

🐙Octopi: Object Property Reasoning with Large Tactile-Language Models

Samson Yu[†], Kelvin Lin[†], Anxing Xiao[†], Jiafei Duan[§], and Harold Soh^{†‡}

[†]Dept. of Computer Science, National University of Singapore

[§]University of Washington, [‡]NUS Smart Systems Institute

Contact Authors: samson.yu@u.nus.edu, harold@comp.nus.edu.sg

Abstract—Physical reasoning is important for effective robot manipulation. Recent work has investigated both vision and language modalities for physical reasoning; vision can reveal information about objects in the environment and language serves as an abstraction and communication medium for additional context. Although these works have demonstrated success on a variety of physical reasoning tasks, they are limited to physical properties that can be inferred from visual or language inputs. In this work, we investigate combining tactile perception with language, which enables embodied systems to obtain physical properties through interaction and apply commonsense reasoning. We contribute a new dataset PHYSICLEAR, which comprises both physical/property reasoning tasks and annotated tactile videos obtained using a GelSight tactile sensor. We then introduce OCTOPI, a system that leverages both tactile representation learning and large vision-language models to predict and reason about tactile inputs with minimal language fine-tuning. Our evaluations on PHYSICLEAR show that OCTOPI is able to effectively use intermediate physical property predictions to improve its performance on various tactile-related tasks. PHYSICLEAR and OCTOPI are available at <https://github.com/clear-nus/octopi>.

I. INTRODUCTION

For humans, touch is a crucial sense that provides physical information beyond what vision can provide (e.g., material properties, texture information, temperature), especially during occlusion. This in turn improves our ability to perform physical reasoning [37, 4] and act in our world. Here, we are interested in enabling general purpose robots, specifically those empowered with large-language models (LLMs), to perform similar physical reasoning. While recent work has demonstrated that LLMs and large vision-language models (LVLMs) can provide an impressive level of commonsense and physical reasoning [46, 55, 3, 53], they are limited to either visual or text input modalities, and thus have limited performance in scenarios with visual ambiguity [17].

This paper extends LVLMs to have the sense of touch. We posit that incorporating a tactile modality into LVLMs will enable better physical reasoning in real-world environments. As an example, Fig. 1 illustrates how commonsense knowledge is applied together with tactile information to complete a novel physical task. Here, the robot leverages its tactile inputs together with the LLM’s commonsense knowledge (that ripe avocados are soft) to correctly select the ripe avocado. We use visual-tactile sensors, i.e., the GelSight [60], which provides image frames that reveal physical object properties such as texture and hardness [59]. However, there remains a significant



Fig. 1. Avocado ripeness selection by combining tactile information with commonsense knowledge. Using inputs from its tactile sensor, OCTOPI identifies the left avocado as softer. Using commonsense reasoning, OCTOPI infers that it is ripe and fulfills the user’s request.

domain gap between natural images that typical LVLMs are trained with and the tactile data.

To bridge this gap, we contribute the PHYSICLEAR dataset, which comprises GelSight images on a variety of real world objects, along with object labels and part annotations. PHYSICLEAR complements existing tactile datasets [59, 61, 18, 56, 19] as it provides three physical property annotations, specifically hardness, roughness, and bumpiness, that have been used in prior research [43, 20, 38, 10, 5, 26] and can be potentially inferred from the GelSight data. PHYSICLEAR also includes an training and evaluation suite comprising five reasoning tasks, which can serve as a benchmark for the research community.

Using PHYSICLEAR, we develop OCTOPI (Object Comprehension with Tactile Observations for Physical Intelligence). OCTOPI is a LLaMA-based [49, 50] LVLM (Vicuna [11]) equipped with a CLIP-based [39] tactile encoder, whose representations have been aligned via projection. In experiments, we show that OCTOPI is able to use its tactile modality to predict object properties and reason about scenarios including avocado ripeness.

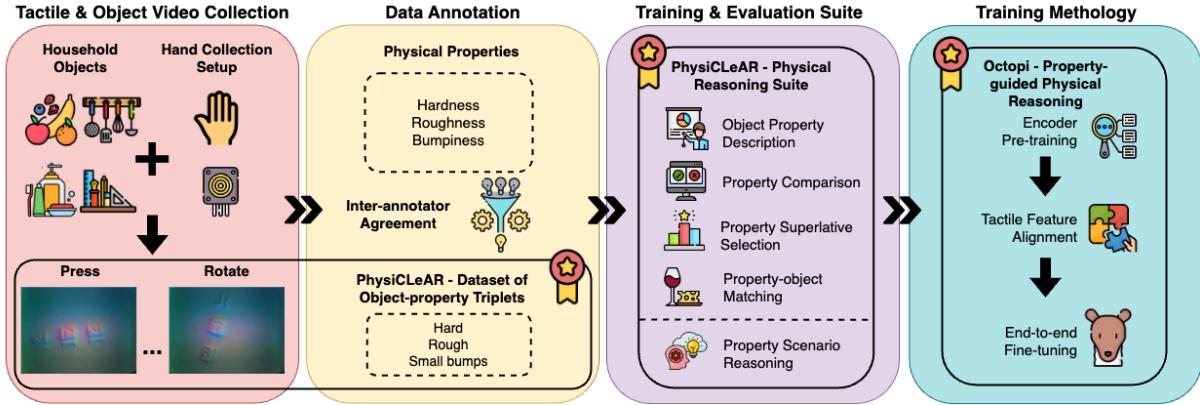


Fig. 2. PHYSICLEAR and OCTOPI (with key contributions starred). We collect tactile videos for everyday household objects by hand with two exploratory procedures: pressing and rotation. The videos are annotated by three annotators for three physical properties: hardness, roughness and bumpiness. PHYSICLEAR leverages the videos and annotations for five language-driven physical description and understanding tasks. OCTOPI is a LVLM fine-tuned on PHYSICLEAR for tactile-grounded physical understanding and reasoning.

TABLE I. Physical Property Details. The physical object properties examined, along with their descriptions and semantic categories.

Property	Description	Categories
Hardness	how easily an object's surface is deformed when pressed	soft, moderately hard, hard
Roughness	how rough an object's surface is	smooth, slightly rough, rough
Bumpiness	the size of bumps on an object's surface	no bumps, small bumps, big bumps

TABLE II. GelSight Dataset Comparisons. PHYSICLEAR provides physical property labels for tactile descriptions and physical reasoning across three physical properties. We further compare against existing datasets across three diversity measures. Property diversity refers to whether there are objects in the dataset that vary across the three properties we selected: hardness, roughness and bumpiness. Object diversity indicates whether there is more than one type of object in the dataset. Material diversity indicates the number of different materials in the dataset.

Dataset	Property Label Availability	Property Diversity	Object Diversity	Material Diversity
Hardness Dataset (2016) [59]	Yes (only hardness)	Yes	Yes	Medium
Clothing Dataset (2018) [61]	Yes	Yes	No (only clothing)	Low
ObjectFolder 2.0 (2022) [18]	No	No (only hard objects)	Yes	Medium
Touch and Go (2022) [56]	No	Yes	Yes	High
ObjectFolder-Real (2023) [19]	No	No (only hard objects)	Yes	Medium
PHYSICLEAR	Yes	Yes	Yes	Medium

Contributions. In summary, this paper makes the following key contributions:

- A new GelSight dataset, PHYSICLEAR, that exhibits property diversity, object diversity, and material diversity for selected physical properties.
- OCTOPI, a framework for physical reasoning that leverages vision-based tactile sensors and the commonsense reasoning capabilities of LLMs.
- An accompanying training and evaluation suite spanning five tasks and baseline results using OCTOPI.

We hope that PHYSICLEAR and OCTOPI will spur research in tactile-enabled physical reasoning for embodied AI systems [14].

II. RELATED WORK

In this section, we briefly review prior work on tactile representation learning with the GelSight sensor, large vision-language models (LVLMs) and language/vision-guided physical reasoning. There has been significant work in tactile-based manipulation and physical reasoning, and we refer

readers desiring information on these topics to relevant survey papers [37, 14, 66, 47, 28].

Tactile Representation Learning with GelSight. Tactile representation learning has advanced significantly in recent years as robotic manipulation often requires more precision beyond what can be provided by vision alone [44]. Among the available tactile sensors, vision-based sensors have gained popularity due to their high-resolution image outputs and versatility. In particular, the GelSight sensor has been used in recent work [32, 59, 60, 24, 61] for inferring physical properties (e.g. hardness, texture and liquid volume) and to manipulate objects [48]. A key benefit of the GelSight is that its image outputs can be easily processed by modern deep learning methods [24]. As a result, popular vision algorithms have been used for tactile representation learning with GelSight [62, 8]. In our work, we exploit recent advances in tactile representation learning to extend the capabilities of LVLMs to reason about vision-based tactile input.

Large Vision-Language Models. Recent advancements in LLMs have spurred a significant increase in efforts to integrate

vision models with LLMs, exemplified by Flamingo [1], BLIP-2 [29], and MiniGPT-v2 [9]. These Large Vision-Language Models (LVLMs) have shown remarkable effectiveness in utilizing web-scale image-text data for image-based reasoning, benefiting a range of applications from robotics [7, 15] to medical imaging [42]. Very recent work involves developing LVLMs that can process video content [30, 36], enabling reasoning over dynamic visual information, or integrate multimodal sensory data [64].

Physical Reasoning with Language and Vision as Context. The exploration of physical reasoning in conjunction with language predates the emergence of LLMs. Early studies focused on assessing model proficiency in physical reasoning. For example, the PIQA [6] benchmark evaluates models on physical common sense, whereas PROST[2] examines their understanding of physical reasoning concepts. Subsequent advancements in language grounding have led to works such CLEVRER [58], PIP [13], SPACE [12] and Phys101 [54], which investigate the acquisition of physical reasoning skills from visual inputs.

In the emerging LLM era, research has focused on object-centric physical reasoning in LLMs. This involves evaluating various LLMs for their physical reasoning capabilities, e.g., NEWTON [52], and employing Vision-Language Models (VLMs) to predict physical properties that are then used to facilitate reasoning, as demonstrated in physically-grounded VLMs [17]. Unlike previous studies that primarily address physical reasoning through the integration of vision and language, OCTOPI stands out as the one of the first models capable of processing tactile images alongside language instructions to enable physical reasoning. There has been very recent work [22] that uses simulated tactile inputs with LLMs, but we focus on real tactile data. Concurrent work [16, 57] also explores real-world tactile data but our work features physical property annotations and a test suite comprising scenario reasoning tasks, and experiments using OCTOPI to evaluate the utility of physical property inference.

III. PHYSICLEAR - TACTILE AND PHYSICAL UNDERSTANDING TRAINING & EVALUATION SUITE

This section describes PHYSICLEAR, which comprises a tactile dataset with physical property and object-part annotations, along with a training and evaluation suite.

A. Physical Property Selection

In this work, we focus on three object properties: *hardness*, *roughness* and *bumpiness*. We list each property's description and categories in Table I. Briefly, *hardness* is characterized by the extent of surface deformation when subjected to pressure; *roughness* pertains to the texture of the surface; and *bumpiness* describes the prominence of surface protrusions. The *hardness* of an object correlates with its compliance and thermal characteristics. In contrast, *roughness* and *bumpiness* are attributes influenced by the surface's friction coefficient [10].

The selection of *hardness*, *roughness*, and *bumpiness* as physical attributes in our research is grounded in their relevance for physical reasoning [43, 20, 38, 10, 5, 26]. Generally, static physical properties of objects are categorized into geometric (e.g., size), material (e.g., hardness), and affective (e.g., comfort) [41]. Our study predominantly addresses material properties, as we deemed geometric and affective properties too challenging to ascertain using the GelSight. The choice of these specific properties was also informed by the data collection methodology [27], tailored to the limitations and strengths of the GelSight sensor, including considerations for its sensitivity and durability.

B. Dataset Collection & Annotation

To facilitate the grounding of our physical reasoning on tactile inputs, we collected a dataset of 74 everyday objects, totalling 408 tactile videos and corresponding videos showing the object as the data was collected. These objects were selected to span across our three selected properties, with variations across object types and materials. Detailed comparisons between PHYSICLEAR and existing GelSight datasets can be found in Table II.

The GelSight data was collected by-hand to mitigate risk of damaging the sensors and due to the challenge of securing different parts of irregularly-shaped objects while performing the required sampling motions. For each selected object, we captured up to seven tactile videos for each distinct region identified by a human evaluator. This process involved a two-step procedure: initially pressing the GelSight sensor against the object to capture pressure readings, followed by rotating the sensor to acquire shear readings. Each video generated from a single GelSight sensor reading constitutes an individual data point within our dataset.

Annotations of the physical properties were carried out by three independent annotators, with the average score used as the final annotation for each data point. Annotators were provided with both the tactile videos and the objects. Each property has three categories, and annotators were given the following guidelines for labeling each property:

- **Hardness:** The label *soft* is for objects that are compressible with little force, *moderately hard* for objects that are compressible with moderate force, and *hard* for objects that are incompressible even with a large pressing force.
- **Roughness:** *smooth* is for objects that present very minimal or no resistance when we slide our finger across its surface, *slightly rough* for objects with slight resistance, and *rough* for objects with significant resistance.
- **Bumpiness:** *no bumps* is for objects with no visible protrusions on its surface, *small bumps* for objects with protrusions smaller than $\approx 1/4$ of the tactile image upon contact, and *big bumps* for objects with protrusions larger than $1/4$.

This process yielded over 1,200 annotations and we observed high inter-annotator agreement scores (ICC3k of 0.894 (hardness), 0.979 (roughness), and 0.792 (bumpiness)). For reference, a score above 0.75 is considered good or excellent

reliability. The dataset was subsequently divided into three distinct subsets (training, validation, and testing) following an 80-10-10 split. This division resulted in 60 objects for training and 7 objects each for validation and testing.

C. Training & Evaluation Suite

PHYSICLEAR’s training and evaluation suite comprises five physical reasoning tasks (Table III). All five tasks use tactile data and natural language instructions as inputs (Table IV). Since the tactile data is in video form, we follow prior LVLM work and represent it as a sequence of frames: X_1, \dots, X_N . We further detail each task’s motivation, setup, evaluation details and whether they are used for training [T] and/or evaluation [E] below:

Object Property Description (OPD) [T, E]. This task addresses property-based description: generating both unstructured and structured descriptions of an object’s hardness, roughness, and bumpiness from tactile videos. It parallels existing image or video captioning tasks, which use natural visuals, and aids the LLM in interpreting tactile signals. The unstructured description provides more complex descriptions in cases where they are visible in the tactile video (e.g. “fibrous structure” for toilet paper or “grains” for a scoop of rice). The language instructions are variants of “Describe the physical properties of $\langle \text{tact_start} \rangle T_1, \dots, T_N \langle \text{tact_end} \rangle$.” The unstructured description is generated using ChatGPT 3.5 and manually cleaned to produce a diverse description of our objects based on our three physical properties. The structured description of an object’s physical properties using our annotations is formatted as: “Overall, it presents a *hardness_label* and *roughness_label* surface with *bumpiness_label*.”

Property Comparison (PC) [T, E]. Given two tactile videos, each of a different object, a specified physical property, and its comparative adjective, determine whether the comparative adjective accurately describes the two videos. From a training perspective, this task helps a model distinguish between the various descriptions of physical properties, thereby aligning its comprehension of physical characteristics with our defined categories of *hardness*, *roughness*, and *bumpiness*. This alignment ability may improve a model’s ability to interpret and reason about the physical world in a manner consistent with human understanding.

Property Superlative Selection (PSS) [T, E]. For three tactile videos, each of a different object, and a specified physical property and its superlative adjective (e.g. hardest for the hardness property), choose the video that the superlative adjective best describes. This task is similar to the PC task and helps the LLM align its physical understanding with that of our physical property descriptions. Furthermore, since prior work has shown that LLMs might perform differently when the polarity of the comparative adjective changes [52], this task seeks to enhance the LLM’s resilience to various comparative descriptions of physical properties.

Property-object Matching (POM) [T, E]. This task requires matching physical properties to objects: given three tactile

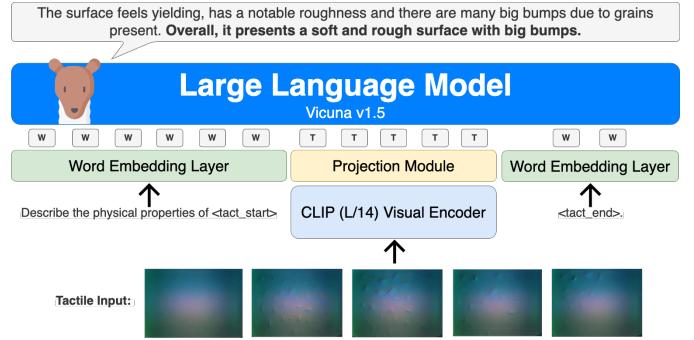


Fig. 3. OCTOPI Framework. Our framework consists of CLIP’s visual encoder, a projection module with two linear layers, and Vicuna v1.5 as the LLM. Language embeddings are derived through tokenization and then Vicuna’s word embedding layer, with $\langle \text{tact_start} \rangle$ and $\langle \text{tact_end} \rangle$ being newly trained word embeddings indicating the start and end of a tactile frame sequence from a single tactile sensor. Tactile frames are fed into the visual encoder followed by the projection module to derive tactile embeddings with the same dimension as the word embeddings.

videos (each featuring a different object) and three specified objects, the goal is to correctly associate each video with an object. This helps to align a model’s existing knowledge of object properties with our haptic perception, as our annotations are based on human touch and serve as the reference for the physical properties and their labels.

Property Scenario Reasoning (PSR) [E]. We provide two tactile videos, each showcasing a different object, along with a real-world scenario that relies on one or more of our defined physical properties. The task is to choose the video that represents the object whose physical properties best meet the scenario’s demands. This approach allows us to assess a model’s physical reasoning capabilities. Details of the scenarios are presented in Table V.

IV. OCTOPI - VISION-LANGUAGE PROPERTY-GUIDED PHYSICAL REASONING

The OCTOPI framework comprises three trained components: 1) tactile input encoder, 2) projection module, and 3) LLM, similar to prior LVLM models [34, 36, 63]. A summary of our overall framework is shown in Fig. 3.

We leverage the capabilities of pre-trained vision models, notably the CLIP [39] visual encoder ViT-L/14, as the foundation for our tactile encoder to derive meaningful feature representations. The encoder’s output is then mapped to the LLM’s word embedding space using a projection module, typically consisting of one or two trainable layers. Our projection module, inspired by LLaVA [34, 33], employs two linear layers with an intermediate GELU activation [21]. Lastly, the LLM serves as the language understanding component in OCTOPI. The performance of the LLM is largely influenced by its pre-training datasets. We utilize the open-source LLaMA-based LLM, Vicuna [11], recognized for its dialogue capabilities.

The inference process is illustrated in Fig. 3. OCTOPI receives an instruction to evaluate the physical properties of

TABLE III. PHYSICLEAR Task Details. PHYSICLEAR contains five physical description and understanding tasks. We give each task's motivation and indicate whether they are used for OCTOPI's training and/or evaluation.

Task	Motivation	Training [T]	Evaluation [E]
Object Property Description (OPD)	LLM learns to describe tactile videos using physical properties.	✓	✓
Property Comparison (PC)	Align LLM physical understanding with physical property descriptions.	✓	✓
Property Superlative Selection (PSS)	Align LLM physical understanding with physical property descriptions.	✓	✓
Property-object Matching (POM)	Align LLM object understanding with physical property descriptions.	✓	✓
Property Scenario Reasoning (PSR)	Evaluate LLM's ability to resolve scenarios using physical properties.	✗	✓

TABLE IV. PHYSICLEAR Task Prompt Setup. We provide the language prompts for each task in PHYSICLEAR. The LLM (i.e. ASSISTANT) responses for each task are ground-truth tactile descriptions, except the last response in *italics* for PSR. This italicized response is generated by OCTOPI-13b, since PSR is not used for training and hence there is no ground-truth response.

Task	Prompt Setup
OPD	USER: Describe the physical properties of $\langle \text{tact_start} \rangle T_1, \dots, T_N \langle \text{tact_end} \rangle$. ASSISTANT: The surface feels somewhat hard, lacks noticeable roughness and there are no textural variations or fine ridges. Overall, it presents a moderately hard and smooth surface with no bumps.
PC	USER: Is the object in the tactile video $\langle \text{tact_start} \rangle T_{1,1}, \dots, T_{1,N} \langle \text{tact_end} \rangle$ covered with bigger bumps than the one in $\langle \text{tact_start} \rangle T_{2,1}, \dots, T_{2,N} \langle \text{tact_end} \rangle$? Describe both objects before answering. ASSISTANT: First object: The surface feels firm and somewhat flexible, is smooth and there are big bumps and ridges. Overall, it presents a hard and smooth surface with big bumps. Second object: The surface feels soft, exhibits no roughness and there are small imperfections and ridges. Overall, it presents a moderately hard and slightly rough surface with small bumps. Conclusion: Yes, the first object is covered with bigger bumps.
PSS	USER: Given three tactile videos: a) $\langle \text{tact_start} \rangle T_{1,1}, \dots, T_{1,N} \langle \text{tact_end} \rangle$, b) $\langle \text{tact_start} \rangle T_{2,1}, \dots, T_{2,N} \langle \text{tact_end} \rangle$, c) $\langle \text{tact_start} \rangle T_{3,1}, \dots, T_{3,N} \langle \text{tact_end} \rangle$. Describe each object and then select the smoothest. ASSISTANT: a) The surface feels yielding, exhibits significant roughness and there are small bumps and ridges present. Overall, it presents a soft and rough surface with small bumps. b) The surface feels somewhat firm and flexible, exhibits no roughness and there are big bumps on the surface. Overall, it presents a moderately hard and smooth surface with big bumps. c) The surface feels soft and yielding, exhibits a slight roughness and there are tiny textural variations due to the fibrous structure. Overall, it presents a soft and slightly rough surface with small bumps. Conclusion: b) is the smoothest.
POM	USER: Given three tactile videos: a) $\langle \text{tact_start} \rangle T_{1,1}, \dots, T_{1,N} \langle \text{tact_end} \rangle$, b) $\langle \text{tact_start} \rangle T_{2,1}, \dots, T_{2,N} \langle \text{tact_end} \rangle$, c) $\langle \text{tact_start} \rangle T_{3,1}, \dots, T_{3,N} \langle \text{tact_end} \rangle$. Describe the object in each video, then match each video to one of the following objects in alphabetical order: 1) a mandarin orange, 2) a cotton ball, 3) a roll of masking tape. ASSISTANT: a) The surface feels firm, exhibits a slight roughness and there are small textural irregularities and fine ridges. Overall, it presents a hard and slightly rough surface with small bumps. b) The surface feels soft and yielding, exhibits no roughness and there are no prominent bumps or ridges. Overall, it presents a soft and smooth surface with no bumps. c) The surface feels soft, exhibits a slight roughness and there are tiny bumps present. Overall, it presents a soft and slightly rough surface with small bumps. Conclusion: a) is a roll of masking tape, b) is a cotton ball and c) is a mandarin orange.
PSR	USER: Describe these two tactile videos: a) $\langle \text{tact_start} \rangle T_{1,1}, \dots, T_{1,N} \langle \text{tact_end} \rangle$, b) $\langle \text{tact_start} \rangle T_{2,1}, \dots, T_{2,N} \langle \text{tact_end} \rangle$. ASSISTANT: a) The surface feels hard and solid, exhibits a lack of roughness and there are no textural variations or fine ridges. Overall, it presents a hard and smooth surface with no bumps. b) The surface feels yielding, has a subtle roughness and there are small irregularities and bumps. Overall, it presents a soft and slightly rough surface with small bumps. USER: In an emergency, which of the objects above can be used to break through thin ice covering a car outside? Select only one most appropriate object for this scenario based on physical property descriptions of the two objects. Use the format 'The most suitable object is x), because xxx'. ASSISTANT: <i>The most suitable object is a) a hard and smooth surface with no bumps, because it is more likely to penetrate the thin ice covering a car exterior and not get stuck.</i>

TABLE V. Scenario Questions with Selected Object and/or Properties. We list the three scenarios we use to evaluate our model's physical reasoning capabilities. For each scenario, we either specify the target object (along with the properties it has) or the target properties to select the target and non-target objects.

Scenario Question	Target Object	Target Properties
Which object is most suitable for removing stains from a non-stick pan without scratching it?	Microfiber cloth	hardness, roughness
Which object would be most easily grippable when wet and slippery?	-	roughness, bumpiness
In an emergency, which of the objects above can be used to break through thin ice covering a car outside?	-	hardness

uncooked rice. The text is tokenized and fed into the LLM's word embedding layer to produce word $[W]$ embeddings. A sequence of five tactile images is processed through the tactile encoder, with the output embeddings sent to the projection module to obtain the final tactile $[T]$ embeddings. Newly trained word embeddings, represented by $\langle \text{tact_start} \rangle$ and $\langle \text{tact_end} \rangle$, mark the beginning and end of the tactile data, respectively. These tactile embeddings are then merged with the word embeddings at designated positions to form the final instruction embeddings for the LLM.

We follow a three-step training methodology: (i) encoder fine-tuning, (ii) tactile feature alignment, and (iii) end-to-end fine-tuning. In the following, we describe each of these steps in greater detail.

A. Encoder Fine-tuning

Existing LViL models take natural images as input and can use CLIP's visual encoder without modification. However, our work involves vision-based tactile inputs, which marks a significant distribution shift from natural images, necessitating

additional fine-tuning to derive useful representations from these inputs.

We fine-tune our visual encoder to obtain useful representations from tactile inputs using multitask physical property classification. We adopt the architecture of ViFi-CLIP [40] so that our visual encoder can be trained on video inputs. In ViFi-CLIP, frame-level embeddings from CLIP’s visual encoder are average-pooled to obtain a video-level representation.

We then append learnable prompts to the pre-trained CLIP visual encoder ViT-L/14 following Visual Prompt Tuning (VPT) [25] and initialize ViFi-CLIP’s visual encoder with the new visual encoder. Specifically, we attach 8 task-specific learnable prompts and a shared linear layer to the input sequence of each Transformer [51] layer in the visual encoder and freeze the entire pre-trained Transformer backbone.

Finally, we add three separate classification heads to ViFi-CLIP, each of which predicts a label for one property (i.e. hardness, roughness or bumpiness), and train all three classification heads simultaneously using the cross-entropy loss. The model achieving the highest combined validation accuracy — correctly predicting all three properties for an object — is selected.

B. Tactile Feature Alignment

We discard the fine-tuned CLIP’s classification layers and use the outputs from its visual encoder as output embeddings. To align the output embeddings from the fine-tuned visual encoder with the LLM, the projection module is trained on language annotations while the encoder and the LLM are frozen. We also fine-tune the embedding layer due to the two new word tokens (i.e. <tact_start> and <tact_end>).

C. End-to-end Fine-tuning

Finally, we used end-to-end fine-tuning to improve the coherence of the LLM’s responses and increase the similarity between its responses and the language annotations. In this stage, only the visual encoder is frozen while the word embedding layer, projection module, and LLM are fine-tuned. We fine-tune the LLM using low-rank adaptation (LoRA) [23] for parameter-efficient fine-tuning.

V. EXPERIMENTAL SETUP

In this section, we evaluate the physical property prediction and reasoning capabilities of our proposed method. We design several experiments to answer the following questions:

- 1) Are our physical property predictions useful for OCTOPI to reason about everyday scenarios?
- 2) Can OCTOPI be used in real robots to help them accomplish tasks using physical reasoning?
- 3) Can OCTOPI’s understanding of the physical properties generalize to the unseen daily life objects?

A. Data Processing

The tactile videos were processed into frames. To focus on a few salient frames for better efficiency, we selected frames that have the top 30% total pixel intensity difference with their

TABLE VI. Results on PHYSICLEAR Physical Understanding Tasks. OCTOPI’s performance on physical understanding tasks improves with object property descriptions (OPD). Performance also increases with larger LLM size, with OCTOPI-13b outperforming OCTOPI-7b across all three tasks.

	Random	OCTOPI-7b	OCTOPI-7b (no OPD)
PC	33.33	48.10	46.51
PSS	33.33	74.67	39.88
POM	16.67	44.39	23.23
	Random	OCTOPI-13b	OCTOPI-13b (no OPD)
PC	33.33	55.06	40.70
PSS	33.33	84.00	39.88
POM	16.67	60.43	18.71

preceding frames. We randomly sampled 5 frames from these salient frames during training and selected 5 frames at uniform intervals from the first salient frame during evaluation. Data augmentation was performed during training in the form of random horizontal and vertical flips with 50% chance for each flip.

B. Training Hyperparameters

Encoder fine-tuning was performed for 30 epochs using the AdamW optimizer [35] with no weight decay, a learning rate of 10^{-3} , batch size of 32, and a cosine annealing learning rate schedule. During tactile feature alignment, the projection module is trained using 8k PHYSICLEAR samples using the AdamW optimizer [35] with no weight decay, a learning rate of 2×10^{-5} , batch size of 16, and a cosine annealing learning rate schedule.

For end-to-end fine-tuning, both the projection module and the LLM’s LoRA parameters are trained using 3k PHYSICLEAR samples using the AdamW optimizer [35] with no weight decay, batch size of 16 and a cosine annealing learning rate schedule. Learning rates of 2×10^{-5} and 2×10^{-4} were used for the projection module and the LLM, respectively. We use a scaling factor of 256, a rank of 128, and a dropout rate of 0.05 for LoRA.

C. Training Requirements

Encoder fine-tuning took 6 hours and required less than 5GB of GPU VRAM. Tactile feature alignment together with end-to-end fine-tuning took 5 hours for OCTOPI-7b and 6.5 hours for OCTOPI-13b. We used 1 NVIDIA RTX A6000 for OCTOPI-7b and 2 NVIDIA RTX A6000s for OCTOPI-13b.

VI. EXPERIMENTAL RESULTS

To address the above questions, we evaluated OCTOPI using (i) accuracy on the physical understanding tasks in PHYSICLEAR’s test set, (ii) accuracy on scenario reasoning tasks, (iii) task success rate on a real robot, and (iv) property prediction accuracy on unseen objects. We tested two versions of OCTOPI, OCTOPI-7b and OCTOPI-13b, which use Vicuna-7b v1.5 and Vicuna-13b v1.5 as their LLMs respectively.

TABLE VII. Results on PHYSICLEAR Scenario Reasoning Tasks. During scenario reasoning, we do not provide ground-truth property descriptions. Our experiments show that leveraging object properties significantly improves scenario reasoning for OCTOPI.

	Random	OCTOPI-7b	OCTOPI-13b
PSR	50.00	69.57	67.39
PSR (w/o OPD)	50.00	63.04	39.13

A. Tactile-grounded Physical Understanding with Object Property Descriptions

During tactile feature alignment and end-to-end fine-tuning, we trained OCTOPI with comparison tasks (i.e. PC, PSS and POM) to align its physical understanding of our physical properties and objects with our labels. We evaluated OCTOPI’s physical understanding with the same single-step prompts used during training and on 500 question-answer pairs in total across the three tasks. The results for physical understanding of unseen test objects are shown in Table VI.

Our results show that both OCTOPI-7b and OCTOPI-13b perform well on all three physical understanding tasks when they are trained to predict property descriptions. Using physical property descriptions, OCTOPI-7b achieves accuracies of 48.10% on PC, 74.67% on PSS and 44.39% on POM. OCTOPI-13b outperforms OCTOPI-7b by 6.96% on PC, 9.33% on PSS and 16.04% on POM. This suggests that OCTOPI’s physical understanding improves significantly with LLM size.

Further, we explored the effect of using physical property descriptions by fine-tuning both OCTOPI-7b and OCTOPI-13b on the physical understanding tasks without intermediate physical property predictions. We found that predictions based on object properties notably improve physical understanding in both OCTOPI-7b and OCTOPI-13b.

B. Scenario Reasoning

We assessed the usefulness of our physical property categories by testing how OCTOPI can reason about everyday scenarios using the physical properties. For reference, the different scenario questions are provided in Table V with the prompts shown in Table IV.

Our results are summarized in Table VII. For both OCTOPI-7b and OCTOPI-13b, including the object property significantly improves performance, which supports our overall hypothesis that leveraging these properties is helpful for these tasks. Interestingly, we observed that the 7b model marginally outperformed the 13b model.

We provide two qualitative examples to show OCTOPI-13b performing commonsense physical reasoning effectively. In the first task, we provide a tactile video of a scoop of uncooked rice and first instruct it to describe the tactile video. We then follow up with an instruction to OCTOPI-13b to determine if the rice is uncooked or cooked. OCTOPI-13b is able to reason that the scoop of rice is uncooked due to its rough surface, as shown in Fig. 4.

Next, we gave OCTOPI-13b two tactile videos corresponding to two different parts of the same toothbrush - the handle and the bristles. It is instructed to describe both objects using

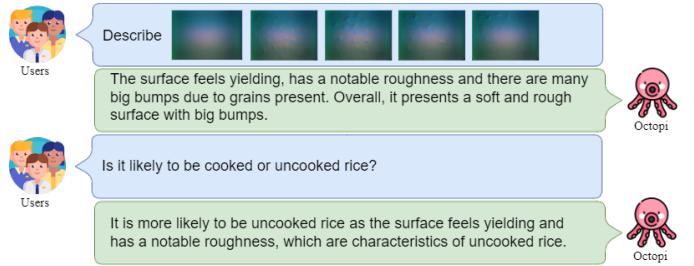


Fig. 4. Rice (Cooked v.s. Uncooked) Reasoning. OCTOPI-13b is prompted to reason about whether a scoop of rice is more likely to be cooked or uncooked based on a tactile video of a scoop on uncooked rice. It reasons about the rice state correctly without being trained to do so.

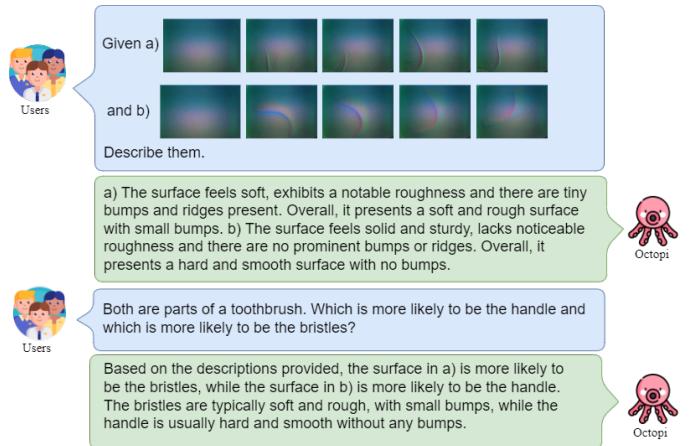


Fig. 5. Toothbrush Part Reasoning. Given a tactile video of a toothbrush’s handle and the same toothbrush’s bristles, OCTOPI-13b is prompted to reason which tactile readings belong to the handle and which belongs to the bristles.

the physical properties. We then instruct it to determine which tactile video belongs to each object part using the physical properties. Fig. 5 shows that OCTOPI-13b is able to reason about the property-object match correctly.

C. Avocado Ripeness Classification

To evaluate OCTOPI’s usefulness as a tactile-grounded physical reasoning system for real world tasks, we integrated two GelSight sensors on a 7-DoF Franka Emika Panda robot and used it for avocado ripeness classification. While ripe avocados generally appear in a shade of brown, their ripeness is difficult to determine using vision alone. At the same time, ripe avocados are softer than unripe ones and thus, tactile sensations can improve classification.

We performed property prediction and ripeness classification evaluations using a set of 10 avocados with 20 tactile samples collected from each avocado (i.e. 200 total samples). During ripeness classification, 100 pairs of avocado samples were selected and OCTOPI was tasked to identify which avocado is riper. Each pair contains avocados in different stages of ripeness. At test time, the Franka robot grasped each avocado once to collect the tactile readings, before passing it to the model. The model gave instructions on which avocado

TABLE VIII. Avocado Property Prediction and Ripeness Classification Results. OCTOPI-13b predicts avocado properties reasonably well with only a pressing motion. For avocado ripeness classification, OCTOPI-13b is able to leverage its commonsense knowledge to use both *hardness* and *bumpiness* properties.

	Random	OCTOPI-13b	PG-InstructBLIP
Property Prediction	3.70	35.50	0.00
Hardness	33.33	57.50	37.50
Roughness	33.33	71.00	3.00
Bumpiness	33.33	64.00	9.50
Ripeness Classification	50.00	63.00	-

to place in the ripe bin, and this was relayed through ROS for execution.

Table VIII shows OCTOPI-13b’s performance on both avocado property prediction and avocado ripeness classification. We split the evaluation into two parts and evaluate the avocado property prediction and ripeness classification separately.

For avocado property prediction, OCTOPI-13b achieves an accuracy of 35.50%, which is significantly higher than that of the random baseline (3.70%). Unlike the PHYSICLEAR dataset, these tactile videos are collected with only pressing and without any rotation. This suggests that OCTOPI-13b’s physical property prediction capability is robust to differences in tactile exploratory procedures.

OCTOPI-13b obtains a ripeness prediction accuracy of 63.00%. To set up this task, we first elicit OCTOPI’s commonsense reasoning for avocado ripeness using our properties by querying “You will be given tactile descriptions that consist of three physical properties: hardness, roughness, bumpiness. Hardness categories in increasing hardness: soft, moderately hard, hard. Roughness categories in increasing roughness: smooth, slightly smooth, rough. Bumpiness categories in increasing bumpiness: no bumps, small bumps, big bumps. Which of these properties help to determine avocado ripeness? Rank them.” OCTOPI-13b responds with “The properties that help determine avocado ripeness are hardness and bumpiness. Roughness is not a reliable indicator of avocado ripeness. In terms of hardness, a ripe avocado will be moderately hard, while an unripe avocado will be hard. In terms of bumpiness, a ripe avocado will have small bumps, while an unripe avocado will have no bumps or small bumps.” This exchange is provided as context to OCTOPI-13b during ripeness classification.

For comparison, we evaluated a physically-grounded vision-language model, PG-InstructBLIP [17], on property predictions of our avocados. PG-InstructBLIP was trained to infer a predetermined set of physical properties from *visual* images of real objects in the EgoObjects dataset [65]. Table VIII shows PG-InstructBLIP’s performance on property prediction for our avocados was poor. Possible reasons for this are that (i) the definitions of the physical properties may not be well-aligned with PHYSICLEAR, and/or (ii) the physical properties of avocados are not clearly apparent using only the visual modality. We could not coax the PG-InstructBLIP model to directly classify avocado ripeness despite trying various prompts; it would always pick the first object.

TABLE IX. Results on PHYSICLEAR Object Property Description Test Set. FT CLIP is the combination of the fine-tuned CLIP visual encoder and the three separate trained classification layers. OCTOPI-7b and OCTOPI-13b perform above the random baseline for object property predictions and have similar performance to the fine-tuned CLIP. OCTOPI-13b performs better than OCTOPI-7b on the prediction task.

	Random	FT CLIP	OCTOPI-7b	OCTOPI-13b
Combined	3.70	57.89	47.37	55.26
Hardness	33.33	86.84	71.05	73.68
Roughness	33.33	76.32	73.68	78.95
Bumpiness	33.33	71.05	81.58	78.95

TABLE X. CLIP Fine-tuning Ablation Results on Object Property Prediction. FT refers to fine-tuned. Using the CLIP fine-tuned on property prediction improves OCTOPI’s performance in property prediction.

OCTOPI-7b (FT CLIP)	OCTOPI-7b (base CLIP)
Combined	47.37
Hardness	71.05
Roughness	73.68
Bumpiness	81.58

OCTOPI-13b (FT CLIP)	OCTOPI-13b (base CLIP)
Combined	55.26
Hardness	73.68
Roughness	78.95
Bumpiness	78.95

D. Object Property Description Prediction

The physical understanding and scenario reasoning capabilities of OCTOPI depends on its initial physical property predictions. We evaluated OCTOPI’s physical property prediction on the PHYSICLEAR test set and show the results in Table IX. Both OCTOPI-7b and OCTOPI-13b perform well above the random baseline for combined and individual property prediction and have similar performance to the fine-tuned CLIP model, indicating that OCTOPI can be used for object property prediction. OCTOPI-13b has a higher combined accuracy (i.e. all three physical properties are correctly predicted for a given object) when compared to OCTOPI-7b, suggesting there are performance gains with larger LLMs for tactile signal grounding (apart from the *bumpiness* property).

VII. ABLATIONS

In this section, we describe ablation studies to examine (i) the impact of the encoder’s learned representations on physical property prediction and (ii) the influence of end-to-end fine-tuning data quantity on physical reasoning. For the following sections, we report test accuracy on unseen objects.

A. Ablation: The Impact of Encoder Fine-tuning

We used vision-based tactile inputs in this work and pre-trained vision foundation models (i.e. CLIP) have shown impressive performance on vision tasks. To test whether additional fine-tuning improves the pre-trained CLIP encoder’s representations for physical property prediction using tactile images, we conducted ablation experiments. We compared the performance of two OCTOPI versions — one trained with the off-the-shelf CLIP encoder and the other trained with the fine-tuned CLIP encoder.

TABLE XI. CLIP Fine-tuning Ablation Results on Physical Understanding Tasks. Using a fine-tuned CLIP improves OCTOPI’s performance in physical understanding tasks for both OCTOPI-7b and OCTOPI-13b.

OCTOPI-7b (fine-tuned CLIP)		OCTOPI-7b (base CLIP)
PC	48.10	30.38
PSS	74.67	42.67
POM	44.39	36.36
PSR	69.57	69.57
OCTOPI-13b (fine-tuned CLIP)		OCTOPI-13b (base CLIP)
PC	55.06	47.47
PSS	84.00	75.33
POM	60.43	57.22
PSR	67.39	45.65

TABLE XII. End-to-end Fine-tuning Physical Property Prediction Result Comparisons. End-to-end fine-tuning with LoRA generally improves physical property prediction accuracies.

OCTOPI-7b (w/ LoRA)		OCTOPI-7b (w/o LoRA)
Combined	47.37	39.47
Hardness	71.05	65.79
Roughness	73.68	76.32
Bumpiness	81.58	71.05
OCTOPI-13b (w/ LoRA)		OCTOPI-13b (w/o LoRA)
Combined	55.26	23.68
Hardness	73.68	36.84
Roughness	78.95	73.68
Bumpiness	78.95	71.05

In Table X, our *Object Property Description* results show that OCTOPI-7b trained with a fine-tuned CLIP encoder outperforms one trained with an unmodified CLIP encoder by 7.90% on combined accuracy. Similarly, OCTOPI-13b with the fine-tuned CLIP visual encoder performs better on the combined, roughness, and bumpiness predictions, with the combined accuracy being 5.26% higher. This suggests that a fine-tuned CLIP generally improves its learned representations for physical property prediction in an end-to-end LVLM.

We further tested both OCTOPI versions on physical understanding tasks with results in Table XI. For OCTOPI-7b, the version trained with a fine-tuned CLIP encoder performs better across the three physical understanding tasks (by 17.72% on PC, 32.00% on PSS, 8.03% on POM). Similarly, OCTOPI-13b with the fine-tuned CLIP encoder has a better performance for physical understanding tasks, which suggests that fine-tuning generally helps physical understanding and physical reasoning performance. Further encoder analysis can be found in Appendix E.

B. Ablation: The Impact of End-to-end Fine-tuning

Table XII shows OCTOPI’s performance on the property prediction task before and after end-to-end fine-tuning with LoRA. For both OCTOPI-7b and OCTOPI-13b, the fine-tuned variants generally performed better. We see sharp improvements for OCTOPI-13b with improvements across the properties. Our results suggest that end-to-end fine-tuning improves physical property prediction accuracy. Similar to the property prediction task, we observed that fine-tuning with LoRA also improves OCTOPI’s performance on physical understanding tasks (Table XIII).

TABLE XIII. End-to-end Fine-tuning Physical Understanding Result Comparisons. End-to-end fine-tuning for physical understanding tasks significantly improves physical understanding for both OCTOPI-7b and OCTOPI-13b.

Rand.	OCTOPI-7b (LoRA)	OCTOPI-7b (w/o LoRA)
PC	33.33	48.10
PSS	33.33	74.67
POM	16.67	44.39
PSR	50.00	69.57
Rand.	OCTOPI-13b (LoRA)	OCTOPI-13b (w/o LoRA)
PC	33.33	55.06
PSS	33.33	84.00
POM	16.67	60.43
PSR	50.00	67.39

VIII. CONCLUSION AND DISCUSSION

In this work, we extended large vision-language models (LVLMs) to process and describe tactile inputs using physical properties. We proposed a tactile dataset called PHYSICLEAR, comprising data from vision (Camera) and tactile (GelSight) sensors collected from everyday objects, along with physical property annotations. We also present OCTOPI, a large tactile-language model trained using datasets like PHYSICLEAR to perform physical property reasoning using tactile inputs.

Our experiments show that OCTOPI is able to describe tactile signals from novel unseen objects and that inferred physical properties can be used for physical reasoning and robot task completion in scenarios with visual ambiguity. We studied the impact of different components in OCTOPI, and found that using a task-specific visual encoder that is fine-tuned on our labels improves performance significantly across all tasks. This suggests that improvements to the visual encoder will yield benefits. In addition, parameter-efficient LLM fine-tuning consistently improved performance.

Our work opens up future work in tactile robotics. We are currently working on tactile encoder improvement and more diverse exploratory procedures to obtain additional physical properties. It would also be interesting to combine different datasets (e.g., those using other tactile sensors [44, 45]), along with other modalities such as robot proprioception. We plan to also perform physical understanding alignment with object images and LLM fine-tuning with additional physical understanding data [52, 31].

ACKNOWLEDGEMENTS

This research is supported by the National Research Foundation, Singapore under its Medium Sized Center for Advanced Robotics Technology Innovation.

REFERENCES

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35: 23716–23736, 2022.

- [2] Stéphane Aroca-Ouellette, Cory Paik, Alessandro Roncone, and Katharina Kann. [Prost: Physical reasoning of objects through space and time](#). *arXiv preprint arXiv:2106.03634*, 2021. URL <https://arxiv.org/pdf/2106.03634.pdf>.
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. [Qwen-vl: A frontier large vision-language model with versatile abilities](#). *arXiv preprint arXiv:2308.12966*, 2023. URL <https://arxiv.org/pdf/2308.12966.pdf>.
- [4] Anton Bakhtin, Laurens van der Maaten, Justin Johnson, Laura Gustafson, and Ross Girshick. [PHYRE: A New Benchmark for Physical Reasoning](#). 2019. URL <https://arxiv.org/pdf/1908.05656.pdf>.
- [5] Wouter M. Bergmann Tiest. Tactual perception of material properties. *Vision Research*, 50(24):2775–2782, 2010. ISSN 0042-6989. doi: <https://doi.org/10.1016/j.visres.2010.10.005>. URL <https://www.sciencedirect.com/science/article/pii/S0042698910004967>. Perception and Action: Part I.
- [6] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. [Piqa: Reasoning about physical commonsense in natural language](#). In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020. URL <https://arxiv.org/pdf/1911.11641.pdf>.
- [7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. [Rt-2: Vision-language-action models transfer web knowledge to robotic control](#). *arXiv preprint arXiv:2307.15818*, 2023. URL <https://arxiv.org/pdf/2307.15818.pdf>.
- [8] Guanqun Cao, Jiaqi Jiang, Danushka Bollegala, and Shan Luo. [Learn from Incomplete Tactile Data: Tactile Representation Learning with Masked Autoencoders](#). In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10800–10805. IEEE, 2023. URL <https://arxiv.org/pdf/2307.07358.pdf>.
- [9] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhosseiny. [Minigpt-v2: large language model as a unified interface for vision-language multi-task learning](#). *arXiv preprint arXiv:2310.09478*, 2023. URL <https://arxiv.org/pdf/2310.09478.pdf>.
- [10] X. Chen, Fei Shao, Cathy Barnes, Tom Childs, and Brian Henson. [Exploring Relationships between Touch Perception and Surface Physical Properties](#). *International Journal of Design*, 3:67–76, 08 2009. URL <https://arxiv.org/pdf/1704.03822.pdf>.
- [11] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#), March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- [12] Jiafei Duan, Samson Yu, and Cheston Tan. [Space: A simulator for physical interactions and causal learning in 3d environments](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2058–2063, 2021.
- [13] Jiafei Duan, Samson Yu, Soujanya Poria, Bihan Wen, and Cheston Tan. [PIP: Physical Interaction Prediction via Mental Simulation with Span Selection](#). In *European Conference on Computer Vision*, pages 405–421. Springer, 2022. URL http://phys101.csail.mit.edu/papers/phys101_bmvc.pdf.
- [14] Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and Cheston Tan. [A survey of embodied ai: From simulators to research tasks](#). *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2):230–244, 2022.
- [15] Jiafei Duan, Yi Ru Wang, Mohit Shridhar, Dieter Fox, and Ranjay Krishna. [Ar2-d2: Training a robot without a robot](#). *arXiv preprint arXiv:2306.13818*, 2023.
- [16] Letian Fu, Gaurav Datta, Huang Huang, William Chung-Ho Panitch, Jaimyn Drake, Joseph Ortiz, Mustafa Mukadam, Mike Lambeta, Roberto Calandra, and Ken Goldberg. [A touch, vision, and language dataset for multimodal alignment](#). *arXiv preprint arXiv:2402.13232*, 2024.
- [17] Jensen Gao, Bidipta Sarkar, Fei Xia, Ted Xiao, Jiajun Wu, Brian Ichter, Anirudha Majumdar, and Dorsa Sadigh. [Physically grounded vision-language models for robotic manipulation](#). *arXiv preprint arXiv:2309.02561*, 2023. URL <https://arxiv.org/pdf/2309.02561.pdf>.
- [18] Ruohan Gao, Zilin Si, Yen-Yu Chang, Samuel Clarke, Jeannette Bohg, Li Fei-Fei, Wenzhen Yuan, and Jiajun Wu. [Objectfolder 2.0: A multisensory object dataset for sim2real transfer](#). In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10598–10608, 2022. URL <https://arxiv.org/pdf/2204.02389.pdf>.
- [19] Ruohan Gao, Yiming Dou, Hao Li, Tanmay Agarwal, Jeannette Bohg, Yunzhu Li, Li Fei-Fei, and Jiajun Wu. [The ObjectFolder Benchmark: Multisensory Learning With Neural and Real Objects](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17276–17286, June 2023. URL <https://arxiv.org/pdf/2306.00956.pdf>.
- [20] Yang Gao, Lisa Anne Hendricks, Katherine J Kuchenbecker, and Trevor Darrell. [Deep learning for tactile understanding from visual and haptic data](#). In *2016 IEEE international conference on robotics and automation (ICRA)*, pages 536–543. IEEE, 2016. URL <https://arxiv.org/pdf/1511.06065.pdf>.
- [21] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [22] Yining Hong, Zishuo Zheng, Peihao Chen, Yian Wang, Junyan Li, and Chuang Gan. [Multiply: A multisensory object-centric embodied large language model in 3d world](#). *arXiv preprint arXiv:2401.08577*, 2024.
- [23] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu

- Chen. [Lora: Low-rank adaptation of large language models](#). *arXiv preprint arXiv:2106.09685*, 2021. URL <https://arxiv.org/pdf/2106.09685.pdf>.
- [24] Hung-Jui Huang, Xiaofeng Guo, and Wenzhen Yuan. [Understanding dynamic tactile sensing for liquid property estimation](#). *arXiv preprint arXiv:2205.08771*, 2022. URL <https://arxiv.org/pdf/2205.08771.pdf>.
- [25] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pages 709–727. Springer, 2022.
- [26] Jiaqi Jiang and Shan Luo. Robotic perception of object properties using tactile sensing. In *Tactile Sensing, Skill Learning, and Robotic Dexterous Manipulation*, pages 23–44. Elsevier, 2022.
- [27] Roberta Klatzky and Catherine L Reed. Haptic exploration. *Scholarpedia of Touch*, pages 177–183, 2016.
- [28] Mark H. Lee. Tactile sensing: New directions, new challenges. *The International Journal of Robotics Research*, 19(7):636–643, 2000. doi: 10.1177/027836490001900702. URL <https://doi.org/10.1177/027836490001900702>.
- [29] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. [Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models](#). *arXiv preprint arXiv:2301.12597*, 2023. URL <https://dl.acm.org/doi/10.5555/3618408.3619222>.
- [30] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. [Videochat: Chat-centric video understanding](#). *arXiv preprint arXiv:2305.06355*, 2023. URL <https://arxiv.org/pdf/2305.06355.pdf>.
- [31] Lei Li, Jingjing Xu, Qingxiu Dong, Ce Zheng, Xu Sun, Lingpeng Kong, and Qi Liu. Can language models understand physical concepts? In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11843–11861, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.726. URL <https://aclanthology.org/2023.emnlp-main.726>.
- [32] Rui Li and Edward H. Adelson. Sensing and recognizing surface textures using a gelsight sensor. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1241–1247, 2013. doi: 10.1109/CVPR.2013.164.
- [33] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.
- [34] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. [Visual instruction tuning](#). *arXiv preprint arXiv:2304.08485*, 2023. URL <https://arxiv.org/pdf/2304.08485.pdf>.
- [35] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [36] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. [Video-ChatGPT: Towards Detailed Video Understanding via Large Vision and Language Models](#). *arXiv preprint arXiv:2306.05424*, 2023. URL <https://arxiv.org/pdf/2306.05424.pdf>.
- [37] Andrew Melnik, Robin Schiewer, Moritz Lange, Andrei Muresanu, Mozghan Saeidi, Animesh Garg, and Helge Ritter. [Benchmarks for Physical Reasoning AI](#). *arXiv preprint arXiv:2312.10728*, 2023. URL <https://arxiv.org/pdf/2312.10728.pdf>.
- [38] Matthew Purri and Kristin Dana. [Teaching cameras to feel: Estimating tactile physical properties of surfaces from images](#). In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVII 16*, pages 1–20. Springer, 2020. URL <https://arxiv.org/pdf/2004.14487.pdf>.
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. [Learning transferable visual models from natural language supervision](#). In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. URL <https://arxiv.org/pdf/2103.00020.pdf>.
- [40] Hanoona Rasheed, Muhammad Uzair Khattak, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Fine-tuned clip models are efficient video learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6545–6554, 2023.
- [41] Maki Sakamoto and Junji Watanabe. [Exploring Tactile Perceptual Dimensions Using Materials Associated with Sensory Vocabulary](#). *Frontiers in Psychology*, 8, 2017. URL <https://api.semanticscholar.org/CorpusID:14038261>.
- [42] Mehmet Saygin Seyfioglu, Wisdom O Ikezogwo, Fatemeh Ghezloo, Ranjay Krishna, and Linda Shapiro. [Quilt-LLaVA: Visual Instruction Tuning by Extracting Localized Narratives from Open-Source Histopathology Videos](#). *arXiv e-prints*, pages arXiv–2312, 2023. URL <https://arxiv.org/pdf/2312.04746.pdf>.
- [43] Kuniyuki Takahashi and Jethro Tan. [Deep visuo-tactile learning: Estimation of tactile properties from images](#). In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8951–8957. IEEE, 2019. URL <https://arxiv.org/pdf/1803.03435.pdf>.
- [44] Tasbolat Taunyazov, Weicong Sng, Hian Hian See, Brian Lim, Jethro Kuan, Abdul Fatir Ansari, Benjamin CK Tee, and Harold Soh. [Event-driven visual-tactile sensing and learning for robots](#). *arXiv preprint arXiv:2009.07083*, 2020. URL <https://arxiv.org/pdf/2009.07083.pdf>.
- [45] Tasbolat Taunyazov, Luar Shui Song, Eugene Lim, Hian Hian See, David Lee, Benjamin CK Tee, and Harold Soh. Extended tactile perception: Vibration sensing through tools and grasped objects. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1755–1762. IEEE, 2021.
- [46] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricu,

- Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. *Gemini: a family of highly capable multimodal models*. *arXiv preprint arXiv:2312.11805*, 2023. URL <https://arxiv.org/pdf/2312.11805.pdf>.
- [47] Johan Tegin and Jan Wikander. Tactile sensing in intelligent robotic manipulation—a review. *Industrial Robot: An International Journal*, 32, 02 2005. doi: 10.1108/01439910510573318.
- [48] Stephen Tian, Frederik Ebert, Dinesh Jayaraman, Mayur Mudigonda, Chelsea Finn, Roberto Calandra, and Sergey Levine. *Manipulation by feel: Touch-based control with deep predictive models*. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 818–824. IEEE, 2019. URL <https://arxiv.org/pdf/1903.04128.pdf>.
- [49] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. *Llama: Open and efficient foundation language models*. *arXiv preprint arXiv:2302.13971*, 2023.
- [50] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. *Llama 2: Open foundation and fine-tuned chat models*. *arXiv preprint arXiv:2307.09288*, 2023. URL <https://arxiv.org/pdf/2307.09288.pdf>.
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [52] Yi Ru Wang, Jiafei Duan, Dieter Fox, and Siddhartha Srinivasa. *NEWTON: Are Large Language Models Capable of Physical Reasoning?*
- [53] Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. *Vary: Scaling up the Vision Vocabulary for Large Vision-Language Models*. *arXiv preprint arXiv:2312.06109*, 2023. URL <https://arxiv.org/pdf/2312.06109.pdf>.
- [54] Jiajun Wu, Joseph J Lim, Hongyi Zhang, Joshua B Tenenbaum, and William T Freeman. *Physics 101: Learning Physical Object Properties from Unlabeled Videos*. In *BMVC*, volume 2, page 7, 2016. URL http://phys101.csail.mit.edu/papers/phys101_bmvc.pdf.
- [55] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. *Next-gpt: Any-to-any multimodal llm*. *arXiv preprint arXiv:2309.05519*, 2023. URL <https://arxiv.org/pdf/2309.05519.pdf>.
- [56] Fengyu Yang, Chenyang Ma, Jiacheng Zhang, Jing Zhu, Wenzhen Yuan, and Andrew Owens. *Touch and go: Learning from human-collected vision and touch*. *arXiv preprint arXiv:2211.12498*, 2022. URL <https://arxiv.org/pdf/2211.12498.pdf>.
- [57] Fengyu Yang, Chao Feng, Ziyang Chen, Hyoungseob Park, Daniel Wang, Yiming Dou, Ziyao Zeng, Xien Chen, Rit Gangopadhyay, Andrew Owens, et al. Binding touch to everything: Learning unified multimodal tactile representations. *arXiv preprint arXiv:2401.18084*, 2024.
- [58] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B Tenenbaum. *Cleverer: Collision events for video representation and reasoning*. *arXiv preprint arXiv:1910.01442*, 2019. URL <https://arxiv.org/pdf/1910.01442.pdf>.
- [59] Wenzhen Yuan, Mandayam A. Srinivasan, and Edward H. Adelson. Estimating object hardness with a gelsight touch sensor. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 208–215, 2016. doi: 10.1109/IROS.2016.7759057.
- [60] Wenzhen Yuan, Siyuan Dong, and Edward H. Adelson. Gelsight: High-resolution robot tactile sensors for estimating geometry and force. *Sensors*, 17(12), 2017. ISSN 1424-8220. doi: 10.3390/s17122762. URL <https://www.mdpi.com/1424-8220/17/12/2762>.
- [61] Wenzhen Yuan, Yuchen Mo, Shaoxiong Wang, and Edward H Adelson. *Active clothing material perception using tactile sensing and deep learning*. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4842–4849. IEEE, 2018. URL <https://arxiv.org/pdf/1711.00574.pdf>.
- [62] Ben Zandonati, Ruohan Wang, Ruihan Gao, and Yan Wu. *Investigating Vision Foundational Models for Tactile Representation Learning*. *arXiv preprint arXiv:2305.00596*, 2023. URL <https://arxiv.org/pdf/2305.00596.pdf>.
- [63] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.
- [64] Yiyuan Zhang, Kaixiong Gong, Kaipeng Zhang, Hongsheng Li, Yu Qiao, Wanli Ouyang, and Xiangyu Yue. *Meta-transformer: A unified framework for multimodal learning*. *arXiv preprint arXiv:2307.10802*, 2023. URL <https://arxiv.org/pdf/2307.10802.pdf>.
- [65] Chenchen Zhu, Fanyi Xiao, Andrés Alvarado, Yasmine Babaei, Jiabo Hu, Hichem El-Mohri, Sean Chang, Roshan Sumbaly, and Zhicheng Yan. *Egoobjects: A large-scale egocentric dataset for fine-grained object understanding*. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- [66] Yixin Zhu, Tao Gao, Lifeng Fan, Siyuan Huang, Mark Edmonds, Hangxin Liu, Feng Gao, Chi Zhang, Siyuan Qi, Ying Nian Wu, et al. Dark, beyond deep: A paradigm shift to cognitive ai with humanlike common sense. *Engineering*, 6(3):310–345, 2020.