

A mobile application for keyword search in real-world scenes

Shrinivas Pundlik^{1,2}, Anikait Singh^{1*}, Gautam Baghel^{1*}, Vilte Baliutaviciute¹, Gang Luo^{1,2}

¹Schepens Eye Research Institute of Mass Eye & Ear, Boston MA

²Department of Ophthalmology, Harvard Medical School, Boston MA

*Contributed to this work while at Schepens Eye Research Institute.

Corresponding Author:

Shrinivas Pundlik (shrinivas_pundlik@meei.harvard.edu)

Schepens Eye Research Institute, 20 Staniford Street, Boston MA 02114

1 **Abstract**

2 We frequently engage in activities in our daily lives that involve visual search centered on keywords, such
3 as searching for the due date on a utility bill, or finding calories in a food product. Keyword search in a
4 cluttered environment is difficult in general, and even more challenging for people with low vision. While
5 magnification can help in reading for low vision people, it does not facilitate efficient visual search due to
6 the constriction of the field of view. The motivating observation for this study is that, in a large number of
7 visual search tasks, people know what are they looking for (i.e., they know the keywords), they just do
8 not know where to find them in the scene. We have developed a mobile application that allows the users
9 to input keywords (by voice or by typing), uses an optical character recognition (OCR) engine to search
10 for the provided keyword in the scene captured by the smartphone camera, and zooms in on the instances
11 of the keyword detected in the captured images, to facilitate efficient information acquisition. In this
12 paper we describe the development and evaluation of various aspects of the application, including
13 comparing the various mainstream OCR engines that power the app, and an evaluation study comparing
14 the app to the conventional optical magnifier vision aid. Normally sighted adults, while wearing blur
15 glasses to lower their visual acuity, performed keyword searches for a series of items ranging from easy to
16 difficult with the app and with a handheld magnifier. While there was no difference in the search times
17 between the two methods for the easier tasks, the app was significantly faster than the magnifier for the
18 difficult tasks.

19 **Keywords**

20 Low-vision aid; mobile application; optical character recognition (OCR); timed instrumental activities of
21 daily living (TIADL) tasks

22 **1. Introduction**

23 Visual search is an important, frequently performed visual task in daily life, such as looking for a street
24 name when walking, or finding the calorie content of a food product. Performing visual search in a
25 cluttered environment can be very demanding, even for those with normal vision.[1, 2] It is even more

26 challenging for people with visual impairments.[3-7] Among the various aspects of daily life that are
27 negatively affected due to visual impairment, the limitations in performing visual search are one of the
28 key challenges that hampers efficient information acquisition.[8, 9] As vision impairment can run the
29 gamut from moderate loss of visual acuity (moderate low-vision) to complete blindness (no light
30 perception), the challenges faced in performing daily tasks vary greatly between patients with different
31 conditions, and between individuals. [10, 11]

32 Visual search can be considered a spot reading task, which is the act of quickly locating and acquiring a
33 specific piece of information from a scene. People with moderate to severe vision loss (up to the level of
34 legal blindness), such as those who have central vision loss, often rely on magnification for reading or
35 discerning the details of their surroundings.[12, 13] However, increased magnification leads to a smaller
36 field of view, which makes searching for a particular detail within the scene highly inefficient. For
37 example, when searching for a particular ingredient on a product label, the information of interest is only
38 a small portion of the label (Figure 1a). However, one usually needs to go through the list from top to
39 bottom and end-to-end. When using magnification, this may take visually impaired people an even longer
40 time due to the reduction of the field of view caused by magnification. On the other hand, the visual field
41 restriction for people with peripheral vision loss (PVL), such as those with retinitis pigmentosa (RP),
42 greatly affects their visual search ability despite having good visual acuity.[14, 15] They may not have
43 problems in discerning the scene details, but often have difficulty in knowing where to look for the
44 required information and may need help in locating the targets (Figure 1b). Thus, visual search related
45 challenges could arise due to a variety of reasons, including inability to locate the targets, inability to
46 discern target details, or both. Therefore, a vision aid assisting in visual search can help people with
47 vision loss in various daily life tasks.

48 To tackle the visual search related challenges in visually impaired users, dedicated devices [16] and
49 services [17] have been developed that leverage artificial or human intelligence. Considering the growing
50 prevalence of mobile devices in the general as well as visually impaired population [18-21], mobile

51 applications aiding in visual search specifically targeting visually impaired users are also now
 52 available.[22, 23] The main idea behind many of these vision aids is to perform one or multiple functions
 53 including object detection, optical character recognition (OCR), and/or scene categorization using
 54 computer vision, and then provide some feedback to the user via a predefined tags or descriptions.
 55 Performing generic object detection can be challenging in the real world (for example, product
 56 identification based on barcode or appearance) and the predefined categories of object classes may be too
 57 restrictive to cover the rich variety of objects encountered in everyday situations. Comparatively, OCR is
 58 a much more well-defined problem with mature and established technologies available to tackle it.
 59 Searching with keywords can be intuitive and help narrow down the scope of the search, thereby
 60 improving the odds of obtaining the required information. However, dedicated OCR apps are meant for
 61 document reading instead of detection of text in the scene.[24] Even in applications that can perform
 62 OCR in scene images, the feedback provided to the user is generally not relevant, as the entire text blocks
 63 detected in the scene are continuously read to the user. We have developed a mobile app, Supervision
 64 Search (SVS) [25], that can perform keyword search in scene images so that the users can quickly and
 65 efficiently retrieve the relevant information from their surroundings or from the object of interest.

Calcium Pantothenate)		
Magnesium (as Magnesium Oxide)	10mg	3
Zinc (as Zinc Oxide)	5mg	33
Bacopa (Bacopa Monniera L.) (Stem, Leaf, Flower)	125mg	tt
Arctic Root (Rhodiola Rosea) (Root)	100mg	tt
American Ginseng (Panax Quinquefolius L.) (Root)	75mg	tt
St. John's Wort (Leaf, Flowers)	75mg	tt
Choline (as Choline Bitartrate)	70mg	tt
5-HTP (l-5-Hydroxytryptophan)	25mg	tt
DMAE (as DMAE Bitartrate)	25mg	tt
GABA (Gamma-Aminobutyric Acid)	20mg	tt

(a)



(b)

66 **Figure 1: Different aspects of visual search applied to different tasks.** (a) Rather than reading through the
 67 whole fact sheet, a person often needs to check just a few details about a product, e.g. GABA in this cognitive
 68 supplement. (b) In many navigation situations, visually impaired travelers have difficulty knowing where to
 69 look to find the needed information amongst all the available options. When using magnification to read
 70 distant text, the reduced field of view makes finding the target more difficult in cluttered environments.

71 The key arguments in support of development of SVS app are: i) searching for items represented by text
 72 or symbols forms a large part of visual search activities in the daily lives of people with low vision, ii) in

73 most cases, people already know what they are searching for (i.e., the keywords are already known), and
74 iii) if the keywords are located in the scene, then the information related to those keywords is available in
75 the general vicinity (spatially) of the detected keywords. Searching with keywords, for instance using
76 Google, is ubiquitous when the information is in digital form. The SVS app merely generalizes the same
77 approach in real-world scenes. By leveraging powerful OCR engines, the SVS app searches for user input
78 keywords, then highlights and zooms in on the found instances of the keyword in the captured scene
79 image to facilitate quick and easy retrieval of information related to the keyword.

80 In this paper, we describe the concept of keyword search in natural images and its relevance in various
81 daily life activities, detail the design of the SVS app to target different requirements of low vision users,
82 and present results of its preliminary evaluation: both of the underlying algorithms, and of the app by
83 human subjects with simulated visual acuity loss. The goal of this work is to test whether the approach of
84 keyword search in natural scene images can be utilized to design a vision aid for visually impaired people,
85 and whether the SVS app can provide additional benefits compared to conventional visual aids.

86 **2. Keyword Search in Scene Images**

87 For developing a visual search assistance application for visually impaired people, we needed to identify
88 methods for performing the search, as well as define methods for seamless interaction with the user. For
89 performing the keyword search, we relied on established optical character recognition (OCR)
90 technologies instead of developing custom algorithms from scratch.

91 **2.1 Optical Character Recognition Engine**

92 An OCR engine lies at the heart of a keyword search application, as its capabilities and limitations
93 essentially shape much of the usability of the search application. Quite simply, an OCR engine processes
94 the input image to detect text regions, recognize the characters, group them into words, and output strings
95 of text with associated metadata, such as its location in the image (coordinates of the bounding box for a
96 word), and possibly its orientation with respect to predefined axis, among other data. The main
97 application for OCR technology has been to digitize printed documents; and in the space of assistive

98 technologies for visually impaired people, OCR is extensively used in applications for assistance in
99 reading printed material and documents.[24] However to be truly useful in general scenarios, OCR needs
100 to work for text embedded anywhere in the scene, and not just documents. Thus, scene text recognition is
101 more challenging and cannot be handled very well by applications focused on document OCR processing.
102 Therefore more sophisticated OCR engines are needed.

103 Fortunately, there are various commercially available 3rd party OCR engines, available for incorporation
104 within end-user applications, which have been trained using advanced machine learning techniques to
105 work in highly demanding real-world images containing text. Four different OCR Engines that can be
106 implemented via an API were evaluated to determine the accuracy and usability of each engine: Google
107 Machine Learning Kit (ML Kit) [26], Microsoft Azure's Cognitive Services (Azure OCR) [27],
108 ABBYY's Real Time Recognition SDK (ABBY RTR) [28], and Amazon's Rekognition SDK
109 (Rekognition) [29].

- 110 • *Google Machine Learning Kit* (version 16.0.0) — This service provides machine learning models
111 for text recognition that are both native to the device and cloud-based. In this study however, we
112 only used the native version of the API that is freely available. This engine returns the words that
113 are found in the image along with the coordinates of the bounding box of each word.
- 114 • *Microsoft Azure's Cognitive Services* (version 2.0) — This cloud-based service provides image
115 processing algorithms to identify content present in the image. Pictures of the content are sent to
116 the service to be processed, and a JSON file is returned with the detected text, a bounding box for
117 each word, and its orientation angle. There is a limitation on the maximum image size (4MB) that
118 can be sent to the service. Being a cloud-based engine exclusively, it requires a network
119 connection to work and is not particularly well-suited for real-time applications.
- 120 • *ABBYY's Real Time Recognition SDK* (version 1.0.7.56 free demo version) — A native model
121 similar to Google's Machine Learning Kit, this OCR Engine can also process frames of a live
122 video stream in real time. A confidence level for detection needs to be set as an operating

123 parameter for accuracy of detection of a word within the image. This engine returns the location
124 of the line in which the word is present within the image, but the exact bounding box for a
125 particular word was not available for the version used in this work.

- 126 • *Amazon’s Rekognition SDK* (version 2.6) — This is a cloud-based service that returns the
127 detected text, the confidence level of each word, and the “geometry” of the word, which contains
128 a polygon that surrounds the word. A limitation of this OCR service (the version that was
129 available for this work) was that the image size sent to the cloud server was restricted to 5MB,
130 and it only processed up to a maximum of 50 words within the image.

131 **2.2 Evaluation of OCR Engines**

132 The accuracy and robustness of the above 4 OCR engines in finding a specific search query was evaluated
133 in a variety of natural images captured in offices, train stations, stores, and on streets in downtown Boston
134 (Figure 2). For each image, a set of most relevant keywords were identified to be tested by the OCR
135 engines. An image could have more than one keywords associated with it. Keywords were selected such
136 that they would make sense in a realistic application. For example, in an image of an intersection, the
137 likely search query would be the street name present in the image; or for a food product, the likely search
138 queries could be any of the nutritional factors. A total of 203 keywords for 117 images were identified.
139 High resolution images were captured using Samsung Galaxy S7, S8, and Google Pixel smartphones;
140 although the images were not of exact same size (Mean \pm SD diagonal was 4939 \pm 258 pixels). Each image
141 was tested at 3 different zoom levels — 1x, 2x, and 4x. If an OCR engine was not able to find the
142 keyword at the lowest zoom level, then next zoom level was tested. Images were zoomed such that the
143 search keyword was present near the center of the image. The Lanczos interpolation was implemented, as
144 it allows for detailed upsampling of the image and preservation of small text. The true location of the
145 keyword within the captured images was determined manually for comparison of the spatial accuracy of
146 the keyword localization. For each keyword, the detected text and the location of each word was saved
147 (with the exception of ABBYY Real Time Recognition service, where it was not possible to obtain the

148 location of the detected keyword from the SDK version that we were using for testing).

149 The performance of the OCR engines was evaluated in terms of the number of keywords successfully

150 found, and the zoom level at which they were found. Particularly, we were interested in determining the

151 cumulative detection success and the success in detection at the base scale (no zoom or 1x condition) for

152 performance considerations. For both of these outcomes, the proportions of successfully detected

153 keywords were compared between the 4 OCR engines. The closeness of the location of the detected

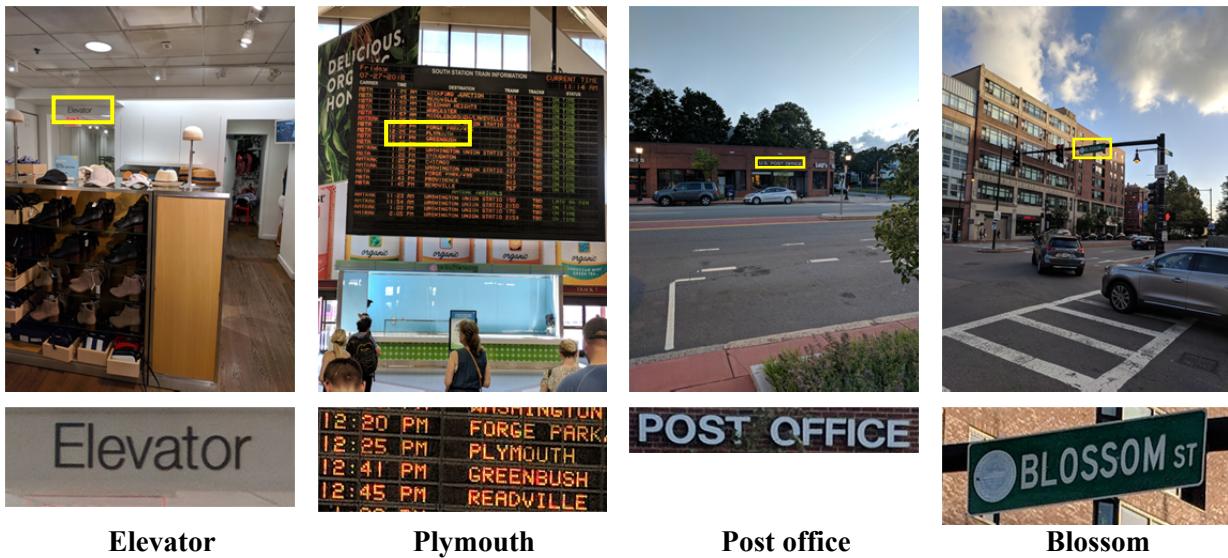
154 keyword in the image to the ground truth location was also determined. The distance (Euclidean) between

155 the detected word and the ground truth location was calculated and normalized with respect to the image

156 dimensions. For determining the localization accuracy, only successfully detected single keyword

157 instances were considered. Log-transformed keyword distances were compared between the 3 OCR

158 engines, excluding ABBYY RTR for which we did not have data.



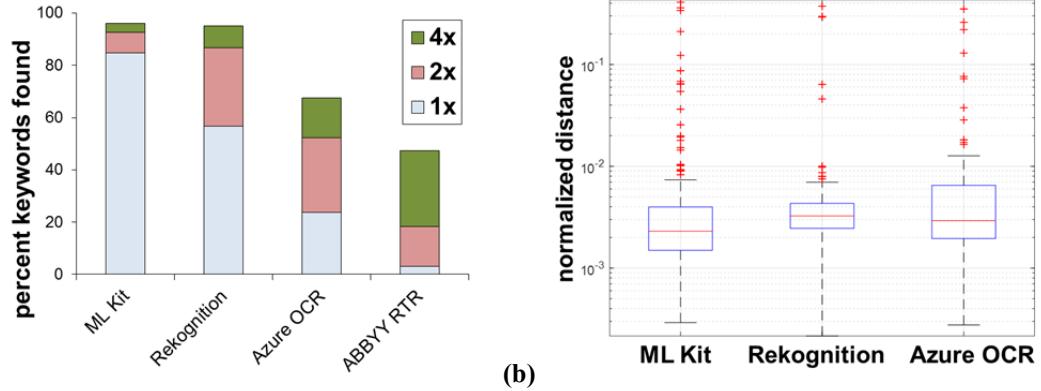
159 Figure 2: Examples of natural images (captured with a mobile device) used for testing the OCR engines.
 160 Actual images used for testing are shown in the top row and the search keyword associated with each
 161 image is listed in the bottom row. The location of the keyword in the images is highlighted in the images with
 162 an overlaid yellow rectangle (for the purpose of illustration, does not reflect OCR output). The highlighted
 163 patch is zoomed in the middle row below the images. The keywords for each image were chosen to represent
 164 realistic usage in way finding scenarios.

165 2.3 Results

166 Out of the 203 keywords tested, only 3 were not detected by any of the OCR engines. A 4-sample test for
167 equality of proportions showed that the overall successful search rate differed significantly between the

168 OCR engines ($\chi^2 = 187.5$, df = 3, p < 0.001), with ML Kit and Rekognition having a significantly higher
169 overall success rate than Azure OCR and ABBYY RTR (p < 0.001 for all multiple pairwise comparison
170 with Bonferroni correction). The proportion of successfully detected keywords was significantly higher
171 for Azure OCR compared to ABBY RTR (p < 0.001). There was no difference in the overall success rate
172 between ML Kit and Rekognition. Successful search rate at 1x zoom level differed significantly between
173 the 4 OCR engines ($\chi^2 = 325.16$, df = 3, p < 0.001), with ML Kit being the most successful in detecting
174 keywords without needing to zoom in, followed by Rekognition, Azure OCR, and ABBYY RTR. All
175 pairwise comparisons were significant (Bonferonni correction). While the overall successful detection
176 rate was not significantly different between ML Kit and Rekognition, ML Kit detected significantly more
177 keywords at the base zoom level (Figure 3a). The keyword search success rates at different image zoom
178 levels for the 4 OCR engines are shown in Figure 3a and summarized in Table 1.

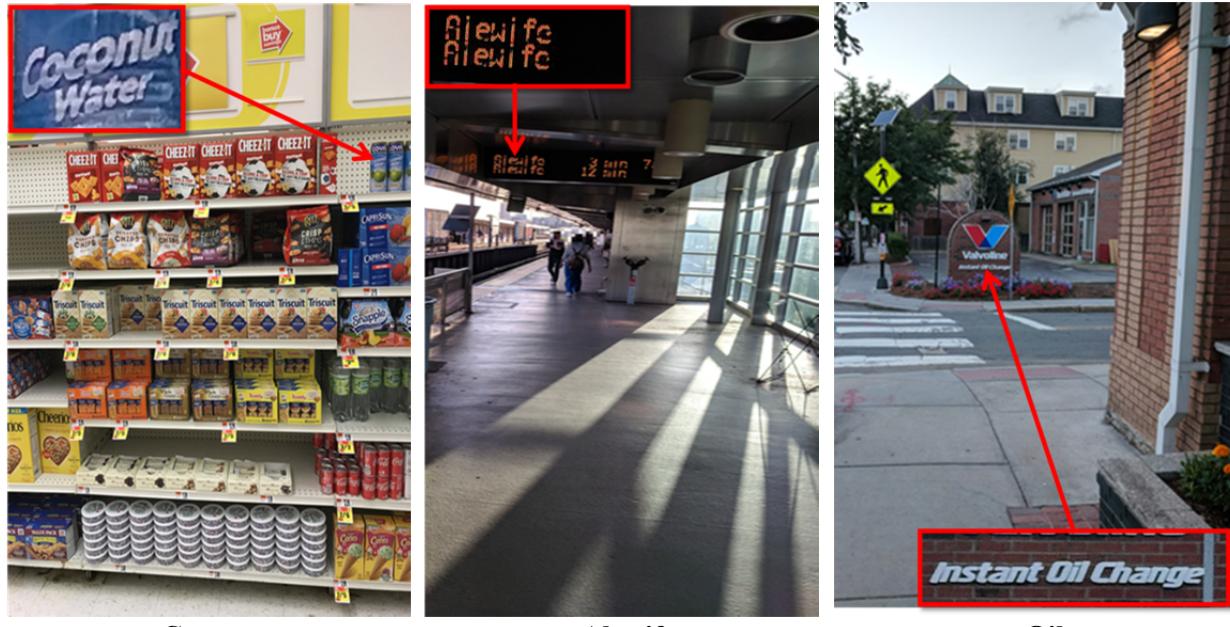
179 The keyword localization accuracy was evaluated by computing the distances between the keywords
180 correctly found in the image and the ground truth locations, based on 164, 163, and 113 cases out of 203
181 for ML Kit, Rekognition, and Azure OCR, respectively. The keyword localization errors were
182 significantly different when comparing the three engines (Kruskal-Wallis $\chi^2 = 10.54$, df = 2, p = 0.005),
183 with pairwise comparisons revealing a significantly larger localization error for ML Kit compared to
184 Rekognition (p = 0.003 with Bonferroni correction). There was no significant difference between
185 Rekognition and Azure OCR. Relative to the size of the images, the keyword localization errors were
186 between 0.2% and 0.4% of the image dimensions. Some examples of keyword search failures are shown
187 in Figure 4.



188
189
190
191
192
193
194
Figure 3: OCR testing results in natural images captured by a mobile device. (a) Comparison of the
percentage of keywords successfully found by the 4 OCR engines tested in our study (from a total of 203
keywords over 117 images) at different image zoom levels. Cumulative percentage of detected keywords for
ML Kit and Rekognition were significantly higher than others. ML Kit was also significantly better than the
other 3 in successful detection at base zoom level (1x). (b) Comparison of the average normalized distance to
the detected keyword in the image from the ground truth location (manually marked for each keyword).
Keyword localization error was significantly different between Rekognition and ML Kit.

195 Table 1: Successful keyword search rates at different zoom levels for the 4 OCR engines.

	ML Kit	Rekognition	Azure OCR	ABBYY RTR
1x	172 (84.7%)	115(56.6%)	48 (23.6%)	6 (3.0%)
2x	16 (7.9%)	61 (30%)	58 (28.6%)	31 (15.3%)
4x	7 (3.4%)	17 (8.4%)	31 (15.3%)	59 (29.1%)
Total	195 (96.1%)	193 (95.1%)	137 (67.5%)	96 (47.3%)



197
198
199
200
Figure 4: Some examples where the OCR engines failed to detect the keywords. The keywords associated with
the images are shown below the pictures. The location of each keyword within an image is indicated by a red
arrow, and the keyword region is shown in the inset for clarity. Font type variations, background clutter, and
image degradation are some of the main reasons for OCR failure in natural images.

201 **2.4 Discussion**

202 Our evaluation of the OCR engines with natural images shows that ML Kit performs the best in terms of
203 proportion of keywords successfully recognized, as well as the proportion of keywords recognized at the
204 base image scale without needing to zoom in. Not having to zoom-in for successful recognition can
205 potentially save processing time and maintain the maximum possible the field of view of the scene,
206 thereby improving the usability of the search paradigm in complex real world scenarios. Even though ML
207 Kit was found to be statistically significantly worse at localization of keywords but overall just slightly
208 (localization error of 1.4% of the average image size) compared to Rekognition(error of 1% of the
209 average image size), the differences are mainly due to outliers. One of the main reasons for this is OCR
210 mistakes, such as the merging of the keyword with neighboring word (for example, when searching for
211 keyword ‘pharmacy’, the localization is affected by the preceding text ‘CVS’. So while the keyword is
212 found, the localization error with respect to the ground truth location can be higher. While highly precise
213 localization of the detected keyword in the image is an important consideration, a small localization error
214 can be tolerated when searching for information associated with the keyword. As previously reported, the
215 0.2 to 0.4% localization error amounts to less than 20 pixel difference, which is negligible considering the
216 information associated with the keyword will still be in the neighborhood of the searched keyword.

217 One of the main reasons for search failure, particularly for OCR engines other than the ML Kit, was the
218 small size of the text in the image, which was resolved in most cases by increasing the scale. Other
219 reasons were image degradation and presence of non-standard fonts (such as too artistic, oriented in
220 different directions, or composed of dots as seen in the LED display on train station in Figure 4). Another
221 factor that could be responsible for search failure is OCR errors, where the recognized text from the
222 image contains the keyword but OCR mistake such as merging, typos, or special characters prevents its
223 successful identification.

224 There was a large disparity in the capabilities of the OCR engines to detect slanted text. ML Kit and
225 ABBYY RTR were limited to detection of text in a relatively narrow range of orientations ($\approx \pm 30^\circ$).

226 Azure OCR was able to detect text in all possible orientations due to its ability to switch the axes of the
227 images based on the device orientation. However, it assumed that all the text within the image had similar
228 orientations and thus could not handle multiple orientations of text within the same search image.
229 Rekognition was able to handle text in multiple orientations.

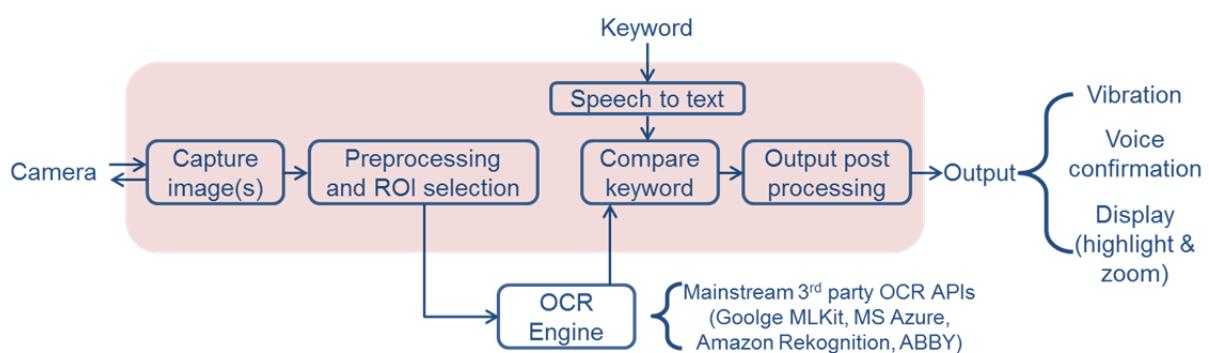
230 From the usage-standpoint, each OCR engine had its own strengths and limitations. Rekognition and
231 Azure OCR were cloud-based engines and thus required network connection. Given the lag in
232 transmitting the input image and retrieving the detected text within the image from the cloud, real-time
233 implementation with these OCR engine is not practically feasible. Moreover, the version of Rekognition
234 SDK tested in this study had an upper limit on the number of words it could detect in a given image.
235 Thus, despite being relatively successful in detecting almost all the keywords in our outdoor test images,
236 it could not handle situations where there was a lot of text in the image. The ABBY RTR engine could be
237 implemented natively (not requiring cloud-based processing), but the demo version we tried in this study
238 was clearly inferior to other OCR engines in scene text detection. Comparatively, ML Kit worked
239 natively on the device (did not require cloud-based processing), was capable of providing real-time text
240 detection output, and showed robust performance in challenging images, particularly detecting text
241 without needing to zoom-in. Based on the experimentation with the OCR engines, we chose ML Kit for
242 implementing the Android version of the SVS app with the capability for real-time scene text detection.

243 **3. The Supervision Search (SVS) Mobile Application**

244 **3.1 Description of SVS app**

245 The operational concept of the SVS app is shown in Figure 5. The scene is captured by the device's
246 camera and the image(s) are processed by the OCR engine to detect all the text present in the image(s).
247 The keyword input by the user is searched for within the detected text, and depending on the usage mode,
248 the user is informed via a combination of various methods: vibration of the device, voice confirmation,
249 and displaying the zoomed-in and highlighted instance of the keyword. The keyword can be input via
250 speech (standard voice input in the mobile device) or via keypad. The underlying OCR engine can be

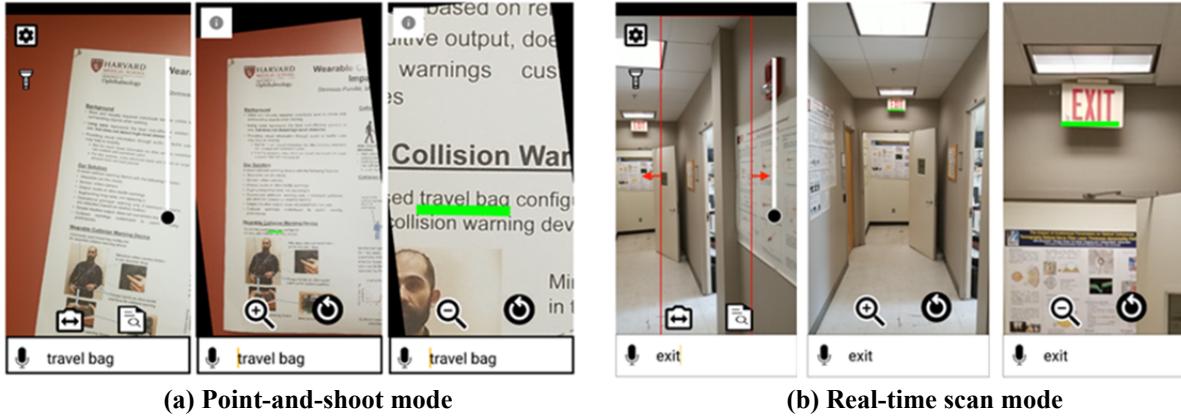
251 implemented using any of the mainstream APIs, including the 4 OCR engines discussed above. The
 252 application interface remains the same irrespective of the choice of the OCR engine, although the choice
 253 of OCR engine affects the performance and other operational parameters (such as availability of a
 254 network if the OCR API is cloud-based). Currently, the iOS version of the SVS app is based on Azure
 255 OCR. The Android version of the SVS app has been implemented using the ML Kit due to its native
 256 processing, real-time performance, and its overall robustness. The app functionality explained below
 257 pertains to the Android version.



259 **Figure 5: Overview of the operational steps in performing keyword search in natural images using a mobile
 260 device. The shaded region represents the block diagram of the Supervision Search app.**

261 There are two modes of operation for the SVS app: point-and-shoot mode and real-time scan mode
 262 (Figure 6). The point-and-shoot mode is ideal for searching a keyword within an object. On the other
 263 hand, the real-time scan mode is more suited for searching in a wide field of view, for example: localizing
 264 the exit door. In point-and-shoot mode, a single high resolution picture of the scene is captured with the
 265 rear camera of the mobile device and processed by the OCR engine to detect the presence of the keyword
 266 input by the user. In the real-time scan mode, the user presses and holds on the screen while moving the
 267 device in a scanning pattern over the target area (for example, holding the mobile device in front and
 268 scanning from side to side). Each video frame acquired by the camera is grabbed from the image buffer
 269 and processed by the OCR engine to detect the presence of the keyword in the scene in real-time. In order
 270 to improve the localization precision of the keyword, a scanning slice is introduced to restrict the field of
 271 view of the camera, thereby reducing the search area (searching only within a narrow strip of each

272 frame's image). This scanning slice was set as the central 50 percent region of the screen (horizontally).
 273 By reducing the region of interest for search, users can orient themselves more accurately with the
 274 detected target and reduce the negative impact of motion blur from the movement of the device when
 275 scanning.



276 **Figure 6: The two operational modes of the SVS app.** Screenshots of the input screen, the search notification
 277 screen, and the zoomed in location of the found keyword are shown for each mode. (a) In the point-and-shoot
 278 mode, the user takes a picture of a scene and the search result is highlighted and zoomed in for better
 279 viewing. In the real-time scan mode, there is a region of interest at the middle of the screen. When the
 280 keyword is found within this region of interest, then the app provides an indication. Thus the direction of
 281 camera pointing is loosely related to the spatial location of the keyword in the scene. This helps in wayfinding,
 282 for example- in this case, the location of the exit sign.

283 To minimize the impact of inevitable errors in OCR and possible typos in provided keywords, the SVS
 284 app does not require an exact match to keywords. Instead, the matching is based on a modified
 285 Levenshtein Distance.[30] When comparing the keyword with each word in the detected text, this
 286 algorithm determines the number of insertions, deletions, and substitutions that are required to change the
 287 detected word to the search keyword (Figure 7). Thus, with more edits, the probability that the candidate
 288 word is the search query diminishes. The net number of edits in each word is then divided by the length of
 289 the smaller of two words being compared to generate a distance value. This normalization allows the
 290 comparison algorithm to be dynamic — allowing for more edits in the detected word with a longer search
 291 query. Finally, the distance of each word in the detected text with the keyword is compared with a
 292 threshold. If the distance is lower than the tolerance, then the given word in the detected text is considered
 293 a match to the input keyword. This threshold value was set at 0.3 for this study. Such an approach for

294 string comparison is useful in dealing with OCR inaccuracies, as well as tolerating minor differences in
295 the input keyword and the actually present word (such as handling plurals or other minor variations or
296 typos made by the speech to text).

* I N S P I R A T I O N A L
| | | | | | | | | | | |
P E R S P I R A T O R Y * *
i s s s s s d d

297
298 **Figure 7: An example showing how the modified Levenshtein distance is calculated when comparing two**
299 **words. In this graphic, ‘i’ represents an insertion, ‘d’ represents a deletion and s represents a substitution.**
300 **Punctuation and plural words are ignored in the algorithm. Hyphenated words are regarded as multiple**
301 **separated words. The distance between inspirational and perspiratory is 8. Due to fact that perspiratory is a**
302 **shorter word — with 12 characters— the modified algorithm would return a value of 0.67. This number is**
303 **greater than the tolerance value in the app and thus these words are not similar.**

304 When the searched keyword is found, it is highlighted in the image by a flashing green bar to make it
305 more visible to the user. At the same time, voice output from the device indicates the number of instances
306 of the keyword that were found in the image. In addition, the app allows the user to directly zoom in on
307 the highlighted keyword to retrieve any necessary information related to it that is present in its immediate
308 neighborhood. In the case of multiple instances of the keyword being present, the app allows the user to
309 zoom in on each instance sequentially. The user also has the option to search for another keyword within
310 the captured image. The highlighting and zoom parameters are updated accordingly. In the real-time scan
311 mode, the app indicates the successful detection of the keyword in the image (in real-time) by vibration.
312 Thus, as the user is scanning with the mobile device, instant feedback is received when the keyword is
313 found. This helps in associating the presence of the keyword within the captured image with its actual
314 location in the scene via proprioception.

315 The native ML Kit OCR engine is restricted to detecting text with an orientation of $\pm 30^\circ$ with respect to
316 the horizontal. Moreover, the exact orientation of the detected text is not returned by this engine. Thus to
317 handle the rotated text (unidirectional rotation only), the average orientation of the entire text is computed

318 from the bounding box information for each word in the detected text. After computation of the average
319 orientation, the image rotation is cancelled before displaying it on the screen for the convenience of the
320 user.

321 **3.2 Evaluation of the SVS App**

322 Preliminary evaluation of the point-and-shoot mode of the SVS app was conducted with human subjects
323 to determine whether the concept of keyword search can provide any benefit in performing visual search
324 activities in daily life. We recruited 6 adults with normal vision (best corrected visual acuity 20/20 or
325 better and no other known vision disorders) from our institute for performing the app evaluation. During
326 the experiment, they wore blur glasses (blurring filter attached to a no-power lens) that reduced their
327 visual acuity to the level of 20/100 – 20/125. The study followed the tenets of the Declaration of Helsinki
328 and informed consent was obtained from all the study participants. The protocol was approved by the
329 institutional review board at the Massachusetts Eye and Ear Infirmary.

330 The evaluation task was a variation of the timed instrumental activities of daily living (TIADL) task.[31,
331 32] The TIADL tasks consist of a sample of visual activities routinely performed in daily life such as
332 reading ingredients of food products, instructions on medicines, finding phone number in a directory, and
333 using tools, among others. In our study, we curated tasks that required keyword search, increased the
334 overall number of tasks to be performed, and changed the complexities of the tasks such that they varied
335 from simple to complex (difficult). While task difficulty is a relative term, in this study difficulty refers to
336 the time required to complete the search (find the keyword). Briefly, these included finding specific
337 information from the items such as finding due date on a utility bill, or whether the food product
338 contained nuts, or finding lowest priced clothing item from a catalogue. The items were classified into 5
339 categories: printed sheets (flyers, utility bills etc.), documents (booklets, brochures etc.), restaurant
340 menus, food products, and other household items (Table 2). Overall, 50 items were selected and split
341 evenly in two conditions where the subjects searched for the information either using the SVS app or
342 using handheld optical magnifiers (4x and 12x). Pilot testing was done to make the overall grouping of

343 the items balanced in the two conditions (items in both groups were at about same difficulty level).

344 **Table 2: Modified TIADL tasks used for evaluation study. A total of 50 items were identified and were**
345 **divided equally in two sets to be tested in two conditions: while using the SVS app and while using optical**
346 **magnifiers.**

Item Category	Description	Typical Information Requested	Number
Printed sheets	- Utility bills - Bank Statements - Flyers - Forms (Tax etc.)	Due date, amount due, amount related to a particular transactions, reward balance, contact info on the flyer	15
Documents	- Telephone directory - Brochures - Catalogues - Manuals	Telephone number of a particular business or a person, price of an item in the catalogue, specifications listed in the manual	12
Restaurant Menus		Price of items containing a specific kind of food (for example, chicken, salmon etc.)	6
Food products		Nutritional information, ingredients	5
Other common items	- Electronics - Hardware - Stationary	Battery capacity, memory, photo resolution, setting time, voltage	12

347

348 The search task administration to the subjects occurred by presenting an item (each trial consisted of a
349 different item) and asking them for specific information contained within the item. The task question was
350 framed in a manner such that the intended search keyword was either part of the question or was implicit
351 in the question. However, they keyword was not directly provided to the subjects, and they were
352 supposed to come up with the keyword on their own when searching for the requested information. For
353 example, if the keyword to be searched in a credit card statement was reward points, the subjects were
354 asked to report what was the reward point balance for the month. The examiner would time the task
355 (using stopwatch) after presenting the item and communicating the keyword for it. The timing stopped
356 when the subject responded with the correct answer (the trial continued if the answer was incorrect). For
357 the SVS app condition, timing of the trial included keyword input time, and any subsequent errors due to
358 inputting an incorrect keyword or speech recognition errors by the voice input functionality of the
359 smartphone. When using optical magnifiers, the subjects were free to choose from 4x and 12x magnifiers
360 depending upon the text size in the item. SVS app was run on a Samsung Galaxy S8 smartphone. Subjects
361 were trained to use the SVS app and optical magnifier for searching with some practice items before the

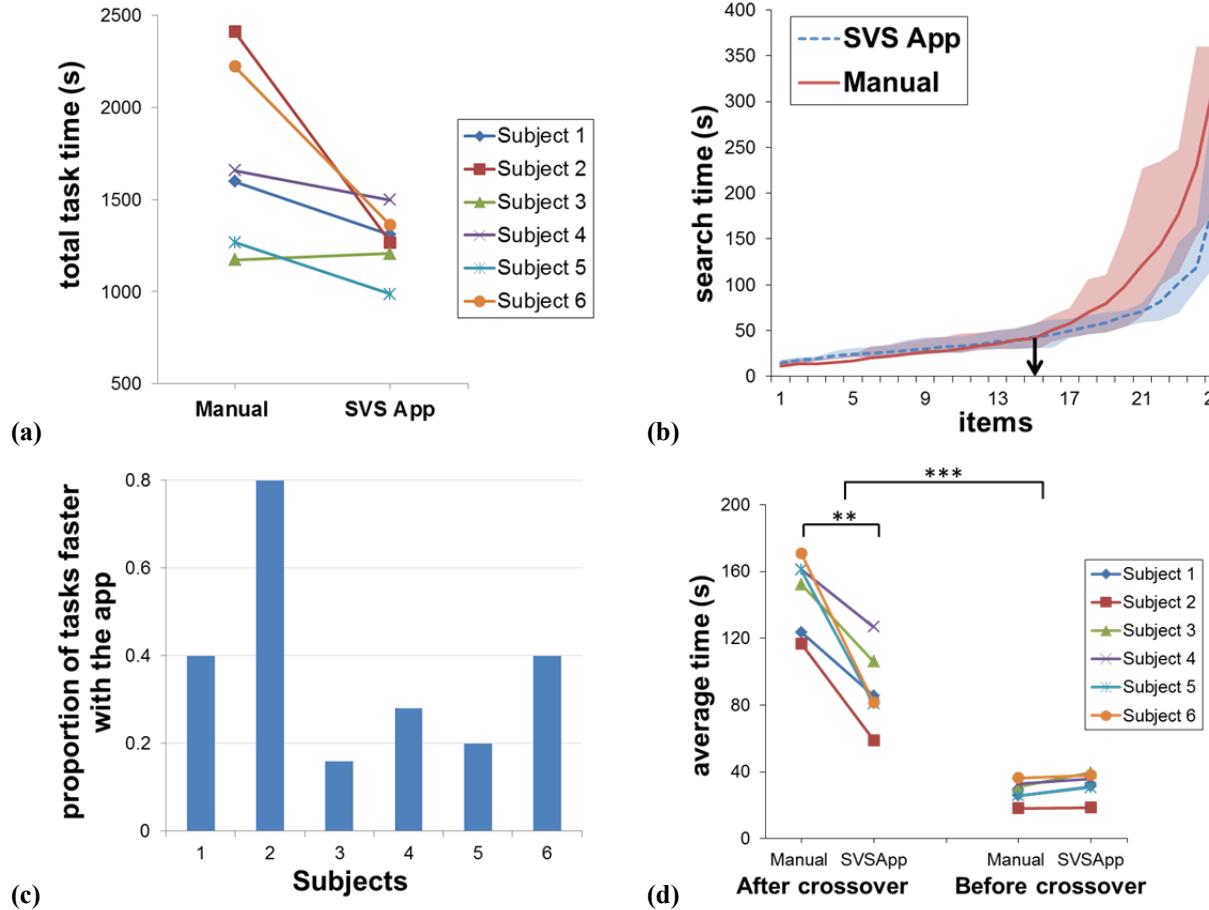
362 actual trials. The order of device use and task group was counter-balanced. The experiment was
363 conducted in a well lit room (standard office room lighting). The total study time for a subject was about
364 2.5 hours including task training (\approx 30 minutes) and a break of \approx 10 minutes between the two conditions.
365 An upper limit was placed on the search time for each trial (360 seconds) and if the subject did not find
366 the required information within this time the trial was ended and the search time for that trial was set as
367 360 s. For determining the benefit of the app in visual search, we compared the cumulative task time
368 between the two conditions. Cumulative task time was the sum of search times for individual trials in a
369 given condition. Since the individual tasks varied a lot in terms of their difficulty within the conditions,
370 we also performed further fine grained comparisons by considering task difficulties. The search times for
371 a subject were sorted from fastest to slowest for each condition. The average sorted task time was
372 compared between the two conditions. In the sorted task list, a crossover point was determined for each
373 subject after which the tasks were increasingly difficult to perform manually compared to the app. Tasks
374 sitting at a higher order in the list compared to the crossover point were therefore considered difficult to
375 perform manually using the magnifier. Comparisons of the median task time before and after crossover
376 point were done between the two conditions.

377 **3.3 Results**

378 Cumulative task completion times between SVS app and optical magnifier were not significantly
379 different, but the difference approached significance level of 0.95 (paired t-test: $t = 2.4$, $df = 5$, $p = 0.059$)
380 (Figure 8a). The cumulative task completion time reduced with the app for 5 out of 6 subjects (average \pm
381 std. for 6 subjects: manual method = 29 ± 8 minutes, with app = 21 ± 3 minutes). For the one remaining
382 subject, cumulative task time with the app was only about half a minute more than manual searching.
383 However, the cumulative task time does not fully inform us about the underlying data given the task
384 variety and the variable difficulty level of the tasks, which is evident with the skewed average sorted
385 search times (Figure 8b). It can be seen that search with magnifier was almost as fast as the SVS app for
386 the items before the crossover point for the two curves (easier tasks). However, over the remaining items

387 beyond the crossover point, the SVS app appeared to be substantially faster. The crossover point varied
388 between subjects (Figure 8c), indicating that the relative proportion of the search tasks for which the app
389 was faster compared to manual search varied between subjects. Averaged over 6 subjects, about 37%
390 tasks were faster with the app.

391 We then compared the average search times for tasks before and after crossover point for each subject
392 with and without the app using repeated measures ANOVA. There were significant effects of the search
393 method ($F(1,5) = 30.5$, $p = 0.003$), crossover point ($F(1,5) = 231.4$, $p < 0.001$), and their interaction
394 ($F(1,5) = 47.8$, $p < 0.001$) on the average search time. Post hoc tests with Bonferroni correction showed
395 that searching with app was significantly faster than manual searching for the more difficult tasks (tasks
396 beyond the crossover point) (average of 6 subjects - without app: 148 ± 22 seconds; with app: 90 ± 23
397 seconds, $p=0.009$). The difference in the average search times between the two methods before the
398 crossover point (easier tasks) was not statistically significant (without app: 28 ± 6 seconds; with app: 32 ± 8
399 seconds, $p=0.12$) (Figure 8d).



400
401 Figure 8: Evaluation of SVS app with normally sighted subjects wearing blur glasses (n=6). (a) The
402 cumulative task time (sum of individual trial times) reduced when using the SVS app in all but one subject
403 (#3). In some subjects (for example #2 and #6), there was a drastic reduction compared to manual search with
404 a magnifier. (b) Average sorted search time across subjects with and without the SVS app. The shaded region
405 shows the limits (max. & min.) of the response for the respective conditions. The plot shows that the search
406 tasks were of variable difficulty levels ranging from simple to difficult. Difficult tasks took longer to complete.
407 The crossover point (denoted by downward pointing arrow on the chart) can be determined from these
408 curves: the point after which search with the optical magnifier takes longer compared to the SVS app. (c)
409 Chart showing the proportion of the tasks that were faster with the SVS app for each subject. These
410 correspond to the tasks that sit at a higher order than the crossover point for the given subject. (d) Average
411 task time before (easy) and after (difficult) the crossover point between the two conditions. For the difficult
412 tasks, the median search time with the app decreased significantly. However, there was no statistically
413 significant difference in search time between the two apparatuses for easier tasks. ** denote p < 0.01, ***
414 denote p < 0.001.

414 3.4 Discussion

415 The evaluation results show that SVS app can lead to an overall saving in the search time, particularly for
416 complex search tasks that are deemed to be difficult. On an average, there were savings of ≈ 7.5 minutes
417 (26% reduction) with the app compared to search with the optical magnifier for cumulative task time. But
418 the cumulative search time does not tell the entire story since the specific tasks used in this study were

419 somewhat arbitrary. As Figure 8b indicates, the search times increase exponentially as the tasks get more
420 complex. This is true for both with and without app conditions. It can be expected that including more
421 difficult tasks would lead to more time saving, or vice versa. Therefore, the absolute time saving in
422 cumulative search time with the app is not very meaningful in the sense of generalization. Previous
423 studies using the TIADL methodology tested with a limited number of tasks, for instance 17 in Owsley et
424 al. 2001 [31], 5 in Taylor et al. 2014 [33], and 3 in Wittich et al. [34]. It is arguable that too few tasks may
425 not be able to capture the variety of situations encountered by people in their daily living.

426 In this evaluation study, we included 50 tasks with a wide range of difficulties, in an attempt to
427 generalize the study findings. A key consideration is not simply the number of tasks, but the difficulty
428 range. In Figure 8b, we can see that the search time with magnifier appeared to increase more rapidly
429 compared to the SVS app, as the difficulty level of the tasks increased. This finding will remain
430 unchanged, i.e. it is generalizable, no matter what tasks are included, as long as the range of task
431 difficulty is sufficiently broad. To perform more detailed analysis, we employed the method of separating
432 the tasks around the point where the sorted search time curves for the two apparatuses cross over (the
433 crossover point), essentially separating search tasks into two categories: easy vs. difficult. By separating
434 the tasks around the crossover point allowed us to understand the benefit of the app across broad range of
435 task difficulty as well as for different individuals. The results clearly indicate that for the easy tasks the
436 search time was almost the same for both the apparatuses, but for the difficult tasks the SVS app helped in
437 reducing the task time by a large margin (close to 40% reduction as compared to the optical magnifier).

438 The search time reduction with the SVS app is not seen across the board for all items because the app
439 requires some fixed amount of setup time, including inputting the search keyword and taking the picture.
440 Thus, the SVS app may not be time-saving when the tasks are easy (for example, when searching in less
441 cluttered objects). It should be noted that the search time for the SVS were inclusive of instances of
442 failure of keyword detection due to the following reasons: subjects inputting a keyword that was not
443 present in the item (or inputting an incorrect keyword), failure of OCR engine due to image capture issues

444 (such as blur or specular reflections off of the surface of the objects), failure of search due to incorrect
445 pointing of the camera (keyword out of the field of view of the camera), and incorrect speech recognition
446 (voice input returning incorrect keyword). On the contrary, searching with the optical magnifier did not
447 require any overhead time: the subjects could start the search the moment they heard the question for the
448 given task. The SVS app was able to reduce the overall search time in complex tasks despite these
449 limitations, errors, and inaccuracies.

450 The human subject evaluation was currently limited to the point-and-shoot mode of the SVS app. The
451 ability of the OCR engines to successfully search in images of outdoor scenes, as shown in this paper,
452 lends a lot of promise to the utility of real-time scan mode. Future work should evaluate the real-time scan
453 mode for the purpose of wayfinding. The human subject study sample was also limited to normally
454 sighted individuals with simulated visual acuity loss, limiting validity and generalizability. However, their
455 long search time in manual search condition indicates that the blur lenses somewhat replicated the
456 difficulty visually impaired people face. Evaluation of the app in visually impaired subjects will be future
457 work.

458 **4. General Discussion & Conclusions**

459 In this work we have shown that keyword search in real-world scenes can be implemented in a vision-
460 assistive mobile application, and it can help make keyword-based searches faster and easier in real-world
461 visual search tasks when visual acuity is lowered. Our specific keyword search offers a direct and realistic
462 chance of returning a successful hit, and consequently providing the users with the related information
463 without overloading them with irrelevant information. Use of the mobile platform is justified due to the
464 increasing popularity of mobile devices, even in the elderly population.[18, 21] Furthermore, vision aids
465 can be made cheaper and more accessible if delivered via mobile platform. In the support of
466 implementing keyword-based vision search as a vision aid, there is evidence that spot reading is highly
467 prevalent in people using mobile app-based smart vision aids.[35] Spot reading is required for quick
468 information gathering tasks and thus, a keyword searching application can help in spot reading and

469 information acquisition.

470 One limitation of this approach is that the user needs to input a keyword that is present in the scene (or a
471 close variation of it), which may not always be the case. This was observed in our human subject testing
472 of the SVS app, where the subjects entered keyword that was not present in the scene and thus they had to
473 retry with a different keyword. We did not explicitly provide the subjects with the keywords they were
474 supposed to search, and instead expected the subjects to come up with their own relevant keyword for the
475 given task. This study design was more realistic in simulating the daily life tasks where subjects are
476 generally aware of likely keywords but do not explicitly know which keyword they should search for.
477 Although we have designed some tolerance for typos introduced by the OCR engines or minor changes in
478 the keyword depending on the context (for example, plurals ending with an ‘s’), our approach cannot deal
479 with synonyms of a keyword. In the future, we will deal with the issue by introducing a smarter search
480 that can suggest context-based synonyms of the input keyword if a direct match is not found.

481 In conclusion, this work describes a novel vision aid implemented on a mobile device that can help make
482 visual search related daily life tasks easier to perform by visually impaired people. Impaired vision is a
483 complex condition and loss of visual acuity is only one aspect of it. Still, the preliminary results shown in
484 this paper support the approach of keyword searches in real-world scenes for aiding visual search. The
485 SVS app can be potentially beneficial to the people with low vision in complex search tasks, though
486 further investigation is required.

487 **Acknowledgments**

488 This work was supported in part by the Innovations in Technology Low Vision Research Award from
489 Research to Prevent Blindness.

490 **References**

491 [1] J. M. Wolfe, "Guided Search 2.0 A revised model of visual search," *Psychonomic Bulletin &*
492 *Review*, vol. 1, pp. 202-238, 1994.

- 493 [2] M. P. Eckstein, "Visual search: A retrospective," *Journal of Vision*, vol. 11, 2011.
- 494 [3] M. MacKeben and D. Fletcher, "Target search and identification performance in low vision
495 patients," *Investigative Ophthalmology and Vision Science*, vol. 52, pp. 7603-7609, 2011.
- 496 [4] N. Smith, D. Crabb, and D. Garway-Heath, "An exploratory study of visual search performance
497 in glaucoma," *Ophthalmic & Physiological Optics*, vol. 31, pp. 225-232, 2011.
- 498 [5] T. Kuyk, L. Liu, and P. Fuhr, "Feature search in persons with severe visual impairment," *Vision
499 Research*, vol. 42, pp. 3224-3234, 2005.
- 500 [6] P. Satgunam and G. Luo, "Does Central Vision Loss Impair Visual Search Performance of Adults
501 More than Children?," *Optometry and Vision Science*, vol. 95, pp. 443-451, 2018.
- 502 [7] G. Luo, P. Satgunam, and E. Peli, "Visual search performance of patients with vision impairment:
503 Effect of JPEG image enhancement," *Ophthalmic & Physiological Optics*, vol. 32, pp. 421-428,
504 2012.
- 505 [8] E. Brady, M. R. Morris, Y. Zhong, S. White, and J. P. Bigham, "Visual Challenges in the
506 Everyday Lives of Blind People," in *CHI 2013: Changing Perspectives, Session: Design for the
507 Blind*, 2013, pp. 2117-2126.
- 508 [9] S. West, G. Rubin, A. Broman, B. Muñoz, K. Bandeen-Roche, and K. Turano, "How Does Visual
509 Impairment Affect Performance on Tasks of Everyday Life?," *Arch Ophthalmol.*, vol. 120, pp.
510 774-780, 2002.
- 511 [10] A. Colenbrander, "Assessment of functional vision and its rehabilitation," *Acta Ophthalmologica*,
512 vol. 88, pp. 163-173, 2010.

- 513 [11] R. Massof, L. Ahmadian, L. Grover, J. Deremeik, J. Goldstein, C. Rainey, C. Epstein, and G.
514 Barnett, "The Activity Inventory: An Adaptive Visual Function Questionnaire," *Optometry and*
515 *Vision Science*, vol. 84, pp. 763-774, 2007.
- 516 [12] N. X. Nguyen, M. Weismann, and S. Trauzettel-Klosinski, "Improvement of reading speed after
517 providing of low vision aids in patients with age-related macular degeneration," *Acta*
518 *Ophthalmologica*, vol. 87, pp. 849–853, 2009.
- 519 [13] S. Smallfield, K. Clem, and A. Myers, "Occupational Therapy Interventions to Improve the
520 Reading Ability of Older Adults With Low Vision: A Systematic Review," *The American*
521 *Journal of Occupational Therapy*, vol. 67, pp. 288-295, 2013.
- 522 [14] S. Haymes, A. Johnston, and A. Heyes, "Relationship between vision impairment and ability to
523 perform activities of daily living," *Ophthalmic & Physiological Optics*, vol. 22, pp. 79-91, 2002.
- 524 [15] J. Szlyk, W. Seiple, G. Fishman, K. Alexander, S. Grover, and C. Mahler, "Perceived and actual
525 performance of daily tasks: relationship to visual function tests in individuals with retinitis
526 pigmentosa," *Ophthalmology*, vol. 180, pp. 66-75, 2001.
- 527 [16] Orcam. *MyEye 2* <https://www.orcam.com/en/myeye2/>.
- 528 [17] Aira. <https://aira.io/>.
- 529 [18] Pew Research Center. (2018 Mobile Fact Sheet <https://www.pewinternet.org/fact-sheet/mobile/>.
530 *Internet and Technology*.
- 531 [19] V. Braimah, J. Robinson, R. Chun, and W. M. Jay, "Usage of accessibility options for the
532 iPhone/iPad in a visually impaired population," presented at the The Association for Research in
533 Vision and Ophthalmology, 2014.

- 534 [20] M. D. Crossland, R. S. Silva, and A. F. Macedo, "Smartphone, tablet computer and e-reader use
535 by people with vision impairment," *Ophthalmic and Physiological Optics*, vol. 34, pp. 552-557,
536 2014.
- 537 [21] J. Morris, J. Mueller, M. L. Jones, and B. Lippincott, "Wireless Technology Use and Disability:
538 Results from a National Survey," *Journal on Technology and Persons with Disabilities*, I.
539 *Barnard et al. (Eds): Annual International Technology and Persons with Disabilities Conference*,
540 pp. 70-80, 2013.
- 541 [22] Microsoft. *Seeign AI* <https://www.microsoft.com/en-us/seeing-ai>.
- 542 [23] Google, "Lookout" <https://www.blog.google/outreach-initiatives/accessibility/lookout-discover-your-surroundings-help-ai/>."
- 544 [24] American Federation for the Blind. Read Printed Text with Your Smartphone
545 [https://www.afb.org/blindness-and-low-vision/using-technology/accessible-identification-
546 systems-people-who-are-blind-3. *Blindness and Low Vision*.](https://www.afb.org/blindness-and-low-vision/using-technology/accessible-identification-systems-people-who-are-blind-3)
- 547 [25] SuperVision Search.
548 <https://play.google.com/store/apps/details?id=edu.harvard.meei.supervisionsearch>.
- 549 [26] Google. *ML Kit* <https://firebase.google.com/docs/ml-kit/recognize-text>.
- 550 [27] Microsoft Azure. *Vision API (OCR)* [https://azure.microsoft.com/en-us/services/cognitive-
551 services/computer-vision/#text](https://azure.microsoft.com/en-us/services/cognitive-services/computer-vision/#text).
- 552 [28] ABBYY. *Mobile Capture* <https://www.abbyy.com/en-us/mobile-capture-sdk/>.
- 553 [29] Amazon. *Rekognition* <https://aws.amazon.com/rekognition/>.

554 [30] C. D. Manning, P. Raghavan, and H. Schütze, "Edit distance. Single-Link and Complete-Link
555 Clustering

556 " in *Introduction to Information Retrieval*, ed: Cambridge University Press, 2008.

557 [31] C. Owsley, G. McGwin Jr., M. Sloane, B. Stalvey, and J. Wells, "Timed instrumental activities of
558 daily living tasks: relationship to visual function in older adults," *Optometry and Vision Science*,
559 vol. 78, pp. 350-359, 2001.

560 [32] C. Owsley, M. Sloane, G. McGwin Jr., and K. Ball, "Timed Instrumental Activities of Daily
561 Living Tasks: Relationship to Cognitive Function and Everyday Performance Assessments in
562 Older Adults," *Gerontology*, vol. 48, pp. 254-265, 2002.

563 [33] J. Taylor, R. Bambrick, M. Dutton, R. Harper, B. Ryan, R. Tudor-Edwards, H. Waterman, C.
564 Whitaker, and C. Dickinson, "The p- EVES study design and methodology: a randomised
565 controlled trial to compare portable electronic vision enhancement systems (p- EVES) to optical
566 magnifiers for near vision activities in visual impairment," *Ophthalmic & Physiological Optics*,
567 vol. 34, pp. 558-572, 2014.

568 [34] W. Wittich, J. Jarry, E. Morrice, and A. Johnson, "Effectiveness of the Apple iPad as a Spot-
569 reading Magnifier," *Optometry and Vision Science*, vol. 95, pp. 704-710, 2018.

570 [35] G. Luo, "Use of Mobile Magnifier App - an Analytics Study," *Optometry and Vision Science*, vol.
571 In revision, 2019.

572

573