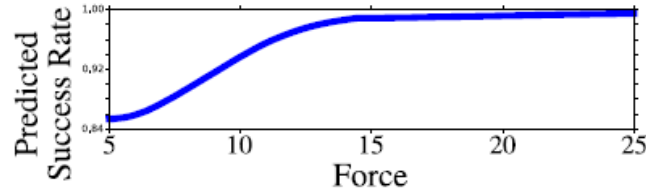
TABLE II
DETAILED GRASPING RESULTS USING DIFFERENT POLICIES FOR THE "EASY" AND "HARD" TEST OBJECTS

	Objects											Average grasp success	
		215g	160g	40g	125g	125g	65g	135g	30g	380g	140g		10g
"Easy" set	Methods	% grasp success (# success / # trials)											
	Vision only	76% (38/50)	70% (7/10)	60% (6/10)	50% (5/10)	50% (5/10)	90% (9/10)	40% (4/10)	60% (6/10)	90% (9/10)	10% (1/10)	100% (10/10)	63.2%
	Tactile + Vision	95% (95/100)	100% (10/10)	100% (10/10)	100% (10/10)	90% (9/10)	100% (10/10)	90% (9/10)	100% (10/10)	80% (8/10)	90% (9/10)	90% (9/10)	94.0%
	Cylinder fitting	90% (18/20)	90% (18/20)	80% (16/20)	55% (11/20)	100% (20/20)	100% (20/20)	90% (18/20)	75% (15/20)	35% (7/20)	20% (4/20)	100% (20/20)	75.9%
	Objects											Average grasp success	
		230g	120g	195g	50g	70g	85g	38g	165g	65g	340g		110g
"Hard" set	Methods	% grasp success (# success / # trials)											
	Vision only	60% (6/10)	80% (8/10)	30% (3/10)	30% (3/10)	80% (8/10)	40% (4/10)	60% (6/10)	50% (5/10)	50% (5/10)	50% (5/10)	20% (2/10)	50%
	Tactile + Vision	80% (8/10)	100% (10/10)	50% (5/10)	80% (8/10)	90% (9/10)	70% (7/10)	100% (10/10)	40% (4/10)	60% (6/10)	80% (8/10)	60% (6/10)	73.6%
	Cylinder fitting	95% (19/20)	100% (20/20)	35% (7/20)	100% (20/20)	90% (18/20)	15% (3/20)	90% (18/20)	85% (17/20)	15% (3/20)	15% (3/20)	95% (19/20)	66.8%

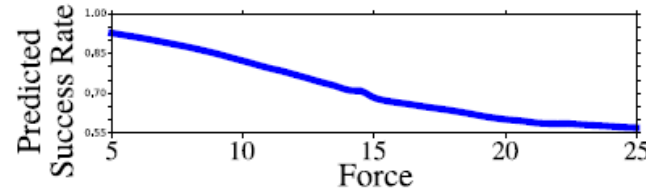
Experimental Results

We now examine qualitatively: what strategies has our model learned and what behaviors does it produce?

1. Grasping Force



(a) Stable grasp.



(b) Unstable grasp.

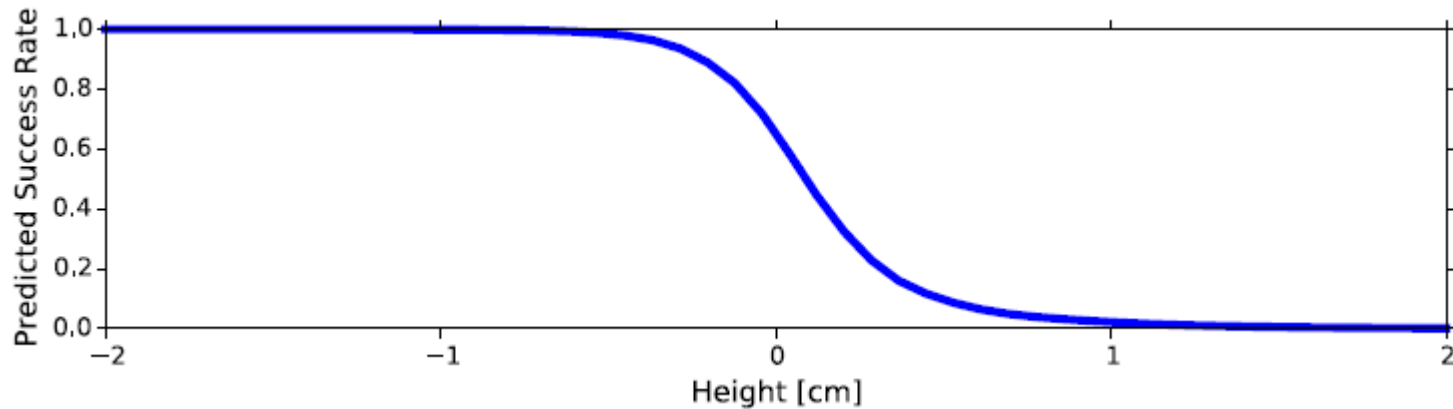
The model learned that, when stably in contact with the object, there is a correlation between force applied and success rate. However, for unstable grasps, the model learned that increasing the grasp force might misplace the object and result in an unsuccessful grasp.



Experimental Results

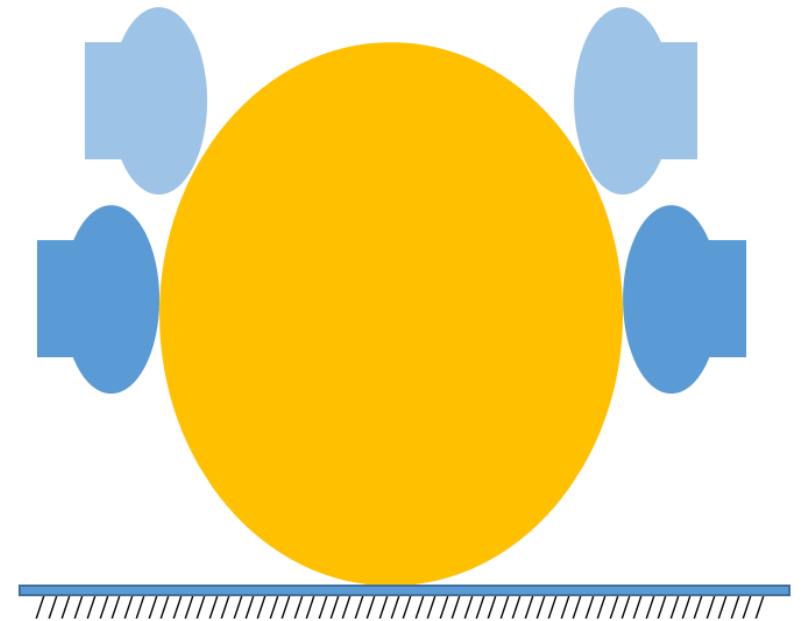
We now examine qualitatively: what strategies has our model learned and what behaviors does it produce?

1. *Grasping Force*
2. *Height and Center-of-Mass*



Example of predicted grasp success rate varying the height of the fingers.

The model learned that decreasing the height of the fingers generally increases the success rate.

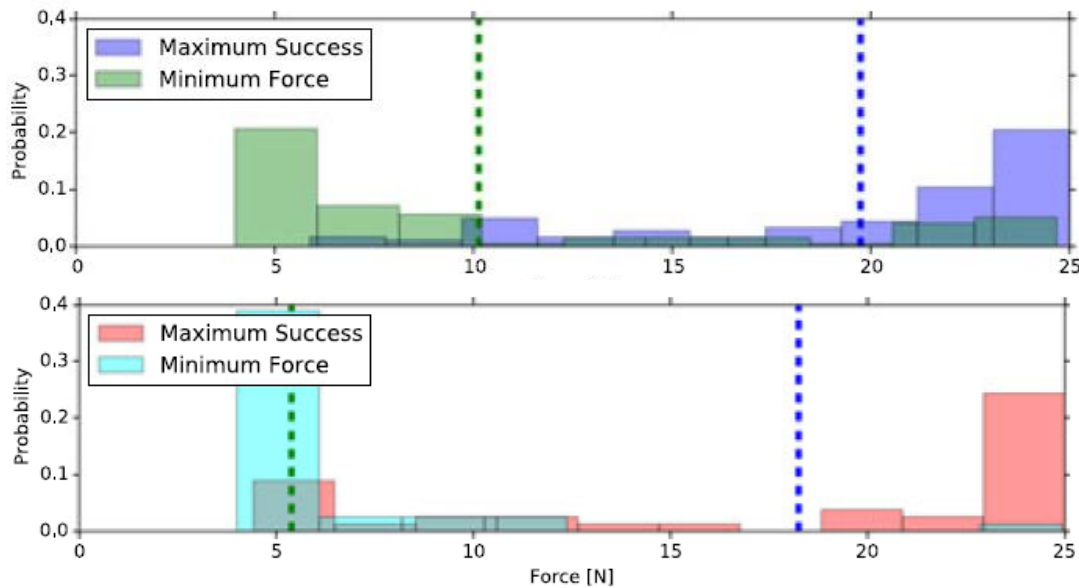


Experimental Results

We now examine qualitatively: what strategies has our model learned and what behaviors does it produce?

1. *Grasping Force*
2. *Height and Center-of-Mass*
3. *Minimum Force Grasp*

We typically do not use the strongest grasp possible, but rather employ the minimum amount of contact force, out of consideration for energy consumption and object fragility



Histogram and mean of the forces applied in the successful grasps.

(a) Although the success rates for the two Tactile+Vision policies are similar (95% maximum success vs 94% minimum force), the mean force applied is significantly reduced when using the minimum force policy (10 vs 20 N).

(b) The success rates for the Vision only policies is lower at 76%, but again the mean force applied is significantly reduced when using the minimum force policy (6 vs 18 N).

Discussion of Results

Strengths and weaknesses:

- **Minimize the applied forces while maintaining a high success rate**
 - Exploit the visuo-tactile action-conditioned model to minimize the applied forces while maintaining a high success rate.
- **Compliant, small and irregular objects.**
 - Based on these experiments, the largest improvements in performance of the model seem to happen in the presence of compliant objects, and objects where it is difficult to visually ascertain a good grasp, such as small or irregular objects.

Discussion of Results

Strengths and weaknesses:

- **Complex:**
 - The controller requires both visual and tactile operation and is complicated to design
- **High-bandwidth:**
 - The input information required by the controller is visual and tactile information.
 - Both visual and tactile information are large matrix and high frequency data

Critique / Limitations / Open Issues

Limitations and solutions:

- The action-conditioned model only makes single-step predictions, and does not perform information-gathering actions
- We can use reinforcement learning models based on policy gradient algorithms to enable robots to acquire more information through perception and exploration while executing actions
- Only consider relatively coarse actions
- A model using fine-grained actions could more delicately manipulate the object before the grasp, and potentially react to slippage during the lift-off, however, it also increases the computing resource consumption and bandwidth requirements.

Prospect for Future Work

Future work that could be built on?

1. Action-conditional model and tactile sensing give better **adaptivity** and **evolvability** than traditional vision-based grasping.
2. The ability of **data gathering** actions during grasping helps model evolution, similar to what Tesla is doing to enhancing their auto-driving algorithm.
3. Work of **fine-grained** and **continuous** control loop could be built onto the action-conditional model for future work. (potentially using torque control)

Prospect for Future Work

Future work that could be built on?

4. Use proper optimization algorithms to **reduce computational cost** during grasping. (for example, applying PSO to optimize decision making process when grasping)*

5. Introduce “**human in the loop**” into training dataset, and invent an interactive method for human to “teach”, which may be helpful for the trained model to refer to human’s tricks to grasp things.*

Extended Readings

Some later papers that go further from this: [Dimensions](#)

[A Recognition Method for Soft Objects Based on the Fusion of Vision and Haptics](#)

Teng Sun, Zhe Zhang, Zhonghua Miao, Wen Zhang
2023, Biomimetics - Article

[Vision-Tactile Fusion Based Detection of Deformation and Slippage of Deformable Objects During Grasping](#)

Wenjun Ruan, Wenbo Zhu, Kai Wang, Qinghua Lu, Weichang Yeh, Lufeng Luo, Caihong Su, Quan W...
2023, Cognitive Systems and Information Processing - Chapter

[Heuristic grasping of convex objects using 3D imaging and tactile sensing in uncalibrated grasping scenarios](#)

Augusto Gómez Eguíluz, I. Rañó
2022, Expert Systems with Applications - Article

Publication metrics

[About](#)

Dimensions Badge



150 Total citations
97 Recent citations

43 Field Citation Ratio
n/a Relative Citation Ratio

Altmetric



Blogs (1)
Twitter (10)
Patents (5)
Reddit (1)
Video (1)
Mendeley (355)



Document history

2018-07-19 Published online

Funded by

Nvidia (United States)

Sources of this paper:

[More Than a Feeling: Learning to Grasp and Regrasp Using Vision and Touch | IEEE Journals & Magazine | IEEE Xplore](#)

AncoraSIR.com

Siyuan Wong & 2022.3.9

More Than a Feeling: Learning to Grasp and Regrasp using Vision and Touch



SUSTech
Southern University
of Science and Technology

Summary

- **Main problem:**

- Learning grasping and re-grasping skills using visual and tactile senses.

- **Why significant and hard?**

- Multiple senses help robots grasp better; Using iterative method is more adaptive.
 - Hardware limitations and the difficulty of multi-variable data fusion

- **Key limitations of prior work:**

- Poor model universality and transferability, limited data fusion skill.

- **Key insight:**

- Multiple-sensing and proper data fusion give better ability to grasp
 - Action-conditional models help robots to "think about the future"

Thanks !



AncoraSIR.com



SUSTech
Southern University
of Science and Technology