

# UltraTac: Integrated Ultrasound-Augmented Visuotactile Sensor for Enhanced Robotic Perception

Junhao Gong<sup>1\*</sup>, Kit-Wa Sou<sup>1\*</sup>, Shoujie Li<sup>1†</sup>, Changqing Guo<sup>1</sup>, Yan Huang<sup>1</sup>, Chuqiao Lyu<sup>1</sup>, Ziwei Song<sup>1</sup>, Wenbo Ding<sup>1,2†</sup>

**Abstract**— Visuotactile sensors provide high-resolution tactile information but are incapable of perceiving the material features of objects. We present UltraTac, an integrated sensor that combines visuotactile imaging with ultrasound sensing through a coaxial optoacoustic architecture. The design shares structural components and achieves consistent sensing regions for both modalities. Additionally, we incorporate acoustic matching into the traditional visuotactile sensor structure, enabling the integration of the ultrasound sensing modality without compromising visuotactile performance. Through tactile feedback, we can dynamically adjust the operating state of the ultrasound module to achieve more flexible functional coordination. Systematic experiments demonstrate three key capabilities: proximity sensing in the 3–8 cm range ( $R^2 = 0.99$ ), material classification (average accuracy: 99.20%), and texture-material dual-mode object recognition achieves 92.11% accuracy on a 15-class task. Finally, we integrate the sensor into a robotic manipulation system to concurrently detect container surface patterns and internal content, which verifies its promising potential for advanced human-machine interaction and precise robotic manipulation.

## I. INTRODUCTION

Visuotactile sensors (e.g., GelSight [1], DIGIT [2], GelSlim [3]) capture high-resolution optical images of surface deformations to reveal texture, contact location, and force distribution. Nonetheless, they acquire data only during contact and cannot detect material properties [4], [5]. This limitation challenges robotic systems requiring proximity sensing [6]. Moreover, the inability to ascertain material characteristics, including internal cavities and content identification, compromises safety and efficiency in dynamic environments, particularly when handling delicate or unfamiliar objects [7].

To address these limitations, the integration of multiple sensing modalities is explored in various contexts. The work presented in [8], [9] demonstrates a novel approach to proximity sensing by employing selective membrane

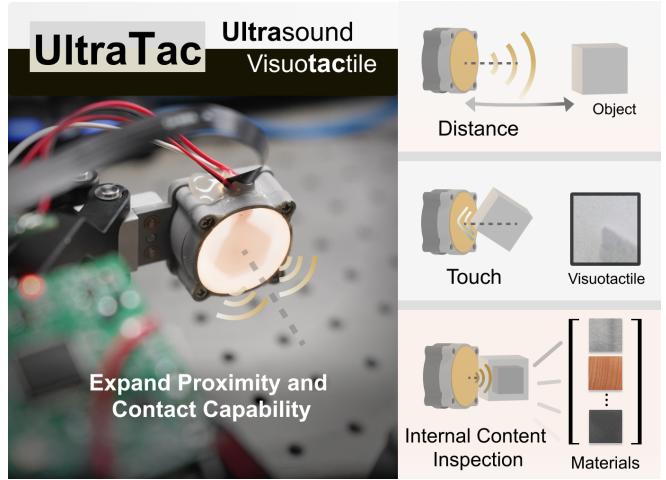


Fig. 1. Overview of UltraTac: an integrated sensor combining visuotactile and ultrasound sensing.

transparency and integrated internal light sources, thereby enabling the accurate detection of the distance between objects and sensors prior to physical contact. However, this approach encounters challenges in dark environments or those with complex lighting conditions. Furthermore, the integration of both sensing modalities using a single camera with time-division multiplexing exacerbates the complexities associated with light source control and data acquisition. Song et al. [10] employ a flexible triboelectric sensor to achieve simultaneous material and texture recognition, while Lee et al. [11] employ hybrid triboelectric, piezoelectric, and piezoresistive operating mechanisms to simultaneously distinguish between materials and textures. However, both approaches face a common challenge: extracting two distinct modalities from a unified signal may lead to crosstalk, where interference between the modalities degrades signal fidelity and classification performance.

To overcome the limitations of traditional visuotactile sensors, ultrasound sensing adds value by measuring distances without contact and analyzing materials beneath the surface. Operating in a proximity sensing mode, ultrasound measures distances before contact is made and provides material features by analyzing reflected acoustic waves [12]. Recent studies show that ultrasound sensors detect objects at distances of several centimeters with high precision [13], making them valuable for proximity sensing. Additionally, when in contact with objects, ultrasound penetrates surfaces

\*These authors contributed equally to this work.

This work was supported by National Key R&D Program of China (No.2024YFB3816000), Shenzhen Key Laboratory of Ubiquitous Data Enabling (No. ZDSYS20220527171406015), Guangdong Innovative and Entrepreneurial Research Team Program (2021ZT09L197), and Tsinghua Shenzhen International Graduate School-Shenzhen Pengrui Young Faculty Program of Shenzhen Pengrui Foundation (No. SZPR2023005), and Meituan Academy of Robotics Shenzhen.

†Corresponding author: Shoujie Li (lsj20@mails.tsinghua.edu.cn), Wenbo Ding (ding.wenbo@sz.tsinghua.edu.cn)

<sup>1</sup>Shenzhen Ubiquitous Data Enabling Key Lab, Shenzhen International Graduate School, Tsinghua University, Shenzhen 518055, China.

<sup>2</sup>RISC-V International Open Source Laboratory, Shenzhen 518055, China.

This paper has supplementary downloadable material available at: ultra-tac.junhaogong.top

to reveal internal structures and material properties that remain invisible to optical methods [14].

In this paper, we introduce UltraTac, an innovative integrated sensor that effectively tackles these challenges by seamlessly merging visuotactile imaging with ultrasound sensing via a coaxial optoacoustic design (Fig. 1). Our study presents three major contributions:

- **Novel integrated sensor architecture:** We develop a compact sensor that unifies visuotactile imaging with ultrasound through a coaxial design. By using a translucent, ultra-thin elastomer and an annular lead zirconate titanate (PZT) transducer with an engineered impedance matching stack, we achieve optimal acoustic coupling while preserving optical clarity for simultaneous tactile and ultrasound sensing.
- **Enhanced sensing capabilities via ultrasound augmentation:** Our system provides two capabilities unavailable to conventional tactile sensors: proximity sensing and contact-based material classification. These modes can be dynamically interchanged based on tactile feedback, enabling adaptive perception across diverse interaction scenarios.
- **Extensive experimental validation and application demonstrations:** Through systematic experiments, we validate the sensor's effectiveness in proximity detection, material classification, and dual-modal object recognition tasks, demonstrating its potential for robotic manipulation, content inspection, and other applications requiring object assessment.

## II. RELATED WORK

In recent years, multimodal visuotactile sensors have significantly advanced robotic perception. Comprehensive reviews by Luo et al. [15] and Dahiya et al. [16] highlight the benefits of multimodal approaches in addressing the limitations of individual sensing technologies.

Proximity sensing enables early anticipation of interactions, enhancing performance in complex manipulation tasks. CompdVision integrates a compound-eye imaging system that employs far-focus stereo units for external depth estimation and near-focus units for tactile deformation tracking, enabling accurate pre-contact object localization [17]. Similarly, the multimodal fingertip sensor by SaLoutos et al. combines embedded pressure sensors with time-of-flight proximity modules to measure contact forces and detect approaching objects [18]. Moreover, Li et al.'s M<sup>3</sup>Tac extends sensor capability by fusing visible, near-infrared, and mid-infrared imaging, thereby achieving high-resolution proximity sensing, precise three-dimensional reconstruction, and temperature measurement [19].

For material classification specifically, researchers demonstrate that combining different sensing modalities can improve discrimination between visually similar materials [20], [21]. Song et al. [22] combine tactile and temperature sensing for enhanced material recognition, while Sou et al. [23] integrate mechanoluminescence with tactile sensing for event-trigger perception. Abderrahmane et al. [24] show how

multiple sensing modalities improve object recognition even for previously unseen items.

Non-destructive testing and subsurface analysis are critical areas where traditional tactile sensing falls short. Current ultrasound methods for internal inspection [25] show promise but typically require specialized equipment separate from tactile systems. Capturing both surface properties and internal characteristics within a single compact sensor represents an advancement for applications ranging from quality control to medical diagnostics [26].

Nonetheless, seamlessly integrating tactile and ultrasound sensing into a unified sensor architecture remains an open challenge. Prior attempts to combine these modalities often yield bulky systems with separate sensing elements, limiting their practical application in robotic end-effectors [27]. Moreover, maintaining optical clarity for visuotactile sensing alongside efficient acoustic coupling for ultrasound necessitates innovative solutions.

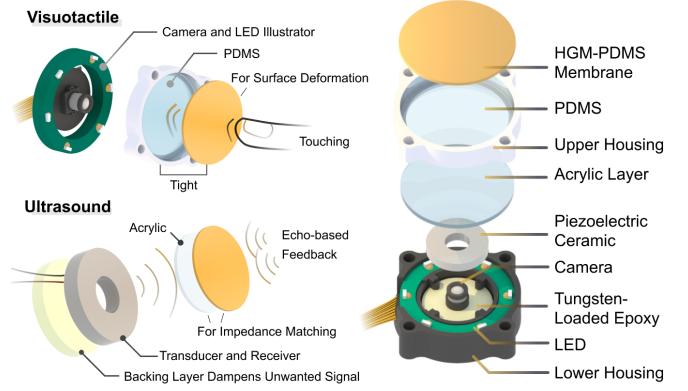


Fig. 2. Structure illustration showing the visuotactile and ultrasound modules with exploded component view.

## III. DESIGN AND IMPLEMENTATION

Traditional visuotactile sensors integrate a camera, acrylic substrate, flexible deformation layer, and light-blocking membrane for optical imaging and tactile detection. However, this structure impedes ultrasound propagation due to acoustic impedance mismatches between layers. Our design challenge was to create a structure to enable effective ultrasound transmission while preserving visuotactile capabilities.

### A. Design Principles

UltraTac delivers dual-modal perception by integrating surface texture information from camera imaging through an elastomer membrane with material properties detected via ultrasound, as shown in Fig. 2. We optimize the sensor's mechanics, materials, electronic systems, and algorithms to achieve effective dual-modal sensing capabilities.

Mechanically, UltraTac is built around a coaxial optoacoustic structure, ensuring that the camera and the ring-shaped piezoelectric ceramic (PZT) share the same central axis. This layout preserves the optical path for visuotactile sensing while simultaneously providing an unobstructed

route for ultrasound transmission and reception. This structural integration ensures consistent sensing regions for both modalities, enhancing the spatial alignment of visuotactile and ultrasound sensing within the sensor.

From a materials perspective, we introduce several unconventional materials in the visuotactile domain to achieve effective acoustic impedance matching. In parallel, the electronic systems have been designed as a highly integrated processing unit that supports both camera and ultrasound functionalities. Furthermore, from an algorithmic standpoint, the system employs a touch-triggered dual-pathway approach that dynamically adapts sensor functionality; when no contact is detected, ultrasound performs proximity sensing, whereas upon contact it switches to material classification while visual data simultaneously processes texture recognition.

TABLE I  
ACOUSTIC IMPEDANCE OF DIFFERENT MATERIALS

Materials	Acoustic Impedance (MRayls)
Tungsten	100 [28]
Epoxy	3.4 [28]
PZT	25-35 [28]
Acrylic	3.2 [29]
PDMS	1.1 [30]
Hollow Glass Microsphere	0.2 [31]
Air	0.000415 [29]

### B. Acoustic Matching

In this work, we employ acoustic impedance matching principles based on two key formulas. The first is the single-layer matching formula [29]:

$$Z_m = \sqrt{Z_1 \cdot Z_2}, \quad (1)$$

where  $Z_1$  and  $Z_2$  are the acoustic impedance of the two media and  $Z_m$  is the optimal impedance of the matching layer. The second formula is the quarter-wavelength thickness condition:

$$d = \frac{\lambda_m}{4} = \frac{c_m}{4f}, \quad (2)$$

where  $d$  is the matching layer thickness,  $\lambda_m$  is the wavelength in the layer,  $c_m$  is the speed of sound in the layer, and  $f$  is the operating frequency. These formulas, together with the data in Table I, will be used for acoustic matching.

From a materials standpoint, we aim to enhance acoustic transmission efficiency while preserving optical transparency. Typically, the light-blocking layer in sensors consists of flexible membranes infused with metal powders or dyes. However, as shown in Table I, these membranes exhibit significantly higher acoustic impedance than air, leading to substantial reflection of ultrasound signals at the interface [29]. To mitigate this loss, we incorporate hollow glass microspheres (HGM) into the dyed membrane, thereby reducing the impedance mismatch between the PDMS layer and air. HGM is widely employed to lower acoustic impedance, with its effectiveness dependent on the volume fraction and particle size [31]. In our work, HGM is mixed with dyed

PDMS at a 1:1 volume ratio to function both as a light-blocking layer for imaging and as an acoustic matching layer at the PDMS-air interface. The composite layer's optimal impedance was determined using the single-layer matching formula (Equation (1)), ensuring enhanced transmission efficiency.

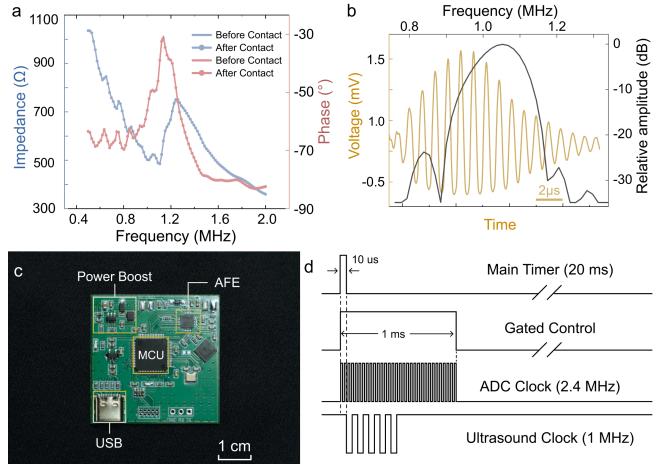


Fig. 3. Ultrasound transducer characteristics. (a) Impedance and phase response across frequency. (b) Time-domain signal pulse and frequency response. (c) The PCB of the ultrasound module. (d) Timing diagram of ultrasound module.

Under the elastomer layers, an acrylic substrate is typically used to support the deformation of the elastomer. To satisfy the acoustic impedance matching conditions, rather than altering the material type, we adjust the acrylic thickness to 0.7 mm, calculated based on the propagation speed of 1 MHz ultrasound in acrylic and the quarter-wavelength formula (Equation (2)). This optimized thickness allows the acrylic to function as a single-layer matching medium between the PZT and PDMS.

Additionally, an ultrasound sensor requires a backing layer to suppress undesired acoustic reflections. Based on Equation (1), tungsten-loaded epoxy is widely used in ultrasound applications as a backing-layer material because its acoustic impedance can be tuned by adjusting the ratio of tungsten powder to resin [28]. In our design, we employ tungsten-loaded epoxy prepared at a 3:2 volume ratio.

We conduct an acoustic analysis of the entire sensor. Fig. 3(a) presents the impedance and phase variations over the frequency range of 500 kHz to 2 MHz. The minimum impedance corresponds to the sensor's resonant frequency (1.05 MHz), while the maximum impedance is observed at the anti-resonant frequency (1.21 MHz). The marked curve represents the impedance response under applied pressure. The impedance characteristics remain largely unchanged before and after compression, indicating that the deformation of the flexible layer above the PZT does not affect the ultrasound functionality. This provides a fundamental physical basis for the integration of the tactile sensing modality. Fig. 3(b) shows the echo signal obtained when the sensor is placed on a container filled with water. The yellow waveform represents the time-domain signal, while the black waveform

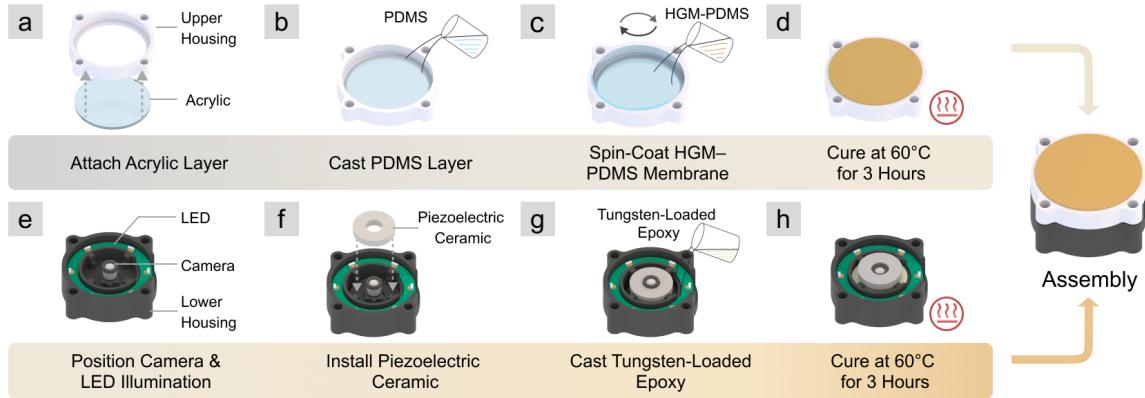


Fig. 4. Fabrication process of the sensor showing parallel assembly paths. (a-d) Upper housing with acrylic layer attachment, PDMS casting, HGM-PDMS spin-coating, and curing. (e-h) Lower housing with camera and LED positioning, piezoelectric ceramic installation, tungsten-epoxy casting, and final assembly of both components.

corresponds to the frequency-domain signal obtained after performing a Fourier transform on the time-domain data. From the information presented in the figure, we observe that the sensor's central frequency is 1.05 MHz.

### C. Electronic Systems

**Ultrasound Module:** To integrate ultrasound functionality into robotic systems, a compact and highly integrated ultrasound transmission and reception system is required. To address this need, a dedicated PCB (Fig. 3(c)) was designed ( $4\text{ cm} \times 4\text{ cm}$ ) that requires only a single USB cable for both power and communication. The PCB, featuring a power boost circuit, generates excitation voltages up to 40 V for ultrasound signals and achieves a dynamic sensing refresh rate of approximately 50 Hz, meeting the real-time demands of human–machine interaction. The system architecture comprises a 144 MHz microcontroller unit (MCU) and a highly integrated analog front-end (AFE). The MCU dynamically adjusts both the number of transmitted ultrasound pulses and the gain in the reception chain, enabling flexible modulation of ultrasound functionality.

The timing diagram of the ultrasound acquisition system is shown in Fig. 3(d). The MCU operates with a 20 ms cycle for both transmission and acquisition, utilizing a 2.4 MHz ADC to capture the echo envelope, with all triggers controlled by a hardware timer. During the idle period within the 20 ms cycle, the MCU transmits data and updates the operational status of the ultrasound system. Specifically, during the operation of the ultrasound distance measurement mode, the AFE is configured to emit five pulses per cycle to minimize the blind zone, while the reception chain gain is set to an amplification factor of 55.5 dB at 1 mV. Conversely, when the ultrasound content recognition mode is activated, the pulse count is increased to 20 in order to enhance recognition accuracy. The transition between these two operational modes is governed by haptic feedback; when the tactile modality detects contact between the sensor and an object, a contact trigger signal is transmitted to the MCU, thereby initiating the switch in ultrasound functionality.

**Optical Module:** The optical system consists of a micro-camera with a 1.88 mm focal length, positioned at the center of the annular piezoelectric ceramic. The camera lens, with a 6 mm diameter, is precisely integrated within the ceramic ring. A circular LED board is mounted around the outer perimeter of the piezoelectric ceramic, ensuring uniform illumination without interfering with ultrasound wave propagation. Experimental evaluations indicate that the imaging area closely approximates a 15 mm diameter circle, aligning with the dimensions of the ultrasound annular ring. This spatial consistency suggests a high degree of congruence between the optical and ultrasound sensing regions.

### D. Fabrication

The fabrication process follows a systematic eight-step procedure as illustrated in Fig. 4, with parallel assembly paths for the upper and lower sensor components.

For the upper housing assembly (Fig. 4 (a-d)), the process begins by attaching a thin acrylic layer (0.7 mm) to the bottom of the ring-shaped upper housing using cyanoacrylate adhesive 502. This acrylic layer serves dual functions as both an acoustic matching element and a deformation platform for the elastomer. A PDMS elastomer layer with a mass ratio of 30:1 (base to curing agent) is then cast into the housing structure, creating the base elastic layer for tactile sensing. Following this, an HGM-PDMS mixture (volume ratio 1:1) is applied using spin-coating at 3000 rpm for 30 s to form a thin, uniform flexible membrane over the PDMS layer. This specialized membrane functions as both the light-blocking layer for tactile imaging and acoustic matching layer for ultrasound transmission. The upper assembly is then heat-cured at 60 °C for three hours.

For the lower housing assembly (Fig. 4 (e-h)), the camera module is precisely positioned at the center of the lower housing, with the LED illumination ring arranged around it to provide uniform lighting for tactile imaging while maintaining the central optical path. The annular PZT transducer is installed around the camera module, ensuring alignment with the central optical axis to maintain the coaxial design. A tungsten-loaded epoxy backing (volume ratio 3:2) is cast

using a soft-tip syringe to fill the cavity behind the PZT transducer, creating the acoustic backing layer necessary for ultrasound performance. This lower assembly is also heat-cured at 60 °C for three hours. After fabricating and curing both the upper and lower assemblies, they are aligned and combined to form the coaxial sensor module.

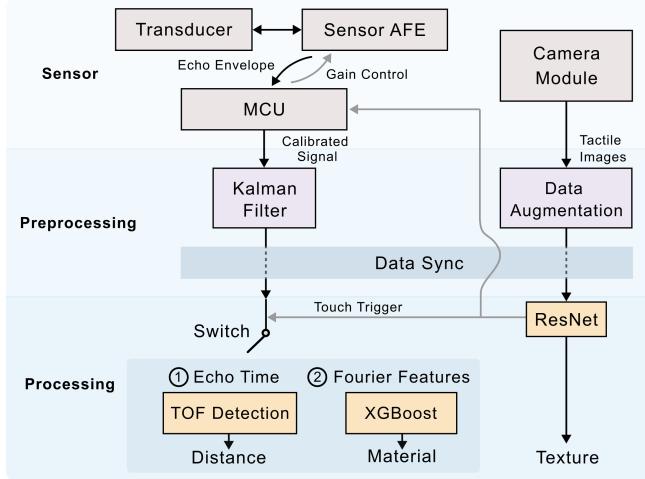


Fig. 5. Overview of the three-level ultrasound–visuotactile fusion pipeline.

### E. System Pipeline

The pipeline shown in Fig. 5 is organized into three hierarchical levels. At the sensor level, the ultrasound path consists of a transducer connected to a sensor AFE, which extracts echo envelopes and communicates with the MCU for gain control. Parallel to this, the camera module captures tactile information.

At the preprocessing stage, ultrasound signals are Kalman filtered for noise reduction, while tactile images are augmented. A critical data sync mechanism aligns 50 Hz ultrasound data with 30 Hz visual data by pairing each ultrasound measurement with its nearest camera frame.

At the processing level, a dual-pathway approach is employed upon touch detection. A touch trigger signal from the camera path activates a switch that routes ultrasound data to either Time-of-Flight (ToF) detection for calculating object distance using echo time during non-contact phases, or XGBoost for material classification using Fourier features during contact. Simultaneously, a ResNet18 model processes tactile images for texture recognition [32].

## IV. EXPERIMENTAL EVALUATION

To validate the sensor’s capabilities, we conducted four systematic experiments: proximity detection, material classification, dual-modal classification, and internal content inspection.

### A. Proximity Detection Experiment

We exploit the propagation and reflection characteristics of ultrasound in air by measuring the ToF of the ultrasound signal in the time-domain. This approach allows the sensor to measure the distance to an object, providing valuable

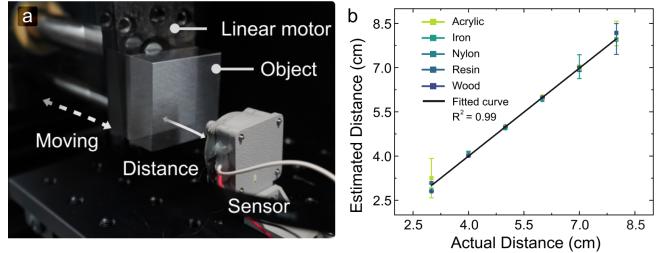


Fig. 6. Proximity detection experiment. (a) Experimental setup. (b) Comparison of estimated vs. actual distances.

information before contact to assist with tasks like collision avoidance and trajectory planning. In our processing method, during sensor initialization, an echo signal (after noise removal) is recorded as a reference frame. Once the distance measurement function is activated, the current frame is subtracted point-by-point from the reference frame. The time corresponding to the voltage peak in the resulting difference frame is identified as the echo time, which is then used to calculate the predicted distance using the method described above.

We evaluate the sensor’s proximity sensing capabilities using the linear motor stage setup shown in Fig. 6(a), where objects of five different materials (acrylic, iron, nylon, resin, and wood) are mounted at distances ranging from 3.0 cm to 8.0 cm from the sensor.

Fig. 6(b) demonstrates the system’s measurement accuracy, with estimated distances strongly correlating with actual distances ( $R^2 = 0.99$ ). The error bars indicate measurement consistency across all five materials, despite their distinct acoustic impedance. The linear relationship between estimated and actual distances remains robust throughout the tested range, with average estimation errors below  $\pm 0.5$  cm. This material-independent proximity detection capability enables reliable distance estimation for robotic manipulation tasks, supporting collision avoidance and approach trajectory planning.

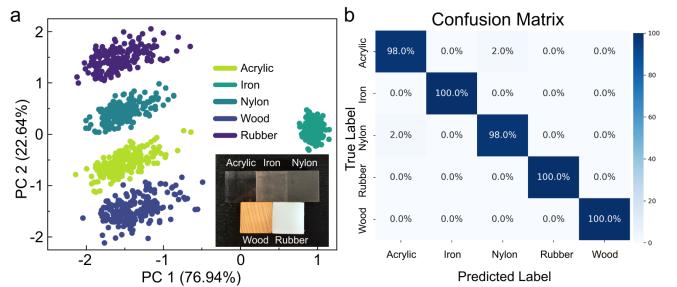


Fig. 7. Material classification experiment. (a) PCA of the spectral features for five materials. The inset shows photographs of the corresponding physical samples. (b) Confusion matrix with 99.20% average accuracy.

### B. Material Classification Experiment

In this experiment, we explore material classification using ultrasound with the aim of capturing additional dimensions of information and overcoming the limitations of traditional

tactile sensors, which are confined to detecting only surface characteristics. To minimize interference from properties unrelated to the acoustic characteristics, experiments are conducted using uniform blocks with consistent size and thickness. For each echo signal frame, a Fourier transform is performed to extract spectral features such as spectral contrast, spectral kurtosis, spectral skewness, and spectral entropy. Fig. 7(a) demonstrates the clear separability achieved through Principal Component Analysis (PCA), with PC1 and PC2 explaining 76.94% and 22.6% of variance respectively across five distinct materials (Acrylic, Iron, Nylon, Rubber, and Wood).

The confusion matrix in Fig. 7(b) confirms high classification accuracy: Iron, Rubber, and Wood achieve perfect classification (100%), while Acrylic (98%) and Nylon (98%) show minimal confusion with each other (2% error rate). These results validate ultrasound's effectiveness for reliable material discrimination, capturing information inaccessible to conventional tactile sensors.

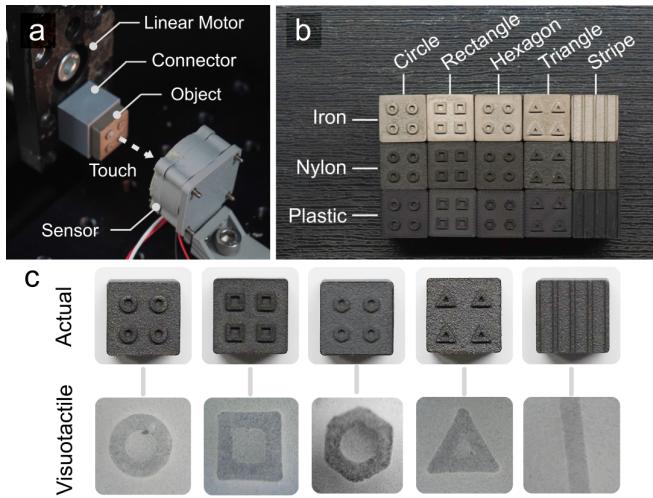


Fig. 8. Dual-modal experiment. (a) Experimental setup. (b) Test objects with patterns across three materials. (c) Actual objects and their corresponding visuotactile results.

### C. Dual-Modal Classification Experiment

The dual-modal assessment combines tactile and ultrasound sensing to enhance object classification by leveraging both surface pattern recognition and material classification. As illustrated in Fig. 8(a), the experimental setup comprises a linear motor that facilitates horizontal movement with a cube-shaped test object attached via a mechanical connector. During contact between the sensor and the object, the system concurrently acquires tactile images and ultrasound echo data. To develop the classification models, we collected 200 samples for each of five materials at intervals of 100 ms, capturing both modalities. The dataset was divided using an 8:2 train-test split to ensure robust model evaluation. For the ultrasound data, each frame was transformed via Fourier analysis to extract features as inputs for XGBoost material classification. For the visuotactile data, image frames were preprocessed with data augmentation operations such as

resizing, rotation, and noise injection before being input to ResNet model for surface pattern classification. Both training processes were executed on an RTX4070 platform.

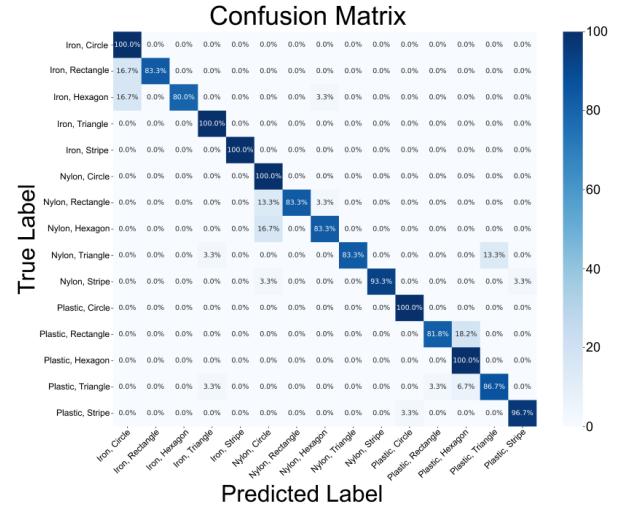


Fig. 9. Confusion matrix for dual-modal experiment with 92.11% average accuracy.

To evaluate the overall system performance, we tested it with various objects as depicted in Fig. 8(b). These objects feature diverse surface patterns (circle, rectangle, hexagon, triangle, stripe) applied to materials including iron, nylons, and plastic. Fig. 8(c) presents a comparison between the actual test objects and the corresponding sensor readings, thereby demonstrating the system's capability to accurately capture and classify surface geometry and material characteristics.

The classification results in Fig. 9 demonstrate the effectiveness of this approach in a 15-class classification task, achieving an average accuracy of 92.11%. The confusion matrix shows strong performance across all materials and shapes, with minor misclassifications occurring primarily between shapes of the same material. This high accuracy confirms that the integration of tactile and ultrasound sensing enables robust discrimination of both geometric features and material properties.

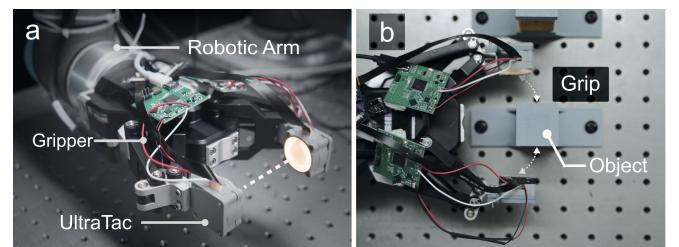


Fig. 10. Internal contents inspection experiment. (a) Robotic gripper setup. (b) Gripper grasping test object.

### D. Internal Content Inspection Experiment

The internal contents inspection experiment demonstrates UltraTac's capability for automated quality control by identifying container surface patterns and internal contents without

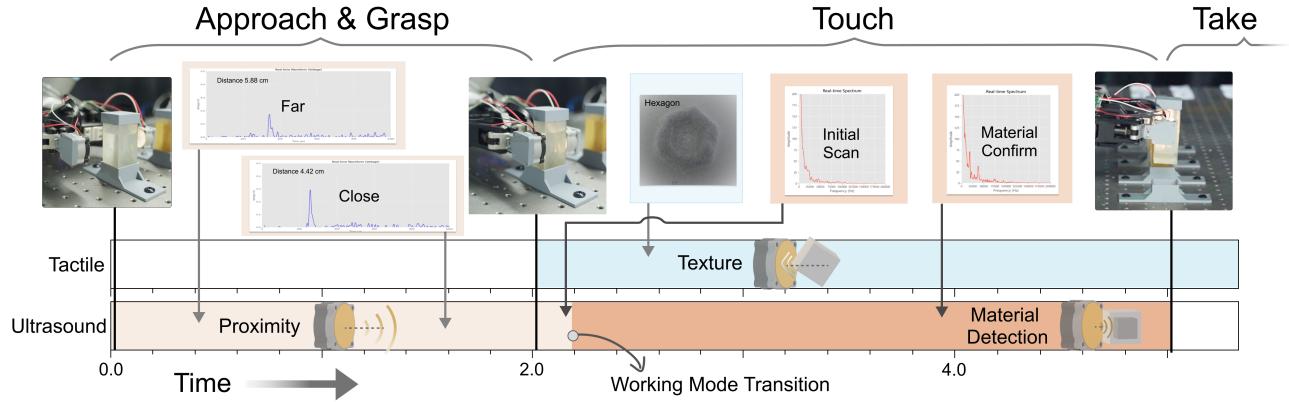


Fig. 11. Timeline of the system operation during internal contents inspection: Approach & Grasp (ultrasound distance measurement), Touch (tactile texture recognition and ultrasound material detection after working mode transition), and Take (transport execution).

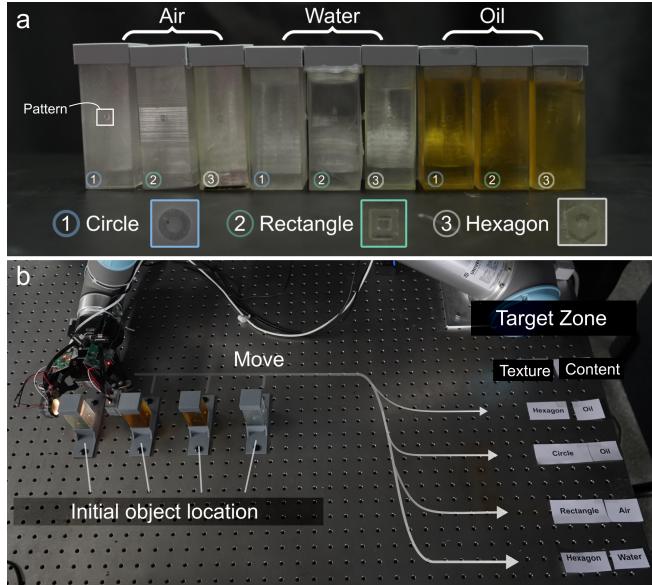


Fig. 12. Test set and workflow. (a) Test containers with varied contents and patterns. (b) Workflow showing transport from initial positions to target zones sorted by texture and content.

opening packages. As shown in Fig. 10(a), two UltraTac sensors are mounted on a robotic gripper attached to a robotic arm for validation. Fig. 10(b) shows the gripper interacting with test objects during grasping, with sensors acquiring both tactile and ultrasound data.

The inspection process follows a structured timeline as illustrated in Fig. 11 and comprises three distinct phases: approach and grasp, touch, and take. In the approach and grasp phase (0-2 s), the ultrasound modality provides proximity sensing to ensure an optimal grasp position by employing ToF calculations to accurately determine the distance. Upon contact at 2 s, the sensor system undergoes a working mode transition. In the subsequent touch phase, the tactile modality captures the surface textures, as exemplified by the hexagonal pattern, while the ultrasound modality switches to a material detection mode to analyze the internal contents by leveraging the spectral characteristics of the echo signals. Finally,

during the take phase, the robotic arm executes the transport operation, delivering the container to its designated location based on both its surface pattern and internal content.

We created a test set of nine containers featuring various internal contents (air, water, oil) and surface patterns (circle, rectangle, hexagon), as illustrated in Fig. 12(a). Each content type includes all three surface patterns, which poses challenges for discrimination: identical patterns with different contents (for example, hexagon-patterned containers containing different substances) and identical contents with different patterns (such as various patterns all containing oil). Fig. 12(b) presents the workflow in which containers are transported from their initial positions to target zones, sorted according to both texture and content type (operation videos and results are listed on the project website). The successful placement of each container underscores the efficacy of the sensor system in concurrently identifying surface features and internal properties, thereby enhancing the overall performance of the robotic application.

## V. LIMITATIONS AND FUTURE WORKS

Ultrasound sensing complements visuotactile perception by adding proximity detection and internal inspection, but it has limits: echoes from targets closer than 3 cm fall into a blind zone due to pulse duration and receiver recovery, and HGM fillers—whose particles are larger than conventional ones—degrade imaging resolution. Future improvements such as increasing excitation voltage, reducing transducer size, or optimizing transducer placement may alleviate the need for strict acoustic matching while improving overall sensor performance.

## VI. CONCLUSIONS

In this paper, we introduced UltraTac, an integrated ultrasound-augmented visuotactile sensor based on a novel coaxial optoacoustic architecture that combines high-resolution tactile imaging with ultrasound proximity sensing. By aligning a micro-camera and a ring-shaped PZT transducer along a common optical axis and employing innovative acoustic matching techniques within a flexible,

optically transparent membrane, UltraTac overcomes the key limitations of conventional visuotactile sensors, which cannot capture material properties or perform proximity sensing. Our dual-pathway signal processing pipeline further enables dynamic switching between proximity detection and material classification modes based on tactile feedback.

Systematic experimental evaluations demonstrate that UltraTac delivers robust performance across multiple sensing tasks. Proximity detection experiments revealed a strong linear correlation ( $R^2 = 0.99$ ) between estimated and actual distances, while material classification based on Fourier-derived spectral features achieved 99.2% accuracy across various materials. Dual-modal object recognition experiments confirmed the sensor's capability with a 92.11% accuracy in a 15-class task, and integration into a robotic gripper showcased its practical utility for automated quality control by successfully identifying both container surface patterns and internal contents without opening packages. Future work will focus on further miniaturization, refinement of signal processing algorithms, and integration of additional sensing modalities to enhance tactile perception in unstructured, real-world environments, paving the way for safer, more efficient, and intelligent human–machine interaction systems.

## REFERENCES

- [1] W. Yuan, S. Dong, and E. H. Adelson, "GelSight: High-resolution robot tactile sensors for estimating geometry and force," *Sensors*, vol. 17, no. 12, p. 2762, 2017.
- [2] M. Lambeta, P.-W. Chou, S. Tian, B. Yang, B. Maloon, V. R. Most, D. Stroud, R. Santos, A. Byagowi, G. Kammerer et al., "DIGIT: A novel design for a low-cost compact high-resolution tactile sensor with application to in-hand manipulation," *IEEE Robotics and Automation Letters*, vol. 5, no. 3, pp. 3838–3845, 2020.
- [3] D. Ma, E. Donlon, S. Dong, and A. Rodriguez, "Dense tactile force estimation using GelSlim and inverse FEM," in 2019 International Conference on Robotics and Automation (ICRA), 2019, pp. 5418–5424.
- [4] Z. Kappassov, J.-A. Corrales, and V. Perdereau, "Tactile sensing in dexterous robot hands—Review," *Robotics and Autonomous Systems*, vol. 74, pp. 195–220, 2015.
- [5] S. Wang, J. Wu, X. Sun, W. Yuan, W. T. Freeman, J. B. Tenenbaum, and E. H. Adelson, "3D shape perception from monocular vision, touch, and shape priors," in 2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2022, pp. 2143–2150.
- [6] R. Li, R. Platt, W. Yuan, A. ten Pas, N. Roscup, M. A. Srinivasan, and E. Adelson, "Localization and manipulation of small parts using GelSight tactile sensing," in 2014 IEEE/RSJ International Conference on Intelligent Robots and Systems, 2014, pp. 3988–3993.
- [7] K. Shimonomura, "Tactile image sensors employing camera: A review," *Sensors*, vol. 19, no. 18, p. 3933, 2019.
- [8] F. R. Hogan, M. Jenkin, S. Rezaei-Shoshtari, Y. Girdhar, D. Meger, and G. Dudek, "Seeing Through your Skin: Recognizing Objects with a Novel Visuotactile Sensor", in 2021 IEEE Winter Conference on Applications of Computer Vision (WACV), Waikoloa, HI, USA: IEEE, Jan. 2021, pp. 1217–1226.
- [9] F. R. Hogan, J.-F. Tremblay, B. H. Baghi, M. Jenkin, K. Siddiqi, and G. Dudek, "Finger-STS: Combined Proximity and Tactile Sensing for Robotic Manipulation", *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 10865–10872, Oct. 2022.
- [10] Z. Song et al., "A flexible triboelectric tactile sensor for simultaneous material and texture recognition", *Nano Energy*, vol. 93, p. 106798, Mar. 2022.
- [11] G. Lee et al., "Fingerpad-Inspired Multimodal Electronic Skin for Material Discrimination and Texture Recognition", *Advanced Science*, vol. 8, no. 9, p. 2002606, May 2021.
- [12] B. Zhu, T. Geng, G. Jiang, Z. Guan, Y. Li, and X. Yun, "Surrounding object material detection and identification method for robots based on ultrasonic echo signals", *Applied Bionics and Biomechanics*, vol. 2023, no. 1, p. 1998218, 2023.
- [13] I.-J. Cho, H.-K. Lee, S.-I. Chang, and E. Yoon, "Compliant ultrasound proximity sensor for the safe operation of human friendly robots integrated with tactile sensing capability", *Journal of Electrical Engineering and Technology*, vol. 12, no. 1, pp. 310–316, 2017.
- [14] Pavlin CJ, Sherar MD, Foster FS. Subsurface ultrasound microscopic imaging of the intact eye. *Ophthalmology*. 1990 Feb;97(2):244-50.
- [15] S. Luo, J. Bimbo, R. Dahiya, and H. Liu, "Robotic tactile perception of object properties: A review," *Mechatronics*, vol. 48, pp. 54–67, 2017.
- [16] R. S. Dahiya, G. Metta, M. Valle, and G. Sandini, "Tactile sensing—From humans to humanoids," *IEEE Transactions on Robotics*, vol. 26, no. 1, pp. 1-20, 2010.
- [17] L. Luo et al., "CompdVision: Combining Near-Field 3D Visual and Tactile Sensing Using a Compact Compound-Eye Imaging System", Mar. 15, 2024, arXiv: arXiv:2312.07146. Accessed: Jul. 01, 2024.
- [18] A. SaLoutos, E. Stanger-Jones, M. Guo, H. Kim, and S. Kim, "Design of a Multimodal Fingertip Sensor for Dynamic Manipulation", in 2023 IEEE International Conference on Robotics and Automation (ICRA), London, United Kingdom: IEEE, May 2023, pp. 8017–8024.
- [19] S. Li et al., "M3Tac: A Multispectral Multimodal Visuotactile Sensor With Beyond-Human Sensory Capabilities", *IEEE Transactions on Robotics*, vol. 40, pp. 4484–4503, 2024.
- [20] D. Xu, G. E. Loeb, and J. A. Fishel, "Tactile identification of objects using Bayesian exploration", in 2013 IEEE international conference on robotics and automation, 2013, pp. 3056–3061.
- [21] W. Yuan, C. Zhu, A. Owens, M. A. Srinivasan, and E. H. Adelson, "Shape-independent hardness estimation using deep learning and a GelSight tactile sensor," in 2017 IEEE International Conference on Robotics and Automation (ICRA), 2017, pp. 951–958.
- [22] Z. Song et al., "SATac: A Thermoluminescence Enabled Tactile Sensor for Concurrent Perception of Temperature, Pressure, and Shear," 2024 IEEE International Conference on Robotics and Automation (ICRA), Yokohama, Japan, 2024, pp. 5680–5686.
- [23] K.-W. Sou, W.-S. Chan, K.-C. Lei, Z. Wang, S. Li, D. Peng, and W. Ding, "A Bio-Inspired Event-Driven Mechanoluminescent Visuotactile Sensor for Intelligent Interactions," *Advanced Functional Materials*, 2024.
- [24] Z. Abderrahmane, G. Ganesh, A. Crosnier, and A. Cherubini, "Haptic zero-shot learning: Recognition of objects never touched before," *Robotics and Autonomous Systems*, vol. 105, pp. 11–25, 2018.
- [25] Q. Mao, Z. Liao, J. Yuan, and R. Zhu, "Multimodal tactile sensing fused with vision for dexterous robotic housekeeping", *Nature Communications*, vol. 15, no. 1, p. 6871, 2024.
- [26] G. Davis, R. Nagarajah, S. Palanisamy, R. A. R. Rashid, P. Rajagopal, and K. Balasubramaniam, "Laser ultrasonic inspection of additive manufactured components", *The International Journal of Advanced Manufacturing Technology*, vol. 102, pp. 2571–2579, 2019.
- [27] E. Kerr, T. M. McGinnity, and S. Coleman, "Material classification based on thermal and surface texture properties evaluated against human performance," in IEEE International Conference on Autonomous Robot Systems and Competitions (ICARSC), 2018, pp. 168–173.
- [28] V. T. Rathod, 'A Review of Acoustic Impedance Matching Techniques for Piezoelectric Sensors and Transducers', *Sensors (Basel)*, vol. 20, no. 14, p. 4051, Jul. 2020.
- [29] L. E. Kinsler, A. R. Frey, A. B. Coppens, and J. V. Sanders, *Fundamentals of acoustics*. John wiley & sons, 2000.
- [30] D.-S. Lin, X. Zhuang, S. H. Wong, M. Kupnik, and B. T. Khuri-Yakub, "Encapsulation of Capacitive Micromachined Ultrasonic Transducers Using Viscoelastic Polymer", *Journal of Microelectromechanical Systems*, vol. 19, no. 6, pp. 1341–1351, Dec. 2010.
- [31] X. Xu, L. Zhang, H. Guo, X. Wang, and L. Kong, 'Acoustic characterization of transmitted and received acoustic properties of air-coupled ultrasonic transducers based on matching layer of organosilicon hollow glass microsphere', *Micromachines*, vol. 14, no. 11, p. 2021, 2023.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, 'Identity mappings in deep residual networks', in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV* 14, 2016, pp. 630–645.