

000

001

002

003

Audi-Exchange: AI-Guided Hand-based Actions to Assist Human-Human Interactions for the Blind and the Visually Impaired

004

005

006

007

008

009

010

011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035

036

037

038

039

040

041

042

043

044

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090

091

092

093

094

095

096

097

098

099

0100

0101

0102

0103

0104

0105

0106

0107

0108

0109

0110

0111

0112

0113

0114

0115

0116

0117

0118

0119

0120

0121

0122

0123

0124

0125

0126

0127

0128

0129

0130

0131

0132

0133

0134

0135

0136

0137

0138

0139

0140

0141

0142

0143

0144

0145

0146

0147

0148

0149

0150

0151

0152

0153

0154

0155

0156

0157

0158

0159

0160

0161

0162

0163

0164

0165

0166

0167

0168

0169

0170

0171

0172

0173

0174

0175

0176

0177

0178

0179

0180

0181

0182

0183

0184

0185

0186

0187

0188

0189

0190

0191

0192

0193

0194

0195

0196

0197

0198

0199

0200

0201

0202

0203

0204

0205

0206

0207

0208

0209

0210

0211

0212

0213

0214

0215

0216

0217

0218

0219

0220

0221

0222

0223

0224

0225

0226

0227

0228

0229

0230

0231

0232

0233

0234

0235

0236

0237

0238

0239

0240

0241

0242

0243

0244

0245

0246

0247

0248

0249

0250

0251

0252

0253

0254

0255

0256

0257

0258

0259

0260

0261

0262

0263

0264

0265

0266

0267

0268

0269

0270

0271

0272

0273

0274

0275

0276

0277

0278

0279

0280

0281

0282

0283

0284

0285

0286

0287

0288

0289

0290

0291

0292

0293

0294

0295

0296

0297

0298

0299

0300

0301

0302

0303

0304

0305

0306

0307

0308

0309

0310

0311

0312

0313

0314

0315

0316

0317

0318

0319

0320

0321

0322

0323

0324

0325

0326

0327

0328

0329

0330

0331

0332

0333

0334

0335

0336

0337

0338

0339

0340

0341

0342

0343

0344

0345

0346

0347

0348

0349

0350

0351

0352

0353

0354

0355

0356

0357

0358

0359

0360

0361

0362

0363

0364

0365

0366

0367

0368

0369

0370

0371

0372

0373

0374

0375

0376

0377

0378

0379

0380

0381

0382

0383

0384

0385

0386

0387

0388

0389

0390

0391

0392

0393

0394

0395

0396

0397

0398

0399

0400

0401

0402

0403

0404

0405

0406

0407

0408

0409

0410

0411

0412

0413

0414

0415

0416

0417

0418

0419

0420

0421

0422

0423

0424

0425

0426

0427

0428

0429

0430

0431

0432

0433

0434

0435

0436

0437

0438

0439

0440

0441

0442

0443

0444

0445

0446

0447

0448

0449

0450

0451

0452

0453

0454

0455

0456

0457

0458

0459

0460

0461

0462

0463

0464

0465

0466

0467

0468

0469

0470

0471

0472

0473

0474

0475

0476

0477

0478

0479

0480

0481

0482

0483

0484

0485

0486

0487

0488

0489

0490

0491

0492

0493

0494

0495

0496

0497

0498

0499

0500

0501

0502

0503

0504

0505

0506

0507

0508

0509

0510

0511

0512

0513

0514

0515

0516

0517

0518

0519

0520

0521

0522

0523

0524

0525

0526

0527

0528

0529

0530

0531

0532

0533

0534

0535

0536

0537

0538

0539

0540

0541

0542

0543

0544

0545

0546

0547

0548

0549

0550

0551

0552

0553

0554

0555

0556

0557

0558

0559

0560

0561

0562

0563

0564

0565

0566

0567

0568

0569

0570

0571

0572

0573

0574

0575

0576

0577

0578

0579

0580

0581

0582

0583

0584

0585

0586

0587

0588

0589

0590

0591

0592

0593

0594

0595

0596

0597

0598

0599

0600

0601

0602

0603

0604

0605

0606

0607

0608

0609

0610

0611

0612

0613

0614

0615

0616

0617

0618

0619

0620

0621

0622

0623

0624

0625

0626

0627

0628

0629

0630

0631

0632

0633

0634

0635

0636

0637

0638

0639

0640

0641

0642

0643

0644

0645

0646

0647

0648

0649

0650

0651

0652

0653

0654

0655

0656

0657

0658

0659

0660

0661

0662

0663

0664

0665

0666

0667

0668

0669

0670

0671

0672

0673

0674

0675

0676

0677

0678

0679

0680

0681

0682

0683

0684

0685

0686

0687

0688

0689

0690

0691

0692

0693

0694

0695

0696

0697

0698

0699

0700

0701

0702

0703

0704

0705

0706

0707

0708

0709

0710

0711

0712

0713

0714

0715

0716

0717

0718

0719

0720

0721

0722

0723

0724

0725

0726

0727

0728

0729

0730

0731

0732

0733

0734

0735

0736

0737

0738

0739

0740

0741

0742

0743

0744

0745

0746

0747

0748

0749

0750

0751

0752

0753

0754

0755

0756

0757

0758

0759

0760

0761

0762

0763

0764

0765

0766

0767

0768

0769

0770

0771

0772

0773

0774

0775

0776

0777

0778

0779

0780

0781

0782

0783

0784

0785

0786

0787

0788

0789

0790

0791

0792

0793

0794

0795

0796

0797

0798

0799

0800

0801

0802

0803

0804

0805

0806

0807

0808

0809

0810

0811

0812

0813

0814

0815

0816

0817

0818

0819

0820

0821

0822

0823

0824

0825

0826

0827

0828

0829

0830

0831

0832

0833

0834

0835

0836

0837

0838

0839

0840

0841

0842

0843

0844

0845

0846

0847

0848

0849

0850

0851

0852

0853

0854

0855

0856

0857

0858

0859

0860

0861

0862

0863

0864

0865

0866

0867

0868

0869

0870

0871

0872

0873

0874

0875

0876

0877

0878

0879

0880

0881

0882

0883

0884

0885

0886

0887

0888

0889

0890

0891

0892

0893

0894

0895

0896

0897

0898

0899

0900

0901

0902

0903

0904

0905

0906

0907

0908

0909

0910

0911

0912

0913

0914

0915

0916

0917

0918

0919

0920

0921

0922

0923

0924

0925

0926

0927

0928

0929

0930

0931

0932

0933

0934

0935

0936

0937

0938

0939

0940

0941

0942

0943

0944

0945

0946

0947

0948

0949

0950

0951

0952

0953

0954

0955

0956

0957

0958

0959

0960

0961

0962

0963

0964

0965

0966

0967

0968

0969

0970

0971

0972

0973

0974

0975

0976

0977

0978

0979

0980

0981

0982

0983

0984

0985

0986

0987

0988

0989

0990

0991

0992

0993

0994

0995

0996

0997

0998

0999

1000

Audi-Exchange: AI-Guided Hand-based Actions to Assist Human-Human Interactions for the Blind and the Visually Impaired

Daohan Lu^{1,4} and Yi Fang^{1,2,3}

¹ Multimedia and Visual Computing Lab, New York University, New York, United States.

² Tandon School of Engineering, New York University, New York, United States.

³ Department of Electrical and Computer Engineering, New York University, Abu Dhabi, United Arab Emirates.

⁴ College of Arts and Science, New York University, New York, United States.

Abstract. Vision loss or low vision poses significant challenges to blind-or-visually-impaired (BVI) individuals when interacting with humans and objects. Although many apps and assistive devices can help them better interact with the environment and objects, the current state of assistive technology leaves human-human interaction needs of the BVI largely unaddressed. Because of this, we introduce a new wearable mobile assistive platform, named Audi-Exchange, to address part of the problem. Developed with mobile-optimized computer vision and audio engineering techniques, Audi-Exchange facilitates a specific area of human-human interaction by helping the BVI user accurately locate another person's hand with spatial audio in order to pass objects over to or receive objects from the other person. Audi-Exchange differs from existing academic and commercial assistive technologies in that it is intuitive to use and non-intrusive when worn. We conduct several experiments to investigate Audi-Exchange's effective as an assistive human-human interaction tool and discover encouraging results.

Keywords: Assistive Technology, Human-Human Interaction, Visually Impaired

1 Introduction

In the workplace or daily life, people need to interact with other people for various social and working needs. However, people who are blind or visually impaired (BVI) have fewer or lesser tools available to them for interacting with another person, namely those that rely on a clear vision, such as sensing eye contact or reading body language. This leads to a harder time when they interact with other people, especially with those who are not familiar with the appropriate practices for interacting with BVI individuals. This can lead to strong real-world impacts. In the workplace, for instance, visual impairment leads to lower efficiency when interacting and communicating with coworkers and customers, which makes employment more difficult. Thus, we wish to investigate

the feasibility and effectiveness of developing an assistive device to collect important visual cues in the environment and relays this information to the BVI user to help the BVI user interact with other sighted or BVI people. In this paper, we focus on facilitating human-human interaction that involves handing over of objects. The choice is due to the common situation in which a BVI person needs to exchange objects with another person in casual settings as well as in the workplace. For instance, a BVI person may need to accept a credit card from someone or offer a cup of water to someone. To help the BVI with these types of actions, we introduce the Audi-Exchange wearable mobile assistive platform. When Audi-Exchange is activated, it collects images from a camera and uses an efficient convolution neural network (CNN) to detect another person's hand within the camera's field of view. When a hand is detected, Audi-Exchange tracks its on-camera location and computes its corresponding 3D direction based on the camera's optical parameters. Next, the direction of the target hand is relayed to the BVI user through stereo headphones as an audio tone, which is processed by a head-related transfer function (HRTF) to appear as coming from the computed direction. In this fashion, we leverage computer vision and audio engineering to augment or substitute impaired human vision to allow the BVI to "hear" how they should move when handing over or receiving an object from another person. We assess Audi-Exchange's effectiveness and speed by evaluating it on a proof-of-concept system consisting of a personal computer and a camera through several experiments. Lastly, we show that the system could be ported to work with mainstream smartphone hardware with a few modifications to the algorithm and discuss the next step for Audi-Exchange.

2 Assistive Technology Landscape

When reviewing academic sources, we found an abundance of research that aimed to develop assistive devices for the blind and visually impaired. A large portion falls into the category of electronic travel aids (ETAs), which are devices that gather information, including nearby objects and obstacles, about the surrounding environment via dedicated sensors and transfer it to the user [12,21]. Some examples include an RFID-based indoor navigation system [16], the smart cane [24], the Path Force Feedback Belt [19], and Substitute Eyes (a hand-worn ultrasonic obstacle warning device) [11]. Another category that assistive devices fall into is vision substitution, which is using computer vision to capture images of the surroundings and transforming the raw visual information into a VI-friendly form, such as haptics or sounds. Vision substitution devices are more general-purpose as the information gathered by the system is less condensed compared to ETAs and relayed in a more direct and visual form to the VI, with the exception being text recognition sometimes is performed on images to read text printed in the surroundings to the user. For instance, FingerReader reads printed text to the user with a hand-worn camera [22]. Silicon Eyes informs the user of the color of objects nearby in addition to facilitating outdoor navigation [20]. Overall, however, even though many assistive devices developed in the academic

community exist that cater to different aspects of BVI assistance, most have significant drawbacks. Mainly, they can be cumbersome due to being equipped with complex sensors and can be intrusive due to needing to be worn at various locations around the body to effectively gather information. Also, few papers mention estimated production costs for the the assistive devices that were proposed, which leads to uncertainty as to whether these devices can be efficiently built and widely distributed.

There also exists a number of commercial assistive devices for the VI. A major portion are visual enhancement tools that are designed to be used by the low vision and not fully blind (perhaps because products made for the low vision enjoy a larger market compared to those made for the fully blind). The basic white cane is arguably the most widely used [10,14] and affordable tool useful to both the low vision and fully blind, but its function is limited by its short range as a direct extension of physical touch. Advancement in digital imaging sensors led to a variety of vision enhancement tools that process images of the surroundings using camera sensors and image processing techniques (such as zooming in or boosting contrast). These include IrisVision [7], Acesight [5], and eSight [6]. Vision enhancement can also be accomplished with certain accessibility apps running on Google Glasses [1] and Microsoft Hololens [2]. More advanced solutions exist (such as the OrCam MyEye 2.0 [8]) where visual information gathered by the camera sensors is condensed into audio notifications for text reading, identifying objects, and recognizing faces. Even though commercial solutions are generally more comfortable to use, they are generally very costly. The vision augmentation devices mentioned above come in the price range of around \$2,000-4,000 USD, with OrCam MyEye 2.0 [8] being the most expensive, costing \$4,250 at the time of writing. Despite having being somewhat common, few hardware-based assistive solutions made by researchers or commercial firms are designed to specifically cover the VI's human-human interaction needs.

On the other hand, purely software-based solutions can be run on common smartphones and are more affordable, with Microsoft Seeing AI and BlindSquare [4,9] being popular free vision enhancement and BVI navigation apps. However, software-based solutions generally also have a set of significant shortcomings as they tend to (1) lack on-board visual processing, which leads to functional reliance on online computing, or (2) like the hardware-based solutions discussed previously, they only facilitate interactions with the environment while the VI's human interaction needs are largely ignored. In summary, despite innovations in mobile information gathering and visual computing that enabled many innovative hardware and software assistive solutions for the BVI, there are still only few that are designed specifically to help the user have smoother, richer interactions with other people.

3 Motivation and Merit

Because existing assistive technology for the BVI is often costly to obtain, uncomfortable to use, and does not address human interaction needs (see fig. 1), we

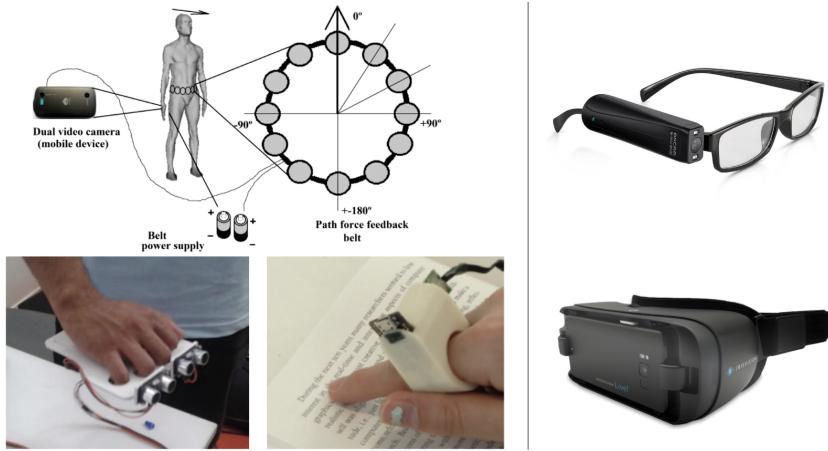


Fig. 1. Existing assistive device solutions for the BVI can be cumbersome, intrusive (first column), or costly (second column). Also, few devices are designed to address human-human interaction needs. The three images in the first column read from left to right and then top to bottom are: the Path Force Feedback Belt [19], Substitute Eyes [11], and FingerReader [22]. The two images in the second column read from top to bottom are: OrCam MyEye 2 [8] and IrisVision Live! [7], which sell for \$4,250 and \$2,950 respectively at the time of writing.

wish to design a new assistive device to address these problems. A 2020 comprehensive literature review on wearable device design reveals that “comfort”, “intuitiveness”, and “mobility” are the three most desired properties of a wearable device in terms of user experience [13]. Comfort corresponds to the minimization of discomfort or pain, so the device should have a small form factor, be worn on convenient locations, and have good thermal dissipation [13]. For “intuitiveness”, the human-computer interaction pattern should be easy to learn and it should resemble familiar interactive patterns used in everyday life. For “mobility”, the device should be lightweight and not obstruct the user’s physical activity. These three properties are not exclusive to health-sighted individuals and are arguably more important for blind or visually impaired individuals because they expect a higher degree of usability and reliability from assistive wearable devices. As a consequence, we design a new assistive device that specifically aims to facilitate human-human interactions while maximizing comfort, intuitiveness, and mobility when worn. We think that good user experience is critical to the actual rate of adoption in the BVI community and the social impact that results from it.

Named Audi-Exchange, we propose an assistive platform that uses real-time onboard mobile computer vision and audio engineering to help the BVI better handle hand-based human-human interactions. From the user’s perspective, Audi-Exchange sends spatial audio cues through a pair of headphones to help the user locate another person’s hand to hand over or receive an object. For example, it helps a blind cashier reach over the counter to receive a credit card from

a customer. From a technical standpoint, Audi-Exchange uses a smartphone's camera and processor and mobile-optimized computer vision algorithms to track the hand motion of a person next to the BVI user and relay this information in real-time to the BVI user through the stereo headphones via spatial audio. In this manner, Audi-Exchange is designed to be responsive, energy-efficient, intuitive, comfortable, and mobile. Unlike assistive technologies that rely on dedicated sensors and processors, all parts of Audi-Exchange can run on mainstream smartphones using onboard mobile CPU/GPU/Neural Processor as the "brain", the camera as the sensory input, and a pair of wired or wireless headphones as the sensory output. Regular headphones be optionally replaced with bone conduction headphones for better perception of ambient sounds.

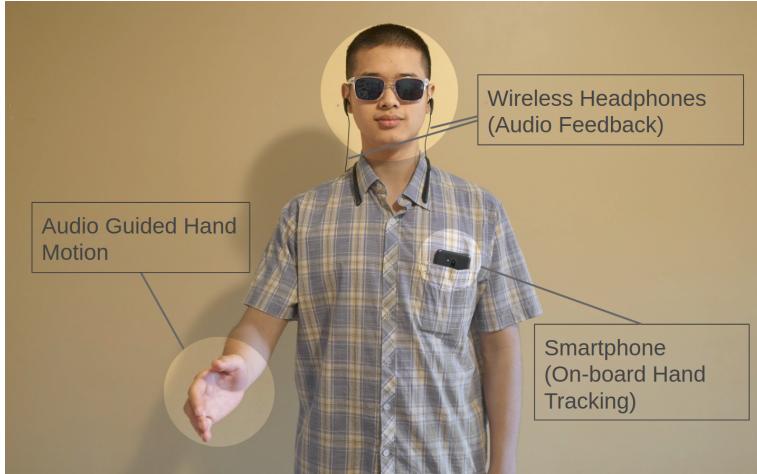


Fig. 2. Audi-Exchange Concept Design. By tracking the interactive partner's hand on a smartphone and sending spatial-audio-based interactive cues via wireless headphones, Audi-Exchange helps the BVI with hand-related human-human interactions in a lightweight and comfortable package.

4 Method

4.1 Hand Tracking

To implement Audi-Exchange, it is necessary to have a fast and accurate hand-tracking algorithm for the computer program to determine the location of the interactive partner's hand on an image. Hence, we base the hand tracking functional unit of Audi-Exchange on an efficient hand tracking algorithm [23] utilizing MobileNet (a mobile-optimized computer vision architecture) [15] and SSD (a computationally efficient object detection method) [17]. Running on a 2014



Fig. 3. Some images depicting the performance of the hand tracking algorithm by [23], which is the one we adopt for hand tracking in this paper. Image credit to [23].

13" Macbook Pro with a 2.6GHz dual-core Intel Core i5, we found the hand tracking algorithm to be able to process an average of 13 frames per second (FPS) with an input image size of 320 by 180 pixels. Such a speed is sufficient for tracking a non-fast-moving hand in real-time. Note that modern smartphones are typically equipped processors that can outperform the 2014 Macbook Pro by significant margins and many even come with dedicated neural processors that further speed up computer vision and machine learning applications, as is the case with the latest Apple iPhone (dubbed the "Neural Engine") and Google Pixel (dubbed the "Neural Core"). As a result, the 13-FPS speed that we recorded represents a bottom-line performance figure when hand tracking is run on older, less capable devices. We show in fig. 3 some test images of the hand tracking results as published by the author [23], while our test results are detailed in sec. 5.1.

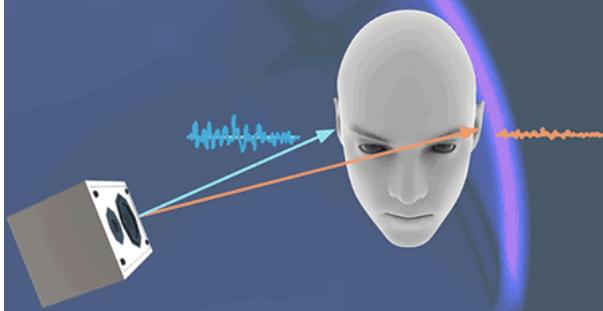
4.2 3D Audio

After a hand is detected, we wish to transmit this information to the BVI user in the form of an audio signal. To this end, we develop the 3D Audio module, which utilizes the OpenAL API [3] to embed 3D location cues into an arbitrary sound source. The method in which OpenAL (and other similar APIs like Resonance Audio [18]) accomplishes this is through the head-related transfer function (HRTF). The HRTF is a transfer function that takes in an artificial source of sound, denoted as $x(f)$, which is a function of frequency, and outputs separate channels of audio to the left (denoted $X_L(f)$) and right (denoted $X_R(f)$) ears of the listener to create directional sound [25]. When given HRTFs $H_L(f, \theta, \phi)$ and $H_R(f, \theta, \phi)$ and the source audio $x(f)$, the output $X_L(f)$ and $X_R(f)$ can be

270 computed as:

$$\begin{aligned} X_L(f) &= H_L(f, \theta, \phi) * x(f) \\ X_R(f) &= H_R(f, \theta, \phi) * x(f) \end{aligned} \quad (1)$$

271 where θ, ϕ are known constants representing the azimuth and elevation of the
 272 simulated audio source. Cartesian (x, y, z) coordinates can be used too, though
 273 only the radial direction will be taken into account by the HRTFs with the
 274 distance simply determining volume. In this way, original audio signals are filtered
 275 to carry directional cues to the listener. We utilized this fact to help the user
 276 locate a virtual audio source whose location corresponds to the actual location
 277 of the hand being tracked. When given a 3D location, the 3D Audio module
 278 plays a tone processed by the HRTF so that it appears as coming from the given
 279 location when the user hears it through stereo headphones. By relaying visual
 280 information to the user through sound, 3D Audio serves as a crucial part that
 281 makes it possible for the device to communicate visual information to a BVI
 282 user. Based on our testing, we found the 3D Audio module to consume negligi-
 283 ble computational resources when it processes and plays spatial audio, making
 284 it compatible with the mobile-optimized design of Audi-Exchange.



300 **Fig. 4.** A graphic illustrating the Interaural time differences (ITD)
 301 used in 3D audio software to help humans determine the horizontal location of low-frequency sounds.
 302 Image source: Resonance Audio [18]

303 304 4.3 Audi-Exchange

305 Finally, we combine the hand tracking and 3D audio modules to build the Audi-
 306 Exchange prototype. It works as follows: in sequential order, a camera captures
 307 an image, which is sent to the hand tracking module. After the hand location
 308 is determined, a virtual sound source is created with 3D audio to reflect the
 309 horizontal location of the detected hand. As the user hears the 3D sound through
 310 headphones, they acknowledge the location of the partner's hand in the real
 311 world. Then, the user moves their own hand to the perceived location in an
 312 attempt to touch the other person's hand. The process repeats until the action
 313 is completed as determined by the user. The flow of information from the BVI
 314

315 user's perspective is illustrated in 5. By offering auditory feedback in response
 316 to visual interactive cues (hand motion), Audi-Exchange repairs or augments
 317 hand-eye coordination for VBI user to efficiently hand objects over to or receive
 318 objects from another person. Although the Audi-Exchange test system used in
 319 this paper is set up for laptops and desktop computers, it is simple to port the
 320 program to mobile devices by migrating to open-source mobile-based computer
 321 vision and 3D audio algorithms (see sec. 6).



331 **Fig. 5.** Audi Exchange relays information about the location of the interactive partner's
 332 hand to the BVI user to guide the user toward handing over to or receiving an object
 333 from another person, who may be VIB or sighted.

337 5 Experiments

339 5.1 Hand Tracking Experiment



352 **Fig. 6.** Overview of the four test cases. From left to right: 1. Regular Lighting / Empty
 353 Hand, 2. Regular Lighting / Holding Pen, 3. Regular Lighting / Holding Cup, and 4.
 354 Poor Lighting / Empty Hand. We choose frames where the subject and hand are in
 355 the center of the image for clarity. Zoom in for better details.

358 In this experiment, we wish to find how reliably a mobile hand tracking program
 359 can pick up the interactive partner's hand. We designed four different test scenes

360 to investigate some realistic settings under which Audi-Exchange would be used.
 361 They are categorized by the lighting conditions and the object that the partner's
 362 hand is holding. The categories are 1. Regular Lighting / Empty Hand, 2. Regular
 363 Lighting / Holding Pen, 3. Regular Lighting / Holding Cup, and 4. Poor Lighting
 364 / Empty Hand. We show one image from each test case in fig. 6. The first two
 365 test scenes happen most frequently, with the interactive partner shown on the
 366 camera trying to receive an object by extending an empty hand or trying to
 367 give a small object with their hand. The third test scene exemplifies situations
 368 in which the partner offers a large object that occludes part of their hand when
 369 seen from the camera's perspective. The fourth test scene aims to see how well
 370 the hand can be tracked when the partner's hand is not properly lit and light is
 371 coming from behind the partner. All test scenes are filmed indoors with warm-
 372 colored LED lighting. A smartphone camera is used captures a person facing the
 373 camera from head-height, which simulates the BVI user's egocentric view. The
 374 person in view moves the extended hand around throughout the video to create
 375 natural motion blur and location variations. The smartphone camera used has
 376 a horizontal field of view of approximately 90° and an aspect ratio of 16 by 9,
 377 with the shutter speed set to 1/50 of a second, ISO to 1200, and white balance
 378 to "Auto". All images are resized to 320 by 180 pixels before being sent to the
 379 hand tracker.

380
 381 To quantitatively evaluate tracking performance, we record the number of
 382 true positives, false positives, and false negatives for each test scene (see tab.
 383 1). A true positive is defined as when the hand offering or receiving an object
 384 is detected or when both hands of the partner are detected (which is a trivial
 385 issue, as the hand we're interested in will be elevated and extended forward,
 386 meaning it can be distinguished from the other hand by picking the topmost
 387 hand). For similar reasons, exactly two closely overlapping bounding boxes on
 388 the correct hand is also considered a true positive, as the issue can be trivially
 389 solved by treating it as if two hands were detected. A false positive is defined as
 390 when the hand of interest is not detected, but an irrelevant object or the hand
 391 not of interest is. In the "False Positive (Wrong Hand)" column, we show the
 392 number of false positive frames in which the false positive is due to detecting
 393 only the irrelevant hand and no other objects. Lastly, a false negative is defined
 394 by not detecting any object or hand in a frame. Note that in all test cases, all
 395 frames include the hand of interest, so the true negative statistic is not available
 396 (or, equivalently, it has 0 instances). For each test case, there are 50 randomly
 397 selected frames from a video sequence. A frame has an equal chance of having
 398 the person and hand appear in the middle, left, and right parts of the image. A
 399 human annotator determines whether a frame is considered a true positive, false
 400 positive, false positive – wrong, or false negative by observing the visual output
 401 of the hand tracking algorithm. The annotator moves on the next frame when
 402 ready. The quantitative result is shown in tab. 1. Select qualitative hand tracking
 403 outputs are shown in fig. 7. We include a video segment in the supplementary
 404 materials demonstrating the hand tracking software used in this experiment.

	True Positive	False Positive	False Positive (Wrong Hand)	False Negative
Test 1	46	3	2	1
Test 2	47	2	0	1
Test 3	21	22	15	7
Test 4	0	50	0	0

Table 1. The recorded instances of true positives, false positives, false positives due to detecting the irrelevant hand, and false negatives for each hand tracking test scene. A total of 50 frames were examined for each test scene. Tests 1-4 correspond to Regular Lighting / Empty Hand, Regular Lighting / Holding Pen, Regular Lighting / Holding Cup, and Poor Lighting / Empty Hand, in this order.

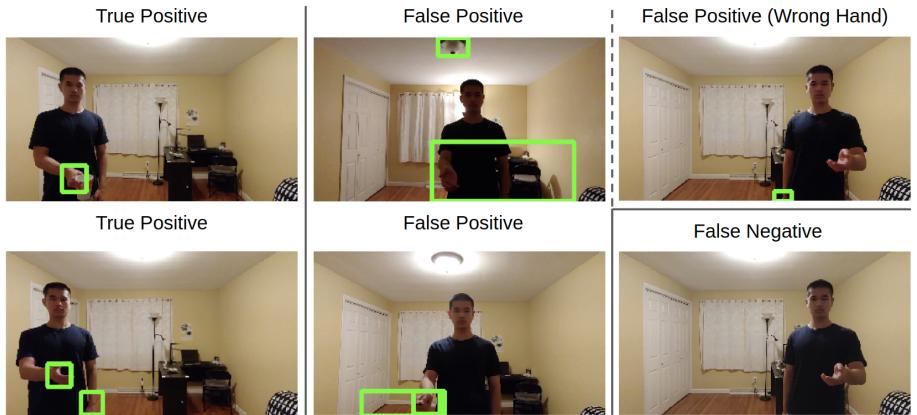
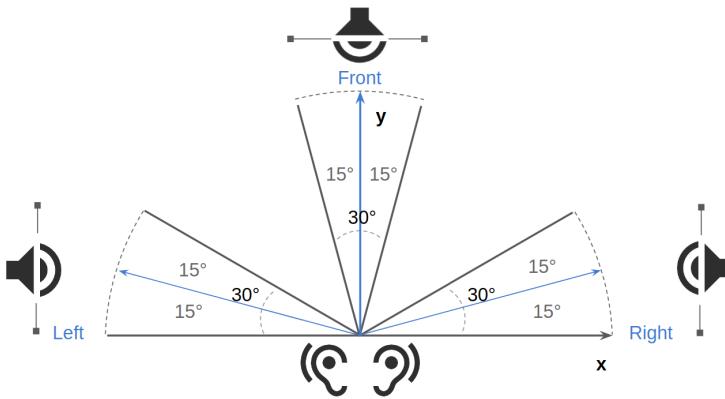


Fig. 7. Some examples of the visual output of the hand tracking algorithm. Each is annotated as True Positive, False Positive, False Positive (Wrong Hand), or False Negative Classes. Note that a “False Positive (Wrong Hand)” image counts toward False Positive also. Zoom in for more details.

We found the hand tracking algorithm to be robust under regular lighting conditions for tracking both the receiving and giving hand poses in the ‘Empty hand’ and ‘Holding pen’ test scenes. Even when a significant portion of the hand is occluded by an object, which is the case for the ‘Holding Cup’ test scene, the algorithm was still able to detect the hand of interest 42% of the time. If the hand of interest was not detected, the other hand is detected 51.7% of the time, which to some extent still reflects the location of the hand of interest as the two hands usually appear at nearby locations. In poor lighting, however, the algorithm fails to reliably track the hand. Specifically, the unlit overhead light in the room was picked up as a hand in all 50 frames tested. Another finding worth mentioning was that tracking performance was insensitive to the on-camera location of the hand. Hands appearing on the edges of a frame are tracked with approximately the same accuracy as those near the center. Based on our testing, we believe the mobile-optimized hand-tracking algorithm [23] selected for this experiment to

450 be sufficiently reliable to locate hands in real-time when the hands are properly
 451 lit and exposed.

453 5.2 3D Audio Experiment



470 **Fig. 8.** An illustration of the setup for the 3D Audio Experiment. “Speaker” and “ear”
 471 icons credit to Google.

473 In this experiment, we wish to find how accurately a human can perceive the
 474 location of a computer-generated audio signal through stereo headphones. We
 475 set up this experiment so that a computer program first randomly selects one of
 476 three horizontal directions: front, left, and right. When a direction is selected,
 477 an exact angle is further generated by drawing from a uniform distribution over
 478 $\pm 15^\circ$ from the mean angle of that direction (15° for “right”, 90° for “front”,
 479 and 165° for “left”). A graphical depiction of this process is shown in fig. 8.
 480 Then, the program creates a virtual 3D audio source with an azimuth of the
 481 angle previously chosen, an elevation of 90° (parallel to the ground), and a
 482 radial distance of $1m$. The virtual audio source plays a 440Hz sine wave for 1.0
 483 second, physically transmitted through a pair of headphones worn over the ears
 484 of the test subject. After the signal plays, the test subject enters the perceived
 485 location of the audio source on the computer program, after which the next
 486 virtual audio source is created and the process repeats. We gather 50 data points,
 487 each consisting of the source direction generated by the computer and the source
 488 direction perceived by the human. The confusion matrix is visualized below (see
 489 fig. 9). We include a video segment in the supplementary materials demonstrating
 490 the test program and the computer-generated 3D audio. If desired, the reader
 491 can experience a demo of the 3D audio experiment by playing the video segment
 492 with stereo headphones on.

493 As the data in fig. 9 shows, a human is fully capable of accurately discerning
 494 the direction of a computer-generated 3D audio signal when they are separated

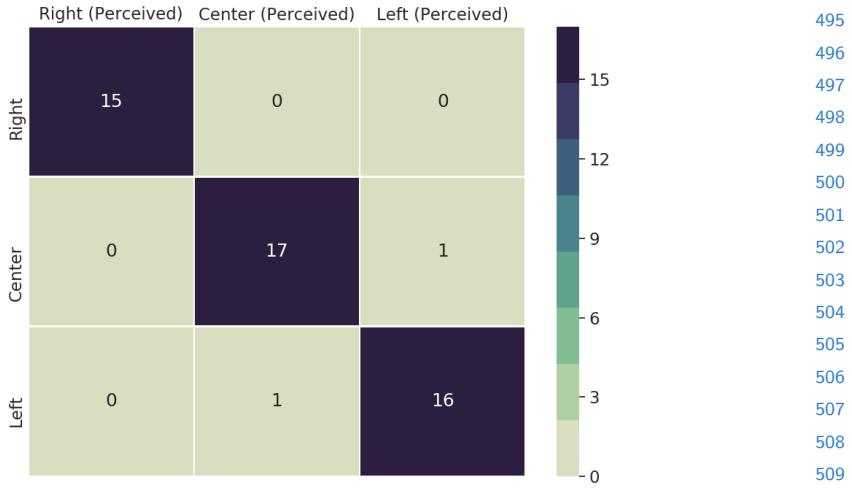


Fig. 9. Confusion matrix for the 3D Audio experiment. Rows are indexed by computer-generated directions while columns are indexed by human-perceived directions.

by different azimuth angles. In only 2 out of 50 trials have the human subject incorrectly determined the direction of the generated signal, which translates to an accuracy of 96%. This provides the insight that the 3D Audio module is a reliable way to encode spatial information into a sound to allow a listener to determine its virtual location. Spatial sounds are also highly intuitive, as almost all people with regular hearing use their two ears to locate sound sources in real life for various purposes and 3D Audio is designed to simulate this experience.

5.3 Audi-Exchange Experiment

The Audi-Exchange Experiment integrates both Hand Tracking and 3D Audio to serve as the Audi-Exchange assistive device prototype (see fig. 10). This experiment uses the “Regular Lighting / Empty Hand” images that were previously used in sec. 5.1, but we further split them into three classes, front, left, and right, based on the location of the hand appearing in the image. In each trial, an image is randomly drawn from the three directional classes and fed into the hand detector. After that, a virtual audio source is created at ground level and 1m in front of the user with a horizontal location being proportional to that of the detected hand such that the virtual horizontal location is within $[-1m, 1m]$. Then, the virtual location is normalized to have a distance of 1m so that all audio signals will have equal volume regardless of the location. In other words, if $x, y \in [0, 1]$ is the horizontal and vertical location of the hand on an image,

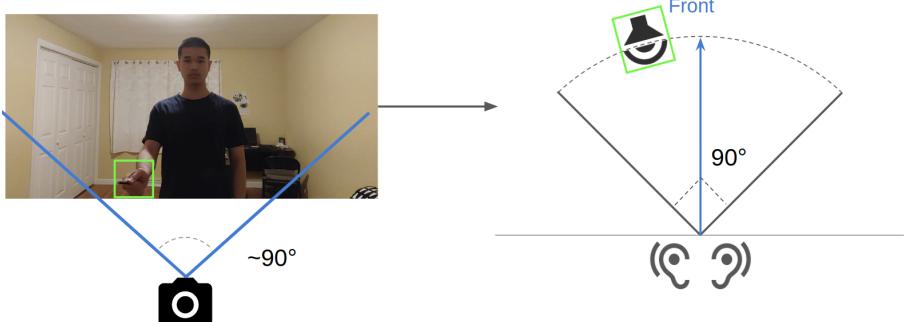


Fig. 10. Illustration of the Audi-Exchange Experiment. The hand-tracking and 3D audio algorithms translate visual interaction cues of an extended hand into a 3D audio signal for the BVI user. Camera image shown was cropped to better fit the page. “Camera”, “speaker”, and “ear” icons credit to Google

then the location of the virtual sound source (x', y', z') is computed as:

$$\begin{aligned} (\hat{x}, \hat{y}, \hat{z}) &= ((x - 0.5) * 2, 1, 0) \\ (x', y', z') &= \frac{(\hat{x}, \hat{y}, \hat{z})}{\|(\hat{x}, \hat{y}, \hat{z})\|} \end{aligned} \quad (2)$$

where positive x' is the right, positive y' is the front, and positive z' is the top. Note that the vertical location of the hand is disregarded because we think horizontal location is the most important piece of information to facilitate hand-based interactions. We place the virtual audio source at (x', y', z') as computed in eq. 2, which plays a 440Hz sine wave sound for 1.0 second through headphones. Then, the test subject enters the perceived direction of the hand based solely on sound and proceeds to the next trial when ready. If the hand tracker fails to detect the hand, the program automatically records a “not detected” response. For each trial, we record the ground-truth direction of the hand and the perceived direction of the hand. The data is showed in the confusion matrix in fig. 11. The Audi-Exchange Experiment is more challenging than the independent 3D Audio and Hand Tracking experiments for two reasons. (1) Hand tracking is not perfectly accurate, so the computed image location of the tracked hand can be off. (2) When mapping the image location into the 3D location for the virtual audio source, the virtual location spans only 90° horizontally as opposed to 180° (the case for the 3D Audio Experiment) in an attempt to reproduce the actual azimuth direction of the hand in the real world, so there are smaller perceptive differences among sounds coming from different virutal directions. Despite these challenges, the test subject has correctly located the hand n% of the time based on sound alone. The high tested effectiveness of the Audi-Exchange prototype strongly indicates that Audi-Exchange would be feasible and effective as an assistive platform when used by BVI individuals.

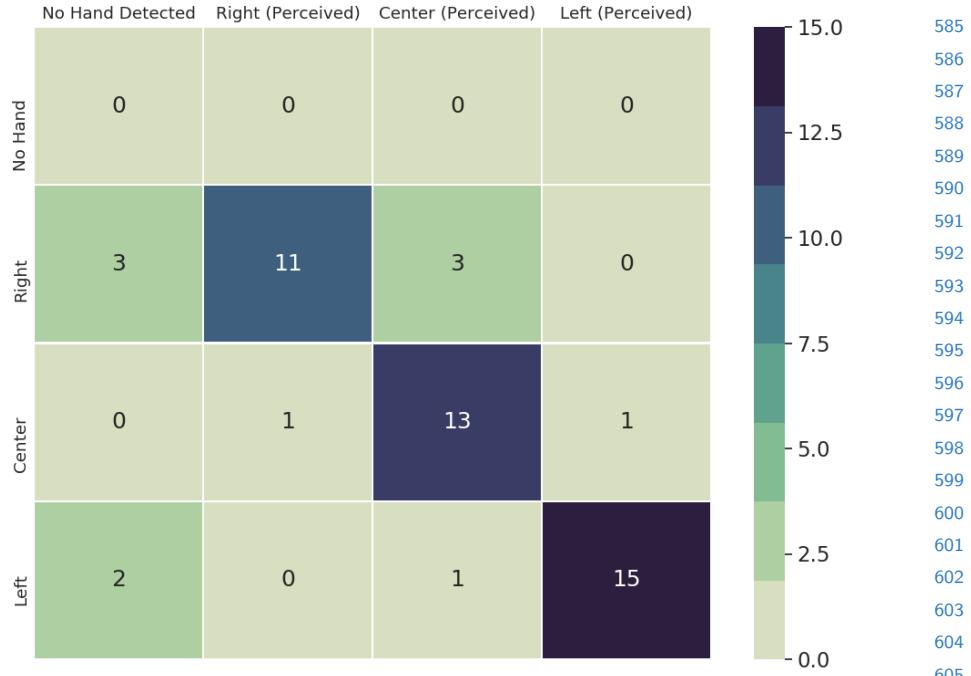


Fig. 11. Confusion matrix for the Audi-Exchange Experiment. Rows are indexed by ground-truth directions while columns are indexed by human-perceived directions.

6 Porting to Mobile Devices

Although we tested the Audi-Exchange prototype on a desktop computer, it can be updated and ported to mobile computing platforms (e.g. smartphones) with relative ease. The current hand tracking algorithm is already based on a neural network architecture designed for mobile usage [15] can be run in efficiently on smartphones. While the OpenAL [3] library used in this paper is mostly used for desktop applications, newer and more powerful spatial open-source sound APIs exist that target mobile applications specifically. Resonance Audio is an example [18]. Hence, we expect porting the Audi-Exchange prototype to smartphones to be a simple process, from where we can continue to develop the Audi-Exchange assistive technology.

7 Conclusion

We have developed a working version of Audi-Exchange, a mobile assistive technology that guides the BVI when handing over and receiving objects from another person with 3D sound cues. We have implemented Audi-Exchange with a mobile-optimized neural network (i.e. MobileNet + SSD [15,17]) and OpenAL [3]. We tested the system with images of a person using their hand to offer and

receive objects captured in front of a smartphone camera and demonstrated that under proper lighting, the hand can be reliably tracked and the user can accurately determine its location by listening to a corresponding computer-generated 3D audio signal through stereo headphones. This gives us confidence that blind- or-visually-impaired individuals can use Audi-Exchange to better complete object exchanges with their hand. As the next step, we plan to construct more in-depth experiments as well as update Audi-Exchange to further improve and validate Audi-Exchange.

630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674

638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674

675 References

- 677 1. Google glass, <https://www.google.com/glass/start/> 677
- 678 2. Microsoft hololens, <https://www.microsoft.com/en-us/hololens/> 678
- 679 3. Openal: Cross platform 3d audio, <https://openal.org/> 679
- 680 4. Seeing ai, <https://www.microsoft.com/en-us/ai/seeing-ai> 680
- 681 5. Acesight - low vision eletronic glasses (July 2020), <https://www.acesight.com/> 681
- 682 6. esight 3 — low vision device for the visually impaired (July 2020), <https://esighteyewear.com/low-vision-device-for-visually-impaired/> 682
- 683 7. Irisvision — wearable low vision glasses for visually impaired (May 2020), <https://irisvision.com/> 683
- 684 8. Orcam myeye 2.0 (July 2020), <https://www.orcam.com/en/myeye2/> 684
- 685 9. Pioneering accessible navigation – indoors and outdoors (May 2020), <https://www.blindsightsquare.com/> 685
- 686 10. Attia, I., Asamoah Brempong, D.: The white cane. its effectiveness, challenges 686 and suggestions for effective use: The case of akropong school for the blind. 689 Journal of Education, Society and Behavioural Science pp. 47–55 (05 2020). 690 <https://doi.org/10.9734/JESBS/2020/v33i330211> 691
- 692 11. Bharambe, S., Thakker, R., Patil, H., Bhurchandi, K.: Substitute 692 eyes for blind with navigator using android. pp. 38–43 (04 2013). 693 <https://doi.org/10.1109/TIIEC.2013.14> 694
- 695 12. Elmannai, W., Elleithy, K.: Sensor-based assistive devices for visually-impaired 695 people: Current status, challenges, and future directions. Sensors **17**(3), 565 (2017) 696
- 697 13. Francés, L., Morer, P., Rodriguez, M., Cazón-Martín, A.: Wearable design 697 requirements identification and evaluation. Sensors **20**, 2599 (05 2020). 698 <https://doi.org/10.3390/s20092599> 699
- 699 14. Hoover, R.: The cane as a travel aid. Blindness pp. 353–365 (1950) 700
- 700 15. Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for 701 mobile vision applications (2017) 702
- 703 16. Kulyukin, V., Gharpure, C., Nicholson, J., Pavithran, S.: Rfid in robot-assisted 703 indoor navigation for the visually impaired. In: 2004 IEEE/RSJ International Conference 704 on Intelligent Robots and Systems (IROS)(IEEE Cat. No. 04CH37566). 705 vol. 2, pp. 1979–1984. IEEE (2004) 706
- 707 17. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, 707 A.C.: Ssd: Single shot multibox detector. Lecture Notes in Computer Science p. 708 21–37 (2016). https://doi.org/10.1007/978-3-319-46448-0_2, http://dx.doi.org/10.1007/978-3-319-46448-0_2 709
- 710 18. Martin Dufour, Aaron McLeran, M.B., et al.: Resonance-audio (2018), <https://github.com/resonance-audio/resonance-audio> 711
- 712 19. Oliveira, J.F.: The path force feedback belt. In: 2013 8th International Conference 712 on Information Technology in Asia (CITA). pp. 1–6. IEEE (2013) 713
- 714 20. Prudhvi, B., Bagani, R.: Silicon eyes: Gps-gsm based navigation assistant for 714 visually impaired using capacitive touch braille keypad and smart sms facility. In: 2013 715 World Congress on Computer and Information Technology (WCCIT). pp. 1–3. 716 IEEE (2013) 717
- 718 21. Roentgen, U.R., Gelderblom, G.J., Soede, M., De Witte, L.P.: Inventory of 718 electronic mobility aids for persons with visual impairments: A literature review. Journal 719 of Visual Impairment & Blindness **102**(11), 702–724 (2008)

- 720 22. Shilkrot, R., Huber, J., Liu, C., Maes, P., Nanayakkara, S.C.: Finger-
721 reader: A wearable device to support text reading on the go. In: CHI
722 '14 Extended Abstracts on Human Factors in Computing Systems. p.
723 2359–2364. CHI EA '14, Association for Computing Machinery, New York,
724 NY, USA (2014). <https://doi.org/10.1145/2559206.2581220>, <https://doi.org/10.1145/2559206.2581220>
- 725 23. Victor, D.: Handtrack: A library for prototyping real-time hand tracking interfaces
726 using convolutional neural networks. GitHub repository (2017), <https://github.com/victordibia/handtracking/tree/master/docs/handtrack.pdf>
- 727 24. Wahab, M.H.A., Talib, A.A., Kadir, H.A., Johari, A., Noraziah, A., Sidek, R.M.,
728 Mutalib, A.A.: Smart cane: Assistive cane for visually-impaired people. arXiv
729 preprint arXiv:1110.5156 (2011)
- 730 25. Wikipedia contributors: Head-related transfer function — Wikipedia, the free
731 encyclopedia. [https://en.wikipedia.org/w/index.php?title=Head-related](https://en.wikipedia.org/w/index.php?title=Head-related_transfer_function&oldid=940911588)
732 _transfer_function&oldid=940911588 (2020), [Online; accessed 18-July-2020]
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- 756
- 757
- 758
- 759
- 760
- 761
- 762
- 763
- 764