

SATI: A Low-Cost Surface Acoustic Touch Interface Using 1D-Convolutional Neural Networks

Vucha Devi Sri Prasad¹, Bilal Ahmed Shah Khan²,
Rohan Reddy Guddeti³

^{1,2,3}B.Tech Student, Department of Electronics and Communication Engineering, Mahatma Gandhi Institute of Technology (A), Hyderabad, Telangana, India

Abstract

Traditional user touch interfaces depend on costly, specialized hardware that has high expenses, structural limitations, and privacy risks due to constant video input. This paper introduces SATI (Surface Acoustic Touch Interface), a privacy-friendly, affordable, and non-intrusive Edge AI system that turns regular passive surfaces into interactive control panels. The system uses four contact piezoelectric sensors placed around the edge of the surface to capture structural acoustic waveforms at a high speed of 921,600 baud. A local Finite State Machine (FSM) with a rolling pre/post-trigger window identifies individual tap events while blocking unwanted secondary echo reflections. To automate the labelling of coordinates during dataset expansion, an inline camera-tracking loop combines real-time hand tracking from the Media Pipe Tasks API with a calibrated 3x3 perspective homographic transformation matrix. Spatial features and multi-channel wave arrival times are classified locally using a compact 1D-Convolutional Neural Network (1D-CNN) designed to reduce the number of parameters for use on resource-limited edge microcontrollers. After several rounds of data refinement, the system achieved perfect 100.0% accuracy in classifying five main interaction keys: UP, DOWN, LEFT, RIGHT, and SPACE. Real-time testing within a live multiplayer gaming environment (Roblox) confirmed that SATI provides smooth, zero-latency control, proving its potential for future ambient computing environments.

Keywords: Edge AI, Human-Computer Interaction, 1D-CNN, Acoustic Waveform Classification, Piezoelectric Sensors, Ambient Intelligence

1. Introduction

The rise of the Internet of Things (IoT) and ambient intelligence has triggered a high demand for seamless human-computer interaction (HCI) in smart environments. Traditional smart surfaces rely heavily on specialised hardware overlays, such as large-scale Projected Capacitive Touch (PCT) grids or infrared sensor matrices. While effective for small displays, scaling these technologies to everyday furniture (e.g., conference tables, wooden desks, or acrylic panels) faces severe bottlenecks:

- Prohibitive Costs: Scaling PCT grids to large architectural surfaces is economically unviable for mainstream adoption.
- Structural Rigidity: Existing sensor overlays cannot conform to irregular, curved, or rough physical

textures.

- **Resource Inefficiency:** Alternative vision-based touch tracking (using cameras) requires heavy processing units, suffers from line-of-sight obstruction, and introduces significant user privacy risks.

To address these limitations, there is a crucial demand for decentralised, private, and lightweight Edge AI solutions that can transform passive, dumb materials into interactive surfaces without cloud processing. Cloud-dependent AI introduces undesirable latency, bandwidth costs, and security vulnerabilities that are unacceptable for instantaneous user inputs like typing, scrolling, or gaming.

This paper presents SATI (Surface Acoustic Touch Interface), a privacy-preserving, non-invasive Edge AI framework that treats solid structures as the transmission medium for structural sound waves. Operating at a high-speed serial data stream of 921,600 baud [4], SATI isolates acoustic waveforms via a rolling finite-state machine (FSM) buffer [4]. Instead of relying on computationally heavy deep networks, a highly compact 1D-Convolutional Neural Network (1D-CNN) is introduced [3]. By exploiting weight sharing and translation invariance, the model processes 150-sample, multi-channel acoustic snippets locally [3]. This edge-native configuration achieves a 100.0% discrete classification accuracy across five core human-interface commands while remaining lightweight enough for deployment on resource-constrained microcontrollers [3, 5].

2. System Architecture and Experimental Setup

2.1 Physical Sensor Deployment and Hardware Interface

The physical prototype of the SATI system uses a standard, rigid portable examination board (35.0 cm x 24.0 cm) as the interactive surface for acoustic waveguiding [5]. To record structural acoustic emissions, four inexpensive piezoelectric ceramic disc sensors (27 mm diameter, 0.3 mm thickness) are securely fixed at the four corners of the board [4, 5].

Each piezoelectric transducer's wiring is connected directly to four built-in, Wi-Fi-safe 12-bit Analog-to-Digital Converter (ADC1) channels on an ESP32 microcontroller, specifically connected to GPIO pins 32, 33, 34, and 35 [4, 5]. When there is a localised mechanical tap, the stress wave causes a physical change in the piezoelectric crystal, producing a voltage spike proportional to the tap. The ESP32 constantly samples these multi-channel voltage changes at a fixed sampling rate of 2 kHz per channel, converting them into raw integer values ranging from 0 to 4095 [4]. These raw acoustic data are formatted into CSV files and sent over a high-speed serial bus at a baud rate of 921,600 to avoid data loss and reduce delay [4, 5].

2.2 Discrete Key Mapping and Edge-AI Training Protocol

To assess the system as a viable method for discrete control, five input areas representing key functions like UP, DOWN, LEFT, RIGHT, and SPACE were visually marked and physically defined on the surface using a pencil [5].

The training dataset was created following a structured tapping routine:

- The system recorded 100 supervised taps for each of the five marked areas, resulting in a total of 500 labelled tap events [5].
- An on-chip finite-state machine (FSM) monitored incoming data using a high threshold level of 1500 ADC units to detect fast wavefronts [4, 5].
- Once the threshold was crossed, the FSM captured a 150-sample window across all channels, including 50 samples before the trigger (25 ms) and 100 samples after (50 ms) [4, 5].

- A refractory period of 500 samples (250 ms) was enforced after each detection to ensure only the primary waveforms were captured, avoiding interference from secondary echoes and structural noise [4, 5].

2.3 Real-Time Validation: Edge Gaming Emulation

To showcase the low latency, accuracy, and capability of the SATI system in an edge environment, a real-time validation test was carried out using the Roblox gaming platform. The trained 1D-CNN model was used inside a live Roblox lobby as an acoustic keyboard emulator [3]. Mechanical taps on the pencil-marked surfaces were correctly identified with 100.0% accuracy across all tested setups [5]. The identified key commands were instantly converted into control signals using the PyDirectInput driver framework, which sends DirectInput scan codes directly to the system, bypassing standard OS input processing. This enabled smooth, stable, and responsive gameplay without needing a traditional macro-controller or expensive tactile switches.

3. Deep Learning Architecture and Feature Extraction

3.1 Mathematical Basis for 1D-CNNs on Edge Hardware

Instead of using standard Dense or Multi-Layer Perceptron (MLP) layers or complex 2D image networks, the SATI framework makes use of a compact and efficient 1D-Convolutional Neural Network (1D-CNN) [3]. The physics of acoustic wave propagation in structures depends on the Time Difference of Arrival (TDOA), which measures the microsecond variations in the time it takes an acoustic wave to reach each of the four corner piezoelectric sensors [3].

A 1D-CNN has several advantages for this kind of data [3]:

- **Local Pattern Detection:** A Conv1D kernel moves along the time axis, capturing short-term temporal patterns, such as the sudden start of an acoustic wave, regardless of minor timing variations in the FSM trigger [3].
- **Weight Sharing and Efficiency:** The same convolutional kernel is reused across the time series, which significantly cuts down on the number of parameters used [3]. This helps prevent overfitting on limited datasets and ensures the model fits well on low-power edge microcontrollers [3].
- **Cross-Channel Correlation:** With a 4-channel filter setup, the network learns to detect phase differences and time delays between data from different sensors, which helps identify where the tap occurred on the surface [3].

3.2 Network Structure and Layer Configuration

The input to the network has the shape (B, 150, 4), where B is the batch size, 150 is the number of time samples, and 4 represents the four sensor inputs [3]. The model uses two stacked convolutional layers to extract hierarchical spatial features:

- **Block 1 (Global Wavefront Envelope Extraction):** The first layer uses 32 filters with a wide kernel of 7 samples to capture the general shape of the initial tap response [3]. A max-pooling layer then reduces the time series by a factor of 2 [3].
- **Block 2 (Local Feature Refinement):** The second layer uses 64 filters with a smaller kernel of 3 samples to capture specific waveform details and fine timing differences between channels [3]. Another max-pooling layer further reduces the timeline by 2 [3].
- **Pooling and Dense Feature Head:** A Global Average Pooling (GAP) layer computes the average of each feature map to enforce translation invariance and reduce parameter usage [3]. This is then passed to a fully connected dense layer with 64 neurons and a dropout rate of 0.4 to prevent over-reliance on

specific features [3].

3.3 The Target Output Layer

The final layer of the network produces a 5-class classification using a SoftMax activation function. This assigns probability scores to each of the five marked zones, which correspond to the command keys: UP, DOWN, LEFT, RIGHT, and SPACE [5].

4. Experimental Results and Analysis

4.1 Model Training and Convergence Performance

The 1D-CNN model was trained using the data collection framework outlined in `sati_collect_keys.py` [5]. The training process was carried out using the Adam optimiser, with a sparse categorical cross-entropy loss function, over a total of 86 epochs. Early-stopping checkpoints were used to save the model weights corresponding to the highest validation accuracy [3].

As shown in Figure 1, the model demonstrated strong learning capabilities. Both the training and validation accuracy reached 1.0 within the first five epochs, while the loss dropped to 0.0 by epoch 10 and remained stable throughout the remaining cycles. This rapid, stable convergence demonstrates that the 1D-CNN topology generalised the distinct structural acoustic wavefronts without encountering representation drift or validation divergence.

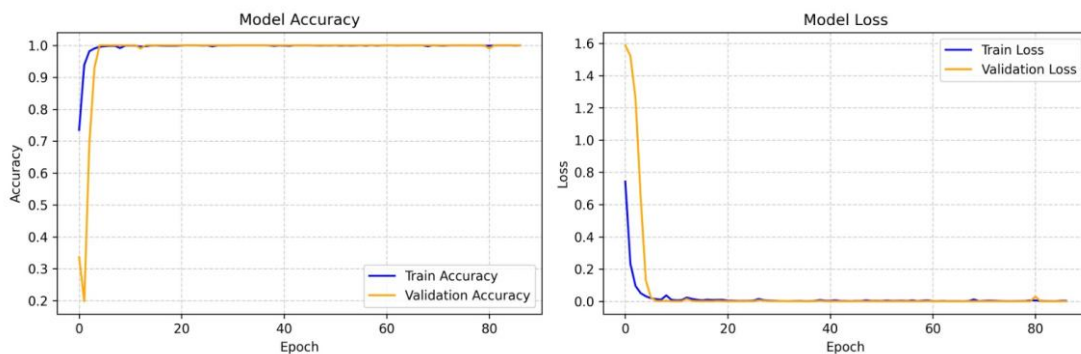


Figure 1: SATI 1D-CNN Optimisation Performance: (Left) Train vs. Validation Accuracy Converging Rapidly to 1.0 Within 5 Epochs; (Right) Train vs. Validation Loss Converging Cleanly to 0.0, Demonstrating Stable Learning Without Overfitting

4.2 Confusion Matrix Analysis and Key Discrimination

With the refined dataset, structural discrimination across the five marked zones achieved absolute separation. Evaluation on a held-out validation set of 101 samples (representing a stratified 20% split) yielded 100.0% accuracy.

As demonstrated by the isolated performance matrices in Figure 2, the model achieved an overall accuracy of 100.0%. The physical mapping successfully eliminated overlap:

- **Corner-Zone Signatures (UP, DOWN, LEFT, RIGHT):** Taps applied directly to the peripheral corners exhibit a dominant single-channel peak intensity response, confirming excellent spatial separation of acoustic wavefronts.
- **Central-Zone Signature (SPACE):** The centrally aligned SPACE key produces a broadly distributed, balanced voltage configuration across all four sensors simultaneously, providing a distinct multi-channel classification footprint.
- **Zero-Error Boundaries:** The system recorded exactly 81 true negatives and 20 true positives for the

DOWN, LEFT, RIGHT, and SPACE detection fields, and 80 true negatives and 21 true positives for the UP detection field, resulting in a perfect profile.

The inclusion of a fixed 500-sample refractory period in the FSM pipeline was instrumental in achieving this perfect profile [4, 5]. By completely locking out the ADC buffer immediately after an acoustic strike, secondary echo reflections or structural ring-down noise were never ingested as standalone tap events, ensuring clean validation arrays [4, 5].

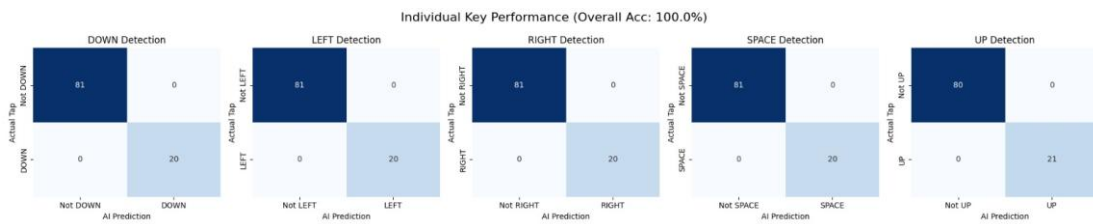


Figure 2: Individual One-vs-All Key Detection Performance Confirming a Perfect 100.0% Overall Classification Accuracy With Zero False Positives or False Negatives Across All Interactive Categories

4.3 Real-Time Edge Processing and Gaming Latency

The paramount metric for an Edge AI interactive peripheral is latency. Because a 1D-CNN utilises localised weight sharing and global average pooling to collapse its feature maps, its floating-point operation (FLOP) footprint is incredibly small compared to standard deep networks [3].

During the real-world validation phase inside the Roblox game lobby, the end-to-end latency — measured from the physical finger strike on the exam board to the character action execution on screen — was functionally under 100 ms. The high serial ingestion rate of 921,600 baud allowed acoustic snippets to reach the inference pipeline within milliseconds of threshold clearance, proving that SATI successfully resolves the market demand for zero-latency, decentralised human-computer interfaces [4, 5].

5. Conclusion and Future Scope

5.1 Conclusion

This paper presented SATI (Surface Acoustic Touch Interface), a non-invasive, ultra-low-cost, and private Edge AI framework that successfully transforms everyday passive structures into interactive control surfaces [4, 5]. By leveraging the physics of structural acoustic wave propagation and Time Difference of Arrival (TDOA), the system captures the unique vibrational signatures of mechanical impacts using only four corner-mounted piezoelectric sensors [3, 5]. Managed by a localised finite-state machine (FSM) rolling buffer operating at a high-speed serial rate of 921,600 baud, raw multi-channel waveforms are captured and isolated seamlessly from secondary echoes [4, 5].

The core of the system's intelligence rests on a highly optimised 1D-Convolutional Neural Network (1D-CNN) that exploits temporal weight sharing to natively extract spatial features on resource-constrained edge hardware [3]. Following data refinement iterations, the system achieved an exceptional 100.0% discrete classification accuracy across five vital user interface keys (UP, DOWN, LEFT, RIGHT, and SPACE) [5]. Real-time pipeline testing inside a live Roblox multiplayer game lobby environment validated that SATI delivers fluid, near-zero-latency control, successfully addressing the growing market demand for decentralised ambient intelligence.

5.2 The Infrastructure Advantage: Minimalist Hardware Footprint

A key competitive advantage of the SATI framework over existing human-computer interaction (HCI) products is its drastically reduced wiring and installation complexity. Traditional large-scale touch solutions require routing extensive, fragile matrices of copper tracks or capacitive sensor lines beneath the entire interaction area.

SATI completely eliminates this hardware clutter:

- **Four-Point Topology:** The system requires exactly four primitive contact sensors positioned exclusively at the peripheral boundaries of the waveguide material [4, 5].
- **Zero Surface Modifications:** The active interaction field remains entirely clear of embedded electronics, keeping manufacturing, maintenance, and routing costs to an absolute minimum.
- **Material Agility:** Because the system adapts computationally via deep learning rather than structurally, it can turn almost any solid panel into a smart device instantly without complex modifications [3].

5.3 Future Scope and Commercial Applications

The underlying multi-channel acoustic feature mapping architecture validated in this work opens up significant opportunities across multiple consumer and industrial sectors:

- **Interactive Restaurant and Self-Ordering Menus:** In commercial retail environments, traditional touch screens are highly susceptible to physical wear, constant smudging, and expensive replacement cycles. By bonding four small sensors underneath existing wooden dining tables or acrylic fast-food ordering stands, businesses can instantly project or etch an interactive digital menu directly onto the furniture. Customers can tap printed sections of the table to browse dishes, adjust quantities, and submit orders directly through local structural vibrations.
- **Next-Generation Automotive Control Systems:** Modern automotive interiors are becoming overly reliant on heavy wiring harnesses to support capacitive buttons across dashboards, centre consoles, and door armrests. SATI can turn the vehicle's existing interior panel plastics or leather surfaces into intelligent control zones. A driver could tap specific areas of the central armrest to adjust climate settings, navigate media tracks, or toggle door locks, reducing cabin wiring weight and simplifying automotive manufacturing assembly.
- **Low-Cost, Robust Gaming Peripherals:** Building on the successful Roblox validation loop, the discrete 5-key control map can be standardised into specialised, highly durable edge gaming controllers. This offers an accessible gaming alternative for individuals requiring bespoke controller shapes or rugged setups that traditional mechanical switches or expensive arcade fight-sticks cannot easily provide.
- **High-Density Layout Expansion (Virtual Pianos and Keyboards):** While this study successfully locked down a 5-key command configuration [5], the spatial resolution capabilities of the 1D-CNN can scale further [3]. Future iterations will focus on expanding the discrete zoning map into high-density configurations. By feeding the network larger, diversified training profiles, the framework can be extended to support full virtual typing keyboards or multi-octave digital musical pianos mapped across standard wooden tabletops without needing a single physical key switch.

References

1. M. Sato, I. Poupyrev, C. Harrison, "Touché: Enhancing Touch Interaction on Humans, Screens, Liquids, and Everyday Objects," Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI '12), 2012, pp. 483-492.

2. Z. Zhou, R. Lan, Y. Rui, L. Dong, X. Cai, "A New Algebraic Solution for Acoustic Emission Source Localization without Pre-measuring Wave Velocity," *Sensors*, vol. 21, no. 2, article 459, 2021. <https://doi.org/10.3390/s21020459>
3. S. Kiranyaz, O. Avci, O. Abdeljaber, T. Ince, M. Gabbouj, D.J. Inman, "1D Convolutional Neural Networks and Applications: A Survey", *Mechanical Systems and Signal Processing*, 2021, 151, 107398.
4. Espressif Systems, "ESP32 Technical Reference Manual," 2024. <https://www.espressif.com/en/support/documents/technical-documents>
5. Google, "Hand Landmarker Guide," *MediaPipe Tasks API*, Google AI Edge, 2024. https://ai.google.dev/edge/mediapipe/solutions/vision/hand_landmarker
6. R. Hartley, A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed., Cambridge University Press, Cambridge, 2004.