

Individual Differences in Face Looking Behavior Generalize from the Lab to the World

Matthew F. Peterson, Jing Lin, Ian Zaun, and Nancy Kanwisher
Massachusetts Institute of Technology
Department of Brain and Cognitive Sciences

Correspondence concerning this article should be addressed to Matthew F. Peterson, Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA 02139. Email: mfpeters@mit.edu

Abstract

1
2
3 Recent laboratory studies have found large, stable individual differences in the location people
4 first fixate when identifying faces, ranging from the brows to the mouth. Importantly, this
5 variation is strongly associated with differences in fixation-specific identification performance
6 such that an individual's recognition ability is maximized when looking at their preferred
7 location (Mehouder, Arizpe, Baker, & Yovel, 2014; Peterson & Eckstein, 2013). This finding
8 suggests that face representations are retinotopic and individuals enact gaze strategies that
9 optimize identification, yet the extent to which this behavior reflects real-world gaze behavior is
10 unknown. Here, we used mobile eye-trackers to test whether individual differences in face-gaze
11 generalize from lab to real-world vision. In-lab fixations were measured with a speeded face
12 identification task, while real-world behavior was measured as subjects freely walked around the
13 MIT campus. We found a strong correlation between the patterns of individual differences in
14 face-gaze in the laboratory and real-world settings. Our findings support the hypothesis that
15 individuals optimize real-world face identification by consistently fixating the same location and
16 thus strongly constraining the space of retinotopic input. The methods developed for this study
17 entailed collecting a large set of high-definition, wide field-of-view natural videos from head-
18 mounted cameras and the viewer's fixation position, allowing us to characterize subject's
19 actually-experienced real-world retinotopic images. These images enable us to ask how vision is
20 optimized not just for the statistics of the "natural images" found in web databases, but of the
21 truly natural, retinotopic images that have landed on actual human retinae during real-world
22 experience.

23
24 *Keywords:* mobile eye tracking, eye movements, face recognition, natural systems, retinal image
25 statistics

1 **Introduction**

2 The crux of the problem of visual recognition is the ability to appreciate that an object is
3 the same across the very different images it casts on the retina due to changes in position, size,
4 lighting, and viewing angle, to name a few (DiCarlo & Cox, 2007). Recent work suggests that
5 for the case of face recognition, position invariance is achieved in part by behavior rather than by
6 computation: people fixate a consistent and stereotyped position on the face, thus minimizing
7 variability in the retinal position of face images (Gurler, Doyle, Walker, Magnotti, &
8 Beauchamp, 2015; Mehoudar et al., 2014; Peterson & Eckstein, 2012). In particular, robust
9 individual differences are found in the precise location where people make their first saccade into
10 the face, with a continuous distribution ranging from the brows to the mouth. These differences
11 are robust over time, task, face familiarity, and variation in low-level properties such as color,
12 size, and contrast (Gurler et al., 2015; Mehoudar et al., 2014; Or, Peterson, & Eckstein, 2015;
13 Peterson & Eckstein, 2012, 2013). Most importantly, face recognition performance drops by
14 nearly 20% when faces are presented at another subject's preferred looking position if it differs
15 from one's own (Or et al., 2015; Peterson & Eckstein, 2013). This work suggests that the
16 representations that underlie face recognition are retinotopically specific, with position
17 invariance largely attained not by cortical computations (Riesenhuber & Poggio, 1999; Serre,
18 Wolf, Bileschi, Riesenhuber, & Poggio, 2007) but by looking behavior. However, all of this
19 work has been conducted in laboratory settings, with eye movements monitored as subjects
20 performed tightly controlled tasks in which photographs of faces are presented at a fixed distance
21 while head and body movements are restricted by a chinrest.

22 The lab-testing situation differs from real-world face viewing in a number of respects, yet
23 few studies have investigated real-world gaze on faces in non-clinical populations (Einhäuser et
24 al., 2009; Macdonald & Tatler, 2013, 2015). In the lab, visual stimulation is limited to a centrally
25 presented computer screen, whereas real-world faces generate a wide array of retinal images of
26 unpredictable sizes and positions anywhere in the visual field. In the world, unlike the lab, retinal
27 stimulation is determined not only by eye movements, but also by head direction and body
28 orientation. Further, real-world vision is dynamic and interactive, with goals shifting moment to
29 moment, rather than fixed by task instructions. Perhaps most importantly, in the real world the
30 face we are looking at is often looking back at us, engendering a social context associated with
31 tasks, signals, actions, and behavioral consequences that are distinct from the lab. Given the

1 dramatic differences between these conditions, it is important to know whether the consistent
2 individual differences in face-looking behavior documented in previous lab studies are also
3 found in everyday real-world vision. Here we asked this question by measuring each subject's
4 preferred face fixation position in the lab with the same methods used previously, and then by
5 sending them off for a walk around the MIT campus while wearing a mobile eye tracker. This
6 design enabled us to monitor where individuals fixated on faces that came into view during
7 naturalistic real-world vision. If position invariance for face recognition is indeed solved in large
8 part by looking behavior (rather than computation), then individual differences in preferred face-
9 fixation positions measured in lab should generalize to real-world behavior. Failure to find this
10 result would suggest that the prior results reflect a special case, and would cast doubt on the
11 hypothesis that position invariance in face recognition is solved by eye movements. A failure to
12 generalize would also call into question the extent to which face recognition behavior measured
13 in the lab should be applied to our understanding of how the brain processes faces during normal
14 operation.

15 Beyond answering whether face-fixation behavior observed in the lab generalizes to the
16 world, the present study will enable us to make a first foray into a broader research program of
17 characterizing what might be called "retinal image statistics". Most prior studies of natural image
18 statistics use photographs from the web that likely represent a biased sample of the images
19 people actually see in everyday life. First, these photos reflect situations in which someone used
20 a camera to select and frame a small portion of the visual world at a specific moment. The
21 criteria for the photographer's selection likely differ from the criteria viewers use to select
22 saccade targets. Second, most photographs are thrown away, and the ones that survive and get
23 posted on the web are a nonrandom sample, less likely to be marred by the occlusions, blur, bad
24 lighting, or other factors that reduce the intelligibility or attractiveness of the image but are
25 common in real-world contexts. Third, and perhaps most importantly, images on the web do not
26 come with information about where viewers were fixating. Fixation position matters enormously,
27 because acuity declines sharply from the fovea toward the periphery, meaning that only a few
28 degrees of the world around fixation are seen with high resolution. For all these reasons the
29 standard web-photo-based analyses of natural image statistics do not represent an unbiased
30 sample of the visual information that reaches the brain. Because our mobile eye tracking study
31 records both the image seen by the subject, and the subject's eye position on that image, our

- 1 study provides a collection of experienced images with the fixation point on each, a necessary
- 2 first step in a broader study of the statistics of experienced natural retinal images.

1 Methods

2 The study was run in two Stages. In Stage I, participants identified celebrity faces
3 presented on a computer screen while their eye movements were monitored. Each subject was
4 categorized into one of three groups according to where they tended to fixate on the faces. A
5 subset of these subjects from each group were later recalled to participate in Stage II, in which
6 they wore a mobile eye tracker to monitor their gaze while they walked around natural
7 environments.

8

9

1 **Methods (Stage I: In Lab)**

2 *Participants*

3 70 participants were recruited using flyers and departmental subject lists (40 MIT
4 students and 30 from the Cambridge community; 48 female; age: mean=28.0, min=18, max=62).
5 Subjects received \$20 for participation, gave informed consent, and had normal or corrected to
6 normal vision. The study was approved by the MIT Committee on the Use of Humans as
7 Experimental Subjects.

8 9 *Eye tracking*

10 The right eye of each participant was tracked using an SR Research EyeLink 1000
11 Desktop Mount sampling at 1,000 Hz. A nine-point calibration and validation were run at the
12 beginning of the session and after every 40 trials with a mean error of no more than 0.5° visual
13 angle. Saccades were classified as events where eye velocity was greater than 22°/sec and eye
14 acceleration exceeded 4000°/sec².

15 16 *Stimuli and display*

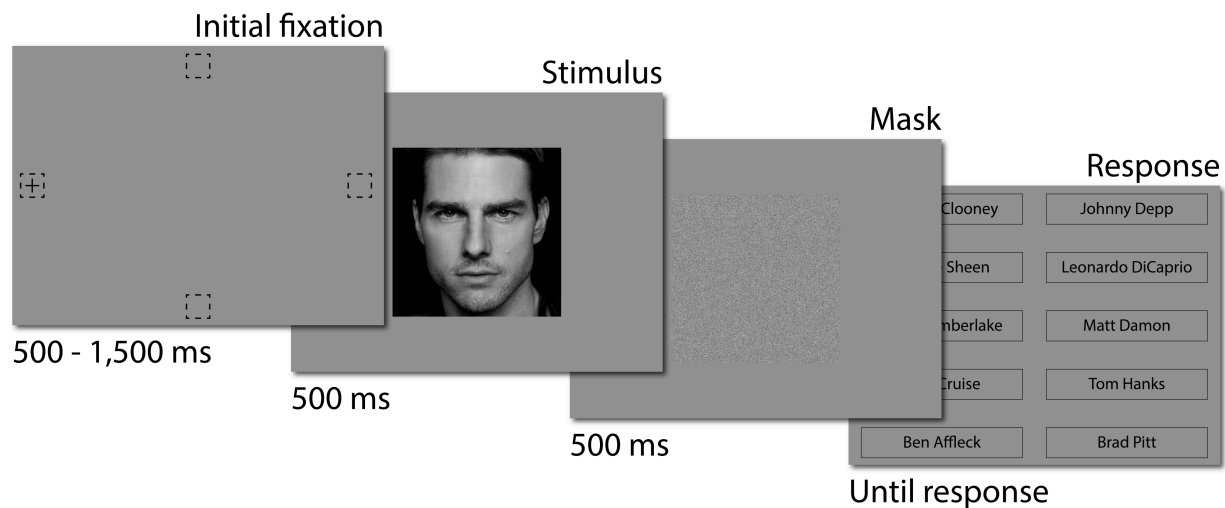
17 Stimuli were 160 frontal view images of 80 well-known Caucasian celebrities (e.g., Tom
18 Cruise, Jennifer Lawrence) acquired using Google image search (two different images per
19 celebrity, 40 male and 40 female). Images were converted to grayscale, rotated to an upright
20 orientation, scaled so that the center of the eyes and center of the mouth were in the same
21 position and separated by 6.0°, cropped from the top of the head to the chin (463 by 463 pixels or
22 16.9°) and contrast energy normalized. All stimuli were presented on a 17-inch CRT monitor
23 with a resolution of 1024 by 768 pixels and refresh rate of 85 Hz. Subjects sat 50 cm from the
24 monitor, with each pixel subtending 0.036°.

25 26 *Procedure*

27 Participants saw each of the 160 images in random order. Following the procedure used
28 in our earlier studies (Peterson & Eckstein, 2012, 2013), a trial began with a fixation cross
29 located 10° from the center of the monitor at either the left, right, top, or bottom edge of the
30 screen (location randomly selected). The subject fixated the center of the cross and pressed the
31 spacebar when ready. After a random, uniformly distributed delay between 500 and 1500 ms, the

1 cross disappeared and the randomly sampled face image was displayed at the center of the
 2 monitor. Note that in an earlier control experiment we found that the pattern of individual
 3 differences in preferred fixation behavior on centrally presented faces were conserved when
 4 faces were presented at unpredictable locations (Peterson & Eckstein, 2013). During the delay
 5 period the subject was required to maintain fixation at the cross, with a deviation of more than
 6 1.0° resulting in an error message and restarting of the trial. The face image remained visible for
 7 500 ms, during which eye movements were allowed, and was then replaced with a 500 ms high
 8 contrast white noise mask. A response screen then appeared consisting of two columns of five
 9 names each (the correct name of the face they had just seen and nine randomly sampled foils of
 10 the same gender, positions randomized). The subject used the mouse to click on the name they
 11 thought was correct after which the correct answer was highlighted for 500 ms before
 12 commencing the next trial (Figure 1).

13
 14



15
 16 Figure 1. In-lab famous face identification paradigm (Stage I).

17
 18

19 *Analysis*

20 Identification performance was quantified as the proportion of trials with a correct
 21 identification (*PC*). Individual's face fixation behavior was quantified by computing the mean
 22 location of the first into-face fixation (i.e., the location at the end of the first into-image saccade

1 as defined above Methods: Eye tracking) across the 160 image presentations. We then defined an
2 individual's Relative Fixation metric, γ , as the distance of their mean fixation upwards from the
3 mouth relative to the total distance between the mouth and eyes:

4
$$\gamma = \frac{y_{fixation} - y_{mouth}}{y_{eyes} - y_{mouth}} \quad \text{Eqn. 1}$$

1 **Methods (Stage II: Real World)**

2 *Participants*

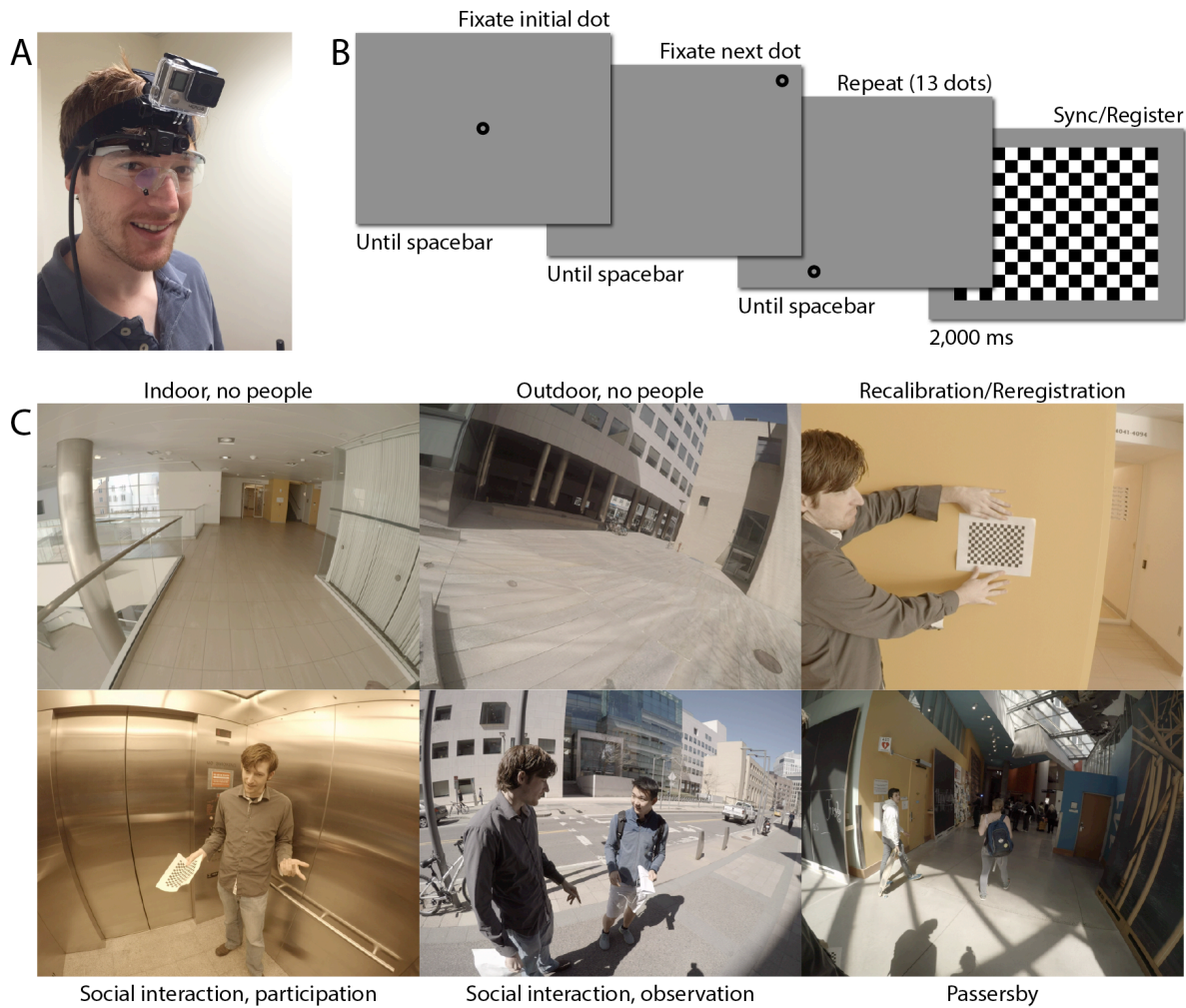
3 “Looking groups” were defined before the current study based on independent data from
4 250 subjects who had participated in similar face identification studies at the University of
5 California, Santa Barbara (Or et al., 2015; Peterson & Eckstein, 2012, 2013, 2014). As in the
6 current study, the previous work measured the mean location of subjects’ first into-face fixation.
7 Inter-individual variation was found to be large and consistent along the vertical dimension,
8 ranging from the eyebrows to the mouth (Gurler et al., 2015; Mehoudar et al., 2014; Or et al.,
9 2015; Peterson & Eckstein, 2013). Using these data, we defined criteria to categorize people into
10 three looking groups: Upper Lookers (UL) were the 15% of the sample who looked highest up
11 on the face, Lower Lookers (LL) were the 15% who looked lowest, and Middle Lookers (ML)
12 were everybody in between. We used these predefined criteria to categorize the original 70
13 subjects from Stage I of the current study into looking groups based on the average location of
14 their first into-face fixation from Stage I. The current sample yielded 11 ULs (15.7%), 45 MLs
15 (64.3%), and 14 LLs (20.0%). For each looking group, we recalled the 10 subjects with the
16 highest calibration scores, as measured by the EyeLink, to participate in Stage II (10 ULs, 10
17 MLs, and 10 LLs; Figure 5). As with Stage I, subjects received \$20 for participation, provided
18 informed consent, and had normal or corrected to normal vision. The study was approved by the
19 MIT Committee on the Use of Humans as Experimental Subjects.

20

21 *Procedure*

22 Subjects were told only that we were interested in assessing everyday, natural visual
23 experience. Critically, we did not mention any specific interest in faces or people. Subjects were
24 first fitted with the mobile eye tracker glasses and GoPro camera (Figure 2A) before initial
25 calibration, validation, and registration (see below and Figure 2B). The experimenter then
26 accompanied the subject for 8-12 minutes around the lab and nearby hallways of the Brain and
27 Cognitive Sciences Building and the Stata Center across the street, engaging in conversation
28 aimed toward making them feel comfortable with the apparatus. Subjects were then instructed to
29 walk unaccompanied across campus walkways, courtyards, a long hallway, and a busy city street
30 to a pre-designated location (12-15 minutes). The experimenter met the subjects at the location
31 and accompanied them back to the calibration room (5 minutes), concluding the study (25-30

1 minutes total). Each subject followed a similar path that exposed them to a representative sample
 2 of environmental settings (indoor locations like hallways, rooms, corridors, etc., and outdoor
 3 locations like streets, yards, etc.) and social contexts (no people, engaged in one-on-one
 4 interaction, watching others interact, etc.; Figure 2C). Subjects were all run at a similar time of
 5 day to maximize the between-subject consistency of environmental and social conditions.
 6
 7



8
 9 Figure 2. Real world eye tracking paradigm (Stage II). (A) Subjects were fitted with a pair of Applied
 10 Science Laboratory (ASL) eye tracking glasses. A supplemental GoPro camera enhanced the quality and
 11 field of view of the recorded video of the subject’s visual environment. (B) Calibration (moving dot) and
 12 ASL-to-GoPro video synchronization and registration (checkerboard) were automated and standardized
 13 across participants. (C) Each subject walked a similar route through the uncontrolled environments

1 around the MIT campus. Routes and times were chosen to ensure that a variety of locations and social
2 settings were sampled.

3
4

5 *Real world eye tracking: Overview*

6 Measuring and analyzing eye movements in unconstrained real world environments poses
7 multiple challenges. Here, we detail a standardized framework that allows the experimenter to
8 reliably collect and analyze accurate data. The framework focuses on standardized routines that
9 maximize the consistency, precision, and retention of data, while avoiding possible subject-
10 specific and task-specific biases. It also allows for frequent validation across time, a critical
11 aspect as data from mobile eye trackers can be marred by subject/apparatus motion and changing
12 environmental (e.g., lighting) and eye (e.g., pupil size) states that can dramatically compromise
13 initial calibration. Finally, the framework develops a combination of automatic algorithms and
14 novel crowdsourcing techniques for analysis and interpretation.

15
16

16 *Apparatus*

17 Real world gaze direction was measured at 60 samples per second with a pair of Applied
18 Science Laboratory (ASL) Mobile Eye-XG Eye Tracking Glasses. The ASL tracker uses two
19 cameras to estimate fixation position relative to the central region of the visual world in front of
20 the wearer (Figure 2A). The first camera, termed the scene camera, rests on the top rim of the
21 glasses and records video at 60 frames per second (fps), with a field of view (FOV) spanning 64°
22 horizontally and 48° vertically (640 by 480 pixels). The scene camera was adjusted to align the
23 center of its FOV with that of the subject's. The second camera, termed the eye camera, records
24 an infrared (IR) image of the subject's right eye reflected off a partially IR-reflective coated lens
25 that protrudes from the main lens. This allows the eye camera to detect both the subject's pupil
26 and the corneal reflection of a pattern of three dots produced by an IR emitter (with one dot
27 selected as the primary). The position and orientation of both the eye camera and the IR-
28 reflective lens were adjusted for each subject so that the pupil was centered in the eye camera's
29 FOV and the three IR dots were near the pupil center when the subject looked straight ahead.
30 The eye camera lens was then focused to maximize pupil and IR dot sharpness.

1 To improve upon the scene camera's FOV, resolution, and image sensor quality (contrast
2 sensitivity, temporal properties, etc.), subjects wore a supplementary GoPro Hero4 Black camera
3 (FOV spanning 110° horizontally and 90° vertically; 2704 by 2028 pixels; 30 fps). The GoPro
4 was positioned just above the eye tracker glasses and adjusted so that its FOV center aligned
5 with that of the ASL's (Figure 2A). A substantial fisheye distortion was present at the extreme
6 edges of the GoPro FOV. However, the fixations analyzed in the study were mainly restricted to
7 the central region where distortion was minimized.

8 9 *Calibration*

10 The ASL estimates gaze position by learning the mapping between specific locations in
11 the world (in x-y coordinates relative to the scene camera) and the displacement vector from the
12 pupil center to the primary IR dot registered by the eye camera. To minimize head movements
13 during calibration, subjects placed their heads on a chin rest located 42 cm from an 18" CRT
14 monitor centered in the subject's FOV with a resolution of 1024 by 768 pixels (spanning 50°
15 horizontally and 37.5° vertically). To maximize calibration accuracy and reliability, subjects
16 completed a standardized calibration task written in MATLAB and PsychToolbox 3.0.10 (a
17 JavaScript version was also developed, see Supplementary Material). Subjects first fixated on a
18 centrally presented black dot (outer radius 1.0°) with a small gray circular center (inner radius
19 0.15°). When the subject was confident they were fixating steadily as close to the dot center as
20 possible, they pressed the spacebar. The dot then relocated randomly to one of twelve positions
21 arranged in a 4 x 3 grid, spaced 14.0° apart horizontally and 15.8° vertically (spanning 42.0° by
22 31.6° ; Figure 2B). The subject would then fixate the new dot location and again press the
23 spacebar, proceeding through the 13 locations (12 grid plus initial central). After all dots were
24 fixated, an image of the entire array appeared, during which s/he was instructed to look at the
25 center of each dot, starting from the upper left and moving left to right and row by row for post
26 hoc validation.

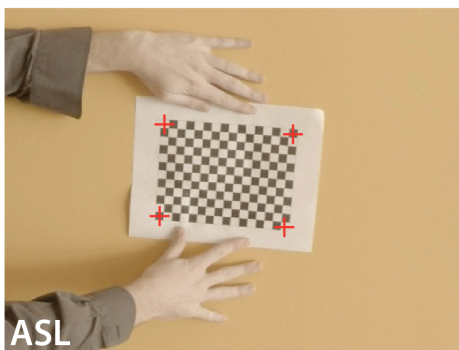
27 This data was used after the testing session for manual calibration using ASL's EyeXG
28 software. Independent raters viewed the scene camera video in slow motion (8 fps) with an
29 image of the pupil and displacement vector from the eye camera superimposed. For each
30 calibration dot transition event, the raters waited for the subject's eye to move and stabilize on
31 the new location as ascertained by an abrupt shift in the overlaid pupil/displacement vector. The

1 rater used a mouse to manually select the location of the center of the current calibration dot on
 2 the scene camera image (Figure 3A). The ASL EyeXG software then computed a function that
 3 mapped the displacement vectors (eye camera) to the dot locations (scene camera) for the 13
 4 calibration dots for each subject.

5
 6

Manual

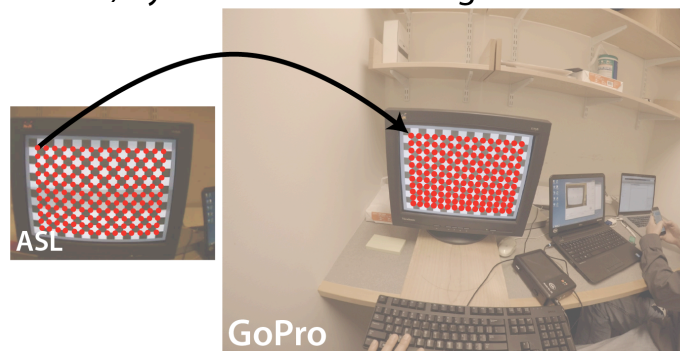
A) Calibration



C) Recalibration

Automatic

B) Synchronization & Registration



D) Fixation Detection & Remapping

7
 8 Figure 3. Post-processing of eye tracking and video data. (A) A subject-specific function that estimates
 9 gaze direction is learned by registering the location of each calibration dot (relative to the ASL scene
 10 camera) to the position of the pupil center and corneal reflection (from the eye camera). (B) The vertices
 11 of the post-calibration checkerboard pattern are automatically detected in both scene recordings, allowing
 12 for automatic synchronization and coordinate-registration between videos. (C) Data quality was validated
 13 every three minutes by having the subject fixate the corners of a checkerboard pattern. (D) Saccade and
 14 fixation events were automatically detected and their spatial coordinates mapped to the high-resolution,
 15 wide field of view GoPro video.

16

1

2 *Gaze location and fixation event detection*

3 Subjects' gaze location (in x - y coordinates) relative to the scene camera image for each
4 valid frame was estimated by the ASL EyeXG software using the mapping function learned
5 during calibration (Figure 3D). Frames were defined as invalid if the corneal reflection was lost
6 during saccades, blinks, large eccentricity fixations, or extreme external IR illumination and
7 were not included in the analysis. Across all subjects, $67.3 \pm 3.4\%$ (mean \pm standard error of the
8 mean) of frames were classified as valid, with no significant difference in the percentage of valid
9 frames between looking groups (ULs: $69.3 \pm 5.5\%$, MLs: $67.3 \pm 7.3\%$, LLs: $65.4 \pm 5.4\%$; $p = 0.91$).

10 Fixations were defined by the automated ASL algorithm as events where six or more
11 consecutive samples (100 ms) were measured within 1° of the sample group centroid. Fixation
12 events were terminated when three consecutive samples measured greater than 1° from the
13 fixation centroid or when pupil data was lost for 12 or more samples (200 ms; Figure 3D). To
14 check the accuracy of this automated algorithm, we re-analyzed the data using two well-
15 validated methods for categorizing fixation events in noisy eye tracking data with significant
16 flicker (intermittent loss of pupil contact; Holmqvist et al., 2011; Wass, Smith, & Johnson,
17 2012).

18 First, we re-analyzed all data following a modified version of the fixation-detection
19 algorithm for unreliable eye tracking data described in (Wass et al., 2012). The procedure was as
20 follows: 1) Samples labeled as missing data, or with out-of-range coordinates (x more than 32°
21 and/or y more than 24° from the scene camera center), were labeled as invalid; 2) Valid data was
22 smoothed with a bilateral filter to reduce small within-fixation jitter while preserving large
23 saccadic displacements (Durand & Dorsey, 2002; Frank, Vul, & Johnson, 2009; Stampe, 1993);
24 3) The mean absolute deviation (MAD) in gaze position was calculated within a 6 sample (100
25 ms) sliding window; 4) Windows with a MAD less than $50^\circ/\text{sec}$ were classified as potential
26 fixations, with consecutive qualifying windows concatenated into longer potential fixations; 5)
27 Potential fixations separated by less than 9 consecutive invalid samples (150 ms) were
28 concatenated if they were displaced by less than 1° , with invalid samples assigned the mean
29 position of the preceding potential fixation; 6) Potential fixations were labeled as valid fixations
30 if they were immediately preceded and followed by a likely saccade event (MAD of 3

1 preceding/succeeding samples greater than $100^\circ/\text{sec}$) and displaced from the mean of the
2 preceding and succeeding potential fixations by at least 1° .

3 Second, the authors hand-coded fixation events for a random sample of data through
4 visual inspection of gaze-position vs time plots (Holmqvist et al., 2011; Wass, Smith, & Johnson,
5 2012) and fixation-overlaid scene camera video (Figure 6). Fixations were defined as epochs of
6 relatively stable gaze position preceded and succeeded by abrupt shifts in gaze.

7 The two validation procedures were in good agreement with the automatic ASL
8 algorithm (see Supplementary Figure S1 for an example of fixations detected by each
9 procedure). Both of the validation procedures detected fewer, and longer, fixations than the ASL,
10 largely due to the merging of shorter fixations interrupted by brief false saccade events into
11 longer fixations (Supplementary Figure S1). Fixation position did not differ across procedures,
12 and, critically, the positions of on-face fixations were unaffected, with strong across subject
13 correlations between the ASL procedure's γ values and both the filtering procedure ($r = .95, p <$
14 $.01$; Supplementary Figure S2) and hand-coding ($r = .89, p < .01$).

15

16 *Synchronization and registration*

17 The ASL EyeXG software outputs an estimated gaze location for each frame in x - y
18 coordinates relative to the ASL native scene camera, but ultimately we wanted to map these
19 fixation coordinates to the higher resolution, larger FOV GoPro video. To do this, we presented a
20 16 by 12 checkerboard pattern on the monitor immediately after validation (Figure 3B). After the
21 fact, we implemented an automatic routine in Matlab that searched for the first frame in the
22 native scene camera video in which a 16 by 12 checkerboard pattern could be detected. The time
23 in the video was recorded and the coordinates of the checkerboard vertices (192 points)
24 automatically detected (Figure 3B). The same was done with the GoPro video. The video streams
25 were then synchronized by aligning the checkerboard onset times. Then, we computed the
26 projective linear transform matrix, T , that mapped the 192 vertex points from ASL to GoPro
27 coordinates with the minimum mean-square error. The transform matrix was then used to map
28 gaze coordinates for each frame and each fixation event from the ASL video to the GoPro
29 (Figure 3D).

30

31 *Recalibration and reregistration*

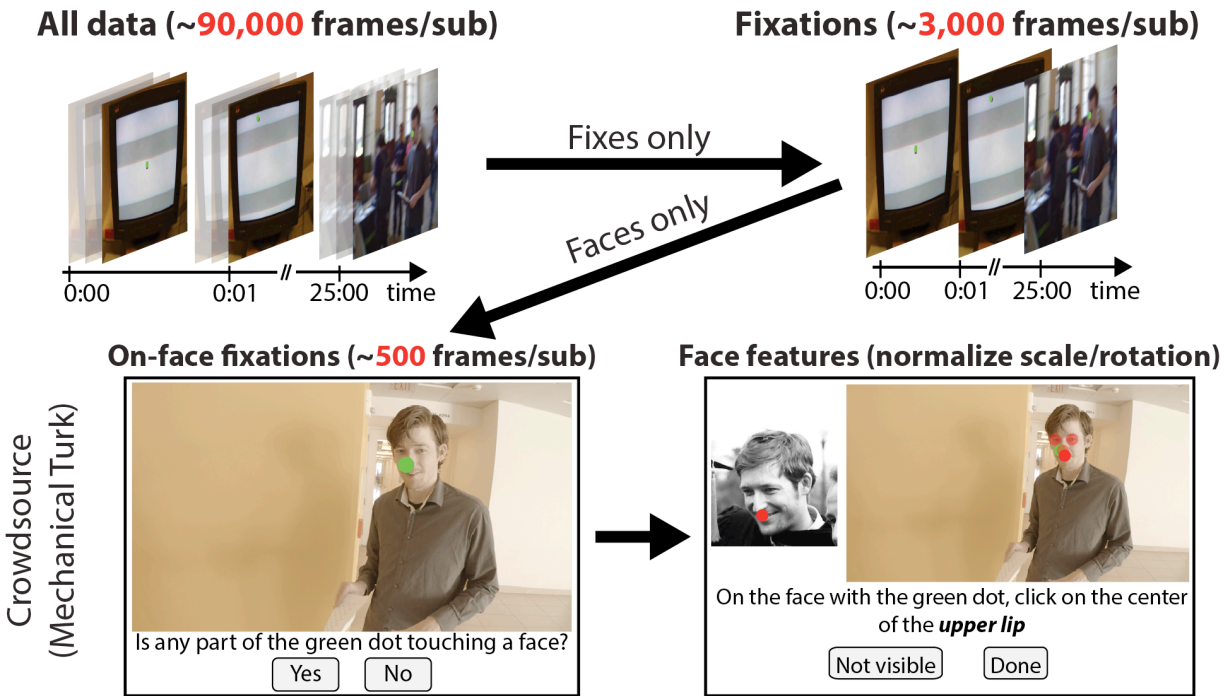
1 To ensure data validity over the course of the study, subjects regularly performed a
2 recalibration and reregistration routine. Every three minutes, the subject was instructed to stop
3 and hold at arm's distance a calibration/registration checkerboard pattern centered at eye level.
4 While keeping the head steady, the subject would fixate, in turn, the extreme upper-left, upper-
5 right, lower-left, and lower-right corners of the checkerboard for two seconds each before
6 resuming their walk (Figure 3C). Similar to the initial calibration, independent raters viewed
7 each recalibration at 8 frames per second. For each of the four corner fixations, the raters waited
8 until the subject's eye moved and stabilized on the new location indicated by a sudden shift and
9 stabilization of the overlaid pupil/displacement vector. The rater selected the location of the
10 center of the current recalibration target on the scene camera image (Figure 3C), which the ASL
11 EyeXG software used to augment the displacement vector to gaze location mapping function.
12 Similarly, the 16 by 12 checkerboard pattern and its corresponding vertices were automatically
13 detected in both videos and any necessary adjustments to the transform matrix were applied.

14

15 *Analysis: automatic fixation event filtering*

16 On average, we obtained 24.2 minutes (87,165 frames) of data per subject (Figure 4A).
17 For this study, we were interested only in the fixation location targeted by saccades. This
18 information is contained completely in the image and gaze position corresponding to the first
19 frame of each detected fixation event. This allowed us to greatly reduce our data set by
20 automatically selecting, for each fixation, a single video frame and eye position for further
21 analysis (average of 3,023 frames/subject; Figure 4B).

22



1
2 Figure 4. Analysis and interpretation of fixation data. The current study is concerned only with the
3 locations of distinct fixation events, greatly reducing the amount of data to be analyzed (from 60 to
4 around 2 samples/second). Since only on-face fixations were relevant here, data was further refined with
5 the help of human raters on Mechanical Turk. Finally, human raters were again enlisted to determine the
6 location of the on-face fixations relative to the eyes and mouth.

7
8 *Analysis: crowdsourcing face-fixation events*

9 One of the primary difficulties with studies conducted outside traditional laboratory
10 environments is the decreased ability to control subjects' sensory input. In the lab, the
11 experimenter precisely determines the spatial and temporal characteristics of visual stimulation.
12 Thus, the position (x, y) of gaze at some time (t) unambiguously maps to known stimulus
13 properties. Unconstrained environments do not provide this level of control, as the
14 spatiotemporal properties of the visual stimulus are not known a priori. This situation makes
15 measurements of gaze timing and position necessary but not sufficient for mapping to
16 meaningful stimulus properties. The difficulty of this mapping is determined by the stimulus
17 properties the experimenter is interested in, the quality of the visual recording, and the
18 complexity of the visual environment.

1 In this study we are interested in how people look at faces. This goal requires the ability
2 to reliably determine whether a fixation is on a face given only the recorded video image and the
3 associated x - y gaze position. While advances in algorithms and computing resources have led to
4 impressive gains in automatic face detection within complex images (Phillips & O’Toole, 2014;
5 Taigman, Yang, Ranzato, & Wolf, 2014), the combination of high resolution video and
6 unconstrained environmental uncertainty poses a serious challenge to even the most advanced
7 computer face detection systems. In this type of scenario, humans remain the gold standard for
8 face detection accuracy. However, this advantage comes at a cost of processing capacity: an
9 individual can accurately detect faces only up to a certain speed.

10 To maximize accuracy and throughput, we developed a simple crowdsourcing algorithm
11 using Amazon Mechanical Turk. By drawing on the judgments of many individuals in parallel,
12 crowdsourcing greatly increases the bandwidth of human-based face recognition. Turk raters
13 were shown a series of randomly sampled single video frames corresponding to fixation onsets
14 as described in the previous section. For each image (trial), a bright green dot was overlaid at the
15 measured fixation location, and the rater responded whether any portion of the green dot was
16 touching a face (Figure 4C). To ensure raters were real humans who understood and were
17 actively attending to the task, each image was rated by multiple people. If the first two raters
18 agreed, the response was taken as truth and the image was removed from the rating pool. If the
19 first two raters did not agree, the image was shown to a third tie-breaking rater. Individual raters’
20 performance was monitored by calculating their miss (responding No Face when two separate
21 raters responded Face) and false alarm rates (responding Face when two other raters responded
22 No Face). For online quality assurance, each trial had a 1 in 30 chance of being a probe. The
23 probe set was a mixture of 80 author-verified images and an expanding set of images that had
24 already been successfully rated by two other raters (who had not themselves been excluded
25 because of low concordance with other raters), with author-verified images more likely to be
26 sampled on earlier trials. If the rater disagreed with the consensus, they would be given a
27 warning message. Raters were allowed two mistakes; a third disqualified them from further
28 participation and all of their rating data was discarded from final analyses. Post hoc manual
29 verification by the authors of a random sample of rated images revealed no false positives or
30 negatives.

31

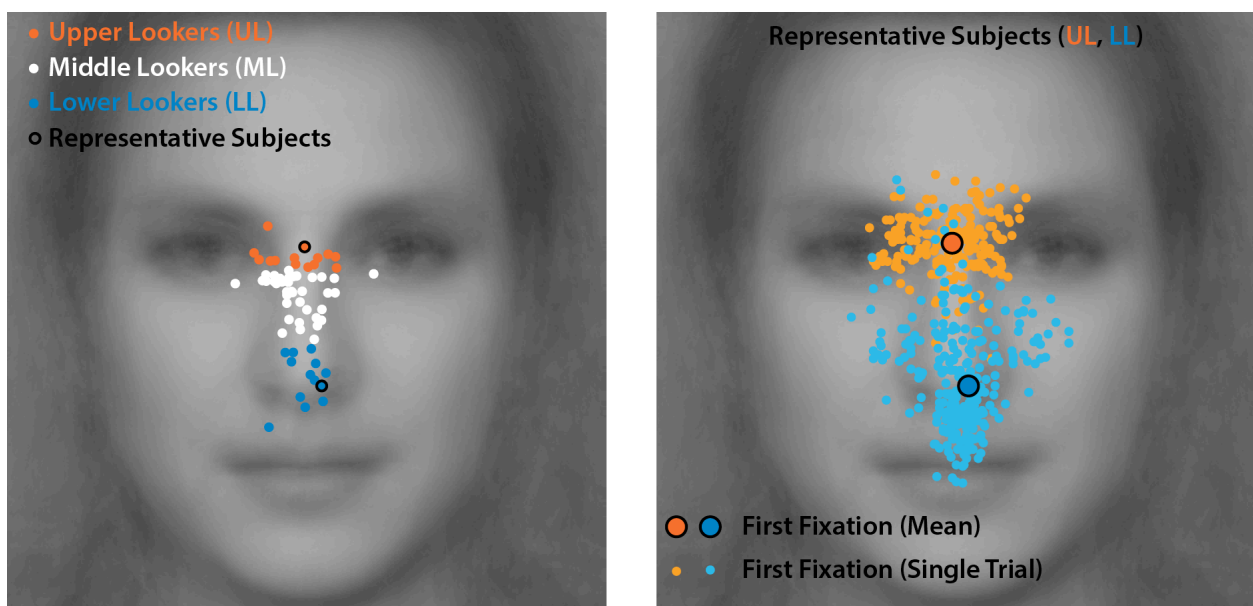
1 *Analysis: crowdsourcing face-fixation location*

2 To quantitatively compare within face fixation location between the laboratory and the
3 real world, we need to compute the Relative Fixation metric, γ (see Equation 1 in the Analysis
4 section of Methods for Stage I). In the lab, this calculation is simple, as the position of the eyes
5 and mouth are set and known by the experimenter. For the mobile section, we need to estimate
6 these locations on the video frames where faces could be present at any combination of location,
7 pose and size. We again turned to crowdsourcing with a second Mechanical Turk task. Raters
8 were shown random frames that were determined from the first Turk task to have on-face
9 fixations (again signified by a green dot). If the rater determined that the image was originally
10 misclassified as face-present in the first Turk task, a No Face option was available that recycled
11 the image back to the previous Turk task pool. Otherwise, raters were first asked to rotate the
12 image until the face with the dot on it was upright and then clicked on the center of one of the
13 visible eyes and the center of the upper lip (the upper lip was chosen so as to minimize the
14 variability in estimated mouth position due to plastic changes arising from talking, expressions,
15 etc.; Figure 4D). γ was then computed as before (Equation 1). Each image was scored by two
16 raters. If the raters disagreed by more than ten degrees of rotation and/or more than 10% of the
17 eye-to-mouth distance, a third rater scored the image and the two most similar ratings were
18 averaged. After the fact, manual verification of a random sample of rated images showed good
19 agreement by the raters and no systematic biases.

1 Results

2 *In lab initial face fixation behavior*

3 Across subjects, the initial into-face saccade landed on average below the eyes
 4 (mean±standard error of the mean: $\gamma = .757 \pm .025$, $t(69) = 9.86$, $p < .001$) and left of the midline
 5 ($\chi = .041 \pm .014$, $t(69) = 3.02$, $p = .0035$; Figure 5). Consistent with past literature, individuals
 6 varied greatly and consistently in their preferred face-fixation behavior along the vertical
 7 dimension, ranging from the eye brows (max(γ) = $1.11 \pm .061$) to just above the mouth (min(γ) =
 8 $0.17 \pm .065$; Figure 5; Gurler et al., 2015; Mehoudar et al., 2014; Peterson & Eckstein, 2013).
 9



10
 11 Figure 5. Stage I (in-lab) initial-fixation behavior for face identification. On the left, each dot represents
 12 the mean location, across trials, of the initial on-face fixation for one subject. Subjects were categorized
 13 as Upper (orange), Middle (white), or Lower (teal) Lookers according to pre-determined criteria based on
 14 previous work. On the right, fixations for each trial (small dots) and the mean across trials (large dots) for
 15 one UL (orange) and one LL (teal).
 16

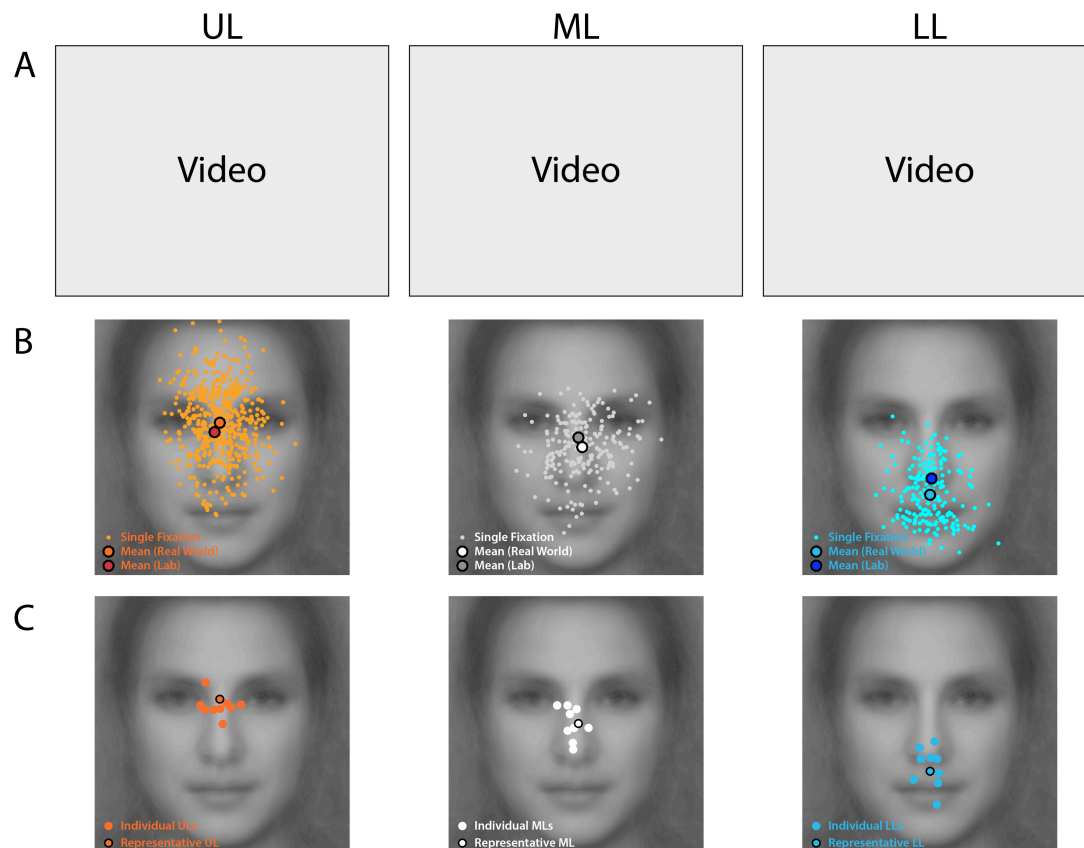
17 An existing independent sample of face-looking behavior ($n = 275$) was used to pre-
 18 define criteria to categorize the current subject sample into three groups. Upper Lookers (UL)
 19 fixate higher on the face than 85% of the total previously-sampled population ($\chi_{UL} = .93$), Lower
 20 Lookers (LL) fixate lower than 85% ($\chi_{LL} = .55$), with Middle Lookers (ML) constituting
 21 everybody else. Using this criteria, 11 of 70 subjects were categorized as ULs (15.7%), 14 as

1 LLs (20.0%), and 45 as MLs (64.3%; Figure 5). The 10 subjects with the best Stage I EyeLink
 2 calibration scores from each group were recalled for the mobile condition, resulting in the
 3 following gamma values (mean±standard deviation) for each group: ULs: $\gamma_{UL} = .995 \pm .098$, MLs:
 4 $\gamma_{ML} = .815 \pm .133$, LLs: $\gamma_{LL} = .326 \pm .183$ (Figure 7A).

6 *Real world face fixation behavior*

7 Subjects' distinctive preferred face fixation behavior can be appreciated in the example
 8 subject videos from each looking group (Figure 6A): Individuals fixated predominantly at their
 9 preferred region, with occasional fixations on other face regions quickly followed by a return to
 10 the preferred region. Most importantly, individuals' preferred real world fixation regions were
 11 consistent with their laboratory fixations (Figure 6B). Grouping subjects according to their in-lab
 12 behavior, the data from real-world viewing showed that ULs ($\gamma_{UL} = .921 \pm .040$) looked
 13 significantly higher than MLs ($\gamma_{ML} = .735 \pm .056$, $t(18) = 2.91$, $p = .005$) who looked significantly
 14 higher than LLs ($\gamma_{LL} = .267 \pm .066$, $t(18) = 5.69$, $p < .001$; Figures 6C, 7A).

15



16

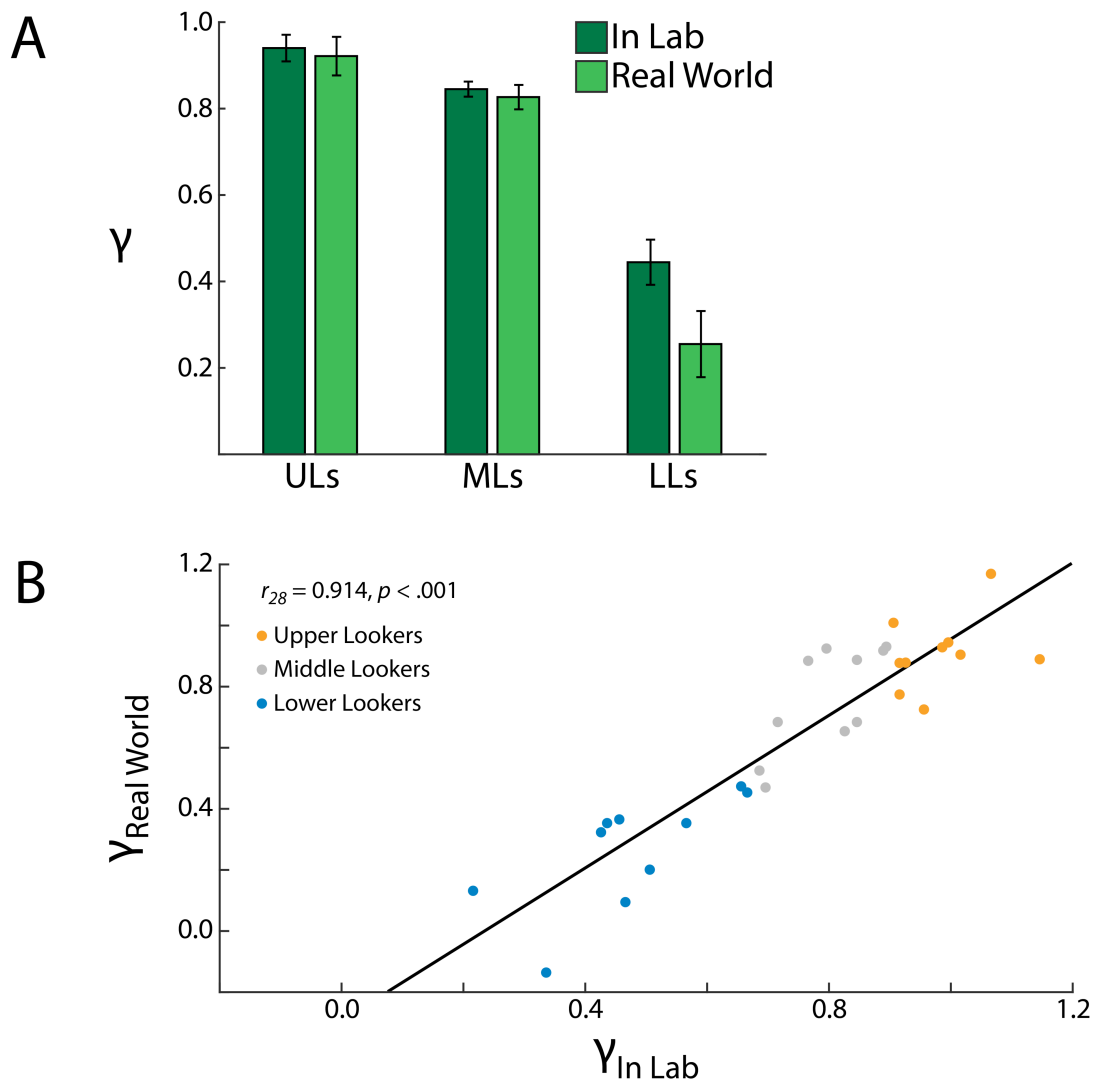
1 Figure 6. Real world face fixation behavior for lab-defined Upper, Middle and Lower Lookers. (A) Each
 2 video is from one representative subject from each group, with the white dot denoting gaze position. (B)
 3 Dots represent individual on-face fixation events for the same subjects as (A). (C) Each dot represents the
 4 mean location across all on-face fixations for a single subject.

5

6 *Relationship between in lab and real world face fixation behavior*

7 A repeated measures two-way ANOVA found significant main effects of looking group
 8 ($F(2,27) = 65.45, p < .001$) and modality (laboratory vs. real world; $F(2,27) = 9.62, p = .004$) on
 9 fixation behavior (γ), but not a significant interaction ($F(2,27) = 0.28, p = .76$; Figure 7A).

10



11

1 Figure 7. Relationship between real-world and in-lab face fixation behavior. (A) The laboratory-measured
2 group differences in the mean location of the initial on-face fixation, from 0 (center of the mouth) to 1
3 (center of the eyes), are also observed under real-world conditions. (B) The conservation of face-gaze
4 patterns between the lab and the world is consistent at the individual level across the range of observed
5 behavior.

6
7 Across the sample, correlational analysis showed that an individual's real-world fixations
8 were strongly predictive of their laboratory behavior ($r(28) = .914, p < .001$; Figure 7B). This
9 relationship was near ceiling given the reliability of the each modality's measurements. For each
10 of 1,000 bootstrap samples, we randomly split each subject's data in half, computed γ for each
11 half, and calculated the correlation between the two halves. The average split-half reliabilities
12 were $r = .996$ and $r = .909$ for the in-lab and real-world measurements, respectively, with an
13 average split-half correlation of $r = .905$ between them (correlation value lower than for the full
14 data set due to smaller sample sizes).

1 **Discussion**

2 Here we tested whether individual differences in face-looking behavior, observed
3 previously only in restricted lab conditions, generalize to the real world. To answer this question,
4 we measured subjects' fixation positions on faces both under controlled laboratory conditions
5 and while they walked around the MIT campus. Our main finding is that face-fixation patterns
6 are remarkably similar in the two situations, with an individual's laboratory fixation behavior
7 strongly predicting their real-world gaze, nearly as well as possible given measurement
8 reliability (Fig. 7). These results demonstrate that the prior lab-based finding of individual
9 differences in face fixation behavior generalizes to real-world vision. They further imply that the
10 superior face recognition performance when an individual fixates their preferred location
11 (Peterson & Eckstein, 2013) both reflects, and optimizes, that person's real-world face
12 recognition behavior. Taken together, these results suggest that real-world face recognition
13 entails two qualitatively distinct stages: face detection in the periphery, and face recognition at
14 the fovea. Finally, the methods developed here provide a rich dataset of images that humans
15 have actually experienced during real-world viewing, including the viewer's fixation position on
16 each image, opening up important new avenues for investigation of the statistics of the images
17 landing on peoples' retinas during natural behavior (Retinal Image Statistics; RIS), and the
18 tuning of human behavior and neural representations to those statistics.

19

20 *Comparing Real-World and In-Lab Eye Movements*

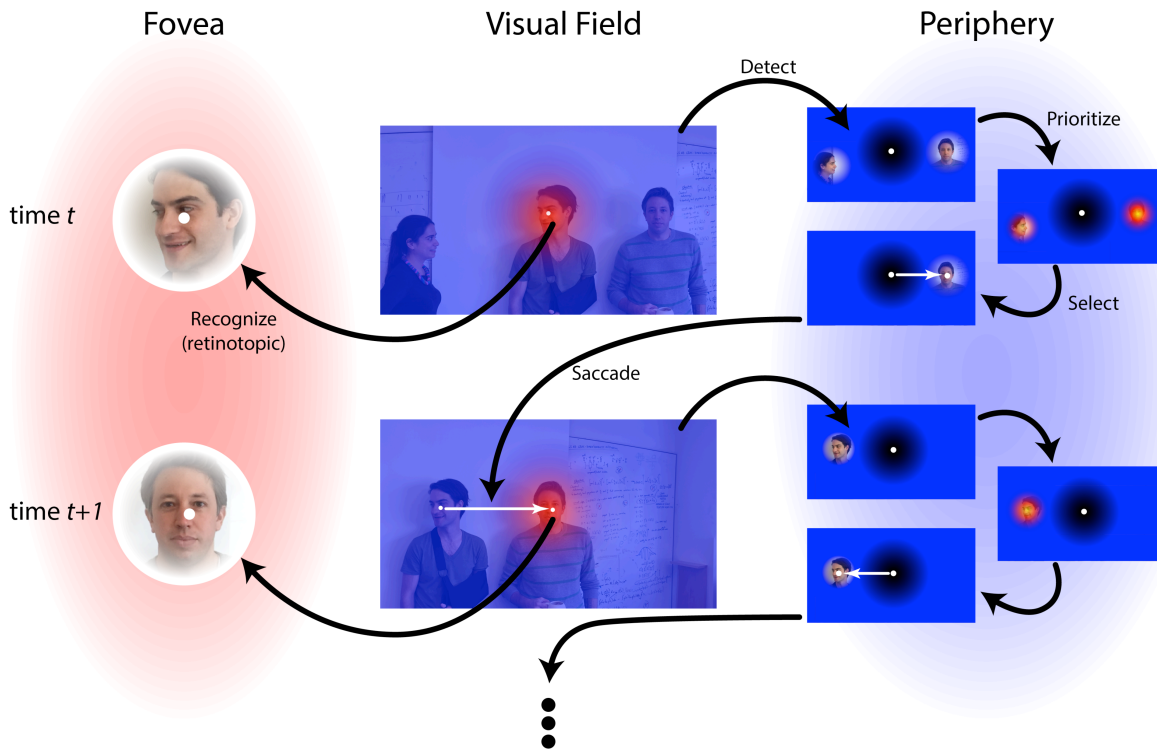
21 The work presented here builds on previous studies that have sought to characterize how
22 people move their eyes in naturalistic real-world environments and how these eye movements
23 relate to those observed under controlled laboratory conditions. Most mobile eye tracking studies
24 have assessed fixation behavior while subjects execute specific tasks, generally within a single
25 location (making tea or sandwiches: Hayhoe, 2000; Hayhoe & Ballard, 2005; Land, Mennie, &
26 Rusted, 1999; driving: Land, 1992; Land & Lee, 1994; visual search: Foulsham, Chapman,
27 Nasiopoulos, & Kingstone, 2014; Mack & Eckstein, 2011; gaze-cueing: Macdonald & Tatler,
28 2013, 2015; social: Einhäuser et al., 2009; Laidlaw, Foulsham, Kuhn, & Kingstone, 2011; Risko,
29 Laidlaw, Freeth, Foulsham, & Kingstone, 2012). A smaller number of studies have assessed eye
30 movements in unconstrained natural environments and behavior (Cristino & Baddeley, 2009;
31 Foulsham, Walker, & Kingstone, 2011; Hart et al., 2009). In general, these studies have assessed

1 coarse statistical trends across groups of subjects (e.g., tendency to fixate the image center in the
2 lab versus a “world-center”, the horizon, outside the lab;). The improved reliability of data
3 collection and efficiency of data analysis provided by the techniques developed here allow for a
4 significant expansion of the type and scope of real-world eye tracking studies (Figs. 2-4).

6 *Peripheral Detection and Foveal Recognition as Distinct Stages of Face Perception*

7 The evidence presented here suggests that real-world face recognition entails a systematic
8 sequence of processing steps in which detection operates in the periphery in parallel with
9 recognition at the fovea (Fig. 8). According to this hypothesis, the detection mechanism
10 continuously monitors for the presence of faces in the visual periphery (Step 1: Detect). Relevant
11 features of peripheral faces that can be computed with adequate precision (e.g., location, size,
12 pose, motion) are then combined to form a retinotopic “face priority map”, which is integrated
13 with other social and non-social priority calculations to form a general attention-guiding priority
14 map (Step 2: Prioritize; Bisley & Goldberg, 2010; Fecteau & Munoz, 2006; Itti, Koch, & Niebur,
15 1998; Koehler, Guo, Zhang, & Eckstein, 2014). Next, the highest priority location is selected for
16 subsequent fixation (Step 3: Select). When a face is selected for the next fixation, the eye
17 movement system exploits the stereotyped T-shaped configuration of facial features to precisely
18 target saccades to the individual’s specific preferred face-fixation position (Step 4: Saccade).
19 This brings the face image to a reliable position on the fovea, where it is processed by
20 specialized recognition mechanisms, shown previously to be highly retinotopically specific (Step
21 5: Recognize; Peterson & Eckstein, 2012). According to this model, face detection and face
22 recognition are fundamentally different processes, with detection occurring for faces in the
23 periphery at a wide range of eccentricities and positions, and recognition proceeding at the fovea
24 on faces that are usually centered at a single stereotyped retinal location. Note that steps 1-4
25 (detection, prioritization, selection, and saccadic targeting of peripheral faces) likely proceed in
26 parallel with Step 5 (recognition of the currently-foveated face).

27



1
 2 Figure 8. Schematic of a parallel peripheral-detection/foveal-recognition model. At any given time, t , the
 3 foveal (red) and peripheral (blue) retinal images are determined by the position of the body, head, and
 4 eyes. Peripheral mechanisms are tuned to image properties that support face detection. Faces likely to
 5 contain important visual information are selected and targeted with eye movements, providing powerful
 6 foveal resources for detailed recognition tasks. The eye movement is precise and individual-specific,
 7 eliminating image translation variance and possibly matching retinotopic face representations.

8
 9 The model of face perception just sketched can be tested using the methods developed in
 10 the current study. In particular, we can use our growing database of natural images our observers
 11 experienced (including their fixation position on those images) to ask: 1) Where do faces land on
 12 the retina in real-world viewing? 2) What are the features of peripherally-viewed faces that guide
 13 selection for saccadic targeting? 3) Is human size invariance for face recognition tuned to the
 14 statistics of retinal face sizes that occur during natural viewing? The general hypothesis, that we
 15 can now test in detail, is that the face detection and face recognition systems are each specifically
 16 tuned for task-specific statistics of experienced natural images.

17

18 *Retinal Image Statistics*

1 More broadly, this work makes possible a richer and more ecologically valid dataset with
2 which to test the core ideas of Natural Systems Analysis (Geisler, 2008): that the computations
3 employed by the visual system are the product of evolutionary optimization for the sensory
4 evidence (i.e., images) and tasks critical for survival. A deep understanding of these systems
5 requires knowledge of the properties of the visual environment in which they operate (i.e.,
6 natural image statistics; Botvinick, Weinstein, Solway, & Barto, 2015; Olshausen & Field, 1996;
7 Simoncelli & Olshausen, 2001; Torralba & Oliva, 2003). While the study of natural image
8 statistics has provided crucial insights into the computations carried out by the visual system, the
9 degree to which these images faithfully represent real-world visual experience is unclear. Prior
10 studies have typically analyzed sets of narrow field-of-view static photographs that have not
11 been selected to reflect everyday visual experience. Critically, these images do not have fixation
12 data, a critical missing element given the radically lower visual acuity in the visual periphery.
13 The framework presented here simultaneously collects high resolution, wide field-of-view video
14 of the visual environment and corresponding eye movements, allowing us to directly measure the
15 retinotopic images people experience in everyday life, which we term Retinal Image Statistics
16 (RIS). This new database should be applicable to myriad problems of vision beyond face
17 perception.

18

19 *Real-world face fixations in impaired populations*

20 Finally, the methods developed here enable us to rigorously measure real-world gaze
21 behavior in populations that may have deficits in face recognition. Fixation behavior may be a
22 prime determinant of successful face recognition, yet how those with possible recognition
23 deficits look at faces in the real world is largely unknown.

24 For example, a deficit in the recognition of faces is frequently reported in Autism
25 Spectrum Disorder (ASD). While the findings in the literature are conflicting, most evidence
26 suggests that face recognition impairments in ASDs are greater under natural viewing conditions
27 (e.g., static vs. dynamic, computer images vs. real faces; Jemel, Mottron, & Dawson, 2006;
28 Weigelt, Koldewyn, & Kanwisher, 2012). The literature is also conflicting on the question of
29 whether ASDs and TDs differ in the way they look at faces, but avoidance of faces in general
30 and eyes in particular apparently becomes more pronounced with increasing naturalism (Gharib,
31 Adolphs, & Shimojo, 2014; Klin, Jones, Schultz, Volkmar, & Cohen, 2002; Speer, Cook,

1 McMahon, & Clark, 2007). Most importantly, few studies have measured gaze behavior on faces
2 in natural viewing in ASD (Magrelli et al., 2013; Vabalas & Freeth, 2015), and none have done
3 so on a large scale during normal behavior in unconstrained environments. Overall, the evidence
4 suggests that any differences in face perception between ASDs and TDs should be greatest under
5 these conditions, which we can test in the future using the methods developed here.

6 Another disorder that may be informed by tests of real-world gaze behavior is
7 developmental prosopagnosia (DP), a lifelong deficit in face recognition in the absence of known
8 neurological damage (Behrmann & Avidan, 2005; Duchaine & Nakayama, 2006; Zhang, Liu, &
9 Xu, 2015). The few studies that have examined face-looking behavior in DPs have incorporated
10 small sample sizes (often a single patient) and laboratory viewing conditions (Barton, Radcliffe,
11 Cherkasova, & Edelman, 2007; Bate, Haslam, Tree, & Hodgson, 2008; Pizzamiglio et al., 2015;
12 Schmalzl, Palermo, Green, Brunsdon, & Coltheart, 2008; Schwarzer et al., 2006). A natural
13 hypothesis is that some or all of the deficits in face recognition in DPs result from suboptimal
14 and/or inconsistent looking behavior on faces, which could disrupt the normal development of
15 face representations and/or the ability to enact eye movement strategies that reliably constrain
16 retinotopic input.

17 Finally, it is of great interest to understand how, why, and when individuals acquire their
18 distinct face gaze behavior. One possibility is that retinotopic tuning of face representations is
19 present at birth, with location tuning varying across the population. This account holds that
20 individuals learn fixation strategies that are optimized for their specific tuning. A second, more
21 likely, possibility is that face representations are not strongly tuned to position at birth. Rather,
22 individuals vary, for whatever reasons, in where they look on faces. This early, retinotopic visual
23 experience might then guide the learning and development of the basic structure of face
24 representations. This situation could create a positive feedback scenario, such that the
25 performance advantage for fixating a specific region provides an incentive to maintain this
26 looking behavior. On this hypothesis, any early disruption of face-looking behavior could lock in
27 a self-reinforcing cycle of suboptimal face representations and suboptimal face-looking behavior,
28 providing a possible account of developmental prosopagnosia and/or face deficits in ASD. This
29 hypothesis could also account for the lifelong face perception deficits in individuals treated early
30 in life for bilateral or left (but not right) lateralized congenital cataracts that deprive face-
31 selective regions in the right hemisphere of patterned visual input for a brief period after birth

1 (Le Grand, Mondloch, Maurer, & Brent, 2001, 2003). A final possibility is that although face
2 representations are retinotopically specific, the general ability to encode new faces is not itself
3 tuned to an individual's particular fixation preference. Rather, consistently fixating the same
4 position causes most face memories to be encoded relative to the individual's specific preferred
5 gaze location. According to this hypothesis, the stability of an individual's specific face-fixation
6 behavior optimizes recognition by matching the retinotopic position of the current face to the
7 retinotopic positions of previously encoded faces. This matching hypothesis predicts that
8 individuals should identify new faces best when they are trained and tested at the same fixation
9 position. Critically, there should be no correlation between individual differences in preferred
10 fixation position and the fixation position during learning that leads to maximum recognition
11 performance during test.

1 Conclusion

2 In sum, we found that individual differences in face fixation behavior reported previously
3 in the lab generalize to real-world viewing. These findings suggest a distinction between two
4 components of face perception: detection of faces in the periphery, and recognition of faces in
5 the fovea. These findings also suggest possible causes of lifelong deficits in face perception in
6 developmental prosopagnosia, autism, and congenital cataracts. Finally, the methods developed
7 here make possible the large-scale collection natural images as seen by humans, including the
8 critical information of fixation position on each image, a dataset that may open up important new
9 constraints on natural systems analysis (Geisler, 2008).

10

11

1 Acknowledgments

2 We would like to thank Jason Fischer for his helpful comments on this manuscript. This work
3 was supported by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC
4 award CCF – 1231216. The authors have no financial or proprietary interests.

5

1 References

- 2 Barton, J. J. S., Radcliffe, N., Cherkasova, M. V., & Edelman, J. A. (2007). Scan patterns during
3 the processing of facial identity in prosopagnosia. *Experimental Brain Research*, *181*(2),
4 199–211. <http://doi.org/10.1007/s00221-007-0923-2>
- 5 Bate, S., Haslam, C., Tree, J. J., & Hodgson, T. L. (2008). Evidence of an eye movement-based
6 memory effect in congenital prosopagnosia. *Cortex*, *44*(7), 806–819.
7 <http://doi.org/10.1016/j.cortex.2007.02.004>
- 8 Behrmann, M., & Avidan, G. (2005). Congenital prosopagnosia: face-blind from birth. *Trends in*
9 *Cognitive Sciences*, *9*(4), 180–187. <http://doi.org/10.1016/j.tics.2005.02.011>
- 10 Bisley, J. W., & Goldberg, M. E. (2010). Attention, Intention, and Priority in the Parietal Lobe.
11 *Annual Review of Neuroscience*, *33*, 1–21. <http://doi.org/10.1146/annurev-neuro-060909->
12 [152823](http://doi.org/10.1146/annurev-neuro-060909-152823)
- 13 Botvinick, M., Weinstein, A., Solway, A., & Barto, A. (2015). Reinforcement learning, efficient
14 coding, and the statistics of natural tasks. *Current Opinion in Behavioral Sciences*, *5*, 71–
15 77. <http://doi.org/10.1016/j.cobeha.2015.08.009>
- 16 Cristino, F., & Baddeley, R. (2009). The nature of the visual representations involved in eye
17 movements when walking down the street. *Visual Cognition*, *17*(6/7), 880–903.
18 <http://doi.org/10.1080/13506280902834696>
- 19 DiCarlo, J. J., & Cox, D. D. (2007). Untangling invariant object recognition. *Trends in Cognitive*
20 *Sciences*, *11*(8), 333–341. <http://doi.org/10.1016/j.tics.2007.06.010>
- 21 Duchaine, B. C., & Nakayama, K. (2006). Developmental prosopagnosia: a window to content-
22 specific face processing. *Current Opinion in Neurobiology*, *16*(2), 166–173.
23 <http://doi.org/10.1016/j.conb.2006.03.003>

- 1 Durand, F., & Dorsey, J. (2002). Fast Bilateral Filtering for the Display of High-dynamic-range
2 Images. In *Proceedings of the 29th Annual Conference on Computer Graphics and*
3 *Interactive Techniques* (pp. 257–266). New York, NY, USA: ACM.
4 <http://doi.org/10.1145/566570.566574>
- 5 Einhäuser, W., Schumann, F., Vockeroth, J., Bartl, K., Cerf, M., Harel, J., ... König, P. (2009).
6 Distinct Roles for Eye and Head Movements in Selecting Salient Image Parts during
7 Natural Exploration. *Annals of the New York Academy of Sciences*, 1164(1), 188–193.
8 <http://doi.org/10.1111/j.1749-6632.2008.03714.x>
- 9 Fecteau, J. H., & Munoz, D. P. (2006). Saliency, relevance, and firing: a priority map for target
10 selection. *Trends in Cognitive Sciences*, 10(8), 382–390.
11 <http://doi.org/10.1016/j.tics.2006.06.011>
- 12 Foulsham, T., Chapman, C., Nasiopoulos, E., & Kingstone, A. (2014). Top-down and bottom-up
13 aspects of active search in a real-world environment. *Canadian Journal of Experimental*
14 *Psychology/Revue Canadienne de Psychologie Expérimentale*, 68(1), 8–19.
15 <http://doi.org/10.1037/cep0000004>
- 16 Foulsham, T., Walker, E., & Kingstone, A. (2011). The where, what and when of gaze allocation
17 in the lab and the natural environment. *Vision Research*, 51(17), 1920–1931.
18 <http://doi.org/10.1016/j.visres.2011.07.002>
- 19 Frank, M. C., Vul, E., & Johnson, S. P. (2009). Development of infants' attention to faces during
20 the first year. *Cognition*, 110(2), 160–170. <http://doi.org/10.1016/j.cognition.2008.11.010>
- 21 Geisler, W. S. (2008). Visual Perception and the Statistical Properties of Natural Scenes. *Annual*
22 *Review of Psychology*, 59(1), 167–192.
23 <http://doi.org/10.1146/annurev.psych.58.110405.085632>

- 1 Gharib, A., Adolphs, R., & Shimojo, S. (2014). “Don’t Look”: Faces with Eyes Open Influence
2 Visual Behavior in Neurotypicals but not in Individuals with High-Functioning Autism.
3 *Journal of Vision, 14*(10), 681–681. <http://doi.org/10.1167/14.10.681>
- 4 Gurler, D., Doyle, N., Walker, E., Magnotti, J., & Beauchamp, M. (2015). A link between
5 individual differences in multisensory speech perception and eye movements. *Attention,*
6 *Perception, & Psychophysics, 77*(4), 1333–1341. [http://doi.org/10.3758/s13414-014-](http://doi.org/10.3758/s13414-014-0821-1)
7 0821-1
- 8 Hart, B. M. ’t, Vockeroth, J., Schumann, F., Bartl, K., Schneider, E., König, P., & Einhäuser, W.
9 (2009). Gaze allocation in natural stimuli: Comparing free exploration to head-fixed
10 viewing conditions. *Visual Cognition, 17*(6-7), 1132–1158.
11 <http://doi.org/10.1080/13506280902812304>
- 12 Hayhoe, M. (2000). Vision Using Routines: A Functional Account of Vision. *Visual Cognition,*
13 7(1-3), 43–64. <http://doi.org/10.1080/135062800394676>
- 14 Hayhoe, M., & Ballard, D. (2005). Eye movements in natural behavior. *Trends in Cognitive*
15 *Sciences, 9*(4), 188–194. <http://doi.org/10.1016/j.tics.2005.02.009>
- 16 Holmqvist, K., Nyström, M., Andersson, R., Dewhurst, R., Jarodzka, H., & Weijer, J. van de.
17 (2011). *Eye Tracking: A comprehensive guide to methods and measures*. OUP Oxford.
- 18 Itti, L., Koch, C., & Niebur, E. (1998). A Model of Saliency-Based Visual Attention for Rapid
19 Scene Analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence,*
20 20(11), 1254–1259. <http://doi.org/10.1109/34.730558>
- 21 Jemel, B., Mottron, L., & Dawson, M. (2006). Impaired Face Processing in Autism: Fact or
22 Artifact? *Journal of Autism and Developmental Disorders, 36*(1), 91–106.
23 <http://doi.org/10.1007/s10803-005-0050-5>

- 1 Klin, A., Jones, W., Schultz, R., Volkmar, F., & Cohen, D. (2002). Visual Fixation Patterns
2 During Viewing of Naturalistic Social Situations as Predictors of Social Competence in
3 Individuals With Autism. *Archives of General Psychiatry*, 59(9), 809.
4 <http://doi.org/10.1001/archpsyc.59.9.809>
- 5 Koehler, K., Guo, F., Zhang, S., & Eckstein, M. P. (2014). What do saliency models predict?
6 *Journal of Vision*, 14(3), 14. <http://doi.org/10.1167/14.3.14>
- 7 Laidlaw, K. E. W., Foulsham, T., Kuhn, G., & Kingstone, A. (2011). Potential social interactions
8 are important to social attention. *Proceedings of the National Academy of Sciences*,
9 108(14), 5548–5553. <http://doi.org/10.1073/pnas.1017022108>
- 10 Land, M. F. (1992). Predictable eye-head coordination during driving. *Nature*, 359(6393), 318–
11 320. <http://doi.org/10.1038/359318a0>
- 12 Land, M. F., & Lee, D. N. (1994). Where we look when we steer. *Nature*, 369(6483), 742–744.
13 <http://doi.org/10.1038/369742a0>
- 14 Land, M. F., Mennie, N., & Rusted, J. (1999). The Roles of Vision and Eye Movements in the
15 Control of Activities of Daily Living. *Perception*, 28(11), 1311–1328.
16 <http://doi.org/10.1068/p2935>
- 17 Le Grand, R., Mondloch, C. J., Maurer, D., & Brent, H. P. (2001). Neuroperception: Early visual
18 experience and face processing. *Nature*, 410(6831), 890–890.
19 <http://doi.org/10.1038/35073749>
- 20 Le Grand, R., Mondloch, C. J., Maurer, D., & Brent, H. P. (2003). Expert face processing
21 requires visual input to the right hemisphere during infancy. *Nature Neuroscience*, 6(10),
22 1108–1112. <http://doi.org/10.1038/nn1121>

- 1 Macdonald, R. G., & Tatler, B. W. (2013). Do as eye say: Gaze cueing and language in a real-
2 world social interaction. *Journal of Vision*, 13(4), 6. <http://doi.org/10.1167/13.4.6>
- 3 Macdonald, R. G., & Tatler, B. W. (2015). Referent expressions and gaze: Reference type
4 influences real-world gaze cue utilization. *Journal of Experimental Psychology: Human*
5 *Perception and Performance*, 41(2), 565–575. <http://doi.org/10.1037/xhp0000023>
- 6 Mack, S. C., & Eckstein, M. P. (2011). Object co-occurrence serves as a contextual cue to guide
7 and facilitate visual search in a natural viewing environment. *Journal of Vision*, 11(9), 9.
8 <http://doi.org/10.1167/11.9.9>
- 9 Magrelli, S., Jermann, P., Noris, B., Ansermet, F., Hentsch, F., Nadel, J., & Billard, A. (2013).
10 Social orienting of children with autism to facial expressions and speech: a study with a
11 wearable eye-tracker in naturalistic settings. *Frontiers in Psychology*, 4.
12 <http://doi.org/10.3389/fpsyg.2013.00840>
- 13 Mehoudar, E., Arizpe, J., Baker, C. I., & Yovel, G. (2014). Faces in the eye of the beholder:
14 Unique and stable eye scanning patterns of individual observers. *Journal of Vision*, 14(7),
15 6. <http://doi.org/10.1167/14.7.6>
- 16 Olshausen, B. A., & Field, D. J. (1996). Natural image statistics and efficient coding. *Network:*
17 *Computation in Neural Systems*, 7(2), 333–339. <http://doi.org/10.1088/0954->
18 [898X_7_2_014](http://doi.org/10.1088/0954-898X_7_2_014)
- 19 Or, C. C.-F., Peterson, M. F., & Eckstein, M. P. (2015). Initial eye movements during face
20 identification are optimal and similar across cultures. *Journal of Vision*, 15(13).
21 <http://doi.org/10.1167/15.13.12>

- 1 Peterson, M. F., & Eckstein, M. P. (2012). Looking just below the eyes is optimal across face
2 recognition tasks. *Proceedings of the National Academy of Sciences*, *109*(48), E3314–
3 E3323. <http://doi.org/10.1073/pnas.1214269109>
- 4 Peterson, M. F., & Eckstein, M. P. (2013). Individual Differences in Eye Movements During
5 Face Identification Reflect Observer-Specific Optimal Points of Fixation. *Psychological*
6 *Science*, *24*(7), 1216–1225. <http://doi.org/10.1177/0956797612471684>
- 7 Peterson, M. F., & Eckstein, M. P. (2014). Learning optimal eye movements to unusual faces.
8 *Vision Research*, *99*, 57–68. <http://doi.org/10.1016/j.visres.2013.11.005>
- 9 Phillips, P. J., & O’Toole, A. J. (2014). Comparison of human and computer performance across
10 face recognition experiments. *Image and Vision Computing*, *32*(1), 74–85.
11 <http://doi.org/10.1016/j.imavis.2013.12.002>
- 12 Pizzamiglio, M. R., Luca, M. D., Vita, A. D., Palermo, L., Tanzilli, A., Dacquino, C., & Piccardi,
13 L. (2015). Congenital prosopagnosia in a child: Neuropsychological assessment, eye
14 movement recordings and training. *Neuropsychological Rehabilitation*, *0*(0), 1–40.
15 <http://doi.org/10.1080/09602011.2015.1084335>
- 16 Riesenhuber, M., & Poggio, T. (1999). Hierarchical models of object recognition in cortex.
17 *Nature Neuroscience*, *2*(11), 1019–1025. <http://doi.org/10.1038/14819>
- 18 Risko, E. F., Laidlaw, K. E. W., Freeth, M., Foulsham, T., & Kingstone, A. (2012). Social
19 attention with real versus reel stimuli: toward an empirical approach to concerns about
20 ecological validity. *Frontiers in Human Neuroscience*, *6*.
21 <http://doi.org/10.3389/fnhum.2012.00143>

- 1 Schmalzl, L., Palermo, R., Green, M., Brunsdon, R., & Coltheart, M. (2008). Training of familiar
2 face recognition and visual scan paths for faces in a child with congenital prosopagnosia.
3 *Cognitive Neuropsychology*, 25(5), 704–729. <http://doi.org/10.1080/02643290802299350>
- 4 Schwarzer, G., Huber, S., Grüter, M., Grüter, T., Groß, C., Hipfel, M., & Kennerknecht, I.
5 (2006). Gaze behaviour in hereditary prosopagnosia. *Psychological Research*, 71(5),
6 583–590. <http://doi.org/10.1007/s00426-006-0068-0>
- 7 Serre, T., Wolf, L., Bileschi, S., Riesenhuber, M., & Poggio, T. (2007). Robust Object
8 Recognition with Cortex-Like Mechanisms. *IEEE Transactions on Pattern Analysis and*
9 *Machine Intelligence*, 29(3), 411–426. <http://doi.org/10.1109/TPAMI.2007.56>
- 10 Simoncelli, E. P., & Olshausen, B. A. (2001). Natural Image Statistics and Neural
11 Representation. *Annual Review of Neuroscience*, 24(1), 1193–1216.
12 <http://doi.org/10.1146/annurev.neuro.24.1.1193>
- 13 Speer, L. L., Cook, A. E., McMahon, W. M., & Clark, E. (2007). Face processing in children
14 with autism: Effects of stimulus contents and type. *Autism*, 11(3), 265–277.
15 <http://doi.org/10.1177/1362361307076925>
- 16 Stampe, D. M. (1993). Heuristic filtering and reliable calibration methods for video-based pupil-
17 tracking systems. *Behavior Research Methods, Instruments, & Computers*, 25(2), 137–
18 142. <http://doi.org/10.3758/BF03204486>
- 19 Taigman, Y., Yang, M., Ranzato, M., & Wolf, L. (2014). DeepFace: Closing the Gap to Human-
20 Level Performance in Face Verification. In *2014 IEEE Conference on Computer Vision*
21 *and Pattern Recognition (CVPR)* (pp. 1701–1708).
22 <http://doi.org/10.1109/CVPR.2014.220>

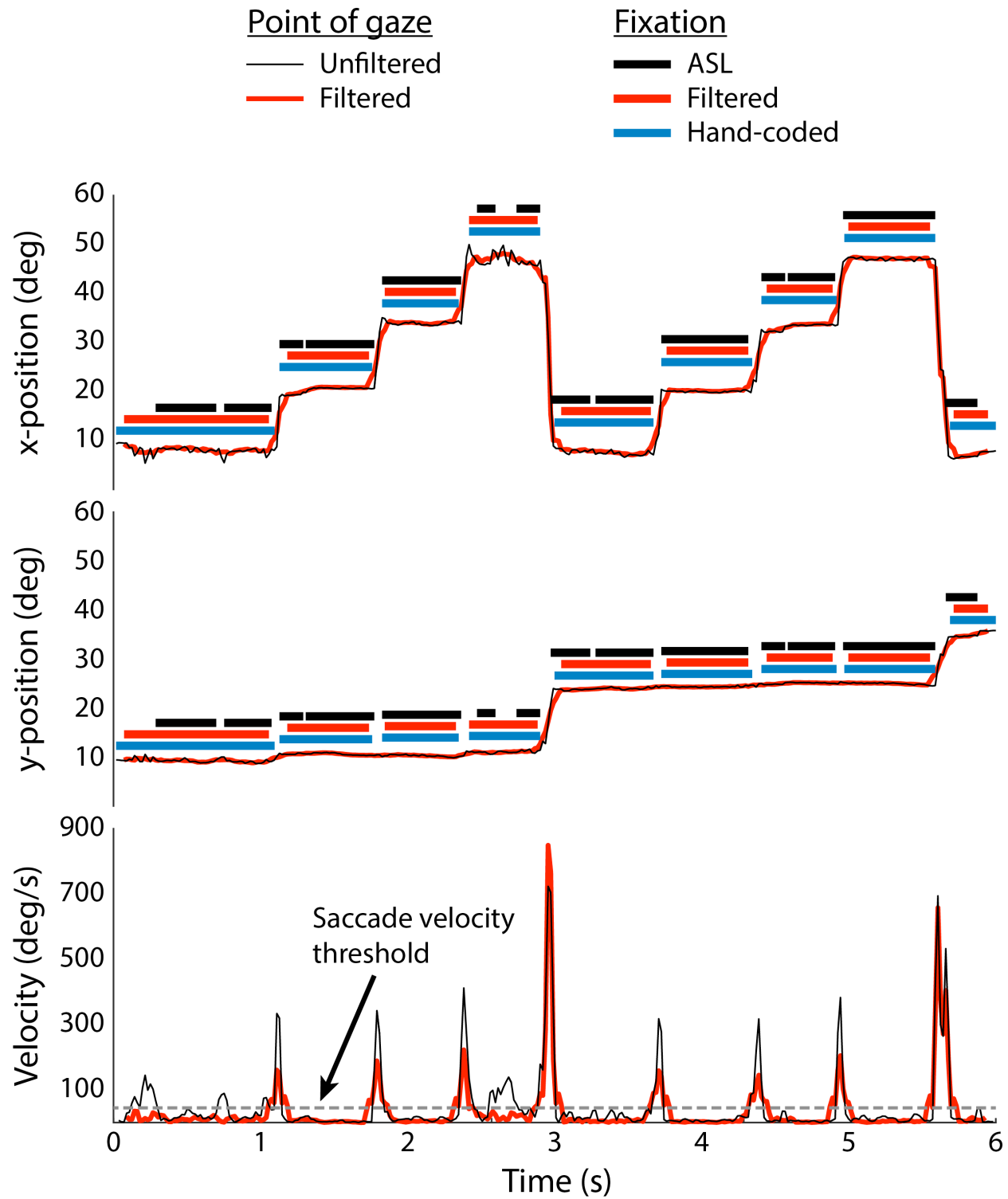
- 1 Tatler, B. W., & Vincent, B. (2008). Systematic tendencies in scene viewing. *Journal of Eye*
2 *Movement Research*, 2(2), 1–18.
- 3 Torralba, A., & Oliva, A. (2003). Statistics of natural image categories. *Network: Computation in*
4 *Neural Systems*, 14(3), 391–412. http://doi.org/10.1088/0954-898X_14_3_302
- 5 Vabalas, A., & Freeth, M. (2015). Brief Report: Patterns of Eye Movements in Face to Face
6 Conversation are Associated with Autistic Traits: Evidence from a Student Sample.
7 *Journal of Autism and Developmental Disorders*, 1–10. [http://doi.org/10.1007/s10803-](http://doi.org/10.1007/s10803-015-2546-y)
8 015-2546-y
- 9 Wass, S. V., Smith, T. J., & Johnson, M. H. (2012). Parsing eye-tracking data of variable quality
10 to provide accurate fixation duration estimates in infants and adults. *Behavior Research*
11 *Methods*, 45(1), 229–250. <http://doi.org/10.3758/s13428-012-0245-6>
- 12 Weigelt, S., Koldewyn, K., & Kanwisher, N. (2012). Face identity recognition in autism
13 spectrum disorders: A review of behavioral studies. *Neuroscience & Biobehavioral*
14 *Reviews*, 36(3), 1060–1084. <http://doi.org/10.1016/j.neubiorev.2011.12.008>
- 15 Zhang, J., Liu, J., & Xu, Y. (2015). Neural Decoding Reveals Impaired Face Configural
16 Processing in the Right Fusiform Face Area of Individuals with Developmental
17 Prosopagnosia. *The Journal of Neuroscience*, 35(4), 1539–1548.
18 <http://doi.org/10.1523/JNEUROSCI.2646-14.2015>

19

20

1 **Supplementary Figures**

2

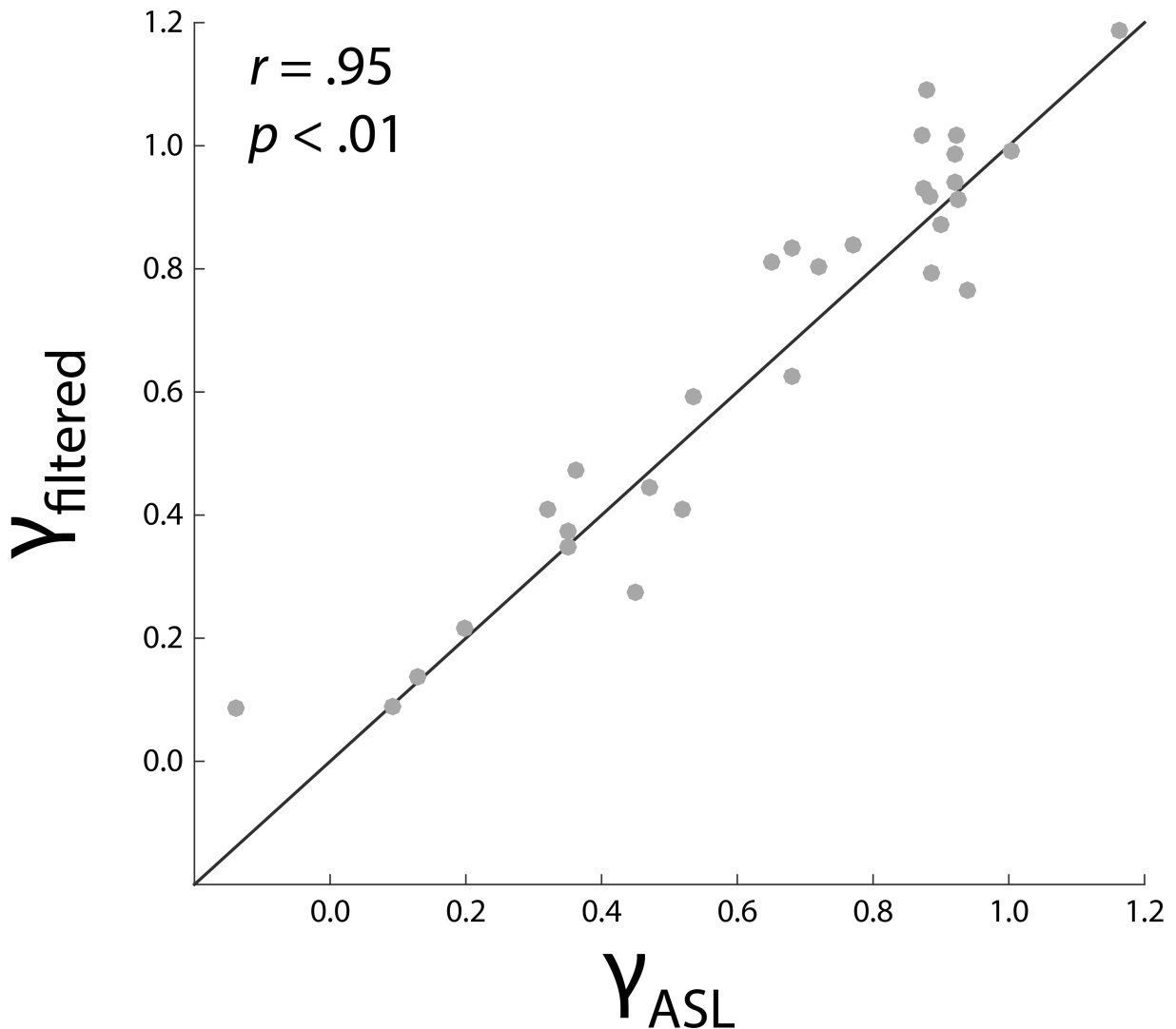


3

4 Figure S1. Example of fixation validation for a six second segment from one subject's data. Thin black
5 lines in the top and middle plots are the raw x and y gaze coordinates, respectively, with thin red lines

1 denoting the gaze position after bilateral filtering and interpolation. Above the traces, bars represent the
2 times of individual fixation events detected by the ASL algorithm (black), by the new filtering procedure
3 (red), and by hand (blue). The bottom plot shows instantaneous velocity for the raw (black) and filtered
4 (red) data, with the dotted grey line denoting the threshold used for saccade detection.

5
6



7

8 Figure S2. Correspondence between each subject's mean face-fixation location (γ) according to fixations
9 detected by the ASL algorithm (x-axis, γ_{ASL}) and the new filtering procedure (y-axis, $\gamma_{filtered}$).