

# TOUCH BEYOND VISION: A SURVEY OF VISION-TACTILE-LANGUAGE MODELS IN EMBODIED INTELLIGENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Embodied intelligence increasingly leverages multimodal perception—particularly vision and language—to support rich interaction with the physical world. Yet the tactile modality remains under-explored, despite its essential role in human perception and manipulation. In this survey, we systematically review research at the intersection of vision, tactile sensing, and language, which we refer to as Vision-Tactile-Language (VTL) models. We provide (i) a historical context tracing the shift from vision-centric embodied systems to multisensory agents, (ii) foundational aspects of tactile sensing and representation, (iii) methods for integrating vision and touch, (iv) emerging architectures that incorporate language alongside vision and touch, (v) applications in embodied robotics, (vi) current challenges and open problems, and (vii) a forward-looking outlook toward tactile foundation models. We conclude by arguing that touch closes a key gap in embodied AI, enabling truly grounded perception, reasoning and action.

## 1 INTRODUCTION

Embodied artificial intelligence (AI) aims to endow agents with the ability to perceive, reason, and act in physical environments. Historically, embodied learning has been dominated by vision-centric paradigms, relying on visual perception for navigation, manipulation, and object understanding, and later extended to include language for high-level semantic reasoning and instruction following (Ma et al., 2024; Driess et al., 2023). However, in contrast to human perception—which integrates vision, audition, and touch—most AI systems still lack tactile feedback, a crucial channel for understanding contact, force, and texture (Pirozzi, 2020).

Humans rely heavily on tactile cues to perceive material properties, shape compliance, and micro-slip during manipulation. Without such feedback, robots must infer contact indirectly through vision, which often fails under occlusion, lighting variation, or complex object geometries (Yuan et al., 2023). The absence of touch limits the robustness of manipulation, sim-to-real transfer, and adaptability in contact-rich tasks. Recent reviews of tactile sensing for robotics emphasize that, despite decades of progress in tactile hardware and signal processing, the integration of tactile sensing with multimodal learning remains at an early stage (Zhang & Luo, 2025; Pirozzi, 2020).

In parallel, vision-language foundation models such as CLIP (Radford et al., 2021), Flamingo (Alayrac et al., 2022), and PaLM-E (Driess et al., 2023) have demonstrated remarkable abilities to align perception with linguistic reasoning. These successes have inspired a new question: *if models can learn from vision and language jointly, can we extend this framework to include touch?* The emerging research direction of Vision-Tactile-Language (VTL) models explores precisely this—combining tactile perception with visual and linguistic understanding to create more grounded and physically aware agents.

This paper presents a comprehensive survey of recent advances in vision-tactile-language modeling for embodied intelligence. Our goals are three-fold:

1. To systematically review tactile sensing technologies, datasets, and representation methods in the context of embodied AI.

- 054 2. To analyse how vision and touch are integrated, and how language supervision further  
055 enhances representation learning.  
056  
057 3. To identify major challenges—including data scarcity, sensor heterogeneity, cross-modal  
058 alignment, and benchmarking—and to outline future directions toward tactile foundation  
059 models.

060 Concretely, our contributions are:

- 061  
062 • We trace the historical evolution from vision-centric to multisensory embodied systems.  
063 • We summarise tactile sensing fundamentals and discuss existing tactile datasets and repre-  
064 sentation paradigms.  
065 • We review recent architectures that align or fuse vision and touch, as well as VTL models  
066 integrating language grounding.  
067 • We discuss open challenges and propose potential directions for scalable and generalisable  
068 multisensory AI.  
069

070 The remainder of this paper is structured as follows. Section 2 outlines the historical trajectory  
071 from vision-dominant perception to multisensory embodiment. Section 3 discusses the foundations  
072 of tactile sensing. Section 4 reviews vision–tactile integration. Section 5 surveys emerging VTL  
073 architectures. Section 6 highlights applications in robotic manipulation and multimodal learning.  
074 Section 7 analyses challenges and open problems, and Section 8 concludes with future outlooks.  
075

## 076 2 HISTORICAL CONTEXT: FROM VISION-CENTRIC TO MULTISENSORY 077 EMBODIMENT 078

079 The trajectory of embodied intelligence has historically been shaped by advances in visual percep-  
080 tion. Early research in robotics and AI focused primarily on vision-based sensing for navigation,  
081 recognition, and scene understanding (Anderson et al., 2018; Savva et al., 2019). Visual sensors  
082 such as RGB and depth cameras provided rich environmental information, enabling breakthroughs  
083 in object detection, visual navigation, and embodied instruction following. The dominance of visual  
084 data was partly due to its accessibility and the existence of large-scale datasets such as ImageNet and  
085 COCO, which catalyzed deep learning in computer vision. However, visual modalities alone lack  
086 direct physical grounding: they perceive geometry and appearance but not the forces, compliance,  
087 or friction that govern interaction with objects.  
088

089 **Emergence of Language in Embodied AI.** The integration of language understanding marked the  
090 second major milestone. With the advent of large-scale pretraining, models such as CLIP (Radford  
091 et al., 2021), Flamingo (Alayrac et al., 2022), and PaLM-E (Driess et al., 2023) connected vision  
092 to natural language supervision, bridging low-level perception with high-level semantics. This gave  
093 rise to embodied vision–language–action frameworks (Ma et al., 2024; Sharma et al., 2022), where  
094 agents follow textual commands (“pick up the red cup”) or answer visual questions about their sur-  
095 roundings. Yet even with such advances, these systems remain predominantly exteroceptive—they  
096 perceive the world visually and linguistically, but cannot “feel” it.

097 **Early Exploration of Tactile Sensing.** Parallel to the evolution of visual perception, tactile sens-  
098 ing research developed largely within the robotics community. Early tactile arrays and force sensors  
099 provided contact and pressure information for robotic hands, with notable examples such as BioTac,  
100 and DIGIT sensors (Lambeta et al., 2020). These technologies enabled fine-grained measurements  
101 of contact geometry, shear forces, and slip detection, offering insight into material and shape prop-  
102 erties. Reviews on tactile sensing (Pirozzi, 2020; Zhang & Luo, 2025) highlight significant progress  
103 in sensor fabrication, data collection, and representation learning. However, the tactile modality  
104 long remained siloed from mainstream machine learning due to the lack of large-scale datasets and  
105 standardized formats.

106 **From Vision–Tactile Fusion to Multisensory Embodiment.** Recent works have started to bridge  
107 this divide by combining tactile feedback with visual perception (Yuan et al., 2023; Gao et al.,

2023). More recently, the rise of multimodal foundation models (Alayrac et al., 2022; Driess et al., 2023; Chen et al., 2023) has inspired a shift from pairwise modality fusion to unified multisensory representation learning.

The evolution of embodied AI has thus progressed from vision-only perception (Dosovitskiy et al., 2021; He et al., 2016) to vision–language understanding (Radford et al., 2021; Li et al., 2022) and now toward multisensory embodiment, where tactile feedback provides the missing channel of physical grounding (Huang et al., 2023; Li et al., 2024). This historical trajectory motivates the remainder of this paper: exploring how tactile sensing can be systematically incorporated into modern multimodal learning frameworks to achieve more robust, human-like embodied intelligence.

### 3 FOUNDATIONS OF TACTILE PERCEPTION IN AI

The tactile modality provides embodied agents with a means to directly sense the physical properties of their environment. Unlike vision, which captures the appearance of objects, or language, which conveys semantic abstraction, touch delivers *grounded contact information*—forces, textures, stiffness, temperature, and deformation. This section surveys the foundations of tactile sensing in artificial intelligence, covering sensor technologies, datasets, and learned representations that enable tactile perception to interface with machine learning frameworks.

#### 3.1 TACTILE SENSORS AND DATA ACQUISITION

Modern tactile sensing technologies have evolved considerably over the past decade. Early tactile sensors were resistive or capacitive arrays that measured local pressure changes, such as the Takktile and Weiss sensors. More recent designs exploit optical or vision-based principles, where a camera observes deformations in an elastomeric surface illuminated by LEDs (Lambeta et al., 2020; Donlon et al., 2018). These *visuotactile sensors*, exemplified by DIGIT (Lambeta et al., 2020), and Omni-Tact (Padmanabha et al., 2020), provide high-resolution “tactile images” of contact geometry and can infer surface normals, slip, and local force distributions. Such sensors have enabled a transition from analog contact signals to structured visual-like data that are well-suited for convolutional and transformer-based neural architectures.

Beyond optical methods, other designs include piezoresistive (Pirozzi, 2020), capacitive, magnetic, and triboelectric sensors. Hybrid sensors combine multiple modalities (e.g., force and vibration) to capture both static and dynamic contact cues. Recent reviews emphasize the emergence of compact, inexpensive, and soft tactile sensors that can be mounted on robotic grippers and humanoid hands (Zhang & Luo, 2025; Pirozzi, 2020). However, despite rapid hardware progress, tactile data remain challenging to collect at scale due to calibration requirements, sensor fragility, and limited throughput compared to visual data collection.

#### 3.2 TACTILE DATASETS

Tactile datasets have historically been limited in both scale and diversity compared with vision and language corpora, but recent efforts are rapidly expanding their scope. Early tactile and visuotactile datasets such as Tacto (Murli et al., 2021) and TouchSim (Saal & Bensaïa, 2017) focused primarily on low-level perception and control, generating synthetic tactile data or simulating skin mechanics for tasks like grasp stability and slip detection (Calandra et al., 2017). However, tactile data remain highly dependent on physical embodiment—the same object can yield distinct signals under varying sensor geometries, contact forces, or material properties—making cross-platform generalization difficult.

To address these limitations, emerging large-scale multimodal corpora integrate tactile, visual, and linguistic modalities to connect low-level sensations with high-level semantics. Touch100k (Cheng et al., 2024) represents a major milestone in this direction, comprising over 100,000 touch–language–vision triplets with fine-grained natural-language annotations and a curriculum-based pre-training strategy (TLV-Link) that aligns tactile, visual, and linguistic representations. Such resources bridge the gap between sensory grounding and semantic reasoning, enabling scalable pretraining, zero-shot tactile understanding, and laying the foundation for next-generation embodied multimodal models.

### 3.3 TACTILE REPRESENTATIONS AND LEARNING METHODS

A central challenge in tactile perception for AI is how to represent high-dimensional and heterogeneous signals in a way compatible with learning architectures. Existing approaches can be broadly categorized as follows:

**Tactile-as-Image Representations.** Visuotactile sensors yield image-like data capturing surface deformations, which can be processed using convolutional neural networks (CNNs) or vision transformers (ViTs) (Lambeta et al., 2020). These methods treat tactile frames analogously to visual images, allowing the use of pretrained backbones for feature extraction. However, unlike visual images, tactile images encode local geometry and contact intensity rather than global appearance.

**Signal and Time-Series Representations.** For non-visual tactile arrays, the data are often treated as temporal signals—pressure, vibration, or force sequences. Recurrent or temporal convolutional networks are used to model temporal dynamics during contact events (Calandra et al., 2017). Such methods excel in recognizing slip, compliance, and dynamic interactions.

**Latent Embeddings and Cross-Modal Learning.** Recent research employs self-supervised and contrastive learning to align tactile embeddings with visual or linguistic spaces (Gomes et al., 2022; Zheng & Yuan, 2021; Gao et al., 2023). Early work demonstrated that self-supervised tactile representation learning can capture material and geometric properties without manual labels (Yuan et al., 2017; Calandra et al., 2018). More recent visuotactile frameworks exploit cross-modal contrastive objectives to map tactile and visual signals into a shared latent space, enabling bidirectional prediction and multimodal reasoning (Gomes et al., 2022; Gao et al., 2023).

Tactile perception in AI has thus evolved from simple pressure sensing to high-dimensional visuotactile imaging and multimodal representation learning (Dong et al., 2021; Li et al., 2024). Although data collection and standardization remain bottlenecks, the combination of tactile sensing with deep representation learning has paved the way for large-scale, transferable tactile understanding, forming the basis for the vision–tactile–language models discussed in Section 5.

## 4 VISION–TACTILE INTEGRATION: BRIDGING PERCEPTION MODALITIES

While tactile sensing provides local physical interaction information, vision supplies global spatial context and object appearance. Integrating these two modalities allows embodied agents to reason about both what an object *looks like* and what it *feels like*. This section reviews the main approaches for fusing vision and touch, covering architectural designs, learning objectives, and representative applications.

### 4.1 ALIGNMENT APPROACHES

The central challenge of vision–tactile integration lies in bridging heterogeneous sensory streams: tactile signals are typically localized, high-frequency, and contact-dependent, while visual data are global and spatially dense. Several approaches have been proposed to align these modalities at different levels of abstraction:

**Feature-level alignment.** Contrastive learning methods train vision and tactile encoders such that embeddings of corresponding contact events are close in latent space, while mismatched pairs are pushed apart (Calandra et al., 2017). This objective encourages cross-modal correspondence without explicit supervision, analogous to CLIP-style visual–text alignment (Radford et al., 2021).

**Cross-attention and multimodal fusion.** Beyond alignment, multimodal transformers and cross-attention architectures explicitly fuse visual and tactile tokens. Each modality is processed by its own encoder (e.g., a Vision Transformer for images, a CNN for tactile data), and information is integrated via cross-attention layers that learn inter-modal dependencies (Gao et al., 2023). Such architectures allow mutual conditioning—visual cues guide tactile feature weighting and vice versa—improving material recognition and manipulation stability.

**Geometry- and contact-based alignment.** For robotic manipulation, accurate geometric correspondence between visual and tactile frames is critical. Some works align the two modalities by estimating contact poses in the visual frame (Yuan et al., 2023; Gomes et al., 2022). Others use 3D reconstruction to integrate tactile depth into object models, effectively “completing” occluded visual surfaces using touch. These approaches bridge spatial grounding between external and contact-based sensing.

## 4.2 CROSS-MODAL LEARNING AND REPRESENTATIONS

Vision–tactile integration not only enhances perception but also enables knowledge transfer across modalities. Existing approaches can be broadly grouped into three paradigms. First, *vision supervision for tactile representation* leverages large-scale visual data to regularize tactile encoders and compensate for the scarcity of tactile datasets. Second, *tactile augmentation for visual recognition* uses touch to disambiguate visually similar yet physically distinct objects (e.g., sponge vs. rubber block); when fused with vision, tactile cues help correct uncertainty under occlusion or lighting variation and improve robustness in manipulation (Zhang & Luo, 2025). Third, several works pursue *shared latent spaces and generative modeling*, where vision and touch are jointly embedded in a unified representation that supports bidirectional generation—predicting tactile signals from vision (“imagine what it feels like”) or reconstructing appearance from touch (“see by feeling”) (Gao et al., 2023; Gomes et al., 2022). Such cross-modal representations lay the foundation for tactile simulation, data augmentation, and more general embodied multimodal learning.

## 4.3 REPRESENTATIVE WORKS

A variety of systems illustrate the effectiveness of vision–tactile integration across different architectures and objectives. Early efforts focused on explicit model-based fusion, where tactile and visual signals are combined through spatial or feature-level alignment. Robot Synesthesia (Yuan et al., 2023) introduces visuotactile sensing for in-hand manipulation using point-cloud-based fusion of tactile and visual inputs, while the VisuoTactile Transformer (Gao et al., 2023) employs cross-attention mechanisms to integrate modalities for material classification and slip detection.

Building on these approaches, recent studies pursue *unified tactile representation learning* that generalizes across sensors, materials, and interaction types. UniTouch (Yang et al., 2024) aligns tactile signals from heterogeneous vision-based sensors with visual, linguistic, and even auditory modalities through shared encoder alignment, establishing one of the first unified tactile representation frameworks. AnyTouch (Feng et al., 2025) further bridges static (image) and dynamic (video) tactile data, introducing a unified static–dynamic representation that enhances cross-sensor generalization and robust visuotactile reasoning.

## 4.4 CHALLENGES IN VISION–TACTILE INTEGRATION

Despite significant progress, several challenges remain open. First, *data imbalance* between visual and tactile samples leads to biased models, since tactile data are harder to collect. Second, *temporal synchronization* between asynchronous visual and tactile streams can degrade alignment during dynamic interactions. Third, *domain heterogeneity* across sensors and lighting conditions limits model transferability. Finally, there is no widely adopted benchmark for evaluating vision–tactile fusion, making quantitative comparison difficult (Zhang & Luo, 2025). These issues motivate the development of unified multimodal benchmarks and pretraining strategies that we discuss further in Section 5.

## 5 VISION–TACTILE–LANGUAGE MODELS (VTL MODELS)

Building upon the success of vision–language foundation models, recent research has begun exploring the inclusion of tactile sensing as a third modality to achieve multisensory embodiment. Vision–Tactile–Language (VTL) models aim to jointly learn representations that integrate visual appearance, tactile feedback, and linguistic semantics, enabling agents to perceive, reason, and act with physical grounding.

## 270 5.1 MOTIVATION AND CONCEPTUAL OVERVIEW

271  
272 Language provides a natural supervisory signal for describing physical interactions: phrases such  
273 as “soft fabric,” “rough metal,” or “the object slipped” correspond directly to tactile events. VTL  
274 models thus seek to map tactile signals and visual observations into semantic spaces where linguist-  
275 tic concepts can act as anchors. Incorporating tactile input allows grounding of language in the  
276 material and force-based properties of the physical world—an essential component of embodied  
277 understanding.

278 Conceptually, a VTL model extends the standard vision–language paradigm by introducing a tactile  
279 encoder that processes visuotactile images, pressure maps, or force sequences, and aligns its latent  
280 representations with those of vision and text encoders through joint pretraining objectives. This  
281 multimodal alignment allows queries such as “what does this surface feel like?” or “find an object  
282 that feels smooth and looks shiny” to be answered through joint reasoning.

## 283 5.2 MODEL ARCHITECTURES

284  
285 Most existing VTL systems adopt a three-branch architecture: a vision encoder, a tactile encoder,  
286 and a language encoder, each producing modality-specific embeddings that are fused via cross-  
287 modal transformers or contrastive alignment layers. Representative architectural patterns include:

288  
289 **Contrastive Tri-Modal Alignment.** Inspired by CLIP (Radford et al., 2021), models such as tac-  
290 tile, and language embeddings using triplet or multi-modal contrastive losses. Given paired visual  
291 frames, tactile readings, and natural language descriptions, the model learns a shared embedding  
292 space where semantically related samples are close together. This design supports zero-shot re-  
293 trieval across modalities (e.g., “find the tactile reading corresponding to this caption”).

294  
295 **Cross-Attention Fusion Transformers.** Other architectures employ cross-attention layers to al-  
296 low modalities to attend to each other dynamically. For instance, the VisuoTactile Transformer (Gao  
297 et al., 2023) extends to tri-modal settings, where tactile embeddings influence both vision and lan-  
298 guage representations for multimodal reasoning.

299  
300 **Unified Multisensory Encoders.** Recent research explores unified multimodal encoders that pro-  
301 cess both visual and tactile inputs as spatial tokens. Such encoders—often transformer-based—learn  
302 inter-modality dependencies implicitly, leveraging large-scale pretraining and masked modeling ob-  
303 jectives (Xie & Correll, 2025). This direction mirrors the trend toward “foundation models” in  
304 robotics and perception.

## 305 5.3 TRAINING STRATEGIES

306  
307 Training VTL models requires balancing heterogeneous data sources and objectives. Several strate-  
308 gies have emerged to effectively align visual, tactile, and linguistic modalities.

309  
310 **Joint Pretraining.** Paired visual, tactile, and textual samples are used for multimodal pretraining  
311 with objectives such as contrastive alignment, masked modeling, and captioning. Contrastive losses  
312 (e.g., InfoNCE) encourage cross-modal correspondence (Radford et al., 2021; Li et al., 2024), while  
313 masked token prediction and language-conditioned decoding leverage textual supervision (Alayrac  
314 et al., 2022; Driess et al., 2023). Such joint objectives enable models to predict linguistic descrip-  
315 tions of tactile sensations or infer contact states from text-based cues, as demonstrated in Touch-  
316 GPT (Huang et al., 2023).

317  
318 **Two-Stage Training.** When tactile data are limited, a two-stage paradigm is often adopted: first,  
319 vision–tactile encoders are pretrained using self-supervised contrastive learning (Zheng & Yuan,  
320 2021; Gao et al., 2023), then a language module is introduced through instruction tuning or caption-  
321 ing datasets (Li et al., 2022; Tu et al., 2025). This decoupled approach stabilizes optimization and  
322 improves transferability across domains.

323  
**Multitask and Reinforcement Learning.** For embodied and robotic applications, multimodal  
pretraining can be integrated with reinforcement learning, where tactile and visual feedback jointly

inform policy optimization. Language serves as a high-level goal specification, while touch provides fine-grained corrective signals for grasping or insertion tasks (Calandra et al., 2018; Lee et al., 2019; Driess et al., 2023). This combination bridges high-level reasoning with low-level control, paving the way for tactile-informed embodied agents.

#### 5.4 BENCHMARKS AND EVALUATION

Compared with vision–language evaluation tasks such as visual question answering or image captioning, benchmarks for vision–tactile–language (VTL) systems remain nascent. Existing studies typically evaluate four categories of tasks: (i) *multimodal retrieval*, where models retrieve vision–tactile pairs given a language query (e.g., “a soft rubber object”); (ii) *haptic captioning*, which generates natural language descriptions of tactile experiences; (iii) *cross-modal prediction*, involving the synthesis of tactile readings from visual input or vice versa (Gao et al., 2023); and (iv) *manipulation policy learning*, which integrates language-guided touch feedback to predict or optimize grasp success (Calandra et al., 2017).

Despite these efforts, standardized large-scale benchmarks for VTL remain unavailable. Data scarcity, heterogeneous sensor modalities, and inconsistent linguistic annotations hinder reproducibility and comparison. Ongoing community initiatives aim to construct unified visuotactile–language corpora and open-source evaluation suites—akin to COCO or LVIS—for advancing large-scale and standardized assessment in this emerging field (Zhang & Luo, 2025; Xie & Correll, 2025).

#### 5.5 REPRESENTATIVE VTL SYSTEMS

Recent systems demonstrate the expanding scope of vision–tactile–language (VTL) modeling, ranging from core tri-modal architectures to generative and reasoning-oriented frameworks. ForceFM (Xie & Correll, 2025) provides one of the earliest overviews and prototypes of touch-aware foundation models, emphasizing large-scale tactile–language data alignment and scalable pre-training strategies. Building on alignment-based foundations, newer studies explore generation and higher-level reasoning. TextToucher (Tu et al., 2025) enables fine-grained *text-to-touch* generation for controllable tactile synthesis conditioned on linguistic prompts, while Octopi (Yu et al., 2024) grounds large tactile–language models on tactile videos to infer object properties such as material, stiffness, and contact dynamics. The latter also releases the PhysiCLeAR dataset and a standardized language-evaluation benchmark, supporting quantitative assessment of tactile understanding.

#### 5.6 DISCUSSION

The emergence of VTL models marks a paradigm shift from purely visual-language perception to multisensory foundation models. By unifying the three major perceptual axes—vision, touch, and language—these systems can achieve more grounded reasoning, semantic generalization, and physical adaptability. Nonetheless, significant research gaps remain in dataset scale, modality synchronization, and open benchmarking. Bridging these challenges will be key to realizing tactile foundation models capable of zero-shot tactile reasoning and cross-modal generation.

### 6 APPLICATIONS AND EMERGING TRENDS

The integration of tactile sensing with vision and language has opened new frontiers for embodied intelligence. By enriching perception with contact-based feedback and semantic grounding, VTL systems enable agents to interact with their environments more safely, robustly, and intelligently. This section reviews major application areas where tactile information has proven beneficial and highlights emerging trends that leverage multimodal fusion in robotics and AI.

#### 6.1 ROBOTIC MANIPULATION AND GRASPING

Robotic manipulation is perhaps the most natural domain for tactile sensing. While visual perception informs object identity and geometry, tactile feedback provides crucial cues about contact state,

grip stability, and surface compliance. Integrating both modalities improves performance in grasp planning, slip detection, and in-hand manipulation.

Early work demonstrated that tactile sensing improves grasp success prediction compared to vision-only systems (Calandra et al., 2017). Using DIGIT sensors, robots can detect micro-slip and adjust grasp force in real time (Lambeta et al., 2020). Vision–tactile fusion further enhances robustness by allowing grasp corrections even under visual occlusion (Yuan et al., 2023). Recent VTL models extend these capabilities with linguistic grounding: for instance, a robot can interpret the command “grasp the softest sponge” or “pick the rough cylinder” and use tactile cues to discriminate between visually similar items.

Beyond simple grasping, tactile feedback plays an essential role in fine manipulation tasks such as peg insertion, surface polishing, and fabric folding. Touch-guided policies and contact-rich reinforcement learning (Cao et al., 2022) have shown that combining tactile feedback with visual context significantly accelerates skill acquisition and reduces physical trial failures.

## 6.2 SIMULATION-TO-REAL TRANSFER

One of the persistent challenges in robotics is the sim-to-real gap: policies trained in simulation often fail in the real world due to unmodeled contact dynamics, frictional variance, and material compliance. Tactile sensing mitigates this problem by providing *direct physical grounding* during real-world interaction. By incorporating tactile feedback into policy learning, robots can adapt on-line to deviations between simulated and actual dynamics (Murali et al., 2021; 2022).

Moreover, simulated tactile data can augment real datasets. Frameworks such as Tacto (Murali et al., 2021) and TouchSim (Saal & Bensmaia, 2017) allow large-scale synthetic data generation for pretraining tactile encoders, similar to how synthetic imagery bootstraps visual models. Domain randomization techniques are used to bridge sensor-specific variations, improving cross-platform generalization. Combined with visual and linguistic supervision, such methods pave the way toward foundation models that generalize across domains (Xie & Correll, 2025).

Beyond simulation pretraining, recent work shows that language can serve as a unifying interface to fuse heterogeneous sensors into generalist visuomotor policies. FuSe (“Beyond Sight”) fine-tunes pretrained policies to incorporate touch (and audio) via language grounding, yielding more than 20% improvement in real-world success on contact-rich manipulation tasks (Jones et al., 2025).

## 6.3 INTERACTIVE AND GROUNDED LEARNING

VTL systems enable a new paradigm of interactive embodied learning, where agents learn through continuous feedback from both humans and the environment. Language serves as a high-level communication channel, while vision and touch provide perceptual grounding. For example, a human could instruct a robot to “press harder until you feel resistance” or “rub the surface to check if it’s smooth,” linking linguistic concepts with tactile experiences.

Such systems have potential applications in human–robot collaboration, haptic teleoperation, and assistive technologies. In rehabilitation or prosthetics, for instance, VTL-based interfaces could provide language-mediated tactile feedback to improve intuitiveness and user control (Dahiya & Valle, 2019).

Beyond teleoperation, MimicTouch collects multi-modal human *tactile* demonstrations to directly teach contact-rich manipulation skills, mitigating modality mismatch (latency, occlusion) common in vision-only setups and improving policy learning from grounded touch cues.

## 6.4 EMERGING TRENDS

Several research trends are emerging at the intersection of vision, touch, and language:

- **Tactile Foundation Models.** Inspired by visual-language pretraining, researchers are building large-scale tactile datasets and encoders capable of general tactile reasoning (Xie & Correll, 2025).

- 432 • **Generative Touch Simulation.** Diffusion and generative adversarial models are being  
433 applied to synthesize tactile images or force patterns conditioned on visual or linguistic  
434 prompts, enhancing data diversity (Gao et al., 2023).
- 435
- 436 • **Integration with 3D and Physics Engines.** Combining tactile data with 3D scene un-  
437 derstanding (e.g., NeRFs or physics-based rendering) enables full-scene physical reason-  
438 ing (Yuan et al., 2023).
- 439
- 440 • **Multisensory Human–AI Interaction.** Haptic-language interfaces are emerging where  
441 users communicate tactile expectations through speech, and the agent responds both ver-  
442 bally and physically (Dahiya & Valle, 2019).

443 The applications of VTL systems extend far beyond robotics. By combining perception, language,  
444 and touch, embodied agents gain the capacity for nuanced interaction—grasping, exploring, de-  
445 scribing, and learning through experience. These capabilities form a critical step toward physically  
446 grounded AI, where multimodal fusion enables both cognitive understanding and safe, adaptive ma-  
447 nipulation of the real world.

## 448 7 CHALLENGES AND OPEN PROBLEMS

449

450 Despite the rapid progress of multisensory learning, significant challenges remain before Vi-  
451 sion–Tactile–Language (VTL) models can reach the maturity and generality of vision–language  
452 systems. This section highlights the key bottlenecks and open research questions across data, model,  
453 and evaluation dimensions, as well as ethical and safety considerations for embodied tactile AI.

### 454 7.1 DATA SCARCITY AND SENSOR DIVERSITY

455

456 A fundamental limitation in VTL research is the scarcity of large-scale, standardized tactile datasets.  
457 Unlike vision or language corpora, tactile data are inherently physical, requiring contact interactions  
458 to collect (Zhang & Luo, 2025). This process is slow, labor-intensive, and sensor-dependent. Even  
459 small variations in sensor design—such as gel stiffness, camera placement, or surface coating—can  
460 drastically alter tactile readings (Lambeta et al., 2020). As a result, tactile datasets are fragmented  
461 across research groups, each using different calibration pipelines and data formats.

462 To overcome this challenge, the community has begun developing simulation frameworks such as  
463 Tacto (Murali et al., 2021) and TouchSim (Saal & Bensmaia, 2017) to generate synthetic tactile data,  
464 analogous to how synthetic imagery has boosted vision research. However, bridging the *reality gap*  
465 between simulated and real tactile signals remains an open problem (Murali et al., 2022). Stan-  
466 dardization of data collection protocols, metadata annotation (e.g., contact force, material type), and  
467 multimodal synchronization will be crucial for scalable tactile datasets and pretraining.

### 468 7.2 CROSS-MODAL REPRESENTATION ALIGNMENT

469

470 Aligning vision, tactile, and language modalities poses significant representational and temporal  
471 challenges. Vision captures global appearance, touch captures localized contact dynamics, and lan-  
472 guage encodes symbolic abstractions. Synchronizing and fusing such heterogeneous information  
473 streams is nontrivial. Temporal misalignment is common in manipulation, where tactile signals  
474 lag or lead visual observations due to actuation delays (Yuan et al., 2023). Spatial correspondence  
475 between tactile contact regions and visual frames also depends on accurate pose estimation.

476 Furthermore, tactile signals differ statistically from visual data—while visual embeddings are high-  
477 dimensional and dense, tactile signals are often sparse or low-frequency. Naïvely aligning these  
478 modalities may bias the joint representation toward visual features. Recent approaches use con-  
479 trastive or cross-attention architectures (Gao et al., 2023), yet robust tri-modal alignment strategies  
480 that respect both temporal and spatial consistency are still lacking. Future work may benefit from  
481 explicit geometric priors or physics-informed alignment mechanisms.

486 7.3 GENERALIZATION ACROSS EMBODIMENTS  
487

488 VTL models often overfit to specific hardware configurations or robot embodiments. A model  
489 trained on a particular tactile sensor or manipulator may fail when deployed on another platform with  
490 different mechanical compliance, friction, or camera placement (Pirozzi, 2020; Zhang & Luo, 2025).  
491 Achieving cross-embodiment generalization requires disentangling sensor-specific noise from in-  
492 trinsic contact representations. Potential directions include meta-learning across multiple robots,  
493 sensor-domain adaptation, or pretraining on simulated sensor ensembles (Xie & Correll, 2025).  
494

495 7.4 BENCHMARKING AND EVALUATION  
496

497 Unlike vision–language research, which benefits from standardized datasets (e.g., COCO, VQA),  
498 the VTL field lacks shared benchmarks. Existing works evaluate on disparate tasks—material clas-  
499 sification, slip detection, captioning—using inconsistent metrics. This fragmentation hampers direct  
500 comparison and reproducibility. Community-wide efforts are needed to define open evaluation pro-  
501 tocols for:

- 502 • Cross-modal retrieval and captioning tasks.
- 503
- 504 • Manipulation benchmarks integrating tactile success rates.
- 505
- 506 • Multisensory generalization tests across sensors and domains.

507 In addition, there is a need for unified embodied evaluation environments, where agents can interact  
508 physically and semantically with real or simulated objects under multimodal feedback (Savva et al.,  
509 2019).  
510

511 7.5 COMPUTATIONAL AND MODELING CHALLENGES  
512

513 Tactile data streams are high-frequency and high-dimensional, especially for visuotactile sensors  
514 that output RGB images at 60–120 Hz. Processing these signals jointly with video and language re-  
515 quires substantial computational resources. Efficient architectures for multi-stream fusion—possibly  
516 through hierarchical encoders or adaptive sampling—remain an active area of research. Moreover,  
517 large-scale tri-modal pretraining demands novel data-efficient objectives to prevent overfitting given  
518 the limited tactile data availability.  
519

520 7.6 ETHICAL, SAFETY, AND SOCIETAL CONSIDERATIONS  
521

522 As embodied AI systems increasingly interact with the physical world, safety and ethics become  
523 central concerns. Tactile-enabled robots can exert forces on humans or fragile objects, introducing  
524 new risks compared to purely visual systems. Unintended excessive pressure or unsafe contact could  
525 cause damage or harm. Therefore, tactile feedback must be integrated with strict safety constraints  
526 and transparent control policies (Dahiya & Valle, 2019).

527 From a societal perspective, the rise of VTL models also raises issues of data privacy and bias.  
528 Tactile datasets collected from human–robot interaction may inadvertently encode sensitive biomet-  
529 ric or material information. Furthermore, linguistic supervision may import semantic biases (e.g.,  
530 associating textures or shapes with subjective adjectives). Responsible development of tactile AI  
531 requires robust consent protocols, anonymization of tactile-human data, and interpretability tools  
532 for contact-based decision making.  
533

534 7.7 SUMMARY  
535

536 In summary, the challenges for VTL research span multiple levels—from sensing hardware to mul-  
537 timodal representation and societal implications. Addressing these obstacles will require interdis-  
538 ciplinary collaboration across robotics, computer vision, haptics, and natural language processing.  
539 Overcoming them is essential to unlock the next generation of embodied AI agents that can truly  
see, feel, and reason about the physical world.

## 8 FUTURE OUTLOOK

Vision–Tactile–Language (VTL) modeling represents a crucial step toward building physically grounded artificial intelligence. Yet the field remains at an early stage, with most studies focusing on small-scale prototypes or domain-specific applications. Looking ahead, we envision several transformative research directions that could shape the future of multisensory embodied AI.

### 8.1 TOWARDS TACTILE FOUNDATION MODELS

The remarkable generalization ability of vision–language foundation models suggests the potential for analogous tactile foundation models. Such models would be pretrained on large-scale multi-modal datasets encompassing visual, tactile, and linguistic modalities, enabling zero-shot tactile reasoning and open-ended physical understanding. Recent works like ForceFM (Xie & Correll, 2025) outline early frameworks for tri-modal pretraining, aligning touch with semantics and visual context. However, achieving foundation-level scale requires advances in tactile data standardization, sensor simulation, and multimodal annotation.

In the near future, self-supervised pretraining on simulated tactile data (Murali et al., 2021; Saal & Bensmaia, 2017) combined with real-world fine-tuning could yield scalable tactile representations similar to ImageNet pretraining in vision. Moreover, continual learning frameworks that allow models to update tactile knowledge from ongoing embodied experience—“feeling to learn”—will be essential for lifelong adaptation.

### 8.2 GENERATIVE TOUCH MODELS

Generative modeling represents an emerging frontier for tactile AI. Generative Touch Models aim to synthesize tactile signals conditioned on visual or linguistic input, effectively allowing AI systems to “imagine what something feels like.” For example, given an image of a fabric and the caption “smooth silk,” the model could generate corresponding tactile patterns or pressure maps (Gao et al., 2023). Such synthetic tactile data can support low-cost data augmentation, simulation training, and interactive visualization.

Diffusion models and cross-modal transformers have already demonstrated success in text-to-image and text-to-audio synthesis; extending these paradigms to touch will require modeling spatial–temporal contact dynamics and high-frequency force signals. Recent explorations of visuotactile diffusion networks and generative adversarial tactile models suggest promising results for haptic rendering and virtual reality applications (Gomes et al., 2022; Yuan et al., 2023).

### 8.3 INTEGRATION WITH 3D AND PHYSICS-BASED REPRESENTATIONS

Tactile sensing naturally complements 3D geometry and physics-based reasoning. While vision provides global 3D shape, touch offers fine-grained local surface geometry and material cues. Integrating tactile feedback into 3D reconstruction pipelines enables reconstruction of occluded surfaces and estimation of friction or compliance. Combining tactile data with differentiable physics simulators could enable embodied agents to predict future contact forces, energy transfer, or deformation, enhancing physical reasoning.

Future embodied systems may fuse 3D NeRF representations with tactile textures, producing unified “touch-aware” scene representations. Such models would allow queries like “what regions of this object feel rough or soft?” or “how does pressing here deform the surface?”—bridging perception, simulation, and reasoning in one cohesive framework.

### 8.4 MULTISENSORY HUMAN–AI COLLABORATION

The long-term vision for VTL models extends beyond robotics into multisensory human–AI interaction. By grounding communication in shared perceptual spaces, robots and humans could exchange not only words and visuals but also sensations. For instance, a collaborative robot could report “the surface feels slippery” or “contact pressure exceeds safe threshold,” while humans could instruct using tactile language such as “press gently until you feel texture.”

594 This capability is especially relevant in healthcare, teleoperation, and assistive robotics (Dahiya &  
595 Valle, 2019). Integrating tactile feedback with natural language dialogue and visual context could  
596 enable intuitive, trustable human–AI partnerships. Future interfaces may even provide bidirectional  
597 haptic feedback, allowing humans to “feel what the robot feels” in real time through wearable de-  
598 vices.

## 600 8.5 BEYOND VTL: TOWARD VTLA AND OMNI-VTLA MODELS

601  
602 Recent work extends the perception-oriented VTL paradigm toward Vision–Tactile–Language–  
603 Action (VTLA) models that couple tactile perception with language-conditioned control policies.  
604 TLA integrates tactile cues and language instructions to learn contact-rich manipulation strate-  
605 gies (Hao et al., 2025), while VTLA adds visual input and preference learning for insertion tasks,  
606 achieving high success rates and strong sim-to-real transfer (Zhang et al., 2025). Going further, Om-  
607 niVTLA unifies diverse tactile sensors through a semantic-aligned tactile Vision Transformer and  
608 introduces the large-scale ObjTac dataset with over 135k tri-modal samples (Cheng et al., 2025).  
609 Together, these models move beyond perception to embodied decision-making, representing a step  
610 toward universal multisensory foundation models.

## 611 8.6 ETHICAL AND SOCIETAL CONSIDERATIONS

612  
613 As tactile AI becomes pervasive, ethical design and inclusive datasets will be critical. Touch is an  
614 intimate modality: its misuse in human–robot interaction raises privacy and safety concerns (Dahiya  
615 & Valle, 2019). Ensuring transparency, consent, and fairness in tactile data collection will be funda-  
616 mental for socially acceptable deployment. Collaborations between engineers, ethicists, and policy-  
617 makers will be needed to define standards for responsible tactile AI.

## 618 8.7 SUMMARY

619  
620 In summary, the future of embodied intelligence lies in bridging perception, language, and physi-  
621 cal interaction. VTL models are paving the way toward multisensory foundation models that can  
622 both understand and generate tactile experiences. By uniting vision, touch, and language through  
623 scalable learning frameworks and ethical design, researchers can move closer to realizing physically  
624 grounded, adaptive, and communicative artificial agents that truly “see, feel, and understand” the  
625 world.

## 626 9 CONCLUSION

627  
628 This survey provides a comprehensive overview of recent advances at the intersection of vision,  
629 tactile sensing, and language in embodied intelligence. We traced the historical trajectory from  
630 vision-centric perception to multisensory embodiment, highlighting how tactile sensing completes  
631 the perceptual triad that underpins human interaction with the physical world. We summarized  
632 the foundations of tactile sensing technologies, datasets, and representation learning, followed by  
633 an in-depth discussion of vision–tactile integration and the emergence of Vision–Tactile–Language  
634 (VTL) models. Our review reveals a growing convergence between robotics, computer vision, and  
635 natural language processing, where tactile perception is evolving from a niche sensory channel into  
636 a key component of multimodal reasoning and embodied understanding, opening new avenues for  
637 contact-rich manipulation, interactive human–robot communication, and grounded learning.

638  
639 Despite this progress, several challenges persist: the scarcity and heterogeneity of tactile  
640 datasets (Zhang & Luo, 2025), difficulties in cross-modal alignment (Gao et al., 2023), limited cross-  
641 embodiment generalization (Xie & Correll, 2025), and the absence of standardized benchmarks for  
642 evaluating VTL models. Addressing these issues will require coordinated efforts across haptics,  
643 perception, and machine learning communities. Looking forward, we envision a new generation  
644 of tactile foundation models that unify visual, tactile, and linguistic modalities under scalable self-  
645 supervised frameworks. Such models could reason about and generate tactile experiences, simulate  
646 physical interactions, and communicate them through natural language, ultimately enabling robots  
647 that both perceive and *feel*. The integration of touch, vision, and language thus represents a pivotal  
step toward truly grounded artificial intelligence—systems that can not only see and talk about the

648 world but also *understand* it through direct physical experience. We hope this survey serves as both  
649 a roadmap and a call to action for advancing toward that multisensory future.

## 651 REFERENCES

- 652 Jean-Baptiste Alayrac et al. Flamingo: A visual language model for few-shot learning. In *NeurIPS*,  
653 2022.
- 654 Peter Anderson, Qi Wu, Damien Teney, Joel Bruce, Anton van den Hengel, et al. Vision-and-  
655 language navigation: Interpreting visually-grounded navigation instructions in real environments.  
656 In *CVPR*, 2018.
- 657 Roberto Calandra, Andrew Owens, Mukundan Upadhyaya, Wenzhen Yuan, Justin Lin, Edward H.  
658 Adelson, and Sergey Levine. The feeling of success: Does touch sensing help predict grasp  
659 outcomes? In *Conference on Robot Learning (CoRL)*, 2017.
- 660 Roberto Calandra, Andrew Owens, Dinesh Jayaraman, Justin Lin, Wenzhen Yuan, Jitendra Malik,  
661 Sergey Levine, and Edward H. Adelson. More than a feeling: Learning to grasp and regrasp using  
662 vision and touch. In *IEEE Robotics and Automation Letters (RA-L)*, volume 3, pp. 3300–3307,  
663 2018. URL <https://ieeexplore.ieee.org/document/8460901>.
- 664 Yu Cao, Xin Li, Tong Wu, and Yuke Zhu. Contact-rich manipulation via deep reinforcement learning  
665 with tactile feedback. *IEEE Transactions on Robotics*, 2022.
- 666 Guanglu Chen, Zhe Li, Yuntao Zhang, et al. Shikra: Unifying vision-language understanding  
667 and generation with instruction-finetuned multimodal large language models. *arXiv preprint*  
668 *arXiv:2306.02522*, 2023.
- 669 Ning Cheng, Changhao Guan, Jing Gao, Weihao Wang, You Li, Fandong Meng, Jie Zhou, Bin  
670 Fang, Jinan Xu, and Wenjuan Han. Touch100k: A large-scale touch-language-vision dataset for  
671 touch-centric multimodal representation. *arXiv preprint arXiv:2406.03813*, 2024.
- 672 Zhengxue Cheng, Yiqian Zhang, Wenkang Zhang, Haoyu Li, Keyu Wang, Li Song, and Hengdi  
673 Zhang. Omnivtla: Vision–tactile–language–action model with semantic-aligned tactile sensing.  
674 *arXiv preprint arXiv:2508.08706*, 2025.
- 675 Ravinder S. Dahiya and Maurizio Valle. Human–robot interaction through tactile sensing: A survey.  
676 *IEEE Transactions on Robotics*, 35(3):703–724, 2019.
- 677 Siyuan Dong, Wenzhen Yuan, and Edward H. Adelson. Gelsight imagenet: Large-scale visuotactile  
678 dataset for representation learning. *arXiv preprint arXiv:2107.14152*, 2021.
- 679 Elliot Donlon, Siyuan Dong, Moju Liu, Jie Li, and Edward H. Adelson. Gelslim: A high-resolution,  
680 compact, robust, and calibrated tactile-sensing finger. In *IEEE/RSJ IROS*, 2018.
- 681 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
682 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-  
683 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at  
684 scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- 685 Danny Driess et al. Palm-e: An embodied multimodal language model. *arXiv preprint*  
686 *arXiv:2303.03378*, 2023.
- 687 Ruoxuan Feng, Jiangyu Hu, Wenke Xia, Tianci Gao, Ao Shen, Yuhao Sun, Bin Fang, and Di Hu.  
688 Anytouch: Learning unified static–dynamic representation across multiple visuo-tactile sensors.  
689 In *International Conference on Learning Representations (ICLR)*, 2025.
- 690 Tianyu Gao, Wei Xu, Ruolin Chen, and Shan Luo. Visuotactile transformer: Cross-attention fusion  
691 for multimodal material recognition. In *IEEE/RSJ IROS*, 2023.
- 692 Pedro Gomes, Hugo Araujo, and Alexandre Bernardino. Learning visuotactile 3d object recon-  
693 struction via self-supervised cross-modal learning. *IEEE Robotics and Automation Letters*, 7(2):  
694 2765–2772, 2022.

- 702 Peng Hao, Chaofan Zhang, Dingzhe Li, Xiaoge Cao, Xiaoshuai Hao, Shaowei Cui, and Shuo  
703 Wang. Tla: Tactile-language-action model for contact-rich manipulation. *arXiv preprint*  
704 *arXiv:2503.08548*, 2025. doi: 10.48550/arXiv.2503.08548.
- 705  
706 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
707 nition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
708 *(CVPR)*, 2016.
- 709 Zhiyong Huang, Ziyu Gao, Wei Chen, Bin Fang, Hao Zhang, et al. Touchgpt: A  
710 vision–tactile–language model for multisensory embodied intelligence. *arXiv preprint*  
711 *arXiv:2308.10901*, 2023.
- 712  
713 Joshua Jones, Oier Mees, Carmelo Sferrazza, Kyle Stachowicz, Pieter Abbeel, and Sergey Levine.  
714 Beyond sight: Finetuning generalist robot policies with heterogeneous sensors via language  
715 grounding. *arXiv preprint arXiv:2501.04693*, 2025.
- 716 Mike Lambeta, Po-Wei Chou, Stephen Tian, et al. Digit: A novel design for a low-cost compact  
717 high-resolution tactile sensor with application to in-hand manipulation. In *IEEE Robotics and*  
718 *Automation Letters*, volume 5, pp. 3838–3845, 2020.
- 719  
720 Michelle A. Lee, Yuke Zhu, Krishnan Srinivasan, Silvio Savarese, Li Fei-Fei, and Animesh Garg.  
721 Touch and go: Learning to manipulate with tactile sensing and reinforcement learning. In *Pro-*  
722 *ceedings of the Conference on Robot Learning (CoRL)*, 2019.
- 723  
724 Junnan Li, Dongxu Li, Silvio Savarese, and Steven CH Hoi. Blip: Bootstrapping language-image  
725 pre-training for unified vision-language understanding and generation. In *International Confer-*  
726 *ence on Machine Learning (ICML)*, 2022.
- 727  
728 Yutong Li, Weihao Wang, Jing Gao, and Bin Fang. Vt-lm: Vision–tactile–language model for  
729 multimodal representation learning. *arXiv preprint arXiv:2402.04387*, 2024. URL <https://arxiv.org/abs/2402.04387>.
- 730  
731 Y. Ma et al. A survey on vision-language-action models for embodied ai. *arXiv preprint*  
732 *arXiv:2405.14093*, 2024.
- 733  
734 Aditya Murali, Roberto Calandra, and Wenzhen Yuan. Tactile transfer: Learning to transfer from  
735 simulated to real tactile perception. *IEEE Robotics and Automation Letters*, 7(3):6003–6010,  
736 2022.
- 737  
738 Aditya Murali et al. Tacto: A fast, flexible, and open-source simulator for tactile sensors. In  
739 *Conference on Robot Learning (CoRL)*, 2021.
- 740  
741 Ananya Padmanabha, Stephen Tian, Po-Wei Chou, et al. Omnitact: A multi-directional high-  
742 resolution touch sensor. In *IEEE Robotics and Automation Letters*, volume 5, pp. 5881–5888,  
743 2020.
- 744  
745 Salvatore Pirozzi. Tactile sensors for robotic applications. *Sensors*, 20(24):7009, 2020.
- 746  
747 Alec Radford et al. Learning transferable visual models from natural language supervision. In  
748 *ICML*, 2021.
- 749  
750 Hannes P. Saal and Sliman J. Bensmaia. Touchsim: Simulating tactile signals from the whole hand  
751 with millisecond precision. *PNAS*, 114(28):E5693–E5702, 2017.
- 752  
753 Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, et al. Habitat: A platform for embodied  
754 ai research. In *ICCV*, 2019.
- 755  
756 Pratyusha Sharma, Abhinav Gupta, and Saurabh Kumar. Learning generalizable vi-  
757 sion–language–action models for embodied tasks. In *NeurIPS*, 2022.
- 758  
759 Jiahang Tu, Hao Fu, Fengyu Yang, Hanbin Zhao, Chao Zhang, and Hui Qian. Texttoucher: Fine-  
760 grained text-to-touch generation. In *Proceedings of the AAAI Conference on Artificial Intelli-*  
761 *gence*, volume 39, pp. 7455–7463, 2025.

756 Weipeng Xie and Nikolaus Correll. Towards forceful robotic foundation models: A literature survey.  
757 *arXiv preprint arXiv:2504.11827*, 2025.  
758

759 Fengyu Yang, Chao Feng, Ziyang Chen, Hyouneseob Park, Daniel Wang, Yiming Dou, Ziyao Zeng,  
760 Xien Chen, Rit Gangopadhyay, Andrew Owens, and Alex Wong. Binding touch to everything:  
761 Learning unified multimodal tactile representations. In *Proceedings of the IEEE/CVF Conference*  
762 *on Computer Vision and Pattern Recognition (CVPR)*, pp. 26340–26353, June 2024.

763 Samson Yu, Kelvin Lin, Anxing Xiao, Jiafei Duan, and Harold Soh. Octopi: Object property rea-  
764 soning with large tactile-language models. In *Robotics: Science and Systems (RSS)*, 2024.  
765

766 Wenzhen Yuan, Siyuan Dong, and Edward H. Adelson. Gelsight: High-resolution robot tactile  
767 sensors for estimating geometry and force. *Sensors*, 17(12):2762, 2017.

768 Y. Yuan, H. Che, Y. Qin, et al. Robot synesthesia: In-hand manipulation with visuotactile sensing.  
769 *arXiv preprint arXiv:2312.01853*, 2023.  
770

771 Chaofan Zhang, Peng Hao, Xiaoge Cao, Xiaoshuai Hao, Shaowei Cui, and Shuo Wang. Vtla:  
772 Vision–tactile–language–action model with preference learning for insertion manipulation. *arXiv*  
773 *preprint arXiv:2505.09577*, 2025.

774 J. Zhang and S. Luo. Recent advances and challenges of tactile sensing for robotics: From sensors  
775 to data, representation and applications. *ScienceDirect*, 2025.  
776

777 Yilin Zheng and Wenzhen Yuan. Self-supervised learning for tactile representation with tacti-  
778 lessl. *IEEE Robotics and Automation Letters (RA-L)*, 6(4):8333–8340, 2021. URL <https://ieeexplore.ieee.org/document/9507732>.  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809