

# Strata: Fine-Grained Acoustic-based Device-Free Tracking

Sangki Yun<sup>\*</sup>, Yi-Chao Chen, Huihunag Zheng, Lili Qiu, and Wenguang Mao  
The University of Texas at Austin  
{sangki, yichao, huihuang, lili, wmao}@cs.utexas.edu

## ABSTRACT

Next generation devices, such as virtual reality (VR), augmented reality (AR), and smart appliances, demand a simple and intuitive way for users to interact with them. To address such needs, we develop a novel acoustic based device-free tracking system, called Strata, to enable a user to interact with a nearby device by simply moving his finger. In Strata, a mobile (*e.g.*, smartphone) transmits known audio signals at inaudible frequency, and analyzes the received signal reflected by the moving finger to track the finger location. To explicitly take into account multipath propagation, the mobile estimates the channel impulse response (CIR), which characterizes signal traversal paths with different delays. Each channel tap corresponds to the multipath effects within a certain delay range. The mobile selects the channel tap corresponding to the finger movement and extracts the phase change of the selected tap to accurately estimate the distance change of a finger. Moreover, it estimates the absolute distance of the finger based on the change in CIR using a novel optimization framework. We then combine the absolute and relative distance estimates to accurately track the moving target. We implement our tracking system on Samsung Galaxy S4 mobile phone. Through micro-benchmarks and user studies, we show that our system achieves high tracking accuracy and low latency without extra hardware.

## Keywords

Acoustic Tracking; Gesture Recognition; Channel Impulse Response

## 1. INTRODUCTION

**Motivation:** Smart appliances, Virtual Reality (VR), and Augmented Reality (AR) are all taking off. The availability of easy-to-use user interface is the key to their success. Smart TVs are still cumbersome to navigate through menus. Many smart appliances require users to manually launch smartphone applications and click through, which is even more cumbersome than actually turning on/off switches. VR and AR are expected to hit \$150 billion

<sup>\*</sup>Sangki Yun is currently a research staff at Hewlett Packard Labs.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MobiSys'17, June 19-23, 2017, Niagara Falls, NY, USA

© 2017 ACM. ISBN 978-1-4503-4928-4/17/06...\$15.00

DOI: <http://dx.doi.org/10.1145/3081333.3081356>

by 2020. They provide immersive experience, and open the doors to new ways of training, education, meeting, advertising, travel, health care, emergency responses, and scientific experiments. However, the current user interfaces of VR/AR are rather limited. Our vision is to develop a device-free user interface (UI) so that a user can freely move his or her hand to control game consoles, VR, AR, and smart appliances.

Directly tracking hand or finger movement without any device is appealing due to convenience. We have interviewed VR, AR, game developers and users. They all prefer device-free based user interface (UI) (*i.e.*, controlling devices directly using hands without holding anything) since it is cumbersome to hold a device outside the view and it is more natural to play games and interact with VR/AR objects using hands directly.

**Challenges:** There have been considerable work on enhancing the user interface of VR/AR devices. Google Daydream [8] and HTC VIVE [37] provide controllers for motion tracking. Device-free tracking is more convenient for VR users. Microsoft HoloLens [11] enables vision based device-free gesture tracking using camera and depth sensor fusion. However, it costs more than \$3,000 in the market. Vision based tracking incurs significant energy and computational cost [30]. Also, its tracking accuracy highly depends on the background color and lighting condition.

Besides vision techniques, radio and acoustic signals based object tracking has received significant research attention. However, enabling accurate device-free tracking using radio and acoustic signals is particularly challenging. In such a case, reflected signal has to be used for tracking. Reflected signal is much weaker than directly received signal (*e.g.*, in free space, the directly received signal attenuates by  $1/d^2$  whereas the reflected signal attenuates by  $1/d^4$ , where  $d$  is the distance between the device and target to be tracked). Moreover, it is more difficult to handle multiple reflection paths in device-free tracking. In device-based tracking, one may rely on finding the first arriving signal since the straight-line path between the sender and receiver is shortest. In comparison, in device-free tracking, the path of interest is not the shortest, which makes it even harder to distinguish which path should be used for tracking.

In particular, there has been considerable amount of work on device-free motion tracking using RF signals (*e.g.*, [13, 2, 41, 32]) and acoustic signals (*e.g.*, [21, 40]). However, they either require specialized hardware, or provide insufficient accuracy. For example, WiDeo [13] is a recent WiFi-based tracking scheme and pushes the tracking accuracy to within a few centimeters, but it requires full-duplex wireless hardware. WiTrack [2] tracks user position with 10-13 cm error by using customized hardware that transmits chirp signal through 1.67 GHz bandwidth. WiDraw [32] achieves 5 cm tracking accuracy with the support of 25 WiFi APs around

the user. Therefore, they require complicated hardware while providing insufficient accuracy to enable motion-based UI for VR/AR users. Millimeter-wave based tracking schemes, such as Google Soli [17] or mTrack [41] achieve higher accuracy, but require 60 GHz antenna arrays that are not widely available yet. The acoustic based tracking schemes, such as AAMouse [44] or CAT [18], achieve around 1 cm or lower median error on commodity devices. However, they are device based schemes and assume the first arriving signal should be used for tracking, which does not hold for reflected signals. RF-IDraw [38] achieves high tracking accuracy, but it requires the user to hold the RFID tag to be tracked.

Recently, two state-of-the-art acoustic tracking schemes (*i.e.*, FingerIO [21] and LLAP [40]) enable device-free tracking only using the existing speaker and microphones in mobile device. Although they achieve higher accuracy than the other schemes, their tracking accuracy is still limited due to multi-path and other movements. For example, based on our extensive experiment, LLAP [40] achieves 0.7 cm distance estimation error in 1D tracking, but its trajectory error in 2D space increases to 1.9 cm. When there are people moving around the tracking object, the accuracy further degrades. The accuracy of FingerIO is even lower than LLAP.

Therefore, despite significant progress, achieving highly accurate and responsive device-free tracking on commodity hardware remains an open challenge. According to personal communication with game and application developers, sub-centimeter level accuracy and within 16 ms response time are required in order to provide good user experience. This is especially challenging to achieve using a commodity device, such as a smartphone, given its limited processing power and lack of special hardware.

**Our approach:** Built on the existing work, we develop a new device-free tracking using acoustic signals. In our system, a mobile device (*e.g.*, smartphone) continuously transmits inaudible acoustic signals. The signals are reflected by nearby objects, including a moving finger, and arrive at the microphone on the same mobile. The mobile analyzes the received signal to estimate the channel, based on which it estimates the distance change and absolute distance to locate the finger. Due to the small wave-length of acoustic signals, it is promising to derive the distance change based on the phase. Phase is also more robust to imperfect frequency response of a speaker. However, like many wireless signals, audio signals go through multiple paths to reach the receiver (*e.g.*, due to reflection by different objects). Such multipath propagation poses significant challenges for phase-based tracking. To address the challenge, we estimate channel impulse response (CIR) in the time-domain. The estimate gives the channel coefficient of each channel tap. We then select an appropriate channel tap and use the phase of the selected tap to estimate the distance change of a finger.

To further derive the absolute distance, we develop a novel framework to estimate the absolute distance of the path reflected by the moving finger during a few consecutive intervals such that its changes match with the changes in the CIR during these intervals and the distance changes between the intervals match with the phase measurement. Inferring the absolute distance serves two purposes: (i) it allows us to get the initial absolute distance so that we can translate the subsequent distance change into a new absolute distance, and (ii) it can be used to improve the tracking accuracy and alleviate error accumulation in subsequent intervals by combining it with the relative distance change.

We implement our approach on Samsung Galaxy S4 mobile phone, which has one speaker and two microphones, and enable real-time tracking of the user's moving finger. Using extensive evaluation, we show our system has three distinct features: (i) high accuracy: within 0.3 cm distance tracking error, 1.0 cm 2D tracking error, and

0.6 cm drawing error in a 2D space; (ii) low latency: we can update the position every 12.5ms, and (iii) easy to deploy: with a software app installation, a smartphone can track a nearby finger movement without extra hardware.

**Paper outline:** The rest of this paper is organized as follows. We describe our approach in Section 2, and evaluate its performance in Section 3. We review related work in Section 4, and conclude in Section 5.

## 2. OUR APPROACH

In this section, we present a fine-grained acoustic-based device-free tracking.

### 2.1 Overview

We use the phase change of the acoustic channel to estimate the distance change. This allows us to achieve high accuracy because the acoustic wavelength is very short. For example, the wavelength is 1.9 cm in 18 KHz audio frequency. Only 1 mm movement causes the reflected path length to change by 2 mm, which results in  $0.21\pi$  phase change, large enough to detect.

However, in practice, due to multi-path propagation (*i.e.*, a signal traverses multiple paths before arriving at the receiver), the impact of a moving target on the overall channel can be very complicated, and varies across environments. In order to address this challenge, we use the phase from the estimated channel impulse response (CIR) rather than the raw received signal. CIR is a characterization of all signal traversal paths with different delays and magnitudes [27]. Specifically, it is a vector of channel taps where each channel tap corresponds to multi-path effects within a specific delay range. By focusing on the phase change of certain channel taps whose delays are close to the target delay range, we can effectively filter out the phase change incurred by the movement of objects outside a certain range as determined by the number of taps being used.

The phase change only gives us the distance change. We need to know the absolute distance at some point in order to translate the distance change into an absolute distance for tracking. Moreover, using the distance change alone incurs error accumulation, since the distance at a given time is estimated as the sum of all previous distance changes plus the initial position, each of which may incur an error. To address both issues, we develop a technique to estimate the absolute distance, which is used to get the initial position and also enhance the tracking accuracy by combining it with the distance change over time.

Putting together, our overall system consists of the following steps, which we will elaborate in this section.

1. Estimate channel impulse response (CIR) (Section 2.3);
2. Identify the channel tap corresponding to the target, and track the phase change of the selected tap in CIR to estimate the distance change (Section 2.4);
3. Estimate the absolute distance based on CIR (Section 2.5);
4. Combine the absolute distance with the relative distance to get a more accurate distance estimate, and track the target's position based on the distance to different landmarks (*e.g.*, microphones) (Section 2.6).

Our system implementation uses 18-22 KHz frequency and 48 KHz sampling rate. We can easily support other bandwidths and sampling rates.

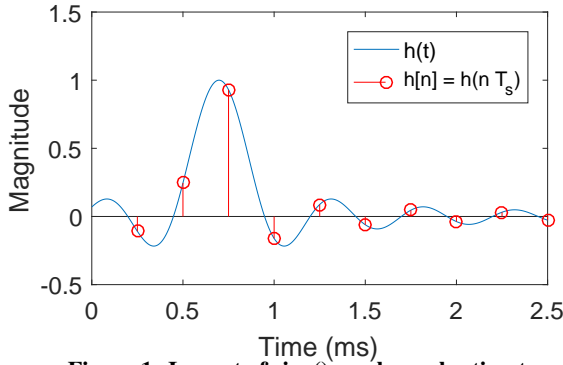


Figure 1: Impact of sinc() on channel estimate.

## 2.2 Background

A wireless signal, including an audio signal, travels through a straight line from the transmitter to the receiver in free space. In reality, due to obstacles in the environment, a single transmitted signal will reach the receiver via multiple paths (e.g., paths going through different reflectors). Therefore, the received signal is a superposition of multiple signals with different delays. The received signal via multipath is traditionally modeled as the following Linear Time-Invariant (LTI) system [35]. Suppose the channel has  $L$  paths and the received signal from path  $i$  has delay  $\tau_i$  and amplitude  $a_i$  determined by the travel distance of the path and reflectors. Then, the received signal  $y(t)$  is the summation of  $L$  signals, as shown below:

$$y(t) = \sum_{i=1}^L a_i x(t - \tau_i) = \sum_{i=1}^L a_i e^{-j2\pi f_c \tau_i} s(t - \tau_i) = h(t) * x(t), \quad (1)$$

where  $s(t)$  and  $x(t)$  are the transmitted baseband and passband signals at time  $t$ , respectively, and  $h(t)$  is the channel impulse response.  $h(t) = \sum_{i=1}^L a_i e^{-j2\pi f_c \tau_i} \delta(t - \tau_i)$ , where  $\delta(t)$  is Dirac's delta function [22].

The channel estimate from the received baseband symbols is a discrete output of  $h(t)$  sampled every  $T_s$  [35], which is

$$h[n] = \sum_{i=1}^L a_i e^{-j2\pi f_c \tau_i} \text{sinc}(n - \tau_i W), \quad (2)$$

where  $\text{sinc}(t) = \frac{\sin(\pi t)}{\pi t}$ . Conventionally,  $h[n]$  is called the  $n$ -th channel tap, because CIR is regarded as a discrete-time filter in LTI system. Note that sinc function decays over time, so the impact of delayed signal on the measured  $h[n]$  is small when the difference between  $nT_s$  and  $\tau_i$  are sufficiently large. However, if they are relatively close, movement is captured in multiple channel taps due to the signal dispersion effect of *sinc* function. This is illustrated in simulation result in Figure 1, where the channel is affected by one reflected signal at 30 cm (i.e.,  $\tau = 0.697$  ms). From the figure, we can observe that not only the closest channel tap from  $\tau$  (i.e.,  $h[3]$ ), but also the other nearby channels are affected by a single reflected signal.

## 2.3 Estimating Channel Impulse Response

**Single-carrier communication channel:** To estimate the channel, we design data communication for acoustic channel. An important design decision is whether we should use single-carrier or multi-carrier (e.g., OFDM) communication. OFDM is widely used in modern wireless communication due to its efficiency and robustness to Inter Symbol Interference (ISI) caused by multipath. However, it yields channel estimation in frequency domain, while chan-

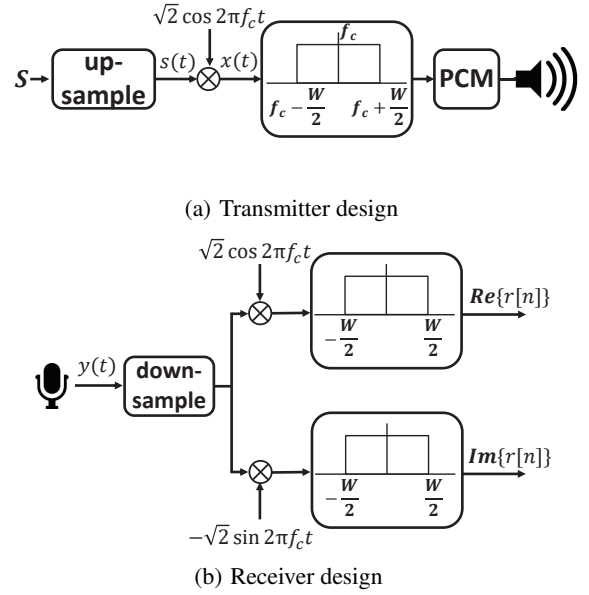


Figure 2: Transmitter and Receiver system diagram.

nel estimate is often more useful in time domain for tracking and localization purpose [19]. So we need to transform channel coefficients from frequency domain to time domain, but this process requires additional computation due to IFFT. Therefore, we use a single-carrier based communication system to directly estimate the channel in time domain without extra processing.

**Transmission signal design:** A transmitter sends a known training sequence for channel estimation. Let  $\mathbf{S} = \{s[1], \dots, s[K]\}$  denote the training sequence, where  $K$  is the length of the sequence. It can be any random bits. We choose 26-bit GSM training sequence because it is known to have good properties for synchronization and channel estimation [25] and widely used in single carrier communication. We modulate  $\mathbf{S}$  to BPSK symbols, where bits 0 and 1 are mapped to baseband symbols 1 and -1, respectively.

Figure 2(a) illustrates signal generation and transmission process. To transmit a modulated symbol over the inaudible frequency band, we first need to reduce the signal bandwidth so that it does not exceed the maximum allowed bandwidth of the inaudible band. Let  $f_s$  and  $B$  denote the sampling rate and the channel bandwidth, respectively. To limit the bandwidth of the transmitted symbol, we upsample the symbol at a rate of  $\frac{f_s}{B}$ , which is done by zero padding and low-pass filtering to smooth discontinuity [22]. Finally, we up-convert the signal to transmit it over the inaudible band. Let  $f_c$  denote the center frequency of the passband. We change the frequency of the signal by multiplying  $\sqrt{2} \cos(2\pi f_c t)$  to the baseband signal:  $x(t) = \sqrt{2} \cos(2\pi f_c t) s(t)$ , where  $s(t)$  and  $x(t)$  are upsampled baseband and passband signals, respectively. Since BPSK only has real parts,  $x(t) = \sqrt{2} e^{-j2\pi f_c t} s(t)$ .

To remove noise outside the transmission band, we perform band-pass filtering on  $x(t)$  with pass-band from  $f_c - \frac{B}{2}$  Hz to  $f_c + \frac{B}{2}$  Hz. The generated passband signals are transmitted through the smartphone speaker. Since the transmitted training sequence is always fixed, it can be generated offline and saved as a format of 16-bit Pulse Coded Modulation (PCM) in a Waveform Audio (WAV) file, which can be played by any mobile device that supports it (e.g., smartphone or smart watch).

We refer to the training sequence as a frame. Between frames, we insert a fixed gap (i.e., zero symbols) to avoid inter-frame inter-

ference. The gap should be sufficiently long so that the delayed signal from the previous frame does not interfere with the new frame. However, it should be as short as possible to provide low latency. Our study shows 24 zero symbols between frames are sufficient. As a result, a frame has 50 symbols. Given that the baseband symbol interval is  $\frac{1}{B} = 0.25$  ms, each frame lasts 12.5 ms. So we update new channel estimate and the target's position every 12.5ms, which is below 16 ms required for providing seamless user experience.

**Receiver design:** Figure 2(b) illustrates the signal reception and baseband conversion process. The received passband signal  $y(t)$  arriving at the microphone is converted into a baseband symbol  $r[n]$  using the following down-conversion process:  $y(t)$  is multiplied by  $\sqrt{2}\cos(2\pi f_c t)$  and  $-\sqrt{2}\sin(2\pi f_c t)$  to get the real and imaginary parts of the received baseband symbol, respectively. We perform low-pass filtering and down-sampling to select a signal every symbol interval. This gives us the following baseband signal<sup>1</sup>:

$$\begin{aligned} r[n] &= \sqrt{2}\cos(2\pi f_c t)y(t) - j\sqrt{2}\sin(2\pi f_c t)y(t) \\ &= \sqrt{2}e^{-j2\pi f_c t}y(t), \end{aligned}$$

where  $t$  is the time that the  $n$ -th baseband symbol is sampled (*i.e.*,  $t = n \times T_s$ , where  $T_s$  is a symbol interval).

**Frame detection:** After passband-to-baseband signal conversion, the receiver detects the first symbol of the received frame by energy detection and cross-correlation. We first detect the rough beginning of the frame based on energy detection: if the magnitude of three consecutive symbols is higher than the threshold  $\sigma$ , we treat the first symbol as the beginning of the frame symbols. Our implementation uses  $\sigma = 0.003$ . This value needs to be carefully selected depending on the phone and the volume of the speaker. Then we find more precise starting point based on cross-correlation. Specifically, we find the sample that gives the maximum cross-correlation magnitude between the received and transmitted frames. Note the frame detection procedure is only necessary at the beginning of tracking. Since the frame interval is fixed, once a frame is detected, the subsequent frames can be determined by adding a constant frame interval.

**Channel estimation:** Next we estimate the channel based on the received frame and the known training sequence. There are several existing channel estimation algorithms in single carrier communication system. We use least-Square (LS) channel estimation since it is cheap to compute on a mobile. We mainly focus on the algorithm implementation, and refer readers to [25] for the fundamental theory behind it.

For LS channel estimation, one needs to decide the reference length  $P$  and the memory length  $L$ , where  $L$  determines the number of channel taps we can estimate and  $P + L$  is the training sequence length. Increasing  $L$  allows us to estimate more channel taps but reduces the reliability of estimation. Our implementation uses  $P = 16$  and  $L = 10$ , which implies we can track movement up to 50 cm away (see Section 2.4.1). One can easily adapt  $P$  according to the environment.

Let  $\mathbf{m} = \{m_1, m_2, \dots, m_{L+P}\}$  denote the training sequence.

A circulant training matrix  $\mathbf{M} \in \mathbb{R}^{P \times L}$  is:

$$\mathbf{M} = \begin{bmatrix} m_L & m_{L-1} & m_{L-2} & \dots & m_1 \\ m_{L+1} & m_L & m_{L-1} & \dots & m_2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ m_{L+P} & m_{L+P-1} & m_{L+P-2} & \dots & m_{P+1} \end{bmatrix}.$$

Let  $\mathbf{y} = \{y_1, y_2, \dots, y_{L+P}\}$  denote the received training sequence. The channel is estimated as

$$\hat{\mathbf{h}} = (\mathbf{M}^H \mathbf{M})^{-1} \mathbf{M}^H \mathbf{y}_L, \quad (3)$$

where  $\mathbf{y}_L = \{y_{L+1}, y_{L+2}, \dots, y_{L+P}\}$ .

Given the pre-computed  $(\mathbf{M}^H \mathbf{M})^{-1} \mathbf{M}^H$ , the computational cost of the channel estimation is only the matrix-to-vector multiplication, which is  $O(P \times L)$ . Considering  $P = 16$  and  $L = 10$ , the channel estimation complexity is low enough to implement on a mobile.

To further improve the channel estimation accuracy, unlike in the traditional digital communication, which picks one out of every  $r$  samples during downsampling, where  $r$  is the upsampling rate, we use the average of the first  $l$  in the  $r$  samples to estimate the channel every sampling interval.

## 2.4 Tracking Phase Change

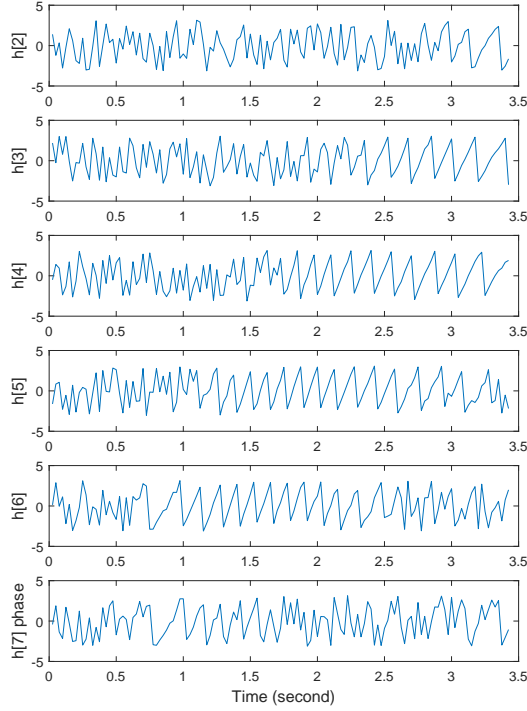
### 2.4.1 Overview

Next we track the phase change based on CIR estimates. We study the impact of the reflected signal using the following experiment. We move a small aluminum ball ( $< 1$  cm diameter) attached to a long and thin wood stick. The person moving the stick is over 1 m away from the ball (outside the range of CIR taps). The ball is initially 30 cm away from the smartphone and moved 20 cm towards the phone. Figure 3 shows the phase of multiple channel taps while the ball is moving towards the smartphone. How to observe the phase of the moving object will be explained in detail in Section 2.4.2. The result shows that the phase rotates in multiple taps, which indicates the path length is changing. This change is caused by the moving ball. As shown in Figure 3, even though only a single small object moves, the phase rotation is observed in multiple taps. We repeat the experiments multiple times, and find that the phase rotation is observed approximately in 3 consecutive taps and the reflected signal with delay  $\tau$  affects the three  $h[n]$  that have the smallest  $|\tau - nT_s|$ . As a result,  $h[n]$  can be approximated as:

$$h[n] \approx \sum_k a_k e^{-j2\pi f_c \tau_k}, \quad (n - \frac{3}{2})T_s < \tau_k < (n + \frac{3}{2})T_s. \quad (4)$$

In other words, each channel tap  $h[n]$  contains the phase and magnitude of the reflected signals whose delays are between  $(n - \frac{3}{2})T_s$  and  $(n + \frac{3}{2})T_s$ . The path length changes with the delay  $\tau_k$  according to  $d_k = \tau_k V_c$ , where  $d_k$  is the travel distance of the path  $k$  and  $V_c$  is the propagation speed of the audio (*i.e.*,  $V_c \approx 340$  m/s). Assuming that the speaker and microphone are closely located, the distance from the microphone to the reflecting object is approximately half of the travel distance. Therefore,  $h[n]$  indicates the object's distance from microphone is between  $(n - \frac{3}{2})\frac{T_s V_c}{2}$  and  $(n + \frac{3}{2})\frac{T_s V_c}{2}$ . Given  $T_s = 0.25$  ms,  $\frac{T_s V_c}{2} = 4.25$  cm and each tap captures objects across 12.75 cm range. This enables us to filter out the movement of objects outside the target range. For example, if we want to track the movement of a finger within 50 cm from the mobile, we can limit the channel taps to the first 10 taps to filter out the movement outside 50 cm. This is because the 10th tap may contain information from objects up to around 12th taps away, which gives  $12 * 4.25 = 51$  cm.

<sup>1</sup>Conventionally,  $(\cdot)$  and  $[\cdot]$  notations are used to represent analog and digital signals, respectively. Here every signal is digital because a mobile app cannot access the analog signal. We use  $(\cdot)$  and  $[\cdot]$  notations to distinguish upsampled signal with rate  $f_s$  from the downsampled signal with rate  $B$ .



**Figure 3: Phase change in multiple channel taps while moving a ball.**

Next we track the phase change using the CIR estimate. While the CIR vector captures the channel with different propagation distances, it is challenging to extract the phase change caused by the target movement based on CIR since multiple paths with similar distances are mixed within each channel tap. To address the issue, we decompose the problem into the following two steps: (i) if we know which channel tap is affected by the moving target, how to extract the phase change caused by the target's movement, and (ii) how to determine which channel tap is affected by the target. Below we present our approaches to address both issues.

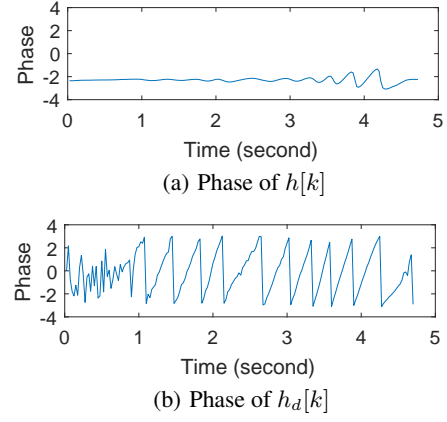
#### 2.4.2 Estimate Phase Change

We assume the  $k$ -th channel tap is affected by the target's movement. In order to observe the phase change of the moving target, we compare the two consecutive channel measurements. Taking difference between the two consecutive channels effectively removes dominant static reflections. Let  $L_k$  denote the number of paths observed in  $h[k]$ . Suppose the  $L_k$ -th path is the path reflected from the moving finger, while the other  $L_k - 1$  paths remain the same during the two consecutive channel measurement periods  $t - 1$  and  $t$ . Then,

$$h[k]^{t-1} = \sum_{i=1}^{L_k} a_i e^{-j2\pi f_c \tau_i(t-1)},$$

$$h[k]^t = \sum_{i=1}^{L_k-1} a_i e^{-j2\pi f_c \tau_i(t)} + a_{L_k} e^{-j2\pi f_c (\tau_{L_k}(t-1) + \tau_d(t))},$$

where  $h[k]^t$  is the  $k$ -th channel tap estimated from the  $t$ -th frame and  $\tau_d(t)$  is the delay difference caused by the target movement between the  $t$ -th and  $(t-1)$ -th frame intervals (i.e.,  $\tau_d(t) = \tau_{L_k}(t) -$



**Figure 4: Phase of the channel impulse responses while a finger is moving.**

$\tau_{L_k}(t-1)$ ). By taking their difference, we get

$$h_d[k]^t = a_{L_k} (e^{-j2\pi f_c (\tau_{L_k}(t-1) + \tau_d(t))} - e^{-j2\pi f_c \tau_{L_k}(t-1)}), \quad (5)$$

where  $h_d[k]^t = h[k]^t - h[k]^{t-1}$ . Equation 5 assumes that  $a_{L_k}$  associated with a propagation path is constant over two consecutive measurements due to a very small distance change in a 12.5 ms interval. From the angle of  $h_d[k]^t$ , we observe the phase rotation caused by the change of  $\tau_{L_k}(t)$ .

$$\begin{aligned} \angle(h_d[k]^t) &= \angle(e^{-j2\pi f_c (\tau_{L_k}(t-1) + \tau_d(t))} - e^{-j2\pi f_c \tau_{L_k}(t-1)}) \\ &= \angle(e^{-j2\pi f_c \tau_{L_k}(t-1)} (e^{-j2\pi f_c \tau_d(t)} - 1)) \\ &= \angle(e^{-j2\pi f_c \tau_{L_k}(t-1)}) + \frac{\angle(e^{-j2\pi f_c \tau_d(t)})}{2} + \frac{\pi}{2}, \quad (6) \end{aligned}$$

where  $\angle(X)$  is the phase of the complex number  $X$ . We can prove  $\angle(e^{-j2\pi a} - 1) = \frac{\angle(e^{-j2\pi a})}{2} + \frac{\pi}{2}$  geometrically. The proof is omitted in the interest of brevity. Equation 6 assumes  $\angle(h_d[k]^t)$  is smaller than  $\pi$ .

Figure 4 (a) and (b) show the phases of  $h[k]$  and  $h_d[k]$ , respectively, while a user is moving his finger towards the speaker and microphone. In the collected trace, we conjecture  $h[k]$  includes the finger movement related path between 1.0 second and 4.6 second. In Figure 4(a), the phase is very stable and the change by the finger movement is not clear because the majority portion of  $h[k]$  contains signals from the static paths. After removing the impact of the static paths, we can observe clear phase rotation due to the finger movement from  $h_d[k]$ .

From the phase difference between  $h_d[k]^{t+1}$  and  $h_d[k]^t$ , we get the phase rotation caused by the delay difference  $\angle(e^{-j2\pi f_c \tau_d(t)})$ , and eventually the travel distance of the finger during the measurement interval using the relation between the phase change. Note that  $\tau_d(t) = \tau_{L_k}(t) - \tau_{L_k}(t-1)$ . Using Equation 6, we represent the phase difference as

$$\begin{aligned} \angle(h_d[k]^{t+1}) - \angle(h_d[k]^t) &= \angle(e^{-j2\pi f_c \tau_{L_k}(t)}) \\ &\quad - \angle(e^{-j2\pi f_c \tau_{L_k}(t-1)}) + \frac{1}{2} (\angle(e^{-j2\pi f_c \tau_d(t+1)}) - \angle(e^{-j2\pi f_c \tau_d(t)})) \\ &= \angle(e^{-j2\pi f_c \tau_d(t)}) + \frac{1}{2} (\angle(e^{-j2\pi f_c \tau_d(t+1)}) - \angle(e^{-j2\pi f_c \tau_d(t)})). \end{aligned}$$

By solving the above equation, we can calculate  $\angle(e^{-j2\pi f_c \tau_d(t)})$ . Without prior, we can simply assume  $\tau_d(t+1) = \tau_d(t)$ . Once we get the phase rotation, we can calculate the distance change based on  $d^t = \lambda \times \angle(e^{-j2\pi f_c \tau_d(t)}) / 2\pi$ , where  $d^t$  is the distance

change of the dynamic path at time  $t$ , and  $\lambda$  is the wavelength of the audio signal. This relationship holds as long as the phase change is smaller than  $\pi$ , which holds for our finger movement speed and interval duration.

### 2.4.3 Finding Channel Tap Corresponding to the Target

Section 2.4 assumes we already know which tap to use for tracking the finger movement. This section describes how to find the right tap that includes the path reflected from the finger among multiple possible taps. Note that as mentioned in Section 2.4.1, the phase rotation by the finger movement is observed in multiple taps rather than in a single tap. Therefore, we just need to select one of these taps.

The channel taps can be classified as dynamic taps (*i.e.*, those that includes dynamic paths) and static taps (*i.e.*, those that do not). The right taps should be dynamic taps, since we are interested in tracking finger movement. If all taps are static taps, it means the finger does not move and its position does not need to be updated.

**Criterion 1:** Compared to the static taps, the dynamic taps have relatively larger variation of the channel over time. Therefore, we develop the following test to identify dynamic paths in the tap  $k$ :

$$M_1[k]^t = \frac{|h[k]^t - h[k]^{t-1}|}{|h[k]^t|} > \sigma_t,$$

which compares the normalized difference in the magnitude of two consecutive channels with a threshold  $\sigma_t$ . We use  $\sigma_t = 0.05$ .

**Criterion 2:** While the above condition distinguishes between the dynamic and static taps, the noise in the channel estimation might cause the classification error. Therefore, we add another criterion based on the following observation: the phase rotation of static tap  $k$ , denoted as  $h_d[k]$ , is very unstable because all static paths are removed during the differentiation process and the remaining value in  $h_d[k]$  may contain random noise. In comparison, if  $k$  is a dynamic tap, the phase rotation of  $h_d[k]$  is much more stable because  $h_d[k]$  includes the dynamic path and its phase change over the measurement interval is stable. This is evident from Figure 3 and 4, which show the phase changes when the dynamic path is not included in the channel tap.

Based on this observation, we develop the following criterion to select the final tap. We measure the stability of phase change, which is defined as the phase change difference over the last three measurements. Specifically, we find the maximum phase change over the three periods:  $M_2[k]^t = \max_{i=t-2, t-1, t} f(i)$ , where  $f(i) = |\angle(e^{-j2\pi f_c \tau_d^k(i)}) - \angle(e^{-j2\pi f_c \tau_d^k(i-1)})|$ . We select the tap with the smallest maximum phase change (*i.e.*, the smoothest tap) from all taps that satisfy the criterion 1 as the final tap.

## 2.5 Estimating Absolute Distance

So far, we have focused on tracking the distance change of the finger by observing the phase. We need an absolute initial distance at some point in order to get the distance over time. Therefore, we develop a method to estimate the absolute distance based on the channel coefficients.

How to accurately estimate the absolute distance based on the channel estimate is an open problem. WiDeo [13] and Chronos[36] cast this problem as non-linear optimization problems, which search for the parameters associated with each channel tap such that the sum across all taps matches the overall channel measurement. This is effective when the channel is sparse. [13, 36] show this approach achieves decimeter-level accuracy in WiFi. We have tried this approach in our context, and found the accuracy is poor when applied

to acoustic signals due to many multipaths, which results in many unknowns and a severely under-constrained system.

**Basic framework:** We develop a new formulation to address the under-constrained problem. It is motivated by the following important observation. We do not need to reconstruct complete channel profile in order to track a moving finger. Instead, we just need to reconstruct the delay and magnitude of the path that is reflected by the finger. Since a finger is small, it is reasonable to assume only one path is reflected by the finger. Therefore, we can take the difference between the two consecutive CIR estimates, which will cancel out all static paths and reduce the unknowns to the number of the parameters associated with the path that is reflected by the finger. Recall Equation 2 models the impact of reflection on the estimated channels. We take the difference between the two consecutive CIR, which removes all static paths and only keeps the path reflected by the moving finger:

$$h_d[n] = a(e^{-j2\pi f_c(\tau+\tau_d)} \text{sinc}(n - (\tau + \tau_d)W) - e^{-j2\pi f_c \tau} \text{sinc}(n - \tau W)), \quad (7)$$

where  $a$  and  $\tau$  are the amplitude and delay of the signal reflected from the finger, respectively<sup>2</sup>. Based on the measured  $h_d[n]$ , our goal is to find  $\tau$  and  $a$  that minimize the difference between the measured and estimated CIR change, where the estimated CIR is derived from the mathematical model using Equation 7.

To further improve the accuracy, we minimize the difference between the measured and estimated CIR change over multiple consecutive CIR measurements. Our implementation considers the past 3 CIR measurements (*i.e.*, the CIR change from  $t-2$  to  $t-1$  and from  $t-1$  to  $t$ ). We still compute the coordinate every interval except that we use the most recent 3 CIR measurements to construct the optimization problem. Therefore, we can improve the accuracy while maintaining 12.5 ms tracking interval.

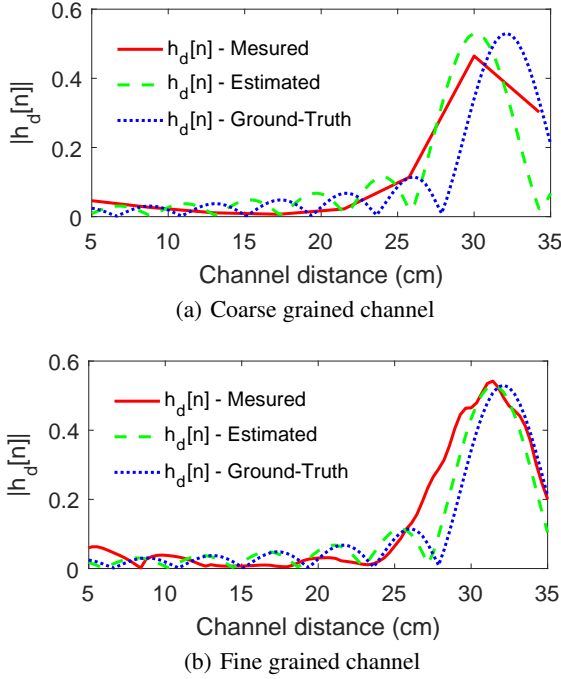
Putting together, we solve the following optimization problem:

$$\min_{\tau, \tau_d^{est}(i), a} : \sum_i \sum_{n=1}^L \left[ h_d^{t,t+i}[n] - a(e^{-j2\pi f_c(\tau+\tau_d^{est}(i))} \text{sinc}(n - (\tau + \tau_d^{est}(i))W) - e^{-j2\pi f_c \tau} \text{sinc}(n - \tau W)) \right]^2 + \alpha \sum_i |\tau_d^{est}(i) - \tau_d(i)| \quad (8)$$

where  $h_d^{t,t+i}[n]$  denotes the measured CIR change in the  $n$ -th tap from the  $t$ -th measurement to the  $t+i$ -th measurement,  $\tau_d^{est}(i)$  is the inferred delay change from the  $t$ -th measurement to the  $t+i$ -th measurement, and  $\tau_d(i)$  is the delay change derived from the phase measurement in Section 2.4.  $a$ ,  $\tau$ , and  $\tau_d^{est}(i)$  are unknowns. The first term in the objective captures the fitting error between the measured CIR change versus the CIR change derived from the absolute distance based on  $a$ ,  $\tau$  and  $\tau_d^{est}(i)$ , and the second term captures the fitting error between the inferred delay change versus the delay change measured from the phase.  $\alpha$  captures the relative weight between the two terms, and set to 100 in our evaluation due to more accurate distance change measurement.

Compared with the channel decomposition in [13, 36], which tries to find  $\tau$ 's and  $a$ 's associated with all paths in the channel, our scheme finds  $\tau$  and  $a$  associated with the path reflected by the moving finger. This has two benefits: (i) it reduces the number of unknowns and improves the accuracy while reducing computation cost, (ii) it removes all static paths and helps reduce the error.

<sup>2</sup>Note that Equation 5 is based on Equation 4, which is an approximation to Equation 2.



**Figure 5: Measured, re-generated, and ground-truth channel differences with coarse and fine grained channels based absolute distance estimation.**

Moreover, we leverage the information across multiple measurements to further enhance the accuracy. Once we get  $\tau$ , we can easily calculate the absolute distance as the product of delay ( $\tau$ ) and sound propagation speed.

**Enhancement:** While the above approach is useful to find the absolute distance, its accuracy is limited by the resolution of the estimated channel. Since the channel bandwidth is limited to 4 KHz, the baseband symbol interval is 0.25 ms, which is translated into the channel resolution of 4.25 cm. When we try to find  $\tau$  using  $h[n]$  sampled every 4.25 cm, the error increases due to the coarse resolution.

To enhance the accuracy, we exploit the over-sampled signals to achieve finer resolution. For the over-sampled signals received between  $h[k]$  and  $h[k+1]$ , we estimate the channel using the same method in Section 2.3. These samples are 0.0208 ms apart (*i.e.*,  $\frac{1}{F_s}$  ms), which corresponds to the distance resolution of 3.5 mm. As a result, the resolution of the channel estimation is limited not by the baseband symbol rate but by the sampling rate of the audio signal, which is 12 times higher! With this fine-grained channel estimation, we find  $\tau$  using the same optimization framework.

Figure 5 shows the impact of coarse-grained and fine-grained channel estimation on the absolute distance estimate. In this experiment, we record the audio signal while a user is moving his finger, and estimate the absolute distance using the measured channels. For ease of interpretation, we represent the x-axis of Figure 5 as the distance corresponding to the delay of the channel. The red lines in the figures show measured channel difference. The green and blue lines correspond to the channel differences by plugging in the estimated and ground-truth delays into Equation 7, respectively. The ground-truth finger distance is 32 cm. As we can see, the channel difference under the coarse channel estimation deviates from the ground-truth, and has 2 cm error. In comparison, the channel difference from the fine grained estimation is close to the ground-truth, and has only 0.4 cm error.

## 2.6 Tracking the moving finger in 2D space

In Section 2.4 – Section 2.5, we present an approach to estimate the distance to the target. This allows us to track in 1D space, which is already useful for some applications. Next we describe how to track a finger in 2D space by leveraging two microphones on a smartphone (*e.g.*, Samsung S series).

**Combine relative and absolute distance:** We use both absolute distance and distance change estimated from the phase to track the target. At the beginning, we use the absolute distance to get the initial position of the target. Afterwards, we can get two distance estimates:  $d_k^p$  estimated from the phase and  $d_k^c$  estimated from the channel difference, where  $d_k^p = d_{k-1} + \Delta d_k^p$  and  $\Delta d_k^p$  is computed as a function of the phase change. We then combine the two distance estimates using a weighting factor  $\beta$ :  $d_k = (1 - \beta)d_k^p + \beta d_k^c$ .  $\beta$  is set to 0.1 due to more accurate phase change measurement. In section 3.3.3, we extensively evaluate the various  $\beta$  values.

**Estimate coordinate:** Given the phone form factor, we know the relative locations of the speaker and microphones. Suppose the speaker is located at the origin (0,0), and the two microphones are located at  $(x_1, 0)$  and  $(x_2, 0)$ , respectively. We assume they are all aligned in the same Y-axis. The finger should be on the ellipse whose foci are  $(x_1, 0)$  and  $(x_2, 0)$  and the total distance to the foci are  $d_1$  and  $d_2$ , respectively. Using two microphones as landmarks, we can track the finger by finding the intersection of the two ellipses. There are multiple intersections between two ellipses when they overlap, and we select the one closer to the previous position.

## 3. PERFORMANCE EVALUATION

### 3.1 Experiment setup

To evaluate performance and conduct user study, we implement an Android app that processes the audio signal and tracks the finger movement in real-time. We use Samsung Galaxy S4 with Android 5.1.1 as a tracking device. It has one rear speaker and two microphones at the top and bottom of the phone with 14 cm separation. The mobile app plays audio file generated as explained in Section 2.3, and analyzes the audio signal received from the microphones to track the finger in real-time. The speaker's volume is set to 80% of the maximum. We use inaudible audio frequency between 18 KHz and 22 KHz for transmission, and set the tracking interval to 12.5 ms. To convert passband to baseband, we implement an infinite-impulse response (IIR) low-pass filter.

To collect the ground-truth of the finger movement, we let the user move the finger on the top of a smart tablet, shown in Figure 6. It collects the touch event of the screen and generates the ground-truth trajectory of the finger movement, which is compared with the position estimated by our approach. This setup is only needed for collecting the ground truth to quantify the accuracy, and our scheme lets users freely draw in the air. In user study, we show simple shapes, such as a triangle, diamond, and circle on the tablet screen, and ask the user to trace the shapes. The average distance of the shapes is 22.3 cm. For data collection and user study, we recruit 5 users: four are men and one is woman. They are college students between 21 to 29. This age-group is among the most active VR/AR users. The data is collected at a student office in the department building surrounded by desks, cubicles, walls and small objects nearby, and people are allowed to move around the tracking device.



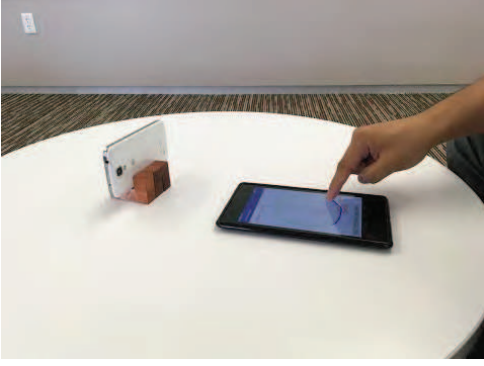


Figure 6: Testbed setup for the performance evaluation.

### 3.2 Baseline Schemes

We compare our scheme with the following two acoustic signal based device-free tracking schemes:

**Low Latency Acoustic Phase (LLAP):** LLAP [40] tracks the finger movement by observing the phase change of the reflected signal. In this scheme, the transmitter continuously sends sine waves and the receiver tracks the moving target based on the phase change of the received waves. We implemented LLAP closely following [40]. We generate the transmission signal  $\sin(2\pi f_c t)$  using MATLAB, stored as 16-bit PCM-format WAV file, and transmitted through the speaker. At the receiver side, the audio signal from the microphone is first multiplied by  $\sqrt{2} \cos(2\pi f_c t)$  and  $-\sqrt{2} \sin(2\pi f_c t)$ , and goes through a low-pass filter to get the real and imaginary parts of the received signal, respectively. The phase of the moving finger is tracked by Local Extreme Value Detection (LEVD) algorithm in [40] that filters out static signals and tracks the phase change by the finger movement. It improves the accuracy of the phase estimation by averaging multiple received signals. We take the phase after averaging all received signals during 12.5 ms so that both Strata and LLAP have the same tracking delay. Also, we used the multiple frequencies to address the multi-path fading as proposed in [40]. We did not implement its absolute distance estimation algorithm in LLAP for initialization. Instead, we used the ground-truth for the initial position estimation.

The main difference between Strata and LLAP is that the former separately measures the phase change of the signals with different delays while the latter measures the phase change caused by all of the surrounding objects. Since multiple body parts may move together while the user moves the finger, using the combined phase change yields less accurate finger tracking. The problem is even worse when there are other moving objects and people nearby, which can be common in practice. Our evaluation results of its 1D distance estimation error is similar to what is presented in [40]. Interestingly, [40] did not evaluate the 2D tracking error.

**Cross-correlation based tracking:** This scheme tracks a moving object by measuring the change in the cross-correlation of the received signal, called *echo profile* proposed by FingerIO [21]. It transmits an OFDM symbol every interval and locates the moving finger based on the change of the echo profiles of the two consecutive frames. Specifically, every interval it compares the difference between the echo profiles of the two consecutive frames and filters out the points whose differences are smaller than a pre-defined threshold. Among the remaining unfiltered points, it selects the point closest to the beginning of the interval since the direct path is the shortest path. We carefully followed the description of [21] in our implementation, but achieved much lower accuracy than [21] perhaps due to additional optimization used but omitted in [21]. We identified a few ways to improve [21]. First, we use larger band-

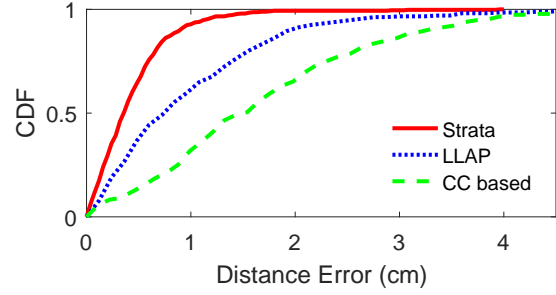


Figure 7: Impact of distance on tracking error.

width and longer FFT size for OFDM symbol, both of which help to improve the tracking accuracy. We select FFT size of 256 and 6 KHz bandwidth (16 – 22 KHz) for transmission while the other approaches including ours still use 4 KHz bandwidth. Second, we find the difference between the two consecutive profiles is small as the finger is moving, which makes it challenging to select an appropriate threshold for filtering. Instead, we select the point that has the maximum cross correlation difference. Our results show this helps to improve accuracy and we can now roughly track the position of the moving finger, but sometimes detect a random location due to noise. Therefore, we further filter out the positions that are too far from the previous position. Our evaluation picks the maximum peak that is within  $\pm 10$  cm away from the previous position. 10 cm is a loose bound to tolerate the error in estimating the previous position. After getting the distance estimation, we use the same algorithm introduced in Section 2.6 to locate the finger in a 2D plane.

### 3.3 Experimental Result

#### 3.3.1 Phase based Tracking

**Tracking accuracy in 1D:** We first evaluate the accuracy of estimating distance change. In this experiment, the user initially places the finger 20 cm away from the mobile phone, and moves 10 cm towards it. We repeat the experiment 200 times for each scheme and collect the CDF of the distance error. As shown in Figure 7, the median errors of Strata, LLAP, and cross-correlation based tracking (*i.e.*, CC based) are 0.3 cm, 0.7 cm, and 1.5 cm, respectively. Their 90th percentile errors are 0.8 cm, 1.8 cm, and 3.2 cm, respectively. The performance of LLAP is similar to the result presented in [40] when the finger is a moving object and the initial distance is 20 cm (*e.g.*, see Figure 11(c) in [40]). Note that in [40], the authors mostly used a hand rather than a finger to evaluate the tracking performance. Figure 7 indicates that Strata can accurately track the movement by observing the phase rotation from CIR. Its median tracking error is less than half of LLAP. Since Strata separately tracks the phase in the CIR vector, it can effectively filter out the measurement noise from the user’s body movement. In comparison, LLAP cannot since all signals are mixed up. The accuracy of the cross-correlation based tracking (even after optimization) is lower than both phase-based tracking schemes.

**Tracking accuracy in 1D with other moving objects:** One advantage of Strata is that it can distinguish movements at different distances. This allows us to easily filter out the interference incurred by the movement of the surrounding objects in the phase tracking. By limiting the channel taps to the first 10, we can ignore the phase change caused by the moving object with the distance larger than 50 cm.

To validate its effectiveness, we perform the distance tracking experiment with a person moving in the background while another



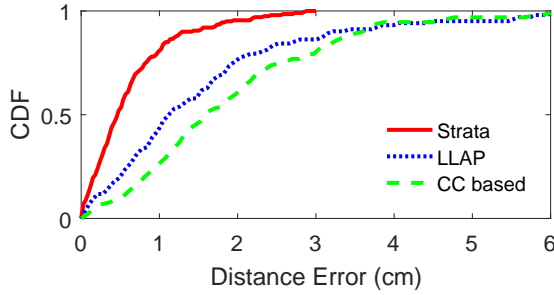


Figure 8: CDF of the distance tracking error with interrupting user.

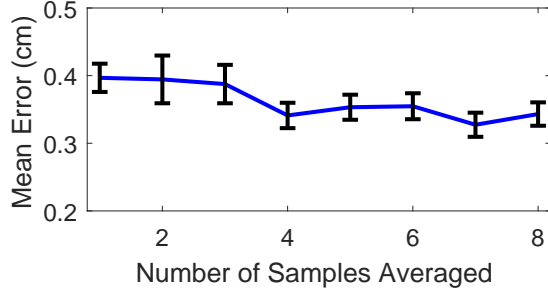


Figure 9: Tracking error with a varying number of samples averaged.

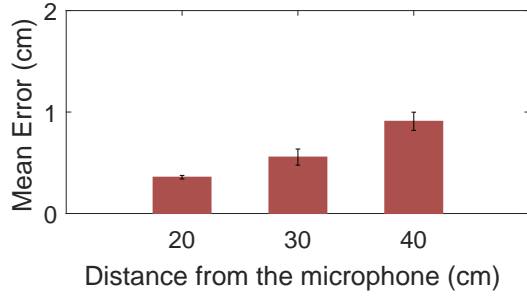


Figure 10: Impact of distance on tracking error.

user is moving his finger. The background user is approximately 60 cm away from the mobile phone. Figure 8 plots CDF of the distance tracking error. In Strata, the tracking accuracy is almost not affected by the background user: the median error increases by only 1 mm over no background moving user.

For the cross-correlation based tracking, we set it to focus on the phase change within the range between 0 to 40 cm. Even if there is a change in echo-profile in the distance of the moving user, we ignore it. As a result, it can effectively avoid the interference and achieves similar median tracking error as before: 1.6 cm. In comparison, LLAP incurs considerable degradation due to the moving user. The median error increases from 0.8 cm to 1.2 cm. Since it does not have a mechanism to distinguish the phase change caused by different objects, the background movement significantly degrades its tracking accuracy. LLAP [40] proposes to combine the phase tracking results of different frequencies to mitigate the multipath fading, but we find it is not sufficient to cancel out the large fading caused by moving people near the object to be tracked.

**Varying the number of samples:** Figure 9 is the average 1D tracking error with 95% confidence interval when various number of samples are used for averaging downsampled signals as explained in Section 2.3. The result shows that averaging 4 samples reduces the error from 0.39 cm to 0.35 cm compared to simple downsampling. Moreover, it reduces the variance. The 90th percentile error

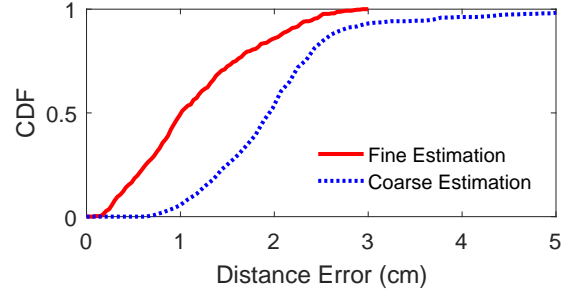


Figure 11: CDF of the absolute distance error.

decreases from 1.1 cm to 0.8 cm. Using more than 4 samples does not lead to significant additional improvement. So we use 4 samples to reduce the computation cost.

Besides the number of samples, there are a few parameters that should be carefully selected for accurate tracking: 1) the frame detection threshold  $\sigma$ , 2) the dynamic channel tap selection threshold  $\sigma_l$ , and 3) weighting factor  $\beta$ . We extensively address  $\beta$  in the next section, and skip detailed microbenchmark of  $\sigma$  and  $\sigma_l$ . These parameters require calibration on different phones, but are relatively easy to get after a few trials.

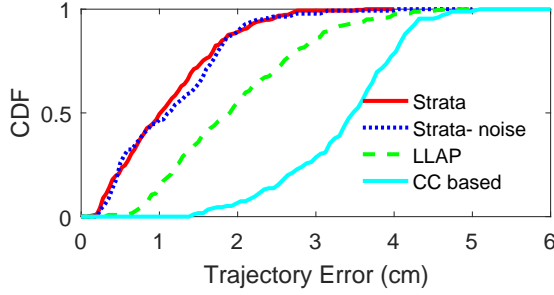
**Impact of distances:** We further vary the distance of the moving target. Specifically, we set the initial distance of the finger from the microphone to 20 cm, 30 cm, and 40 cm, move it 10 cm towards the mobile phone, and measure the distance error. Figure 10 shows the average error with 95% confidence interval. The mean tracking errors are 0.35 cm, 0.55 cm, and 0.9 cm, respectively, as the distance changes from 20 cm to 30 cm to 40 cm. In all cases, the average error is within 1 cm.

### 3.3.2 Estimating Absolute Distance

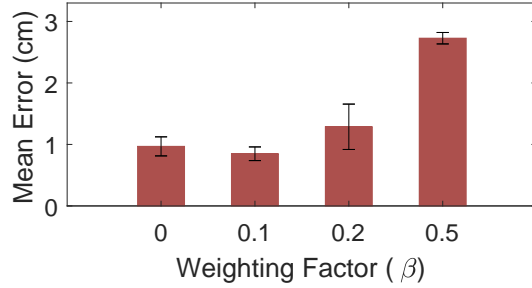
Strata can not only track the distance change from the phase, but also estimate the absolute distance using the channel difference as explained in Section 2.5. To evaluate the accuracy of the absolute distance estimation, we collected the audio and ground-truth data while users are tracing the shapes on the tablet where initial position of the finger is 20 cm away from the phone. The distance error is measured by calculating the difference between the estimated distance of the finger and microphones versus the ground-truth distance. Figure 11 shows the CDF of the distance error in 200 collected traces. The median and the 90th percentile errors are 1.0 cm and 2.1 cm, respectively. Note that the fine-grained channel estimation significantly improves the accuracy. In terms of the median error, the error reduction from the coarse channel estimation is 48%. We also implement the other schemes that exploits CIR to detect the distance of the target (*i.e.*, WiDeo [13] and Chronos [36]) and evaluate their absolute distance estimation accuracy, but they perform poorly for acoustic signals: their median tracking error is larger than 10 cm, so we do not include the result in Figure 11.

### 3.3.3 Combining Relative and Absolute Distance Estimation

**Comparison with the other acoustic tracking schemes:** Finally, we evaluate the tracking error in 2D plane. Given the finger's distance from the left and right microphones, Strata and the baseline schemes track the finger position in a 2D plane as explained in Section 2.6. As shown in the previous evaluation results, the phase-based distance tracking is more reliable than the absolute distance estimation. Therefore, we set  $\beta$  of Strata to 0.1 to give a higher weight to the estimate based on phase based tracking. For



**Figure 12: CDF of the trajectory errors with Strata, LLAP, and cross-correlation based tracking.**

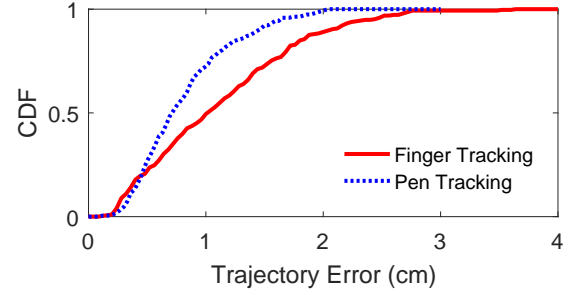


**Figure 13: 2D trajectory errors with various weighting factors.**

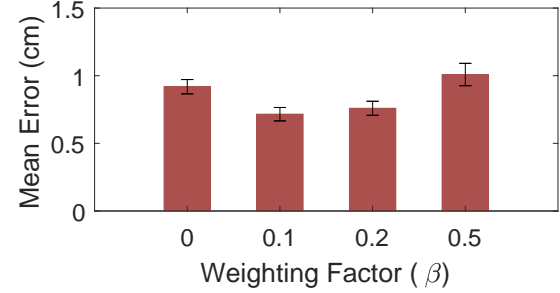
LLAP, we provide the ground-truth initial position and let it track the finger starting from the ground-truth initial position. We collect the finger trajectory tracked by the three schemes as well as the ground-truth trajectory while the users are following the shapes on the tablet screen. The initial distance is 20 cm. The trajectory error is calculated by averaging the difference in the distance between the estimated and ground-truth positions from all samples. Figure 12 shows the CDF of the trajectory errors of the three schemes after repeating the experiment 200 times for each scheme. The median errors of Strata, LLAP, and cross-correlation based tracking are 1.01 cm, 1.9 cm, and 3.47 cm, respectively. The 90th percentile errors are 2.05 cm, 3.19 cm, and 4.18 cm, respectively. Similar to 1D tracking result, Strata yields the tracking error close to half of that in LLAP, and cross-correlation based tracking shows the highest error among them. Note that the trajectory tracking error in 2D is not reported in [40].

Figure 12 also shows the tracking experiment result in noisy environment. Using a PC with an external speaker, we played EDM music [9] near the tracking device with the distance approximately 1m. We played it with median volume in the speaker, where the sound pressure level near the phone was 70 dB. The result shows noise has little impact on the tracking accuracy of Strata. Compared to the other result without artificial noise, the distribution of the tracking error is quite similar. The median error increases by 1 mm, but the 90th percentile error is hardly affected by the noise. Similar findings were reported in several different acoustic signal based tracking schemes (e.g., [44, 40, 33, 18]).

Figure 13 shows the average trajectory error with various weighting factors on the absolute distance estimation (i.e.,  $\beta$ ). The result shows that when we set  $\beta$  to small values, such as 0.1, it improves the accuracy over the phase-based tracking alone. However, increasing  $\beta$  more than 0.1 tends to increase the error because the absolute distance error is less accurate than the phase based tracking. Note that even when  $\beta = 0$ , the absolute distance estimate is still useful for getting the initial position.



(a) CDF of the 2D trajectory errors.



(b) 2D trajectory errors with various weighting factors.

**Figure 14: Strata tracking evaluation when pen is tracked instead of finger.**

**Tracking pen:** So far, we have focused on tracking the finger movement. Our CIR based tracking can effectively nullify the tracking noise from surrounding reflectors when they are sufficiently separated. However, the reflections from the hand and arm movement can affect the finger movement tracking, since the hand and arm are close to the finger and their reflections are observed in the same channel tap of the CIR. To understand the impact of hand and arm movement on the tracking accuracy, we conduct the following experiment. We track a pen attached to a 60 cm thin stick instead of a finger. The user can remotely move it so that his body movement does not affect tracking. We wrap the pen with aluminum foil so that it can generate touch event to the capacitive sensing based touch screen in the tablet, and collect the ground-truth of the pen trajectory using the same method as before.

Figure 14 shows the pen tracking result. Its median and 90th percentile errors are 0.71 cm and 1.44 cm, respectively, which are 0.3 cm and 0.6 cm reduction from the corresponding finger tracking, respectively. The error decreases since the impact of hand and arm movement is removed. Figure 14(b) shows the mean tracking errors with various weighting factors. Interestingly, when  $\beta = 0$ , the finger and pen tracking have similar tracking errors. In this case, we only use the phase to track the movement. Since the finger and hand have roughly the same movement, phase-based tracking is more robust to the arm and hand movement. In comparison, since the hand and arm have different distances from the finger, their movement has a larger impact on the absolute distance estimation. In Figure 14(b), we observe the mean pen tracking error reduces from 0.92 cm to 0.7 cm by increasing  $\beta$  from 0 to 0.1. Unlike the finger tracking, further increasing  $\beta$  does not degrade the tracking error of a pen. These results confirm our intuition that part of the absolute distance estimation error comes from the hand and arm movement.

**User study:** We evaluate how accurately users can draw shapes with real-time feedback. We implement a JAVA program that re-

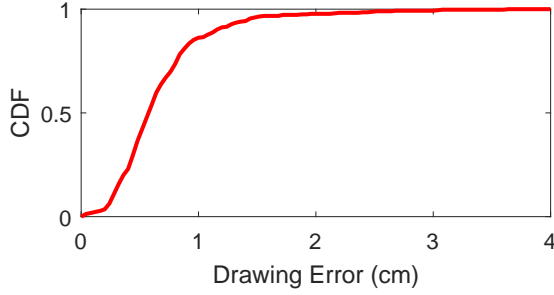


Figure 15: CDF of the drawing error in user study.

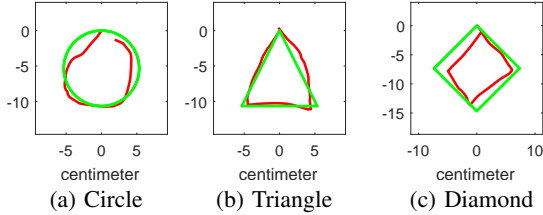


Figure 16: Median drawing error samples drawn by finger tracking.

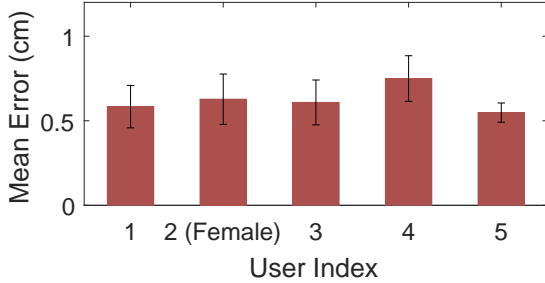


Figure 17: Mean drawing errors of 5 users.

ceives the tracking result from the mobile phone and visualizes a pointer on a PC screen controlled by the finger movement in real-time. Then we asked the user to move the pointer to trace the shapes (e.g., triangle, diamond, and circle) on the screen. The average surface area of them are 31.1 cm. The quality of the drawings is quantified using the drawing error, which is also used in [40]. For each point that Strata estimates using the acoustic signals, we find the closest point on the original shape to compute the distance and compute the distance over all points.

Figure 15 shows the CDF of the drawing errors using 200 trajectories collected from 5 users. The median and 90th percentile drawing errors are 0.57 cm and 1.14 cm, respectively. When the users draw with real-time feedback, they can compensate portion of the tracking error by moving his finger towards the desired position. As a result, the drawing error tends to lower. Figure 16 shows sample drawings of the three shapes under median drawing errors. As we can see, the traced trajectories are close to the original shapes. Figure 17 shows the mean drawing errors of 5 users. The user 2 corresponds to a female user while the other users are male. The result shows that users achieve similar accuracy.

**Tracking time:** Every 12.5ms, Strata performs low-pass filtering, passband to baseband conversion, channel estimation, and eventually tracks the finger movement. In our Samsung Galaxy S4, the average processing time is 2.5 ms. According to [40], the processing time of LLAP at each interval is 4.3 ms when the same device is used. For the down-conversion, it already takes 3.5 ms. We expect that since LLAP uses multiple sine waves in different bands

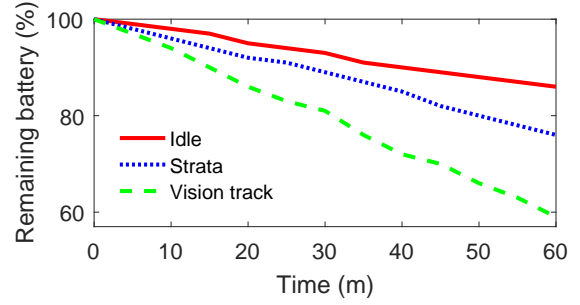


Figure 18: Energy consumption evaluation.

during the down-conversion process, it needs to filter the signal for every band they use to receive the sine waves (9 in our implementation). On the other hand, Strata performs filtering only once for the whole band. As a result, Strata spends less than 1 ms for the down-conversion.

**Energy consumption:** Finally, we evaluate the energy efficiency of Strata. In this experiment, we first fully charge the battery of the phone, turn on the tracking app and let it track the finger movement, and measure the battery consumption during 1 hour. For comparison, we also evaluate the energy consumption of the vision based tracking scheme. Since there is no particular finger tracking app in the app market, we downloaded a popular camera app in Google Play store, *Face Camera* [10], which tracks human face in camera and adds some visual effects on it. We believe finger tracking and face tracking are not quite different in terms of energy consumption. While running *Face Camera* app, we only use camera based track and do not use other features, such as taking picture or adding visual effect. We also evaluate the energy consumption of the idle phone as baseline. In all of these cases, we keep the screen turned on assuming active use of the device.

Figure 18 shows the energy consumption during one hour. While the phone is not actively running apps, it consumes 14% of its battery in an hour. When the acoustic finger tracking is used, it spends 8% more battery. On the other hand, vision based tracking consumes 27% additional energy, which is 3 times more than the energy consumption of the acoustic tracking. The result demonstrates Strata can achieve energy efficient finger tracking and is more energy efficient than vision based tracking.

## 4. RELATED WORK

We classify existing work into (i) device-free tracking and (ii) device-based tracking.

### 4.1 Device-Free Tracking and Gesture Recognition

**Vision-based:** Vision has been widely used in object tracking and gesture recognition. Among them, Microsoft Kinect [1] has been commercially very successful. It uses depth sensor as well as camera to recognize a user's movement. LeapMotion [16] is another vision based object tracking device on the market. Microsoft HoloLens [11] provides vision and depth sensing based AR capability. While it is expected that it can provide fine grained finger-level gesture recognition based on rich amount of signal resources [34], it requires a high-end device that costs over \$3,000.

To increase the user base for VR/AR, Google provides a \$10 Google cardboard [6] and many other companies sell VR headsets within \$100 (e.g., Samsung Gear), which transforms a smartphone into an AR/VR device. It is expected that smartphone-based VR/AR will likely dominate in the future due to widely available

smartphones. Since existing smartphones have limited computation power, energy, and sensing capability (e.g., no depth sensors), audio based tracking used in Strata is more attractive than vision to support smartphone-based VR by using the existing speaker and microphone on the phone and significantly reducing energy consumption. For example, [30] reports acoustic sensing consumes 20% energy compared to vision-based object recognition, and our system is even more energy efficient than vision-based approach as shown in Section 3.3.3.

**RF-based:** Recently, RF based device-free object tracking schemes in smart home and office environment have received significant attention [2, 32, 13, 41, 17]. WiSee [24] is a pioneering work that uses WiFi signals to recognize 9 gestures in several environments. WiTrack [2] applies Frequency Modulated Continuous Wave (FMCW) to track a user's location with 10-13 cm error in the  $x$  and  $y$  coordinates, and 21 cm in the  $z$ -dimension. It uses customized hardware that can sweep the channel in 1.7 GHz bandwidth. WiDraw [32] estimates angle of arrival (AoA) using CSI, and achieves a median tracking error of 5 cm using 25 WiFi access points (APs). The requirement of such a large number of APs significantly limits the applicability of WiDraw. In comparison, we only require 1 speaker and 2 microphones on one machine to achieve higher accuracy. WiDeo [13] tracks human motion based on reflected WiFi signals with 7 cm tracking error. It implements on WARP using 4 antennas with full-duplex capable transceivers, which is not readily available on the market. VitalRadio [3] also uses specialized hardware to monitor the breathing rate by observing the phase change of the RF signal. mTrack [41] and Soli [17] use 60 GHz signals for gesture recognition. While 60GHz is promising, it requires significant extra hardware for sending, receiving, and processing signals in real-time. [39] proposes Soli gesture recognition scheme with the support of Convolutional Neural Networks (i.e., deep learning), which is too expensive to run on mobile device.

**Acoustic-based:** Both LLAP [40] and FingerIO [21] track the finger movement using the reflected audio signal from a mobile phone. LLAP develops a phase based tracking while FingerIO uses OFDM symbol based movement detection. Our evaluation in Section 3.2 shows Strata out-performs both schemes because we extract the path associated with the finger movement and track its phase change instead of using the mixed signals. The authors of [7] develop a system, called UltraHaptics, to provide haptic feedback based on acoustic radiation force from a phased array of ultrasonic transducers. This requires customized hardware while Strata supports tracking in software. ApneaApp [20] uses FMCW to track heart-beat by looking at periodic patterns. Gesture recognition performs pattern matching and requires significant training data. In comparison, continuous tracking is more challenging due to the lack of training data or patterns to match against.

## 4.2 Device-based Tracking and Recognition

**IMU-based:** IMU sensors have been commonly used for motion tracking. Several works have reported that accelerometer has large error and its error increases rapidly over time due to double integration over time [38, 44]. Gyroscope has pretty good accuracy, but is not easy to use, since users have difficulty in how much to rotate in order to control certain displacement movement [44].

**Acoustic-based:** Audio is attractive for localization and tracking due to its slow propagation speed, which helps improve the accuracy. One line of research uses acoustic signals to estimate distance. For example, BeepBeep [23] develops a novel approach that allows the sender and receiver with unknown clock offsets to measure the one-way propagation delay. SwordFight [45] addresses

several practical challenges in using audio signals for tracking (e.g., quickly detecting the signal, reducing computation overhead, and accounting for measurement error during movement). Another line of research uses acoustic signals for localization. Cricket [31] uses both audio and RF to achieve median error of 12 cm with 6 beacon nodes. Swadloon [12] exploits the Doppler shift of audio signal for fine-grained indoor localization. Its error is around 50 cm. [15] uses chirp-based ranging to achieve localization error within 1 m. AAMouse [44] uses the Doppler shift to estimate the velocity, from which it computes the distance to localize the mobile. Its median error is 1.4 cm. Recently, CAT [18] develops a distributed FMCW and combines it with the Doppler shift estimation to achieve 7 mm error. The third line of audio based approaches target gesture recognition. DopLink [4] and Spartacus [33] use the Doppler shift of the audio signal between the two mobile phones for gesture recognition. DopLink [4] detects if a device is being pointed by another device. Spartacus [33] uses the Doppler shift to determine the device's moving direction and pairs it with another device moving in the same direction. However, these all require the user to hold the device, which have different applications from Strata that tracks hand movement without device.

**RF-based:** WiFi has been widely used for localization and tracking (e.g., [5, 26, 28, 36]). Many WiFi based localization uses received signal strength and their accuracy is limited due to obstacles and multipath. More recently, several works use channel state information (CSI), which reports the RSS on each OFDM subcarrier group, to provide accurate location distinction [28]. ArrayTrack [42] and SpotFi [14] further use phase of the received signal to enhance the accuracy. For example, ArrayTrack achieves a median error of 23 cm using 16 antennas. RF-IDraw [38] achieves 3.7 cm tracking error. Tagoram [43] and MobiTagbot [29] achieve higher tracking accuracy than the above RF-based schemes by assuming the RFID tag or reader moves at a constant speed, so they have different applications from Strata.

## 5. CONCLUSION

This paper develops a novel device-free acoustic tracking system that achieves 1.0 cm median tracking error. Through this process, we gain several important insights: (i) phase-based tracking is effective for acoustic signals due to its small wavelength, but we need to use the channel estimate from an appropriate tap (instead of the overall channel) to achieve high accuracy, (ii) it is hard to estimate the absolute distance by directly decomposing the channel since acoustic channel is not sparse, and our formulation based on the change removes static paths, which significantly reduces the number of unknowns and improves accuracy, and (iii) combining both distance change and absolute distance helps further improve the accuracy. Moving forward, we are interested in further enhancing the accuracy and developing applications on top of device-free tracking.

## 6. REFERENCES

- [1] Microsoft X-box Kinect. <http://xbox.com>.
- [2] F. Adib, Z. Kabelac, D. Katabi, and R. Miller. Wittrack: Motion tracking via radio reflections off the body. In *Proc. of NSDI*, 2014.
- [3] F. Adib, H. Mao, Z. Kabelac, D. Katabi, and R. C. Miller. Smart homes that monitor breathing and heart rate. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 837–846. ACM, 2015.
- [4] M. T. I. Aumi, S. Gupta, M. Goel, E. Larson, and S. Patel. Doplink: Using the doppler effect for multi-device

- interaction. In *Proc. of ACM UbiComp*, pages 583–586, 2013.
- [5] P. Bahl and V. N. Padmanabhan. Radar: An in-building RF-based user location and tracking system. In *INFOCOM 2000. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, volume 2, pages 775–784, 2000.
- [6] Google Cardboard. <https://www.microsoft.com/microsoft-hololens/en-us>.
- [7] T. Carter, S. A. Seah, B. Long, B. Drinkwater, and S. Subramanian. Ultrahaptics: Multi-point mid-air haptic feedback for touch surfaces. In *Proceedings of UIST*, 2013.
- [8] Google Daydream. <https://vr.google.com/daydream/>.
- [9] Edm mix 2017 - best remixes of popular music. <https://www.youtube.com/watch?v=IpHXpQ5sWZI&t=1038s>.
- [10] Face camera - snappy photo. <https://play.google.com/store/apps/details?id=com.fotoable.snapfilters>.
- [11] Microsoft HoloLens. <https://www.microsoft.com/microsoft-hololens/en-us>.
- [12] W. Huang, Y. Xiong, X.-Y. Li, H. Lin, X. Mao, P. Yang, and Y. Liu. Shake and walk: Acoustic direction finding and fine-grained indoor localization using smartphones. In *Proc. of IEEE INFOCOM*, 2014.
- [13] K. Joshi, D. Bharadia, M. Kotaru, and S. Katti. Wideo: Fine-grained device-free motion tracing using RF backscatter. In *Proc. of NSDI*, 2015.
- [14] M. Kotaru, K. Joshi, D. Bharadia, and S. Katti. Spotfi: Decimeter level localization using wifi. In *ACM SIGCOMM*, volume 45, pages 269–282. ACM, 2015.
- [15] P. Lazik and A. Rowe. Indoor pseudo-ranging of mobile devices using ultrasonic chirps. In *Proc. of ACM SenSys*, pages 99–112, 2012.
- [16] Leap motion. <https://www.leapmotion.com/>.
- [17] J. Lien, N. Gillian, M. E. Karagozler, P. Amihoud, C. Schwesig, E. Olson, H. Raja, and I. Poupyrev. Soli: ubiquitous gesture sensing with millimeter wave radar. In *Proc. of SIGGRAPH*, 2016.
- [18] W. Mao, J. He, and L. Qiu. Accurate audio tracker. In *Proceedings of ACM MobiCom*, Oct. 2016.
- [19] A. T. Mariakakis, S. Sen, J. Lee, and K.-H. Kim. Sail: Single access point-based indoor localization. In *Proceedings of the 12th annual international conference on Mobile systems, applications, and services*, pages 315–328. ACM, 2014.
- [20] R. Nandakumar, S. Gollakota, and N. Watson. Contactless sleep apnea detection on smartphones. In *Proc. of ACM MobiSys*, 2015.
- [21] R. Nandakumar, V. Iyer, D. Tan, and S. Gollakota. Fingerio: Using active sonar for fine-grained finger tracking. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 1515–1525. ACM, 2016.
- [22] A. V. Oppenheim, R. W. Schaffer, J. R. Buck, et al. *Discrete-time signal processing*, volume 2. Prentice-hall Englewood Cliffs, 1989.
- [23] C. Peng, G. Shen, Y. Zhang, Y. Li, and K. Tan. BeepBeep: a high accuracy acoustic ranging system using COTS mobile devices. In *Proc. of ACM SenSys*, 2007.
- [24] Q. Pu, S. Gupta, S. Gollakota, and S. Patel. Whole-home gesture recognition using wireless signals. In *Proc. of ACM MobiCom*, 2013.
- [25] M. Pukkila. Channel estimation modeling. *Nokia Research Center*, 2000.
- [26] A. Rai, K. K. Chintalapudi, V. N. Padmanabhan, and R. Sen. Zee: zero-effort crowdsourcing for indoor localization. In *Proc. of ACM MobiCom*, 2012.
- [27] T. S. Rappaport et al. *Wireless communications: principles and practice*, volume 2. Prentice Hall PTR New Jersey, 1996.
- [28] S. Sen, J. Lee, K.-H. Kim, and P. Congdon. Avoiding multipath to revive inbuilding wifi localization. In *Proceeding of the 11th annual international conference on Mobile systems, applications, and services*, pages 249–262. ACM, 2013.
- [29] L. Shangguan and K. Jamieson. The design and implementation of a mobile rfid tag sorting robot. In *Proceedings of ACM MobiCom*, pages 31–42. ACM, 2016.
- [30] K. G. Shin and Y.-C. Tung. Real-time warning for distracted pedestrians with smartphones, Sept. 25 2015. US Patent App. 14/865,262.
- [31] A. Smith, H. Balakrishnan, M. Goraczko, and N. Priyantha. Tracking moving devices with the cricket location system. In *Proc. of ACM MobiSys*, 2005.
- [32] L. Sun, S. Sen, D. Koutsonikolas, and K. Kim. Withdraw: Enabling hands-free drawing in the air on commodity wifi devices. In *Proc. of ACM MobiCom*, 2015.
- [33] Z. Sun, A. Purohit, R. Bose, and P. Zhang. Spartacus: spatially-aware interaction for mobile devices through energy-efficient audio sensing. In *Proc. of ACM Mobisys*, pages 263–276, 2013.
- [34] J. Taylor, L. Bordeaux, T. Cashman, B. Corish, C. Keskin, T. Sharp, E. Soto, D. Sweeney, J. Valentin, B. Luff, A. Topalian, E. Wood, S. Khamis, P. Kohli, S. Izadi, R. Banks, A. Fitzgibbon, and J. Shotton. Efficient and precise interactive hand tracking through joint, continuous optimization of pose and correspondences. In *Proc. of SIGGRAPH*, 2016.
- [35] D. Tse and P. Viswanath. *Fundamentals of wireless communication*. Cambridge university press, 2005.
- [36] D. Vasisht, S. Kumar, and D. Katabi. Decimeter-level localization with a single wifi access point. In *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*, pages 165–178, 2016.
- [37] Vive. <http://www.htcvive.com>.
- [38] J. Wang, D. Vasisht, and D. Katabi. RF-IDraw: virtual touch screen in the air using RF signals. In *Proc. of ACM SIGCOMM*, 2014.
- [39] S. Wang, J. Song, J. Lien, I. Poupyrev, and O. Hilliges. Interacting with soli: Exploring fine-grained dynamic gesture recognition in the radio-frequency spectrum. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology (UIST)*, pages 851–860. ACM, 2016.
- [40] W. Wang, A. X. Liu, and K. Sun. Device-free gesture tracking using acoustic signals. In *Proceedings of ACM MobiCom*, pages 82–94. ACM, 2016.
- [41] T. Wei and X. Zhang. mTrack: high precision passive tracking using millimeter wave radios. In *Proc. of ACM MobiCom*, 2015.
- [42] J. Xiong and K. Jamieson. Arraytrack: A fine-grained indoor location system. In *Proc. of NSDI*, pages 71–84, 2013.
- [43] L. Yang, Y. Chen, X.-Y. Li, C. Xiao, M. Li, and Y. Liu. Tagoram: Real-time tracking of mobile RFID tags to high

precision using cots devices. In *Proc. of ACM MobiCom*, 2014.

[44] S. Yun, Y. chao Chen, and L. Qiu. Turning a mobile device into a mouse in the air. In *Proc. of ACM MobiSys*, May 2015.

[45] Z. Zhang, D. Chu, X. Chen, and T. Moscibroda. Swordfight: Enabling a new class of phone-to-phone action games on commodity phones. In *Proc. of ACM MobiSys*, 2012.