

Localizing Acoustic Objects on a Single Phone

Hongzi Zhu^{ID}, Member, IEEE, Yuxiao Zhang, Zifan Liu, Xiao Wang^{ID}, Student Member, IEEE,
Shan Chang^{ID}, Member, IEEE, and Yingying Chen^{ID}, Fellow, IEEE

Abstract—Finding a small object (e.g., earbuds, keys or a wallet) in an indoor environment (e.g., in a house or an office) can be frustrating. In this paper, we propose an innovative system, called *HyperEar*, to localize such an object using only a single smartphone, based on enhanced time-difference-of-arrival (TDoA) measurements over acoustic signals issued from the object. One major challenge is the hardware limitations of a Commercial-Off-The-Shelf (COTS) phone with a short separation between the two microphones and the low sampling rate of such microphones. *HyperEar* enhances the accuracy of TDoA measurements by virtually increasing distances between microphones through sliding the phone in the air. *HyperEar* requires no communication for synchronization between the phone and the object and is a low-cost and easy-to-use system. We evaluate the performance of *HyperEar* via extensive experiments in various indoor conditions and the results demonstrate that, for an object of 7 m away, *HyperEar* can achieve a mean localization accuracy of about 15 cm when the object in normal indoor environments.

Index Terms—Object finding, time difference of arrival, indoor environment, smartphone, acoustic source localization.

I. INTRODUCTION

IT IS often the case that people find it is very frustrating to find their personal objects, such as earbuds, keys and wallets, in an indoor environment. With the proliferation of mobile devices (e.g., smartphones and tablets), it would be appealing to use one such mobile device to localize such small personal objects. To this end, a tag issuing inaudible acoustic signals can be attached to a private object and can be localized by the smartphone of the object owner.

Such an application pose four rigid requirements to a system as follows: 1) *far operational distance*: as a tag may be far from the mobile device of a user, especially for a large indoor environment, the system should work well at a long distance. 2) *good localization accuracy*: object finding calls for dm- or even cm-level localization accuracy. 3) *excellent user*

Manuscript received October 1, 2020; revised April 1, 2021; accepted May 12, 2021; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor K. Lee. Date of publication May 24, 2021; date of current version October 15, 2021. This work was supported in part by the National Natural Science Foundation of China under Grant 61772340, Grant 61672151, and Grant 61972081; and in part by the DHU Distinguished Young Professor Program. (Corresponding author: Shan Chang.)

Hongzi Zhu, Yuxiao Zhang, Zifan Liu, and Xiao Wang are with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: hongzi@sjtu.edu.cn).

Shan Chang is with the School of Computer Science and Technology, Donghua University, Shanghai 201620, China (e-mail: changshan@dhu.edu.cn).

Yingying Chen is with the Department of Electrical and Computer Engineering, Rutgers University, Piscataway, NJ 08854 USA (e-mail: yingche@scarletmail.rutgers.edu).

Digital Object Identifier 10.1109/TNET.2021.3080820

experience: the system should be easy to use and minimize the involvement of users; otherwise, users would get bored and quit using the system. Last but not least, 4) *low deployment cost*: the system should rely on devices that are cheap or users already have.

In the literature, existing techniques are insufficient to meet all above requirements. A number of pioneering sound source positioning systems [1]–[6] heavily rely on wireless communication for synchronization. For example, in Cricket [1], ultrasound is used to localize the speaker via Time-of-Flight (ToF) and the speaker and microphones are synchronized through radio frequency (RF) signals. Beep-Beep [4] also needs wireless communication for phone-to-phone ranging. Another large portion of existing work [7]–[12] is based on dedicated hardware such as microphone arrays deployed in specific locations in advance. For example, an 8-microphone-array is used on a mobile robot to localize a sound source in 3D [7]. Stefanakis *et al.* [12] propose a system to use a 4-microphone-array to estimate the angle-of-arrival (AoA) of multiple sound source. Several recent studies [13], [14] use smart devices to localize a sound source at mm-level accuracy. For example, keystrokes can be snooped with a smartphone placed by a keyboard [13]. vTrack [14] uses two or three microphones of a smartphone to localize and track a speaker near the phone and can achieve an accuracy of 2.3mm on 0.26m×0.2m regions. These systems can achieve very accurate localization but only work in a short range. WalkieLokie [15] is one recent work that also considers to localize a remote sound source with one single smart device. However, it needs users to continuously walk. As a result, there is no existing solution, to the best of our knowledge, to localizing an small acoustic object with one single smartphone.

In this paper, we propose an indoor acoustic source localization system, called *HyperEar*, which uses a Commercial-Off-The-Shelf (COTS) smartphone to localize a remote speaker in an indoor environment. *HyperEar* can be deployed as a software (e.g., an app) without any hardware modification. The basic idea of *HyperEar* is get highly accurate Time Difference of Arrival (TDoA) measures of inaudible acoustic beacons with a smartphone. More specifically, to find the target speaker, a user, holding his/her phone, first selects an appropriate direction and then slides the phone in the air along the direction back and forth for several times. For each onboard microphone, one TDoA can be obtained before and after one slide. With at least two microphones and two corresponding TDoA measures, triangulation can be conducted

to estimate the relative location of the phone with respect to the speaker.

Two main challenges lie in the HyperEar scheme. First, with the short distance between both microphones and their limited sampling rate, how to obtain accurate TDoA measurements is not straightforward. For instance, given the small size of a phone, two onboard microphones are separated at most at a distance of about 13-15cm. If the sampling frequency is 44.1KHz as adopted by most smartphones and the sound speed is 343m/s, there would be about only 40 distinguishable TDoA measurements, dividing the space into 40 regions. Consequently, locations in the same TDoA region cannot be discriminated. As such regions expand quickly as the distance from the phone increases, it leads to huge location ambiguity when the speaker is far away from the phone.

To deal with this challenge, HyperEar incorporates two techniques to increase the density of TDoA regions. First, leveraging the uneven distribution of TDoA regions, the phone is rolled around its z -axis to find the optimal direction, where the speaker stays in the area with densest regions (i.e., with least location ambiguity). Second, the phone is slid along the identified direction in the air and TDoAs are calculated based on the estimated sliding distance of each microphone. In this way, the number of measurable TDoAs is no longer restricted by the size of the phone but the sliding distance of the phone. Therefore, by increasing the sliding distance, the density of TDoA regions can also be increased.

The second challenge is how to accurately estimate the sliding distance of phone with low-end onboard sensors. With unstable hand operations, it is unlikely to achieve perfect movements of the phone. Furthermore, the ever-changing posture of the phone and error-prone acceleration readings result in unexpected rotation and displacement estimation errors.

To tackle this challenge, we leverage an appealing feature of HyperEar, which allows slight phone rotation and displacement from ideal slides. The reason is that the TDoA measurement error, due to unstable sliding movements, occurring on one microphone is about the same amount as that occurring on the other microphone, and they are cancelled with each other during the triangulation calculation. To deal with low-quality acceleration readings, HyperEar leverages the fact that the true velocity of the phone at the starting and at the ending of a sliding movement should be zero, and adopts a linear model to remove the accumulative errors caused by taking the integral of acceleration readings over time.

HyperEar can be used for 2D localization and can be easily extended for 3D localization as well. It has many advantages as follows. First, it reaches the minimum deployment cost with a cheap speaker and the user's smartphone. Second, it requires no synchronization or communication between the speaker and the phone. Last but not least, it is easy to use and has very light workload and elastic requirements for users. We have implemented HyperEar on two types of smartphones, and evaluated the performance of HyperEar through extensive experiments in two indoor environments with different noise types and levels. The results show that, for a speaker of 7m away, HyperEar can achieve a mean localization accuracy of

about 15cm in a normal indoor environment and 37cm in a noisy shopping mall.

We highlight our main contributions made in this paper as follows:

- We propose a novel scheme, called HyperEar, to passively estimate highly accurate TDoA of acoustic beacons issued from a distant speaker using a smartphone.
- We implement HyperEar on two Android-based smartphones, i.e., i.e., Samsung Galaxy S4 and Samsung Galaxy Note3, which demonstrates the feasibility of the proposed scheme.
- We conduct a systematic evaluation that shows the high accuracy of HyperEar. The results demonstrate the efficacy of the HyperEar design.

The remainder of this paper is organized as follows. We present the system model, basic concepts of TDoA-based localization and the restrictions of conducting acoustic localization on a smartphone in Section II. Section III introduces the architecture of HyperEar. Acoustic signal processing is elaborated in Section IV. Section V describes how to estimate the phone displacement when the user sliding the phone in the air, given the noisy inertial sensor readings. The procedures of acoustic source localization based on triangulation is introduced in Section VI. We discuss the practical issues that may be encountered in Section VII. Section VIII presents the performance evaluation and experiment results. Section IX compares HyperEar with related work. Finally, we present concluding remarks of our work and summarize the directions for future work in Section X.

II. RESTRICTIONS OF ACOUSTIC LOCALIZATION ON SMARTPHONES

In this section, we introduce the system model and analyze key factors that affect the localization accuracy.

A. System Model

To be practical, we consider a system with only two components:

- **Objects with one Single Speaker:** a speaker can be attached to a target personal object and periodically plays an acoustic signal. This signal can be audible or inaudible to human ears, depending on the specific application scenarios.
- **One Single Smartphone:** HyperEar exploits two microphones and the inertial sensors, i.e., the accelerometer and the gyroscope, embedded in a smartphone. Note that it does not need any synchronization or data communication between the phone and the speaker.

B. Basic Concepts of TDoA-Based Localization

As we do not know when a particular sound signal is emitted from the speaker, the TDoA of this sound signal can be measured with two microphones. In essence, a TDoA measurement reveals geometry information about the direction of an incoming sound. Consequently, establishing the exact

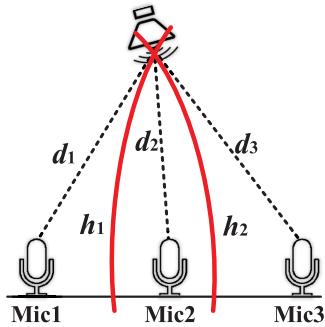


Fig. 1. Localization with triple microphones.

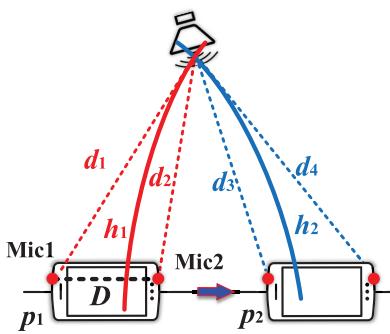


Fig. 2. Naive localization scheme with a two-microphone phone.

position of this sound in 2D normally requires triangulation, e.g., at least two distinct direction measurements. For example as illustrated in Figure 1, let d_1 and d_2 denote the distance between the speaker and two microphones Mic1 and Mic2, respectively. Suppose Δt_1 is the TDoA measured over Mic1 and Mic2, and then the distance difference Δd_1 from the speaker to this pair of microphones can be represented as

$$\Delta d_1 = d_1 - d_2 = \Delta t_1 \times S \quad (1)$$

where S is the velocity of sound. All possible positions satisfying Δd_1 lie on the half hyperbola h_1 as illustrated by the left red curve in Figure 1. Similarly, with Mic2 and a third microphone Mic3, another half hyperbola h_2 , illustrated by the right red curve, can be estimated and used to derive the relative location of the speaker (i.e., the intersection of h_1 and h_2) with respect to these microphones.

As most COTS smartphones are associated with only two microphones, one naive solution, as illustrated in Figure 2, can be used to localize a phone if the location of the speaker is known. First, a half hyperbola h_1 can be calculated at a position p_1 with the onboard microphones Mic1 and Mic2. Then, the phone can be moved to another position p_2 and a second half hyperbola h_2 can be obtained. As a result, the relative location between the speaker and the phone can also be determined if the moving direction as well as the distance from p_1 to p_2 are known.

C. Challenges for Locating Remote Objects

Several hardware limitations of a COTS smartphone make the above naive solution very challenging when the speaker is far from the phone.

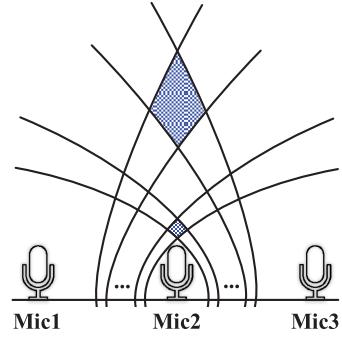


Fig. 3. Location ambiguity increases for far objects.

Limited Sampling Rate. When recording, the sound is digitized using the Analog to Digital Converter (ADC) of a microphone at a fixed sampling rate f_s . Therefore, the resolution of TDoA measurements is restricted by the f_s . Though current state-of-the-art audio hardware on smartphones supports sampling rate of up to 192kHz, operating system usually limits this to 44.1kHz, which means the resolution of TDoA measurements is about 0.023ms or the resolution of distance difference Δd (i.e., the distance interval between two adjacent hyperbolas) is about 7.78mm at $S = 343$ m/s.

Near Separation Between Microphones. The number of distinguishable hyperbolas also depends on the range of possible distance difference Δd , which is bounded by the distance between the two microphones on the phone. As can be inferred from Figure 2, the range of Δd is $[-D, D]$, where D is the distance between the two microphones. Given a sampling rate f_s and D , the number of distinguishable hyperbolas N can be calculated as

$$N = \lfloor 2Df_s/S \rfloor. \quad (2)$$

For example, the distance between the two microphones of a Samsung Galaxy S4 is 13.66cm. With a sampling rate of 44.1kHz, this yields only 35 measurable hyperbolas. As illustrated in Figure 3, with the limited number of distinguishable hyperbolas, the density of hyperbolas drops dramatically as the distance from the microphones increases. While adopting naive TDoA localization schemes can achieve cm- or mm-level accuracy for very near objects, it is very challenging to accurately localize a far sound source. For instance, the localization error of the above naive scheme can reach up to 18.6cm and 266.7cm when the sound source is located at 1m and 5m away from a Samsung Galaxy S4 smartphone, respectively.

Low-End Inertial Sensors. In the naive scheme, onboard inertial sensors can be used to estimate the information of the moving direction and distance of the phone from p_1 to p_2 . Deriving accurate motion information with sampled and error-prone sensor readings, especially under the condition that phone movements are carried out by untrained users, is very difficult.

D. Key Observations on TDoA Measurements

With regard to measuring TDoAs on smartphones, we have two main observations. First, as shown in Figure 4(a), it is

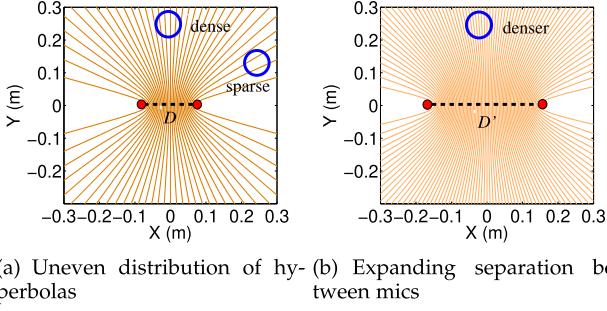


Fig. 4. Two key observations about TDoA measurements.

clear to find that the distribution of hyperbolae over space is quite uneven, with the central areas having a denser distribution of hyperbolae than other sideward areas. Second, if we increase the separation between two microphones from D to D' , as shown in Figure 4(b), the number of hyperbolae will also increase according to (2), leading to a higher density of hyperbolae at remote locations.

III. OVERVIEW OF HYPEREAR

We propose an innovative scheme, called *HyperEar*, to solve the problem of indoor smartphone localization with a single remote speaker. The core idea is to expand the TDoA measurement range by moving the phone in the air so that the number of distinguishable hyperbolae is increased, reducing the location ambiguity at a far distance from the phone. Moreover, the motion of the phone is tracked and estimated by processing the noisy inertial sensor readings. As depicted in Figure 5, the structure of HyperEar incorporates six main components.

Acoustic Signal Preprocessing (ASP). ASP performs three functions to improve the accuracy of TDoA measurements. First, ambient sound being out of the frequency band of the sound signal is filtered out. Second, interpolation is carried out to achieve sub-sample resolution. Third, sampling frequency offset (SFO) errors between the speaker and each of the two microphones are corrected.

Speaker Direction Finding (SDF). The key function of SDF is to find the direction of the speaker. It needs to detect sound signals at both microphones and measure the TDoA for each signal. Based on the relationship of previous TDoA measurements, it gives instructions to help a user find the direction of the speaker.

Motion Signal Preprocessing (MSP). The motion information of the phone is required in HyperEar. This component preprocesses the accelerometer and gyroscope readings. First, it removes high-frequency noise from both signals. It then segments the movement of the phone based on the power level of the acceleration signal.

Phone Displacement Estimation (PDE). In HyperEar, the phone is required to move in both horizontal and vertical directions. Therefore, the displacements of the phone in each direction are estimated. This component takes the segmented acceleration signals as input to estimate the moving speed and distance of the phone along some direction.

2D TDoA Localization (TTL). The key function of this component is to estimate the distance between the speaker

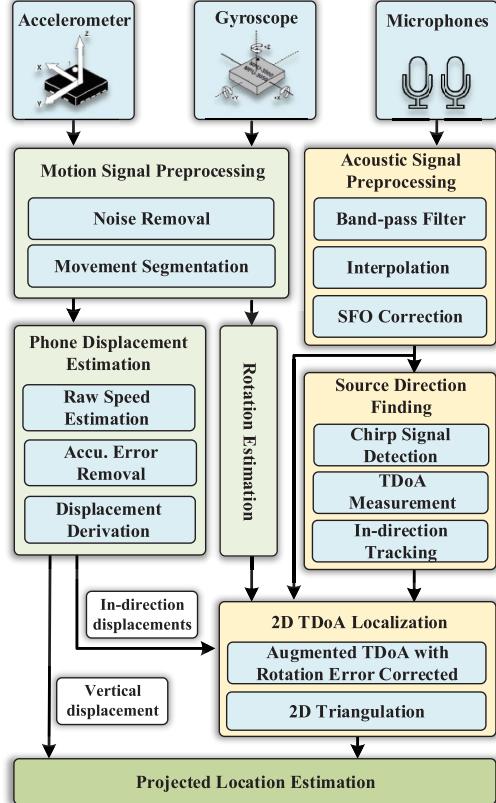


Fig. 5. System architecture of HyperEar.

and the phone. Instead of measuring a TDoA based on two microphones at the same position, it measures a TDoA based on two positions at the same microphone. It then integrates the information of the TDoA measurements, the motion estimation, and the direction of the speaker to perform 2D triangulation.

Projected Location Estimation (PLE). PLE tackles indoor smartphone localization in 3D scenarios. Instead of estimating the relative location of the speaker in 3D space, HyperEar calculates the projected location of the speaker on the floor map. In this way, it does not need to know the height information of the speaker and that of the phone.

For a user, to find the target speaker, he/she needs to hold his/her phone and slowly rolls the phone along its z -axis in order to find the direction of the speaker according to the instructions of SDF. Then, HyperEar will indicate the user to slide the phone back and forth for several times in the direction which is perpendicular to the speaker direction. During this procedure, MSP and PDE are conducted to estimate the displacement of each slide, based on which accurate TDoA are estimated for 2D localization. Finally, the user repeats the sliding operations along the same direction but at a different heights so that PLE is conducted to estimate a projected location in 3D scenarios.

IV. ACOUSTIC SIGNAL PROCESSING

A. Sampling Frequency Offset Correction

The sampling frequencies of the Digital to Analog Converter (DAC) of the speaker and the ADCs of microphones

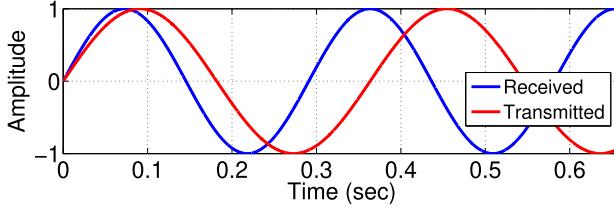


Fig. 6. Illustration of time shift of received signal due to SFO.

exhibit an offset due to non-synchronized clocks. As illustrated in Figure 6, this can cause the received signal after ADC a time shift with respect to the transmitted signal. Sampling frequency offset (SFO) leads to inaccurate TDoA measurements for two reasons: 1) it affects signal detection due to the distortion of received signals; 2) the time shift makes time measurement incorrect.

To tackle this error, we embed a sinusoid signal of a fix and non-overlapping frequency in the sound signal. Specifically, let f_s and f_m denote the sampling rate of the speaker's DAC and that of one microphone's ADC, respectively, f_0 denote the frequency of the embedded sinusoid signal, and f'_0 denote the frequency of the sinusoid signal received at this microphone. Therefore, the time shift factor f_s/f_m can be derived as f_0/f'_0 . As a result, given any smartphone, we can calibrate the SFOs between the sound source and each microphone of the phone, by scaling the duration of received samples according to the corresponding time shift factor.

B. Signal Detection and TDoA Measurement

In HyperEar, the speaker periodically plays a chirp signal, in which the frequency linearly increases with time, for its good cross correlation property. We adopt a similar method introduced in BeepBeep [4] to detect signals at each microphone. To detect the signal, a sliding window of the recorded audio signal at each microphone is correlated with a reference chirp signal. Specifically, let $x[n]$ and $y_t[n]$ denote the reference chirp signal and the recorded signal within the sliding window starting at time t , respectively. We calculate the zero-normalized cross-correlation $\rho_{xy}(t)$ as

$$\rho_{xy}(t) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\sigma_x \sigma_{y_t}} (x[i] - \mu_x)(y_t[i] - \mu_{y_t}) \quad (3)$$

where n is the number of samples in $x[n]$ and $y_t[n]$, μ_x is the average of $x[n]$ and σ_x is the standard deviation of σ_x .

We take the Min-Max feature scaling of $\rho_{xy}(t)$ and consider the first peak that has the prominence larger than a threshold as the presence of the reference chirp signal. An empirical threshold of 0.7 is used in this work.

When the i th signal is detected at both microphones of a phone, the TDoA of the signal is measured as $t_{Mic1}^i - t_{Mic2}^i$, where t_{Mic1}^i and t_{Mic2}^i represent the timestamps of the i th signal detected at Mic1 and at Mic2, respectively.

C. In-Direction Position Tracking

The purpose of finding the direction of the speaker is two-fold. First, if the direction of the speaker is known, we can roll

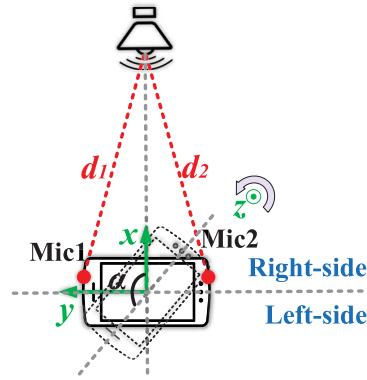


Fig. 7. Illustration of phone rotation and an in-direction position.

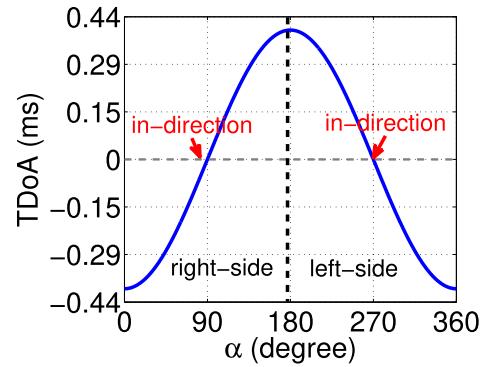


Fig. 8. Relationship between α and the distance difference $d_1 - d_2$.

the phone along its z -axis so as to make the speaker sit in the dense area of hyperbolas as depicted in Figure 4(a). Second, this direction information is required in the triangulation calculation. The SDF can find the direction of the speaker based on the fact that, when the speaker is aligned with the x -axis of the phone, the TDoA measured on the phone should be zero.

As depicted in Figure 7, when the phone is slowly rolled along its z -axis, the measured TDoA varies in the range of $[-D/S, D/S]$. Let $\alpha \in [0^\circ, 360^\circ]$ denote the angle between the direction of the speaker and the positive direction of y -axis of the phone. The speaker is considered on the right-side of the phone when $\alpha \in [0^\circ, 180^\circ]$ and on the left-side when $\alpha \in [180^\circ, 360^\circ]$. When $t_{Mic1}^i - t_{Mic2}^i = 0$, the direction of the speaker is found, i.e., $\alpha = 90^\circ$ means that the speaker locates in the positive direction of x -axis and $\alpha = 270^\circ$ means that the speaker locates in the negative direction of x -axis. In both cases, the phone is set at a so-called *in-direction* position and we stop rolling the phone. Figure 8 depicts the relationship between measured TDoAs and α obtained with a Galaxy S4.

In practice, it is unnecessary to search the whole TDoA range in order to find the *in-direction* position between the phone and the speaker. When the TDoA tends to approach zero, we slow down the rolling speed until a small enough TDoA is measured; otherwise, we change the rolling direction and continue with the search. In this way, an *in-direction* position can be effectively and accurately located.

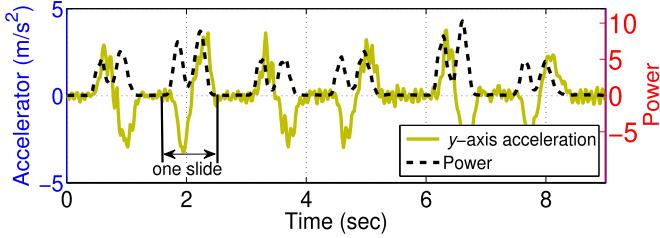


Fig. 9. Segmenting movements based on power of acceleration.

V. PHONE DISPLACEMENT ESTIMATION

In HyperEar, the TDoA measurement range is expanded by incorporating phone movements. Therefore, accurate motion information of the phone is key to both TDoA measurement and the final triangulation calculation. The PDE can obtain accurate motion estimation through comprehensive signal processing on raw inertial sensor readings.

A. Motion Signal Preprocessing

1) *Noise Removal*: We first use gravimeter to cancel the gravity to get linear acceleration data. With low-end inertial sensors, the acceleration and angular speed signals along each axis of the phone contain high-frequency noise. We remove such high frequency noise by passing each signal through a low pass filter. In this work, we use a moving average (SMA) filter, which is the unweighted mean of the previous n samples. We empirically choose the value of n to be 4 to achieve -3dB cut-off frequency at 15Hz with the sampling rate of the accelerometer and gyroscope being 100Hz.

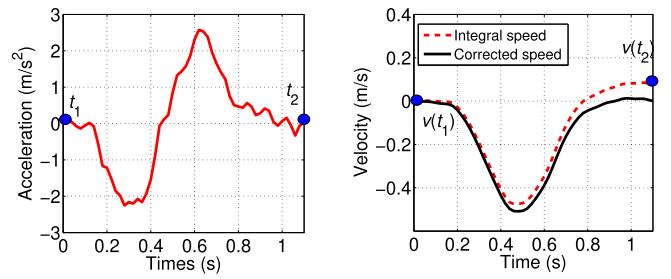
2) *Movement Segmentation*: In HyperEar, phone movements only consist of simple sliding operations along a given direction. In order to determine the starting and ending points of each slide, we first examine the power levels of the acceleration signals. Particularly, we calculate the power levels of the acceleration signal along y -axis by averaging the accumulative square of the signal amplitude in a sliding time window as shown below

$$P(t) = \frac{1}{W} \sum_{n=t}^{t+W} a(n)^2, \quad (4)$$

where W is the length of the time window and $a(\cdot)$ is the amplitude of the acceleration signal. We empirically take the length of the sliding window as 4 samples (i.e., 40ms with the sampling rate of 100Hz). Figure 9 illustrates the power levels of y -axis acceleration signal when the phone is slid back and forth along its y -axis. We consider that a slide starts when the power levels exceeds a threshold and stops when the power levels goes below the threshold for m samples. An empirical threshold of 0.2 and $m = 8$ are used in this work.

B. Sliding Velocity and Displacement Estimation

An intuitive way to estimate the moving speed of the phone along one axis is to calculate the integral of linear acceleration along that axis over time. For example, Figure 10(a) illustrates the y -axis linear acceleration of one slide as a function of



(a) Acceleration of a typical slide (b) Linear accumulative errors of integral speeds

Fig. 10. Illustration of speed estimation with accumulative errors removed.

time. The integral speed is depicted as the dashed curve in Figure 10(b). It can be seen that at the end of the slide, the integral speed drifts apart from the ground truth (i.e., zero at the ending point of a slide).

As studied in our prior work [16], the above accumulative error of integrals is approximately a linear function of time. Given the fact that the true velocity at both ends of a slide is zero, the linear model of errors can be derived and utilized to infer accurate moving speed. In specific, let t_1 and t_2 respectively denote the starting and ending points of a slide, and $v(t_1)$ and $v(t_2)$ respectively denote the integral velocity values at t_1 and t_2 , as illustrated in Figure 10. For linear accumulative errors, the slope of the linear model can be estimated as

$$\text{err}_a = \frac{v(t_2)}{t_2 - t_1}. \quad (5)$$

The instant velocity between t_1 and t_2 , therefore, can be corrected as $v^*(t) = v(t) - \text{err}_a \times (t - t_1)$, where $v(t)$ is the integral velocity at time t . For example, the solid curve as shown in Figure 10(b) illustrates the corrected speed.

Given the corrected sliding velocity $v^*(t)$, the displacement between any two time instants during a slide can be derived by taking the integral of $v^*(t)$ over time.

VI. LOCALIZATION THROUGH TRIANGULATION

A. 2D Localization Based on Augmented TDoAs

Assume that the speaker and the phone co-locate in the same horizontal plane. Without losing generality, we take the case as depicted in Figure 11, where the speaker locates on the right side of the phone. Suppose that the phone hears a signal at position p_1 with the corresponding timestamps at Mic1 and Mic2 being t_1 and t_3 , respectively, and hears the next n th signal at position p_2 with the corresponding timestamps at Mic1 and Mic2 being t_2 and t_4 , respectively. The TDoA of Mic1 at p_1 and p_2 can be measured as $\Delta t'_1 = t_2 - t_1 - n \times T$, where T is the period of signals and n is the number of detected chirp signals according to the algorithm described in Subsection IV-B. Similarly, the TDoA of Mic2 at p_1 and p_2 can be measured as $\Delta t'_2 = t_4 - t_3 - n \times T$.

If Cartesian coordinates are introduced such that the origin is the center of the two positions of Mic1 and the x -axis is the

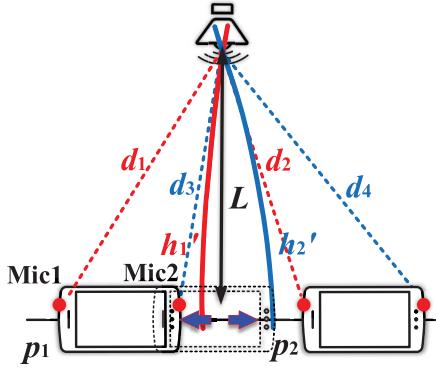


Fig. 11. Augmented TDoA Measurements.

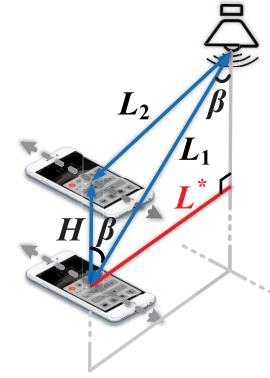


Fig. 12. Projected location on map.

inverse y -axis of the phone, then the half hyperbola h'_1 can be represented as

$$\sqrt{(x - \frac{D'}{2})^2 + y^2} - \sqrt{(x + \frac{D'}{2})^2 + y^2} = \Delta t'_1 \times S, \quad (6)$$

where D' is the estimated sliding distance between p_1 and p_2 , and (x, y) is the coordinates of the speaker. Accordingly, the half hyperbola h'_2 can be represented as

$$\sqrt{(x - D - \frac{D'}{2})^2 + y^2} - \sqrt{(x - D + \frac{D'}{2})^2 + y^2} = \Delta t'_2 \times S, \quad (7)$$

where D is the distance between Mic1 and Mic2 on the phone.

The intersection of h'_1 and h'_2 , i.e., the solution of (6) and (7), is the relative location of the speaker in this coordinate system. In particular, we are interested in the distance L as illustrated in Figure 11, which is the y coordinate of the solution.

B. Projected Location Estimation in 3D

In practice, it is hard to know the height relationship between a target speaker and the phone. Fortunately, for indoor localization applications, the location of the smartphone on a floor map is concerned. With this condition, the 2D localization based on augmented TDoAs can be extended to more general cases where the speaker and the phone have different heights and do not share a common horizontal plane.

In specific, the phone is required to slide on two horizontal planes with different heights as illustrated in Figure 12. The scheme for phone displacement estimation can also be used to estimate the height change between the two horizontal planes (e.g., the H as depicted in Figure 12) by conducting the same signal processing procedure on z -axis acceleration readings. Given the estimation of L_1 , L_2 , and H , the angle β can be calculated as

$$\beta = \arccos \frac{H^2 + L_1^2 - L_2^2}{2 \cdot H \cdot L_1}. \quad (8)$$

Therefore, the projected distance L^* can be calculated as $L_1 \times \sin(\beta)$.

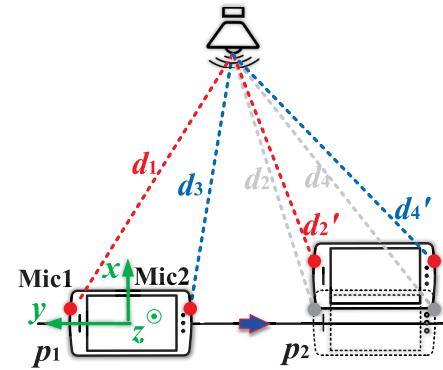


Fig. 13. TDoA errors due to displacement from the expected moving path.

VII. DISCUSSION

We discuss the major design issues encountered in implementing HyperEar as follows.

A. Limitations of Hand Operations

As hand operations are unstable, it is inevitable that the posture of the phone unexpectedly changes and/or the position of the phone deviates from the expected moving path. Figure 13 illustrates an example of phone displacement in x -axis direction while sliding the phone along its y -axis. In this case, the measured distance differences at both microphones (i.e., $d_1 - d'_1$ and $d_3 - d'_4$) differ from the expected values (i.e., $d_1 - d_2$ and $d_3 - d_4$). When the speaker is far enough from the phone, the measured TDoA at each microphone deviates from the expected value by a similar amount, making the 2D TDoA localization little affected. An analogous case is the displacement in z -axis.

In contrast, when there is a rotation around the z -axis of the phone as illustrated in Figure 14, such deviation is not symmetry on both microphones and should be corrected. This error is inevitable, especially when sliding the phone over a long distance. In HyperEar, for each slide, the rotation around z -axis is estimated by calculating the integral of the z -axis angular speed signal. Given the estimated moving distance D' , the distance between two microphones D , and the estimated rotation angle, new half hyperbola equations for

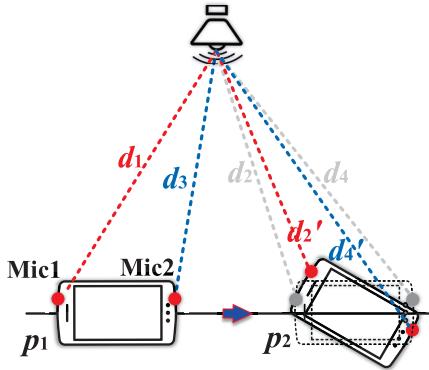


Fig. 14. TDOfA errors due to posture change of the phone.

both microphones can be established and solved. Cases such as rotation around x - or y -axis have little impact on TDOfA localization and are omitted in HyperEar design.

B. Multipath Effects and NLoS Conditions

In indoor environments, reflected signals from multipaths overlap with the signal from the Line-of-Sight (LoS) path. Such signal combination can cause the maximum peak of cross-correlation appear at some NLoS path, which is a little bit lagged with respect to the LoS path. To handle the multipath effects in HyperEar, instead of finding the peak with the maximum correlation value, we locate the first peak that has the prominence larger than a threshold as stated in Subsection IV-B.

In the case where the LoS path is blocked, HyperEar cannot find the audio source in one time. In this situation, with the proposed signal detection algorithm, a user would be guided along the direction of the first multipath with sufficient received power. If the object is not found, the user can operate HyperEar at the new location until the object is found.

C. Doppler Shift

In the design of HyperEar, a user needs to slide his/her phone in order to localize a target speaker. It is possible that a Doppler shift caused by such sliding movements may affect the TDOfA estimation precision and the final localization accuracy. The change in frequency can be calculated as

$$\Delta f = \frac{\Delta v}{v_0 f_s} \quad (9)$$

where Δv is the velocity of the phone relative to the speaker; v_0 is the velocity of sound in the air and f_{signal} is the frequency of emitted acoustic signal. It should be noted that the design of HyperEar makes the Doppler shift caused by sliding movements negligible because of three following reasons. First, Δv is very small because the sliding movements are along the tangential direction of the acoustic signals found through the in-direction position tracking as described in Subsection IV-C. Second, even if the found direction is not that perfect and a small Δv exists, the sliding speed comparing to v_0 is negligible. Third, as a fix frequency f_0 is embedded

in f_{signal} for SFO correction as described in Subsection IV-A, any possible Doppler shift will be corrected along with SFO errors.

D. Power Consumption of Battery Powered Speakers

As a target speaker may be battery powered, for example, earbuds, a tiny cheap speaker installed on wallets or keys, it is essential for HyperEar to have an ultra low power consumption. It is impractical for such a speaker to emit chirp signals all the time at a high frequency. A better design could be that the speaker only starts to emit signals when triggered, e.g., being still for a long period of time detected by a low-power accelerometer. For example, a 59mm×38mm sound module [17] with 3 200mAh LR44 button batteries and a 0.5w 8Ω speaker costs less than five dollars. If the speaker emits a 100ms chirp signal at a frequency of 5Hz, the module can continuously work for over 12,200 seconds. If a user takes one minute to find an object with HyperEar, this allows object finding for 200 times before batteries die. Moreover, the emission frequency should also be adaptive to the remaining battery.

E. Potential Application Scenarios

In addition to finding small objects, HyperEar can be used in many application scenarios. For example, with the prevalence of audio broadcasting systems deployed, HyperEar can be used to navigate a user in public buildings (e.g., office buildings, shopping malls, parking lots, hotels, and railway stations). In this application, each speaker identifies itself by playing unique chirp signals of different frequency. A user can download the digital map of the building and all chirp signals in advance. When a speaker is identified with HyperEar, the user can be localized on the digital map. For another example, a person can use HyperEar to help find his/her car in a large underground parking lot, where the car plays inaudible chirp signals after being locked. In such applications, the key issue of applying HyperEar is the operation range. There are two ways to extend the operation range as follows. First, more powerful speakers can be used to generate chirp signals of larger volume. Second, chirp signals of lower frequency can be used. With auto correlation, a pre-defined chirp signal can be accurately identified even with ambient noise. From the results presented in Subsection VIII-B, HyperEar works even when the operation distance is 30 meters, which should be sufficient in most indoor scenarios.

VIII. PERFORMANCE EVALUATION

A. Methodology

We have implemented HyperEar as an application on two models of smartphones, i.e., Samsung Galaxy S4 and Samsung Galaxy Note3. Both phones run Android 5.0 with two separated microphones that support 16-bit 44.1kHz sampling rate and stereo recording. The distance between the two microphones is 13.66cm and 15.12cm for Samsung Galaxy S4 and Samsung Galaxy Note3, respectively. A cheap desktop speaker with 2W RMS power and 150Hz-20kHz frequency

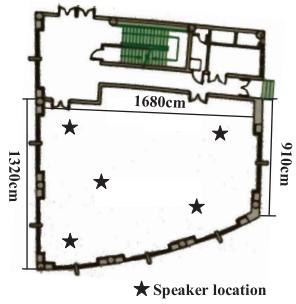


Fig. 15. The meeting room and selected speaker locations for localization experiments.



Fig. 16. One of the volunteers sliding an experimental phone in the meeting room.

response is used. The speaker is mounted on a tripod for both 2D and 3D localization, and connected to a laptop which keeps playing chirp signals on every 200ms.

We evaluate HyperEar in two indoor environments:

- *A meeting room.* The room is approximately 17m×13m with a front stage and ten rows of seats arranged in theatre fashion, as shown in Figure 15. We randomly select five positions in the room to set the speaker. For each position, we let ten volunteers, four females and six males with stature ranging from 160cm to 187cm, to operate the experimental smartphones according to the specific requirements of each experiment. An example operation of HyperEar is illustrated in Figure 16. We also control the noise level of the room by asking volunteers to keep quiet or to chat.
- *A shopping mall.* To set the speaker, we randomly select five positions in the corridor of a shopping mall, which is 95m×16.5m with shops open on both sides. We conduct experiments during off-peak hours when there is background soft music and busy hours when the place is crowded and with advertisement broadcasting.

We evaluate the HyperEar system using the metric of *accuracy* defined as the projected Euclidean distance from the estimated location and the ground truth location of the speaker on the floor map.

B. Signal Design and Frequency Selection

In HyperEar, the speaker requires to periodically emit a chirp signal. To design an appropriate signal to use, however, is non-trivial. On one hand, ambient noise could interfere

with sound signals if their frequency bands are overlapped, resulting in inaccurate TDoA measurements. As reported in previous studies [18], [19], the ambient noise in crowded spaces like cafes is relatively high at frequencies lower than 2KHz. Furthermore, in the application scenario of indoor localization, it is better to adopt sound waves with frequency higher than 17kHz, that are unobtrusive to the human ear [19], [20]. On the other hand, the performance of smartphone microphones degenerates significantly at high frequency bands. To identify appropriate frequencies to use, we conduct extensive field experiments. Specifically, we generate linear chirp signals of 150-2150Hz, 2150-4150Hz, 4150-6150Hz, 6150-8150Hz, 8150-10150Hz, 10150-12150Hz, 12150-14150Hz and 14150-16150Hz, respectively. We play each chirp signal with a desktop speaker of 2W root-mean-square (RMS) power and control the volume level to be about 60dB, measured with a COTS noise monitor. We use an experimental smartphone, i.e., Samsung Galaxy S4, to record and detect each chirp signal at different operation distances from the speaker in a corridor of the shopping mall with the presence of white noise recorded with levels of 40dB, 50dB and 60dB, respectively.

Figure 17 plots the zero-normalized cross-correlation $\rho_{xy}(t)$ obtained at different operation distances and different noise levels. In this experiment, chirp signals can be successfully identified in all cases. It can be seen that, in general, as the operation distance and noise level increase, the value of cross-correlation decreases. Moreover, choosing frequencies ranging from 2150Hz to 6150Hz can achieve better correlation performance even when the operation range is over 30 meters and signal to noise ratio is 0 dB. In the current implementation, we choose to use a 2000-6410Hz linear chirp signal of 100ms and let the speaker play the selected signal on every 200ms (i.e., $T = 200$ ms). The cubic spline interpolation method is applied to the acoustic signal recorded at each microphone to achieve sub-sample time resolution, leading to more accurate TDoA measurements given the limited ADC sampling rate of the microphones.

C. Effectiveness of Sliding

In this experiment, we first examine the effectiveness of sliding the phone in improving the accuracy of TDoA measurements and 2D localization. For each speaker position in the meeting room, we randomly select five testing positions that are at a distance of 5m away from the speaker. To eliminate the impact of unstable hand operations, in this experiment, we mount each phone on a level slide ruler when sliding. In particular, for each selected position, the ruler is set so that it has the same height as the speaker and the phone is in-direction when it moves to the center of the ruler. We then slide the phone on the ruler for 50 times with different distances, ranging from 10cm to 60cm with an interval of about 10cm. For each slide, the sliding distance is estimated with PDE (see Section V) and the distance from the speaker to the ruler is estimated with 2D Localization (see Subsection VI-A).

Figure 18 plots the phone displacement estimation errors. It can be seen that there are residual displacement estimation

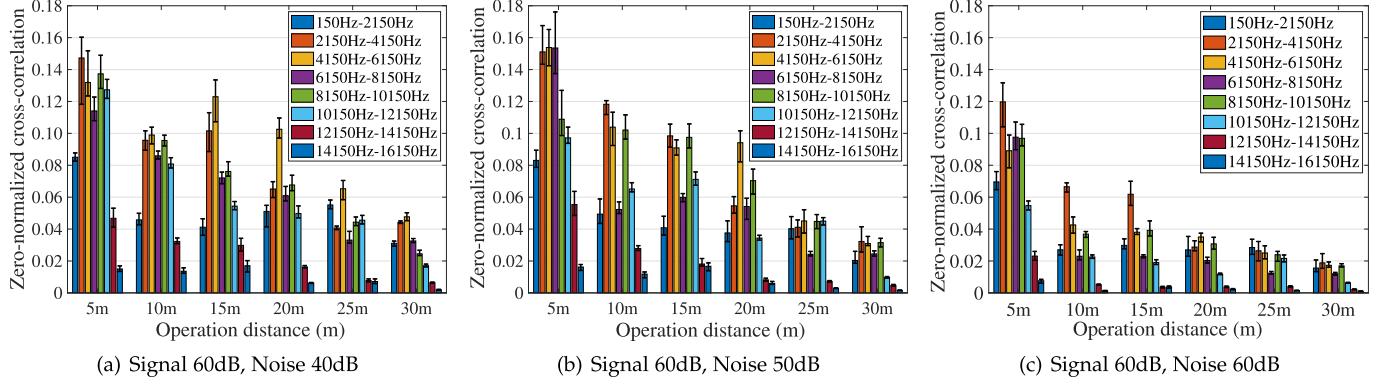


Fig. 17. Signal frequency selection with respect to various environmental noise levels and operation distances.

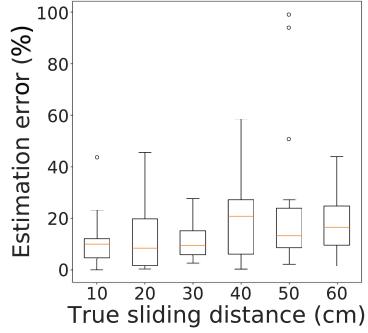


Fig. 18. Displacement estimation errors, S4 on slider ruler.

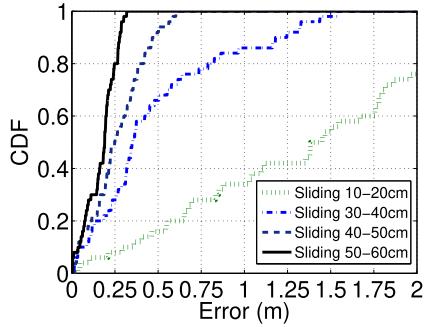


Fig. 19. CDF of 2D localization errors, S4 on slider ruler.

errors after conducting the algorithm as proposed in Subsection V-B to eliminate the linear accumulative error of integrals. Moreover, the average displacement estimation errors slight increase as the sliding distance increases. For instance, the average displacement estimation errors are less than 10% when the sliding distance is less than 30cm and such errors can reach up to 15% when the sliding distance is over 40cm. Given the restriction of human kinesiology, the effective sliding distance is also restrained.

Figure 19 depicts the cumulative distribution function (CDF) of 2D localization errors over all slides on the S4. We also have similar result on the Note3. It is obvious to see that increasing the sliding range will greatly reduce the 2D localization errors. For example, the average localization error is 18cm when sliding range is 50-60cm comparing to the value

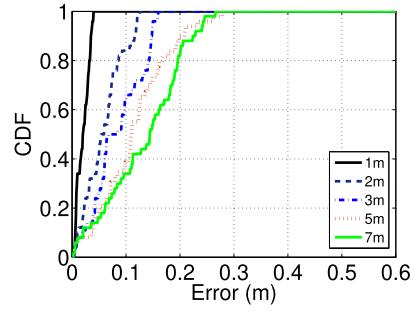


Fig. 20. CDF of 2D localization errors under different operational ranges, S4 on slide ruler.

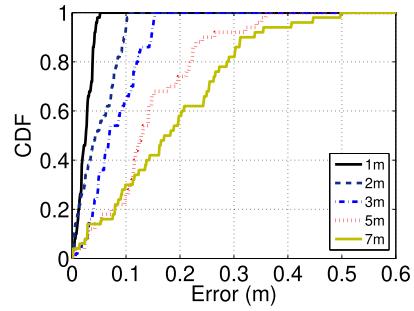


Fig. 21. CDF of 2D localization errors under different operational ranges, Note3 on slide ruler.

of 142cm when the range is 10-20cm. In practice, while it is ideal to slide the phone for longer distances, it is hard for a user to stably control the operation as the sliding distance increases. In HyperEar, slides with an estimated distance over 50cm and z -axis rotation angle less than 20° are automatically selected for use.

D. Impact of Speaker Distance

We then study the effective operation distance of HyperEar. We take a similar setting as in the above experiment except that each time we slide the phone in a range of 50-60cm on the slide ruler and change the distance between the speaker and the slide ruler from 1m to 7m at an interval of 1m.

Figure 20 and Figure 21 plot the CDFs of 2D localization errors for each distance with the S4 and the Note3,

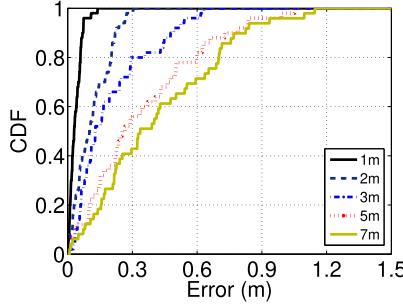


Fig. 22. CDF of 2D localization errors under different operational ranges, S4 in hand.

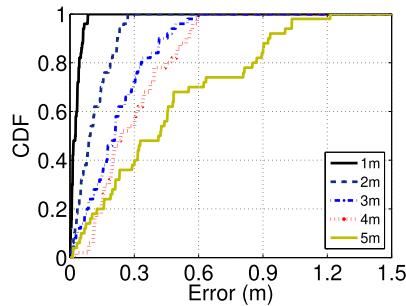


Fig. 23. CDF of 2D localization errors under different operational ranges, Note3 in hand.

respectively. It can be seen from both figures that: 1) as the speaker distance increases, the 2D localization accuracy gradually decreases; 2) HyperEar achieves better performance on the S4 than on the Note3. For example, the mean and 90%-precision accuracy for 1m distance are 2.0cm and 3.5cm, respectively, and the corresponding values are 14.4cm and 22.3cm when the distance is 7m.

E. Impact of Hand Operations

In this experiment, we repeat the above 2D localization experiment except that all sliding operations are conducted by our volunteers with one single hand as illustrated in Figure 16. In specific, for each location of the speaker and each testing position, we ask each volunteer to use each experimental phone to first find the direction of the speaker, and then to slide the phone at a similar height as the speaker over a distance of about 50cm for ten times.

When the speaker distance is 7m, the mean error of source direction finding is 1.8° with the S4. Figure 22 plots the CDFs of 2D localization errors for each speaker distance with the S4. It can be seen that comparing to Figure 20, the performance of HyperEar degrades notably (by more than two times). For example, the mean accuracy drops from 2.6cm to 3cm for 1m operational distance and from 14.4cm to 33.6cm when the distance is 7m. The results on the Note3 are similar as shown in Figure 23. This degradation is mainly because of inaccurate TDoA measurements caused by two main factors. First, the posture of the phone is varying during hand operations, which makes acceleration readings inaccurate, leading to

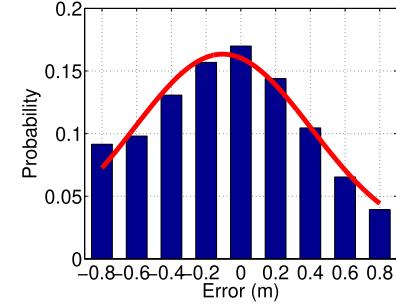


Fig. 24. The distribution of 2D localization errors follows a Gaussian.

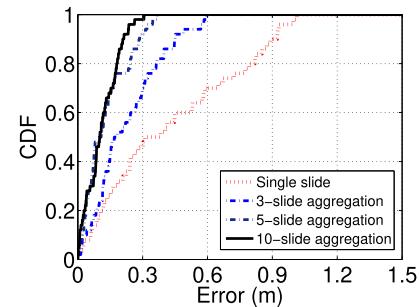


Fig. 25. Aggregating multiple slides improves 2D localization accuracy, S4 in hand.

inaccurate displacement estimates. Second, though the rotation errors due to posture change are remedied during TDoA measurements, the angle estimation is not accurate and therefore residual rotation errors exist.

We further examine the distribution of 2D localization errors. Figure 24 depicts the histogram of localization errors when the speaker distance is 7m. We find that the localization errors can be well captured by a Gaussian distribution with a negative mean. The reason of this negative mean is that, if there is an error in direction finding, the estimated distance L as illustrated in Figure 11 will be always smaller than the ground truth. With Gaussian error distribution, we can reduce 2D localization errors (as the mean is not zero) by averaging multiple measurements. In principle, more measurements lead to a better error removal, but also reduce the ease of use. Figure 25 plots the CDFs of 2D localization errors, obtained when different number of measurements of the same volunteer are averaged, over all volunteers and all position settings for a speaker distance of 7m on the S4. It can be seen that averaging over multiple measurements is effective and averaging over three to five slides seems to be a good tradeoff between localization accuracy and operational workload. For example, the median errors for single slide, aggregation over 3 slides, 5 slides, and 10 slides are 34.6cm, 17.5cm, 10.6cm, and 10.1cm, respectively.

F. 3D Localization

We then examine the performance of the full-version HyperEar smartphone localization system in 3D scenarios. We change the height of the speaker to 0.5m and randomly

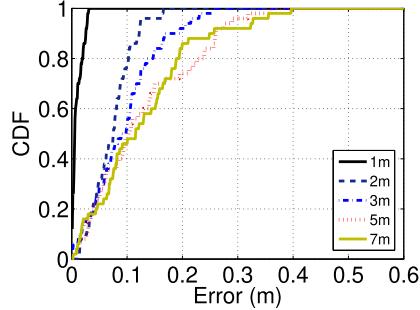


Fig. 26. CDF of 3D localization errors with 5-slide aggregation, S4 in hand.

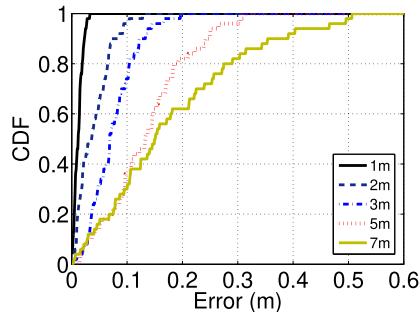


Fig. 27. CDF of 3D localization errors with 5-slide aggregation, Note3 in hand.

select 5 speaker positions in the meeting room. For each speaker position, we change the distance of randomly-selected testing positions from 1m to 7m. For each position of the speaker and each testing position, we ask each volunteer to use each experimental phone to first find the direction of the speaker, slide the phone at one customized height for five times, change to another customized height, and slide the phone for five times again.

Figure 26 and Figure 27 plot the CDFs of 3D localization errors for different speaker distances on the S4 and the Note3, respectively. It can be seen that HyperEar can achieve accurate localization in 3D scenario. For example, over a distance of 7m, the mean and 90%-precision localization accuracy on the S4 is 15.8cm and 25.2cm, respectively, and the corresponding values on the Note3 are 19.4cm and 37.5cm, respectively.

G. Different Indoor Environments

We consider the impact of different indoor environments to the performance of HyperEar. In this experiment, we mount the speaker on the tripod for the ease of deployment at five randomly selected positions in both environments as described in the methodology. For each speaker position, five testing positions that are 7m away from the speaker are selected. For each speaker position and each testing position, we ask each volunteer to perform 3D localization using HyperEar with both experimental phones. We conduct the experiment with different types of noise and control the volume of the speaker so that different signal-to-noise ratio (SNR) values are studied.

Figure 28 plots the CDFs of 3D localization errors on the S4. It can be seen that HyperEar performs stably in the

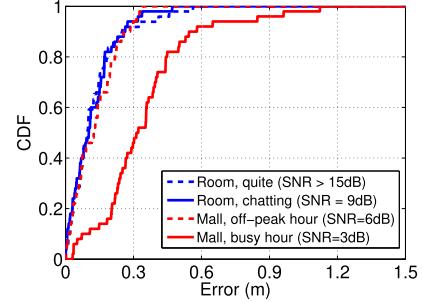


Fig. 28. CDFs of 3D localization errors in two indoor environments with different noise types and levels.

meeting room where the background noise is mainly voice. Recall that we choose a 2-6.4kHz linear chirp signal while human voice is normally lower than 2kHz, which will be filtered out and has little effect to the localization performance. We observe similar results for the shopping mall environment in off-peak hours when there is background music. Though the frequency band of the background noise in the shopping mall overlaps with that of our chirp signal, HyperEar can achieve good performance when SNR is higher than 6dB. When in busy hours, the background noise level dramatically changes over time, making the average SNR low and affecting the localization performance. In the worst case, the mean accuracy is 37.2cm at a distance of 7m.

IX. RELATED WORK

We classify existing sound source ranging and localization systems into two categories based on system complexity. We further divide each category into two subcategories, i.e., communication-based and communication-free.

A. Dedicated Hardware Based

1) *Communication-Based*: In Cricket [1], a listener uses ultrasound to measure the ToF of a beacon with the help of RF signals and calculates its location with three ToF measurements. The system has been reported to perform 100% accurately on $1.2m \times 1.2m$ regions. Wang *et al.* [8] implemented a robot navigation system with a distributed array of 24 microphones. As the robot speaks, TDoA of each pair of microphones is measured and used to calculate the robot's localization. The accuracy can reach about 7cm within about 3m range. Zheng *et al.* [6] realize a 2D localization system with multiple microphones. The system combines GPS information and wireless sensor network and can achieve a localization accuracy of 10cm within a $11m \times 11m$ area.

2) *Communication-Free*: Valin *et al.* [7] use an array of no less than four microphones to find the direction of an acoustic source. The system can reach about 1.7° error when the source is 3-5m away to the array of 8 microphones. Zhang *et al.* [9] present a unified maximum likelihood framework of sound source localization and beamforming. The system uses a microphone array and can achieve AoA accuracy of 6° on a $7m \times 6m \times 2.5m$ region. Stefanakis *et al.* [12]

use a 4-microphone-array to estimate the AoA of multiple sound sources. The system uses perpendicular cross-spectra algorithm to derive AoA of signal and count the number of sources. The AoA accuracy is 2° within 1.3m range and the counting accuracy is 94.3% for 3 sources.

B. COTS Mobile Device Based

1) *Communication-Based*: BeepBeep [4] is a high accuracy ranging system between two COTS devices. The basic idea of BeepBeep is for each phone to emit a chirp signal, capture two signals (i.e., one from itself and one from another phone), and calculate the relative distance with information exchanged through wireless communication. This method can achieve an accuracy of 5cm within 10m. Qiu *et al.* [5] propose a method based on BeepBeep to realize 3D localization between two smartphones. It can reach an accuracy of tens of centimeters within 5m.

2) *Communication-Free*: Liu *et al.* [13] use a smartphone to snoop keystrokes, reaching mm-level audio ranging. Key strokes are grouped based on TDoA. Acoustic features of key strokes are further used to differentiate keys. vTrack [14] uses two or three microphones of a smartphone and combines TDoA, AoA and power level information to localize and track a speaker near the phone. After that, it uses doppler-effect to track the movement of the speaker. vTrack can achieve an accuracy of 2.3mm on a $0.26m \times 0.2m$ region. Although these approaches can achieve extremely high accuracy with single device, they can only work in a very short range. Shake and Walk [21] uses one single microphone to find the direction of a speaker. The basic idea is to detect the frequency change caused by the doppler effect when user moves the phone. It can achieve less than 3° error in 32m distance. Walkielokie [15] is the most related work with HyperEar. It also uses one single smart device to localize a remote sound source. The system requires a user to walk and uses the doppler effect to calculate relative distances of the speaker. Walkielokie can achieve sub-meter accuracy within a range of tens of meters. However, this approach needs the user to continuously walk.

X. CONCLUSION AND FUTURE WORK

In this paper, we have proposed HyperEar, an indoor object finding system based on a single smartphone. HyperEar overcomes hardware limitations posed by a phone and can achieve 15cm accuracy on average for a desktop speaker of 7m away in normal indoor environments. HyperEar minimizes the system deployment cost by relying on cheap or existing devices, which paves the way for wide application of HyperEar.

The current implementation of Hyper has three main limitations, which direct the way of our future work. First, the system adopts a linear chirp sound signal that is audible to the human ear. While this may be helpful in the application of object finding, constantly broadcasting such sounds in public places will be annoying. In the future, we will examine to use inaudible sound signals and investigate the impact of signal distortion due to frequency selectivity of smartphone microphones. Second, HyperEar estimates the phone displacements using inertial sensor readings. In the future, we will study more

accurate schemes, such as stereo vision techniques, to tracking the motion of the phone. Third, performing sliding operation for multiple times in order to find an acoustic target may be unpleasant to use, especially in NLoS conditions. This may limits the application of HyperEar in complex environments. We will investigate the user experience of HyperEar in various scenarios to improve the operability of the system.

REFERENCES

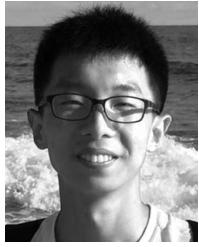
- [1] N. B. Priyantha, A. Chakraborty, and H. Balakrishnan, "The cricket location-support system," in *Proc. ACM MobiCom*, 2000, pp. 32–43.
- [2] A. Harter, A. Hopper, P. Steggles, A. Ward, and P. Webster, "The anatomy of a context-aware application," *J. Wireless Netw.*, vol. 8, nos. 2–3, pp. 187–197, Mar. 2002.
- [3] C. V. Lopes, A. Haghagh, A. Mandal, T. Givargis, and P. Baldi, "Localization of off-the-shelf mobile devices using audible sound: Architectures, protocols and performance assessment," *ACM SIGMOBILE Mobile Comput. Commun. Rev.*, vol. 10, no. 2, pp. 38–50, Apr. 2006.
- [4] C. Peng, G. Shen, Y. Zhang, Y. Li, and K. Tan, "BeepBeep: A high accuracy acoustic ranging system using cots mobile devices," in *Proc. ACM SenSys*, 2007, pp. 1–14.
- [5] J. Qiu, D. Chu, X. Meng, and T. Moscibroda, "On the feasibility of real-time phone-to-phone 3D localization," in *Proc. ACM SenSys*, 2011, pp. 190–203.
- [6] X. Zheng, S. Yang, N. Jin, L. Wang, M. L. Wymore, and D. Qiao, "DivA: Distributed Voronoi-based acoustic source localization with wireless sensor networks," in *Proc. IEEE INFOCOM*, Apr. 2016, pp. 1–6.
- [7] J.-M. Valin, F. Michaud, J. Rouat, and D. Létourneau, "Robust sound source localization using a microphone array on a mobile robot," in *Proc. IEEE/RSJ IROS*, vol. 2, Oct. 2003, pp. 1228–1233.
- [8] Q. H. Wang, T. Ivanov, and P. Aarabi, "Acoustic robot navigation using distributed microphone arrays," *Inf. Fusion*, vol. 5, no. 2, pp. 131–140, Jun. 2004.
- [9] C. Zhang, D. Florêncio, D. E. Ba, and Z. Zhang, "Maximum likelihood sound source localization and beamforming for directional microphone arrays in distributed meetings," *IEEE Trans. Multimedia*, vol. 10, no. 3, pp. 538–548, Apr. 2008.
- [10] J. Xiong and K. Jamieson, "ArrayTrack: A fine-grained indoor location system," in *Proc. USENIX NSDI*, 2013, pp. 1–14.
- [11] X. Alameda-Pineda and R. Horraud, "A geometric approach to sound source localization from time-delay estimates," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 22, no. 6, pp. 1082–1095, Jun. 2014.
- [12] N. Stefanakis, D. Pavlidi, and A. Mouchtaris, "Perpendicular cross-spectra fusion for sound source localization with a planar microphone array," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 25, no. 9, pp. 1821–1835, Sep. 2017.
- [13] J. Liu, Y. Wang, G. Kar, Y. Chen, J. Yang, and M. Gruteser, "Snooping keystrokes with mm-level audio ranging on a single phone," in *Proc. ACM MobiCom*, 2015, pp. 142–154.
- [14] S. Chung and I. Rhee, "vTrack: Virtual trackpad interface using mm-level sound source localization for mobile interaction," in *Proc. ACM UbiComp*, 2016, pp. 41–44.
- [15] W. Huang *et al.*, "Walkielokie: Sensing relative positions of surrounding presenters by acoustic signals," in *Proc. ACM UbiComp*, 2016, pp. 438–450.
- [16] H. Han *et al.*, "SenSpeed: Sensing driving conditions to estimate vehicle speed in urban environments," in *Proc. IEEE INFOCOM*, Apr. 2014, pp. 727–735.
- [17] *Sound Module*. Accessed: May 10, 2021. [Online]. Available: https://detail.tmall.com/item.htm?spm=a230r.1.14.13.77b049b9oyvNZ5&id=624879861872&cm_id=140105335569ed55e27b&abbucket=2
- [18] R. Nandakumar, K. K. Chintalapudi, V. Padmanabhan, and R. Venkatesan, "Dhwani: Secure peer-to-peer acoustic NFC," in *Proc. ACM SIGCOMM*, 2013, pp. 63–74.
- [19] Q. Wang, K. Ren, M. Zhou, T. Lei, D. Koutsikopoulos, and L. Su, "Messages behind the sound: Real-time hidden acoustic signal capture with smartphones," in *Proc. ACM MobiCom*, 2016, pp. 29–41.
- [20] W. Wang, A. X. Liu, and K. Sun, "Device-free gesture tracking using acoustic signals," in *Proc. ACM MobiCom*, 2016, pp. 82–94.
- [21] W. Huang *et al.*, "Shake and walk: Acoustic direction finding and fine-grained indoor localization using smartphones," in *Proc. IEEE INFOCOM*, Apr. 2014, pp. 370–378.



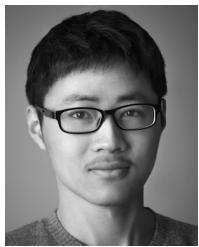
Hongzi Zhu (Member, IEEE) received the Ph.D. degree in computer science from Shanghai Jiao Tong University in 2009. He was a Post-doctoral Fellow with the Department of Computer Science and Engineering, The Hong Kong University of Science and Technology, in 2009, and the Department of Electrical and Computer Engineering, University of Waterloo, in 2010. He is currently a Professor with the Department of Computer Science and Engineering, Shanghai Jiao Tong University. His research interests include vehicular networks, mobile sensing, and mobile computing. He received the Best Paper Award from IEEE Globecom 2016.



Shan Chang (Member, IEEE) received the B.S. degree in computer science and technology and the Ph.D. degree in computer software and theory from Xi'an Jiaotong University in 2004 and 2013, respectively. From 2009 to 2010, she was a Visiting Scholar with the Department of Computer Science and Engineering, The Hong Kong University of Science and Technology. She was also a Visiting Scholar with BBC Research Laboratory, Electrical and Computer Engineering Department, University of Waterloo, from 2010 to 2011. She is currently a Professor with the Department of Computer Science and Technology, Donghua University. Her research interests include security and privacy in wireless networks and machine learning algorithms.



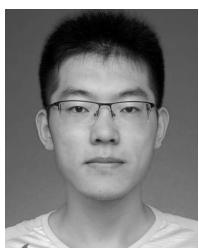
Yuxiao Zhang received the B.S. degree from the Department of Computer Science and Engineering, Shanghai Jiao Tong University in June 2017. He is currently pursuing the master's degree with Yale university. His research interests include pervasive computing, wireless networks, and mobile sensing.



Zifan Liu received the B.S. degree from the Department of Computer Science and Engineering, Shanghai Jiao Tong University in June 2018. His research interests include mobile sensing and computer vision.



Yingying Chen (Fellow, IEEE) received the Ph.D. degree in computer science from Rutgers University, New Brunswick, NJ, USA, in 2007. She is currently a Tenured Professor of electrical and computer engineering with Rutgers University and a member of the Wireless Information Network Laboratory (WINLAB). Prior to joining Rutgers, she was a Professor with the Department of Electrical and Computer Engineering, Stevens Institute of Technology, Hoboken, NJ, USA. Her research interests include cyber security and privacy, mobile and pervasive computing, and mobile healthcare. She has published over 80 journal and refereed conference articles in these areas. She was a recipient of the NSF CAREER Award and Google Faculty Research Award. She also received an NJ Inventors Hall of Fame Innovator Award. She is the recipient of Best Paper awards from IEEE CNS 2014 and ACM MobiCom 2011. She also received the IEEE Outstanding Contribution Award from the IEEE New Jersey Coast Section each year during 2005-2009. Her research has been reported in numerous media outlets including MIT Technology Review, Fox News Channel, The Wall Street Journal, and National Public Radio. She is an Editorial Board of the IEEE TRANSACTIONS ON MOBILE COMPUTING, IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, and IEEE NETWORK.



Xiao Wang (Student Member, IEEE) received the B.S. degree from the Department of Computer Science and Engineering, Shanghai Jiao Tong University in 2019. He is currently pursuing the master's degree with UM-SJTU Joint Institute, Shanghai Jiao Tong University. His research interests include mobile computing and big data analysis.