



**"It's not a tumor": A framework for capitalizing on individual diversity to boost target detection**

Journal:	<i>Psychological Science</i>
Manuscript ID	PSCI-17-1243.R2
Manuscript Type:	Research Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Corbett, Jennifer; Brunel University London, College of Life and Health Sciences, Division of Psychology; Brunel University London, Center for Cognitive Neuroscience Munneke, Jaap; Brunel University London, College of Life and Health Sciences, Division of Psychology; Brunel University London, Center for Cognitive Neuroscience
Keywords:	Visual Search, Signal detection, Wisdom of crowds

SCHOLARONE™  
Manuscripts

Only

1  
2  
3 Averaging diversity boosts detection  
4  
5  
6  
7  
8  
9  
10  
11

12 **“It’s not a tumor”:**  
13 **A framework for capitalizing on individual diversity to boost target detection**  
14

15 Jennifer E. Corbett<sup>1,2</sup> \* & Jaap Munneke<sup>1,2</sup>

16 <sup>1</sup> *Brunel University London; College of Life and Health Sciences, Division of Psychology*

17 <sup>2</sup> *Brunel University London; Center for Cognitive Neuroscience*  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43

44 **\*Address correspondence to:**

45 Jennifer E. Corbett  
46 Brunel University London  
47 Division of Psychology, MJ-122  
48 Kingston Lane, UB8 3PH  
49 Uxbridge, London  
50 Phone: +44 756 309 3475  
51 jennifer.e.corbett@gmail.com  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Abstract**

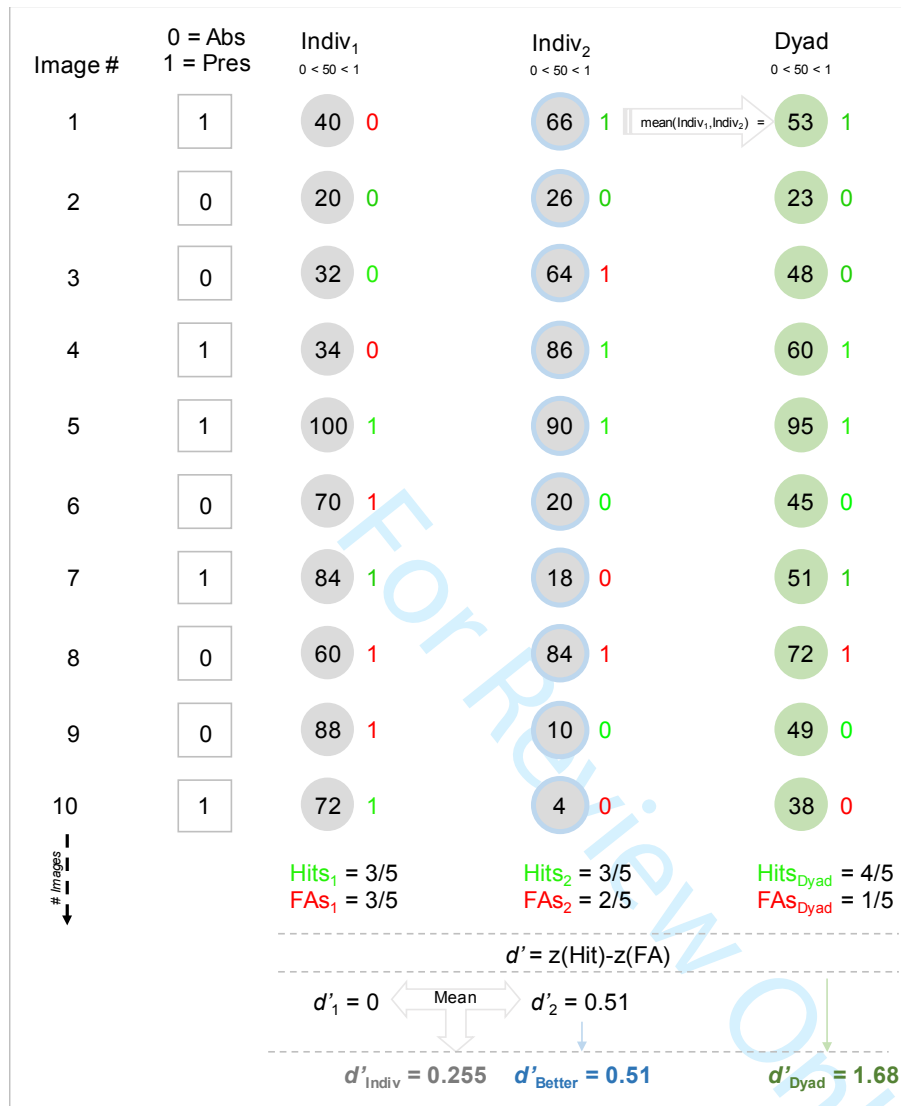
Even experts routinely miss infrequent targets, such as weapons in baggage scans or tumors in mammograms because the visual system is not equipped to notice the unusual. To date, limited progress has been made towards improving human factors that mediate such critical diagnostic tasks. Here we present a novel framework for pairing individuals' estimates to increase target detection. Using a "wisdom of crowds" approach capitalizing on the visual systems' ability to efficiently combine information, we demonstrate how averaging two, non-interacting individuals' continuous estimates of whether a briefly presented image contained a pre-specified target can significantly boost detection across a range of tasks. Furthermore, we show how pairing individuals' estimates to maximize decorrelated patterns of performance in one task can optimize performance on a separate task. Results make significant advances towards combating severe deficits in target detection, with straightforward applications for maximizing performance within limited pools of observers.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Despite spectacular developments in image processing and automated search algorithms (Liu et al., 2017), humans still perform the majority of critical visual search tasks such as breast cancer and airport baggage screenings. In stark contrast to a growing class of computer vision algorithms which input the entire image in parallel (Itti & Koch, 2001; Serre, Oliva, & Poggio, 2007) the human visual system is only capable of processing a few objects in detail at any given moment (Luck & Vogel, 1997). To circumvent these capacity restrictions, the visual system acts like a “statistician,” averaging out redundant and specific information, “filling in” details based on spatiotemporal context, and leaving us with a representation of abstract, statistical regularities to efficiently represent the maximum amount of information (Attneave, 1954; Baddeley, 1997; Barlow, 1961; Kersten, 1987; Olshausen & Field, 1996; Simoncelli & Olshausen, 2001). However, relying heavily on current and prior context means that we are designed not to notice targets that are out of the ordinary (Simons & Chabris, 1999). For example, even experts are prone to missing infrequent targets, such as weapons in baggage scans or tumors in routine mammograms (Wolfe, Horowitz, & Kenner, 2005).

For over a decade since Wolfe and colleagues’ initial discovery that we miss approximately 30% of rare targets (Wolfe et al., 2005), researchers have struggled to develop new methods to improve the detection of unusual and important targets (Van Wert, Horowitz, & Wolfe, 2009; Wolfe et al., 2007; Wolfe, Brunelli, Rubinstein, & Horowitz, 2013; Wolfe et al., 2005). Such attempts have mainly been aimed at increasing individual observers’ “vigilance,” for example, by interspersing bursts of frequent targets into otherwise rare target searches (Wolfe et al., 2007, 2013). Instead of trying to augment the performance of a system equipped to function optimally within its biological and computational constraints, our approach is grounded in observations abstracted from the visual system’s inherent structure and function. Although it

1  
2  
3 may seem that a system with restricted capacity for detail is destined for failure in complex  
4 search tasks, these apparent weaknesses can be translated into strengths if we approach this  
5 problem from the understanding that the visual system capitalizes on regularities by sampling the  
6 most unique, decorrelated aspects of visual input. For example, individual neurons coding for the  
7 most discrepant information in an image (Barlow, 1961) mimic the “wisdom of crowds”  
8 approach of asking multiple individuals for different ground truth guesstimates (Galton, 1907;  
9 Hogarth, 1978). Both cases take advantage of regression towards the mean, such that the more  
10 independent the samples, the more information that is accrued from averaging, and therefore the  
11 greater the reduction in error. Along these lines, studies of the “crowd within” have focused on  
12 similar methods for averaging guesses from a single individual to improve performance on a  
13 range of cognitive and perceptual tasks (Corbett, Fischer, & Whitney, 2011; Herzog & Hertwig,  
14 2009; Vul & Pashler, 2008). All of these findings point to the intriguing and potentially life-  
15 saving possibility of even greater benefits of independence from combining two estimates from  
16 *different* observers. To these ends, we developed and tested the novel framework illustrated in  
17 Figure 1 for pairing individuals’ continuous estimates of target presence to increase detection  
18 performance across a range of tasks.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

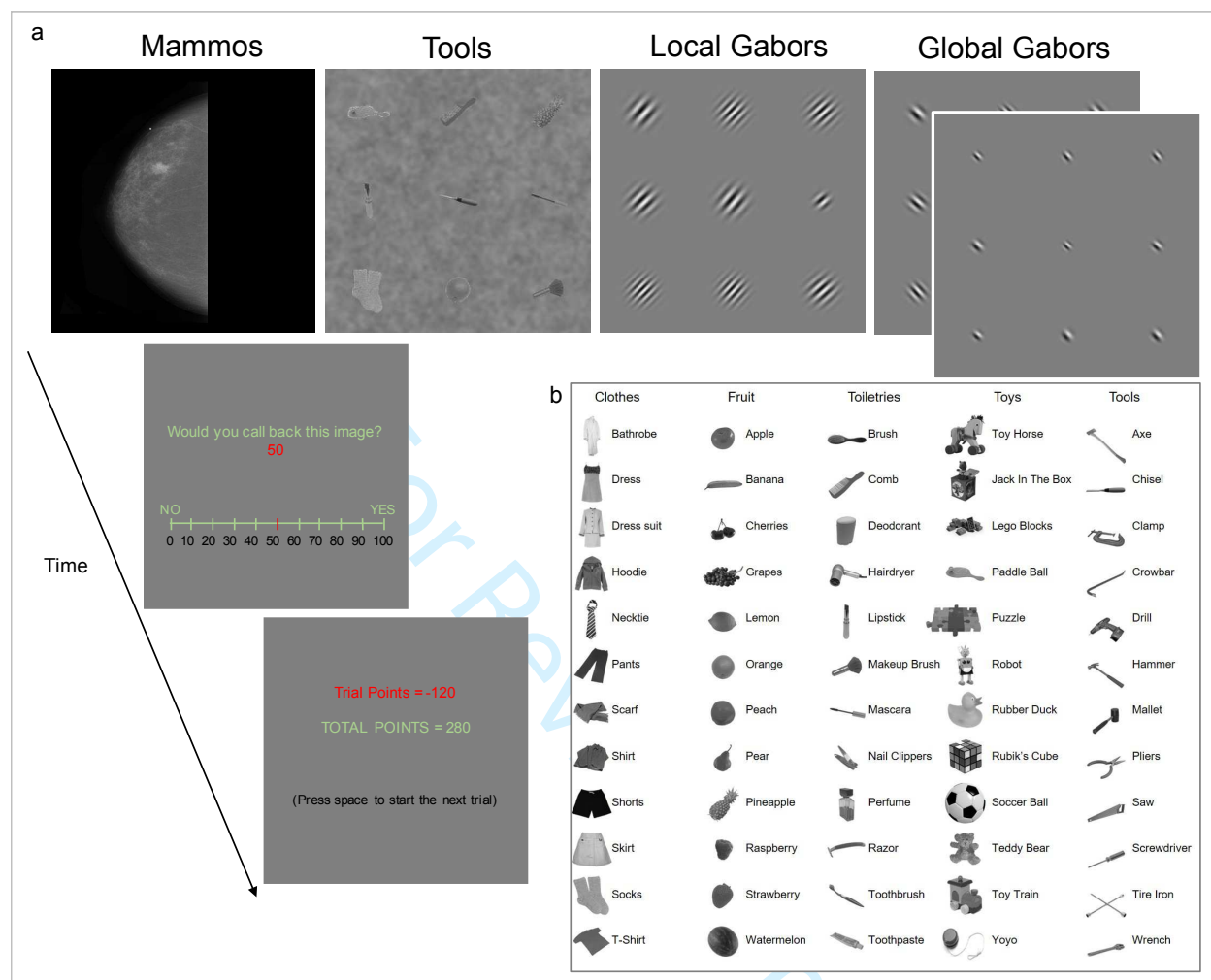


**Figure 1. Framework for improving performance by averaging individuals' estimates.** For each image, individuals provide a continuous estimate (from 0=No, to 100=Yes) about whether a pre-specified target is present. Individuals perform the task independently. Each row represents hypothetical responses to a single image. The first column corresponds to the specific image number (Image #), and the number inside the square in the second column indicates whether a target is present (1) or absent (0) in that image. Numbers inside gray circles in the third and fourth columns represent the corresponding continuous responses of the two observers in the hypothetical dyad, and red and green 0's and 1's to the right of these estimates represent how that continuous response translates into a binary response. Responses less than 50 are scored as 0=Absent, and responses greater than 50 are scored as 1=Present. Green numbers indicate correct responses and red numbers indicate incorrect responses, such that green 1's correspond to Hits and red 1's correspond to False alarms (FAs). Numbers inside green circles in the final column represent the dyad's average estimate (the average of Indiv<sub>1</sub> and Indiv<sub>2</sub>'s continuous responses), and the red and green 0's and 1's to the right of these estimates represent the dyad's corresponding binary response. Each individual's  $d'$  ( $d'_1$  and  $d'_2$ ) is calculated by subtracting their normalized FA rates (FAs<sub>1</sub> and FAs<sub>2</sub>) from their respective normalized Hit rates (Hits<sub>1</sub> and Hits<sub>2</sub>), and the dyad's  $d'$  ( $d'_{\text{Dyad}}$ ; green text) is calculated by subtracting the dyad's normalized FA rate (FAs<sub>Dyad</sub>) from the dyad's normalized Hit rate (Hits<sub>Dyad</sub>). The average  $d'$  for the dyad ( $d'_{\text{Indiv}}$ ; gray text) is calculated as the average of the two individuals'  $d'$ s (the mean of  $d'_1$  and  $d'_2$ ). Indiv<sub>2</sub> (gray circles outlined in blue) had a higher  $d'$  than Indiv<sub>1</sub> ( $d'_2 > d'_1$ ), and is therefore used as the  $d'$  of the better observer in the dyad ( $d'_{\text{Better}}$ ; blue text). Overall, averaging individuals' estimates results in lower error rates (Hits and FAs) and therefore a higher rate of target detection ( $d'$ ) than the average of the individuals' detection rates or the detection rate of the better individual in the dyad.

1  
2  
3 Using this framework, we conducted two experiments to test the predictions that: 1)  
4 pairing individuals' estimates will result in better performance compared to the performance of  
5 either individual (Experiments 1 & 2), and 2) pairing estimates from individuals with the most  
6 decorrelated performance on one task will result in the greatest improvements on a separate task  
7 (Experiment 2).  
8  
9  
10  
11  
12  
13  
14  
15  
16

17 In a first experiment, we investigated how averaging independent estimates from two  
18 different individuals might improve target detection. We presented naïve observers with  
19 unilateral mammography x-rays verified as having either an abnormality (a tumor) or consisting  
20 of normal, healthy tissue. Observers' task was to rate the extent to which they thought a tumor  
21 was present on a pseudo-continuous scale ranging from 0 (No) to 100 (Yes), excluding 50.  
22 Participants completed two "Mammos" conditions in which either 5% or 50% of the total  
23 number of trials contained tumor present images (see Methods and Figure 2a for more details).  
24 Given the noted benefits of combining independent information (Barlow, 1961; Corbett et al.,  
25 2011; Galton, 1907; Herzog & Hertwig, 2009; Hogarth, 1978; Vul & Pashler, 2008), we  
26 predicted that the average of two observers' estimates would result in higher  $d'$  values (a signal  
27 detection measure that accounts for both types of possible miss and false alarm errors) than  
28 either of their individual estimates.  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Methods



**Fig. 2. a. Trial Sequence.** On each trial, observers were presented with a search display for 500 ms, followed by a response slider to indicate the extent to which they thought the image should be called back for containing a pre-specified target. Feedback about the points earned on the trial and the total points accumulated over the entire duration of the experimental block was displayed after a response was made. *Mammos (Exp. 1)*: Observers were presented with a unilateral mammogram and their task was to decide whether a tumor was present. They completed two conditions in which either 5% or 50% of the trials contained tumor present images (5% Mammos and 50% Mammos) in a blocked design. *Tools (Exp 2)*: In the Tools task, observers were presented with a display of nine objects and their task was to decide if a tool target was present. They completed two conditions (blocked) with a tool target present on either 5% or 50% of trials (5% Tools and 50% Tools). *Local Gabors (Exp 2)*: On each trial in the Local Gabors condition, observers were presented with an image of nine Gabors, and their task was to decide whether or not to call the image back for containing a target Gabor that was larger or smaller than all the other Gabors. *Global Gabors (Exp 2)*: On each trial in the Global Gabors condition, observers were presented with two successive displays of nine Gabors, temporally separated by a 1000 ms blank screen, and their task was to decide whether to call back the sequence for having a “target” difference in global mean size between the first and second images. In both Gabors conditions, 50% of trials contained a target. **b. Objects used in the Tools tasks.** As in Wolfe and colleagues (Wolfe et al., 2005), gray-scaled objects (Brodeur, Dionne-Dostie, Montreuil, & Lepage, 2010; Konkle, Brady, Alvarez, & Oliva, 2010) from five different stimulus categories (including tools) were used to create the stimulus displays in the Tools task. Each category consisted of 12 exemplars.



1  
2  
3       *Participants:* There were 25 participants in Experiment 1, and 23 participants in  
4  
5 Experiment 2. All participants were students at Bilkent University, participated in only one of the  
6  
7 experiments, had normal or corrected-to-normal vision, and volunteered to take part in the  
8  
9 experiment either for course credit or money. In each experiment, the participant with the highest  
10  
11 points won an additional 50 Turkish Lira (see *Points/Feedback*). Seven participants were  
12  
13 dismissed from each experiment for failing to pass a pre-established performance threshold  
14  
15 during practice blocks (see *Points/Feedback*). Only the data from the remaining 18 observers in  
16  
17 Experiment 1 (11 female, mean age=20.56 years, age range=19-25), and 16 observers in  
18  
19 Experiment 2 (14 female, mean age=21.25 years, age range=18-25) were included in subsequent  
20  
21 analyses. We had to restrict our sample size to 16 participants in Experiment 2, as this was the  
22  
23 maximum number of individuals for which we could generate all possible paired combinations  
24  
25 given computational limitations (see *Decorrelated Dyads*). All procedures and protocols were in  
26  
27 accordance with the guidelines of Bilkent University's ethical review board.  
28  
29  
30  
31  
32  
33  
34

35       *Task & Trial sequence:* Participants were instructed to act like radiologists (Experiment  
36  
37 1) or airport baggage screeners (Experiment 2), and estimate whether they detected the presence  
38  
39 of a target item from a pre-specified category (e.g., a tumor, a tool) in a briefly presented image.  
40  
41 On each trial in Experiment 1, participants were presented with a unilateral mammogram for 500  
42  
43 ms. On each trial in Experiment 2, participants were presented with an image of nine objects for  
44  
45 500 ms. In both experiments, the target screen was followed by a response slider ranging from 0-  
46  
47 100 (adapted from Evans, Georgian-Smith, Tambouret, Birdwell, & Wolfe, 2013; Evans,  
48  
49 Haygood, Cooper, Culpan, & Wolfe, 2016). Participants used the mouse to adjust the slider to  
50  
51 indicate whether to call back the image based on whether they detected a target (Figure 2a). To  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 implement a modified version of a two-alternative forced-choice task (2AFC), we restricted the  
4 slider responses so that participants could choose any value from 0-100 except 50. After they  
5 finished adjusting the slider, participants were presented with feedback on their points for the  
6 given trial and their total points for the block until they pressed the space bar to continue to the  
7 next trial.  
8  
9  
10  
11  
12  
13  
14  
15  
16

17 *Stimuli & Procedure:* All participants performed the tasks individually and  
18 independently, in a dimly lit room. Participants had no knowledge that their performance would  
19 be paired with another individual offline after the experiment, nor did they receive any feedback  
20 that had any relationship to the performance of any other participant. An HP PC was used to  
21 present stimuli at a viewing distance of 57 cm, on a 21.3” NEC LCD monitor with a 60-Hz  
22 refresh rate and a resolution of 1600 x 1200 pixels. MATLAB (Version 2016a) and the  
23 Psychophysics Toolbox (PTB-3; Brainard, 1997; Pelli, 1997) controlled all presentation, timing,  
24 response functions, and data collection. See the Supplementary Material available online for  
25 details of the training and practice blocks associated with each experiment.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 Experiment 1: Each participant completed two experimental conditions (5% and 50%  
41 target prevalence) in a single, one-hour experimental session. The two conditions were run in  
42 separate blocks, and the order of blocks was counterbalanced over participants. In each  
43 condition, observers were presented with a unilateral mammogram centered in the middle of a  
44 black screen. Mammograms were acquired from the University of South Florida Digital  
45 Database for Screening Mammography (DDSM; Heath et al., 1998; Heath, Bowyer, Kopans,  
46 Moore, & Kegelmeyer, 2001), and presented in native pixel resolution (subtending  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 approximately  $10^\circ \times 18^\circ$  of visual angle). There were 400 unique images in each condition. In the  
4  
5 5% condition, a tumor was present in 20 of the 400 images. In the 50% condition, a tumor was  
6  
7 present in 200 of the 400 images. All participants saw the same set of 400 images in the 5%  
8  
9 condition and another set of 400 images in the 50% condition. The order of the images in each  
10  
11 set was randomized, and each participant saw each image only once throughout the experiment.  
12  
13 See the Supplementary Material available online for a detailed description of the mammogram  
14  
15 stimuli and the corresponding supplementary image database.  
16  
17  
18  
19  
20  
21

22 Participants completed two stages of training before beginning the experimental blocks.  
23  
24 First, they were presented with 50 of the target present mammograms and 50 of the target absent  
25  
26 (healthy) mammograms, in random order. Participants were not required to make a response, but  
27  
28 only to study the images so that they became familiarized with typical target present and absent  
29  
30 images. Each image remained on the screen until the participant pressed the space bar. If the  
31  
32 image contained a target, the original mammogram was re-presented with an overlay showing  
33  
34 the region identified as abnormal circled in red. The 100 trials in the second training stage were  
35  
36 identical to trials in the main 50% experimental condition, except: 1) the training images were  
37  
38 presented simultaneously with the response slider and remained on the screen until the  
39  
40 participant responded, and 2) if the image contained a target, the image was re-presented with the  
41  
42 corresponding overlay of the abnormality. None of the images used in the training and practice  
43  
44 blocks were used in the experimental blocks.  
45  
46  
47  
48  
49  
50

51 Experiment 2: Each participant completed four experimental conditions (Local Gabors,  
52  
53 Global Gabors, 5% Tools, and 50% Tools; Figure 2) in a single experimental session lasting  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 approximately 2 hours. The two Gabors conditions and the two Tools conditions were run in  
4 immediate succession, but the order of these two types of conditions and the order of the  
5 individual conditions within each type were counterbalanced over observers. In each condition,  
6 observers were presented with nine stimuli arranged equidistantly in a 3 x 3 grid, subtending  
7 approximately 16° x 16° of visual angle and centered in the middle of a gray screen.  
8  
9

- 10 • *Local Gabors*: Participants were presented with a total of 288 images, with a target  
11 defined as a Gabor that was smaller or larger than the other eight Gabors present in  
12 50% of the images. For each image, the size, spatial frequency, contrast, and  
13 orientation of the nine individual Gabors were pseudo-randomly assigned from two  
14 possible sets of nine values per dimension (chosen based on the results of previous  
15 pilot studies). The small size set ranged in diameter from 1.9° to 2.1° in 0.02° steps,  
16 and the large size set ranged from 3.4° to 3.6°, also in 0.02° steps. The low spatial  
17 frequency set ranged from 1.2 cycles per degree (cpd) to 2 cpd in 0.1 cpd steps, and  
18 the high spatial frequency set ranged from 6 cpd to 6.8 cpd in 0.1 cpd steps. The low  
19 contrast set ranged from 73% to 81% contrast relative to the gray background, and the  
20 high contrast set ranged from 83% to 91% contrast. The left-tilted orientation set  
21 ranged from -49° to -41° of tilt from vertical, and the right-tilted set ranged from 41°  
22 to 49° of tilt from vertical.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

45 144 target absent images were generated by factorially combining the 2 Sizes x 2  
46 Spatial frequencies x 2 Contrasts x 9 Locations, with two repetitions of each of the  
47 combinations. The set of leftward or rightward orientations was randomly selected for  
48 the nine Gabors in each image. For each target absent image, a corresponding target  
49 present image was constructed from an identical set of parameters except that the  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Gabor representing the mean size of the set in the target absent condition was  
4 replaced by a Gabor representing the mean size of the opposite set in the target  
5 present condition. With the exception of the target size and location in target present  
6 images, the individual values in each dimension were randomly assigned to the nine  
7 possible locations in each image.  
8  
9  
10  
11  
12  
13

- 14  
15  
16  
17 • *Global Gabors*: Participants were presented with a total of 256 pairs of successive  
18 images of nine Gabors each. The target was present in 50% of the pairs, and was  
19 defined as a difference between the mean size of the items in the first image of nine  
20 Gabors and the mean size of the nine Gabors in the second image. There were nine  
21 possible sets of nine individual Gabor sizes, constructed by multiplying a base set  
22 ranging from  $1.3^\circ$  to  $1.7^\circ$  in  $0.05^\circ$  steps by a constant value from 0 to 8. In target  
23 present pairs, the sizes of the Gabors in the first or second image (determined  
24 randomly for the pair) were taken from sets 1, 2, 3, 4, 6, 7, 8, or 9. The sizes of the  
25 Gabors in the other image were taken from set 5, for a total of eight possible mean  
26 size differences between the two images in each pair. In 50% of target absent pairs,  
27 the sizes of both images were taken from the same set (“repeated pairs”; 1, 2, 3, 4, 6,  
28 7, 8, or 9). In the other half of target absent pairs, the sizes of Gabors in both images  
29 were taken from set 5 (“set 5 pairs”). This manipulation guarded against the tendency  
30 to call back the pair of images on a given trial as having a target difference present if  
31 either of the images in the pair contained the set 5 sizes. The same spatial frequencies,  
32 contrasts, and orientations were used as in the Local Gabors condition. The entire set  
33 of 256 image pairs was constructed from two repetitions of the full factorization of  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 the 8 possible Mean size differences x 2 Present/Absent conditions x 2 Spatial  
4 frequencies x 2 Contrasts x 2 Absent pair types (repeated pairs/set 5 pairs). Each  
5 target present image had a corresponding target absent image with all of the same  
6 parameters except that both target absent images contained the same nine sizes (either  
7 repeated pairs or set 5 pairs). As in the Local Gabors condition, the set of leftward or  
8 rightward orientations for the nine Gabors was randomly selected for each pair of the  
9  
10 128 target present images, and the same orientation set was used for the  
11 corresponding target absent pair of images. In all images in all sequential pairs, the  
12 individual values in each dimension were randomly assigned to the nine possible  
13 locations.  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

- 29 • *5% Tools*: There were a total of 360 images, each comprised of nine objects. A single  
30 tool target was present in 5% of the images, and the other eight items were drawn  
31 from one of four categories (Clothing, Fruit, Toiletries, and Toys; Figure 2b). There  
32 were 12 objects per category (including Tools), for a total of 60 possible items. The  
33 images for individual items were taken from open-source databases published in  
34 previous work (Brodeur et al., 2010; Konkle et al., 2010). Each item subtended  $2.5^\circ \times$   
35  $2.5^\circ$  (the average size of the individual Gabors). Individual items were presented at  
36 60% transparency, each item in each image was tilted a random  $\pm 15^\circ$  from vertical in  
37  $1^\circ$  steps, and a pink noise mask was added to the entire image in attempt to better  
38 equate the task to a real-world baggage screening task (Wolfe et al., 2005). There  
39 were 18 target present images (two per location) and 342 target absent images. For  
40 each target present image, there were 19 corresponding target absent images, each  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 containing the same nine non-tool objects drawn randomly from one of the four non-  
4 tool categories. Although the same nine items were used in the target absent images,  
5 the individual items were placed randomly in the nine locations in each image. In the  
6 corresponding target present image, a tool was randomly selected to replace the non-  
7 tool item in the target location, with the locations of the eight other non-tool items  
8 randomly determined.

- 19 • *50% Tools*: A tool target was present in half of the 360 total images. All parameters  
20 were identical to those used in the 5% Tools condition, except that there were 180  
21 target present images (20 per location) and 180 target absent images, with target  
22 absent images containing nine randomly selected items from the four non-tool  
23 categories and the corresponding target present image containing a tool randomly  
24 selected to replace one of the non-tool items and presented in the target location.

36 *Points/Feedback*: On each trial, participants earned points for a correct response or lost  
37 points for an incorrect response (adapted from Wolfe et al., 2005). In all 50% conditions,  
38 participants could earn a maximum of 40 points when there was a target present and they  
39 correctly called back the image (i.e., adjusted the slider to a value above 50%; Hit), and 2.5  
40 points when there was no target present and they correctly did not call back the image (Correct  
41 rejection). They could lose a maximum of 40 points if there was a target and they did not call the  
42 image back (Miss), and 75 points if there was no target and they incorrectly called back the  
43 image (False alarm). In the 5% conditions, participants could earn a maximum of 400 points for  
44 a Hit, and a maximum of 2.5 points for a Correct rejection. They could lose a maximum of 400  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 points for a Miss and a maximum of 75 points for a False alarm. We calculated the specific  
4 number of points earned or lost on each trial as a function of the extremeness of the given  
5 response using the following formula:  
6  
7  
8  
9

$$\text{Points} = (x * \text{MaxPoints}) * (1 - ((2 * \text{Response}) / 100))$$

10  
11  
12  
13  
14  
15  
16  
17 where  $x = -1$  for signal present trials,  $x = 1$  for signal absent trials, MaxPoints was the maximum  
18 number of points that could be earned or lost, and Response was the participant's estimate  
19 between 0 and 100 made via the slider adjustment on the trial.  
20  
21  
22  
23

24 In replicating the rare target deficit originally reported by Wolfe and colleagues (2007,  
25 2005), it was necessary to implement a points system that prevented participants from simply  
26 responding 0 on each trial in the 5% prevalence condition (which would result in 95% accuracy),  
27 or simply hovering in the middle of the response range without actually performing the detection  
28 task. Along these lines, the points system allowed us to better equate participants' response  
29 criteria (tendency to say yes or no more frequently) across the two types of conditions (Healy &  
30 Kubovy, 1981; Wolfe & Van Wert, 2010), such that there was a greater penalty for a bias to  
31 respond "YES" in the 50% conditions when targets were frequent, but a greater penalty for a bias  
32 to respond "NO" in the 5% conditions when targets were rare. Participants were informed of the  
33 specific points for each type of response at the start of each block, and required to have a positive  
34 score (points > 0) at the end of every block.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

## 51 **Results**

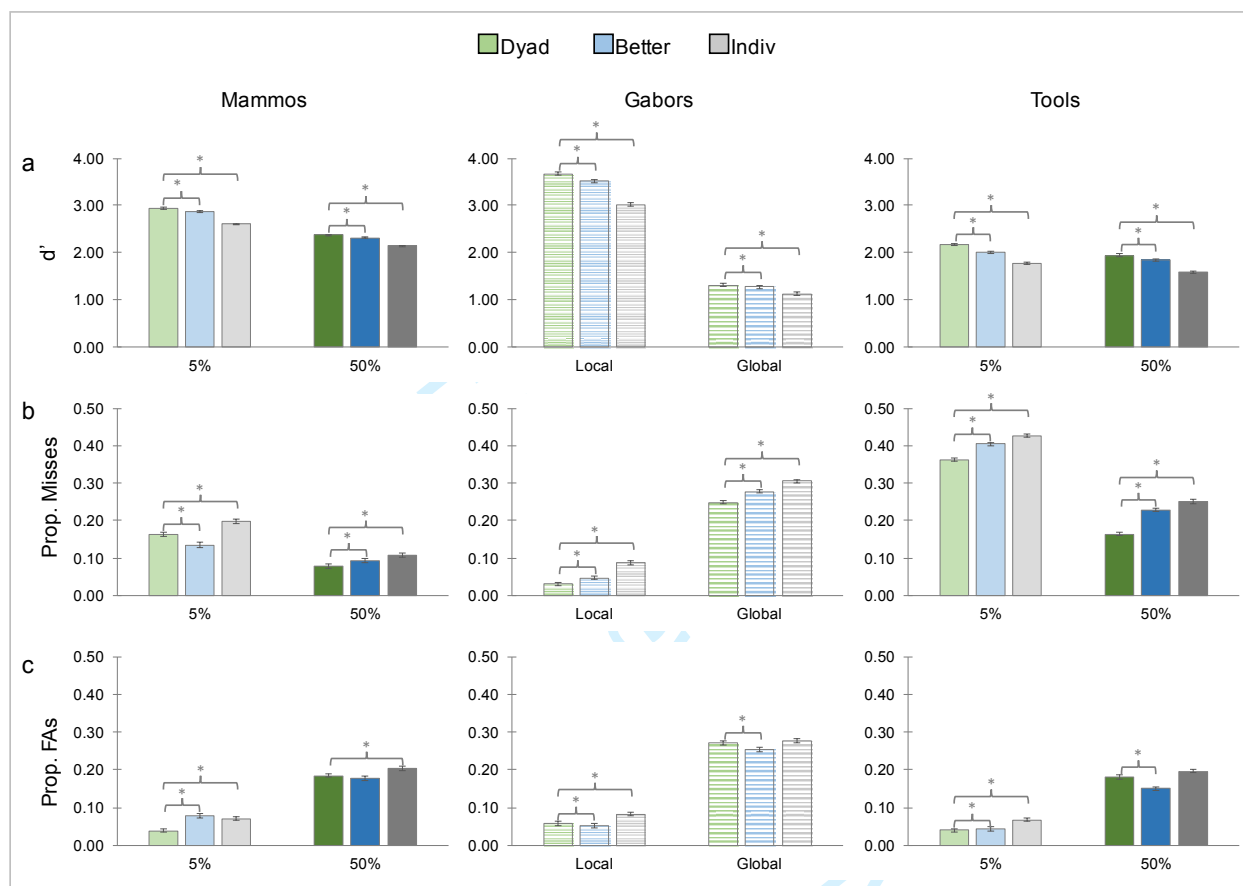
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3        *Dyads versus individuals:* To gauge whether pairing individuals' estimates yielded  
4 improved detection performance, we first paired the 18 (Experiment 1) and 16 (Experiment 2)  
5 participants in  $n-1$  unique arrangements of pairings. Each arrangement included all 18 or 16  
6 participants. Only the pairings of the participants into nine (Experiment 1) and eight (Experiment  
7 2) dyads was different over arrangements, such that each possible dyad was present only once in  
8 the resultant 17 (Experiment 1) and 15 (Experiment 2) arrangements. For example, the dyad  
9 formed by participants 1 and 2 only appeared in one of the 17 possible arrangements of nine  
10 dyads in Experiment 1. In other words, all participants were included in each arrangement of  
11 pairings, but each dyad was unique across all pairings such that the combination of the two  
12 individuals forming a dyad was never repeated (see Table S1 for an illustration of this  
13 procedure). As we planned the statistical comparisons illustrated in Figure 3 and Table S2  
14 between three groups, we conducted our analysis using only these unique arrangements of  
15 pairings based on multiple combinations of the individual participants.  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

33        As illustrated in Figure 3a, averaging individual estimates from two observers improved  
34 signal detection ( $d'$ ) over estimates obtained from the corresponding individuals in both the 5%  
35 Mammos condition (12.89% average improvement) and the 50% Mammos condition (10.81%  
36 average improvement). In both conditions, pairing estimates also improved performance relative  
37 to the performance of the better individual in each dyad. We conducted a follow-up analysis to  
38 explicate how pairing estimates affected both types of miss and false alarm errors. Although  
39 pairing always decreased both misses and false alarms compared to estimates made by the  
40 corresponding individuals, there was a noteworthy trade-off between misses and false alarms in  
41 the 5% Mammos condition, suggesting that estimates made by the better of the two observers in  
42 each dyad were less prone to misses whereas dyad estimates were less likely to result in false  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

alarms (Figure 3b&c). In all reported analyses, paired comparisons refer to two-tailed t-tests that were Bonferroni-corrected for multiple comparisons. Table S2 in the Supplemental Material available online presents the full results of all planned comparisons.



**Figure 3. Dyad versus individual results.** Experiment 1 (Mammos):  $n=9$  dyads; Experiment 2 (Gabors & Tools):  $n=8$  dyads. **a. Detection performance ( $d'$ ):** For each task, dyads' (green) average  $d'$  was significantly higher compared to the average  $d'$  of the corresponding individuals (gray) as well as the  $d'$  of the better of the two individuals (blue). **b. Misses:** Also for each task, the miss rate was significantly lower for dyads compared to individuals, and for dyads compared to the better of the two individuals, with the exception of the 5% Mammos condition. **c. False alarms:** False alarms were significantly lower for dyads versus individuals in all tasks except the Global Gabors and 50% Tools conditions, false alarms were significantly lower for dyads versus the better of the two individuals in the 5% Mammos condition, and false alarms were significantly greater for dyads versus the better of the two individuals in the Global Gabors and 50% Tools conditions. Asterisks represent significant differences, with  $\alpha = 0.0125$  using the Bonferroni correction for four multiple comparisons. Error bars represent the 95% within-subjects confidence intervals for the two-way interaction depicted in each graph (Loftus & Masson, 1994).

*Decorrelated dyads:* After confirming the significant benefits of pairing individuals' independent estimates of the likelihood of tumor presence in mammography images, we sought to generalize this benefit using other types of visual search tasks in Experiment 2. In particular,

1  
2  
3 we first replicated the results of Experiment 1 using a tool detection task in Experiment 2 that is  
4 typically used with naïve observers as a proxy for weapons detection in baggage scans (Fleck &  
5 Mitroff, 2007; Wolfe et al., 2005; Figure 2). Given that previous findings suggest greater benefits  
6 of averaging as estimates become more independent (Corbett et al., 2011; Herzog & Hertwig,  
7 2009; Hogarth, 1978; Vul & Pashler, 2008), we further extended our investigation to examine  
8 whether pairing estimates from observers with the most independent, decorrelated patterns of  
9 performance in separate Gabors tasks could maximize performance in the main Tools search  
10 tasks. Given previous findings that suggest detection might be based in part on a global “gist”  
11 signal (Drew, Evans, Võ, Jacobson, & Wolfe, 2013; Evans et al., 2013, 2016; Rensink, 2004),  
12 we used local and global versions of the Gabors task (Figure 2). This resulted in 4 main  
13 conditions: Local 5%, Global 5%, Local 50%, and Global 50%. We expected the improvement in  
14  $d'$  from averaging dyads' performance in the Tools tasks to increase with their decreasingly  
15 correlated performance in the Gabors tasks.  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34

35 Similar to real-world situations like airport baggage screenings where a pair of security  
36 officers typically work at one of several stations, we sought to optimize overall *group*  
37 performance in the Tools tasks by pairing individuals into eight dyads, such that performance on  
38 the Gabors task was maximally decorrelated over the entire group of 16 individuals. **Importantly,**  
39 **and unlike the analysis used in both Experiments 1 and 2 to compare the performance of dyads**  
40 **and individuals (Figure 3 and Table S2), we only compared dyad performance as a function of**  
41 **correlation on the Gabors tasks and did not restrict this analysis to only one unique occurrence of**  
42 **a particular dyad. Instead, we began by constructing all possible pairings of the 16 individuals**  
43 **using a custom-written recursive MATLAB function, such that one set of the eight dyads was**  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **distinct** relative to any other combination (the only consideration was the uniqueness of the  
4 dyads within a combination relative to the dyads in all previous combinations. Dyads could  
5 repeat over different combinations of pairings, but not within the same combination.). For  $n=8$   
6 possible dyads of the 16 observers, this resulted in  
7  
8  
9  
10  
11

$$\frac{(2n)!}{n!} * \frac{1}{2^n} = 2,027,025 \text{ unique combinations of 16 observers into 8 dyads}$$

12  
13  
14  
15  
16  
17 Using this approach, the number of unique combinations increased exponentially with each  
18 additional dyad, such that standard computational power limited our analysis to a maximum of  
19 16 individuals grouped into eight dyads.  
20  
21  
22

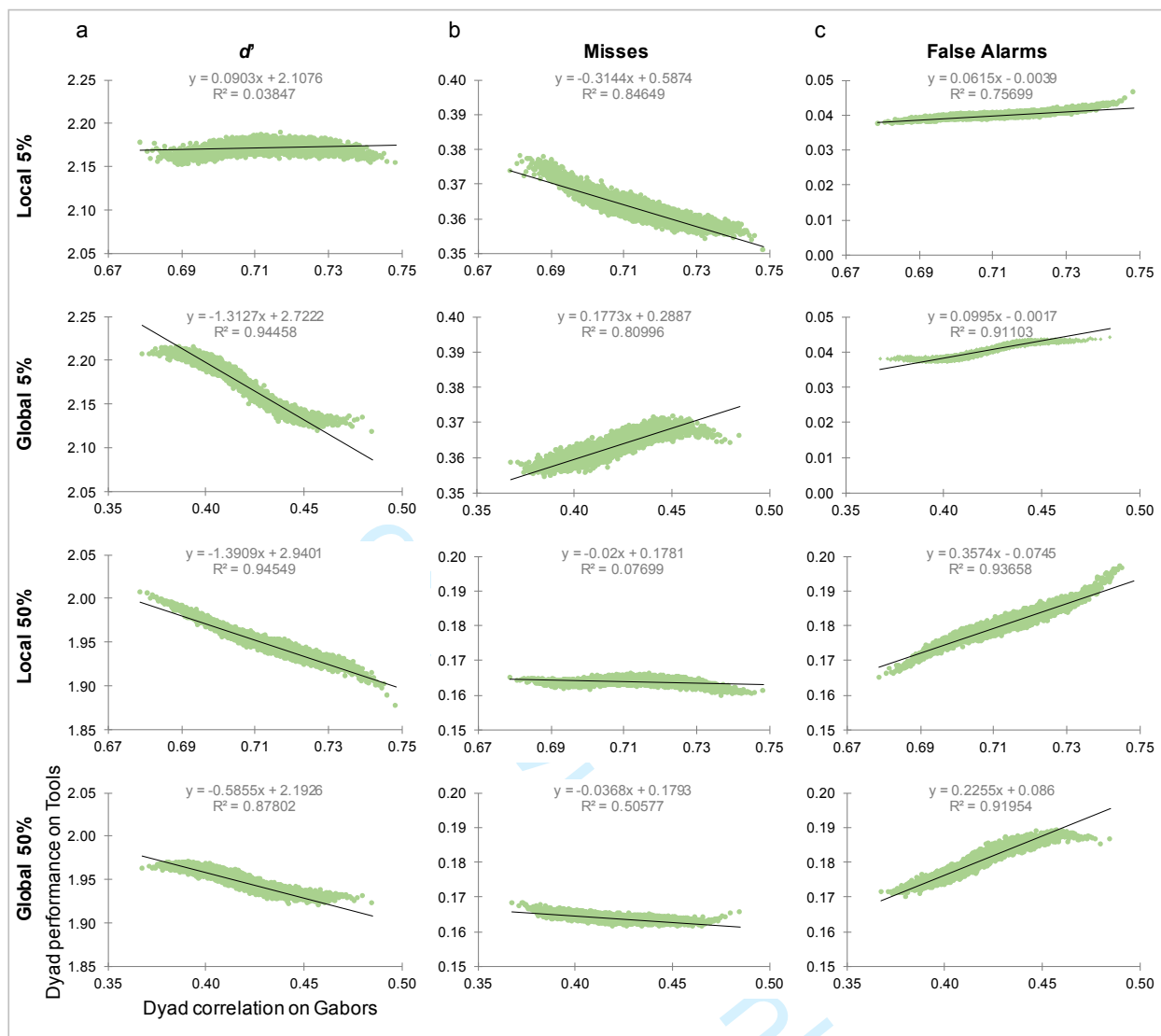
23  
24 Next, separately for the Local and Global Gabors tasks, we calculated the average  
25 correlation between dyads' responses for the individual images (the 288 images in the Local  
26 Gabors condition and the 256 images in the Global Gabors condition) for each set of 2,027,025  
27 combinations. More specifically, the average correlation in responses for each set of eight dyads  
28 was found by summing the eight correlations between each pair of individuals ( $Indiv_1$  and  
29  $Indiv_2$ ), then dividing by the total number of (8) dyads. After finding the average correlation for  
30 each set of combinations, the combinations were sorted from most decorrelated to most  
31 correlated. We then used these combinations on each of the Local and Global Gabors tasks to  
32 compute dyad performance for each of the 2,027,025 combinations in the 5% and 50% Tools  
33 tasks, resulting in four main conditions (Local 5%, Global 5%, Local 50%, and Global 50%).  
34  
35 Finally, we averaged the data from the 2,027,025 combinations over 9,009 bins of 225 pairings  
36 each to facilitate graphical and statistical comparisons.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

51 Although participants' responses were not speeded, their average response times for both  
52 rare (5% Tools:  $M=2020$  ms,  $SD=1332$  ms) and frequent (50% Tools:  $M=1902$  ms,  $SD=1076$  ms)  
53 target detection were similar to those previously reported by Wolfe and colleagues (2005).  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 After first replicating the improved rates of target detection for dyad versus individual  
4 estimates observed in the mammogram task in Experiment 1 for both Gabors and Tools  
5 conditions in Experiment 2 (E2: 16.06% – 22.14% average improvement from averaging; Figure  
6 3 and Table S2), we tested the prediction that pairing estimates from observers such that their  
7 overall performance was maximally decorrelated in the Gabors tasks would result in the greatest  
8 improvements from averaging in the Tools tasks. In support of this hypothesis, we observed a  
9 strong negative correlation ( $r=-0.972$ ), such that  $d'$  values on the 5% Tools task increased as the  
10 average correlation over a given set of pairings in the Global Gabors task decreased (Global 5%).  
11 On the contrary, there was no evidence of a relationship between  $d'$  on the 5% Tools task and the  
12 average correlation over a set of pairings in the Local Gabors task (Local 5%;  $r=0.196$ ).  
13 Improved performance in the 50% Tools task was also well-predicted by observers' decorrelated  
14 performance in the Global (Global 50%;  $r=-0.937$ ) and Local (Local 50%;  $r=-0.972$ ) Gabors  
15 tasks. Taken together, these patterns of results in Figure 4a suggest that pairing estimates from  
16 individuals with decorrelated performance on a simple global detection task can maximize  
17 performance in a separate, more complex rare target detection task, whereas pairing estimates  
18 from individuals with the most decorrelated patterns of performance regardless of the local or  
19 global scale of the simple task may optimize the detection of frequently occurring targets in a  
20 more complex task.

21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45 These conclusions are further supported by the different patterns of improvement in miss  
46 and false alarm rates in the two Tools tasks that resulted from pairing the most decorrelated  
47 observers in the two Gabors tasks. In the 5% Tools condition where dyads showed the lowest  
48 miss and false alarm rates (Figure 3b&c), the more decorrelated the dyads' performance in the  
49 Global Gabors condition, the lower their miss and false alarm rates (Figure 4b&c). In the 50%

1  
2  
3 Tools condition, dyads showed the lowest miss rates (Figure 3b) but higher false alarm rates than  
4 the better of the two individuals in each dyad (Figure 3c). Also in the 50% condition, the more  
5 decorrelated the dyads performed in the Global and especially the Local Gabors conditions, the  
6 greater the improvement in false alarm rates (Figure 4c) with little change in the miss rates  
7 (Figure 4b), such that the greatest improvements from averaging decorrelated dyads were  
8 observed where improvements were most needed. Further evidence of the benefits of capitalizing  
9 on the more fine-grained information available in participants' continuous estimates are  
10 illustrated by the inferior results from the same analyses using participants' dichotomized  
11 estimates in the Gabors tasks instead of their continuous responses (Figure S1).  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



**Figure 4. Decorrelated dyads results.** **a. Detection performance ( $d'$ ):** For all conditions except Local 5%, there was a strong negative correlation between dyads' correlated performance on the Gabors task and their corresponding  $d'$  values on the Tools task. **b. Misses:** Whereas dyads' miss rate in the 5% Tools task decreased with increasingly correlated performance in the Local Gabors task, it decreased with decreasingly correlated performance in the Global Gabors task. There was no strong relationship between dyads' correlated performance in either of the Gabors tasks and dyads' miss rates in the 50% Tools task. **c. False alarms:** In the 5% Tools condition, false alarms slightly increased as dyads' performance became increasingly correlated in both the Local and Global Gabors conditions. In the 50% Tools condition, there was a strong positive correlation between dyads' false alarm rates and correlated performance in both Local and Global Gabors tasks. All plots depict the results of simple regression analyses with all  $F_s(1,9008) > 360$  and all  $p_s < 0.001$ .

## Discussion

The current framework represents a significant advancement towards improving human factors that mediate target detection. Building on parallels between the visual system's ability to

1  
2  
3 capitalize on the variance inherent in the surrounding environment (e.g., Barlow, 1961) and  
4 findings that the “wisdom of crowds” (e.g., Galton, 1907) can be approximated by aggregating  
5 the most conflicting estimates from the same individual (e.g., Corbett et al., 2011; Herzog &  
6 Hertwig, 2009; Vul & Pashler, 2008) the present study demonstrated that averaging estimates  
7 from individuals with the most decorrelated patterns of performance can improve detection  
8 across a range of visual tasks.  
9

10  
11  
12  
13  
14  
15  
16  
17 Several critical differences set the current framework apart from methods used in  
18 previous work (Van Wert et al., 2009; Wolfe et al., 2007, 2013, 2005). First, given evidence that  
19 observers can rapidly detect but not localize changes (Rensink, 2004), and that experts can detect  
20 but not localize the presence of abnormalities in mammography images within a fraction of a  
21 second (Drew et al., 2013; Evans et al., 2013, 2016), we used brief 500 ms stimulus presentations  
22 in attempt to engage the nonselective visual pathway thought to be responsible for the rapid  
23 extraction of global, statistical information (Drew et al., 2013; Evans et al., 2013, 2016; Rensink,  
24 2004; Wolfe, Võ, Evans, & Greene, 2011). Second, instead of having observers provide a binary  
25 present/absent judgment *or a binary judgement followed by a confidence rating*, we required  
26 them to make a continuous response by adjusting a slider to indicate the extent to which they  
27 thought an image should be called back for further inspection because it contained a target. This  
28 allowed us to average individuals’ estimates for each of the unique images, instead of limiting us  
29 to multiplying their average error rates to obtain a measure of paired performance (Wolfe et al.,  
30 2007) *or calculating a separate measure of metacognitive confidence from binary responses and*  
31 *confidence ratings (e.g., Maniscalco & Lau, 2012)*. Third, the current framework does not  
32 require people to interact to reach a joint decision. Human interaction decreases the  
33 independence of individuals’ judgements and results in joint performance that is sub-optimal  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 (e.g., Bang et al., 2017), as it can be no greater than improvements obtained from relying on the  
4 answers of the more confident (Bahrami et al., 2010) and/or accurate (e.g., Erev, Gopher, Itkin,  
5 & Greenspan, 1995) member of each dyad. Instead, observers performed all tasks individually  
6 without interacting or any knowledge that their responses would be combined with those of other  
7 participants. Finally, we used performance on a basic visual task (Gabor) to pair individuals'  
8 estimates on a separate, more complex detection task (Tools). Averaging estimates from  
9 individuals with the most decorrelated patterns of performance in either the Local or Global  
10 Gabor task lead to the greatest improvements in detecting frequent Tools targets. In contrast,  
11 only averaging estimates from individuals with the most decorrelated performance on the Global  
12 Gabor task lead to improvements in detecting rare Tools targets. Whereas previous attempts to  
13 improve rare target detection have interspersed bursts of frequent targets into searches for  
14 infrequent targets (Wolfe et al., 2007, 2013), the present results suggest that the optimal solution  
15 is to pair estimates from observers using different cognitive strategies.  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

33 The most straightforward applications of our results are to begin employing continuous  
34 bounded response scales in safety-critical signal detection tasks and averaging individuals'  
35 estimates. For example, in-line with a recent report by Mayo clinic researchers underscoring the  
36 importance of second opinions in increasing the accuracy of medical diagnoses (Van Such, Lohr,  
37 Beckman, & Naessens, 2017), different radiologists can quickly scan through mammograms and  
38 continuous estimates from the individuals with the most decorrelated patterns can be combined  
39 offline. Importantly, our framework can be applied without requiring any communication  
40 between individuals and without requiring the delay between estimates that is necessary to obtain  
41 similar benefits from aggregating repeated estimates from the same individual (Corbett et al.,  
42 2011; Herzog & Hertwig, 2009; Vul & Pashler, 2008). Therefore, estimates made by an expert at  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 one location (e.g., a transportation security officer at an airport checkpoint) could be paired in  
4  
5 real-time with estimates made by a second person viewing the same images at an offsite location.  
6  
7 In the future, it may even be possible to simulate a decorrelated pattern of performance that can  
8  
9 be averaged online with an existing individual's performance, which could be particularly cost-  
10  
11 effective and useful when personnel or data is sparse. Although we note several possible  
12  
13 applications in human and computer vision, the potential application domain extends broadly  
14  
15 across a diverse range of fields from navigation to engineering to economic forecasting.  
16  
17

18  
19 In conclusion, we have demonstrated a simple yet effective method for capitalizing on the  
20  
21 diversity afforded by averaging continuous estimates of target presence from different  
22  
23 individuals to boost target detection. Not only does this method yield benefits from combining  
24  
25 independent samples in a range of different visual tasks, but our results strongly suggest that  
26  
27 detection in complex tasks can be optimized by pairing estimates from individuals with the most  
28  
29 diverse patterns of performance on a secondary basic detection task.  
30  
31

32  
33  
34  
35 **Acknowledgements:** We thank Bahar Aykut, Miray Kapan, Elif Beyza Koş, Sümeyra Özçalık,  
36  
37 and the members of the Visual Perception and Attention Lab (VPAL) for their help with data  
38  
39 collection. We also thank Aaron Clarke for supplying the code for pre-processing the  
40  
41 mammography images and for his help rating the initial set of mammogram stimuli. Many  
42  
43 thanks to Chris Oriet for helpful comments on earlier versions of this work.  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## References

- 1  
2  
3  
4 Attneave, F. (1954). Some informational aspects of visual perception. *Psychological Review*,  
5  
6  
7 61(3), 183–193. <https://doi.org/10.1037/h0054663>
- 8  
9 Baddeley, R. (1997). The Correlational Structure of Natural Images and the Calibration of  
10  
11 Spatial Representations. *Cognitive Science*, 21(3), 351–372.  
12  
13 [https://doi.org/10.1207/s15516709cog2103\\_4](https://doi.org/10.1207/s15516709cog2103_4)
- 14  
15 Bahrami, B., Olsen, K., Latham, P. E., Roepstorff, A., Rees, G., & Frith, C. D. (2010). Optimally  
16  
17 interacting minds. *Science (New York, N.Y.)*, 329(5995), 1081–1085.  
18  
19 <https://doi.org/10.1126/science.1185718>
- 20  
21 Bang, D., Aitchison, L., Moran, R., Hecce Castanon, S., Rafiee, B., Mahmoodi, A., ...  
22  
23 Summerfield, C. (2017). Confidence matching in group decision-making. *Nature Human*  
24  
25 *Behaviour*, 1(6), 0117. <https://doi.org/10.1038/s41562-017-0117>
- 26  
27 Barlow, H. (1961). Possible principles underlying the transformations of sensory messages. In  
28  
29 W. Rosenblith (Ed.), *Sensory Communication* (pp. 217–234). MIT Press.
- 30  
31 Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, 10(4), 433–436.  
32  
33 <https://doi.org/10.1163/156856897X00357>
- 34  
35 Brodeur, M. B., Dionne-Dostie, E., Montreuil, T., & Lepage, M. (2010). The Bank of  
36  
37 Standardized Stimuli (BOSS), a New Set of 480 Normative Photos of Objects to Be Used  
38  
39 as Visual Stimuli in Cognitive Research. *PLOS ONE*, 5(5), e10773.  
40  
41 <https://doi.org/10.1371/journal.pone.0010773>
- 42  
43 Corbett, J. E., Fischer, J., & Whitney, D. (2011). Facilitating stable representations: Serial  
44  
45 dependence in vision. *PloS One*, 6(1), e16701.
- 46  
47 Drew, T., Evans, K., Võ, M. L.-H., Jacobson, F. L., & Wolfe, J. M. (2013). Informatics in  
48  
49 Radiology: What Can You See in a Single Glance and How Might This Guide Visual  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Search in Medical Images? *Radiographics*, 33(1), 263–274.

4  
5 <https://doi.org/10.1148/rg.331125023>  
6

7  
8 Erev, I., Gopher, D., Itkin, R., & Greenspan, Y. (1995). Toward a Generalization of Signal  
9  
10 Detection Theory to N-Person Games: The Example of Two-Person Safety Problem.  
11  
12 *Journal of Mathematical Psychology*, 39(4), 360–375.

13  
14  
15 <https://doi.org/10.1006/jmps.1995.1034>  
16

17 Evans, K. K., Georgian-Smith, D., Tambouret, R., Birdwell, R. L., & Wolfe, J. M. (2013). The  
18  
19 gist of the abnormal: Above-chance medical decision making in the blink of an eye.  
20  
21 *Psychonomic Bulletin & Review*, 20(6), 1170–1175. [https://doi.org/10.3758/s13423-013-](https://doi.org/10.3758/s13423-013-0459-3)  
22  
23 0459-3  
24

25  
26 Evans, K. K., Haygood, T. M., Cooper, J., Culpan, A.-M., & Wolfe, J. M. (2016). A half-second  
27  
28 glimpse often lets radiologists identify breast cancer cases even when viewing the  
29  
30 mammogram of the opposite breast. *Proceedings of the National Academy of Sciences of*  
31  
32 *the United States of America*, 113(37), 10292–10297.

33  
34  
35 <https://doi.org/10.1073/pnas.1606187113>  
36

37  
38 Fleck, M. S., & Mitroff, S. R. (2007). Rare Targets Are Rarely Missed in Correctable Search.  
39  
40 *Psychological Science*, 18(11), 943–947. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-9280.2007.02006.x)

41  
42 9280.2007.02006.x  
43

44 Galton, F. (1907). Vox populi (The wisdom of crowds). *Nature*, 75(7), 450–451.

45  
46  
47 Healy, A. F., & Kubovy, M. (1981). Probability matching and the formation of conservative  
48  
49 decision rules in a numerical analog of signal detection. *Journal of Experimental*  
50  
51 *Psychology: Human Learning and Memory*, 7(5), 344–354. [https://doi.org/10.1037/0278-](https://doi.org/10.1037/0278-7393.7.5.344)  
52  
53 7393.7.5.344  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 Heath, M., Bowyer, K., Kopans, D., Kegelmeyer, P., Moore, R., Chang, K., & Munishkumaran,  
4  
5 S. (1998). Current Status of the Digital Database for Screening Mammography. In *Digital*  
6  
7 *Mammography* (pp. 457–460). Springer, Dordrecht. [https://doi.org/10.1007/978-94-011-](https://doi.org/10.1007/978-94-011-5318-8_75)  
8  
9 [5318-8\\_75](https://doi.org/10.1007/978-94-011-5318-8_75)  
10  
11
- 12 Heath, M., Bowyer, K., Kopans, D., Moore, R., & Kegelmeyer, P. (2001). The Digital Database  
13  
14 for Screening Mammography. In *Proceedings of the Fifth International Workshop on*  
15  
16 *Digital Mammography* (pp. 212–218). Medical Physics Publishing.  
17  
18
- 19 Herzog, S. M., & Hertwig, R. (2009). The Wisdom of Many in One Mind: Improving Individual  
20  
21 Judgments With Dialectical Bootstrapping. *Psychological Science*, *20*(2), 231–237.  
22  
23 <https://doi.org/10.1111/j.1467-9280.2009.02271.x>  
24  
25
- 26 Hogarth, R. M. (1978). A note on aggregating opinions. *Organizational Behavior and Human*  
27  
28 *Performance*, *21*(1), 40–46. [https://doi.org/10.1016/0030-5073\(78\)90037-5](https://doi.org/10.1016/0030-5073(78)90037-5)  
29  
30
- 31 Itti, L., & Koch, C. (2001). Computational modelling of visual attention. *Nature Reviews*  
32  
33 *Neuroscience*, *2*(3), 194–203. <https://doi.org/10.1038/35058500>  
34  
35
- 36 Kersten, D. (1987). Predictability and redundancy of natural images. *JOSA A*, *4*(12), 2395–2400.  
37  
38 <https://doi.org/10.1364/JOSAA.4.002395>  
39
- 40 Konkle, T., Brady, T. F., Alvarez, G. A., & Oliva, A. (2010). Conceptual distinctiveness supports  
41  
42 detailed visual long-term memory for real-world objects. *Journal of Experimental*  
43  
44 *Psychology: General*, *139*(3), 558–578. <https://doi.org/10.1037/a0019165>  
45  
46
- 47 Liu, Y., Gadepalli, K., Norouzi, M., Dahl, G. E., Kohlberger, T., Boyko, A., ... Stumpe, M. C.  
48  
49 (2017). Detecting Cancer Metastases on Gigapixel Pathology Images. *ArXiv:1703.02442*  
50  
51 *[Cs]*. Retrieved from <http://arxiv.org/abs/1703.02442>  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Loftus, G. R., & Masson, M. E. J. (1994). Using confidence intervals in within-subject designs.

4  
5 *Psychonomic Bulletin & Review*, 1(4), 476–490. <https://doi.org/10.3758/BF03210951>

6  
7 Luck, S. J., & Vogel, E. K. (1997). The capacity of visual working memory for features and

8  
9 conjunctions. *Nature*, 390(6657), 279–281. <https://doi.org/10.1038/36846>

10  
11  
12 Maniscalco, B., & Lau, H. (2012). A signal detection theoretic approach for estimating

13  
14 metacognitive sensitivity from confidence ratings. *Consciousness and Cognition*, 21(1),

15  
16 422–430. <https://doi.org/10.1016/j.concog.2011.09.021>

17  
18  
19 Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive field properties by

20  
21 learning a sparse code for natural images. *Nature*, 381(6583), 607–609.

22  
23 <https://doi.org/10.1038/381607a0>

24  
25  
26 Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: transforming numbers

27  
28 into movies. *Spatial Vision*, 10(4), 437–442. <https://doi.org/10.1163/156856897X00366>

29  
30  
31 Rensink, R. A. (2004). Visual Sensing Without Seeing. *Psychological Science*, 15(1), 27–32.

32  
33 <https://doi.org/10.1111/j.0963-7214.2004.01501005.x>

34  
35  
36 Serre, T., Oliva, A., & Poggio, T. (2007). A feedforward architecture accounts for rapid

37  
38 categorization. *Proceedings of the National Academy of Sciences*, 104(15), 6424–6429.

39  
40 <https://doi.org/10.1073/pnas.0700622104>

41  
42  
43 Simoncelli, E. P., & Olshausen, B. A. (2001). Natural image statistics and neural representation.

44  
45 *Annual Review of Neuroscience*, 24, 1193–1216.

46  
47 <https://doi.org/10.1146/annurev.neuro.24.1.1193>

48  
49  
50 Simons, D. J., & Chabris, C. F. (1999). Gorillas in Our Midst: Sustained Inattentive Blindness

51  
52 for Dynamic Events. *Perception*, 28(9), 1059–1074. <https://doi.org/10.1068/p281059>

- 1  
2  
3 Van Such, M., Lohr, R., Beckman, T., & Naessens, J. M. (2017). Extent of diagnostic agreement  
4 among medical referrals. *Journal of Evaluation in Clinical Practice*.  
5  
6  
7  
8 Van Wert, M. J., Horowitz, T. S., & Wolfe, J. M. (2009). Even in correctable search, some types  
9 of rare targets are frequently missed. *Attention, Perception, & Psychophysics*, *71*(3), 541–  
10 553. <https://doi.org/10.3758/APP.71.3.541>  
11  
12  
13  
14  
15 Vul, E., & Pashler, H. (2008). Measuring the crowd within: probabilistic representations within  
16 individuals. *Psychological Science*, *19*(7), 645–647. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-9280.2008.02136.x)  
17 9280.2008.02136.x  
18  
19  
20  
21 Wolfe, J. M., Brunelli, D. N., Rubinstein, J., & Horowitz, T. S. (2013). Prevalence effects in  
22 newly trained airport checkpoint screeners: Trained observers miss rare targets, too.  
23  
24  
25  
26 *Journal of Vision*, *13*(3), 33–33. <https://doi.org/10.1167/13.3.33>  
27  
28  
29 Wolfe, J. M., Horowitz, T. S., & Kenner, N. M. (2005). Cognitive psychology: rare items often  
30 missed in visual searches. *Nature*, *435*(7041), 439–440. <https://doi.org/10.1038/435439a>  
31  
32  
33 Wolfe, J. M., Horowitz, T. S., Van Wert, M. J., Kenner, N. M., Place, S. S., & Kibbi, N. (2007).  
34 Low target prevalence is a stubborn source of errors in visual search tasks. *Journal of*  
35 *Experimental Psychology. General*, *136*(4), 623–638. [https://doi.org/10.1037/0096-](https://doi.org/10.1037/0096-3445.136.4.623)  
36 3445.136.4.623  
37  
38  
39  
40  
41  
42 Wolfe, J. M., & Van Wert, M. J. (2010). Varying Target Prevalence Reveals Two Dissociable  
43 Decision Criteria in Visual Search. *Current Biology*, *20*(2), 121–124.  
44  
45  
46 <https://doi.org/10.1016/j.cub.2009.11.066>  
47  
48  
49 Wolfe, J. M., Võ, M. L.-H., Evans, K. K., & Greene, M. R. (2011). Visual search in scenes  
50 involves selective and nonselective pathways. *Trends in Cognitive Sciences*, *15*(2), 77–  
51 84. <https://doi.org/10.1016/j.tics.2010.12.001>  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Review Only