



 **Opin vísindi**

This is not the published version of the article / Þetta er ekki útgefna útgáfa greinarinnar

Author(s)/Höf.: Spagnol, S., Hoffmann, R., Herrera Martínez, M., & Unnthorsson, R.

Title/Titill: Blind wayfinding with physically-based liquid sounds.

Year/Útgáfuár: 2018

Version/Útgáfa: Pre-print / Óritrýnt handrit

Please cite the original version:

Vinsamlega vísið til útgefnu greinarinnar:

Spagnol, S., Hoffmann, R., Herrera Martínez, M., & Unnthorsson, R. (2018). Blind wayfinding with physically-based liquid sounds. *International Journal of Human-Computer Studies*, 115, 9-19. doi:<https://doi.org/10.1016/j.ijhcs.2018.02.002>

Rights/Réttur: © 2018 Elsevier Ltd. All rights reserved.

Blind wayfinding with physically-based liquid sounds

Simone Spagnol^{a,*}, Rebekka Hoffmann^{a,b}, Marcelo Herrera Martínez^a,
Runar Unnthorsson^a

^a*Faculty of Industrial Engineering, Mechanical Engineering and Computer Science,
University of Iceland, Dunhagi 5, 107 Reykjavík, Iceland*

^b*Faculty of Psychology, University of Iceland, Sæmundargata 10, 101 Reykjavík, Iceland*

Abstract

Translating visual representations of real environments into auditory feedback is one of the key challenges in the design of an electronic travel aid for visually impaired persons. Although the solutions currently available in the literature can lead to effective sensory substitution, high commitment to an extensive training program involving repetitive sonic patterns is typically required, undermining their use in everyday life. The current study explores a novel sensory substitution algorithm that extracts information from raw depth maps and continuously converts it into parameters of a naturally sounding, physically based liquid sound model describing a population of bubbles. The proposed approach is tested in a simplified wayfinding experiment with 14 blindfolded sighted participants and compared against the most popular sensory substitution algorithm available in the literature - the vOICe (Meijer, 1992) - following a short-time training program. The results indicate a superior performance of the proposed sensory substitution algo-

*Corresponding author

Email addresses: `spagnols@hi.is` (Simone Spagnol), `rebekkah@hi.is` (Rebekka Hoffmann), `marcelo@hi.is` (Marcelo Herrera Martínez), `runson@hi.is` (Runar Unnthorsson)

rithm in terms of navigation accuracy, intuitiveness and pleasantness of the delivered sounds compared to the vOICe algorithm, supporting its usability for the visually impaired community.

Keywords: sensory substitution, sonification, electronic travel aid, physical sound model

1. Introduction

The technique of data sonification is used as an alternative or a complement to data visualization for representing various actions, objects or signals. Sonification can be defined as “*a mapping of numerically represented relations in some domain under study to relations in an acoustic domain for the purposes of interpreting, understanding, or communicating relations in the domain under study*” [40]. Widely accepted sonification techniques include *audification* (i.e., direct playback of data streams as sound waves), *auditory icons* (i.e., discrete environmental sounds), *earcons* (i.e., discrete symbolic sounds), *parameter mapping sonification* between data dimensions and auditory dimensions, and *model-based sonification* (i.e., based on dynamic models of virtual sounding objects) [23, 17].

Sonification is used in very different contexts to represent a great variety of data, ranging from molecular information [19] to geophysical data [16]. Of particular interest are applications in health care, such as in motor rehabilitation systems [1, 39] where task-related auditory information is able to support motor learning and increases attention and engagement levels during rehabilitation tasks. Another widely explored area is that of electronic travel aids [15] and other assistive technologies for visually impaired persons

20 (VIPs) [14], where sonification techniques are designed to substitute visual
21 information [28]. Unfortunately, the majority of the systems exploiting such
22 techniques are still in their infancy and have limited functionalities, small
23 scientific and/or technological value and high cost [15].

24 Available electronic travel aids for VIPs range from simple *obstacle de-*
25 *tectors* with a single range-finding sensor (e.g. ultrasound, infrared), to *envi-*
26 *ronmental imagers* employing data generated from visual representations ac-
27 quired through camera technologies. The most common sonification schemes
28 of obstacle detectors, which only receive range information, are either earcons
29 indicating the presence of an obstacle, or an inversely proportional transform
30 mapping one or more range readings to the loudness and/or pitch of synthetic
31 sounds or musical tones [10]. On the other hand, environmental imagers (i.e.,
32 devices able to deliver a representation of the layout of an environment) allow
33 for greater flexibility in sonification mappings. The most significant exam-
34 ple is provided by the well-known image sonification algorithm used in the
35 vOICe system [30].

36 The vOICe algorithm can be thought of as an inverse spectrogram trans-
37 form, i.e., a time-varying sound whose spectrogram approximately matches
38 an input grayscale image. In particular, the algorithm periodically scans the
39 image from left to right, while associating each row to a different sinusoidal
40 oscillator with fixed frequency (in ascending order from lower to upper rows)
41 and using the brightness of each pixel in turn to control the amplitude of
42 the oscillator. The sound output is then spatialized left to right according to
43 the current scanning point. It has been shown that, following extensive pe-
44 riods of training and exploiting the neural plasticity of the human brain, the

45 vOICe sonification mechanism can lead to effective sensory substitution [31],
46 both in object recognition [53] and spatial learning [35].

47 Although the original vOICe algorithm was designed to sonify 2D grayscale
48 images, its use in blind wayfinding is supported by the observation that a
49 depth map can be directly converted into a grayscale image where brightness
50 corresponds to depth. The use of depth information for the sonification of 3D
51 scenes through either the original vOICe algorithm or slight variations of it
52 has already been proposed and investigated [11, 52]. Furthermore, improve-
53 ments to the pleasantness of sounds (such as using musical tones instead of
54 pure sines) as well as to the spatial feeling and real-time conveyance of the
55 sounds (e.g. presenting independently to each headphone channel simulta-
56 neous scans from the left and right edge to the central column of the image)
57 were proposed [3].

58 The main drawback of most existing sensory substitution devices (SSDs),
59 including the vOICe, is that even though in some cases the conveyed audi-
60 tory information can be successfully interpreted by naïve users, they demand
61 extremely high commitment on the user’s side. A lengthy and strenuous
62 training of up to one year is required in order to enable users to perform
63 most tasks, thus undermining the use of SSDs in everyday life [35]. As
64 Fontana *et al.* point out [18], the prolonged use of SSDs “*leads to the strain*
65 *of the user [...] due to the continuous listening of the same signal at regular*
66 *time intervals. This sound, even if spatialized, produces an unnatural effect*
67 *and causes a progressive fatigue.*” Therefore, the choice of the type of sound
68 as well as the way it is generated should be regarded as a key issue in the
69 design of any sensory substitution algorithm.

70 The current study explores a novel model-based sonification algorithm for
71 translating continuous representations of a dynamic real environment, coded
72 into sequences of depth maps, into auditory feedback. The sensory substitu-
73 tion algorithm we propose is meant to be used for real-time blind wayfinding,
74 with minimum latency between data acquisition and sonification, and with
75 available off-the-shelf hardware technologies. It was designed in an attempt
76 to improve the vOICe algorithm from both an ergonomic and a functional
77 point of view, eventually reducing the required training time, and to be ef-
78 ficiently scalable depending on the available computational resources. The
79 algorithm we propose here directly maps low-order statistics from the raw
80 depth map into the parameters of a physically-based liquid sound model. In
81 this model, physical descriptions of sound events are intentionally simplified
82 to emphasize the most perceptually-relevant timbral features, and to reduce
83 computational requirements as well [4]. The model was specially selected
84 and tuned in order to sound both natural (yet significantly discernible from
85 most daily environmental sounds) and aesthetically pleasant.

86 The remainder of the paper is organized as follows. In Section 2 we
87 describe the generation mechanism of liquid sounds and its use in the design
88 of our *fluid flow* sensory substitution algorithm. In Section 3 we introduce
89 an experiment designed in order to assess the performance and individual
90 preference of the sensory substitution algorithm in a blind wayfinding task.
91 Results are reported in Section 4 and finally discussed in Section 5.

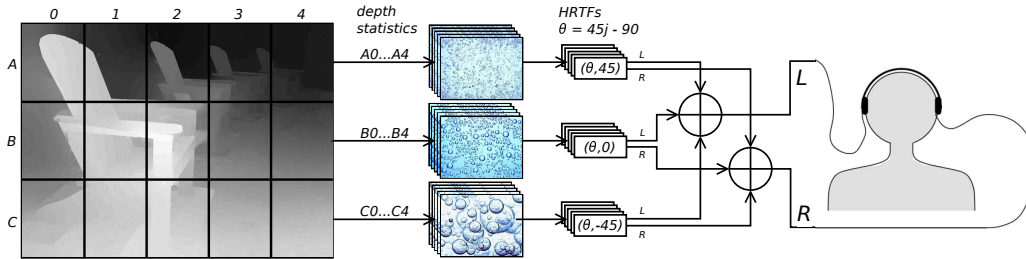


Figure 1: Simplified scheme of the proposed sensory substitution algorithm.

92 2. Sensory substitution with liquid sounds

93 The *fluid flow* sensory substitution algorithm that we propose in this
 94 paper receives a sequence of depth maps as input. Each depth map is di-
 95 vided into 15 equally sized sectors given by the combination of 3 rows and 5
 96 columns. Every sector corresponds to an independent and uncorrelated in-
 97 stance of a liquid sound generator, and its position within the depth map is
 98 spatialized in the frontal hemisphere, allowing for effective source separation.
 99 Figure 1 reports a simplified scheme of the proposed algorithm.

100 2.1. Generation of liquid sounds

101 The building block of the fluid flow algorithm is the *liquid sound gener-*
 102 *ator*. In the physical world, liquid sounds are mostly caused by gas bubbles
 103 trapped inside the liquid rather than by the liquid mass itself. For this reason,
 104 sound is generated through a stochastic process modeling the temporal evo-
 105 lution of a population of bubbles, a synthesis approach previously referred to
 106 as *physically informed sonic modeling by granular synthesis* [57]. The liquid
 107 sound generation algorithm considers individual bubbles to be atomic units
 108 (or *grains*, according to the granular synthesis terminology [37]), synthesized

109 using the well-known physically based Minnaert model [32]. Spherical bub-
 110 bles effectively act as exponentially decaying sinusoidal oscillators: the com-
 111 pressible gas region of the bubble, surrounded by an incompressible liquid
 112 mass, gradually dissipates the energy involved in its creation by a periodic
 113 pulsation, as it would happen in a spring-mass system.

114 Every single bubble k , whose impulse response is

$$i_k(t) = a_k \sin(2\pi f_k^0 t) e^{\zeta_k t} \quad (1)$$

115 is fully defined by means of its radius r_k and depth factor D_k , that uniquely
 116 determine the individual damping factor ζ_k , resonant frequency f_k^0 , and am-
 117 plitude a_k as follows:

$$\zeta_k = \frac{0.13}{r_k} + 0.0072 r_k^{-\frac{3}{2}} \quad f_k^0 = \frac{3}{r_k} \quad a_k = D_k r_k^{\frac{3}{2}} \quad (2)$$

118 Here the depth factor D_k models the lumped effect of the depth of a bubble,
 119 and the effect of different excitation strengths of the bubbles. Bubbles that
 120 are submerged more will be attenuated more. Factor D_k is a dimensionless
 121 number between 0 and 1, where 1 corresponds to a bubble created at the
 122 surface and 0 to a fully submerged bubble.

123 The creation of bubbles is then modeled as a Bernoulli process occurring
 124 at audio rate with success probability $p = 1/\Lambda$, where Λ is the average bubble
 125 rate (bubbles per second). The radius of each successfully produced bubble
 126 k is set to

$$r_k = x_k^{\gamma_r} (r_{MAX} - r_{MIN}) + r_{MIN} \quad (3)$$

127 where $x_k \in [0, 1]$ is a number drawn from a uniform distribution function,
 128 r_{MIN} and r_{MAX} are the minimum and maximum bubble radius values, and

129 γ_r is the radius gamma factor, which allows to increase the ratio of bigger
 130 bubbles relative to smaller bubbles ($0 < \gamma_r < 1$) or *vice versa* ($\gamma_r > 1$).
 131 Similarly, the depth factor D_k is set to

$$D_k = y_k^{\gamma_D} (D_{MAX} - D_{MIN}) + D_{MIN} \quad (4)$$

132 where $y_k \in [0, 1]$ is a number drawn from a uniform distribution function,
 133 D_{MIN} and D_{MAX} are the minimum and maximum depth factor values, and
 134 γ_D is the depth gamma factor, which allows to increase the ratio of bubbles
 135 close to the surface relative to deeper bubbles ($0 < \gamma_D < 1$) or *vice versa*
 136 ($\gamma_D > 1$).

137 Bubble sounds often exhibit a characteristic rise in pitch, especially when
 138 approaching the surface. The phenomenon is mostly caused by the pressure
 139 reduction as the liquid mass above the bubble becomes thinner and thinner.
 140 The effect is modeled in the synthesis algorithm by a global rise factor pa-
 141 rameter ξ . Since bubbles with a rising pitch are created close to the surface,
 142 it seems reasonable to assume they are generally louder than average. This
 143 effect is modeled by a rise cutoff parameter K_ξ . When it is set to a value
 144 $0 < K_\xi < 1$, only bubbles with a depth factor $D_k > K_\xi$ have a nonzero rise
 145 factor ξ . According to the physically based bubble sound model described
 146 in [57], a rising bubble is modeled by making its frequency time-dependent
 147 according to

$$f_k(t) = f_k^0 (1 + \sigma_k t) \quad (5)$$

148 where σ_k is the slope of the frequency rise related to the vertical velocity of
 149 the bubble, modeled as

$$\sigma_k = \xi \zeta_k. \quad (6)$$

150 An implementation of the liquid sound generator described above (*fluid*
151 *flow* module) is included in the Sound Design Toolkit (SDT),¹ an open-source
152 (GPLv2) library of physically based sound synthesis algorithms for Max and
153 Pure Data [4]. In this implementation the stochastic process drives an oscil-
154 lator bank, whose number of voices can be set as a parameter. The size of
155 the oscillator bank defines the polyphony of the algorithm, i.e. the maximum
156 number of bubbles that can be active at the same time. If the maximum num-
157 ber is exceeded, a voice stealing mechanism takes place and the new bubble
158 is assigned to the oscillator that currently has the minimum instantaneous
159 amplitude envelope, resetting all its parameters, base frequency included.
160 Phase alignment allows to avoid audible artifacts during the generation of a
161 new bubble [46].

162 The liquid sound generator is a slightly improved version of the *bubble*
163 *simulator* proposed by van den Doel [57]. The main improvement with re-
164 spect to the van den Doel simulator lies in the use of a single Bernoulli
165 process for a population of bubbles with different radii (i.e., with different
166 base frequencies) rather than 50 Bernoulli processes each set to a fixed base
167 frequency. This strategy allows to represent bubbles of arbitrary size, im-
168 proving the versatility of the algorithm especially with small oscillator banks.

169 2.2. Model-based sonification

170 A global d_{MAX} parameter is defined in order to consider only those points
171 in the depth map whose depth is no greater than this defined parameter.
172 Then, for each sector, two descriptive depth metrics are calculated: *map*

¹<http://soundobject.org/SDT/>

173 *density* and *average depth*. Design choices for mappings between depth map
174 properties and liquid sound features are the following:

- 175 • map density \rightarrow average bubble rate;
- 176 • average depth \rightarrow maximum bubble depth factor.

177 Map density ρ is defined as the number of pixels with depth value no greater
178 than d_{MAX} divided by the total number of pixels in that sector. It is mapped
179 to the *average bubble rate* Λ according to

$$\Lambda = 500\rho^2 \tag{7}$$

180 so that the denser the sector, the more the generated bubbles. The upper
181 limit of 500 bubbles/second was heuristically set following informal investi-
182 gations on the pleasantness and intelligibility of the associated liquid sound.

183 Average depth \bar{d} is defined as the mean depth value (in meters) of all
184 pixels with depth no greater than d_{MAX} in that sector. It is mapped to the
185 *maximum bubble depth factor* D_{MAX} as

$$D_{MAX} = \left(\frac{d_{MAX} - \bar{d}}{d_{MAX}} \right)^2. \tag{8}$$

186 In this way, closer obstacles are transformed in a larger amount of bubbles
187 close to the surface of the water, thus increasing their average loudness and
188 sharpness. As an analogy, it might help to think of the scene as a big aquar-
189 ium seen from above, with the water surface just in front of the observer and
190 all objects producing bubbles.

191 In order to provide a spatial dimension of the depth map, the sound
192 produced by each liquid sound generator is binaurally spatialized by mapping

193 the corresponding depth map sector (R_i, C_j) to the azimuth and elevation
 194 parameters (θ, ϕ) of a generic HRTF filter as follows:

$$\theta = 45j - 90 \quad (9)$$

195

$$\phi = 45 - 45i \quad (10)$$

196 where θ and ϕ are expressed in degrees with respect to the observer according
 197 to a vertical polar coordinate system, $i = 0, 1, 2$ is the row number (top to
 198 bottom), and $j = 0, \dots, 4$ is the column number (left to right). However,
 199 since elevation cues greatly differ from subject to subject [48] and lead to
 200 high variance in vertical localization performance with generic HRTFs [33],
 201 elevation information is redundantly coded into another liquid sound fea-
 202 ture. In particular, sectors belonging to different rows of the depth map are
 203 assigned different bubble radius intervals [r_{MIN}, r_{MAX}] as follows:

$$\begin{aligned} R_0 : r_{MIN} &= 0.2mm, r_{MAX} = 1mm; \\ R_1 : r_{MIN} &= 1mm, r_{MAX} = 5mm; \\ R_2 : r_{MIN} &= 5mm, r_{MAX} = 20mm. \end{aligned} \quad (11)$$

204 Thanks to the inversely proportional relation between bubble radius and res-
 205 onant frequency (see Eq. 2), the above heuristically defined intervals allow
 206 for different characteristic liquid sounds to be produced depending on eleva-
 207 tion, i.e., ranging from light, fizzy sounds for higher elevations (row R_0) to
 208 low, gurgling sounds for lower elevations (row R_2).

209 Other parameters that define the liquid sound generator are kept con-
 210 stant. These include the radius gamma factor ($\gamma_r = 1$), the minimum bub-
 211 ble depth ($D_{MIN} = 0$), the depth gamma factor ($\gamma_D = 1$), the rise factor

212 ($\xi = 0.5$), and the rise cutoff ($K_\xi = 0.5$). Both gamma factors are set to 1
213 in order to preserve the uniform distribution of radius and depth values. On
214 the other hand, the choices for the rise factor and rise cutoff allow for an ad-
215 ditional auditory depth cue. By combining Eq. 8 and Eq. 4 it can be shown
216 indeed that the average depth value at which pitch-rising bubbles start being
217 produced ($D_k > K_\xi$) roughly corresponds to $\bar{d} \approx 0.3d_{MAX}$. This translates
218 at auditory level into a peculiar boiling water sound for close objects, and
219 the closer the object (i.e., the lower the average depth value), the higher the
220 number of pitch-rising bubbles and therefore the clearer the boiling effect.

221 A preliminary version of the fluid flow algorithm was previously pre-
222 sented by the authors in [46]. With respect to the previous version, the
223 main improvements of the algorithm described here lie in the representation
224 of elevation information with different bubble radius values, in using bub-
225 ble depth as a proper physical depth indicator rather than plain amplitude
226 control, and in the use of the rising pitch cue for close objects rather than
227 elevated objects. These design changes were suggested from both test results
228 and informal comments following preliminary experimental trials with offline
229 video sequences [46], that highlighted above all the difficulty of interpreting
230 elevation cues.

231 At the same time, the new mappings provide more meaningful correspon-
232 dences between physical and auditory cues. As a matter of fact, beside the
233 intuitive relationship between physical depth and bubble depth, crossmodal
234 correspondences between pitch (resonant frequency in the bubble model) and
235 elevation are well known in the literature [25] and frequently used in sensory
236 substitution systems (including the vOICE). Furthermore, the boiling ef-

237 fect that gets more and more prominent while approaching an object can be
238 interpreted as an effective natural warning sound [55].

239 **3. Evaluation**

240 The main goal of the experiment presented here is to assess the per-
241 formance and individual preference of the *fluid flow* sensory substitution
242 algorithm in a blind wayfinding task. More in detail, the point-by-point
243 objectives are

- 244 1. to validate the effectiveness of the proposed sounds of giving reliable
245 and distinguishable information in a simplified wayfinding task with a
246 reasonably sized pool of naïve blindfolded participants;
- 247 2. to collect individual judgments about the naturalness, pleasantness and
248 usability of the sounds that are conveyed;
- 249 3. to compare the above results and ratings against those collected using
250 the reference sensory substitution scheme provided through the original
251 vOICe algorithm [30].

252 Our working hypotheses are that: (1) after a short training session, the
253 fluid flow algorithm is able to help participants avoid obstacles in the large
254 majority of the presented cases; (2) performance and completion time are
255 at least comparable to the vOICe algorithm; (3) the individual judgments
256 on the liquid sounds reflect a positive opinion on all the investigated aspects
257 and, in particular, a more positive rating compared to the sounds produced
258 by the vOICe algorithm.

259 *3.1. Sample*

260 Fourteen participants (7F, 7M) participated on a voluntary basis. Ages
261 ranged from 22 to 46 ($M = 30.5$, $SD = 7.2$). All participants spoke fluent
262 English and none of them reported either visual or hearing impairments. All
263 participants gave their informed consent for inclusion before they partici-
264 pated in the study. The study was conducted in accordance with the Decla-
265 ration of Helsinki, and the protocol was approved by the National Bioethical
266 Committee of Iceland (reference number VSN-15-107).

267 *3.2. Experimental setup*

268 The experiment took place in an empty classroom sized $8m$ (length) \times
269 $6.7m$ (width) $\times 3.5m$ (height) inside a building of the University of Iceland.
270 Four pieces of green carpet, sized $4m \times 0.5m$ each, were placed in the middle
271 of the classroom floor in order to delimit a square $3.5m \times 3.5m$ testing area
272 (see Figure 2a). During the whole experiment, to control for confounding
273 effects, windows were kept closed and artificial light was turned on. The
274 absence of any kind of activity in the neighboring classrooms due to sum-
275 mer break guaranteed a quiet environment throughout the testing sessions.
276 The ventilation system of the classroom produced the only significant, yet
277 constant, environmental sound.

278 During the tests, white cardboard boxes were placed in predefined loca-
279 tions of the testing area. The size of a single cardboard box was $0.4m$ (length)
280 $\times 0.4m$ (width) $\times 0.6m$ (height). The number of boxes inside the testing area
281 during each experimental trial ranged from 5 to 8; when less than 8, the un-
282 used boxes were placed along one wall as shown in Figure 2a. Furthermore, a
283 tripod holding a small Bluetooth box speaker (at approximately $1.2m$ height)



Figure 2: Experimental setup. (a) Subject during the experiment. (b) Close up of the equipment.

284 was placed along the end-side of the testing area. The only other significant
285 objects present in the room were a desk and two chairs for the experimenters,
286 all positioned behind the starting point of the participants.

287 Participants wore the following equipment, pictured in Figure 2b: (a) an
288 elastic headband (originally holding a searchlight) with a Structure Sensor
289 camera², a high-performance structured light 3D sensor, tightened to the
290 frontal plastic hold; (b) a pair of open over-ear headphones (AKG K612 Pro)
291 allowing environmental sound to enter the ear; (c) a small backpack carrying
292 a Lenovo Ideapad Y700 laptop running the software to which the camera,
293 headphones and (d) an external battery were connected; (e) a blindfold. In
294 order to ensure regular functioning, the laptop was constantly monitored
295 by an experimenter through a second laptop placed on the desk behind the

²<https://structure.io/>

296 testing area, connected via VPN.

297 Depth maps with a resolution of 640×480 pixels were acquired from
298 the Structure Sensor at a rate of 10 frames per second with the support
299 of an open-source Matlab Wrapper for OpenNI 2.2,³ processed in Matlab,
300 and sonified through the Pure Data software implementing the fluid flow
301 and vOICe algorithms. Depth maps spanned the entire field of view of the
302 Structure Sensor, i.e., 58° horizontal, 45° vertical, and a $0.4m$ to $3m$ depth
303 range. Visual information falling beyond these ranges was therefore not
304 sonified.

305 3.3. Stimuli

306 The sound stimulus conveyed to participants during the experiment was a
307 continuous sonification of the depth data acquired through the Structure Sen-
308 sor, either through the fluid flow algorithm, referred to as FF and described
309 in Section 2, or the vOICe algorithm, referred to as VC and described in
310 the following paragraph. Each algorithm was implemented as a Pure Data
311 patch that constantly receives the depth map statistics data through the
312 OSC (Open Sound Control) protocol. In order to avoid audible artifacts, the
313 incoming depth map statistics values were smoothed with a 100-ms ramp
314 function. In the experiment, the d_{MAX} parameter was set to $3m$ and the
315 number of voices of each liquid sound generator to 32. For the sake of
316 consistency, the level of the sound card was kept constant throughout the
317 experiment for all participants.

³<http://uk.mathworks.com/matlabcentral/fileexchange/42127-matlab-wrapper-for-openni-2-2>

318 The vOICe sensory substitution algorithm was implemented following the
319 specifications from Meijer [30]. The algorithm scans each depth snapshot
320 (resized to 64×64 pixels) from left to right, while associating height (i.e.
321 the vertical coordinate of the pixel) with pitch and depth with loudness.
322 More specifically, every row is associated to an amplitude-controlled oscillator
323 whose fixed frequency exponentially ranges from 500 Hz (bottom row) to 5
324 kHz (top row), while amplitude is inversely proportionally related to the
325 depth value, ranging from 0 for pixels of unknown depth value or where
326 depth is greater than or equal to d_{MAX} , to 1 for pixels of zero depth. The
327 auditory output of the implemented algorithm was compared against the
328 original vOICe software for Windows on a small benchmark set of 10 depth
329 maps from the NYU-Depth Dataset V2⁴ [44], and it was found to never
330 exceed 1 dB of spectral distortion in the 0.5 – 5 kHz range.

331 The generic HRTF filter that we used is provided through the *earplug*~
332 Pure Data binaural synthesis external. The filter renders the angular position
333 of the sound source relative to the subject by convolving the incoming signal
334 with left and right HRTFs from the MIT KEMAR database⁵ [20]. For the
335 sake of consistency, the same HRTF filters were used for both FF and VC.

336 3.4. *Experimental procedure*

337 The experiment was divided in two sessions, each corresponding to a sin-
338 gular sensory substitution algorithm (FF or VC). The two sessions were con-
339 ducted on different days and the order of the sensory substitution algorithms

⁴http://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html

⁵<http://sound.media.mit.edu/resources/KEMAR.html>

340 was randomized and balanced. A single experimental session was composed
341 of three parts presented in the following order: a self-training part, a guided
342 training part, and an experimental test. The purpose of the training was to
343 allow for sufficient interaction with the system and to gain experience with
344 the sonification algorithm prior to the experimental test, where the actual
345 performance data was collected. The duration of the self- and guided train-
346 ing was approximately 10 and 65 minutes, respectively, while the average
347 duration of the experimental test was approximately 40 minutes.

348 *3.4.1. Self-training*

349 Basic information about the sensory substitution algorithm was first pro-
350 vided to participants through a short written description (7 lines) on an
351 experimental sheet, transcribed in the Appendix. Then, participants wore
352 the pair of headphones and freely interacted via keyboard with a simplified
353 demo of the system representing a single virtual object in the field of view
354 of the camera. Participants controlled the azimuth, elevation, distance, and
355 size of the object (see key assignment below), and directly listened to the
356 corresponding sonification:

- 357 • numpads 1 – 9: change the direction of the object on a 3×3 grid: 3
358 azimuths (left, center, right) and 3 elevations (up, middle, down);
- 359 • arrow keys up/down: increase/decrease the distance of the object be-
360 tween $0.5m$ and $3m$, in $0.5m$ steps;
- 361 • keys $+/-$: increase/decrease the size of the object (in terms of % of
362 the occupied area in that sector) from 0% to 100%, in 10% steps.

363 The self-training was designed to introduce participants to the sensory sub-
364 stitution algorithm and the underlying mappings.

365 *3.4.2. Guided training*

366 Participants were equipped with the system (backpack/PC, camera head-
367 band, blindfold, headphones) and then guided through five consecutive train-
368 ing steps as follows.

369 **Step A (3 minutes).** Participants listened interactively to the soni-
370 fication of an empty testing area while being allowed to freely explore the
371 empty room (only being stopped when going too close to an obstacle, e.g.
372 the desk or a wall). Additionally to the floor, at this stage, it was important
373 for participants to listen to and recognize the sonification of walls, ceiling
374 and other fixed objects in the room.

375 **Step B (7 minutes).** One object (made of two or three boxes on top of
376 each other in turn) was placed in the middle of the testing area and partic-
377 ipants were asked to interact with it. Participants were encouraged (guided
378 if necessary) to systematically explore the sonification output in relation to
379 changing their own position, e.g. to (1) go towards/away from the object
380 while facing it, therefore experiencing distance changes, while getting verbal
381 feedback on the current distance; (2) circle the object and stand aside of
382 it while trying to locate it with only head movements; (3) stand $2m$ away,
383 face the object and tilt the head up/down in order to experience elevation
384 changes. At this stage it was important to let participants realize through
385 training that objects closer than $0.4m$ or further than $3m$ were not repre-
386 sented; therefore, participants were invited to explore and experience at what
387 distance the sonification of the object stopped.

388 **Step C (15 minutes).** Participants trained scenes with *a single* object
389 (made of two or three boxes on top of each other) positioned in randomly
390 chosen locations of the testing area within the represented distance range.
391 Pink noise was played on the headphones in order to mask the sound of
392 boxes being moved when preparing the next scene. The participants' task
393 was to first point at the object after head movement only, tell its approximate
394 distance (in meters) and size (2 or 3 boxes), and then to go towards it and
395 touch it. From this step onwards, after successful completion of each scene,
396 participants were invited to temporarily remove the blindfold in order to
397 check the scene they just accomplished.

398 **Step D (20 minutes).** Participants trained scenes with *two* objects
399 (each made of two or three boxes on top of each other) positioned in randomly
400 chosen locations of the testing area within the represented distance range,
401 provided that they were positioned no less than $0.8m$ apart from each other
402 in order to be able to comfortably pass between them. The participants' first
403 task was to point at each object in turn after head movement only and tell
404 again their approximate distance and size. After successful completion of the
405 first task, participants were asked to walk between and past the two objects
406 trying not to touch or collide with them.

407 **Step E (20 minutes).** Participants trained a number of scenes with *two*
408 *or three* objects (randomized), aiming to find their way towards the small
409 speaker placed at a randomly chosen point on the opposite side of the testing
410 area and playing easy-listening pop music [54] at a comfortable level. The
411 obstacles (again 2 or 3 boxes on top of each other) were placed randomly
412 within the testing area, provided that they were positioned no less than

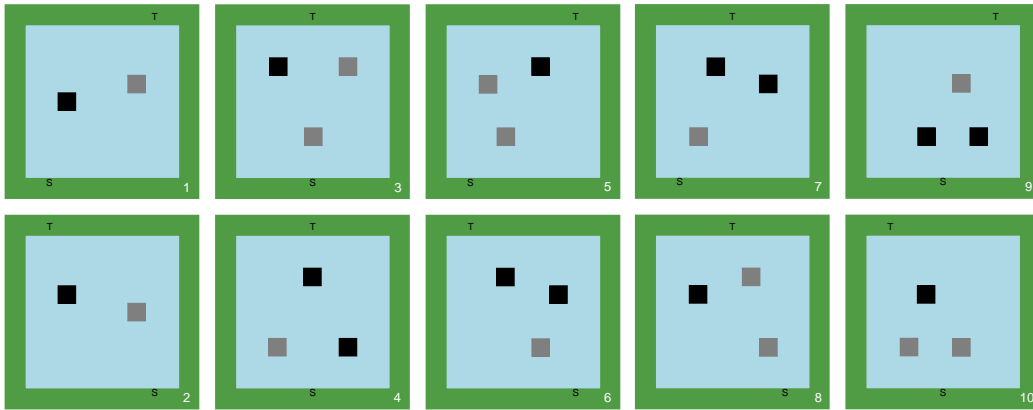


Figure 3: The 10 testing scenes. The 2-box obstacles are depicted as gray squares, and the 3-box obstacles as black squares. The starting and target (end) points are marked with S and T, respectively.

413 0.8m apart from each other (to all sides). Participants were asked to walk
 414 as carefully as possible trying not to touch or collide with the obstacles, to
 415 stay inside the testing area all the time, and to scan the environment before
 416 moving forward. At this stage it was important to tell participants that the
 417 tripod would be represented through sound as well, that they should walk
 418 towards the target without detour (especially when starting on the edges of
 419 the testing area), and that if close to the target, they should try to touch the
 420 target promptly.

421 In order to reduce fatigue, a mandatory 10-minute break was introduced
 422 between Step D and Step E. Participants were invited to take off the system
 423 and relax.

424 3.4.3. *Experimental test*

425 Right after the training, the blindfolded participants tested 10 wayfinding
 426 scenes with two or three objects always positioned within the path towards

427 the target, with a task similar to training step E. However, this time the
428 obstacles (2 or 3 boxes on top of each other each) were not placed randomly
429 within the testing area but in predefined locations, as well as the starting
430 and target (end) points, as shown in Figure 3. The order of the 10 scenes
431 was randomized for each participant and each session. Participants were
432 reminded to walk as carefully as possible, to scan the environment before
433 moving forward, and to walk towards the target without detour. Participants
434 were informed that their goal was to reach the target speaker trying to avoid
435 any collision with obstacles and without leaving the testing area, and that
436 all errors would be counted. For each experimental testing, collected data
437 included:

- 438 • number of collisions with obstacles, while differentiating between mi-
439 nor collisions (i.e., not moving boxes from their position, for instance
440 brushing on them) and major collisions (i.e., boxes moved);
- 441 • number of times the participant left the testing area by treading, even
442 partially, on the carpet (except when in the target’s vicinity);
- 443 • completion time (in seconds, taken with a timer), defined as the time
444 between the moment when the sonification was turned on and the mo-
445 ment when the participant touched the speaker or tripod.

446 After completion of all experimental testing scenes, participants were asked
447 to reply to a questionnaire about the corresponding sensory substitution
448 algorithm by ticking one item in each of three 7-point Likert scales (1 =
449 strongly disagree, 7 = strongly agree):

- 450 1. I feel I could directly understand the meaning of the sounds without
451 training;
- 452 2. I feel that the sounds are pleasant;
- 453 3. I would feel comfortable hearing these sounds on a daily basis.

454 *3.5. Statistical analysis*

455 After an exploratory data analysis on all categories of navigation errors,
456 a more advanced analysis was performed. Due to the dependent, nested
457 structure of the data, and to factor in covariates, linear mixed models with
458 fixed and random effects [36] were fit in R version 3.4.1 (R Development Core
459 Team 2017). The within-subjects design of the current study allowed to sta-
460 tistically control for the differences across participants in every analysis by
461 taking individual variance as random effect into account, which might oth-
462 erwise distort the results. Additionally, training effects might influence the
463 outcome, meaning that participants accomplished more scenes without navi-
464 gation errors when they went through the training and testing procedure for
465 the second time compared to the first time, independent of the sensory sub-
466 stitution algorithm. By randomizing the sequence of the two algorithms, any
467 systematical influence due to training effects was experimentally controlled
468 for. Yet, the training effect might lead to substantial additional variance in
469 the data, which is why it was statistically controlled for by being factored in
470 as random effect into all analyses.

471 *3.5.1. Analysis of performance data*

472 In order to compare the performance between the two sensory substitu-
473 tion algorithms, the probability of passing a scene (meaning the participants

474 did neither collide with any obstacle nor leave the testing area) for each
475 of the two algorithms was calculated, set as outcome variable and fit in a
476 Generalized Linear Mixed Model (GLMM). Due to the categorical nature
477 of the outcome variable, a mixed-effects binomial logistic regression model
478 was performed [22, 24] by executing the the `glmer()` function as part of the
479 `lme4` package in R [5]. For parameter estimation in the GLMM, in order
480 to approximate true likelihood, the Laplace approximation method with an
481 adaptive algorithm using one integration point was performed [7].

482 A model selection process was the first step of the performance analy-
483 sis, in which the improvement of model fits for three different models was
484 compared. Firstly, *Model 0* (a baseline model not containing any fixed pre-
485 dictor but only the random effects of individuals and training) was compared
486 to *Model 1* (with algorithm added as one fixed predictor) in order to deter-
487 mine if taking in algorithm as predictor into the model significantly improves
488 the variance explained by the model. If so, algorithm would have a signif-
489 icant effect on the probability of passing a scene. Secondly, *Model 1* was
490 compared to *Model 2* (with time that was necessary for scene completion
491 added as second fixed predictor, besides algorithm) in order to determine if
492 adding time as predictor significantly improves the variance explained. If so,
493 time would have a significant effect on the probability of passing a scene. A
494 Chi-square distributed Likelihood Ratio Test was performed to determine if
495 the difference between models was significant and therefore select the best
496 model. Finally, the model with the best fit was reported with regression co-
497 efficients, effect direction, confidence intervals and the predictors significance
498 was ascertained with the Wald statistics [58].

499 *3.5.2. Analysis of time data*

500 In the performance analysis described above, the time that participants
501 needed to complete a scene was only indirectly taken into account as pos-
502 sible predictor for passing as scene. However, we were mainly interested in
503 answering the question if the choice of sensory substitution algorithm results
504 in significantly different times (while statistically controlling for training and
505 individual effects). To address this, a subset of data only including passed
506 scenes was created and analyzed with time as continuous outcome variable.
507 This approach was chosen since the occurrence of navigation errors hint at
508 the possibility that scenes were not represented understandably and partic-
509 ipants were not able to interpret the obstacle location, which questions the
510 sense of interpreting failed scenes.

511 A Linear Mixed Model with algorithm as fixed effect and individual dif-
512 ferences and training as random effects was fit using Restricted Maximum
513 Likelihood (REML) [36]. We performed the `lmer()` function as part of the
514 `lme4` package to fit the LMM in R [5], as well as the `lmerTest` package⁶
515 to test if the predictor of the proposed model was significant. The pack-
516 age provides F-test statistics by calculating the degrees of freedom with the
517 Satterthwatie approximation method [41].

518 *3.5.3. Analysis of questionnaire data*

519 We finally investigated for differences in individual questionnaire scores
520 between the two algorithms by running three separate Wilcoxon signed-rank
521 tests, one per questionnaire item (intuitiveness, pleasantness and usability,

⁶<https://CRAN.R-project.org/package=lmerTest>

522 respectively). The choice of the Wilcoxon signed-rank test was due to the
523 within-participants design and to the non-normal distribution of the ques-
524 tionnaire data. Before applying each test, we verified the assumption that
525 the distribution of the differences between the two related groups was sym-
526 metrical in shape by checking that its skew value was between -2 and 2 [27].

527 4. Results

528 The complete individual results from the experiment are reported in Ta-
529 ble 1. In the table, variables C_{MIN} (number of minor collisions), C_{MAJ}
530 (number of major collisions), N_{OUT} (number of times the participant left the
531 testing area), and T_{TOT} (completion time) are aggregated for the 10 scenes.
532 It can be noticed that a lower average number in all types of navigation errors
533 was registered for FF compared to VC.

534 4.1. Performance

535 First, we compared the performance between the two algorithms, FF and
536 VC. To assess whether a scene was reliably and understandably represented
537 by the algorithm, the number of passed scenes was counted. The results show
538 that when using FF, 107 (out of 140) scenes were successfully completed by
539 participants (therefore fulfilling our hypothesis no.1), compared to 77 (out of
540 140) when the same participants used VC.

541 In order to assess whether the higher proportion of passed scenes with
542 FF was statistically significant (on alpha level of .05), the influence of the
543 algorithm on the probability of passing a scene was determined as described
544 in Section 3.5.1 following a model selection process. The results for Model
545 1 and Model 2 are reported in Table 2 with regression coefficients, standard

546 errors, confidence intervals and Wald statistics per predictor. Whereas all
547 models included individual variance and training as random effects, the basic
548 model (Model 0) did not contain any fixed predictors, which is why it is not
549 presented in the table, but served as baseline model for comparison to Model
550 1.

551 According to the Likelihood Ratio Test (LRT), including the predictor of
552 sensory substitution algorithm (Model 1) significantly improved the model
553 fit compared to an empty model without predictors (Model 0), $\chi^2(1, N =$
554 $280) = 20.15, p < .001$. This result indicates that the choice of algorithm,
555 FF or VC, has a significant effect on the outcome variable of performance,
556 meaning that the probability that participants performed a scene without
557 errors was significantly higher when they followed FF compared to VC.

558 In Model 2, time was included as additional fixed predictor to test if it
559 had a significant influence on the performance. We expected that a short
560 completion time, even though at first glance seemingly positive, might in-
561 dicate that participants rushed through the scenes since they were lacking
562 understanding of the scene resulting in collisions. However, including time
563 as predictor (additionally to algorithm) does not significantly improve the
564 model according to the LRT, $\chi^2(1, N = 280) = 2.50, p = .105$, meaning that
565 the completion time is not a predictor for more passed scenes.

566 To summarize, Model 1, only including algorithm as fixed effect while
567 factoring individual variance and training as random effects, explains most of
568 variance in the data. Adding time as predictor does not improve the model fit.
569 The Wald statistics for each fixed predictor of Model 1, reported in Table 2,
570 confirm the significant effect of the sensory substitution algorithm (improving

571 our expectations as stated in hypothesis no.2) and the non-significant effect
572 of time on the probability of passing a scene.

573 *4.2. Time*

574 As shown above, including time as fixed effect to predict if a scene was
575 passed does not significantly improve the model fit, thereby suggesting that
576 if participants completed a scene either quickly or slowly is not related to the
577 fact that the scene was mastered without errors or not.

578 The aim of the detailed time analysis was to investigate if the differ-
579 ent sensory substitution algorithms lead to significantly different completion
580 times. Thus, for the analysis, a subset of data only including passed scenes
581 was created and fit in a LMM with time as continuous outcome variable,
582 algorithm as fixed and individuals and training as random effects, as de-
583 scribed in Section 3.5.2. The resulting parameter is the regression coefficient
584 for the fixed predictor of algorithm (on time as outcome variable), $B = 4.98$
585 $[-1.91, 11.87]$ with $SE = 3.52$, indicating that the choice of algorithm does
586 not influence the time needed for completing the scenes ($F(170, 184) = 2.01$,
587 $p = .159$). In conclusion, using FF does not cause participants to either
588 complete a scene faster or slower, compared to VC.

589 *4.3. Questionnaires*

590 The histograms in Figure 4 report the scores given to each of the 3 ques-
591 tionnaire items. The support in favour of the FF algorithm compared to
592 the VC algorithm was almost unanimous and reflected in all scores, in line
593 with our hypothesis no.3. Intuitiveness FF scores were significantly higher

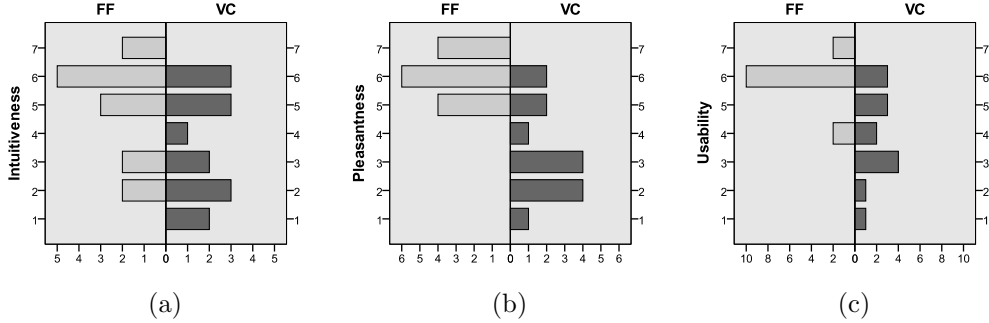


Figure 4: Histograms of questionnaire scores. (a) Intuitiveness. (b) Pleasantness. (c) Usability.

594 ($Z = -2.81, p = .005$) than VC scores (medians: FF = 5.5, VC = 3.5). Sim-
 595 ilarly, usability FF scores were significantly higher ($Z = -2.96, p = .003$)
 596 than VC scores (medians: FF = 6, VC = 4). More interestingly, an over-
 597 whelming difference was found in the pleasantness scores (medians: FF = 6,
 598 VC = 3), according to which participants highly significantly preferred FF
 599 to VC ($Z = -3.2, p = .001$). All the participants judged FF sounds pleas-
 600 ant, while 9 participants out of 14 negatively judged the pleasantness of VC
 601 sounds. Only one participant gave an equal rating to the two types of sounds,
 602 while all other participants gave a higher score to FF sounds.

603 5. Discussion

604 The *fluid flow* sensory substitution algorithm proved to be a usable and
 605 informative sensory substitution scheme for recognizing the location of ob-
 606 stacles in a simplified blind wayfinding task. This conclusion is supported by
 607 the experimental results on a pool of blindfolded sighted participants, who
 608 managed to complete the task in 76% of the proposed scenes. It has to be

609 remarked that the majority of the scenes (see Figure 3) required the par-
610 ticipants to travel through spaces as narrow as 80cm without even brushing
611 against an obstacle. If we apply a minimum tolerance on the committed
612 navigation errors and allow for one minor collision per scene, which in the
613 majority of cases meant that participants recognized the obstacle but did
614 not keep enough distance while walking past it, the percentage of completed
615 scenes grows to 86%.

616 Remarkably, our experimental results indicate a statistically significant
617 superior performance of the fluid flow algorithm compared to the vOICE al-
618 gorithm in terms of obstacle avoidance and navigation accuracy. This finding
619 is supported by qualitative evaluations from the participants collected at the
620 end of each session. For instance, a subset of participants remarked that they
621 preferred to scan the environment themselves by rotating their heads rather
622 than let the algorithm scan at a fixed rate. This remark supports the use
623 of real-time representation of the environment as provided by the fluid flow
624 scheme rather than the vOICE, whose inherently scanning nature combined
625 with head motion results in an unnatural “scan within a scan” not easy to
626 manage for some participants, at least following a short training session. An-
627 other subset of participants reported, following a collision with an obstacle,
628 to have “lost” the obstacle vOICE representation while moving; this issue can
629 also be related to the lack of a real-time feedback for effectively tracking ob-
630 stacles not only during head movement but also during body movement. Due
631 to the high cognitive load on the working memory imposed by the double-
632 scanning with the vOICE algorithm, two participants reported headache after
633 2 hours of training, which did not occur with the real-time presentation used

634 by the fluid flow.

635 On the other hand, one participant deemed the vOICe algorithm to be
636 more convincing in delivering the spatial layout of the obstacles due to the
637 clear left-to-right scanning mechanism. The participant reported that he
638 found the liquid sound representation of obstacles more difficult to separate
639 when there were two or more obstacles in the field of view of the camera, and
640 that he needed head and body movement to resolve the scene layout. This
641 remark may hint at the necessity of a more consistent training with the fluid
642 flow algorithm in static conditions.

643 As reported in the previous section, the time required to complete the
644 scenes was not significantly different between the two algorithms. Two par-
645 ticipants scored exceptionally good performances, completing most scenes
646 without errors and in less than 30 seconds each, independently of the sen-
647 sory substitution algorithm. This results indicates a ceiling effect for certain
648 participants, meaning that the scenes were too easy for them to accomplish
649 and therefore they were not able to differentiate between the two algorithms.
650 The ceiling effects covers potential differences between the algorithms; how-
651 ever, this issue only applied for two out of 14 participants. Some participants
652 were on average both faster and more accurate with the fluid flow algorithm
653 than with the vOICe, while other participants considerably slowed down
654 when using the fluid flow sounds. When asked about the latter behaviour,
655 one participant (at the end of her second session) stated that she had a much
656 better understanding of the scene with the fluid flow sounds and felt like she
657 had more control about her performance than with the vOICe algorithm,
658 and therefore devoted more attention to complete the scene without errors.

659 This conduct is consistent with the fact that prior to the experimental test
660 participants were clearly informed that their task was to minimize navigation
661 errors and not race against time.

662 The proposed algorithm directly receives as input reliable low-level infor-
663 mation conveyed through an off-the-shelf depth sensor, contrary to other sen-
664 sory substitution schemes previously explored by the authors [9, 50, 49, 13]
665 that used obstacle information segmented through computationally heavy
666 image processing techniques. This is a very desirable property in a system
667 that needs to be scalable in order to run on smartphones or embedded sys-
668 tems with low processing power. The scalability of the proposed approach is
669 further supported by the possibility of reducing the resolution of the depth
670 map without considerable loss of information, as well as changing the size
671 of the oscillator bank for each liquid sound generator at the price of sound
672 quality [4]. This would allow for graceful degradation of our rendering ap-
673 proach depending on the available computational resources. Future work will
674 investigate the quality of experience of the sounds produced by the sensory
675 substitution algorithm even in cases of limited computing power.

676 One limitation of the current study lies in the use of a sensor with limited
677 field of view and range information, that disoriented some participants in
678 that the obstacle sonification stopped when getting close enough to it, and
679 required considerable head rotation (both yaw and pitch) for a full scan of
680 the scene. Furthermore, although not directly investigated in this study, the
681 choice of the spatialization technique has an undeniable impact on the spatial
682 perception of sounds, and therefore on the degree of immersion [34] and
683 overall quality of experience. The most effective solution would be the use of

684 individual HRTFs measured on the listener with the addition of head tracking
685 and artificial reverberation [6, 56]. However, obtaining acoustically measured
686 individual HRTF data is only possible with tailored equipment and invasive
687 recording procedures [12]. On the other hand, even though one participant
688 to our study commented that he could “clearly visualize columns of bubbles”
689 where the obstacles were, using non-individual HRTFs is only effective for
690 a limited number of individuals. Different alternative approaches towards
691 HRTF-based spatial rendering were proposed throughout the last decades,
692 ranging from HRTF selection [43, 21] to structural HRTF models [8, 47]
693 and numerical HRTF simulations [26, 59]. Such approaches are expected to
694 progressively bridge the gap between accessibility and accuracy of individual
695 spatial audio [51]. Still, in cases of limited computing power, HRTF rendering
696 can be substituted by constant-power panning [29] to represent horizontal
697 direction at least.

698 Validation with sighted users implies that these results should only be
699 generalized to the visually impaired population with caution. Blind users
700 are generally more adapted to rely on their sense of hearing for orientation
701 and solving daily mobility challenges compared to sighted, e.g. by using
702 echolocation techniques [42]. This might result in even lower training time
703 required for VIPs to successfully apply the fluid flow algorithm. Furthermore,
704 dynamic postural stability is affected by the visual system, which is why the
705 postural stability of sighted individuals with eyes closed has been shown to
706 be superior to that of blind people [2]. This might result in more collisions
707 when VIPs perform the same task compared to sighted people, even when
708 the obstacle is correctly located in the first place. Hence, to control for these

709 possible differences between sighted and blind, similar evaluations of the fluid
710 flow algorithm are currently being carried out within the *Sound of Vision*⁷
711 project, ranging from virtual to complex real world environments [13], re-
712 quired for assessing the usability of the system outside the laboratory.

713 In the final questionnaire, participants reported a clear preference for the
714 fluid flow sounds compared to the vOICE sounds, in terms of intuitiveness,
715 pleasantness, and usability. This result is of great relevance for the integra-
716 tion of the fluid flow sounds in a sensory substitution system for VIPs. Our
717 belief, supported by several participant comments in addition to the ques-
718 tionnaire scores, is that a natural, intuitive, and aesthetically pleasant sonic
719 representation requires little time and effort to be learned while at the same
720 time allowing for longer and less fatiguing practice sessions [45]. In a seminal
721 paper from 2003, yet still as current today as ever, Rocchesso *et al.* [38] assert
722 that “*an aesthetic mismatch exists between the rich, complex, and informa-*
723 *tive soundscapes in which mammals have evolved and the poor and annoying*
724 *sounds of contemporary life in today’s information society*”, recognizing “*the*
725 *need for sounds that can convey information about the environment yet be*
726 *expressive and aesthetically interesting.*” In our view, the use of physically
727 based, natural-sounding liquid sounds perfectly matches this need within the
728 field of sensory substitution.

⁷<https://soundofvision.net/>

729 **Acknowledgments**

730 The authors would like to thank Stefano Baldan for his support with the
731 SDT and all the participants involved in this study. This project has received
732 funding from the European Union’s Horizon 2020 research and innovation
733 programme under grant agreement No. 643636.

734 **Appendix A. Experimental sheet descriptions**

735 **FF.** The system converts the video stream into a liquid streaming sound
736 produced through superposition of bubble sounds. Bubbles simultaneously
737 come from the visible objects direction in space. The bigger the volume
738 occupied by an object in the visible space, the richer the texture of the
739 corresponding streaming sound (i.e., more bubbles produced). The higher
740 the position of the object in the visible space, the fizzier the bubbles sound.
741 The closer an object within the represented distance range, the louder the
742 liquid streaming sound. If the object gets closer than $1m$, bubbles begin to
743 present a characteristic boiling sound.

744 **VC.** The system converts the video stream into a sound made of the
745 superposition of simple tones. The acquired image is scanned in a left to
746 right scanning order, at a rate of one scan per second. Hearing some sound
747 on your left or right thus means having a corresponding object pattern on the
748 left or right side, respectively. During every scan, the higher the pitch, the
749 higher the position of objects in that direction in the visible space. Loudness
750 means distance: the louder the sound, the closer the objects in that direction
751 in the visible space. The bigger the volume occupied by an object in the

752 visible space, the richer (i.e., more simultaneous tones) and the longer the
753 corresponding sound.

754 REFERENCES

- 755 [1] Avanzini, F., Spagnol, S., Rodá, A., De Götzen, A., March 2013. De-
756 signing interactive sound for motor rehabilitation tasks. In: Franinovic,
757 K., Serafin, S. (Eds.), *Sonic Interaction Design*. MIT Press, Cambridge,
758 MA, USA, Ch. 12, pp. 273–283.
- 759 [2] Aydoğ, E., Aydoğ, S. T., Cakci, A., Doral, M. N., May 2006. Dynamic
760 postural stability in blind athletes using the biodex stability system. *Int.*
761 *J. Sports Med.* 27 (5), 415–418.
- 762 [3] Balakrishnan, G., Sainarayanan, G., Nagarajan, R., Yaacob, S., 2008.
763 A stereo image processing system for visually impaired. *Int. J. Signal*
764 *Process.* 2 (3), 136–145.
- 765 [4] Baldan, S., Delle Monache, S., Rocchesso, D., 2017. *The Sound Design*
766 *Toolkit*. Software XIn press.
- 767 [5] Bates, D., Mächler, M., Bolker, B., Walker, S., October 2015. Fitting
768 linear mixed-effects models using lme4. *J. Stat. Softw.* 67 (1), 48 pp.
- 769 [6] Begault, D. R., Wenzel, E. M., Anderson, M. R., October 2001. Direct
770 comparison of the impact of head tracking, reverberation, and individ-
771 ualized head-related transfer functions on the spatial perception of a
772 virtual speech source. *J. Audio Eng. Soc.* 49 (10), 904–916.

- 773 [7] Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R.,
774 Stevens, M. H. H., White, J. S. S., March 2009. Generalized linear mixed
775 models: A practical guide for ecology and evolution. *Trends Ecol. Evol.*
776 24 (3), 127–135.
- 777 [8] Brown, C. P., Duda, R. O., September 1998. A structural model for
778 binaural sound synthesis. *IEEE Trans. Speech Audio Process.* 6 (5),
779 476–488.
- 780 [9] Bujacz, M., Kropidłowski, K., Ivanica, G., Moldoveanu, A., Saitis, C.,
781 Csapó, A., Wersényi, G., Spagnol, S., Jóhannesson, O. I., Unnthórsson,
782 R., Rotnicki, M., Witek, P., July 2016. Sound of Vision - Spatial audio
783 output and sonification approaches. In: Miesenberger, K., Bühler, C.,
784 Penaz, P. (Eds.), *Computers Helping People with Special Needs - 15th*
785 *International Conference (ICCHP 2016)*. Vol. 9759 of *Lecture Notes in*
786 *Computer Science*. Springer Int. Publishing, Linz, Austria, pp. 202–209.
- 787 [10] Bujacz, M., Strumiłło, P., October 2016. Sonification: Review of au-
788 ditory display solutions in electronic travel aids for the blind. *Arch.*
789 *Acoust.* 41 (3), 401–414.
- 790 [11] Capp, M., Picton, P., June 2000. The optophone: an electronic blind
791 aid. *Eng. Sci. Educ. J.* 9 (3), 137–143.
- 792 [12] Cheng, C. I., Wakefield, G. H., April 2001. Introduction to head-related
793 transfer functions (HRTFs): Representations of HRTFs in time, fre-
794 quency, and space. *J. Audio Eng. Soc.* 49 (4), 231–249.

- 795 [13] Csapó, A., Spagnol, S., Herrera Martínez, M., Bujacz, M., Janeczek, M.,
796 Ivanica, G., Wersényi, G., Moldoveanu, A., Unnthórsson, R., May 2017.
797 Usability and effectiveness of auditory sensory substitution models for
798 the visually impaired. In: Proc. 142nd Conv. Audio Eng. Soc. No. 9801.
799 Berlin, Germany, p. 10 pp.
- 800 [14] Csapó, A., Wersényi, G., Nagy, H., Stockman, T., December 2015. A
801 survey of assistive technologies and applications for blind users on mobile
802 platforms: A review and foundation for research. *J. Multimod. User*
803 *Interf.* 9 (4), 275–286.
- 804 [15] Dakopoulos, D., Bourbakis, N. G., January 2010. Wearable obstacle
805 avoidance electronic travel aids for blind: A survey. *IEEE Trans. Syst.*
806 *Man Cybern.* 40 (1), 25–35.
- 807 [16] Dell’Aversana, P., Gabbriellini, G., Amendola, A., January 2017. Sonifi-
808 cation of geophysical data through time-frequency analysis: Theory and
809 applications. *Geophys. Prospect.* 65 (1), 146–157.
- 810 [17] Dubus, G., Bresin, R., December 2013. A systematic review of mapping
811 strategies for the sonification of physical quantities. *PLoS One* 8 (12),
812 28 pp.
- 813 [18] Fontana, F., Fusiello, A., Gobbi, M., Murino, V., Rocchesso, D., Sartor,
814 L., Panuccio, A., 2002. A cross-modal electronic travel aid device. In:
815 *Human Computer Interaction with Mobile Devices*. Vol. 2411 of *Lecture*
816 *Notes in Computer Science*. Springer Berlin Heidelberg, pp. 393–397.

- 817 [19] Garcia-Ruiz, M. A., Gutierrez-Pulido, J. R., July 2006. An overview
818 of auditory display to assist comprehension of molecular information.
819 *Interact. Comput.* 18 (4), 853–868.
- 820 [20] Gardner, W. G., Martin, K. D., June 1995. HRTF measurements of a
821 KEMAR. *J. Acoust. Soc. Am.* 97 (6), 3907–3908.
- 822 [21] Geronazzo, M., Spagnol, S., Avanzini, F., August 2017. Do we need
823 individual head-related transfer functions for vertical localization? The
824 case study of a spectral notch distance metric. Submitted for publication.
- 825 [22] Hartzel, J., Agresti, A., Caffo, B., July 2001. Multinomial logit random
826 effects models. *Stat. Model.* 1 (2), 81–102.
- 827 [23] Hermann, T., Hunt, A., Neuhoff, J. G. (Eds.), November 2011. The
828 Sonification Handbook, 1st Edition. Logos Publishing House, Berlin,
829 Germany.
- 830 [24] Hosmer, D. W., Lemeshow, S., Sturdivant, R. X., August 2013. Logistic
831 regression models for multinomial and ordinal outcomes. In: *Applied*
832 *Logistic Regression*, 3rd Edition. John Wiley & Sons, Inc., New York,
833 NY, USA, pp. 269–311.
- 834 [25] Jamal, Y., Lacey, S., Nygaard, L., Sathian, K., 2017. Interactions be-
835 tween auditory elevation, auditory pitch and visual elevation during mul-
836 tisensory perception. *Multisens. Res.* 30 (3-5), 287–306.
- 837 [26] Katz, B. F. G., November 2001. Boundary element method calculation
838 of individual head-related transfer function. I. Rigid model calculation.
839 *J. Acoust. Soc. Am.* 110 (5), 2440–2448.

- 840 [27] Kim, H.-Y., February 2013. Statistical notes for clinical researchers:
841 Assessing normal distribution (2) using skewness and kurtosis. *Restor.*
842 *Dent. Endod.* 38 (1), 52–54.
- 843 [28] Kristjánsson, A., Moldoveanu, A., Jóhannesson, O. I., Balan, O., Spag-
844 nol, S., Valgeirsdóttir, V. V., Unnthórsson, R., October 2016. Designing
845 sensory-substitution devices: Principles, pitfalls and potential. *Restor.*
846 *Neurol. Neurosci.* 34 (5), 769–787.
- 847 [29] Lee, S.-L., Han, K.-Y., Lee, S.-R., Sung, K.-M., August 2004. Reduction
848 of sound localization error for surround sound system using enhanced
849 constant power panning law. *IEEE Trans. Consum. Electr.* 50 (3), 941–
850 944.
- 851 [30] Meijer, P. B. L., February 1992. An experimental system for auditory
852 image representations. *IEEE Trans. Biomed. Eng.* 39 (2), 112–121.
- 853 [31] Merabet, L., Battelli, L., Obretenova, S., Maguire, S., Meijer, P. B. L.,
854 Pascual-Leone, A., January 2009. Functional recruitment of visual cor-
855 tex for sound encoded object identification in the blind. *Neuroreport*
856 20 (2), 132–138.
- 857 [32] Minnaert, M., 1933. On musical air-bubbles and the sounds of running
858 water. *Phil. Mag.* 16, 235–248.
- 859 [33] Møller, H., Sørensen, M. F., Jensen, C. B., Hammershøi, D., June 1996.
860 Binaural technique: Do we need individual recordings? *J. Audio Eng.*
861 *Soc.* 44 (6), 451–469.

- 862 [34] Nilsson, N., Nordahl, R., Serafin, S., November 2016. Immersion revis-
863 ited: A review of existing definitions of immersion and their relation to
864 different theories of presence. *Hum. Tech.* 12 (2), 108–134.
- 865 [35] Pasqualotto, A., Esenkaya, T., April 2016. Sensory substitution: the
866 spatial updating of auditory scenes mimics the spatial updating of visual
867 scenes. *Front. Behav. Neurosci.* 10 (79).
- 868 [36] Pinheiro, J. C., Bates, D. M., 2000. *Mixed-Effects Models in S and S-*
869 *PLUS. Statistics and Computing.* Springer-Verlag, New York, NY, USA.
- 870 [37] Roads, C., Summer 1988. Introduction to granular synthesis. *Comput.*
871 *Music J.* 12 (2), 11–13.
- 872 [38] Rocchesso, D., Bresin, R., Fernström, M., April–June 2003. Sounding
873 objects. *IEEE Multimedia* 10 (2), 42–52.
- 874 [39] Rosati, G., Oscari, F., Reinkensmeyer, D. J., Secoli, R., Avanzini, F.,
875 Spagnol, S., Masiero, S., June 2011. Improving robotics for neuroreha-
876 bilitation: Enhancing engagement, performance, and learning with au-
877 ditory feedback. In: *Proc. IEEE 12th Int. Conf. Rehab. Rob. (ICORR*
878 *2011).* Zurich, Switzerland, pp. 341–346.
- 879 [40] Scaletti, C., 1994. Sound synthesis algorithms for auditory data repre-
880 sentations. In: Kramer, G. (Ed.), *Auditory Display: Sonification, Audi-*
881 *fication, and Auditory Interfaces.* Vol. 1. Adison-Wesley, Reading, MA,
882 USA, pp. 223–251.
- 883 [41] Schaalje, G. B., McBride, J. B., Fellingham, G. W., December 2002.

- 884 Adequacy of approximations to distributions of test statistics in complex
885 mixed linear models. *J. Agric. Biol. Environ. Stat.* 7 (4).
- 886 [42] Schenkman, B. N., Nilsson, M. E., April 2010. Human echolocation:
887 Blind and sighted persons' ability to detect sounds recorded in the pres-
888 ence of a reflecting object. *Perception* 39 (4), 483–501.
- 889 [43] Seeber, B. U., Fastl, H., July 2003. Subjective selection of non-individual
890 head-related transfer functions. In: *Proc. 2003 Int. Conf. Auditory Dis-*
891 *play (ICAD03)*. Boston, MA, USA, pp. 259–262.
- 892 [44] Silberman, N., Hoiem, D., Kohli, P., Fergus, R., October 2012. Indoor
893 segmentation and support inference from RGBD images. In: *Proc. 12th*
894 *Eur. Conf. on Computer Vision (ECCV'12)*. Florence, Italy, pp. 746–
895 760.
- 896 [45] Singh, A., Piana, S., Pollarolo, D., Volpe, G., Varni, G., Tajadura-
897 Jimnez, A., CdeC Williams, A., Camurri, A., Bianchi-Berthouze, N.,
898 2016. Go-with-the-Flow: Tracking, analysis and sonification of move-
899 ment and breathing to build confidence in activity despite chronic pain.
900 *Hum.-Comput. Interact.* 31 (3–4), 335–383.
- 901 [46] Spagnol, S., Baldan, S., Unnthórsson, R., October 2017. Auditory depth
902 map representations with a sensory substitution scheme based on syn-
903 thetic fluid sounds. In: *Proc. IEEE 19th Int. Work. Multi. Signal Pro-*
904 *cess. (MMSP 2017)*. Luton, UK.
- 905 [47] Spagnol, S., Geronazzo, M., Avanzini, F., March 2013. On the relation

- 906 between pinna reflection patterns and head-related transfer function fea-
907 tures. *IEEE Trans. Audio, Speech, Lang. Process.* 21 (3), 508–519.
- 908 [48] Spagnol, S., Hiipakka, M., Pulkki, V., September 2011. A single-azimuth
909 pinna-related transfer function database. In: *Proc. 14th Int. Conf. Dig-
910 ital Audio Effects (DAFx-11)*. Paris, France, pp. 209–212.
- 911 [49] Spagnol, S., Saitis, C., Bujacz, M., Jóhannesson, O. I., Kalimeri, K.,
912 Moldoveanu, A., Kristjánsson, A., Unnthórsson, R., September 2016.
913 Model-based obstacle sonification for the navigation of visually impaired
914 persons. In: *Proc. 19th Int. Conf. Digital Audio Effects (DAFx-16)*.
915 Brno, Czech Republic, pp. 309–316.
- 916 [50] Spagnol, S., Saitis, C., Kalimeri, K., Jóhannesson, O. I., Unnthórsson,
917 R., October 2016. Sonificazione di ostacoli come ausilio alla deambu-
918 lazione di non vedenti. In: *Proc. XXI Colloquium on Music Informatics
919 (XXI CIM)*. Cagliari, Italy, pp. 47–54.
- 920 [51] Spagnol, S., Wersényi, G., Bujacz, M., Balan, O., Herrera Martínez,
921 M., Moldoveanu, A., Unnthórsson, R., September 2017. Current use
922 and future perspectives of spatial audio technologies in electronic travel
923 aids. Submitted for publication.
- 924 [52] Stoll, C., Palluel-Germain, R., Fristot, V., Pellerin, D., Alleysson, D.,
925 Graff, C., January 2015. Navigating from a depth image converted into
926 sound. *Appl. Bionics Biomech.* 2015.
- 927 [53] Striem-Amit, E., Guendelman, M., Amedi, A., March 2012. Visual acu-

- 928 ity of the congenitally blind using visual-to-auditory sensory substitu-
929 tion. PLoS One 7 (3).
- 930 [54] The High Llamas, March 1996. Hawaii [CD]. CD WOOL 2, Alpaca Park.
- 931 [55] Ulfvengren, P., 2003. Design of natural warning sounds in human-
932 machine systems. Ph.D. thesis, Royal Institute of Technology, Stock-
933 holm, Sweden.
- 934 [56] Välimäki, V., Parker, J. D., Savioja, L., Smith, J. O., Abel, J. S., July
935 2012. Fifty years of artificial reverberation. IEEE Trans. Audio, Speech,
936 Lang. Process. 20 (5), 1421–1448.
- 937 [57] van den Doel, K., October 2005. Physically-based models for liquid
938 sounds. ACM Trans. Applied Perception 2 (4), 534–546.
- 939 [58] Wald, A., November 1943. Tests of statistical hypotheses concerning
940 several parameters when the number of observations is large. Trans.
941 Am. Math. Soc. 54 (3), 426–482.
- 942 [59] Ziegelwanger, H., Majdak, P., Kreuzer, W., July 2015. Numerical calcu-
943 lation of listener-specific head-related transfer functions and sound lo-
944 calization: Microphone model and mesh discretization. J. Acoust. Soc.
945 Am. 138 (1), 208–222.

Table 1: Individual experimental results: number of minor/major collisions (C_{MIN}/C_{MAJ}), number of times the participant left the testing area (N_{OUT}) and total completion time (T_{TOT}), divided by participant and sensory substitution algorithm.

Participant ID	C_{MIN}		C_{MAJ}		N_{OUT}		T_{TOT} [s]	
	FF	VC	FF	VC	FF	VC	FF	VC
01	2	1	0	1	0	0	771	702
02	0	2	0	1	0	0	2451	1453
03	0	6	0	1	0	0	1458	1840
04	2	7	1	5	0	0	909	1292
05	6	4	2	7	1	7	717	1648
06	2	2	1	10	1	0	2172	1138
07	8	8	1	6	2	1	1983	1202
08	0	0	1	0	0	1	963	1506
09	0	1	1	4	0	0	1466	1824
10	0	0	1	0	0	0	230	228
11	3	5	5	13	0	0	822	882
12	0	0	0	0	0	0	407	344
13	1	2	1	0	0	0	410	1021
14	4	9	5	9	0	1	1880	1645
Mean	2	3.4	1.4	4.1	0.3	0.7	1188.5	1194.6
SD	2.5	3.1	1.6	4.4	0.6	1.9	713.4	513.4

Table 2: Results of calculating Generalized Linear Mixed Model for Model 1 and Model 2 including one additional predictor, each with individual variance and training as random effects. The model parameter estimates are calculated basing on Laplace approximation with 1 integration point. Shown are regression coefficients with associated standard errors (SE) and confidence intervals (CI), and Wald statistics (z-value and p -value).

	Predictor	Coeff.	SE	CI [LL,UL]	z-value	p-value
Model 1	Algorithm	-1.30	0.33	[-1.94,-0.66]	-3.96	$p < .001$
Model 2	Algorithm	-1.30	0.33	[-1.95,-0.65]	-3.94	$p < .001$
	Time	-0.01	0.01	[-0.02,0.00]	-1.64	$p = .101$