

Sound2Hap: Learning Audio-to-Vibrotactile Haptic Generation from Human Ratings

Yinan Li

School of Computing and Augmented Intelligence
Arizona State University
Tempe, AZ, USA
yinanli2@asu.edu

Hasti Seifi

School of Computing and Augmented Intelligence
Arizona State University
Tempe, AZ, USA
hasti.seifi@asu.edu

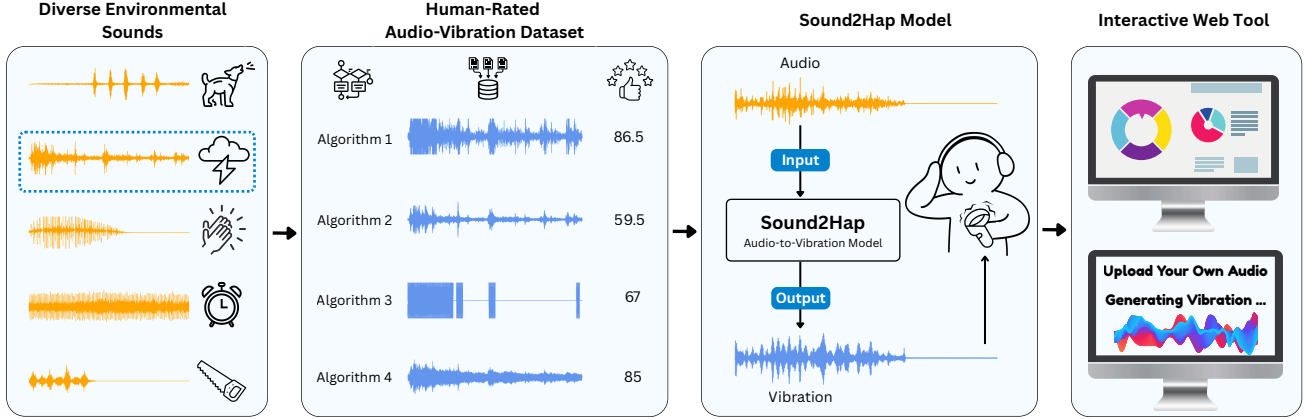


Figure 1: Overview of Sound2Hap. Each diverse environmental sound is first converted into vibrations using four existing signal-processing algorithms. We collect human ratings for each audio-vibration pair, building a dataset of 4,000 rated samples. Using the dataset, we then train the *Sound2Hap* model to generate vibrations directly from sound effects and evaluate it in a user study. In addition, we provide an online tool for visualizing the dataset and generating vibrations using both baseline signal-processing algorithms and our model, enabling designers and researchers to explore or create haptic experiences from environmental sounds.

Abstract

Environmental sounds like footsteps, keyboard typing, or dog barking carry rich information and emotional context, making them valuable for designing haptics in user applications. Existing audio-to-vibration methods, however, rely on signal-processing rules tuned for music or games and often fail to generalize across diverse sounds. To address this, we first investigated user perception of four existing audio-to-haptic algorithms, then created a data-driven model for environmental sounds. In Study 1, 34 participants rated vibrations generated by the four algorithms for 1,000 sounds, revealing no consistent algorithm preferences. Using this dataset, we trained Sound2Hap, a CNN-based autoencoder, to generate perceptually meaningful vibrations from diverse sounds with low latency. In Study 2, 15 participants rated its output higher than signal-processing baselines on both audio-vibration match and Haptic Experience Index (HXI), finding it more harmonious with diverse sounds. This work demonstrates a perceptually validated approach to audio-haptic translation, broadening the reach of sound-driven haptics.

Keywords

Sound-haptic conversion, Audio-haptic conversion, Automatic generation, Multimodal dataset, Haptic Design, Vibrotactile perception

1 Introduction

Haptic feedback about environmental events can enhance realism, immersion, and accessibility in virtual reality (VR), multimedia, and gaming applications [24, 40, 47, 48]. Everyday sounds from nature, animals, and objects carry rich information and emotions, making them a valuable resource for designing haptic effects. For example, in a game or VR experience, vibrations could simulate a passing train, a barking dog, or a crackling fire to enhance immersion and support accessibility for users with visual or hearing impairments [42, 77]. One way to achieve this is through audio-to-vibration mapping, which leverages the shared time-domain structure of both modalities. However, achieving high audio-tactile congruence is challenging due to the differences in the sensory bandwidths of the auditory and tactile systems. While hearing spans 20-20,000 Hz, touch is typically limited to frequencies below 1,000 Hz and has coarser frequency resolution [21, 30]. These disparities make it difficult to translate fine-grained auditory features into perceptually meaningful haptics.

Prior work has introduced various signal-processing techniques to convert audio into vibrotactile patterns, typically by mapping frequency ranges [57, 65] or perceptual features such as loudness, roughness, or pitch to vibration [46, 49]. These approaches have

shown promise in targeted contexts such as enriching music listening or generating strong, discrete effects in games such as explosions or gunfire. However, their reliance on predefined mappings limits scalability and constrains haptic feedback to a narrow set of sensations. As a result, they fall short when applied to more diverse and nuanced audio content. In particular, little is known about how well such methods can translate everyday environmental sounds to vibrations, including their subtle informational and emotional cues.

To explore this, we conducted a study to gather user ratings on the perceptual match between a large, diverse set of audio-vibration pairs (Figure 1). We implemented four existing audio-to-vibration algorithms from the literature [46, 49, 65, 80] and converted 1,000 sound clips from the ESC-50 environmental sound dataset [69] into vibrations. These clips spanned 50 classes, including animal sounds, natural soundscapes, human non-speech sounds, and both domestic and urban environmental noises. In the study, 34 participants held a voice-coil vibration actuator (Haptuator Redesign) between their thumb and index finger to feel vibrations and rated 4,000 audio-vibration pairs for perceptual preference and alignment. Analysis revealed that user ratings varied widely across sound types, showing that fixed algorithmic mappings have limited generalizability. We also observed relationships between sound characteristics and algorithm performance. For example, some methods performed better for impulse or rhythmic sounds, while others were more suitable for broadband or continuous sounds, indicating that each conversion method captures different sound characteristics. These insights motivated the need for a more adaptive, data-driven approach.

Building on these insights, we developed Sound2Hap, a convolutional neural network-based autoencoder model to synthesize vibration signals directly from audio inputs. Using our human-rated dataset from the four existing algorithms, we trained two versions of the model: one using the best-rated vibration for each sound clip (Top-Pair Sound2Hap) and another using all four vibration pairs per clip, weighted by user preferences (Preference-Weighted Sound2Hap). In a user study with 15 participants, we evaluated these models on new, unseen sound clips against a Baseline that applied the best signal-processing method for each sound class. Results showed that both model versions outperformed the Baseline, achieving significantly higher signal ratings and improved scores for Harmony, Discord, and General Score on the Haptic Experience Index (HXI) questionnaire. Between the two model variants, Top-Pair Sound2Hap received more qualitative preference from participants. These findings demonstrate the model’s ability to generalize across diverse environmental sounds while aligning with user preferences. Finally, we provide an interactive web tool that allows researchers and designers to browse our audio-vibration dataset through dynamic visualizations, examine patterns in user ratings across different sounds, and generate new haptic signals for their own audio clips.

Our contributions are as follows:

- A user-rated multimodal dataset with 4,000 paired audio-vibration signals across 50 classes of environmental sounds and 8,000 user ratings. To our knowledge, this is the largest audio-vibration dataset with user ratings.
- Sound2Hap, an audio-to-vibration generative model that converts diverse environmental sound effects into perceptually aligned vibrotactile feedback, trained using two alternative schemes and evaluated for generation time as well as in a user study against signal-processing baselines.
- An interactive web tool that allows designers to explore the user-rated dataset and convert new sound files into vibrations using Sound2Hap and prior signal-processing algorithms.

We make the dataset, audio-to-vibration model, and web tool open-source to support future research in this area. The model is available at GitHub¹. The audio-vibration dataset can be downloaded at Hugging Face². The web tool is available online³.

2 Related Work

The perceptual affinity between audio and vibrations has led to the development of various audio-to-haptic methods in both research and commercial settings. These techniques have been applied to enhance immersion in gaming [8, 43, 89] and movies [9, 86] by synchronizing tactile effects with sound, as well as to improve accessibility by providing vibrotactile substitutes for auditory cues [16, 25, 31, 74]. We categorize these approaches into (a) signal-processing techniques and (b) learning-based neural network models.

2.1 Signal-Processing Methods

Previous sound-to-tactile conversion methods rely primarily on signal-processing algorithms to convert audio signals into haptic feedback. These approaches are computationally efficient with minimal delay, making them well-suited for real-time applications such as music-related haptic experiences. Since the auditory system perceives a wide frequency range (20-20,000 Hz) while tactile sensitivity is limited (< 1 kHz), these systems typically apply filtering and spectral transformations to compress auditory features into the vibrotactile bandwidth.

Some methods map low-level features of the audio signal directly onto vibration waveforms. For systems employing a single vibration actuator, simple spectral transformation techniques have proven effective. Frequency shifting has been used to render high-frequency musical components as haptic effects [17, 65], while dual-band linear resonant actuators have been introduced to better represent treble frequencies [38, 39]. Sung et al. [80] convert audio waveforms to vibrations by centering the vibration frequency at 200 Hz, near the point of highest human tactile sensitivity [30], and introduce variation by mapping low and high intensities to a range of ± 50 Hz around this center frequency. When multiple actuators or actuator arrays are available, richer mappings can be achieved by distributing sound energy across different frequency bands. For example, Karam et al. [44, 45] demonstrate this approach using a bank of bandpass filters to map multiple frequency bands to vibrotactile stimuli across various body locations. However, such a simple bandpass filter cannot effectively capture high-frequency sound features and often fails to selectively emphasize the most informative components of complex audio clips.

¹<https://github.com/Iris1215/Sound2Hap>

²<https://huggingface.co/datasets/yinanli1215/Sound2Hap>

³<https://audiovibration.netlify.app/>

Lee and Choi [49] propose a mapping from the psychoacoustic features of loudness and roughness to the perceived haptic intensity and roughness, enabling selective haptic conversion without machine learning. The mapping controls the type of sound for haptic conversion based on the sound’s perceptual quality. For example, a mapping that emphasizes only perceptually loud and rough sound can render gunshots while ignoring smoother background music, making the conversion selective. Building on this approach, Li et al. [51] introduce additional psychoacoustic features including sharpness, booming, and low-frequency energy, combined with simple threshold detection to improve conversion selectivity. However, these perceptual mapping approaches face classification accuracy challenges in determining when sounds should be converted to vibration, with Lee and Choi [49] reporting up to 20% false positive rates and potential issues with missing low-energy target sounds or generating unnecessary tactile effects for loud, rough human voices.

Signal-processing-based methods are also widely adopted in industry and consumer products. For example, gaming headsets that emphasize low-pitch sounds convert bass frequencies into vibrations [22], and gaming controllers employ sound-driven tactile feedback to enhance immersion [70]. Authoring platforms such as Meta Haptics Studio [58] and bHaptics Studio [6] allow users to convert audio files into vibrotactile signals and fine-tune them through graphical interfaces. Hardware products such as the bHaptics vest further leverage multi-band audio energy to drive a large number of vibration motors distributed across the torso, enabling full-body audio-to-haptic feedback [7].

Most signal-processing methods have distinct advantages and limitations, and are developed primarily for gaming effects or music applications. To assess their efficacy for converting environmental sounds, we implement four distinct methods from the literature and construct a user-rated audio-vibration dataset.

2.2 Learning-Based Methods

Researchers have also explored machine learning (ML) approaches for audio-to-haptic conversion. Yoshida et al. [88] propose VibVid, an ML-based system that learns from audio, video, and acceleration data to automatically generate vibration waveforms, demonstrated on tennis videos. Zhan et al. [93] present a generative cross-modal framework using a Residual U-Net to synthesize tactile signals directly from audio during tool-surface interactions (e.g., sliding or tapping). These works demonstrate the feasibility of generating vibration signals with machine learning, but both are constrained to specific domains, tennis and tool-surface interactions, and therefore do not generalize to the wide variability of everyday environmental sounds that our work addresses. Furthermore, these methods rely on directly collecting paired audio and vibration signals using microphones and accelerometers [78]. However, capturing high-quality vibrations is often infeasible for many environmental sounds. While the growing number of large-scale audio datasets [10, 29, 69] offers a valuable resource for haptic design, these audio datasets lack paired acceleration data, which limits the applicability of prior learning-based approaches for these datasets.

Others use ML models to identify a set of pre-defined sound events from an audio stream, then render preset vibration effects

or employ signal-processing techniques to convert the detected sounds. For instance, Yun et al. [90] employ deep neural networks (DNNs) to detect gunfire sounds in video games and render synchronized motion effects on a motion platform, thereby enhancing 4D gameplay experiences. In a follow-up work, they extend this approach with a random forest classifier to identify suitable sounds for haptic rendering and to selectively produce multimodal tactile stimuli including both vibrations and impacts [91]. More recently, Yun and Choi [89] present a semantic sound-to-haptic conversion system for VR gameplay, where a Long Short-Term Memory (LSTM) model classifies in-game sounds into four key events (gunfire, hits, explosions, reloads) that trigger upper-body vibration patterns on a haptic vest. While these works highlight the potential of machine learning for sound-to-haptic conversion, most focus on classifying pre-defined sound events, which limits the generalizability of the resulting models.

Finally, recent work has developed learning-based approaches to generate haptic signals from other modalities. Generative adversarial networks have been used to synthesize texture vibrations from images or material attributes [5, 11–14, 50, 84, 85]. Heravi et al. [35] propose a model that generates texture signals in real time from user actions such as force and speed, while Faruqi et al. [27] employ a variational autoencoder to design physical textures for 3D-printed objects. More recently, researchers have leveraged large language models (LLMs) to generate haptic signals from free-form inputs, such as generating vibrations from text prompts [52, 60, 72, 80], midair haptic patterns (e.g., rotating lines) from text [79], or thermal feedback from videos [61]. While these approaches facilitate haptic signal design, they do not generate haptics synchronized and aligned with the audio modality.

Our approach learns a mapping from audio waveforms to vibration signals that are perceptually aligned with human experiences, enabling vibration generation from diverse environmental sounds.

3 Overview of Sound2Hap Design and Evaluation Process

Figure 2 illustrates our process for developing an audio-to-vibration model for environmental sounds:

Step 1 - Implement four signal-processing algorithms (Section 4): We selected ESC-50 [69] as a diverse labeled audio dataset for environmental sounds, then implemented four existing audio-to-vibration algorithms from literature to convert 1,000 ESC-50 sound clips into vibration signals, forming the initial audio-vibration dataset.

Step 2 - Collect human ratings (Section 5): In a study with 34 participants, we collected 8,000 ratings on 4,000 audio-vibration pairs, revealing no consistent preference across algorithms and motivating a data-driven approach.

Step 3 - Train the generative model (Section 6): Using the human-rated dataset, we developed Sound2Hap and trained model variants (Top-Pair and Preference-Weighted) to generate vibrations from diverse audio input, aligned with user preferences.

Step 4 - Final user evaluation (Section 7): A second in-person study with 15 participants showed that both Sound2Hap variants outperformed the best signal-processing baselines on perceptual ratings and HXI measures.

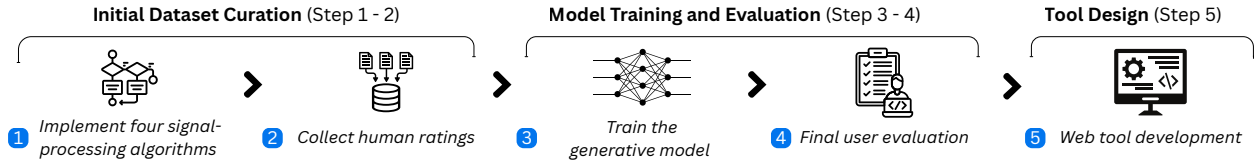


Figure 2: Our overall process for designing and evaluating Sound2Hap.

Step 5 - Web tool development (Section 8): We created an interactive online tool that lets designers explore the dataset and generate vibrations from new sounds using both Sound2Hap and baseline algorithms.

4 Creating An Audio-Vibration Dataset

To build the dataset, we implement four audio-to-vibration conversion algorithms, each with distinct characteristics from the literature, and apply them to a dataset of environmental sound clips.

4.1 Audio Dataset

As our audio source, we use the ESC-50 dataset [69], which contains 2,000 five-second environmental sound clips across five broad categories and 50 detailed semantic classes. The five sound categories include: (1) animal sounds, (2) natural soundscapes and water sounds, (3) human non-speech sounds, (4) interior domestic sounds, and (5) exterior urban environmental noises. We select this dataset due to its broad coverage of real-world environmental sounds, spanning both salient, common events (e.g., dogs barking, laughter) and more subtle, nuanced ones (e.g., brushing teeth, glass breaking). The clips emphasize foreground sound events while minimizing background noise and have clear event labeling, making the dataset well-suited for studying perceptual correspondence between audio and vibration. All sound clips are single-channel (mono), sampled at 44.1 kHz, and stored in 16-bit PCM format.

To ensure consistency across all methods, we apply a normalization procedure to the sound clips. Prior to applying any audio-to-vibration conversion, all input audio signals are peak-normalized with 0 dB headroom and no additional clamping. This step ensures that waveforms are scaled to a consistent amplitude range while avoiding clipping.

4.2 Four Signal-Processing Methods for Audio-to-Vibration Conversion

To create vibrations, we implement four signal-processing approaches based on prior literature and assign the following descriptive labels for ease of reference: Perception-Level Mapping [49], Frequency Shifting [65], Pitch Matching [46], and HapticGen [80]. We select these four methods because they represent distinct signal-processing strategies previously used to convert audio to vibration: directly mapping perceptual features (e.g., roughness), shifting frequency components, matching sounds to a single dominant frequency (i.e., pitch), or mapping intensity and temporal dynamics. Together, they provide a diverse set of baselines for evaluating user preferences. Below, we describe each approach and note any

modifications introduced in our implementations. To ensure comparability across methods, we apply a standardized normalization procedure during vibration post-processing.

4.2.1 Perception-Level Mapping. This method replicates the perceptual framework developed by Lee and Choi [49], which maps the perceptual audio features of loudness and roughness to corresponding vibration intensity and roughness. The algorithm operates on 4096-sample frames to extract audio loudness (La) based on ISO equal-loudness contours⁴ and roughness (Ra) using Vassilakis’s model for spectral peak interactions [87]. These features are then mapped to vibrotactile intensity (Iv) and roughness (Rv) using the equations proposed by the authors. The final vibration synthesis employs two sinusoidal components at fixed frequencies of 175 Hz and 210 Hz, chosen through psychophysical experiments to maximize perceptual control over intensity and roughness, with amplitudes derived from Iv and Rv . Lee and Choi proposed two sets of parameters to convert “game” sounds and “music”. We implement this method using equations for the “game” category, as it is a more relevant category for environmental sounds.

4.2.2 Frequency Shifting. This algorithm is based on the work of Okazaki et al. [65], who propose frequency shifting to compress the broad spectral range of audio signals, especially music, into the limited tactile frequency range (0-1,000 Hz). They sum the original audio signal with one-octave (-12 semitones) and two-octave (-24 semitones) down-shifted versions, followed by a band-pass filter centered at 250 Hz. Their experiments demonstrate high perceived harmony, comfort, and enjoyment, particularly for music containing high-frequency components such as music box sounds. Their implementation uses proprietary software (Hayaemon⁵) for pitch manipulation.

We adapt this approach by developing an open-source Python pipeline with several modifications. We generate one-octave and two-octave down-shifted versions using librosa’s pitch-shifting function⁶, then combine them with the original signal using waveform addition. To address the low-frequency artifacts more effectively, we first apply a 10 Hz high-pass filter to remove frequencies that cause muffled sensations. Then, we apply a 250 Hz Butterworth band-pass filter ($Q = 1.0$, order 4) for the final spectral shaping.

4.2.3 Pitch Matching. This algorithm is based on the work of Kim et al. [46], who propose a crossmodal pitch-matching function to convert short sound events (< 1 second) into vibrotactile signals. Their approach follows two steps. First, they calculate specific loudness values across 24 Bark frequency bands, using ISO 532-1, and

⁴<https://www.iso.org/standard/34222.html>

⁵<http://en.edolfzoku.com/hayaemon2/>

⁶<https://github.com/librosa/librosa>

use these features in a regression model to predict a single vibration frequency that best matches the sound. Second, the sound’s loudness profile over time modulates the amplitude of a sinusoidal vibration, making the vibration mirror the sound’s dynamic intensity. The algorithm is developed using 25 short sounds of glass breaks, gunshots, swords, hitting, and explosions, but its efficacy was not tested with users.

To extend this method to longer and more complex environmental recordings, we implement a time-varying adaptation of the original algorithm. We segment input audio into overlapping 10 ms windows with 50% overlap, following common vibration-processing parameters [53, 80]. For each segment, we then compute both specific loudness (for frequency prediction) and auditory loudness (for amplitude modulation) using ISO 532-1. We apply the regression coefficients from Kim et al.’s perceptual study to each window’s Bark-band features to derive instantaneous vibration frequencies. These time-varying frequency and amplitude parameters are then interpolated across the whole audio duration and synthesized into a continuous sinusoidal waveform using phase accumulation to

ensure smooth frequency transitions. This dynamic approach preserves the core crossmodal pitch matching while extending it to handle longer, frequency-varying sounds.

4.2.4 HapticGen. Based on the importance of temporal rhythm in vibrotactile perception, HapticGen [80] prioritizes intensity and temporal characteristics over spectral frequency details. The algorithm is designed to handle noisy sound clips with overlapping sound effects and complex frequency spectra, where intensity-based mapping provides a more robust haptic signal than spectral analysis. However, its perceptual effectiveness is not evaluated in a study. The vibration signal is synthesized using a sinusoidal numerically controlled oscillator with a base frequency of 200 Hz, selected for its proximity to peak human vibrotactile sensitivity [30]. Audio signals are segmented into 10 ms analysis windows, and short-term Root Mean Square (RMS) energy is computed to drive dynamic frequency modulation within a ± 50 Hz range around the center frequency. This approach enables the system to capture intensity

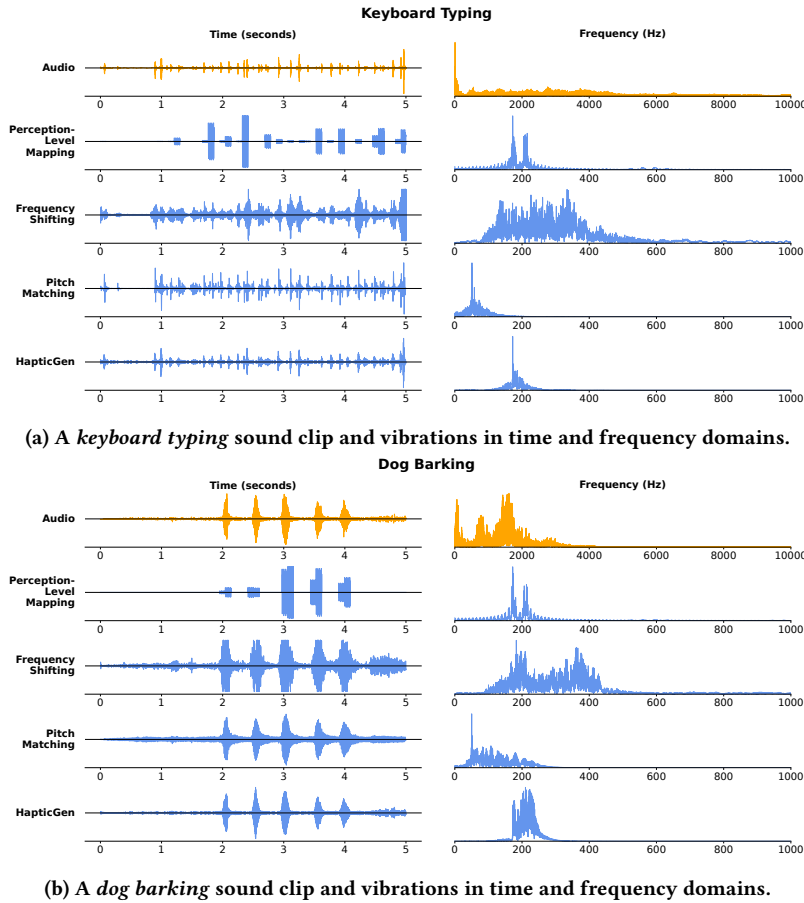


Figure 3: Examples of original audio signals (orange) and converted vibrations (blue), using four signal-processing algorithms. The audio was recorded at 44.1 kHz, 16-bit resolution, and vibrations were rendered at 8 kHz, 16-bit resolution. The y-axis of time-domain plots (left) ranges from -1 to 1. The y-axis of frequency-domain plots (right) shows magnitude, with taller spikes indicating greater energy.

variations while maintaining a frequency range optimized for tactile perception. We use the open-source code⁷ shared by the authors for this method.

4.2.5 Normalization Process for Vibration Signals. After vibration waveform synthesis, all methods apply Root Mean Square (RMS) normalization to ensure consistent perceptual loudness across outputs. For HapticGen, Pitch Matching, and Perception-Level Mapping, the input audio is analyzed in short segments (e.g., 10-100 ms) to generate vibration signals. We estimate segment-wise RMS values and use the maximum value across segments to compute a global normalization factor for the audio. Frequency Shifting generates vibrations from the full audio clip in a single pass, after which we apply a global normalization factor based on the waveform’s overall RMS value.

All final vibration signals are single-channel (mono), saved at an 8 kHz sampling rate using 16-bit PCM encoding. Figure 3 presents the waveforms of two example sound clips alongside their converted vibration signals.

5 User Study 1: Collecting User Ratings for Audio-to-Vibration Techniques

We conducted an in-person user study with 34 participants to collect human preference ratings for vibration signals generated using the above four signal-processing algorithms. Participants rated 4,000 vibrations created from 1,000 sound clips in the ESC-50 dataset, resulting in 8,000 individual ratings.

5.1 Study Methods

5.1.1 Audio-Vibration Stimuli. We selected 20 sound clips from each of 50 audio classes in the ESC-50 dataset (half of the dataset). To ensure acoustic diversity within each class, we extracted comprehensive feature vectors from every clip, including spectral features (centroid, rolloff, and bandwidth), energy metrics (RMS energy, zero-crossing rate, and estimated tempo in beats per minute), 13 Mel-frequency cepstral coefficients (MFCCs), and a 12-dimensional chroma vector. All features were temporally averaged across each sound clip to create static descriptors. For each sound class, we applied K-means clustering ($k=10$) on the feature vectors to divide the acoustic space into clusters. We then sampled 20 clips per class proportionally from these clusters to capture diverse acoustic characteristics. When proportional sampling returned fewer than 20 clips, we randomly selected from the remaining clips to reach the target count. Each sound clip was paired with the four vibrations generated using the four signal-processing algorithms in Section 4.2.

5.1.2 Participants. We recruited 34 participants (24 male, 10 female; mean age = 24 years, $SD = 2.4$) through online advertisements at the authors’ institution. All participants had normal or corrected-to-normal vision and hearing, and no neuropathy or skin injuries on their hands. Four participants reported haptics expertise: all of them with vibrotactile technology, three with force-feedback devices, two with mid-air haptics, and one with electrotactile stimulation. The other 30 participants reported minimal or no prior experience in haptics. Each participant received a \$15 cash reward for their time.

⁷<https://github.com/HapticGen/HapticGen>

The study protocol was approved by the Institutional Review Board (IRB) of the authors’ institution.

5.1.3 Procedure. After collecting informed consent and a background questionnaire, the experimenter introduced the study interface and instructed participants to wear noise-canceling headphones. Participants were asked to hold a voice-coil vibration actuator, namely Haptuator Redesign⁸ from TactileLabs, between the thumb and index finger of their left hand. The actuator was driven by a TPA3116D2 Class D stereo digital audio amplifier board connected to a 12V power supply, with its input linked directly to the PC for waveform playback. We calibrated audio volume, vibration intensity, and amplifier gain during a pilot study and kept them constant across all 34 participants.

Figure 4 shows the study setup and interface. In each trial, participants first listened to a five-second sound clip, followed by feeling four corresponding vibration signals presented in randomized order. Participants could replay both the audio and each vibration signal as many times as needed. After experiencing each vibration, they rated how well it replicated the perceptual effect of the sound on a scale from 0 (no match) to 100 (perfect match). Participants could only proceed to the next clip after playing the sound and four vibrations at least once and providing a rating for each vibration.

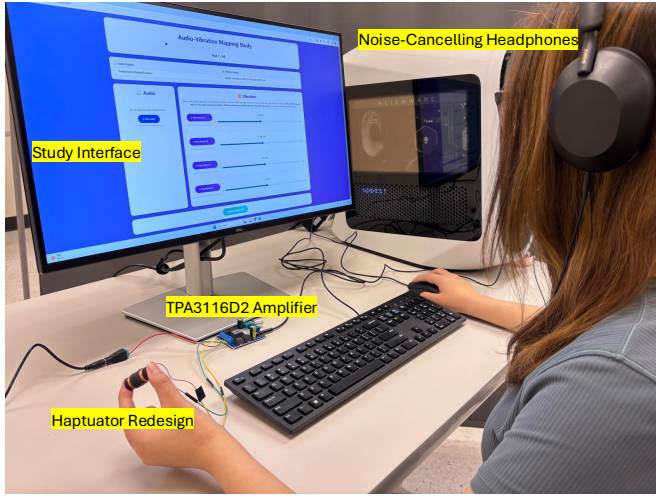
To balance clip distribution across audio classes, maintain participant engagement, and reduce fatigue, we limited each session to 60 trials. The 1,000 sound clips were divided into 17 sets: 16 sets with 60 clips each and one final set with 40 clips. Each of the first 16 sets contained one clip from every class (50 in total), plus 10 additional clips randomly selected from the remaining pool. The final set consisted of the remaining 40 clips. This design ensured that each clip received exactly two ratings, and that participants experienced a wide range of sound classes, helping maintain engagement and reduce fatigue from repeatedly rating similar clips. Each session lasted approximately one hour per participant.

5.2 Results: Insights from User Ratings

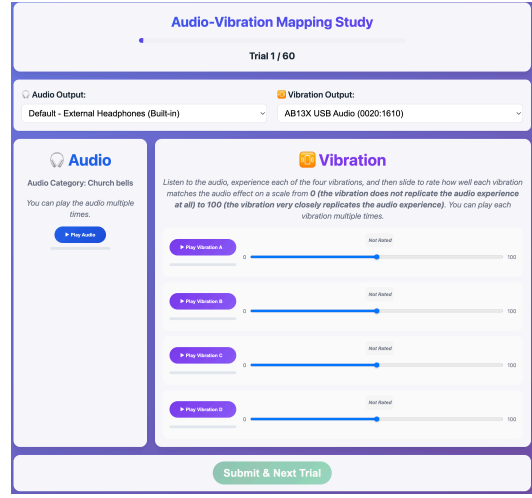
User preferences toward different algorithms varied across sound clips, classes, and categories. Below, we report preferences at the levels of five major categories, 50 sound classes, and 1,000 individual sound clips. For each individual sound clip, the rating is the average between two participants.

Five Categories. Figure 5a shows the distribution of clip-level ratings aggregated within each category. Each box, therefore, summarizes 200 sound-vibration pairs (10 classes \times 20 clips). Ratings for each sound category and algorithm span a wide range, suggesting that each algorithm performs well on some sound clips but poorly on others. For *Natural Soundscapes*, *Human Non-Speech*, and *Interior/Domestic* categories, Pitch Matching received the highest average ratings. In the *Animals* category, Frequency Shifting, Pitch Matching, and HapticGen performed similarly, with no clear winner. For the *Exterior/Urban* category, Frequency Shifting and Pitch Matching also showed comparable ratings. We further averaged user ratings across all sound clips for each algorithm. Overall, Pitch Matching achieved the highest average rating ($Mean = 62.6$, $SD = 22.9$ out of 100), with HapticGen ($Mean = 57.0$, $SD = 23.2$) and Frequency

⁸<https://tactilelabs.com/product/tl-002-14r-haptuator-redesign/>



(a) Study setup.



(b) Study interface.

Figure 4: Study setup and interface for the User Study 1. Participants held a voice-coil vibration actuator (Haptuator Redesign) and rated each vibration’s match to the sound clips on the interface.

Shifting ($Mean = 56.9$, $SD = 24.3$) performing comparably, often close behind. In comparison, Perception-Level Mapping received the lowest overall ratings ($Mean = 31.2$, $SD = 22.9$). Overall, all methods received low or moderate average ratings for the diverse environmental sounds in our study.

Sound Classes. Analyzing trends across the 50 sound classes revealed more nuanced variations (Figure 5b). Table 1 summarizes the winning (highest average rating across its 20 clips) algorithm for each sound class. The results suggest some links between a sound’s acoustic properties and the most effective algorithm. Frequency Shifting proved more effective for continuous, noise-based, or broadband sounds with steady energy or mechanical characteristics, such as rainfall, mechanical drones (e.g., washing machine, vacuum cleaner, airplane), or engines, where conveying the overall frequency texture is more important than preserving fine temporal detail. Meanwhile, HapticGen often excelled with ambient or rhythmic soundscapes, such as sea waves, wind, crickets, and chirping birds, and sounds featuring pronounced echoing rise-and-fall envelopes like church bells and car horns. Finally, Pitch Matching dominated across several sound classes, especially for sequences of short, discrete events such as keyboard typing, door knocks, footsteps, and glass breaking. Both HapticGen and Pitch Matching performed well on animal sounds.

Individual Sound Clips. Clip-level results revealed further insights about the algorithms. Pitch Matching was preferred as the top method for 403 clips, Frequency Shifting for 288 clips, HapticGen for 261 clips, and Perception-Level Mapping for 56 clips. For 8 sound clips, two methods tied as the highest-rated. These results show that user preferences varied depending on the specific sound characteristics. While Perception-Level Mapping had the lowest category-level and class-level average ratings, it still performed better than others for a subset of the clips. Overall, no single existing algorithm received consistently high ratings across all sounds,

underscoring the need for generative models that can capture these nuances and further adapt to the diverse characteristics of audio signals. To this end, we constructed a dataset of 1,000 sound clips, 4,000 vibration signals, and 8,000 human ratings from this study to serve as training data for our model development.

6 Sound2Hap: Audio-to-Vibration Autoencoder Model

We developed Sound2Hap to learn an audio-to-vibration mapping directly from user perception ratings in our dataset of diverse sound effects. This section presents our dataset preparation and augmentation, model architecture and training, and evaluation of Sound2Hap’s vibration generation time.

6.1 Dataset Preparation and Augmentation

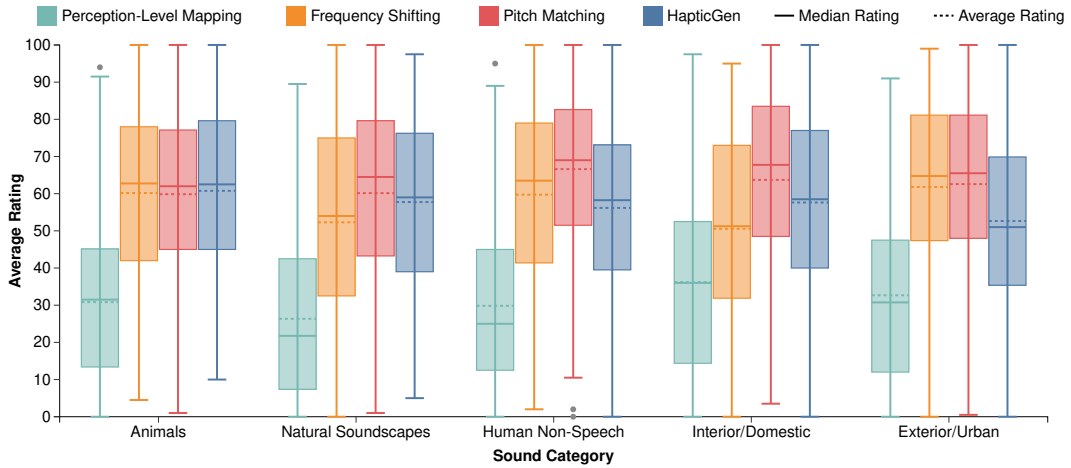
Our dataset contains 1000 sound samples, each paired with four corresponding vibration variants and their associated human preference ratings. We partitioned this dataset into an 80% training set and a 20% validation set. To ensure temporal alignment between audio and vibration signals, we downsampled all audio files from 44.1 kHz to 24 kHz, and upsampled the vibration signals from 8 kHz to 24 kHz using TorchAudio’s FFT-based resampling⁹. During pre-processing, we peak-normalized input audio and normalized all target vibration waveforms to a fixed root mean square (RMS) value. This dual-normalization strategy decoupled the signal’s content from its raw amplitude, producing volume-agnostic audio input and intensity-consistent vibration targets. By standardizing both domains, we preserved the statistical structure of the signals while preventing the concentration of data within a specific range. Such pre-processing was shown to improve cross-modal generation, boosting the adaptability of the model and the quality of learned signal representations [93].

⁹<https://docs.pytorch.org/audio/stable/index.html>

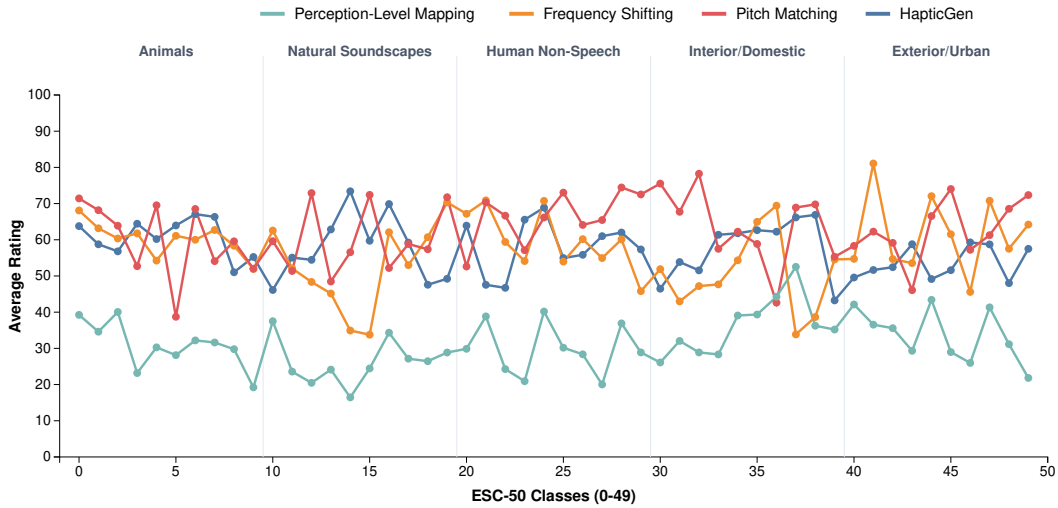
To improve the model’s robustness and generalization, we applied common data augmentation techniques prior to training. For each sound sample, we randomly applied pitch shifting of up to ± 2 semitones [1, 75] and added a small amount of Gaussian noise (up to 0.5% of the signal’s amplitude), each with a 50% probability [59]. This strategy enables the model to learn the core semantic features of the audio, making it robust to small variations in pitch and recording conditions and preventing overfitting. Augmentation was applied only to the training data, while the validation set was kept unchanged to track validation loss reliably, adjust the learning rate, and save the best model for evaluation.

6.2 Model Architecture and Training

We developed an autoencoder-based generative model for audio-to-vibration translation and trained it with two distinct strategies and loss functions (Figure 6). The architecture required an audio encoder to compress complex audio waveforms into a meaningful, low-dimensional latent representation. To achieve this, we adopted the pre-trained encoder and quantizer from the EnCodec model [26], which has demonstrated strong performance in compressing high-fidelity audio in a compact latent space. Our objective was to design a model that not only preserves the temporal and spectral characteristics of source audio but also incorporates subjective human preferences to generate higher-quality haptic feedback. Hence, we



(a) Boxplot of user ratings across the five sound categories for each algorithm. Each box represents 200 data points, one for each sound-vibration pair in the category.



(b) Average user ratings over 50 sound classes per algorithm.

Figure 5: Performance of the four audio-to-vibration algorithms across five overall categories and 50 sound classes.

Algorithm	Animals	Natural Soundscapes	Human Non-Speech	Interior/Domestic	Exterior/Urban
Perception-Level Mapping	-	-	-	-	-
Frequency Shifting	-	10 rain 18 toilet flush	20 crying baby 21 sneezing 24 coughing	35 washing machine 36 vacuum cleaner	41 chainsaw 44 engine 47 airplane
Pitch Matching	0 dog 1 rooster 2 pig 4 frog 6 hen 8 sheep	12 crackling fire 15 water drops 19 thunderstorm	22 clapping 25 footsteps 26 laughing 27 brushing teeth 28 snoring 29 drinking/sipping	30 door knock 31 mouse click 32 keyboard typing 34 can opening 37 clock alarm 38 clock tick 39 glass breaking	40 helicopter 42 siren 45 train 48 fireworks 49 hand saw
HapticGen	3 cow 5 cat 7 insects flying 9 crow	11 sea waves 13 crickets 14 chirping birds 16 wind 17 pouring water	23 breathing	33 door wood creaks	43 car horn 46 church bells

Table 1: Winning algorithm (rows) for each of the 50 sound classes (cells) within the five categories (columns). The number preceding each class name indicates its class number and corresponds to the x-axis in Figure 5b.

focused our design efforts on the decoder and the training objectives specific to the haptic domain.

The model follows a three-stage process, with the encoder and quantizer components frozen during training to function as a fixed audio feature extractor. First, an encoder transforms 24 kHz audio signals into latent representations. We use a CNN encoder with four convolutional blocks (with residual units and downsampling) followed by two Long Short-Term Memory (LSTM) layers. Second, a quantizer applies residual vector quantization to discretize these latent representations into efficient codes. Third, a custom decoder translates these codes into single-channel vibration signals at 24 kHz. We use the SEANet architecture for the decoder, which has a series of residual upsampling blocks and two LSTM layers to capture the long-range temporal dependencies in vibrotactile signals.

In our implementation, we replaced the standard ReLU activations with LeakyReLU to prevent vanishing gradients and improve training stability. Finally, we removed the typical tanh output activation to allow the model to learn natural vibration intensities directly rather than constraining the output amplitudes to a normalized range. We trained two variants of Sound2Hap under distinct loss functions and training strategies, as shown in Figure 6.

6.2.1 Top-Pair Sound2Hap. The first variant, the *Top-Pair Sound2Hap*, is trained using a direct, best-target strategy (Figure 6a). Its objective is to generate an output that precisely matches the single vibration that received the highest human preference rating for a given sound clip.

The training is guided by a composite loss function that compares the generated vibration against the best-rated vibration using four key metrics: (1) Mean squared error (MSE) loss supports basic time-domain waveform reconstruction; (2) Multi-resolution STFT loss compares spectrograms at different resolutions (FFT sizes of 1024, 512, 256); (3) Mel-spectrogram L1 loss emphasizes perceptually relevant frequency bands; (4) Amplitude loss adjusts the

overall intensity of the prediction with the target vibration. See Supplementary Material for the loss function.

6.2.2 Preference-Weighted Sound2Hap. The second variant, the *Preference-Weighted Sound2Hap*, learns from the full spectrum of user ratings rather than a single top choice (Figure 6b). This is achieved through a two-phase training strategy that incorporates a Generative Adversarial Network (GAN).

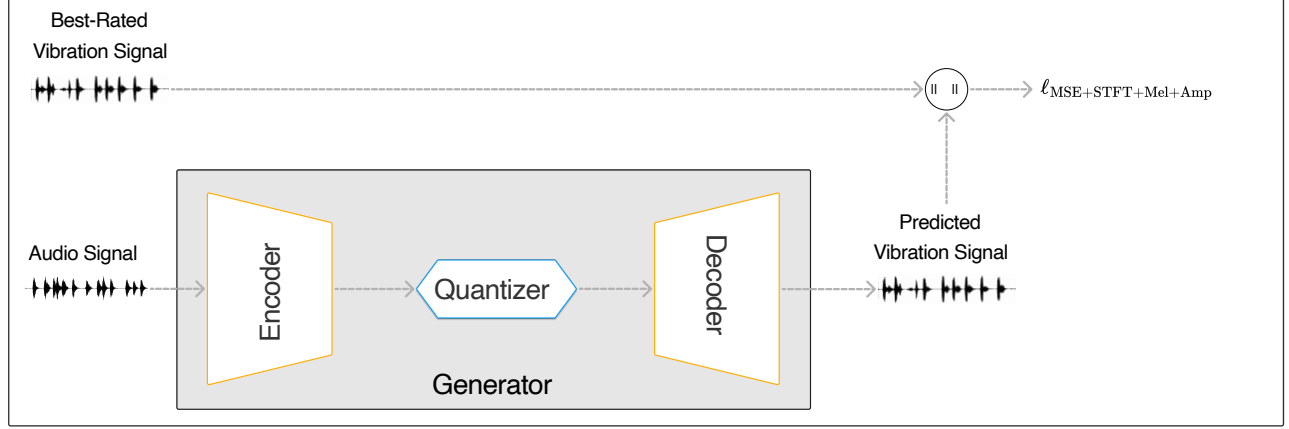
Phase 1: Generator Pre-training. The generator’s decoder is trained using a blended reconstruction loss. For each sound sample, we create a blended target vibration by computing a weighted average of all four reference vibrations, where the weights are the normalized human ratings of each vibration. This approach allows us to generate composite signals that no individual algorithm can produce, thereby enriching the dataset and mitigating limitations in sample size. This is similar to established techniques in audio processing, such as linear time-domain audio mixing used in source-separation training [23, 36] and mixup methods that combine waveforms for regularization [83, 94]. The generator minimizes a loss function between its output and this blended target, combining the MSE, STFT, Mel, and Amp losses (same as Top-Pair Sound2Hap’s losses) with an additional perceptual loss term, which further improves feature-level similarity. In this phase, the blended signals encourage the generator to learn a broader range of vibration patterns based on varied user preferences before being fine-tuned for task-specific refinement in the next phase.

Phase 2: Adversarial Fine-tuning. In the second phase, we introduce a convolutional discriminator and refine the generator in an adversarial setup. The generator’s training objective now combines more terms: (1) Reconstruction loss, which contains the same components as the Phase 1 objective; (2) Adversarial loss that pushes the generator to create outputs that the discriminator classifies as real; (3) Rating loss that rewards the generator for producing vibrations that the discriminator predicts will receive a high user rating.

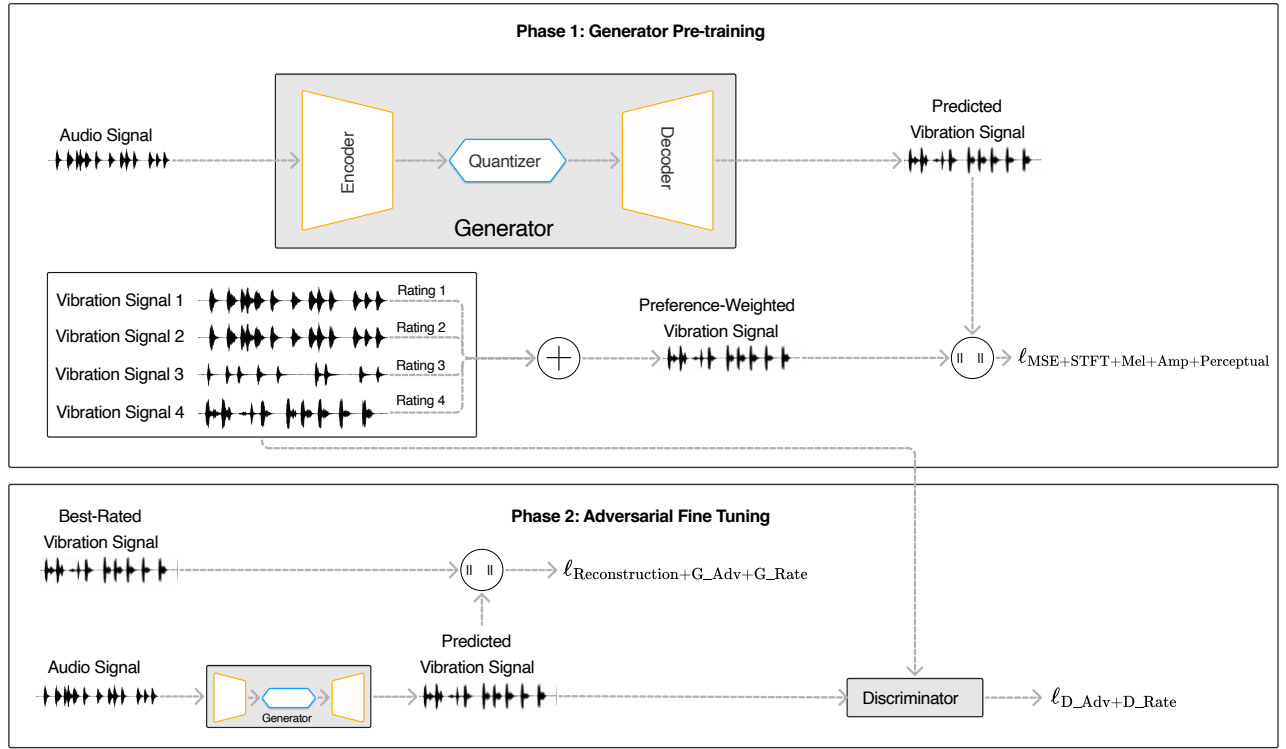
The discriminator is trained on each clip’s four “real” vibration signals from our dataset alongside the generator’s predicted signal, to classify real versus generated vibrations while simultaneously predicting user ratings for real vibrations. See Supplementary Material for the loss function and details of the convolutional discriminator.

6.3 Sound2Hap’s Vibration Generation Time

We evaluate the model’s generation time on sound clips ranging from 1 to 20 seconds in duration. All benchmarks were conducted on 50 randomly selected 5-second, single-channel sound clips (44.1 kHz, 16-bit PCM). Each clip was segmented into five 1-second clips



(a) Top-Pair Sound2Hap



(b) Preference-Weighted Sound2Hap

Figure 6: Sound2Hap model with two training schemes and loss functions: (a) Top-Pair Sound2Hap is trained to directly replicate the best-rated vibration signal; (b) Preference-Weighted Sound2Hap employs a two-phase training strategy: an initial pre-training phase using a blended vibration target, followed by an adversarial fine-tuning phase with a discriminator to refine the generator’s output.

(250 clips total) and into two 2-second clips (ignoring the last second, 100 clips total). To simulate longer inputs, each clip was repeated to create 10-second ($\times 2$) and 20-second ($\times 4$) versions, providing an upper bound for the typical duration of vibration events in user applications. Benchmarks were run on a workstation with an Intel i7-14700HX CPU (16 GB RAM) and an NVIDIA RTX 4060 Laptop GPU (8 GB VRAM), using CUDA acceleration for model execution. The inference times reported below represent averages over all test clips: 250 clips of 1 s, 100 clips of 2 s, and 50 clips each at 5 s, 10 s, and 20 s.

Both Sound2Hap variants maintained vibration generation latencies below 1 second. For the Top-Pair Sound2Hap, inference required 0.44 s for 1 s and 2 s sound clips, 0.48 s for 5 s, 0.52 s for 10 s, and 0.62 s for 20 s. The Preference-Weighted variant showed similar performance: 0.44 s (1 s and 2 s), 0.52 s (5 s), 0.54 s (10 s), and 0.63 s (20 s). These results suggest that Sound2Hap maintains a stable, near-constant overhead for short clips, increasing only modestly with longer durations. In practice, vibration durations are typically limited to a few seconds to prevent overwhelming or numbing the user’s sense of touch.

7 User Study 2: Evaluating the Sound2Hap Generative Model

We conducted an in-person user study with 15 participants to compare vibrations generated by the two Sound2Hap variants against the best-performing algorithm for each sound class from User Study 1 as a baseline. We selected two distinct sound sources, new sound clips from the ESC-50 dataset [69] and the BBC Sound Effects library [10], to test generalizability across different audio datasets.

7.1 Baseline

As a baseline, we used the best-performing signal-processing algorithm for each of the 50 sound classes. This setup mimics a “perfect classifier” that always detects the correct sound class and selects its best-rated method (e.g., Pitch Matching for Dog and HapticGen for Cat sounds). This baseline is motivated by prior work [89, 91], which used a classification approach to identify relevant sounds in an audio track and play corresponding pre-designed vibrations. In practice, we simply assign the best-performing algorithm per class in Study 1 (Table 1) whenever a sound from that class appears. This approach does not involve training or performing actual classification; instead, it assumes perfect knowledge of class membership and access to user preference data per class. By always applying the empirically best-performing method for each class, this baseline represents a practical upper bound on the performance of existing signal-processing methods and serves as a strong baseline for our generative models.

7.2 Study Methods

7.2.1 Audio-Vibration Stimuli. We selected a total of 600 sound clips for evaluation from the ESC-50 dataset and the BBC Sound Effects library to convert to vibrations. Because the Baseline method depends on the best-performing algorithm within each class, we selected clips spanning the same 50 sound classes from both datasets.

These datasets contain no shared sound clips due to their distinct sources¹⁰ and licensing terms.

ESC-50 Dataset. We selected 306 clips from the remaining 1,000 clips in this dataset. We randomly selected six clips from each of the 50 classes. As an exception, we included 12 clips for the “mouse clicking” class since we did not find a matching class in the BBC dataset.

BBC Sound Effects Library. We selected 294 clips from this library. For each ESC-50 class, we searched the BBC database using the exact class name as a keyword, sorting results by relevance. When exact class names did not yield relevant results, we refined our keywords with related terms. For example, using “tap dripping” instead of “water drops”, or “circular saw” instead of “chain saw”, to locate functionally equivalent sounds. For each class, we included the first six clips that (1) contained a clear segment of the target sound, and (2) were free of overlapping sound effects. All selected BBC clips were manually trimmed to the relevant sound segment. To have a uniform clip duration, segments shorter than 5s were padded, and those longer than 5s were trimmed.

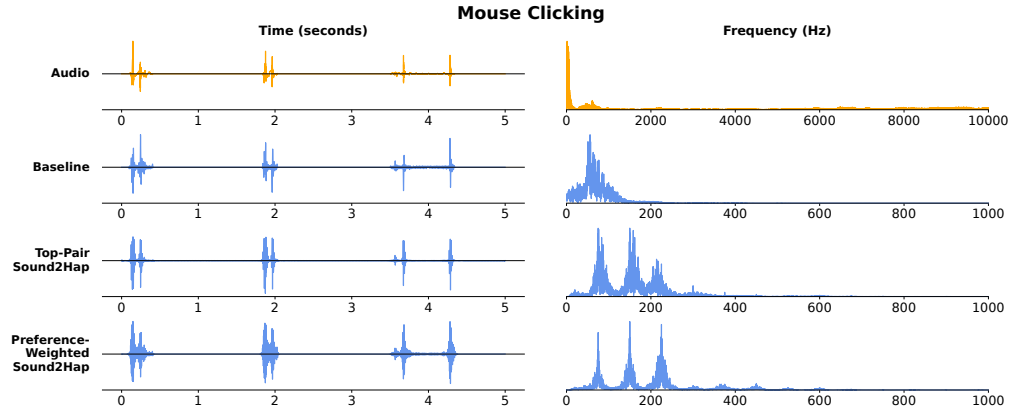
All audio was then converted to mono, resampled to 44.1 kHz, and saved at 16-bit depth. We created three vibrations for each clip using the Baseline method, Top-Pair Sound2Hap, and Preference-Weighted Sound2Hap. Figure 7 shows the audio and vibration signals for two example sound clips.

7.2.2 Participants. We recruited 15 participants (10 male, 5 female; mean age = 23.9 years, SD = 1.9) through online advertisements at the authors’ institution. Participants met the same eligibility criteria as User Study 1. Six participants reported haptics expertise: all with vibrations, two with force-feedback, three with mid-air ultrasound, and one with electrotactile. All other participants indicated minimal or no prior experience with haptic systems. Participants received a \$15 cash reward for their time. The study was conducted under the same IRB approval as User Study 1.

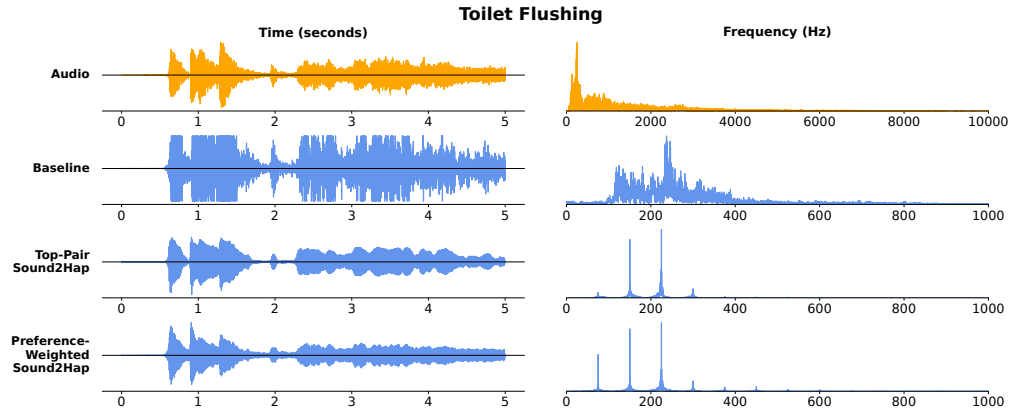
7.2.3 Procedure. The consent form, demographic questionnaires, and introductions to hardware and interface setup followed the same protocol as User Study 1, except that the interface now presented three vibrations. During the main session, each participant evaluated 40 sound clips from 40 distinct sound classes, ensuring broad coverage. For each trial, participants first listened to the sound clip, then experienced three corresponding vibration signals generated by different algorithms (Baseline and two Sound2Hap variants). The vibrations were labeled as generated by Algorithm A, B, and C, with the assignment counterbalanced across participants to reduce order effects. Participants were asked to remember their impressions of each algorithm for subsequent questionnaires and interviews. They could replay the audio and vibrations as many times as needed. After experiencing each vibration, they rated how well it matched the sound using the same 0 to 100 scale as User Study 1. They could proceed only after playing the audio, feeling the vibrations, and entering ratings.

After all trials, participants completed the Haptic Experience Inventory (HXI) questionnaire [76] for each algorithm. The HXI is a validated 20-item questionnaire that measures haptic experience

¹⁰<https://freesound.org/>



(a) A *mouse click* sound clip from the ESC-50 dataset and three generated vibrations in time and frequency domains. The Baseline vibration was generated using Pitch Matching.



(b) A *toilet flush* sound clip from the BBC library and three generated vibrations in time and frequency domains. The Baseline vibration was generated using Frequency Shifting.

Figure 7: Two example audio signals and corresponding vibrations by Baseline, Top-Pair Sound2Hap, and Preference-Weighted Sound2Hap. For the *mouse click* sound, the participant gave ratings of 69 (Baseline), 100 (Top-Pair Sound2Hap), and 79 (Preference-Weighted Sound2Hap) out of 100, and for the *toilet flush* sound, the ratings were 0, 100, and 88 for the three methods, respectively.

across five dimensions: Autotelics, Realism, Involvement, Harmony, and Discord. The 20 items were randomized for each participant and rated on a 7-point Likert scale. At the end, we conducted a 10-minute semi-structured interview about algorithm preferences, factors influencing them, which algorithms suited specific audio categories, and any surprising outputs.

7.3 Results

We present the study results for both quantitative user ratings and qualitative interview insights.

7.3.1 Quantitative User Ratings.

Audio-Vibration Match Ratings. Figure 8 shows the distribution of ratings for the three algorithms. Across the 600 signals, average ratings were 58.28 ($SD = 16.20$) for the Baseline, 76.28 ($SD = 10.14$)

for Top-Pair Sound2Hap, and 73.27 ($SD = 10.05$) for Preference-Weighted Sound2Hap. Assumptions of normality were met, but sphericity was violated; therefore, we applied Greenhouse-Geisser correction. A one-way repeated-measures ANOVA revealed a significant effect of algorithm on signal ratings with $F(1.429, 20.003) =$

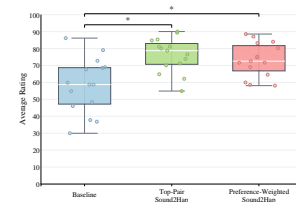


Figure 8: Ratings for audio-vibration match.

10.466, $p = .002$, and $\eta_p^2 = .428$. Pairwise comparisons with Holm-Bonferroni correction showed that both Top-Pair Sound2Hap ($p = .011$) and Preference-Weighted Sound2Hap ($p = .013$) received significantly higher ratings than the Baseline, with no significant difference between the two Sound2Hap variants.

The average ratings on the ESC-50 and the BBC datasets were comparable for each algorithm. The Baseline received the lowest and most varied ratings, with an average of 56.81 ($SD = 33.24$) on ESC-50 and 59.80 ($SD = 30.47$) on BBC. Top-Pair Sound2Hap performed best, with ratings of 76.31 ($SD = 23.27$) on ESC-50 and 76.24 ($SD = 24.22$) on BBC. Preference-Weighted Sound2Hap also performed well, with ratings of 72.65 ($SD = 24.65$) on ESC-50 and 73.93 ($SD = 24.32$) on BBC. These results suggest Sound2Hap’s generalizability across sound datasets.

HXI Factor	F	P	η_p^2
Autotelics	1.444	.251	.093
Involvement	5.202	.027	.271
Realism	4.999	.014	.263
Discord	6.795	.013	.327
Harmony	7.240	.003	.341
General Score	6.007	.017	.300

Table 2: Results of repeated measures ANOVA for five HXI factors and the General Score. Boldface indicates factors that were statistically significant at $p < .05$. All factors except Autotelics exhibited medium effect sizes, suggesting practical significance.

HXI Ratings. Figure 9 shows the distribution of HXI ratings for each algorithm. The HXI factors range from 1 (Strongly Disagree) to 7 (Strongly Agree), and the General Score is the sum of Autotelics, Involvement, Realism, Harmony, and the reverse of the score for Discord (8-Discord). The HXI ratings met the assumption of normality; however, only the Realism and Harmony factors satisfied the sphericity assumption. Therefore, we applied Greenhouse-Geisser corrections for the remaining factors and General Score. Table 2 summarizes the one-way repeated-measures ANOVA results. With the exception of Autotelics, all other factors and the General Score showed significant effects. Pairwise comparisons with Holm-Bonferroni correction revealed several significant effects. On the Realism factor, Preference-Weighted Sound2Hap ($M = 5.75$, $SD = 0.90$) received significantly higher ratings than Baseline ($M = 4.45$, $SD = 1.47$) at $p = .014$. On Discord (lower values indicate better performance), both Top-Pair Sound2Hap ($M = 3.17$, $SD = 1.28$) and Preference-Weighted Sound2Hap ($M = 3.43$, $SD = 0.75$) received significantly lower ratings than Baseline ($M = 4.67$, $SD = 1.21$) at the same $p = .016$. For Harmony, both Top-Pair ($M = 5.52$, $SD = 0.89$) and Preference-Weighted ($M = 5.52$, $SD = 0.83$) were rated significantly higher than Baseline ($M = 4.28$, $SD = 1.30$) at $p = .018$ and $p = .005$, respectively. Finally, for the General Score, both Top-Pair ($M = 27.08$, $SD = 4.12$) and Preference-Weighted ($M = 26.65$, $SD = 3.88$) outperformed Baseline ($M = 21.68$, $SD = 5.19$) at $p = .025$ and $p = .016$. No other pairwise comparisons reached significance.

7.3.2 Qualitative Interview Results.

Algorithm Preferences. When asked about their most-preferred algorithm, eight participants selected Top-Pair Sound2Hap, five preferred Preference-Weighted Sound2Hap, and two chose the Baseline. All participants reported scores between 6 and 9 on a 10-point

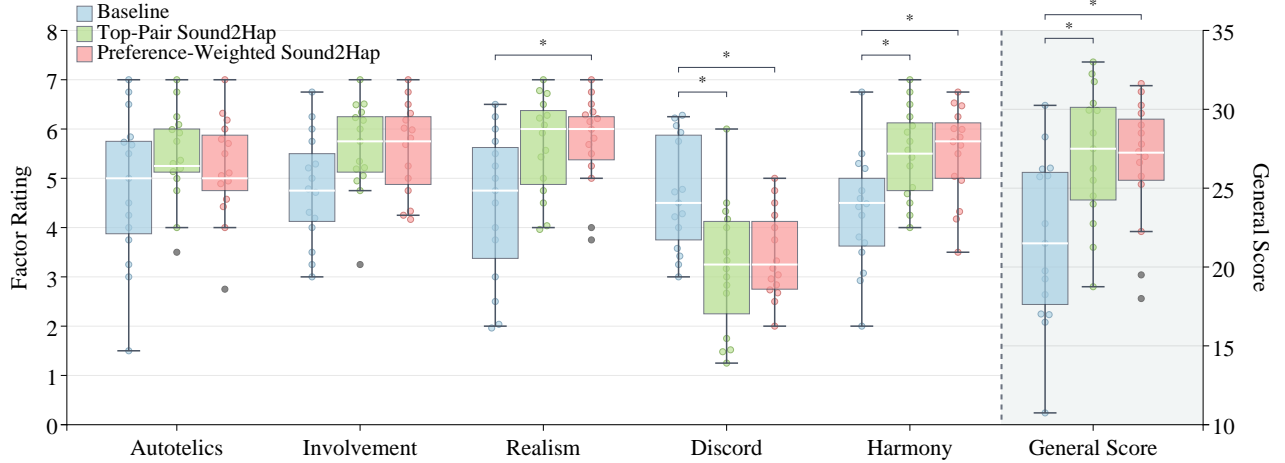


Figure 9: Algorithm ratings on Haptic Experience Index (HXI) questionnaire. Higher values reflect better performance on all factors except Discord, where lower values indicate better performance. * indicates statistical significance after Holm-Bonferroni correction. Autotelics, Involvement, Realism, Discord, and Harmony are rated on a 1 to 7 scale corresponding to the left y-axis, while the General Score, which sums these factors (reversing Discord), has a wider range corresponding to the right y-axis.

scale regarding confidence in their selections. These qualitative preferences aligned with the user ratings.

Factors Shaping Preferences. The primary factors influencing participants’ choices were the accuracy and faithfulness of the vibration patterns to the original audio. Most participants ($n=8$) emphasized the importance of intensity and rhythm, noting that the preferred algorithms better captured the strength and timing of the sound. Accuracy in resembling the pitch was also mentioned ($n=3$), while two participants highlighted overall consistency across the 40 trials. For example, P3 explained, “It [Top-Pair Sound2Hap] was resembling the sound more accurately than the others...the algorithm was also following the breaks and continuing the vibrations according to the intensity of the sound.” Harmony between the audio and haptic feedback was also noted ($n=4$), with participants favoring vibrations that felt most synchronized with the source audio and aligned with their own perceptions and expectations. Top-Pair Sound2Hap was praised for its precision and ability to capture subtle sounds, rhythm, and intensity ($n=8$), though one user felt its vibrations were occasionally too strong. Likewise, Preference-Weighted Sound2Hap was frequently considered the best for its accuracy and consistency ($n=8$), but a few participants ($n=4$) noted it sometimes missed nuances or felt flat. Opinions on the Baseline were the most divided; while some ($n=5$) liked it for capturing small nuances and distinct changes, more users ($n=9$) found it inconsistent, confusing, and unstable compared to the other algorithms, with an intensity that was either too varied or too low.

Unexpected Experiences. When asked if the vibrations ever felt surprisingly different from the audio, the Baseline was the most frequent source of negative surprise ($n=6$). Participants described their vibrations as feeling “very off,” “too continuous,” or “completely different” from what they expected. For example, P4 remarked, “Algorithm B [Baseline] would always seem really different, so I would have to go back and play the audio again [to continue rating].” However, the Baseline also provided a few positive surprises ($n=2$), with one user noting it captured the intensity of a fire crackling in “an amazingly realistic way”. In contrast, surprises from Top-Pair Sound2Hap and Preference-Weighted Sound2Hap were mostly positive ($n=3$), with users appreciating the nuanced sensations for sounds like a ticking clock, a swinging door, and breathing.

Real-World Applicability. Looking toward real-world applications, most participants ($n=13$) expressed a preference for experiencing rhythmic and distinct sounds as vibrations. Popular examples included typing, clapping, thunderstorms, and various alert sounds. Some ($n=5$) also noted that vibrations could be engaging in specific contexts, such as gaming or watching movies, but not during music ($n=2$) and meditation ($n=1$). Conversely, nearly all participants ($n=14$) reported an aversion to converting continuous, loud, or unpleasant noises into vibrations. Helicopters, airplanes, chainsaws, and other jarring sounds were commonly cited as undesirable, as they were perceived as annoying or overwhelming. As one participant (P15) summarized: “Rhythmic sounds are more suitable for being converted into vibrations, whereas continuous noisy sounds feel a bit strange.”

7.3.3 Model-generated Vibrations Comparison with Reference Vibrations. To evaluate why the Sound2Hap variants outperformed

the Baseline, we used 200 sound clips from Study 1’s validation set to analyze how closely the models reproduce vibration patterns. We compare Sound2Hap’s model-generated vibrations to two references: (a) the Baseline vibrations, produced by the winning algorithm at the sound-class level, and (b) the clip-level best vibrations, which received the highest human ratings for each clip in Study 1 (i.e., the winning algorithm for each clip). We generated vibrations for all 200 clips using both Sound2Hap variants and computed root-mean-square error (RMSE) against these two reference signals. Top-Pair Sound2Hap showed an RMSE of 0.335 from the clip-level best vibrations and 0.358 from the Baseline; Preference-Weighted Sound2Hap showed a similar pattern (0.332 vs. 0.357). These results demonstrate that both variants generate signals closer to the clip-level best vibrations than to the Baseline, supporting the conclusion that Sound2Hap captures clip-level perceptual structure beyond what class-level methods can provide.

In our subjective evaluation, we compared Sound2Hap against the class-level Baseline rather than the clip-level best vibrations, as the latter are impractical for real-world use. Identifying the best vibration for a new sound clip requires human ratings across multiple algorithms, so clip-level labels are available only in the curated Study 1 dataset, not in open-world settings (as in Study 2). Class-level knowledge, which reflects which algorithm performs best on average for a sound class, therefore represents the strongest feasible alternative for comparison.

8 An Interactive Web Tool for Audio-Vibration Visualization and Generation

We developed an online tool (Figure 10) to facilitate access to the audio-vibration dataset and algorithms. The tool enables researchers and designers to browse our dataset through an interactive visualization, explore patterns of user ratings for various sounds, and generate new haptic signals for their audio clips. Specifically, the tool includes:

- (1) An *Interactive Visualization View* with a sunburst chart and a line-chart, allows users to explore trends across classes, categories, and conversion algorithms.
- (2) A *Gallery View* at the bottom lets users browse 1,000 sound clips spanning 50 sound classes, together with their 4,000 vibration counterparts and 8,000 human ratings.
- (3) A *Pop-up View* shows each sound clip and corresponding vibrations with playback options, waveforms visualization, and user perception ratings.
- (4) A *Filter Panel* allows users to navigate the dataset efficiently by selecting subsets of sounds, categories, or conversion algorithms.
- (5) A *Generation Environment* lets users upload their own sound clips and generate vibrations using the four signal-processing methods and the two Sound2Hap variants introduced in this work. Generated vibrations are visualized and can be played back or downloaded for further refinement or deployment.

9 Discussion

Below, we first reflect on the result, then discuss the utility of Sound2Hap and outline its limitations and future work.

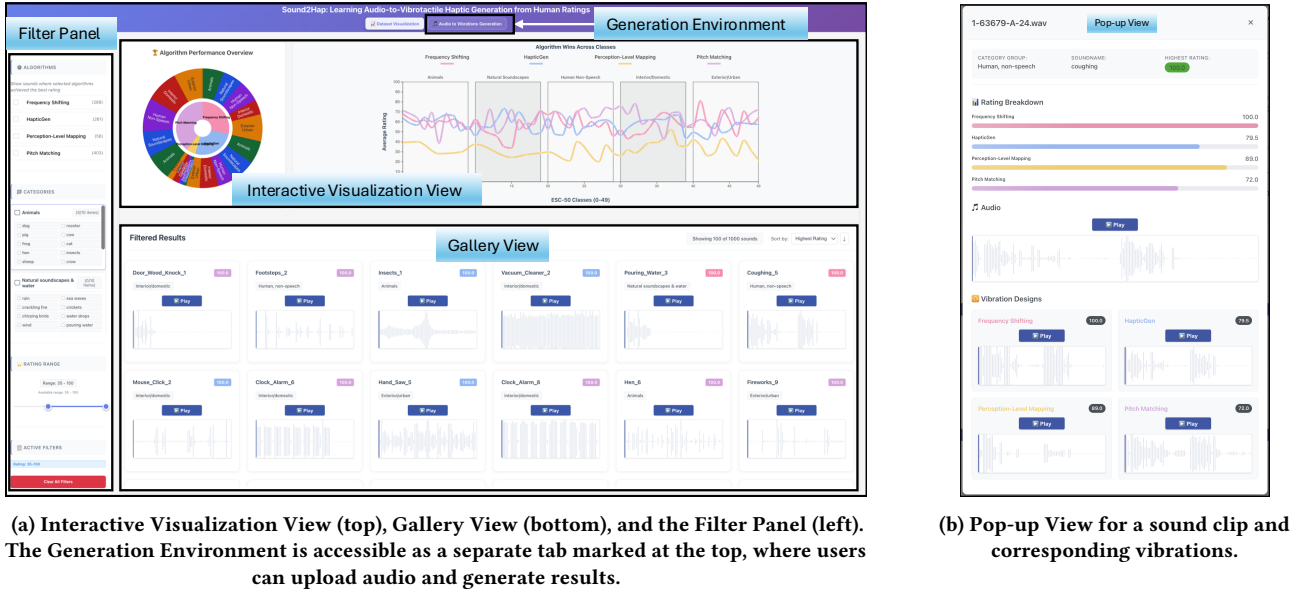


Figure 10: Screenshots of the interactive web tool for audio-vibration visualization and generation.

9.1 Reflecting on Results

User Perception of Audio-to-Vibration Methods. Our work is the first to directly compare user perception of multiple audio-to-vibration conversion methods across diverse sounds. Our results showed that Pitch Matching achieved the highest audio-vibration match ratings (Mean = 62 out of 100), followed by HapticGen ($M = 57$) and Frequency Shifting ($M = 56.9$), while Perception-Level Mapping scored lowest ($M = 31.2$). Although each method performed well for certain clips, their average performances dropped across varied environmental sounds. In Study 2, the Baseline method, which applied the best-performing algorithm for each sound class, achieved a similar average rating ($M = 58$), suggesting that class-level associations cannot fully capture the nuances of individual sounds. In contrast, Sound2Hap achieved significantly higher audio-vibration match ratings ($M = 76$) and better HXI scores for Harmony, Discord, and General Score compared to the Baseline across two datasets and all sound classes. This difference can be partly explained by the training and evaluation setup. The Baseline relies on class-level averages derived from Study 1, where the winning algorithm was identified per sound class (e.g., dog barks) based on mean ratings across clips. While this provides a strong reference, it disregards the substantial clip-level variability observed within each class (e.g., different dog barks). In contrast, Sound2Hap learns directly from clip-level human ratings that encode these finer perceptual nuances, allowing the model to capture subtle relationships between acoustic features and user preference. According to Figure 5a, Frequency Shifting, Pitch Matching, and HapticGen all received comparable mean ratings of around 60 across categories, consistent with the Baseline’s mean ratings in Study 2. However, their highest clip ratings frequently approached 100, indicating that Sound2Hap primarily learns from the top-rated vibration for each audio clip. As a result, the model generalizes more effectively to unseen sounds in Study 2, even within the same class. We further attribute this

improvement to the use of an efficient low-loss audio representation in a compact latent space, as enabled by a pretrained audio encoder. Together, the results support the efficacy and robustness of our learning-based approach to audio-to-vibration conversion.

Comparison between Sound2Hap Variants. Both Sound2Hap variants performed comparably, with users showing a stronger qualitative preference for the Top-Pair model. Both variants outperformed the Baseline in user ratings and achieved similar generation times (~ 0.5 s). Qualitative interview feedback revealed that more participants preferred the Top-Pair Sound2Hap for its precision and ability to capture rhythm and intensity, while describing the Preference-Weighted version as sometimes feeling flat or missing subtle nuances. This perceptual difference reflects the nature of their training targets. The Preference-Weighted Sound2Hap relied on a blended target that combined four vibration signals from different algorithms. This blended approach can introduce destructive interference, smoothing out salient details and producing less distinctive signals. By contrast, the Top-Pair Sound2Hap was trained on a single, high-quality, human-endorsed sample for each sound, providing a cleaner and perceptually coherent learning signal.

Spectral characteristics. As seen in Figure 7, the Sound2Hap variants learned to produce vibrations with distinct frequency spikes rather than a broad spectrum. This is likely a consequence of the training data, where two of the four source algorithms generate signals centered on specific frequencies (e.g., 175 and 210 Hz for Perception-Level Mapping, $200 (\pm 50)$ Hz for HapticGen). Likewise, the Pitch Matching algorithm contributes to this characteristic by predicting a single, dominant frequency for each short audio segment, constrained to a 50 to 400 Hz range. Although this frequency varies smoothly over time, its focus on a single value at any given moment could further encourage the model to generate a narrow frequency spectrum. Thus, the model learned that preferred haptic

signals often contain concentrated energy, with the highest energy frequently centered around 200 Hz, which aligns with the peak of human vibrotactile sensitivity [30]. From user ratings, the model learns that a few frequency spikes can adequately capture the audio nuances while feeling smooth and pleasant to touch. Moreover, this result is consistent with previous research, which demonstrates that the temporal envelope (i.e., rhythmic features) of a haptic signal dominates human perception, proving more salient than the signal’s detailed spectral content [54, 67].

9.2 Utility of Sound2Hap

Web Tool and Dataset. The interactive tool and dataset are designed to support exploration and reuse by researchers and designers. The visualization interface enables users to analyze the dataset by sound categories, classes, or user ratings. For example, a multimodal designer can examine, for each algorithm, the number of sound clips that received the highest ratings within a category (e.g., natural soundscapes) or class (e.g., rain, wind) or filter for highly rated audio-vibration pairs (e.g., scores >90/100) and download them for direct use in their applications. Furthermore, the tool and open-source code enable researchers to generate large-scale vibration datasets with aligned audio and text labels from existing corpora [18, 29, 37], supporting the training of sensory language models that incorporate haptics while reducing the time and cost of manual haptic data collection. By providing access to both the dataset and multiple conversion methods, the tool lowers barriers to creating high-quality, expressive haptic experiences and facilitates further model development.

Multimodal Applications and Research. The model’s perceptually aligned audio-based vibrations support various multimodal use cases in virtual reality, gaming, and interaction design research. For VR and video games, designers can use Sound2Hap to create expressive audio-vibration stimuli and integrate them in their virtual environments to enhance user immersion and engagement. While our results suggest that the model has a fast generation time ($\sim 0.5s$) and is suitable for offline use, real-time conversion (e.g., in VR or gaming) requires stricter latency ($\sim 25 - 75ms$) [2, 64, 91, 92] and would necessitate further model acceleration, such as through knowledge distillation techniques [33]. Beyond applications, high-quality audio-vibration pairs provide congruent stimuli for multimodal interaction research, for example, studying how eye-free feedback modalities affect task performance in XR (e.g., gaze-based object selection) [19]. Furthermore, prior work shows that altering environmental sounds can change self-perception [82], such as footstep sounds modulating perceived body weight [81]. Our approach opens opportunities to investigate whether tactile cues for environmental sounds can produce similar psychological effects.

Moreover, the interchangeable use of audio and vibration cues or sequential playback of the two modalities used in our evaluation also reflects several real-world application scenarios. Interchangeable or substitutive use of modalities is common in alerts, notifications, accessibility cues, and assistive feedback, where either modality may be used in place of the other depending on context. For example, vibration-only notifications in silent mode or vibration replacing audio for accessibility. Sequential playback can support applications that want to notify users in stages, for

example, delivering a vibration alert on a mobile phone before an audio notification to give the user an opportunity to silence the device if needed. At the same time, applications that rely on tightly synchronized audiovisual and haptic stimuli, such as immersive VR or interactive media, would require assessing concurrent playback to ensure perceptual congruence. We discuss this consideration as a limitation and direction for future work in Section 9.3, where we note that evaluating both sequential and concurrent playback will be essential for optimizing Sound2Hap in real-world multimodal environments.

Accessibility Applications. Sound2Hap also holds potential to support accessibility for users with visual or hearing impairments. For instance, perceptually congruent haptics could complement captions for environmental sounds (e.g., “[people laughing]”, “[rain sound]”) to enhance video accessibility for deaf and hard-of-hearing users [3]. For blind and low-vision users, such vibrations may improve video immersion [40] and support clearer audio descriptions [66]. As all our studies involved able-bodied participants, future work should evaluate and refine Sound2Hap with feedback from visually and hearing-impaired users to enhance its perceptual performance and assess its practical utility for accessibility applications.

9.3 Limitations and Future Work

Our work has four main limitations that open up avenues for future research. First, while our dataset represents the largest and most diverse collection of everyday sound events paired with vibrations to date, real-world sounds span an even broader range. Furthermore, Sound2Hap is only tested on sound clips containing a single, salient sound source, while real-world audio often includes overlapping or blended events. Future work can assess the efficacy of Sound2Hap for a wider range of single and overlapping environmental sounds. For overlapping sounds, a promising direction is to develop a modular pipeline that integrates a pre-processing step using audio source separation [55, 63] or saliency detection [71] algorithms. This would allow the system to isolate a foreground event from irrelevant ones (e.g., background noise or speech) and feed the separate and cleaner audio clips to Sound2Hap, expanding its practical applicability. In addition, when audio contains multiple simultaneous sources, the model could be extended to generate multiple vibrations across a vibrotactile array (e.g., a haptic vest), with each vibration representing a distinct sound or source location. Finally, future work could adapt Sound2Hap for other sound domains relevant to interaction design, such as earcons and voice-user-interface cues. Incorporating existing datasets of interface sounds [15, 32] would help generalize the model to create vibrations for symbolic, short-duration sounds commonly used for UI feedback, such as haptic cues that supplement auditory feedback in noisy environments or replace it in silent modes.

Second, our model development and evaluation focus on overall performance across users, but the quantitative ratings and qualitative preferences suggest that users may have individual differences in the algorithm they prefer. Relatedly, in our studies, the participant sample showed a gender imbalance (Study 1: 24 male / 10 female; Study 2: 10 male / 5 female), which may limit the generalizability of the findings. Since prior research suggests that gender

can influence tactile simultaneity thresholds [28], perceived vibration intensity and discomfort [62], and tactile spatial acuity [68], including a gender-balanced sample could reveal user differences in sensitivity to vibration signal variations or synthesis artifacts. Future work should aim to recruit more demographically balanced samples to examine whether gender or other individual factors (e.g., age, music background) impact user preference for audio-to-haptic conversion. Furthermore, future directions could explore adapting personalization approaches from other modalities [4, 56, 73] to further calibrate the model's output to an individual user's perception and preferences.

Third, our evaluation relies on sequential playback of audio and vibrations, allowing participants to focus on tactile nuances without auditory masking. This design also supports unimodal applications, such as content accessibility for users with sensory impairments. Thus, our model is optimized for audio-to-haptic translation when felt in isolation. Yet, many real-world applications rely on concurrent audio-haptic playback, thus our results from isolated vibration perception may not fully apply to them. Prior research shows that synchronizing audio and vibration may alter haptic perception; for instance, a low-frequency sound can cause a high-frequency vibration to be perceived as having a lower frequency [34]. Future work should therefore evaluate the holistic quality of the combined, concurrent stimulus to further optimize haptics for multimodal playback.

Lastly, our evaluation setup focuses on finger-based perception using a voice-coil vibration actuator. This configuration enables precise and high-sensitivity evaluation of vibrotactile signals. However, tactile acuity varies across different body sites. Since the fingertip contains a high density of mechanoreceptors [41], the reported results may represent an upper bound on the perceptual gains of Sound2Hap. Also, other actuator types (e.g., eccentric rotating mass or linear resonant actuators) produce distinct vibration profiles [20]. Future work should therefore examine how well Sound2Hap generalizes across actuator types, body locations, and broader usage contexts.

10 Conclusion

This work introduces Sound2Hap, a data-driven model for generating vibrotactile signals from diverse environmental sounds. By combining large-scale perceptual data collection with two generative model variants, we address the limitations of prior rule-based audio-to-haptics approaches. Our user study shows that Sound2Hap produces signals that align with human perceptual preferences and generalize across sound domains. We release our dataset, models, and an interactive tool to support future research and applications. We envision Sound2Hap enabling scalable haptic content creation and informing new multimodal experiences across XR, accessibility, and entertainment.

Acknowledgments

We thank Mainak Malay Saha for developing the web tool. We also thank the anonymous reviewers, our colleagues, and the study participants for their input on this project. This work was supported by research grants from the National Science Foundation (#2339707) and VILLUM FONDEN (VIL50296).

References

- [1] Olusola O Abayomi-Alli, Robertas Damaševičius, Atika Qazi, Mariam Adedoyin-Olowe, and Sanjay Misra. 2022. Data augmentation and deep learning methods in sound classification: A systematic review. *Electronics* 11, 22 (2022), 3795.
- [2] Bernard D Adelstein, Durand R Begault, Mark R Anderson, and Elizabeth M Wenzel. 2003. Sensitivity to haptic-audio asynchrony. In *Proceedings of the 5th international conference on Multimodal interfaces*. 73–76.
- [3] Sooyeon Ahn, Gyungmin Jin, Gunhyuk Park, and Jin-Hyuk Hong. 2025. Enhancing Video Experiences for DHH Individuals through Sound-Inspired Motion Caption-based Spatiotemporal Tacton. *IEEE Transactions on Haptics* (2025).
- [4] Yuval Alaluf, Elad Richardson, Sergey Tulyakov, Kfir Aberman, and Daniel Cohen-Or. 2024. MyVlm: Personalizing vlms for user-specific queries. In *European Conference on Computer Vision*. Springer, 73–91.
- [5] Yuki Ban and Yusuke Ujitoko. 2018. TactGAN: Vibrotactile designing driven by GAN-based automatic generation. In *SIGGRAPH Asia 2018 Emerging Technologies*. 1–2.
- [6] bHaptics. 2021. bHaptics Studio. <https://www.bhaptics.com/software/studio/>. Retrieved August, 2025.
- [7] bHaptics. 2021. TactSuit X40. <https://www.bhaptics.com/tactsuit/tactsuit-x40>. Retrieved August, 2025.
- [8] Antoine Bourachot, Tifanie Bouchara, and Olivier Cornet. 2023. Impact of an audio-haptic strap to augment immersion in VR video gaming: a pilot study. In *Proceedings of the 18th International Audio Mostly Conference*. 147–153.
- [9] Carmen Branje, Gabe Nespoli, Frank Russo, and Deborah I Fels. 2014. The effect of vibrotactile stimulation on the emotional response to horror films. *Computers in Entertainment (CIE)* 11, 1 (2014), 1–13.
- [10] British Broadcasting Corporation. 2025. BBC Sound Effects. <https://sound-effects.bbcrewind.co.uk/>. Accessed: 2025-08-11.
- [11] Shaoyu Cai, Lu Zhao, Yuki Ban, Takuji Narumi, Yue Liu, and Kening Zhu. 2022. GAN-based image-to-friction generation for tactile simulation of fabric material. *Computers & Graphics* 102 (2022), 460–473.
- [12] Shaoyu Cai and Kening Zhu. 2022. Multi-modal transformer-based tactile signal generation for haptic texture simulation of materials in virtual and augmented reality. In *2022 IEEE International Symposium on Mixed and Augmented Reality Adjunct (ISMAR-Adjunct)*. IEEE, 810–811.
- [13] Shaoyu Cai, Kening Zhu, Yuki Ban, and Takuji Narumi. 2021. Visual-tactile cross-modal data generation using residue-fusion gan with feature-matching and perceptual losses. *IEEE Robotics and Automation Letters* 6, 4 (2021), 7525–7532.
- [14] Guanqun Cao, Jiaqi Jiang, Ningtao Mao, Danushka Bollegala, Min Li, and Shan Luo. 2023. Vis2hap: Vision-based haptic rendering by cross-modal generation. *arXiv preprint arXiv:2301.06826* (2023).
- [15] Zijiang Cao, António Sá Pinto, and Gilberto Bernardes. 2024. Bisaid: Bipolar semantic adjectives icons and earcons dataset. In *Proceedings of the 21st Sound and Music Computing Conference*.
- [16] Doga Cavdir and Ge Wang. 2020. Felt sound: A shared musical experience for the deaf and hard of hearing. In *Proceedings of the International Conference on New Interfaces for Musical Expression*. 176–181.
- [17] Angela Chang and Conor O'Sullivan. 2005. Audio-haptic feedback in mobile phones. In *CHI'05 extended abstracts on Human factors in computing systems*. 1264–1267.
- [18] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. 2020. Vg-sound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 721–725.
- [19] Hyunsung Cho, Naveen Sendhilnathan, Michael Nebeling, Tianyi Wang, Purnima Padmanabhan, Jonathan Browder, David Lindlbauer, Tanya R Jonker, and Kashyap Todi. 2024. SonoHaptics: an audio-haptic cursor for gaze-based object selection in XR. In *Proceedings of the ACM symposium on user interface software and technology (UIST)*. 1–19.
- [20] Seungmoon Choi and Katherine J Kuchenbecker. 2012. Vibrotactile display: Perception, technology, and applications. *Proc. IEEE* 101, 9 (2012), 2093–2104.
- [21] Stanley Coren, Lawrence M Ward, and James T Enns. 2004. *Sensation and perception*. John Wiley & Sons Hoboken, NJ.
- [22] Corsair. 2020. HS60 HAPTIC Stereo Gaming Headset with Haptic Bass. <https://www.corsair.com/us/en/Categories/Products/Gaming-Headsets/USB-Headsets/HS60-HAPTIC/p/CA-9011228-NA>. Retrieved August, 2025.
- [23] Joris Cosentino, Manuel Pariente, Samuele Cornell, Antoine Deleforge, and Emmanuel Vincent. 2020. Librimix: An open-source dataset for generalizable speech separation. *arXiv preprint arXiv:2005.11262* (2020).
- [24] Fabien Danieau, Anatole Lécuyer, Philippe Guillotel, Julien Fleureau, Nicolas Mollet, and Marc Christle. 2012. Enhancing audiovisual experience with haptic feedback: a survey on HAV. *IEEE transactions on haptics* 6, 2 (2012), 193–205.
- [25] Caluã de Lacerda Patata, Saad Hassan, Lloyd May, Michelle M Olson, Toni D'aurio, Roshan L Peiris, and Matt Huenerfauth. 2025. Tactile Emotions: Multimodal Affective Captioning with Haptics Improves Narrative Engagement for d/Deaf and Hard-of-Hearing Viewers. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–17.

- [26] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. 2022. High Fidelity Neural Audio Compression. *arXiv preprint arXiv:2210.13438* (2022).
- [27] Faraz Faruqi, Maxine Perroni-Scharf, Jaskaran Singh Walia, Yunyi Zhu, Shuyue Feng, Donald Degraen, and Stefanie Mueller. 2025. TactStyle: Generating Tactile Textures with Generative AI for Digital Fabrication. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [28] Gina Geffen, Virginia Rosa, and Michelle Luciano. 2000. Sex differences in the perception of tactile simultaneity. *Cortex* 36, 3 (2000), 323–335.
- [29] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. Audio Set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE ICASSP 2017*. New Orleans, LA.
- [30] A Gescheider, Stanley J Bolanowski, and Kathleen R Hardick. 2001. The frequency selectivity of information-processing channels in the tactile sensory system. *Somatosensory & Motor Research* 18, 3 (2001), 191–201.
- [31] Steven Goodman, Susanne Kirchner, Rose Guttman, Dhruv Jain, Jon Froehlich, and Leah Findlater. 2020. Evaluating smartwatch-based sound feedback for deaf and hard-of-hearing users across contexts. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–13.
- [32] Mandip Goswami. 2025. BeepBank-500: A Synthetic Earcon Mini-Corpus for UI Sound Research and Psychoacoustics Research. *arXiv preprint arXiv:2509.17277* (2025).
- [33] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International journal of computer vision* 129, 6 (2021), 1789–1819.
- [34] Waseem Hassan and Kasper Hornbæk. 2024. Audio-Tactile Integration: Concurrent Audio Feedback Can Shift Vibrotactile Frequency Perception. In *International Conference on Human Haptic Sensing and Touch Enabled Computer Applications*. Springer, 74–89.
- [35] Negin Heravi, Heather Culbertson, Allison M Okamura, and Jeannette Bohg. 2024. Development and evaluation of a learning-based model for real-time haptic texture rendering. *IEEE Transactions on Haptics* 17, 4 (2024), 705–716.
- [36] John R Hershey, Zhuo Chen, Jonathan Le Roux, and Shinji Watanabe. 2016. Deep clustering: Discriminative embeddings for segmentation and separation. In *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 31–35.
- [37] Jaesung Huh, Jacob Chalk, Evangelos Kazakos, Dima Damen, and Andrew Zisserman. 2025. Epic-sounds: A large-scale dataset of actions that sound. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2025).
- [38] Inwook Hwang and Seungmoon Choi. 2014. Improved haptic music player with auditory saliency estimation. In *International Conference on Human Haptic Sensing and Touch Enabled Computer Applications*. Springer, 232–240.
- [39] Inwook Hwang, Hyeseon Lee, and Seungmoon Choi. 2013. Real-time dual-band haptic music player for mobile devices. *IEEE transactions on haptics* 6, 3 (2013), 340–351.
- [40] L. Jiang, M. Phutane, and S. Azenkot. 2023. Beyond Audio Description: Exploring 360° Video Accessibility with Blind and Low Vision Users Through Collaborative Creation. In *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI)*. ACM, 1–12.
- [41] Roland S Johansson and Ake B Vallbo. 1979. Tactile sensibility in the human hand: relative and absolute densities of four types of mechanoreceptive units in glabrous skin. *The Journal of physiology* 286, 1 (1979), 283–300.
- [42] Crescentia Jung, Jazmin Collins, Ricardo E Gonzalez Penuela, Jonathan Isaac Segal, Andrea Stevenson Won, and Shiri Azenkot. 2024. Accessible nonverbal cues to support conversations in VR for blind and low vision people. In *Proceedings of the 26th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–13.
- [43] Daeseok Kang, Chang-Gyu Lee, and Ohung Kwon. 2023. Pneumatic and acoustic suit: multimodal haptic suit for enhanced virtual reality simulation. *Virtual Reality* 27, 3 (2023), 1647–1669.
- [44] Maria Karam, Carmen Branje, Gabe Nespoli, Norma Thompson, Frank A Russo, and Deborah I Fels. 2010. The emoti-chair: an interactive tactile music exhibit. In *CHI’10 Extended Abstracts on Human Factors in Computing Systems*. 3069–3074.
- [45] Maria Karam, Frank A Russo, and Deborah I Fels. 2009. Designing the model human cochlea: An ambient crossmodal audio-tactile display. *IEEE Transactions on Haptics* 2, 3 (2009), 160–169.
- [46] Dong-Geun Kim, Jungeun Lee, Gyeore Yun, Hong Z Tan, and Seungmoon Choi. 2023. Sound-to-touch crossmodal pitch matching for short sounds. *IEEE Transactions on Haptics* 17, 1 (2023), 2–7.
- [47] Sang-Youn Kim and Kyu-Young Kim. 2007. Interactive racing game with graphic and haptic feedback. In *International Workshop on Haptic and Audio Interaction Design*. Springer, 69–77.
- [48] Claudia Krogmeier, Christos Mousas, and David Whittinghill. 2019. Human-virtual character interaction: Toward understanding the influence of haptic feedback. *Computer Animation and Virtual Worlds* 30, 3-4 (2019), e1883.
- [49] Jaebong Lee and Seungmoon Choi. 2013. Real-time perception-level translation from audio signals to vibrotactile effects. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 2567–2576.
- [50] Xinwu Li, Huaping Liu, Junfeng Zhou, and FuChun Sun. 2019. Learning cross-modal visual-tactile representation using ensemble generative adversarial networks. *Cognitive Computation and Systems* 1, 2 (2019), 40–44.
- [51] Yaxuan Li, Yongjae Yoo, Antoine Weill-Duflos, and Jeremy Cooperstock. 2021. Towards context-aware automatic haptic effect generation for home theatre environments. In *Proceedings of the 27th ACM symposium on virtual reality software and technology*. 1–11.
- [52] Chungman Lim, Kevin John, Gyungmin Jin, Hasti Seifi, and Gunhyuk Park. 2025. ChatHAP: A Chat-Based Haptic System for Designing Vibrations through Conversation. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–18.
- [53] Chungman Lim and Gunhyuk Park. 2023. Can a computer tell differences between vibrations?: Physiology-based computational model for perceptual dissimilarity prediction. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [54] Chungman Lim, Gunhyuk Park, and Hasti Seifi. 2024. Designing distinguishable mid-air ultrasound tactions with temporal parameters. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–8.
- [55] Xubo Liu, Qiuqiang Kong, Yan Zhao, Haohe Liu, Yi Yuan, Yuzhuo Liu, Rui Xia, Yuxuan Wang, Mark D Plumbley, and Wenwu Wang. 2024. Separate anything you describe. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* (2024).
- [56] Michał Maćkowski, Piotr Brzoza, Mateusz Kawulok, Rafał Meisel, and Dominik Spinczyk. 2023. Multimodal presentation of interactive audio-tactile graphics supporting the perception of visual information by blind people. *ACM Transactions on Multimedia Computing, Communications and Applications* 19, 5s (2023), 1–22.
- [57] Sebastian Merchel, Ercan Altınsoy, and Maik Stamm. 2010. Tactile music instrument recognition for audio mixers. In *Audio Engineering Society Convention 128*. Audio Engineering Society.
- [58] Meta. 2023. Haptics Studio | Oculus Developers. <https://developer.oculus.com/resources/haptics-studio/>. Accessed August, 2025.
- [59] Birger Moell, Jim O’Regan, Shivam Mehta, Ambika Kirkland, Harm Lameris, Joakim Gustafsson, and Jonas Beskow. 2022. Speech data augmentation for improving phoneme transcriptions of aphasic speech using wav2vec 2.0 for the psst challenge. In *4th RAPID Workshop: Resources and Processing of Linguistic, Para-Linguistic and Extra-Linguistic Data from People with Various Forms of Cognitive/Psychiatric/Developmental Impairments, RAPID 2022, Marseille, France, Jun 25 2022*. 62–70.
- [60] Momoka Nakayama, Risako Kawashima, Shintaro Murakami, Yuta Takeuchi, Tatsuya Mori, and Dai Takanashi. 2024. A Method for Generating Tactile Sensations from Textual Descriptions Using Generative AI. In *SIGGRAPH Asia 2024 Posters*. 1–2.
- [61] Seoyong Nam, Minh Chung, Haerim Kim, Eunhae Kim, Taehyeon Kim, and Yongjae Yoo. 2024. Automatic Generation of Multimodal 4D Effects for Immersive Video Watching Experiences. In *SIGGRAPH Asia 2024 Technical Communications*. 1–4.
- [62] Gregory Neely and Lage Burström. 2006. Gender differences in subjective responses to hand–arm vibration. *International journal of industrial ergonomics* 36, 2 (2006), 135–140.
- [63] Aditya Arie Nugraha, Antoine Liutkus, and Emmanuel Vincent. 2016. Multichannel audio source separation with deep neural networks. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 24, 9 (2016), 1652–1664.
- [64] Valeria Ocelli, Charles Spence, and Massimiliano Zampini. 2011. Audiotactile interactions in temporal perception. *Psychonomic bulletin & review* 18, 3 (2011), 429–454.
- [65] Ryuta Okazaki, Hidenori Kuribayashi, and Hiroyuki Kajimoto. 2015. The effect of frequency shifting on audio–tactile conversion for enriching musical experience. *Haptic Interaction: Perception, Devices and Applications* (2015), 45–51.
- [66] Jaclyn Packer, Katie Vizenor, and Joshua A Miele. 2015. An overview of video description: history, benefits, and guidelines. *Journal of Visual Impairment & Blindness* 109, 2 (2015), 83–93.
- [67] Gunhyuk Park and Seungmoon Choi. 2011. Perceptual space of amplitude-modulated vibrotactile stimuli. In *2011 IEEE world haptics conference*. IEEE, 59–64.
- [68] Ryan M Peters, Erik Hackeman, and Daniel Goldreich. 2009. Diminutive digits discern delicate details: fingertip size and the sex difference in tactile spatial acuity. *Journal of Neuroscience* 29, 50 (2009), 15756–15761.
- [69] Karol J. Piczak. 2015. ESC: Dataset for Environmental Sound Classification. In *Proceedings of the 23rd Annual ACM Conference on Multimedia* (Brisbane, Australia, 2015-10-13). ACM Press, 1015–1018. doi:10.1145/2733373.2806390
- [70] PlayStation. 2020. Astro’s Playroom. <https://www.playstation.com/en-us/games/astros-playroom>. Retrieved August, 2025.
- [71] Zuzanna Podwinska, Iwona Sobieraj, Bruno M Fazenda, William J Davies, and Mark D Plumbley. 2019. Acoustic event detection from weakly labeled data using auditory salience. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 41–45.

- [72] Qiaoqiao Ren and Tony Belpaeme. 2025. Touched by chatgpt: Using an llm to drive affective tactile interaction. In *2025 20th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 1563–1567.
- [73] Kimberly Johanna Schelle, Carolina Gomez Naranjo, Martijn Ten Bhömer, Oscar Tomico, and Stephan Wensveen. 2015. Tactile dialogues: Personalization of vibrotactile behavior to trigger interpersonal communication. In *Proceedings of the ninth international conference on tangible, embedded, and embodied interaction*. 637–642.
- [74] Jonathan Isaac Segal, Samuel Rodriguez, Akshaya Raghavan, Heysil Baez, Crescentia Jung, Jazmin Collins, Shiri Azenkot, and Andrea Stevenson Won. 2024. SocialCueSwitch: Towards Customizable Accessibility by Representing Social Cues in Multiple Senses. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–7.
- [75] Nataliya Shakhovska and Ivan Zagorodniy. 2024. Classification of Acoustic Tones and Cardiac Murmurs Based on Digital Signal Analysis Leveraging Machine Learning Methods. *Computation* 12, 10 (2024), 208.
- [76] Tianzheng Shi and Oliver Schneider. 2025. Development and Initial Validation of the Haptic Experience Inventory (HXI). In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–21.
- [77] Francesca Sorgini, Renato Calì, Maria Chiara Carrozza, and Calogero Maria Oddo. 2018. Haptic-assistive technologies for audition and vision sensory disabilities. *Disability and Rehabilitation: Assistive Technology* 13, 4 (2018), 394–421.
- [78] Matti Strese, Clemens Schuwerk, Albert Iepure, and Eckehard Steinbach. 2016. Multimodal feature-based surface material classification. *IEEE transactions on haptics* 10, 2 (2016), 226–239.
- [79] Maciej Stroinski, Kamil Kwarciak, Mateusz Kowalewski, Daria Hemmerling, William Frier, and Orestis Georgiou. 2025. Text-to-Haptics: Enhancing Multi-sensory Storytelling through Emotionally Congruent Midair Haptics. *Advanced Intelligent Systems* 7, 4 (2025), 2400758.
- [80] Youjin Sung, Kevin John, Sang Ho Yoon, and Hasti Seifi. 2025. HapticGen: Generative Text-to-Vibration Model for Streamlining Haptic Design. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–24.
- [81] Ana Tajadura-Jiménez, Maria Basia, Ophelia Deroy, Merle Fairhurst, Nicolai Marquardt, and Nadia Bianchi-Berthouze. 2015. As light as your footsteps: altering walking sounds to change perceived body weight, emotional state and gait. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*. 2943–2952.
- [82] Ana Tajadura-Jiménez, Merle T Fairhurst, and Ophelia Deroy. 2022. Sensing the body through sound. In *The Routledge Handbook of Bodily Awareness*. Routledge, 230–246.
- [83] Yuji Tokozume, Yoshitaka Ushiku, and Tatsuya Harada. 2017. Learning from between-class examples for deep sound recognition. *arXiv preprint arXiv:1711.10282* (2017).
- [84] Yusuke Ujitoko and Yuki Ban. 2018. Vibrotactile signal generation from texture images or attributes using generative adversarial network. In *International conference on human haptic sensing and touch enabled computer applications*. Springer, 25–36.
- [85] Yusuke Ujitoko, Yuki Ban, and Koichi Hirota. 2020. GAN-based fine-tuning of vibrotactile signals to render material surfaces. *IEEE Access* 8 (2020), 16656–16661.
- [86] Shafiq ur Rehman, Muhammad Sikandar Lal Khan, Liu Li, and Haibo Li. 2014. Vibrotactile TV for immersive experience. In *Signal and Information Processing Association Annual Summit and Conference (APSIPA), 2014 Asia-Pacific*. IEEE, 1–4.
- [87] Pantelis N Vassilakis and K Fitz. 2007. SRA: A web-based research tool for spectral and roughness analysis of sound signals. In *Proceedings of the 4th Sound and Music Computing (SMC) Conference*. 319–325.
- [88] Kentaro Yoshida, Seki Inoue, Yasutoshi Makino, and Hiroyuki Shinoda. 2017. VibVid: Vibration estimation from video by using neural network. In *Proceedings of the 27th International Conference on Artificial Reality and Telexistence and 22nd Eurographics Symposium on Virtual Environments*. 37–44.
- [89] Gyeore Yun and Seungmoon Choi. 2025. Real-time Semantic Full-Body Haptic Feedback Converted from Sound for Virtual Reality Gameplay. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–17.
- [90] Gyeore Yun, Hyoseung Lee, Sangyoon Han, and Seungmoon Choi. 2021. Improving viewing experiences of first-person shooter gameplays with automatically-generated motion effects. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–14.
- [91] Gyeore Yun, Minjae Mun, Jungeun Lee, Dong-Geun Kim, Hong Z Tan, and Seungmoon Choi. 2023. Generating real-time, selective, and multimodal haptic effects from sound for gaming experience enhancement. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–17.
- [92] Massimiliano Zampini, Timothy Brown, David I Shore, Angelo Maravita, Brigitte Röder, and Charles Spence. 2005. Audiotactile temporal order judgments. *Acta psychologica* 118, 3 (2005), 277–291.
- [93] Yan Zhan, Xiaoying Sun, Qinglong Wang, and Weizhi Nai. 2023. Method for audio-to-tactile cross-modality generation based on residual U-net. *IEEE Transactions on Instrumentation and Measurement* 73 (2023), 1–14.
- [94] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412* (2017).