

Ultrasound for Gaze Estimation- A Modeling and Empirical Study

Andre Golard¹ and Sachin S. Talathi¹

¹ Facebook Reality Labs, Redmond, WA 98052, USA; agolard@fb.com

* Correspondence: stalathi@fb.com

§ An abbreviated version of this manuscript was published at International Conference on Pattern Recognition workshop on Eye Tracking Techniques, Challenges and Applications, [1]

Abstract: Most eye tracking methods are light-based. As such, they can suffer from ambient light changes when used outdoors, especially for use cases where eye trackers are embedded in Augmented Reality glasses. It has been recently suggested that ultrasound could provide a low power, fast, light-insensitive alternative to camera-based sensors for eye tracking. Here, we report on our work on modeling ultrasound sensor integration into a glasses form factor AR device to evaluate the feasibility of estimating eye-gaze in various configurations. Next, we designed a benchtop experimental setup to collect empirical data on time of flight and amplitude signals for reflected ultrasound waves for a range of gaze angles of a model eye. We used this data as input for a low-complexity gradient-boosted tree machine learning regression model and demonstrate that we can effectively estimate gaze (gaze RMSE error of 0.965 ± 0.178 degrees with an adjusted R^2 score of 90.2 ± 4.6).

Keywords: eye tracking, gaze estimation; ultrasound; CMUT; Machine Learning; Gradient Boosted Regression Trees, Comsol Modeling

Citation: Golard, A.; Talathi, S.S. Ultrasound for Gaze Estimation- A Modeling and Empirical Study. *Journal Not Specified* **2021**, *1*, 0. <https://doi.org/>

Academic Editor: Marco Porta

Received: 19 May 2021

Accepted: 16 June 2021

Published:

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Copyright: © 2021 by the authors. Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Most current eye tracking methodologies use video to capture the position of the iris and/or reflected lights sources – glints [2]. As such, these methods can be affected by ambient light [3], which is particularly true for use cases such as augmented reality with eye glasses. Other light-based methods such as scanning lasers [4], third Purkinje images [5] and directional light sensors [6] can likewise be affected. Speed can also be limited, especially in wearables, where operating a camera at high speed (on the order of 100 Hz or above) would imply high power consumption. At these speeds the camera-based sensors can capture fixations but not other parameters such as saccades, which have been implicated as a markers of schizophrenia spectrum in at-risk mental states [7] as well as other neurological disorders [8]. Fast eye tracking is required for measuring saccades. Current devices capable of measuring saccades are designed for laboratory use, and tend to lack portability [9]. The possibility of using ultrasound for eye tracking has been raised in a patent [10] and there exist studies that use eye-tracking to assist ultrasound procedures [11]. However, to the best of our knowledge, there is no report on experimental study to empirically demonstrate the feasibility for gaze estimation using ultrasound sensors.

A recent paper explored the possibility of using non-contact ultrasound sensors to track fast eye movements [12]. The work focused on the development of a finite element simulation model to investigate the use for ultrasound time of flight data to track fast eye motions. The simulation model is based on a setup made of four transducers positioned perpendicular to the cornea. Distances are measured with each transducer based on the time for it to receive the reflection of its own signal. Given the cornea protrudes, this time changes with the gaze angle. For implementing this simulation setup in any

38 form of glasses-form factor device, the device needs to be precisely positioned relative
39 to the eye. However, we are interested in applications for eye tracking in augmented
40 reality (AR) and virtual reality (VR), where user-specific placement of the sensors is not
41 possible (in AR/VR the eye tracking system will be fixed and the position of the eye will
42 vary from user to user, which means alignment will vary).

43 It is also to be noted that the modeling in [12] was done in the absence of occlusions
44 (such as eyelids). Eye occlusions are known to be problematic for eye tracking systems in
45 general [13]. Furthermore, the authors [12] chose to model standard 40 kHz transducers.
46 While these would be advantageous in terms of minimizing attenuation in air, such
47 a system may be subject to interference from range-finding applications (typically in
48 the 40 - 70 kHz range). Common range finding systems lack the resolution and short
49 distance sensing capabilities required for eye tracking (the typical sensing range would
50 be in meters with a resolution of 1 cm). Another concern for our application of interest is
51 size. Devices would need to fit in glasses frames. Capacitive Micromachined Ultrasonic
52 Transducers (CMUTs) operating at 500 kHz-2 MHz [14] provide the range, resolution
53 and size that is suitable for use in VR and AR devices. This type of transducer has found
54 numerous medical applications in both imaging and therapy [14], which are applications
55 for contact ultrasound.

56 Here, we use the CMUTs for remote sensing as airborne transmitters and receivers.
57 In this mode, the difference in impedance between air and tissue means over 99 percent
58 of the ultrasound signal will be reflected by the eye surface [12]. As such, the size of
59 transducers was a primary concern for our choice of CMUTs for the proposed study
60 and concerns related to test bench size and power consumption did not drive our
61 investigations.

62 In order to systematically investigate the feasibility of near-field ultrasound sensing
63 for eye-tracking with an AR form-factor device, we first did our own finite-element-
64 modeling study using acoustic rays for 1.7 MHz transducers configured on AR glasses.
65 We compared directional and omnidirectional transmit and receive configuration for
66 the sensors to determine where we would expect to see a meaningful signal around
67 glasses frames for a source placed near the glasses branch. We then built a series of table
68 top test bench systems to (a) verify our ability to accurately measure distances in the
69 appropriate range, (b) characterize the transducers, and (c) generate data to be used in
70 a machine learning model to estimate gaze. As such we focus on empirically testing
71 the hypothesis that ultrasound sensors can be used for gaze estimation in the presence
72 of occlusions. We note that in the context of our experiments, gaze is defined by the
73 static orientation of model eye on the goniometer. We demonstrate that ultrasound time
74 of flight and amplitude signals can be leveraged to track gaze in such conditions. In
75 particular, we train a regression model using gradient boosted decision trees to estimate
76 the gaze vector given the set of ultrasound time-of-flight and amplitude signals captured
77 by the CMUT receivers. The nonlinearities introduced by occlusion artifacts make the
78 task of regressing gaze directly from recorded signals non-trivial and we believe that
79 a nonlinear regression model trained on the collected data is best suited to extract the
80 relevant signals for gaze estimation. Results show that the trained model produces a
81 regression R^2 score of 90.2 ± 4.6 % and a gaze RMSE error of 0.965 ± 0.178 degrees.

82 2. Materials and Methods

83 In this section, we describe the set up for acoustic ray tracing modeling, the bench-
84 top experimental setup for data collection, the signal processing steps to extract the
85 ultrasound time of flight and amplitude signals, and the machine learning framework
86 adopted to train a gaze estimation model.

87 2.1. Modeling

88 Ultrasound is modeled as rays released all at once from a single point. Their position
89 is updated at fixed time intervals. We did this so we could trace the path of signals

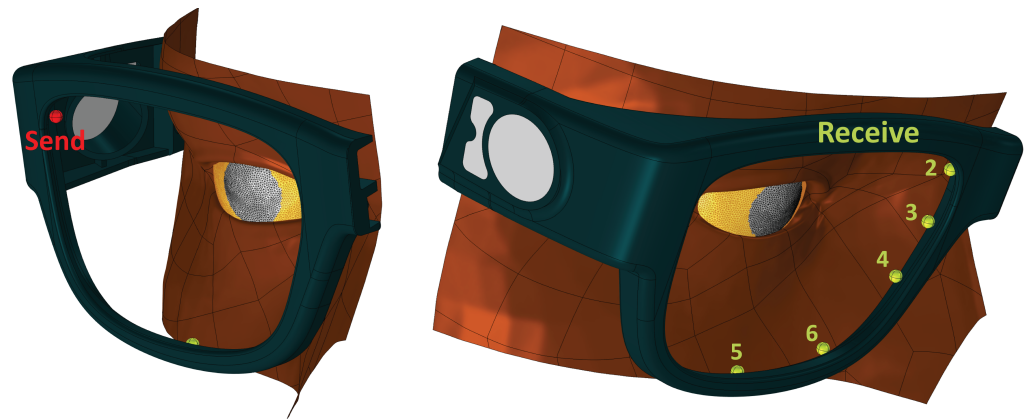


Figure 1. Physical layout for our model. Send location: transducer operating in transmit mode; numbered locations: transducers operating in receive mode. The distance traveled by a ray starting at the send position, reflecting off the cornea, and arriving at receiver 4 is in the 5.28 to 5.35 mm depending on gaze.

90 reaching the receiver and determine if they were reflections from the cornea or the skin
 91 or glasses. We used the acoustic ray tracing features of COMSOL Multiphysics software
 92 (<https://www.comsol.com/release/5.5>). We used a fixed value of 343 m/s for the speed
 93 of sound in air. We accounted for an acoustic wave attenuation in air that corresponds to
 94 1.7 MHz. The absorption attenuation coefficient is 470 (dB/m). To estimate the reflection
 95 of the signal we used the formula $R = \left(\frac{Z_2 - Z_1}{Z_2 + Z_1}\right)^2$ where R is a measure for the fraction
 96 of sound reflected and Z_1 and Z_2 are the impedance of the two media [12]. The acoustic
 97 impedance of a medium are its density times the speed of sound. We used a density of
 98 1 kg/m^3 for air. The densities of the solids range from 911 kg/m^3 for the tear film to
 99 $2,580 \text{ kg/m}^3$ for glass, with $1,051 \text{ kg/m}^3$ for the cornea. The speed of sound in the solids
 100 ranges from 1,450 m/s in fat to 4,500 m/s in glass. Based on these values we estimate
 101 99.87-99.99 percent of the signal will be reflected. These calculations guided our decision
 102 to assume 100 % reflection of ultrasound waves of eye in our modeling.

103 We used a scanned eye surface obtained with an Eye3D scanner (Transfolio, Ma-
 104 rina del Rey, CA, USA). A fit of cornea with a sphere shows a radius of 5.65 mm.
 105 The surface was smoothed, and the mesh size adjusted using Autodesk Meshmixer
 106 (<https://www.meshmixer.com>). We used it to create gaze variants: straight, ± 20 degrees
 107 in the vertical direction, ± 30 degrees horizontal.

108 We used a scanned face and glasses designed in Solidworks to create the eye box
 109 (the space in which rays will propagate). Locations for the transducers are shown
 110 in Figure 1. These positions were arbitrary. (The Comsol model was built by Veryst,
 111 Needham Heights, MA, USA).

112 2.2. Benchtop setup

113 We designed a series of three test benches to evaluate distance measurements, signal
 114 attenuation, transducer directionality, and our ability to estimate gaze.

115 In terms of electronics and data acquisition, all test benches are based on a CMUT
 116 evaluation kit from Fraunhofer IPMS (Dresden, Germany). This test kit is comprised of
 117 CMUT transducers (1.74 MHz), an amplifier, bias-tee, and associated software. These
 118 transducers fit our size and power requirements.

119 We first verified our ability to measure distances, as well as the signal decay due
 120 to attenuation in air given that ultrasound signal attenuation is significant at MHz
 121 frequencies [15]. We used a setup consisting of a pair of transducers aimed at a flat target
 122 attached to a linear translation stage (Test bench 1, Figure 2A).

123 Next we tested the emission properties of the transducers. Our CMUTs are com-
 124 prised of an array of cells connected to a single electrode and a single counter electrode.

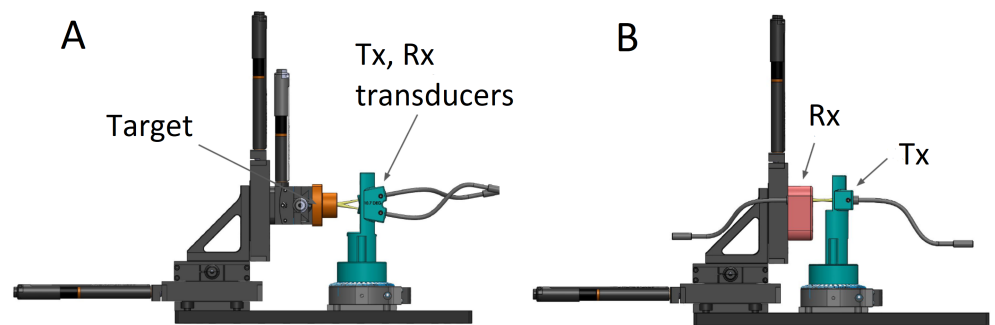


Figure 2. CAD schema for attenuation (A) and directionality (B) test benches. Tx refers to transducer in transmit mode, Rx receive mode. In A the transducers are fixed and the target is moved. In B the Rx is fixed and the Tx rotated.

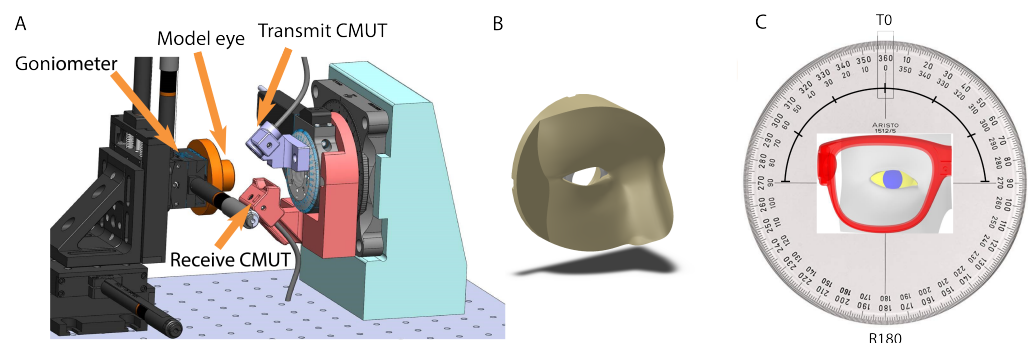


Figure 3. A: CAD schema for the experimental bench-top setup and B: occlusions C. Transducer rotation. The receiver is fixed and the transmitter rotates around an arc. 30 degree steps are shown. We acquired data from -90 to +90 degrees in 10 degree steps.

125 As such they act as a fixed phased array, which is expected to exhibit directionality. We
 126 tested this using a fixed transducer and one on a rotating stage (Test bench 2, Figure 2B).
 127 The Tx transducer was rotated in 1 degree increments and the amplitude of the Rx signal
 128 was recorded.

129 Our third test bench is designed for gaze estimations (Figure 3A). As noted earlier,
 130 we define gaze in terms of the static orientation of model eye on the goniometer. The
 131 transducer side is on the right. We used a pair of transducers (one in transmit mode
 132 and one receiver) mounted on rotating stages to allow us to mimic multiple locations
 133 around a ring (or glasses frame). We acquired data for all transmit and receive locations
 134 covering 360 degrees in 10 degree increments (Figure 3C).

135 On the target side (left part of Figure 3A), a standard sphere on sphere model eye
 136 (cornea radius 7.8 mm, sclera radius 11.925 mm, offset 5.6 mm) was mounted on a
 137 goniometer (Thor Labs). Note these dimensions differ slightly from the scanned eye
 138 used for modeling. This does not affect our findings, see discussion. Gaze angles were
 139 set in one degree increments between ± 5 degrees in both up/down (ϕ) and left/right (θ)
 140 directions.

141 Occlusions (known to affect eye trackers) were added for realism. This is a step
 142 forward from previous modeling which totally ignored occlusions. We did not model or
 143 attempt to integrate eyelashes. Our occlusions consisted of a partial scanned face printed
 144 in flexible material (A40 durometer Polyjet) with a cavity to accommodate the model
 145 eye (Figure 3B). This was mounted in front of and against the model eye and allowed
 146 the eye to move freely.

147 Our test signal consisted of a train of seven oscillations at 1.74 MHz, repeated at
 148 2 kHz. The transmitter was moved to positions around a 180 degree arc opposite the
 149 receiver (-90, -80, . . . , 80, 90), Figure 3C. Fifty runs were recorded for each transducer

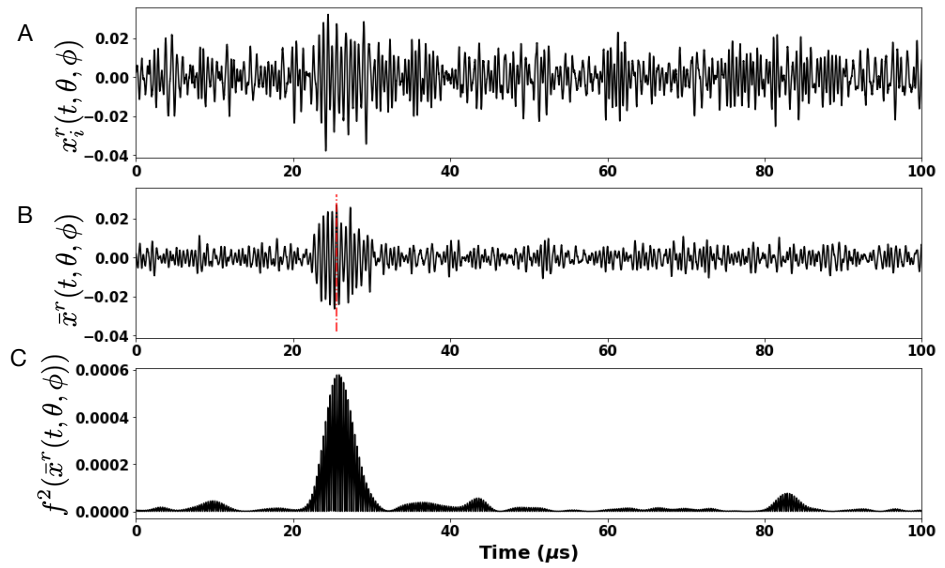


Figure 4. Example of recorded raw time trace of ultrasound sensor signal. The top row (A) shows an example of time trace recorded at the receive Ultrasound CMUT sensor in response to a single burst of test signal. The middle row (B) shows averaged signal computed from the response to a set of 10 bursts of test signal. Finally the last row (C) shows the squared filtered response signal out of a Butterworth filter. The red line indicates the time period of time-to-peak signal detection.

150 position. The series was repeated for all static goniometer positions. The received signal
151 was digitized at 80 MHz.

152 2.3. Data Analysis

153 2.3.1. Feature Engineering

154 The raw signal carries too much noise to allow for accurate peak time and amplitude
155 measurements. Improvements are possible (data not shown). For this proof of concept
156 we used averaging and filtering. In Figure 4A, we show one raw trace, $x_i^r(t, \theta = 0, \phi = 0)$
157 ($i \in [0, 49]$ and $r \in [-90, -80, \dots, 80, 90]$), for the ultrasound signal captured at the
158 receiver, in response to a single test signal emitted by the transmit CMUT transducer.
159 Figure 4B, shows the average of ten traces, defined as $\bar{x}_k^r(t) = 0.1 \sum_{j=1}^{10} x_j(t)$ ($k \in [0, 4]$).
160 The ultrasound time of flight, $\tau_k^r(\theta, \phi)$, and amplitude, $a_k^r(\theta, \phi)$, signal is estimated for
161 each $\bar{x}_k^r(t, \theta, \phi)$ as follows: the signal, $\bar{x}_k^r(t, \theta, \phi)$ is band-pass filtered in the frequency
162 range, [1.6 MHz - 1.9 MHz] using a Butterworth filter of order 4 to generate the filtered
163 version, $f(\bar{x}_k^r)(t, \theta, \phi)$. In Figure 4C, we show the trace for $f^2(\bar{x}_k^r)(t, \theta = 0, \phi = 0)$. The
164 ultrasound time to peak $\tau_k^r(\theta, \phi)$ and the amplitude, $a_k^r(\theta, \phi)$ is obtained by considering
165 a time window of 45 μs around the time instance of peak value for $f^2(\bar{x}_k^r(t, \theta, \phi))$ and
166 finding the first instance of the peak value for $\bar{x}_k^r(t, \theta, \phi)$ within the considered time
167 window. The detected peak value represents the amplitude signal $a_k^r(\theta, \phi)$ and the time
168 to peak, $\tau_k^r(\theta, \phi)$. Thus, for each position $\mathbf{Y} = (\theta, \phi)$ of the model eye on the goniometer,
169 we obtain a set $k=5$ feature vectors $\mathbf{X} \in \mathbb{R}^{36} = \{a^r, \tau^r\}_{r=[-90, -80, \dots, 80, 90]}$ per experimental
170 run. In order to collect sufficiently robust dataset and also to account for changes in
171 day to day environmental fluctuations, we conducted a total of 9 experiments spanning
172 a period of 9 days. In total, for each position, \mathbf{Y} , on the goniometer, we were able to
173 compile a set of 9×5 feature vectors, \mathbf{X} , and our goal for ultrasound based eye tracking
174 is to learn a regression model, $H : \mathbf{X} \rightarrow \mathbf{Y}$; that is, given the ultrasound sensor time of
175 flight and amplitude data, estimate two-dimensional eye gaze coordinates.

176 In Figure 5A and 5B, we plot the distribution of $\tau^r(0, 0)$ and $a^r(0, 0)$ respectively. In
177 the last sub-plot for each of the figures we show how the mean time-of-flight and the
178 mean amplitude signal changes as function of the position of the receiver transducer.

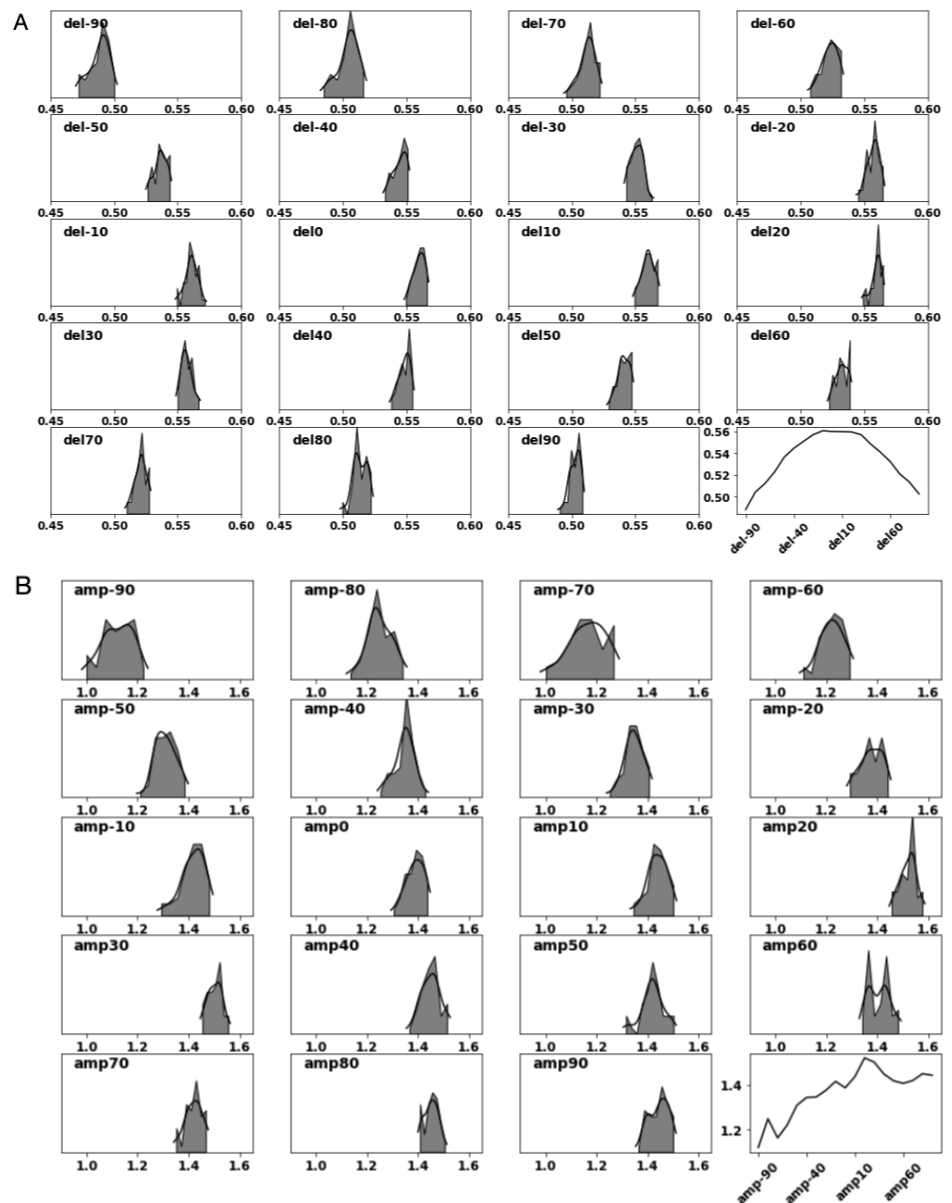


Figure 5. Distribution of the ultrasound time-of-flight (A, del/τ) and amplitude (B, amp/a) signal when the model eye is oriented to gaze angle $\theta = \phi = 0$ degrees

179 It is worth noting that while the time-of-flight signal falls off symmetrically from the
 180 center of gaze, the amplitude signal peaks at receiver $r = 20$, a result of occlusion from
 181 the nose-pad. We also note of the distribution spread for the time-of-flight and the
 182 amplitude signal captured by each receiver, which may be the result of measurement
 183 noise with our test-bench.

184 2.3.2. Gradient Boosted Regression Trees

185 From a machine learning perspective, the task of learning a gaze estimation model
 186 H is categorized as a supervised regression problem. Gradient Boosting Regression
 187 Trees (GBRT) are a powerful class of boosting algorithms for classification and regression
 188 tasks, which combine output from several weak learners into a powerful estimator.
 189 Specifically, GBRT considers additive models of the form: $F_m(x) = F_{m-1}(x) + h_m(x)$,
 190 where h_m are the basis functions modeled as small regression trees of fixed size. For
 191 each boosting iteration, a new boosting tree is added to the GBRT model, F . For our
 192 problem, we train two separate GBRT models to independently estimate the response

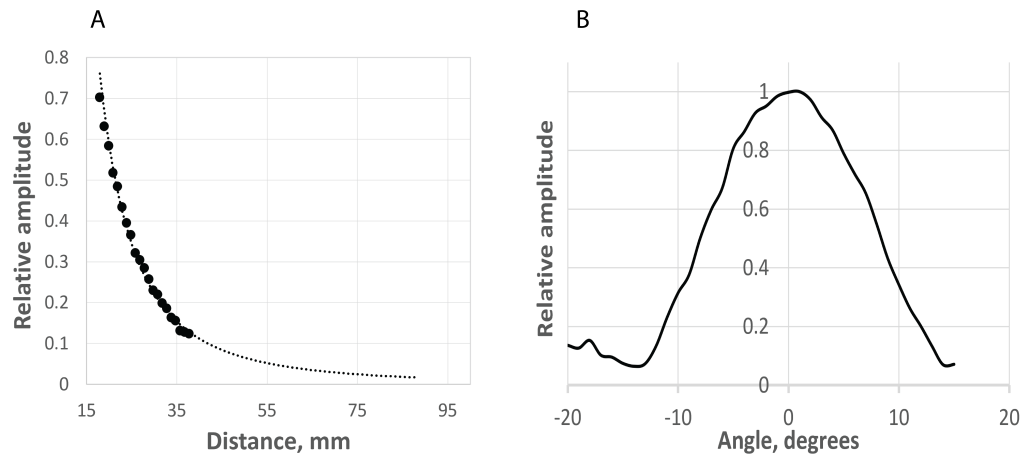


Figure 6. CMUT sensor characterization. A: Attenuation, obtained with Test bench 1, see Figure 2A; B: Directionality, obtained with Test bench 2, see Figure 2B.

193 in the horizontal and vertical dimensions: $\mathbf{Y} = (\theta, \phi)$ as function of the input features,
 194 $\mathbf{X} = (\tau^r, a^r)$. Assuming the GBRT model is comprised of M regression trees with
 195 T_m leaf nodes per regression tree, the GBRT model for each of the gaze regressor is
 196 given as: $F^y(\mathbf{X}, w^y) = w_0^y + \sum_{m=1}^M \sum_{j=1}^{T_m} w_{jm}^y I(\mathbf{X} \in R_{jm}^y)$, where $y = \{\theta, \phi\}$ and R_{jm}^y
 197 represents the j^{th} disjoint partitioning of the input space for the m^{th} regression tree
 198 for the regressor variable, y . The GBRT model weights are estimated from data as
 199 follows: $w^* = w \frac{1}{N} \sum_i^N L(y_i, F(\mathbf{X}_i, w))$ where, L is the squared error loss function. For an
 200 exhaustive description of GBRT, see [16,17].

201 Both the GBRT and linear regression models are trained to minimize the mean-
 202 squared-error between the estimated gaze-vector and the predicted gaze-vector and
 203 we report model performance in terms of root-mean-squared model error on a 5-fold
 204 cross-validation set. In addition we report the adjusted- R^2 as a goodness-of-fit measure
 205 for regression models.

206 3. Results

207 In this section we present findings from our modeling study as well as experiments
 208 conducted using the three benchtop setups described in Section 2.2.

209 We begin by presenting our findings on the CMUT sensor characterization. Data
 210 collected using test bench setup 1 allowed us to investigate the decay characteristics of
 211 the ultrasound signal in air, see Figure 6A. As expected, the ultrasound signal decays
 212 exponentially as a function of distance. An extrapolated fit shows it decays to zero.
 213 The distance axis shows the distance between the pair of transducers and the target
 214 (Figure 2A). Actual travel distance is twice this measurement. The range is similar to the
 215 distances for transducers mounted on eye glasses frames, our use case scenario.

216 Data collected using test bench 2 (Figure 2B) allowed us investigate whether the
 217 CMUT transducers exhibit directionality. Our findings are reported in Figure 6B. The
 218 CMUT transducers indeed exhibit directionality with an emission cone of 10 degrees.
 219 This applies to the transducers in both transmit and receive mode.

220 Based on the above findings we conclude that the strength of ultrasound signal at
 221 the receiver CMUT transducer will depend on two factors: distance and incident angle.
 222 As such, we believe that the amplitude of the ultrasound signal at the receiver contains
 223 relevant information to contribute to our ability to estimate gaze and as shown below,
 224 our findings indeed support this claim.

225 Our modeling study explored two situations: an omnidirectional transducer and
 226 one that mimics the properties of our CMUTs, see Figure 7. 131,072 rays are released
 227 from a point source in each case. The rationale for exploring the two situations is that
 228 while our CMUTs fit our needs, single crystal piezo transducers may provide a robust,

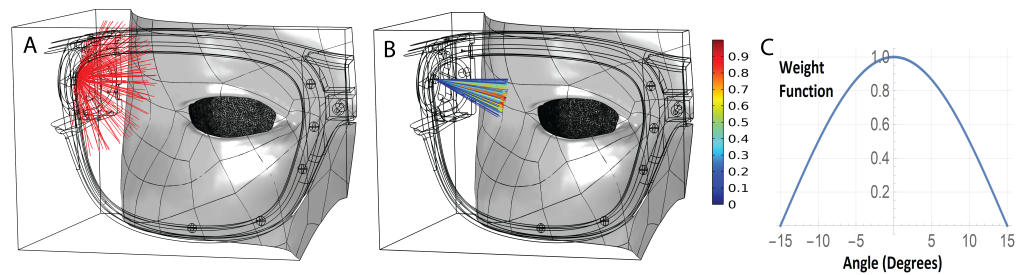


Figure 7. Sensor directionality options. A. Omnidirectional; B. Directional. The color of the ray indicates its intensity; C. Weigh function to mimic the transducer native curve (shown in Figure 6B).

229 inexpensive alternative. They are omnidirectional but can be turned into a directional
 230 device by adding baffles. In terms of size, they would be slightly larger (2.5mm instead
 231 of 1mm in our frequency range).

232 We implemented the sensor native curve by releasing rays with a uniform density
 233 distribution and assigning weight functions to the rays based on rays angle of emission
 234 and reception. The weight function of the rays is $\cos(\min(\alpha*(90/15),90^\circ))$. Alpha is
 235 the angle between the ray and transducer direction.

236 In the directional case we used a similar approach to account for the receiver native
 237 curve. We assigned a similar weight function to the acoustic rays that reach receivers
 238 based on the angle between the incoming acoustic rays and the sensor direction of each
 239 receiver. Therefore, each ray has two weight functions. One weight function is assigned
 240 initially when the ray is released, another weight function is assigned when the ray is
 241 detected by a receiver. The product of the two weight functions is applied. If the angle
 242 between a ray (that reaches a sensor) and the receiver direction is more than 15 degrees
 243 then the ray is not detected (its weight function is zero). If this angle is zero then the
 244 weight function is 1.

245 Figure 8 shows a comparison of the predicted signal at our sensor locations for
 246 directional and omnidirectional transducers. The left and right panels correspond to
 247 thirty degree rotations of the eye to the left (towards the nose) or right. In the case of
 248 omnidirectional transducers the differences between gazes are small. Differences are
 249 more pronounced for directional transducers. Peaks are also better defined with direc-
 250 tional transducers. Late peaks resulting from longer paths due to multiple reflections
 251 are minimized. It is to be noted that such late peaks would be ignored in our analysis,
 252 as we only use the time to peak and peak amplitude for the first peak detected in a
 253 given channel. With the same total number of rays (transducer power), receiver sensors
 254 with a directional transducer have higher signal strength than receiver sensors with
 255 an omnidirectional transducer. We ran the same models for a straight gaze as well as
 256 up/down twenty degree rotations (data not shown), and obtained similar results. Taken
 257 together the directional transducers perform better to resolve gaze.

258 Next we looked at where on the frame we might detect a signal, and why. Figure
 259 9, left, shows signal intensity around the frame. Areas in red have a higher chance of
 260 detecting rays reflected from the eye. Rays reflected off the glasses or skin are ignored.
 261 The center panel shows the path taken by the rays that reach receiver 6. Some of the rays
 262 arrive after multiple reflections from the skin and glasses. The right panel provides a
 263 detailed view of the direction of rays reaching receiver 6 (sphere). Sensor direction is
 264 shown with the solid black line. The majority of these rays will not be detected by the
 265 receiver due to the narrow angle of detection dictated by the receiver native curve. If
 266 we were using omni-directional transducers, rays arriving after two or more reflections
 267 would broaden the signal or create multiple peaks. Directional transducers allow us to
 268 reject unwanted signals before they are counted.

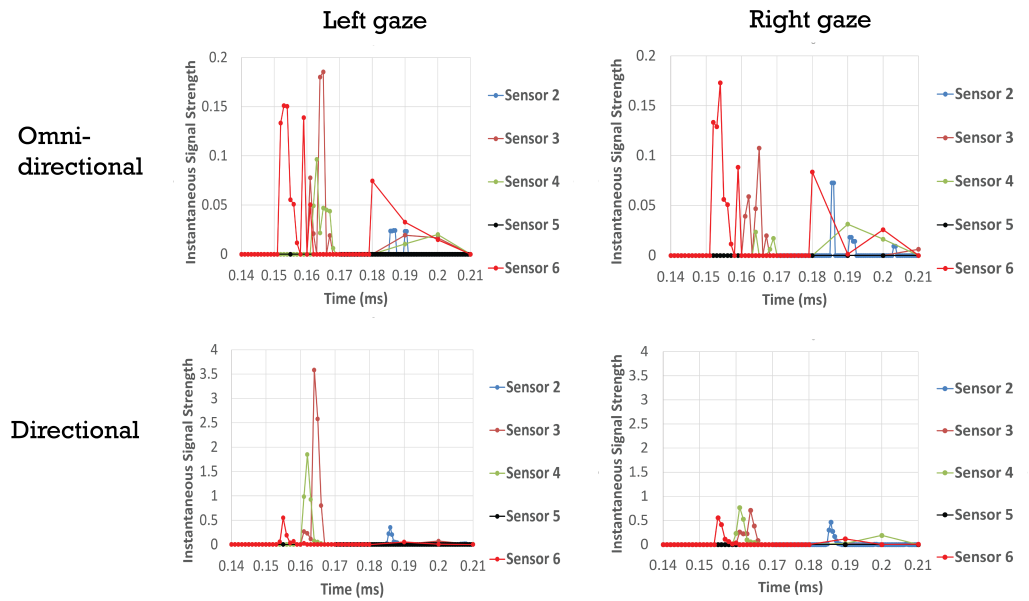


Figure 8. Modeled responses of ultrasound signal measured by the receiver with omni-directional and directional transducers. Left gaze represents a 30 degree rotation away from the nose, Rights gaze is 30 degrees towards the nose. Transducer 1 is not shown, as it operates only in transmit mode, see Figure 1.

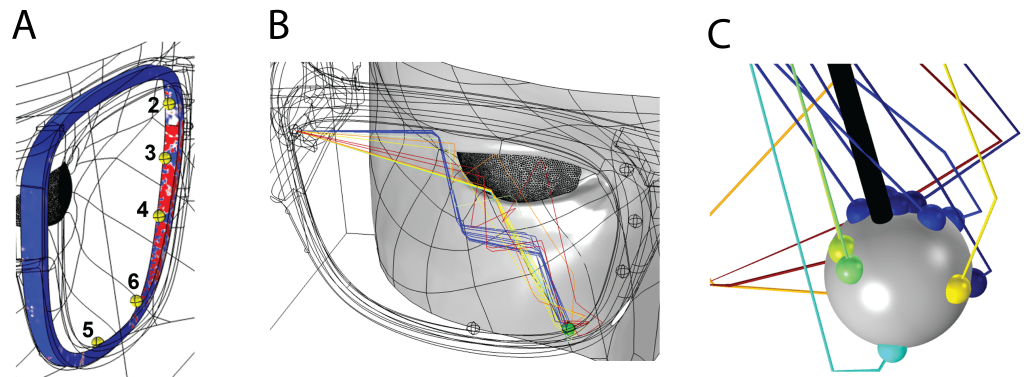


Figure 9. Signal detection: A. Rays reflected from the eye; B. Subset of rays arriving at sensor 6; C. closeup of rays arriving at sensor 6.

269 We next report findings from training a GBRT model on data collected using the
 270 third test bench setup (see Figure 3). For each model eye position on the goniometer,
 271 θ, ϕ , for a fixed receiver transducer position (180 degrees) and for a set of 19 transmit
 272 transducer positions, we fire the ultrasound test signal 50 times, at 2 kHz and record the
 273 raw receiver signal (see Figure 4 top row). In order to increase the strength of ultrasound
 274 response at the receiver we average 10 traces of the raw response signals at a time, to
 275 effectively generate 5 averaged ultrasound response signals, in effect acquiring data at
 276 200 Hz. The averaged response signal is passed through a Butterworth bandpass filter
 277 and we extract two ultrasound signal features: time of flight (τ) and the amplitude at
 278 peak (a), as explained in Section 2.3. In total for each model eye position, we generate a
 279 total of 45 samples for each model eye position on the goniometer over the duration of
 280 the study. For the set of 36 model eye positions, we produce a total of 1620 data samples.

281 We train a GBRT model on these data samples, performing a 5-fold cross-validation
 282 study. The model performance is reported using an adjusted R^2 score [18] and the gaze
 283 RMSE error in degrees. Hyper-parameter search on the GBRT model parameters that
 284 produced the best adjusted R^2 score for 5-fold cross-validation are reported in Table 1.
 285 We obtain gaze RMSE error of 0.965 ± 0.178 and mean adjusted R^2 score of 90.2 % with

Table 1. Hyper-parameters for the trained GBRT model. See [xgboost parameters in sci-kit learn](#) for explanation of these hyper-parameters

Hyper-parameters (XGBoost GBRT Model)	
learning rate	0.0825
max tree depth	5
regression trees	750
min. child weight	23
α regularization	0.01
λ regularization	1

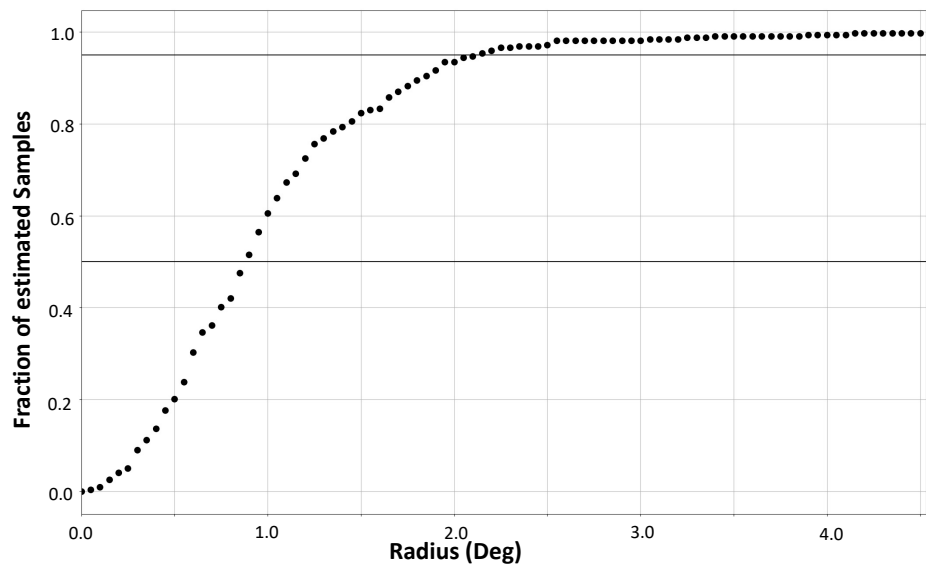


Figure 10. Sensitivity for gaze resolution: fraction of gaze estimates falling within a given radius of corresponding ground truth gaze values.

286 a standard deviation of 4.6, suggesting that almost 90 % of the data fit the regression
 287 model. We also perform a similar analysis using a linear regression model and the results
 288 are reported in Table 2. Nonlinear modeling of the problem through GBRT produced
 289 an improvement in performance for RMSE of $\approx 18\%$ and goodness-of-fit improvement
 290 of $\approx 5.7\%$, in support of our claim that the occlusions introduce nonlinearities in the
 291 ultrasound signals captured by the CMUT receivers, that can be best captured using a
 292 nonlinear regression model.

Table 2. Regression Model for Gaze Estimation. Numbers are presented in terms of mean \pm std. dev

	adjusted- R^2	RMSE
Gradient Boosted Tree	90.2 \pm 4.6	0.965 \pm 0.178
Linear Regression	85.3 \pm 07.6	1.177 \pm 0.236

293 Residuals analysis confirmed that the estimates obtained using the GBRT model
 294 are un-biased (data not shown). In Figure 10, we show the plot of the fraction of GBRT
 295 estimated gaze values that fall within an epsilon-ball of given radius (degrees). We see
 296 that $\approx 50\%$ of estimated gaze values fall within an epsilon ball of radius 0.8 degrees and
 297 $\approx 90\%$ of estimated gaze values fall within an epsilon-ball of radius 2 degrees. Based on
 298 these findings, we conclude that using CMUT ultrasound sensors, we can expect gaze
 299 resolution of up to 2 degrees.

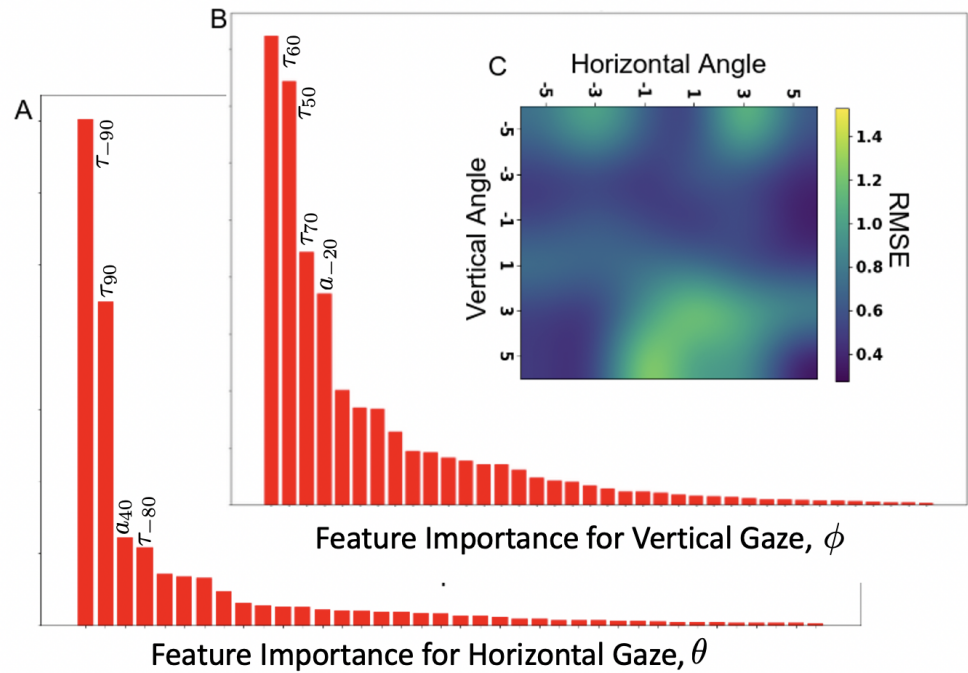


Figure 11. Feature importance and mean accuracy of GBRT models to estimate gaze. A. Horizontal gaze, B. Vertical gaze, C. Error as a function of gaze angle.

300 In Figure 11A and 11B, we show feature importance for the GBRT tree models
 301 trained to estimate the model eye gaze coordinates, θ (horizontal gaze) and ϕ (vertical
 302 gaze). We can see that the top two features for both horizontal and vertical gaze GBRT
 303 model are time of flight ultrasound signal. It has been our observation that while the
 304 time of flight component of ultrasound signal contains dominant information signal to
 305 estimate gaze (95 % contribution to the regression score), the amplitude signal is also
 306 an important contributor for GBRT model to produce an adjusted- R^2 score close to 90
 307 %. In order to test this observation, we trained GBRT model using just the ultrasound
 308 time-of-flight feature and another GBRT model using just the ultrasound amplitude
 309 feature. The findings are: GBRT model trained using time-of-flight features, produces
 310 an adjusted R^2 score of 85.4 ± 5.2 , where as the GBRT model trained using only the
 311 amplitude feature produces an adjusted R^2 score of 78.6 ± 8.2 . In Figure 9C, we show the
 312 mean-RMSE error (across all CV-folds) for the GBRT model. The error is biased towards
 313 the lower half of vertical gaze, primarily resulting from occlusions.

314 4. Discussion

315 This study is the first experimental demonstration of use for ultrasound sensors in
 316 gaze estimation. We show that ultrasonic transducers can effectively produce signals use-
 317 ful to resolve eye gaze, as defined by the static orientation of model eye on a goniometer,
 318 within the range tested, ± 5 degrees in both up/down (θ) and left/right (ϕ) directions.
 319 This range reflects the full deflection of our goniometer. We plan on expanding the range
 320 in future studies.

321 Prior to embarking on our experiments with bench-top setup we conducted ray
 322 tracing modeling. This modeling helped us refine our test bench design, procedures, and
 323 analysis. First, it pointed to the utility of directional over omnidirectional transducers.
 324 Second, it informed us on where, given a source location, we can expect a signal around
 325 the glasses frame. Finally, it provided information on the signals we need to measure
 326 from our bench-top experiments: the time to peak (indicative of distance traveled), and
 327 the amplitude. Due to attenuation in air the amplitude decreases with travel distance.

328 Our modeling indicated that amplitude also carries a signal based on the angle of
329 incidence. This is further evidence for using directional instead of omnidirectional
330 transducers.

331 Our GBRTs show that both amplitude and time of flight contribute to our ability
332 to estimate gaze. This is a new finding as previous modeling work dealt with time of
333 flight alone. As mentioned in our modeling section, two factors contribute to amplitude:
334 attenuation and the incident angle of the incoming sound. One way to compensate for
335 attenuation is to use the time-gain correction built in our amplifier, increasing gain over
336 time to compensate for the signal attenuation with longer distances. When we did this
337 (data not shown) our model performed slightly worse. This indicates that attenuation
338 plays a role in our ability to estimate gaze, and would favor the use of high frequency
339 transducers.

340 For this proof of concept we chose to average ten individual tests prior to filtering
341 the signal and extracting peak and amplitude. This reduces the eye tracking acquisition
342 speed from a maximum of 2kHz to 200Hz, which may not be sufficient to track saccadic
343 eye motion. While this study focused on primarily testing the hypothesis that ultrasound
344 signals can be leveraged to estimate gaze, in future works we will explore avenues to
345 investigate the use for ultrasound in tracking fast eye motion. Specifically, we plan on
346 using a fast-moving model eye coupled with multiple receivers operating at 2kHz. The
347 GBRT models will be adapted so we can test the potential of ultrasound for fast eye
348 tracking to resolve saccades.

349 We are interested in investigating the feasibility for using ultrasound sensors for eye
350 tracking in virtual and augmented reality devices. In addition to sampling speed, power
351 consumption is an important factor to consider. The transducers are very low power, in
352 the milliwatt range. Our current system utilizes a high speed A/D converter. This can
353 be replaced with a low power peak detection circuit. On the compute side, GBRTs are
354 considered low compute. This is in particular true for run time on multi-core machines.
355 Specifically, the run time compute complexity for GBRT models is $O(pn_{\text{trees}}/C)$, where
356 p represent the number of input features and n_{trees} are the number of regression trees
357 and C is the number of compute cores on a given machine. For $n_{\text{trees}}/C \sim 1$, the run
358 time complexity for GBRT is on parity with linear regression models, at $O(p)$.

359 In summary, this study presents data driven proof-of-principle findings to support
360 the claim that ultrasound sensors can be used for gaze estimation.

361 **Author Contributions:** Conceptualization, A.G. and S.S.T.; methodology, A.G. and S.S.T.; valida-
362 tion, S.S.T. ; formal analysis, A.G. and S.S.T.; investigation, A.G. and S.S.T.; resources, A.G. and
363 S.S.T.; data curation, S.S.T.; writing—original draft preparation, A.G.; writing—review and editing,
364 A.G. and S.S.T.; visualization, A.G. and S.S.T.; supervision, S.S.T. and A.G.; project administration,
365 A.G. and S.S.T.; funding acquisition, S.S.T. All authors have read and agreed to the published
366 version of the manuscript.

367 **Funding:** This research was completely funded by Facebook Reality Labs Research, a subsidiary
368 of Facebook, Inc.

369 **Institutional Review Board Statement:** Not applicable, no human data was used in this research

370 **Informed Consent Statement:** Not applicable.

371 **Data Availability Statement:** None

372 **Acknowledgments:** We gratefully thank Robert Cavin and Facebook Reality Labs Research for
373 supporting this work. We also extend our thank you to Alireza Karmani and Nagi Elabbasi from
374 Veryst Engineering LLC. for their assistance in comsol-modeling study.

375 **Conflicts of Interest:** S.S.T is an employee of the funding agency, while A.G is a contract employee
376 hired by the funding agency. Both S.S.T and A.G have been involved in the design of the study,
377 dataset collection, analyses and interpretation of the data, writing of the manuscript, and the
378 decision to publish the results.

References

1. Golard, A.; Talathi, S. Ultrasound for gaze estimation. *Pattern Recognition. ICPR International Workshops and Challenges: Virtual Event, 2021*.
2. Kar, A.; Corcoran, P. A Review and Analysis of Eye-Gaze Estimation Systems, Algorithms and Performance Evaluation Methods in Consumer Platforms. *IEEE Access* **2017**, *5*, 16495–16519.
3. Cheng, D.; Vertegaal, R. An Eye for an Eye: A Performance Evaluation Comparison of the LC Technologies and Tobii Eye Trackers. In *Proceedings of the 2004 Symposium on Eye Tracking Research & Applications, San Antonio, TX, USA, 22–24 March 2004*; p. 61.
4. Sheehy, C.K.; Yang, Q.; Arathorn, D.W.; Tiruveedhula, P.; de Boer, J.F.; Roorda, A. High-speed, image-based eye tracking with a scanning laser ophthalmoscope. *Biomed. Opt. Express* **2012**, *3*, 2611–2622.
5. Cornsweet, T.N.; Crane, H.D. Accurate two-dimensional eye tracker using first and fourth Purkinje images. *J. Opt. Soc. Am.* **1973**, *63*, 921–928.
6. Gramatikov, B.I.; Zalloum, O.H.Y.; Wu, Y.K.; Hunter, D.G.; Guyton, D.L. Directional eye fixation sensor using birefringence-based foveal detection. *Appl. Opt.* **2007**, *46*, 1809–1818.
7. Caldani, S.; Bucci, M.P.; Lamy, J.C.; Seassau, M.; Bendjemaa, N.; Gadel, R.; Gaillard, R.; Krebs, M.O.; Amado, I. Saccadic eye movements as markers of schizophrenia spectrum: Exploration in at-risk mental states. *Schizophr. Res.* **2017**, *181*, 30 – 37.
8. Termsarasab, P.; Thammongkolchai, T.; Rucker, J.C.; Frucht, S.J. The diagnostic value of saccades in movement disorder patients: a practical guide and review. *J. Clin. Mov. Disord.* **2015**, *2*.
9. Gibaldi, A.; Sabatini, S. The saccade main sequence revised: A fast and repeatable tool for oculomotor analysis. *Behav. Res.* **2021**, *53*, 167–187.
10. Scally, B.M.; Perek, D.R. Ultrasound/Radar for Eye Tracking. U.S. Patent 2017/0261610 A1, 2017.
11. Sánchez-Ferrer, M.L.; Grima-Murcia, M.D.; Sánchez-Ferrer, F.; Hernández-Peñalver, A.I.; Fernández-Jover, E.; del Campo, F.S. Use of Eye Tracking as an Innovative Instructional Method in Surgical Human Anatomy. *J. Surg. Educ.* **2017**, *74*, 668–673.
12. Kaputa, D.; Enderle, J. An Ultrasound Based Eye Tracking System. *J. Biomed. Eng. Med. Dev.* **2016**, *1*, 1–4.
13. Hansen, D.W.; Ji, Q. In the Eye of the Beholder: A Survey of Models for Eyes and Gaze. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2010**, *32*, 478–500.
14. Khury-Yacub, B.; Oralkan, O. Capacitive micromachined ultrasonic transducers for medical imaging and therapy. *J. Micromech. Microeng.* **2011**, *21*, 054004–054014.
15. Blackstock, D.T. *Fundamentals of Physical Acoustics*; John Wiley & Sons: Hoboken, NJ, USA, 2000.
16. Hastie, T.; Tibshirani, R.; Friedman, J. *The Elements of Statistical Learning*; Springer: Berlin/Heidelberg, Germany, 2011.
17. Ridgeway, G. Generalized Boosted Models: A guide to gbm package. *Update* **2007**, *1*, 2007.
18. Dodge, Y. *The Concise Encyclopedia of Statistics*; Springer: Berlin/Heidelberg, Germany, 2010.

