# Webcam-Based Hand-Tracked Musical Controller

Purnima Vasistha, Hriday Kharpude, Ana Babnigg
CS 445 Final Project
{purnima4, hridayk2, babnigg2}@illinois.edu

## ABSTRACT

A real-time, webcam-based musical controller is presented that converts simple hand gestures into continuous sound controls. Utilizing MediaPipe Hands and OpenCV, the system estimates about 21 hand landmarks per frame, extracts robust features, and maps them to a lightweight Python audio engine. The index fingertip governs pitch ($x$) and volume ($y$); a thumb↔index pinch is an on/off controller for the sound; a closed fist holds the current pitch; and thumb↔middle and thumb↔ring pinches shift the octave up and down, respectively. Stability within the program is achieved through exponential smoothing and size-normalized distances as well as hysteresis and debouncing. An on-screen HUD displays a current frequency, volume, pinch/hold status, octave, and FPS. The resulting interface is playable and consistent with the stated project objectives.

## I. INTRODUCTION

(Hand)-gestural instruments promise intuitive, low-latency musical control without any needed custom hardware that would not be easily accessible. Inspired by Imogen Heap's Mi.Mu gloves, we transform a single hand into a reliable controller using commodity CV. Our objective is to convert something easily accessible into musically meaningful signals using simple, explainable methods.

## II. MOTIVATION AND IMPACT

This topic was selected to explore a creative, real-time application of computer vision that remains usable outside the class. The initial motivation was to draw on the idea of music controlling, explicitly Imogen Heap's Mi.Mu gloves, and apply that to the broader ideas of Computational Photography. The broader motivation is as follows: (i) demonstrating that robust, musical control can be achieved with only a webcam and open-source libraries, lowering the barrier to entry for gestural instruments; and (ii) providing a compact example of turning noisy landmarks into stable control signals using interpretable techniques (normalization, hysteresis, debouncing, smoothing). The resulting interface supports education, prototyping, and accessibility in human–computer interaction and digital music.

## III. RELATED WORK

Prior work explores gesture-to-sound mapping and interactive audio effects [1], [2]. Modern real-time hand tracking via MediaPipe enables robust landmark estimation with minimal compute.

## IV. APPROACH

### A. Pipeline

1) Capture a webcam frame and mirror it (`cv2.VideoCapture`, `cv2.flip`) so left/right feel natural.
2) Convert to RGB and run MediaPipe Hands (`hands.process`) to obtain 21 hand landmarks if a hand is visible.
3) Extract features: index fingertip $(x, y)$ (landmark 8) and the thumb–index distance (landmarks 4 and 8).
4) Determine gesture states:
   - **Hold (fist/open):** compute a hand openness score (normalized by palm width) and update the state with hysteresis.
   - **Octave shifts:** detect thumb–middle (up) and thumb–ring (down) pinches using size-normalized distances, hysteresis, and a 300 ms debounce; update an integer $\Delta \in [-2, 2]$.
5) Map features to sound: `map_controls` sets base frequency and volume from $x$ and $y$, and uses the thumb–index pinch as a gate.
6) Apply octave: multiply frequency by $2^\Delta$ and clamp it to [50, 2000] Hz.
7) Apply Hold: if the hand just closed to a fist, latch the current frequency; keep it until the hand opens.
8) Smooth frequency and volume with an exponential moving average (EMA, $\alpha = 0.2$).
9) Send `audio_frequency`, `audio_volume`, `audio_pinch` to the audio engine under a lock and draw the HUD.
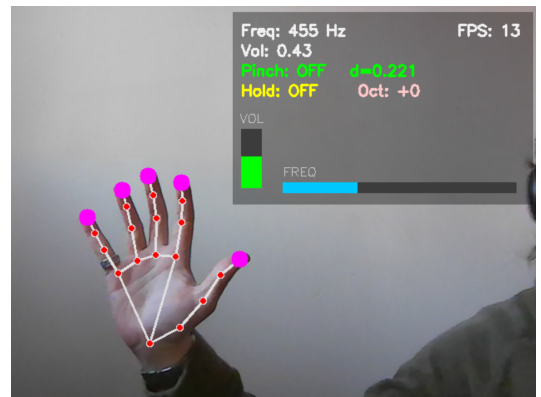


Fig. 1. Shows hand tracking, basic HUD interface with an open hand.

## B. Features and Mapping

Coordinates are normalized to the image: $x \in [0,1]$ (left→right) and $y \in [0,1]$ (top→bottom). The view is mirrored.

- **Pitch:** $f = 200 + 800\,x$ Hz  (linear map from left to right).
- **Volume:** $v = 0.8\,(1-y)$  (higher on screen = louder).
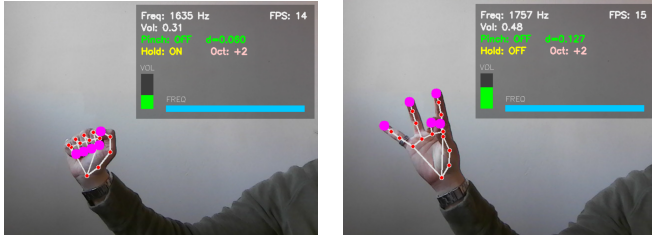- **Gate:** thumb–index pinch is ON if the (unnormalized) distance in image space is $< 0.05$.

An octave offset $\Delta$ from discrete gestures is applied as $f \leftarrow \mathrm{clamp}_{[50,2000]}(f \cdot 2^{\Delta})$.

## C. Gestures

**Smoothing.** An EMA with $\alpha = 0.2$ reduces frame-to-frame jitter while keeping intentional motion responsive.

**Hold (fist).** A hand openness score is the mean distance from the (index, middle, ring, pinky) tips to a palm center, divided by palm width (index MCP↔pinky MCP). Hysteresis thresholds (OPEN_ON=1.05, OPEN_OFF=0.95) avoid flicker: crossing below OPEN_OFF enters FIST (frequency is latched); crossing above OPEN_ON returns to OPEN (live control resumes).

**Octave shift.** Thumb–middle and thumb–ring distances are normalized by palm width and compared to PINCH_ON=0.24 and PINCH_OFF=0.30 with a 300 ms debounce. Each thumb–middle pinch increments $\Delta$ by $+1$ (up to $+2$); each thumb–ring pinch decrements $\Delta$ by $-1$ (down to $-2$).



figureHUD states: Hold ON (left) and Octave +2 (right).

## D. Audio Engine (sound production)

A sounddevice output stream (44.1 kHz, blocksize 512) runs a sine oscillator. In the audio callback, if the gate is OFF the buffer is filled with zeros; otherwise a sinusoid is rendered at audio_frequency with amplitude audio_volume. Shared parameters are updated from the vision thread under a lock to avoid race conditions.

## E. HUD (user interface)

A semi-transparent panel (top-right) displays frequency, volume, pinch state (and distance), hold state, current octave offset $\Delta$, and FPS. A vertical bar shows volume; a horizontal bar shows frequency.



Fig. 2. The HUD shows Pinch: ON and the normalized distance $d$.

## F. Evaluation

Metrics (latency, steadiness, reliability) were captured using an instrumented build of the application with logging (the released code remains unchanged). Here, p90 denotes the 90th percentile and $\Delta_{\mathrm{pp}}$ denotes peak-to-peak range.

- **Latency (CV→audio):** median **15.8 ms** (p90 30.2 ms, $n = 24$). Measures the time from when the vision thread sets the pinch gate to when the audio callback consumes it; well below the 0.2 s target, indicating responsive control.
- **Steadiness (10 s at rest):** frequency $\sigma = \mathbf{10.93}$ Hz, $\Delta_{\mathrm{pp}} = 32.31$ Hz; volume $\sigma = \mathbf{0.002}$, $\Delta_{\mathrm{pp}} = 0.010$. Smaller values indicate steadier control; EMA reduces typical webcam landmark jitter.
- **Reliability (pinch): 4/5** trials successful (preliminary). A success is a clean toggle within a 1 s trial window with a hold $\geq 0.3$ s. As planned, this should be expanded to 30 trials per lighting condition to report percentages.

## V. IMPLEMENTATION DETAILS

Python 3.11; mediapipe 0.10.14; opencv-python 4.10.0.84; numpy 1.26.4; sounddevice. Files: handstuff.py (tracking, mapping, gestures, smoothing, HUD) and audio_engine.py (OutputStream, sine oscillator, shared params). Run python handstuff.py; press q to quit.

**External resources.** MediaPipe Hands (Google), OpenCV, Python SoundDevice (docs linked in References). No external datasets are required; a laptop webcam suffices.

## VI. CHALLENGES AND INNOVATION

**Challenges.** Achieving stable gestures without ML: thresholds must be distance-invariant and non-fluttering. This was addressed via palm-width normalization, hysteresis, and debounce. Audio responsiveness also required a low-latency callback and careful smoothing to avoid sluggishness.

**Innovation.** A compact, fully interpretable pipeline that converts landmarks to controls with low latency; discrete octave control via multi-finger pinches; and a practical "hold" mechanism using openness with hysteresis. The design is extensible for additional aspects such as MIDI output, scale quantization, left-hand XY pad). Innovation is present with the implementation of a moderately complex, real-time interaction system combining multiple robust CV heuristics with an audio callback path and UI; novel integration and tuning rather than reproducing a single static algorithm.

## VII. CONCLUSION

A responsive, single-hand musical controller was demonstrated using only a webcam and common libraries. With simple feature engineering and temporal filtering, the system achieves steady, musically meaningful control (pitch, volume, gate) plus discrete gestures (hold, octave shift).

## ACKNOWLEDGMENTS

## REFERENCES

[1] V. Verfaille, U. Zölzer, and D. Arfib, "Mapping strategies for gestural and adaptive control of digital audio effects," 2006.

[2] S. Matsoukas, A. Zacharakis, and D. Kamarotos, "Towards a real-time mapping of Finger Gesture to Sound," 2012.

[3] Google MediaPipe Hands, https://developers.google.com/mediapipe/solutions/vision/hand_landmarker

[4] OpenCV, https://opencv.org/

[5] Python SoundDevice, https://python-sounddevice.readthedocs.io/

[6] G. D. Nagalapuram, R. S., V. D., D. D., and D. J. Nazareth, "Controlling Media Player with Hand Gestures using Convolutional Neural Network," in *Proc. IEEE Mysore Sub Section International Conference (MysuruCon)*, Hassan, India, 2021, pp. 79–86. doi: 10.1109/MysuruCon52639.2021.9641567.