

# Feel the Bite: Robot-Assisted Inside-Mouth Bite Transfer using Robust Mouth Perception and Physical Interaction-Aware Control

Rajat Kumar Jenamani  
Cornell University  
rj277@cornell.edu

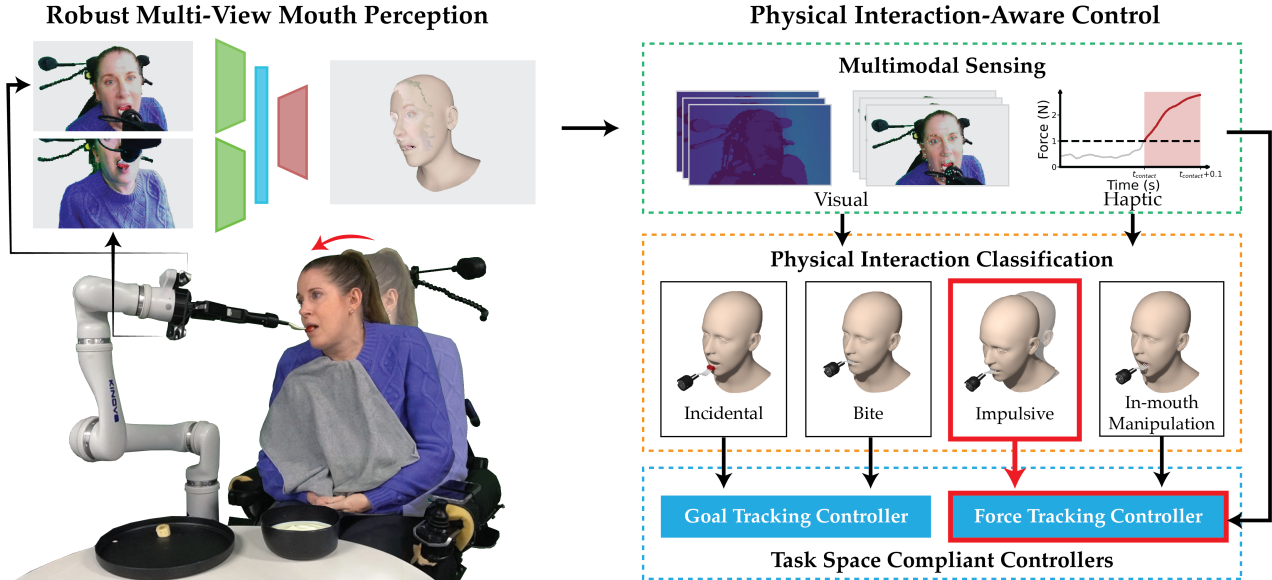
Daniel Stabile  
Cornell University

Ziang Liu  
Cornell University

Abrar Anwar  
University of Southern California

Katherine Dimitropoulou  
Columbia University

Tapomayukh Bhattacharjee  
Cornell University



**Figure 1:** We propose an inside-mouth bite transfer system that uses two key components, (i) robust multi-view mouth perception, and (ii) physical interaction-aware control, to successfully feed care recipients with diverse mobility limitations.

## Abstract

Robot-assisted feeding can greatly enhance the lives of those with mobility limitations. Modern feeding systems can pick up and position food in front of a care recipient’s mouth for a bite. However, many with severe mobility constraints cannot lean forward and need direct inside-mouth food placement. This demands precision, especially for those with restricted mouth openings, and appropriately reacting to various physical interactions – incidental contacts as the utensil moves inside, impulsive contacts due to sudden muscle spasms, deliberate tongue maneuvers by the person being fed to guide the utensil, and intentional bites. In this paper, we propose an inside-mouth bite transfer system that addresses these challenges with two key components: a multi-view mouth perception pipeline robust to tool occlusion, and a control mechanism that employs multimodal time-series classification to discern and react to different

physical interactions. We demonstrate the efficacy of these individual components through two ablation studies. In a full system evaluation, our system successfully fed 13 care recipients with diverse mobility challenges. Participants consistently emphasized the comfort and safety of our inside-mouth bite transfer system, and gave it high technology acceptance ratings – underscoring its transformative potential in real-world scenarios. Supplementary materials and videos can be found at: [emprise.cs.cornell.edu/bitetransfer](https://emprise.cs.cornell.edu/bitetransfer).

## CCS Concepts

- **Human-centered computing** → **Accessibility technologies**;
- **Computer systems organization** → **Robotics**.

## Keywords

Assistive Robots, Physical Human-Robot Interaction, Robot-Assisted Feeding

## ACM Reference Format:

Rajat Kumar Jenamani, Daniel Stabile, Ziang Liu, Abrar Anwar, Katherine Dimitropoulou, and Tapomayukh Bhattacharjee. 2024. Feel the Bite: Robot-Assisted Inside-Mouth Bite Transfer using Robust Mouth Perception and Physical Interaction-Aware Control. In *Proceedings of the 2024 ACM/IEEE International Conference on Human-Robot Interaction (HRI '24)*, March 11–14, 2024, Boulder, CO, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3610977.3634975>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

HRI '24, March 11–14, 2024, Boulder, CO, USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0322-5/24/03...\$15.00

<https://doi.org/10.1145/3610977.3634975>

## 1 INTRODUCTION

Eating is an Activity of Daily Living (ADL) [32]. Losing the ability to independently feed oneself can be devastating, and may elicit feelings of shame and dependence among care recipients [31, 41, 44]. In the United States alone, approximately one million adults with mobility limitations need assistance with eating [46]. Feeding is also one of the most time-consuming ADLs for caregivers [13], placing a significant burden on them [19, 39]. Robot-assisted feeding systems have the potential to enhance the quality of life of care recipients [10] and reduce caregiver burden.

Robot-assisted feeding can be broken down into two stages: bite acquisition and bite transfer [40]. Bite acquisition involves acquiring a bite-sized food item with an appropriate utensil [21, 26–28], whereas bite transfer entails moving it from above the plate to the care recipient’s mouth [7, 24, 43]. In this paper, we focus on bite transfer. Most prior works bring a food item in front of the mouth of a care recipient, who is then expected to lean forward to take the bite. However, for those with severe upper body and neck mobility limitations, leaning forward for an outside-mouth bite transfer can be functionally impossible, necessitating the direct placement of food in their mouths. Even for care recipients who can lean forward, having to repeatedly make this movement can be exhausting. In this work, we present a robot-assisted feeding system that can perform inside-mouth bite transfer and demonstrate its utility for assisting a diverse group of care recipients with severe mobility limitations.

Care recipients who need inside-mouth bite transfers often have complex medical conditions – including limited mouth opening [29], involuntary movements (spasms), and requirement of food transfer at a specific location inside their mouth – which makes feeding them extremely challenging. Involuntary motions, diverse in their type and occurrence, demand alert and adaptable feeding strategies. For example, if an involuntary forward spasm occurs just as the utensil approaches the care recipient’s mouth, the caregiver must quickly retract the utensil to avoid harm. If such a motion happens with the utensil already inside the mouth, the caregiver should comply with the motion, ensuring the utensil moves in harmony with it to avoid injury. Another layer of intricacy emerges when food must be placed precisely within the mouth at a preferred transfer location. In these instances, care recipients may use their tongue to guide the utensil, necessitating that caregivers recognize and act upon this cue. In addition to this, physical interactions also occur when slight errors in sensing/control lead to incidental contacts, and when care recipients intentionally bite down on the food. Consequently, feeding individuals with such complex needs requires precision in perception and control, and identifying and appropriately reacting to various types of physical interactions.

In this paper, we make the following contributions:

**Inside-Mouth Bite Transfer System.** We propose a system that leverages two key components – robust multi-view mouth perception, and physical interaction-aware control – to address the aforementioned challenges and successfully feed care recipients with severe mobility limitations food inside their mouth.

**Robust Multi-View Mouth Perception Method.** Existing methods for mouth perception in robot-assisted feeding rely on a single in-hand camera view, and thus face challenges during inside-mouth bite transfer due to significant occlusion from the feeding

utensil. To address this, we introduce a novel multi-view mouth-perception approach that is robust to tool occlusion (Section 3.1). This method allows for uninterrupted real-time mouth perception throughout the transfer process, enabling detection and adaptation to both voluntary and involuntary head movements.

**Physical Interaction-Aware Control Method.** We propose an interaction-aware controller that uses multimodal sensing for classifying the nature of physical interactions in real-time and reacts accordingly (Section 3.2). Adopting a data-driven approach, we collect a multimodal dataset comprising various types of physical interactions that can occur and train time-series classification models. Based on the detected interaction type, we switch between a goal-tracking controller and a force-tracking controller.

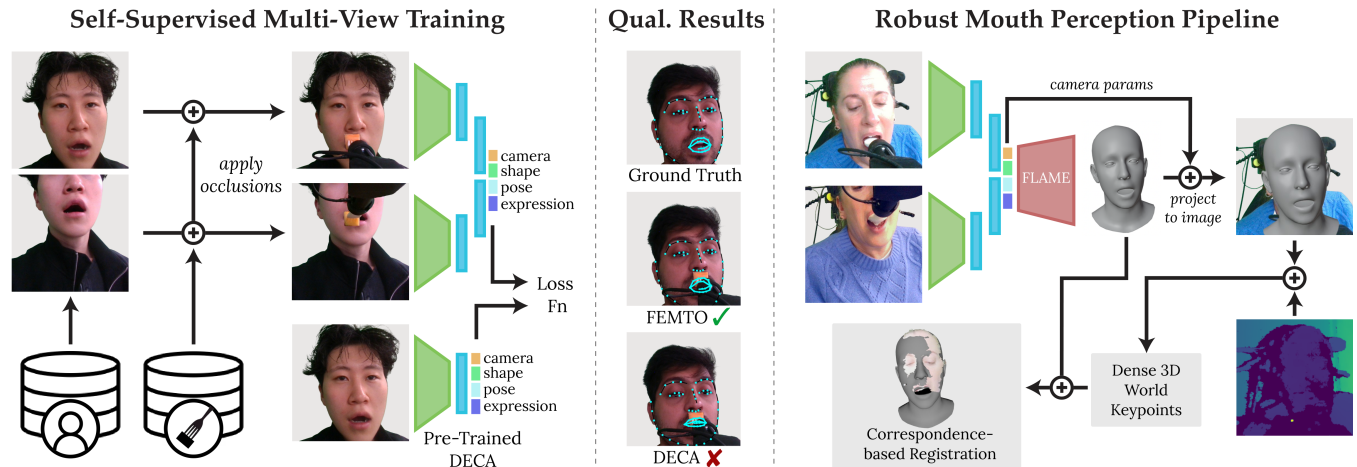
**Evaluation of System Components.** We demonstrate the necessity of both novel components through evaluations with baselines and two ablation studies involving participants without mobility limitations (Section 4).

**Full System Evaluation with Care Recipients.** We demonstrate our system’s efficacy through a user study with 13 care recipients with diverse medical conditions, all necessitating assistance in feeding. Findings suggest users perceive our system as safe and comfortable, and view the technology favorably as measured using a Technology Acceptance Model (TAM) survey [18].

## 2 RELATED WORK

**Robot-Assisted Feeding.** While many commercial feeding systems [3, 4] exist, their limited autonomy and need for manual trajectory programming for bite transfer hinder widespread adoption and retention. Over recent years, various robot-assisted feeding systems with autonomous transfer capabilities have been proposed [7, 11, 24]. However, they assume that care recipients can lean forward, causing the robot to stop at a predetermined distance from their mouth. More recently, a few systems have explored inside-mouth bite transfer [42, 43], but they perceive the user’s mouth pose only once (initially) and do not continuously track it during transfer. This requires users to remain static throughout the entire transfer process, a challenging demand for many individuals with pronounced mobility impairments. These systems also do not consider various types of physical interactions that can arise during inside-mouth bite transfer, such as incidental contacts, impulsive contacts, and in-mouth manipulation. To the best of our knowledge, our work is the first to demonstrate autonomous inside-mouth bite transfer for care recipients having complex medical conditions such as limited mouth opening, involuntary movements, and precise food placement requirements for feeding.

**Mouth Perception.** To effectively feed care recipients with complex medical conditions, it is necessary to accurately perceive their mouth for the whole duration of inside-mouth bite transfer. Contemporary robot-assisted feeding systems use single in-hand cameras, and localize the mouth by either projecting 2D mouth landmarks onto an aligned depth image [43], or fitting a fixed-head model to sparse 3D facial landmarks [24]. However, near the mouth, significant occlusion from the utensil leads to noisy depth data and inaccuracies in single-view perception of 2D landmarks, rendering these methods ineffective for continuous mouth tracking. In this work, we propose a mouth perception pipeline that leverages multiple in-hand cameras and parameterized head models to be robust



**Figure 2: The self-supervised multi-view training of FEMTO (left) enables robust mouth perception under tool occlusion (center). Our perception pipeline (right) leverages FEMTO to accurately estimate mouth pose despite noisy depth due to occlusion.**

to tool occlusion, thus enabling uninterrupted and accurate mouth perception. Among parameterized head models [9, 17, 36], Faces Learned with an Articulated Model and Expressions (FLAME) [36] is of special interest as it separates the representation of identity, pose, and facial expression. Detailed Expression Capture and Animation (DECA) [22], is state-of-the-art for FLAME parameter estimation, and takes a single RGB image as input. However, this reliance on monocular RGB images makes DECA vulnerable to inaccuracies in our use case where significant occlusions around the mouth are present. To overcome this limitation, we propose a novel encoder based on DECA’s architecture that integrates data from multiple cameras, ensuring robust prediction of FLAME parameters.

**Physical Interaction-Aware Control.** Prior work in contact classification distinguishes between incidental collisions and intentional task contacts during collaborative human-robot manipulation [15, 23, 25, 33, 38]. Classical approaches [33] use the spectral norm of external force/torque signals in a specific frequency range for classification. In contrast, several other methods leverage learning-based approaches like Support Vector Machines with time-series features derived from physical contact models [15, 23, 25], or time-series models such as RNNs [38]. While these methods typically use haptic data, we are inspired by works that integrate multiple sensing modalities [20, 35, 37, 45]. To the best of our knowledge, our work is the first to explore multimodal (visual + haptic) sensing for physical human-robot interaction classification. It is also the first to categorize physical interaction types beyond incidental and intentional, by additionally identifying impulsive interactions and further distinguishing intentional interactions into in-mouth manipulation and bite interactions.

### 3 Inside-Mouth Bite Transfer System

Our inside-mouth bite transfer system (Figure 1) consists of a Kinova Gen3 6 DoF robotic arm [2] with a Robotiq 2F-85 gripper [5] grasping a custom-built feeding tool. It uses two Intel RealSense D415 RGBD cameras mounted on the robot’s wrist, one above and one below the utensil, for visual input, and a 6-axis ATI Nano25 F/T sensor [1] for haptic feedback. Our system leverages two key novel components - robust multi-view mouth perception (Section 3.1) and physical interaction-aware control (Section 3.2).

#### 3.1 Robust Multi-View Mouth Perception

We require accurate mouth perception for the entire duration of inside-mouth bite transfer to successfully feed care recipients having small mouth openings, involuntary motions, and requirements of food placements at specific locations inside their mouth. However, real-time mouth perception using a single in-hand camera is challenging as there is significant occlusion from the utensil due to which: (i) state-of-the-art monocular methods such as DECA [22] fail at mouth keypoint detection, and (ii) even if we obtain the mouth keypoints, projecting them on the depth image for pose estimation in the real-world fails as the depth image is noisy in the vicinity of the utensil. We propose a mouth perception pipeline that is robust to this occlusion challenge. Central to this pipeline is a novel method - Face Estimation from Multiple Views under Tool Occlusion (FEMTO). FEMTO uses inputs from two in-hand cameras mounted diametrically opposite to reconstruct a personalized head model for the care recipient, and estimates dense 2D facial keypoints. For the personalized head model, we use the parameterized head model FLAME [36] which is represented by a small number of shape, pose, and expression parameters.

**Model Architecture.** FEMTO’s architecture consists of two parallel encoders; one encoder for the top image and the other for the bottom. These encoders use frozen weights from DECA [22], a single-view model pre-trained on 2 million images from large datasets [12, 14, 47]. The outputs of these encoders are concatenated and further processed through two fully-connected layers. These layers integrate information from both cameras, and, similar to DECA, generate FLAME parameters and an orthographic camera pose for projecting the 3D FLAME mesh into image space.

**Self-Supervised Finetuning on Self-Curated Dataset.** To make FEMTO robust to occlusions, we finetune it on a self-curated visual dataset. We collect about 5000 unoccluded multi-view images with different head poses and facial expressions from 10 participants. Data from 8 participants are used for training FEMTO, and data from the remaining 2 participants are used for comparison against baselines. The latter is detailed in Section 4.1. These images are captured without the utensil in place, and thus DECA can typically generate accurate FLAME parameters and camera pose for these

images. We then synthetically occlude these images by in-painting the utensil (Figure 2). We use various utensils holding distinct food items at this step to ensure that FEMTO can generalize to new utensils and food items. We provide these occluded images to FEMTO, and use annotations generated by DECA for the corresponding un-occluded image as ground truth for supervision. We evaluate all the generated ground truth parameters and manually update them wherever necessary. FLAME’s modular parameterization allows for focusing training on parameters where DECA especially struggles under utensil occlusion, notably jaw pose prediction which is vital for detecting whether the mouth is open or closed.

**Robust Mouth Perception Pipeline.** FEMTO processes the top and bottom images to generate FLAME parameters and the orthographic camera pose of the top RGB image relative to the head model (Figure 2). These FLAME parameters decode into a custom 3D head model, and the generated camera pose can be used to project this model to the top RGB image and obtain dense 2D facial keypoints. We combine these 2D keypoints with the aligned depth image from the top camera to get dense 3D world keypoints. Finally, we use a correspondence-based registration robust to outliers to pose the generated head model to the 3D world keypoints. Despite the depth data around the mouth region being noisy due to occlusion from the utensil, the customized head model enables us to use depth data from other regions of the head for accurate mouth pose estimation. This proposed pipeline runs real-time (5-10Hz) on a system with RTX 3090.

Details on data collection, model training, and correspondence-based registration method used are in the Appendix [6].

### 3.2 Physical Interaction-Aware Control

Our physical interaction-aware control method first classifies the subtle physical interactions during bite transfer and switches appropriate compliant controllers accordingly. We take a data-driven approach to identify the nature of physical-interaction, and collect visual and haptic data for four types of physical interactions:

- (1) Incidental interactions (collisions) between the utensil and user’s mouth that occur outside the mouth due to sensing/control errors as the robot attempts to move inside.
- (2) In-mouth manipulation interactions initiated by the user using their tongue to guide the utensil to a desired transfer location inside their mouth.
- (3) Impulsive interactions due to involuntary spasms occurring while the utensil is inside the mouth.
- (4) Bite interactions that occur when the user takes a bite.

**Data Collection.** We collect a total of 3072 physical interactions (512 interactions X 6 participants without mobility limitations), varying both the utensil and the robot controller. For each physical interaction, we ask participants to perform different variations of the interaction. For example, when taking bites, participants are instructed to use only their teeth around the food item, only their lips around the utensil, or both.

**Multimodal Classification.** For each physical interaction, we analyze visual and haptic data from a 100ms window starting at contact initiation. Visual features include 2D and 3D face keypoints generated from our mouth perception pipeline (Section 3.1). In line with related haptic classification work [23, 25], we augment raw haptic data with time-series features computed over the window,

such as mean, range, kurtosis, Hjorth complexity, and frequency domain. We evaluate four models: Support Vector Machines (SVM) [16], Multi-Layer Perceptron (MLP) [30], Temporal Convolutional Networks (TCN) [34], and Time-Series Transformers (TST) [48].

**Controller Switching.** Once the robot recognizes the nature of physical interaction, it switches to an appropriate controller using an event-driven control method. We use two compliant controllers during inside-mouth bite transfer: a goal-tracking controller and a force-tracking controller. The goal-tracking controller tracks a given goal pose and is ideal for moving the food item to a desired position inside the user’s mouth and moving the utensil outside their mouth as soon as a bite is detected. The force-tracking controller minimizes contact force on the F/T sensor at the end-effector and thus is very compliant and reactive. This is ideal for physical interactions such as impulsive and in-mouth manipulation, where the robot must quickly respond to forces applied to the utensil.

Details on data collection, model training, and controller implementation are provided in the Appendix [6].

## 4 Evaluation of System Components

### 4.1 Mouth Perception Evaluation

We compare the performance of proposed mouth perception method against baselines for accurate 3D mouth keypoint generation, and perform an ablation study to evaluate the necessity of a key feature our mouth perception pipeline enables: continuous, real-time mouth perception for inside-mouth bite transfer.

**Eval. 1: Comparison with Baselines** - We compare the root mean square error (RMSE) of 3D mouth keypoints generated by our novel mouth perception method (FEMTO) against a pre-trained DECA using our proposed pipeline, and contemporary mouth perception methods in robot-assisted feeding [24, 43] (Table 1). For this evaluation, we use data from the 2 test participants in Section 3.1, introducing synthetic utensil occlusion to color and depth images. Ground truth keypoints are generated by applying our mouth perception pipeline with DECA to the unoccluded originals.

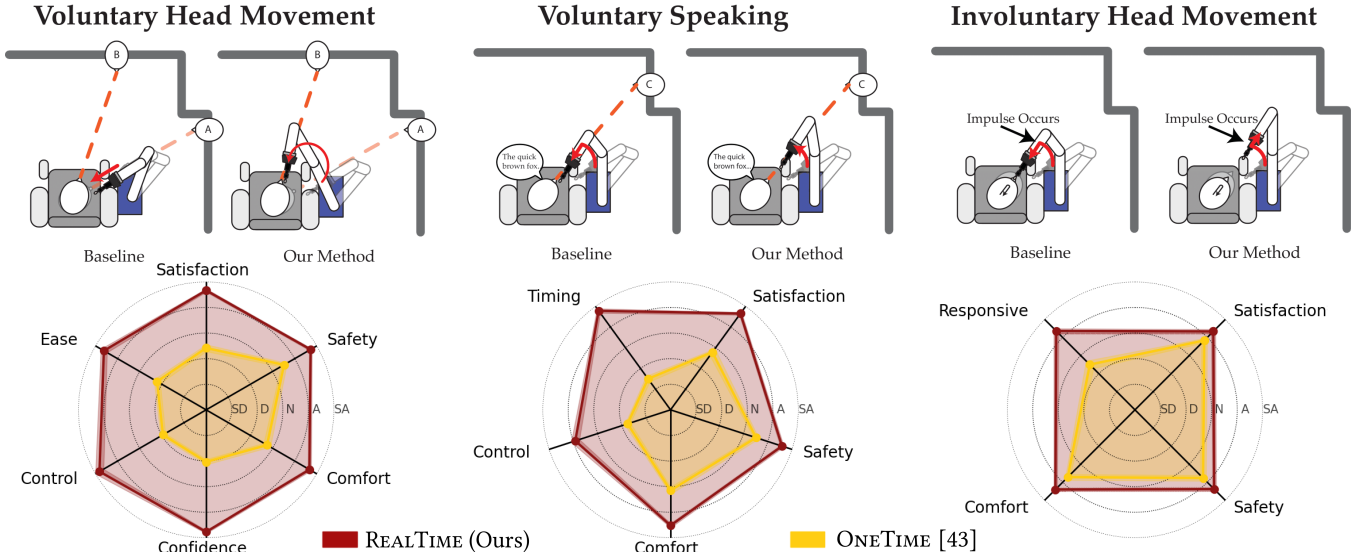
**Table 1: FEMTO outperforms baselines adapted from robot-assisted feeding [24, 43] and head perception [22] literature.**

Method	RMSE (in <i>mm</i> )
Gallenberger et al. 2019 [24]	110.454
Shaikewitz et al. 2023 [43]	82.946
DECA [22]	4.72
<b>FEMTO (Ours)</b>	<b>3.69</b>

**Eval. 2: Necessity of Real-Time Mouth Perception** - We perform an ablation study to evaluate the necessity of continuous, real-time mouth perception for the whole duration of inside-mouth bite transfer. This study examines three scenarios, selected based on feedback from individuals with mobility limitations, that impact their feeding process: (1) voluntary head movement, (2) voluntary speaking, and (3) involuntary head movement (spasm).

**Methods.** In each scenario, we compare two methods: our approach (REALTIME), which employs continuous, real-time mouth perception, and the baseline (ONETIME) [43] which estimates the user’s mouth pose only once when they first open their mouth.

**Procedure.** We conducted the study with 15 participants (11 male, 4 female; ages 20-65) who had no mobility limitations. After participants provide their demographics and complete pre-study



**Figure 3: REALTIME (Ours) enables the robot to track the mouth’s pose and state, which enhances satisfaction, safety, comfort, confidence, control, ease of use, bite timing, and responsiveness across various scenarios in comparison to ONETIME [43]. Response options used across studies: SD (Strongly Disagree), D (Disagree), N (Neutral), A (Agree), SA (Strongly Agree).**

questionnaires, they are seated and strapped into a wheelchair, and asked to simulate three scenarios (Figure 3):

S1. Voluntary Head Movement: Participants initially face person A, and open their mouth to initiate a bite. While the robot is moving towards them to feed, upon an audio cue, they turn to person B as if following a conversation.

S2. Voluntary Speaking: Participants initially face person C, and open their mouth to initiate a bite. While the robot is moving towards them to feed, upon an audio cue, they start reciting a preset script, as if conversing with person C.

S3. Involuntary Head Movement (Spasm): For this scenario, participants are first trained to mimic involuntary impulses, following guidance from an occupational therapist. During the trial, participants initially face person C, and open their mouth to initiate a bite. While the robot is moving towards them to feed, upon an audio cue, participants simulate a sudden forward impulse.

Post-feeding, participants fill out a Likert survey on perceived satisfaction, safety, and comfort [all scenarios], control and ease [S1], control and bite timing [S2], and robot responsiveness [S3]. They are fed a total of 24 times (3 scenarios x 2 methods X 2 utensils x 2 trials), with utensil and method orderings counterbalanced.

**Study Results.** REALTIME consistently outperforms ONETIME in all scenarios (Figure 3). In S1, ONETIME often misses the target by aiming at the initial mouth position, while REALTIME adapts to mouth movements for accurate feeding. In S2, ONETIME disrupts speech by continuing movement, whereas REALTIME is able to perceive the participant’s mouth closing and pauses while they are speaking. In S3, ONETIME maintains its movement towards the participant, even if they spasm. In contrast, REALTIME immediately retracts upon detecting impulse, resuming feeding only when safe.

#### 4.2 Physical Interaction-Aware Control Evaluation

We compare the performance of various time-series classification models for physical interaction classification, and perform an ablation study to evaluate the necessity of physical interaction-aware control for inside-mouth bite transfer.

**Eval. 1: Comparison between Models -** We use data collected with the 6 participants in Section 3.2, and consider two conditions:

- A - Aggregated training and testing datasets from all 6 participants, maintaining an 80:20 split of each participant’s data.
- B - Train on 5 participants, and test on the 6th, novel participant in leave-one-participant-out cross validation fashion.

SVM exhibits superior performance across both conditions (Table 2), which is likely attributable to the low data regime. For all models, combining haptic and visual modalities leads to superior performance compared to using either modality in isolation. This finding strongly advocates for the use of multimodal sensing for classifying physical interactions. The reduction in performance from condition A to B could stem from the unique traits of individual participants affecting physical interaction, aligning with findings from previous studies [15, 38].

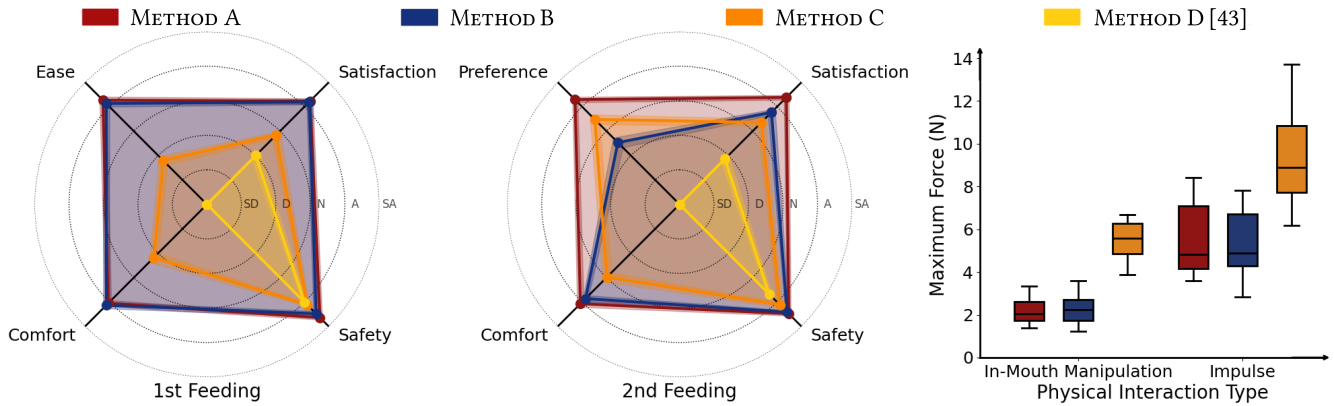
**Performance Improvement with Finetuning.** We investigate finetuning as a solution to the challenge of unique participant traits impacting zero-shot performance of models on novel participants.

**Table 2: Performance (F1 score) of physical interaction classification methods for inside-mouth bite transfer.**

Method	All Participants Aggregated			Novel Participant		
	All	Haptic	Visual	All	Haptic	Visual
SVM [16]	<b>0.903</b>	<b>0.857</b>	<b>0.860</b>	<b>0.872</b>	0.834	0.649
MLP [30]	0.892	0.845	0.842	0.871	<b>0.842</b>	<b>0.712</b>
TCN [34]	0.887	0.787	0.634	0.862	0.784	0.545
TST [48]	0.902	0.831	0.822	0.856	0.819	0.689

True label	Predicted label			
	Incidental	Impulsive	In-Mouth	Bite
Incidental	<b>0.98 ± 0.01</b>	0.01 ± 0.00	0.01 ± 0.00	0.01 ± 0.00
Impulsive	0.08 ± 0.01	<b>0.84 ± 0.01</b>	0.02 ± 0.00	0.06 ± 0.01
In-Mouth	0.01 ± 0.01	0.02 ± 0.01	<b>0.91 ± 0.02</b>	0.05 ± 0.01
Bite	0.04 ± 0.01	0.09 ± 0.02	0.11 ± 0.02	<b>0.76 ± 0.03</b>

**Figure 4: Confusion matrix showing multimodal SVM’s performance for Condition B (Novel Participant).**



**Figure 5: More physical-interaction aware control enhances perceived safety, satisfaction and comfort. This is supported by quantitative data recorded by the force sensor.**

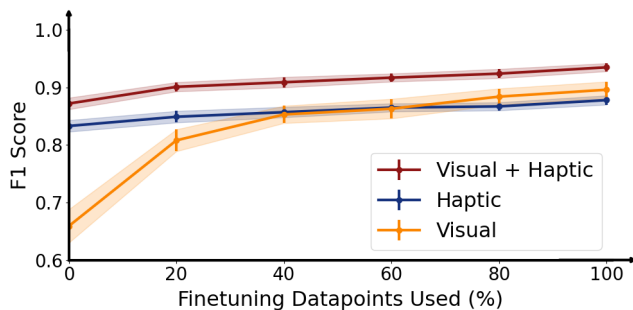
We split the novel participant’s data into an 80:20 ratio, allocating 80% for finetuning and 20% for testing. We evaluate the increase in performance of our best performing model, SVM, with incremental amounts of finetuning datapoints. Results (Figure 6) show steady performance improvement across modalities, with a significant performance improvement even with a small number of data points from the novel participant.

**Eval. 2: Necessity of Physical Interaction-Aware Control -**

We perform an ablation study to assess the importance of physical interaction-aware control for inside-mouth bite transfers, especially for individuals with severe mobility limitations. We consider a scenario where an individual must have food placed on the right side of their mouth due to an inability to bite in the center. When caregivers place food in the center, the person guides them to the correct position using their tongue. Once shown, they expect the caregiver to remember this preference. This person also experiences involuntary head movements, which can occur even while they eat.

**Methods.** We manipulate one within-subject variable: the robot’s level of physical-interaction awareness during feeding. We use four methods that vary in their ability to distinguish between different contact types, prompting the robot to switch to the appropriate controller accordingly. In order of decreasing awareness:

- METHOD A: Incidental vs. Impulsive vs. In-mouth Manipulation vs. Bite classification.
- METHOD B: Incidental vs. Inside Mouth Non-bite (combining impulsive and in-mouth manipulation) vs. Bite classification.
- METHOD C: Non-bite vs. Bite classification.
- METHOD D: Any Contact is Bite [43].



**Figure 6: SVM’s classification performance on novel participants improves steadily with finetuning using their data.**

Methods have access to the ground truth contact type for this study; we also report offline classification accuracy.

**Procedure.** We conducted the study with 14 participants (8 males, 6 females; ages 21-27) who had no mobility limitations. After providing demographics and completing pre-study questionnaires, participants are seated and strapped into a wheelchair. They then undergo training, under an occupational therapist’s guidance, to realistically simulate the actions of the illustrated care recipient. A trial in this study consists of two successive feedings using the same method, in the following sequence:

**1st Feeding:** This feeding focuses on in-mouth manipulation. The robot, initially stationary, begins moving when participants open their mouth. Once the robot moves to the center of their mouth, it says "push," prompting participants to move the food to the right side with their tongue while keeping their head still. When the food is correctly positioned, the robot signals "bite," and participants bite down and wait for the robot to withdraw.

**2nd Feeding:** This feeding focuses on reaction to impulsive physical interactions and remembering preferred bite transfer location. The robot, initially stationary, begins moving when participants open their mouth. Once inside the mouth, the robot says "impulse," prompting participants to simulate an impulse. After the movement, they keep their mouth open for five seconds. If the robot doesn’t place the food on the right side of their mouth within that time, participants have to adjust it with their tongue. When the food is correctly positioned, the robot signals "bite," and participants bite down and remain still until the robot withdraws.

Post-feeding, participants complete a Likert survey assessing satisfaction, safety, and comfort for both feedings. Additionally, the survey evaluates ease of moving the utensil for 1st Feeding and whether the robot automatically moved food to the preferred location for 2nd Feeding. Participants are fed a total of 32 times (2 feedings x 4 methods x 2 utensils x 2 trials), with utensil and method orderings counterbalanced.

**Study Results.** Figure 5 presents the study results. In the 1st Feeding, which focuses on in-mouth manipulation, participants predominantly prefer METHOD A and METHOD B. These methods switch to a force-tracking controller upon detecting in-mouth manipulation, making them easier for participants to manipulate with their tongue compared to METHOD C, which constantly maintains a goal-tracking controller. This preference is quantitatively supported

by the recorded force-torque data. In the 2nd Feeding, centered on reacting to impulsive physical interactions and remembering user preferences, participants rate METHOD A the highest. Both METHOD A and METHOD B transition to a force-tracking controller when impulsive contact is detected, and are perceived as comfortable by participants. However, participants note that METHOD A remembers their preferred transfer location, unlike METHOD B, which cannot distinguish between in-mouth manipulation and impulsive contacts. METHOD C, which continues with the goal-tracking controller during the impulsive motion, is rated as uncomfortable. The quantitative force-torque data show that METHOD C exerts significantly higher maximum force during an impulse as compared to METHOD A and METHOD B. As METHOD C does not respond to the impulsive contacts, it continues to move towards the preferred position even after the impulse. METHOD D [43] retracts from the mouth as soon as any contact is initiated during both 1st Feeding and 2nd Feeding. Consequently, it is ranked as the least satisfying method by the participants.

**Classification Accuracy on Eval. 2 Study Data.** For this analysis, we use our best-performing model, SVM, trained on 3072 physical interaction data points collected in Section 3.2, and test on interaction data from this Eval. 2 study. This evaluation tests the model’s generalization capabilities to more unstructured study settings, compared to the controlled conditions of the prior data collection. SVM achieves an F1 score of 0.719 zero-shot (Table 4). In addition to distribution shifts caused by the unstructured study setting and novel participants, this performance drop could also be due to changes made to the robot system design between the initial data collection and this ablation study, based on end-user feedback. Inline with improvements seen with finetuning in Eval. 1, the model’s performance after further finetuning on 80% of each participant’s data (approximately 10 data points per interaction type) and testing on the remaining 20% of their data, significantly improves to F1=0.772. These results underscore the importance of evaluating finetuning in physical interaction classification models as a means to address various sources of distribution shifts in real-world settings.

All study results are statistically significant (Wilcoxon Signed Rank Test,  $p < 0.05$ ). Details on pilot studies, questions for each measure, and classification results are in the Appendix [6].

**Table 4: SVM’s performance (F1 Score) on Eval. 2 study data.**

	All	Haptic	Visual
Zero-Shot	0.719	0.675	0.326
With Finetuning	0.772	0.697	0.605

## 5 Full System Evaluation with Care Recipients

We evaluate our full system through a user study with 13 individuals with severe mobility limitations, all of whom require assistance with feeding. The objectives of this study were:

- O1. To evaluate the effectiveness and acceptance of our inside-mouth bite transfer system among its intended end users.
- O2. Compare participant preferences between inside-mouth and outside-mouth bite transfer [8] systems, when both options are functionally possible for an individual.

Due to challenges in recruiting many individuals with mobility limitations at any one place, this study took place at three sites: EmPRISE Lab in Ithaca, NY, Columbia University Medical Center in NYC, NY, and a participant’s home (Figure 7) in Taftville, CT.

### 5.1 Methods and Procedures

We considered two systems in this user study: (1) our proposed inside-mouth bite transfer system, and (2) an outside-mouth bite transfer system adapted from Bhattacharjee et al. 2020 [8] that stops at a fixed distance (5 cms) from the user’s mouth.

We informed participants to focus on the interaction from when the robot picks up a food item until they take a bite. Initially, we guided them through practice runs for each method. We used two utensils, a fork for cantaloupe and a spoon for yogurt. A trial involved three successive feedings using the same utensil and bite transfer system, after which we asked questions about their perceived safety and comfort. We evaluated the inside-mouth bite transfer system (O1) with all participants. We compared with the outside-mouth bite transfer system (O2) only with participants who demonstrated the required functional capability during practice runs. Order of the two systems and utensils were counterbalanced. After all trials, a final evaluation phase occurred where we fed participants using inside-mouth bite transfer and asked them to complete a TAM survey [18]. This was followed by a feedback interview to perceive the significance of the key components of our system and identify areas for improvement.

**Table 3: Demographics of user study participants. All participants require assistance with feeding. BiR - Bite Reflex, GaR - Gag Reflex, LR - Limited head/neck ROM, MaS - Masseter Spasticity, TS - Tongue Spasticity**

ID	Age	Gender	Race	Self-described impairment	Impairment time	Challenges with Feeding
P1	44	Female	Caucasian/White	Multiple Sclerosis	25 years	LR, Spasms
P2	33	Female	African American	C3-C4 Spinal Cord Injury	13 years	LR, Spasms
P3	30	Male	Caucasian/White	Arthrogryposis	Since Birth	GaR, LR, MaS, TS
P4	45	Female	Caucasian/White	Multiple Sclerosis	15 years	LR
P5	26	Female	Hispanic	Schizencephaly Quadriplegia	Since Birth	BiR, GaR, LR, MaS, Overbite, TS
P6	27	Male	Caucasian/White	Spinal Muscular Atrophy	Since birth	LR, MaS
P7	49	Female	Hispanic	C4-C5 Spinal Cord Injury	28 years	LR, Spasms
P8	44	Male	Caucasian/White	C4-C5 Spinal Cord Injury	24 years	LR
P9	39	Female	Asian/Pacific Islander	Spinal Muscular Atrophy	Since Birth	LR, MaS
P10	48	Male	African American	C5-C6 Spinal Cord Injury	2.5 years	LR, TS
P11	31	Female	African American	Cerebral Palsy Quadriplegia	Since Birth	LR, TS
P12	25	Male	Caucasian/White	C4-C5 Spinal Cord Injury	6 years	LR
P13	31	Female	African American	Arthrogryposis	Since Birth	LR, TS



**Figure 7: Left: Our inside-mouth bite transfer system feeding a care recipient in their home. Center: Inside-mouth bite transfer challenges: P6 has limited mouth opening, P3 can bite only at left molars, P5's weak jaw control causes incidental contacts, and P1's muscle spasms move her towards the utensil. Right: Our system received high ratings from care recipients for its safety and comfort, and was favorably evaluated in terms of technology acceptance, as indicated by the results of the TAM survey.**

## 5.2 Participants

The study involved 13 participants with diverse ages, genders, self-described impairments, duration of impairments, and care providers (Table 3). Some participants had complex medical conditions, making feeding particularly challenging (Figure 7). For instance, P6 and P9 have Spinal Muscular Atrophy, a condition known to reduce maximum mouth opening and decrease bite force by almost 50% [29]. P3 has an overbite and a sensitive gag reflex, which can be triggered if a utensil presses on their tongue. For solid foods, they prefer the utensil to approach from the left side of their mouth, depositing the food on their left molars. P1, P2 and P7 experience spasms that can happen unexpectedly. P5 is diagnosed with Schizencephaly, which manifests with severe motor limitations (quadriplegia) and poor jaw and tongue control for eating. They have severe spasticity in their tongue, gag and bite reflexes, and an overbite.

All data collection and user studies in this paper were approved by the Cornell University Institutional Review Board.

## 5.3 Findings

**Our inside-mouth bite transfer system is perceived as safe and comfortable.** Figure 7 displays the results. P7 mentioned, "... going in and out, I felt very comfortable, I felt safe, I enjoyed it..." P9 echoed, "... I think I was actually surprised... I liked it, the entrance in and out..." P5's parents were pleasantly surprised, "... I think it worked really much better than I expected. When she moved her mouth, it moved with her... She had fun with it!" P6 emphasized the importance of placement, "The in-the-mouth method worked very well. I was worried before doing it that it would go too deep, and it didn't. I think that making sure it goes in the right place is important, and I think you did very well."

**Users view the technology favorably.** As per Figure 7, users found the inside-mouth bite transfer system to be useful, easy to use, have a positive attitude and intend to use it, and enjoy the process. P2's caregiver remarked, "... it will help a lot of people, especially with her condition... it hurts me sometimes when she can't do certain things, so I'm glad." P6's parents highlighted its potential for social engagement, "[in a social setting,] the caregivers all interact, and the kids are just passive recipients of food. So I think this will involve them more in the eating."

**Even among the few users able to perform outside-mouth bite transfer, some preferred inside-mouth bite transfer.** Out of the 13 users, 9 were unable to lean forward for outside-mouth

bite transfer. Among the 4 users who evaluated both systems, some preferred inside-mouth as they felt that outside-mouth bite transfer requires more physical effort. P10 favored inside-mouth, saying it "requires no movement of the neck." P12 and P4's preferences varied based on their physical state. P12 stated, "If I was having neck pain, maybe I would prefer inside mouth so that I am not moving as much." P11 consistently preferred outside-mouth, noting, "I have the ability to position the food in my mouth myself, with my own comfort."

**Users highlight that continuous, real-time mouth perception and adaptive compliance are essential.** P12 noted, "I could be sitting in front of the TV watching something... I'm ready to take a bite and then it catches my attention, or yeah, talking to someone." P9 added another context: "When a phone rings, and I turn around..." P7 emphasized the importance of device adaptability during spontaneous movements, "On spasms. I mean, if it wiggles with you, that would be perfect. The body can move when there is a spasm, so it could go off a little and if it follows me, that's great."

## 6 DISCUSSION

Our study demonstrates promising results, underlining the potential feeding systems have to help improve the quality of life of care recipients. However, the wide range of disabilities and individual user needs makes crafting universal solutions complex. Participants with unique challenges, like limited mouth opening, stressed the importance of personalized adaptations. For instance, P6 and P9, who have smaller mouth openings, suggested specifically designed utensils, with P9 pointing out, "The spoon was too big." Similarly, opinions on transfer speed varied: P11 described the robot as "a little bit too fast," while P12 found it "too slow." These insights highlight the need for adaptive algorithms and interfaces that are tailored to individual needs, moving away from an one-size-fits-all model. Additionally, while our study focuses on short-term interactions, future work needs to explore long-term usability. Longitudinal studies will offer a deeper insight into the prolonged impact of this technology on users.

## 7 ACKNOWLEDGEMENT

This work was partly funded by NSF IIS #2132846, CAREER #2238792, and DARPA under Contract HR001120C0107. We would like to acknowledge Taylor Sanchez, Megan Sofield, Shubhangi Sinha and Skyler Valdez for their help with the user studies.



## References

- [1] 2024. ATI F/T Sensor. <https://www.ati-ia.com/products/ft/sensors.aspx> (Accessed: 1st January, 2024).
- [2] 2024. Kinova Gen3 6DoF Robotic Arm. <https://www.kinovarobotics.com/product/gen3-robots> (Accessed: 1st January, 2024).
- [3] 2024. Neater Eater Robot. <https://www.neater.co.uk/neater-eater-robotic> (Accessed: 1st January, 2024).
- [4] 2024. Obi. <https://meetobi.com> (Accessed: 1st January, 2024).
- [5] 2024. Robotic 2F-85 Gripper. <https://robotiq.com/products/2f85-140-adaptive-robot-gripper> (Accessed: 1st January, 2024).
- [6] 2024. Supplementary Materials. <https://emprise.cs.cornell.edu/bitetransfer/> (Accessed: 1st January, 2024).
- [7] Suneel Belkhale, Ethan K Gordon, Yuxiao Chen, Siddhartha Srinivasa, Tapomayukh Bhattacharjee, and Dorsa Sadigh. 2021. Balancing Efficiency and Comfort in Robot-Assisted Bite Transfer. *arXiv preprint arXiv:2111.11401* (2021).
- [8] Tapomayukh Bhattacharjee, Ethan K Gordon, Rosario Scalise, Maria E Cabrera, Anat Caspi, Maya Cakmak, and Siddhartha S Srinivasa. 2020. Is more autonomy always better? exploring preferences of users with mobility impairments in robot-assisted feeding. In *2020 15th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 181–190.
- [9] James Booth, Anastasios Roussos, Stefanos Zafeiriou, Allan Ponniah, and David Dunaway. 2016. A 3d morphable model learnt from 10,000 faces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5543–5552.
- [10] Steven W Brose, Douglas J Weber, Ben A Salatin, Garret G Grindle, Hongwu Wang, Juan J Vazquez, and Rory A Cooper. 2010. The role of assistive robotics in the lives of persons with disability. *AJPM&R* (2010).
- [11] Alexandre Candéias, Travers Rhodes, Manuel Marques, Manuela Veloso, et al. 2018. Vision augmented robot feeding. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*. 0–0.
- [12] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. 2018. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*. IEEE, 67–74.
- [13] Adriano Chiò, A Gauthier, A Vignola, Andrea Calvo, Paolo Ghiglione, Enrico Cavallo, AA Terreni, and Roberto Mutani. 2006. Caregiver time use in ALS. *Neurology* 67, 5 (2006), 902–904.
- [14] J Chung, A Nagrani, and A Zisserman. 2018. VoxCeleb2: Deep speaker recognition. *Interspeech 2018* (2018).
- [15] Giovanni Cioffi, Silke Klose, and Arne Wahrburg. 2020. Data-efficient online classification of human-robot contact situations. In *2020 European Control Conference (ECC)*. IEEE, 608–614.
- [16] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20, 3 (1995), 273–297.
- [17] Hang Dai, Nick Pears, William AP Smith, and Christian Duncan. 2017. A 3d morphable model of craniofacial shape and texture variation. In *Proceedings of the IEEE international conference on computer vision*. 3085–3093.
- [18] Fred D Davis, Richard P Bagozzi, and Paul R Warshaw. 1989. User acceptance of computer technology: A comparison of two theoretical models. *Management science* 35, 8 (1989), 982–1003.
- [19] Laura E Dreer, Timothy R Elliott, Richard Shewchuk, Jack W Berry, and Patricia Rivera. 2007. Family caregivers of persons with spinal cord injury: Predicting caregivers at risk for probable depression. *Rehabilitation Psychology* 52, 3 (2007), 351.
- [20] Nima Fazeli, Miquel Oller, Jiajun Wu, Zheng Wu, Joshua B Tenenbaum, and Alberto Rodriguez. 2019. See, feel, act: Hierarchical learning for complex manipulation skills with multisensory fusion. *Science Robotics* 4, 26 (2019), eaav3123.
- [21] Ryan Feng, Youngsun Kim, Gilwoo Lee, Ethan K Gordon, Matt Schmittle, Shivaum Kumar, Tapomayukh Bhattacharjee, and Siddhartha S Srinivasa. 2019. Robot-assisted feeding: Generalizing skewering strategies across food items on a realistic plate. *arXiv preprint arXiv:1906.02350* (2019).
- [22] Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. 2021. Learning an animatable detailed 3D face model from in-the-wild images. *ACM Transactions on Graphics (TOG)* (2021), 1–13.
- [23] Felix Franzel, Thomas Eiband, and Dongheui Lee. 2021. Detection of Collaboration and Collision Events during Contact Task Execution. In *2020 IEEE-RAS 20th International Conference on Humanoid Robots (Humanoids)*. IEEE, 376–383.
- [24] Daniel Gallenberger, Tapomayukh Bhattacharjee, Youngsun Kim, and Siddhartha S Srinivasa. 2019. Transfer depends on acquisition: Analyzing manipulation strategies for robotic feeding. In *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 267–276.
- [25] Saskia Golz, Christian Osendorfer, and Sami Haddadin. 2015. Using tactile sensation for learning contact knowledge: Discriminate collision from physical interaction. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 3788–3794.
- [26] Ethan K Gordon, Xiang Meng, Tapomayukh Bhattacharjee, Matt Barnes, and Siddhartha S Srinivasa. 2020. Adaptive robot-assisted feeding: An online learning framework for acquiring previously unseen food items. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 9659–9666.
- [27] Ethan Kroll Gordon, Amal Nanavati, Ramya Challa, Bernie Hao Zhu, Taylor Annette Kessler Faulkner, and Siddhartha Srinivasa. 2023. Towards General Single-Utensil Food Acquisition with Human-Informed Actions. In *Conference on Robot Learning*. PMLR, 2414–2428.
- [28] Ethan K Gordon, Sumegh Roychowdhury, Tapomayukh Bhattacharjee, Kevin Jamieson, and Siddhartha S Srinivasa. 2021. Leveraging Post Hoc Context for Faster Learning in Bandit Settings with Applications in Robot-Assisted Feeding. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 10528–10535.
- [29] M.W. Granger, P.H. Buschang, G.S. Throckmorton, and S.T. Iannaccone. 1999. Masticatory muscle function in patients with spinal muscular atrophy. *American Journal of Orthodontics and Dentofacial Orthopedics* 115, 6 (1999), 697–702. [https://doi.org/10.1016/S0889-5406\(99\)70296-9](https://doi.org/10.1016/S0889-5406(99)70296-9)
- [30] Simon Haykin. 1994. *Neural networks: a comprehensive foundation*. Prentice Hall PTR.
- [31] Catrine Jacobsson, Karin Axelsson, Per Olov Österlind, and Astrid Norberg. 2000. How people with stroke and healthy older people experience the eating process. *Journal of Clinical Nursing* 9, 2 (2000), 255–264. <https://doi.org/10.1046/j.1365-2702.2000.00355.x>
- [32] Sidney Katz, Amasa B Ford, Roland W Moskowitz, Beverly A Jackson, and Marjorie W Jaffe. 1963. Studies of illness in the aged: the index of ADL: a standardized measure of biological and psychosocial function. *jama* 185, 12 (1963), 914–919.
- [33] Alexandros Kouris, Fotios Dimeas, and Nikos Aspragathos. 2018. A frequency domain approach for contact type distinction in human-robot collaboration. *IEEE robotics and automation letters* 3, 2 (2018), 720–727.
- [34] Colin Lea, Michael D Flynn, Rene Vidal, Austin Reiter, and Gregory D Hager. 2017. Temporal convolutional networks for action segmentation and detection. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 156–165.
- [35] Michelle A Lee, Yuke Zhu, Krishnan Srinivasan, Parth Shah, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Jeannette Bohg. 2019. Making sense of vision and touch: Self-supervised learning of multimodal representations for contact-rich tasks. In *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 8943–8950.
- [36] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. 2017. Learning a model of facial shape and expression from 4D scans. *ACM Trans. Graph.* 36, 6 (2017), 194–1.
- [37] Yunzhu Li, Jun-Yan Zhu, Russ Tedrake, and Antonio Torralba. 2019. Connecting touch and vision via cross-modal prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10609–10618.
- [38] Martina Lippi, Giuseppe Gillini, Alessandro Marino, and Filippo Arrichiello. 2021. A Data-Driven Approach for Contact Detection, Classification and Reaction in Physical Human-Robot Collaboration. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 3597–3603.
- [39] J Lynch and R Cahalan. 2017. The impact of spinal cord injury on the quality of life of primary family caregivers: a literature review. *Spinal cord* 55, 11 (2017), 964–978.
- [40] Rishabh Madan, Rajat Kumar Jenamani, Vy Thuy Nguyen, Ahmed Moustafa, Xuefeng Hu, Katherine Dimitropoulou, and Tapomayukh Bhattacharjee. 2022. Spares: Structuring physically assistive robotics for caregiving with stakeholders-in-the-loop. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 641–648.
- [41] Amal Nanavati, Patricia Alves-Oliveira, Tyler Schrenk, Ethan K Gordon, Maya Cakmak, and Siddhartha S Srinivasa. 2023. Design principles for robot-assisted feeding in social contexts. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot Interaction*. 24–33.
- [42] Daehyung Park, Yuuna Hoshi, Harshal P. Mahajan, Ho Keun Kim, Zackory Erickson, Wendy A. Rogers, and Charles C. Kemp. 2020. Active robot-assisted feeding with a general-purpose mobile manipulator: Design, evaluation, and lessons learned. *Robotics and Autonomous Systems* 124 (2020), 103344. <https://doi.org/10.1016/j.robot.2019.103344>
- [43] Lorenzo Shaikewitz, Yilin Wu, Suneel Belkhale, Jennifer Grannen, Priya Sundaresan, and Dorsa Sadigh. 2023. In-Mouth Robotic Bite Transfer with Visual and Haptic Sensing. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 9885–9895.
- [44] Samantha E. Shune. 2020. An Altered Eating Experience: Attitudes Toward Feeding Assistance Among Younger and Older Adults. *Rehabilitation nursing : the official journal of the Association of Rehabilitation Nurses* (2020).
- [45] Priya Sundaresan, Suneel Belkhale, and Dorsa Sadigh. 2023. Learning visuo-haptic skewering strategies for robot-assisted feeding. In *Conference on Robot Learning*. PMLR, 332–341.
- [46] Danielle M Taylor. 2018. Americans with disabilities: 2014. *US Census Bureau* (2018), 1–32.
- [47] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang. 2019. Racial faces in the wild: Reducing racial bias by information maximization adaptation network. In *Proceedings of the IEEE/CVF international conference on computer*

- vision*. 692–702.
- [48] George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff. 2021. A Transformer-Based Framework for Multivariate Time Series Representation Learning. In *Proceedings of the 27th ACM SIGKDD*

*Conference on Knowledge Discovery & Data Mining (Virtual Event, Singapore) (KDD '21)*. Association for Computing Machinery, New York, NY, USA, 2114–2124. <https://doi.org/10.1145/3447548.3467401>