

## **Chapter Five:**

### **An Update of the Benton Facial Recognition Test**

Ebony Murray, Rachel Bennetts, Jeremy Tree and Sarah Bate

#### **Author Note**

Ebony Murray, Department of Psychology, Bournemouth University, UK.

Rachel Bennetts, Department of Life Sciences, Brunel University London, UK.

Jeremy Tree, Department of Psychology, Swansea University, UK.

Sarah Bate, Department of Psychology, Bournemouth University, UK. SB is supported by a British Academy Mid-Career Fellowship.

Correspondence concerning this report should be addressed to Ebony Murray, Department of Psychology, Faculty of Science and Technology, Poole House, Bournemouth University, Fern Barrow, Poole, BH12 5BB.

Email: [emurray@bournemouth.ac.uk](mailto:emurray@bournemouth.ac.uk)

## Abstract

The Benton Facial Recognition Test (BFRT) is a paper-and-pen task that has traditionally been used to assess face perception skills in neurological, clinical and psychiatric conditions. Despite some criticisms of its stimuli, the task enjoys a simple procedure and is rapid to administer. Further, it has recently been computerised (the BFRT-c), allowing reliable measurement of completion times, while addressing the need for online testing. Here, in response to calls for repeat-screening for the accurate detection of face recognition deficits, we present the BFRT-Revised (BFRT-r): a new version of the BFRT-c that maintains the task's basic paradigm, but employs new, higher quality stimuli that reflect recent theoretical advances in the field. An initial validation study with typical participants indicated that the BFRT-r has good internal reliability and content validity. A second investigation indicated that while younger and older participants achieved similar accuracy scores, completion times were longer in the latter, highlighting the need for age-matched norms when assessing clinical cases. Administration of the BFRT-r and BFRT-c to 31 individuals with developmental prosopagnosia identified 16 cases with impairments in face perception. While these deficits were observed on both tests in eight of the cases, eight others only displayed deficits on one of the two tasks, primarily on the task completion time measure. These findings are discussed in relation to current diagnostic screening protocols for face perception deficits. The BFRT-r is stored in an open repository and is freely available to other researchers.

**Keywords:** face perception; face matching; face recognition; prosopagnosia; Benton; response times

## An Update of the Benton Facial Recognition Test

The Benton Facial Recognition Test (BFRT: Benton & Van Allen, 1968; see Benton et al. (1983) for the formal reference of the test) is a face matching task that is traditionally administered face-to-face using hard copy materials. Participants are simultaneously presented with a target face above an array of six test faces. In the first six trials, one face in the array matches the identity of the target face, and in the final 16 trials, three faces in the array match the identity of the target. The task was originally developed for the assessment of individuals believed to have acquired prosopagnosia (a severe deficit in recognising familiar people from their face) following brain injury (Barton, 2008; Bate & Bennetts, 2015; Van Belle et al., 2011), but has since been widely used to assess face perception skills in a number of neurological, clinical and psychiatric conditions (Annaz et al., 2009; Rabin et al., 2005; Sachse et al., 2014).

Yet, the popularity of the BFRT has reduced in recent years, particularly for the assessment of individuals suspected to have prosopagnosia. At the turn of the century, many more people presented to researchers believing they experience a developmental form of prosopagnosia (Bate et al., 2008; De Luca et al., 2019; Geskin & Behrmann, 2018), prompting a wider individual differences perspective on human face recognition, and the belief that developmental face recognition difficulties may reside on a continuum (Bate & Tree, 2017; Barton & Corrow, 2016). These larger samples of cases have reignited long-standing questions of whether perceptual and mnemonic difficulties are dissociable (De Renzi et al., 1991), and whether subtypes of developmental prosopagnosia (DP) map onto this framework (note that the term “congenital” or “hereditary” prosopagnosia has been used somewhat interchangeably with DP in the literature: e.g. Behrmann & Geskin, 2005; Hasson et al., 2006; Kennerknecht et al., 2007; Palermo et al., 2011). Clearly, to address all these questions, reliable face perception

tasks are required. However, Duchaine and Weidenfeld (2003) reported that when the inner features of the faces in the BFRT were obscured, most typical participants could still achieve a typical score using the hairline and eyebrows alone. Further, Duchaine and Nakayama (2004) found that seven out of 11 DP participants achieved typical scores on the task, again suggesting that external facial cues may be used to aid performance.

Unfortunately, there is also deliberation over alternate tests of face perception, and the field still lacks a reliable task. The most widely used face perception test for the diagnosis of DP is the Cambridge Face Perception Test (CFPT: Duchaine et al., 2007), which presents participants with six morphed faces that are to be organised in order of similarity to a simultaneously presented target face. The task requires proficient use of a computer mouse within a strict time period, and the instructions are complex for online administration, particularly with clinical and older participants (Bate et al., 2018; Bate, Frowd et al., 2019; Bowles et al., 2009). Others query whether morphed faces are unnaturally similar (White et al., 2017), and whether the requirement for similarity judgements initiates higher-level cognitive processes than required for the simplistic identity matching of simultaneously presented naturalistic facial images (Rossion & Michel, 2018). Such simpler face matching tasks are typically found in the forensic face recognition literature (e.g. the Glasgow Face Matching Test: Burton et al., 2010; the Pairs Matching Test: Bate et al., 2018; Bate, Frowd et al., 2019), but are seldom used for the detection of DP due to their low sensitivity to poor performance. Indeed, the chance of responding correctly on all trials is 50% - a score that is within the range achieved by typical perceivers, many of whom find these tasks particularly challenging (e.g. Robertson et al., 2016; Shah et al., 2015). In addition, White and colleagues (2017) reported a response bias in DP participants, where the tendency to respond “different” in a simple same/different face matching task artificially inflated their score on these trials.

Such criticisms led Rossion and Michel (2018) to return to the BFRT, citing advantages in its original paradigm. Despite the external cues to recognition that were highlighted by Duchaine and colleagues (Duchaine & Weidenfeld, 2003; Duchaine & Nakayama, 2006), the BFRT has traditionally been regarded as a difficult test with no ceiling effect (Benton & Van Allen, 1972), that is quick to administer with simple instructions. Importantly, Rossion and Michel point out that decisional response biases are avoided by the task's forced choice procedure (the number of target faces is constant across test sections), making it substantially easier to interpret test scores. Further, Rossion and Michel (2018) highlighted the importance of recording task completion time in addition to accuracy (via a computerised version of the test: the BFRT-c), as a means to detect typical scores that are achieved by compensatory mechanisms. Previous work has adopted this approach when screening for face perception deficits in acquired prosopagnosia, where apparently typical accuracy scores were found to be accompanied by the use of atypical, laboured feature-by-feature matching strategies (e.g. Bukach et al., 2006; Busigny & Rossion, 2010; Delvenne et al., 2004; Farah, 1990; Young et al., 1993).

Another way to address this issue is to administer multiple versions of a face perception task, using rather different facial stimuli. This should prohibit, or at least reduce, the transfer of compensatory strategies that are useful in one version of a task. Indeed, some DPs report the use of particular facial features when images are captured within the same photography session (e.g. no change in skin tone or appearance of the hairline, even when images are cropped), or even consider pictorial cues from the images themselves (Adams et al., 2019). Further, very recent findings also highlight the importance of repeat-testing on key measures of face recognition performance when screening for DP (Bate, Bennetts, Gregory et al., 2019; Murray & Bate, 2020), given issues with task reliability, the occurrence of borderline scores that are

difficult to reconcile, and the possibility that a particular score simply occurred by chance performance (Young et al., 1987).

Yet, no known alternate version of the BFRT exists, and an update of the task using new stimuli is certainly overdue. While the basic paradigm (with the monitoring of completion times) offers a sound means of assessing face perception, the age of the test unsurprisingly lends itself to low image quality. Whilst face recognition can be successful even when images are of low spatial frequency (e.g. Liu et al. 2000), unfamiliar face processing, as is being assessed with the BFRT, benefits from high-quality images (Burton et al., 1999). Moreover, the findings of Duchaine and colleagues (Duchaine & Weidenfeld, 2003; Duchaine & Nakayama, 2004) indicate that extra-facial cues can be used to achieve a typical score on the existing version of the test. While more recent face-processing tasks in the neuropsychological literature have responded to Duchaine and colleagues' criticisms of the BFRT by using tightly controlled images that are captured on the same day and heavily cropped to exclude the external features (e.g. Duchaine & Nakayama, 2006; Biotti et al., 2017; Esins et al., 2016), it has been argued that this procedure actually distances the task from real-world face recognition (Burton, 2013). Rather, variability in facial appearance is a critical feature of everyday face recognition, and should be embraced in, rather than removed from, laboratory tests (Young & Burton, 2017, 2018). In fact, even typical participants struggle to match faces of the same identity when pictured in more "ambient" images that retain the external features of the face, given image-based cues cannot be used as compensatory cues for successful performance (for further discussion, see Burton, 2013).

Here, we introduce a new version of the BFRT-c, the BFRT-revised (BFRT-r), that maintains the format of the original task but employs new, more varied, naturalistic facial images. In Experiment 1, we examine the validity of the BFRT-r in typical participants and

provide norming data for comparison to clinical cases. In Experiment 2, we assess the test's diagnostic utility alongside the BFRT-c in DP.

## Experiment 1

A new face matching task (the BFRT-r) was created that follows the original BFRT paradigm, but is computerised (akin to the BFRT-c) and uses new, more ambient facial images. We initially assessed the psychometric properties of the task and collected norming data from young typical adults. Content validity was assessed using an inverted version of same task.

### Method

#### *Participants*

A total of 165 participants aged 18-35 years took part in Experiment 1. One hundred and nine participants (M age = 24.7 years, SD = 3.5; 55 female) completed the BFRT-c and BFRT-r, but not the inverted version of the BFRT-r (to avoid re-exposure effects). Fifty-six different participants (mean age = 24.9 years, SD = 3.5; 27 male) completed only the inverted version of the BFRT-r. All participants were recruited via the online participant recruitment website *Prolific*, in exchange for a small financial incentive. All were Caucasian and lived within the UK, reported no history of socio-emotional, neurological or psychiatric disorder, and had normal or corrected-to-normal vision. This project was approved by the institutional Research Ethics Committee.

#### *Materials*

*BFRT-c* (Rossion & Michel, 2018): The BFRT-c is the original version of the BFRT, in a computerised format. The test contains a total of 22 trials in which an unfamiliar Caucasian

target face (shown from a frontal viewpoint with a neutral expression) has to be found among a simultaneously-presented array of six Caucasian probe faces, also showing neutral expressions. For the first six trials (half male), the target face has to be found only once within each array, where all faces are shown from a frontal viewpoint, such that the corresponding probe image is very similar to the target image. For the remaining 16 trials (half male), the target face is again presented from a frontal viewpoint. The participant is required to find three images within the six-image array that match the identity of the target. The six faces in each array vary either in terms of head orientation (the second section of the test: eight trials, half female) or lighting (the third section of the test: eight items, half female). Some target faces are repeated: four of the seven female targets appear in two separate sections, one of the seven male targets appears in all three sections, and three male targets are used in two sections. All target identities are also used as distractors in at least one trial of the task.

In each trial, target faces are presented at a slightly different size than those in the array (target faces were 156 x 232 pixels; faces in the array were 201 x 234 pixels, in order to minimise successful matching based on low-level, image-based visual cues: Rossion & Michel, 2018). All images are grayscale and display the overall shape of the face, but are cropped below the chin and beyond the hairline. As in the original version of the task, the order of the trials is not randomised and participants have an unlimited length of time to complete each trial. There is an inter-stimulus interval of 800ms. Information screens at the beginning of each section instruct the participant how many responses to make for each trial, and inform them that response time is recorded.

Participants are required to select their responses by clicking on the appropriate face(s) in each array. For trials that require three responses, participants are able to select faces in any order, but cannot change a response once a face has been selected. The maximum score on the task is 54. Participants can receive one point in each of the six trials that compose Section 1



(where one response is required per trial), and between 0 and 3 points for each of the trials in Section 2 (where three responses are required per trial). Trial completion times are measured to aid data processing (see below), and overall task completion times are monitored for analysis.

*BFRT-r*: The basic paradigm of the BFRT-c is retained, with the same number of trials. However, the facial stimuli are replaced throughout. As gender biases have been shown for the recognition of female but not male faces (e.g. Herlitz & Lovén, 2013; Lovén et al., 2011), we followed the precedent of more recent tests by only using male faces (e.g. the Cambridge Face Memory Test, CFMT: Duchaine & Nakayama, 2006; the CFPT: Duchaine et al., 2007). We initially acquired facial images from a total of 130 Caucasian males (aged 18-34 years:  $M = 21.9$  years,  $SD = 3.2$ ) in exchange for course credit or a small financial incentive. Images were captured within the laboratory, and/or were existing photographs provided by the participant that had been taken within the space of a single year. Thus, different images of the same person had been captured on different days, often months apart, and in many cases, using different cameras. However, images of the same person had all been captured within the same year, preventing any major ageing effects. Blemishes, skin tone and hairstyle varied from image to image, as well as lighting conditions. No image had been manipulated, and all were of sufficiently high quality (no less than 96 DPI). They displayed the target without spectacles. There were variations in viewpoint due to their capture in naturalistic settings.

A unique target was used in each of the 22 trials, and no target was re-used as a distractor. Ten distractor identities were repeated over the 22 trials, but different images of each individual were used where possible; only two images were repeated twice through the test. No distractor identity was repeated in the same array. Distractors were allocated to each trial based on their perceived similarity to the target, as judged by a member of the research team. Pilot testing results supported the judgements made by said member of the research team; the

trials included in the final BFRT-r did not elicit ceiling nor floor effects. In total, the test used images from 76 different individuals.

All images were presented in greyscale. This decision was made based on the pilot testing/materials analysis, which indicated that ceiling effects in the typical population could be achieved when images were in colour. To prevent low-level image matching, target faces were not cropped to exclude any part of the head, hair or ears, but array images were cropped around the hairline (see Figure 1). Target images were larger (166 x 232 pixels) than those in the array (approximately 153 x 200 pixels). As in the BFRT-c, only one of the array faces matched the identity of the target in the first six trials, and three in the remaining 16 trials. In the first 12 trials of the task, all faces are displayed from frontal viewpoints. In the final 10 trials, faces are displayed from frontal, but more naturalistic, viewpoints. The rotation of most faces is small; approximately 10-30 degrees to the left or right. A small number of images (N = 7) are displayed at a larger rotation (less than 45 degrees), but the whole face can be viewed in every photograph (i.e. both eyes are clearly visible; see Figure 1).

As for the BFRT-c, trials were presented in the same order for each participant, with an inter-stimulus interval of 800ms, and responses were made and scored in the same manner. Instructions were identical to the BFRT-c, but additionally informed participants that some images were taken some time apart, and some aspects of the target's appearance (e.g. hairstyle) may have changed during this time. The BFRT-r test materials are available in an open repository: [https://osf.io/vza3m/?view\\_only=404f6d1971924759b126d46cba1d25b7](https://osf.io/vza3m/?view_only=404f6d1971924759b126d46cba1d25b7). A fully programmed version can also be shared with researchers on request, via Testable. The test and its materials are protected by a Creative Commons Attribution-Non Commercial license.

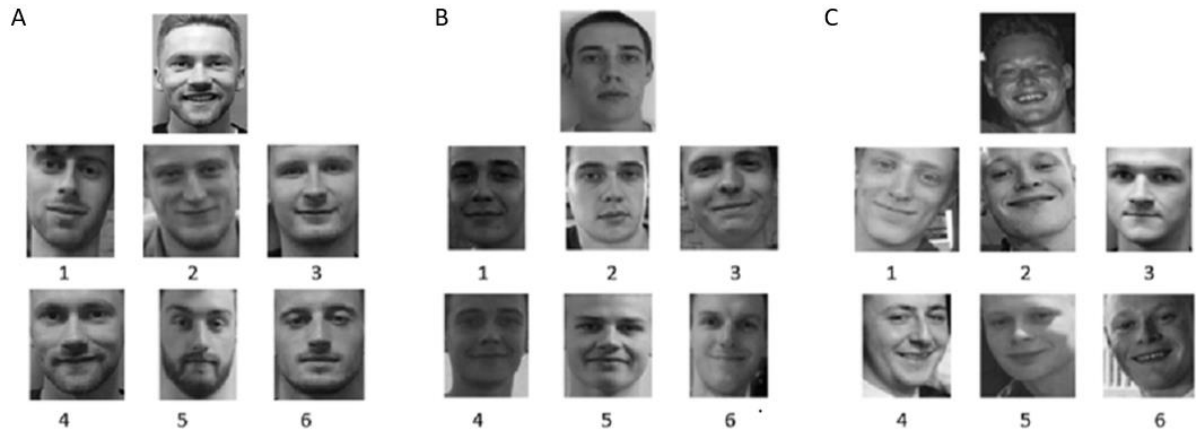


Figure 1. Example trials from the BFRT-r. Panel A shows the trial format for trials 1-6 where all faces are presented from a frontal viewpoint. There is only one correct response (4). Panel B shows the trial format for trials 7-12 where all faces are again shown from a frontal viewpoint, but there are three correct responses (1, 2, 4). Panel C shows the format for trials 13-23, where the target face is shown from a frontal viewpoint and the probe images show some rotation. There are three correct responses (2, 5, 6).

*BFRT-r inverted:* The BFRT-r was also prepared in an inverted format, to assess the content-validity of the test. All stimuli and parameters were identical, with the exception that all images were rotated 180 degrees.

### *Procedure*

All tasks were completed online, using the *Testable* platform ([www.testable.org](http://www.testable.org); see Rezlescu et al., 2020). Participants were required to initially calibrate the tests for screen size, ensuring uniform presentation. The 109 participants that completed the main string of tests completed the BFRT-r first, then the BFRT-c. This enabled us to collect accurate norming data for the new task without introducing practice effects from the repeated use of the same paradigm, or

allowing the possibility of testing fatigue. The 56 participants who only took part in the inverted version of the BFRT-r did not complete any other tests.

### *Data Processing*

As data were collected online, responses were initially screened for task engagement. Each individual's mean response time (and SD) was calculated for each group of trials in each task (i.e. the trials which required one response, and the trials which required three responses), and any responses which were greater than 3 SDs above the mean were removed, as were any responses that were quicker than 150ms. In addition, trials that required three unique responses were screened to ensure correct completion (i.e. to remove trials that had received duplicate responses). If more than 33% of trials for the same test were excluded, the participant's score for the overall task was removed from the final dataset. Overall accuracy scores were also screened for outliers across the dataset, using a three SDs from the mean criterion.

For the participants that completed both the BFRT-r and BFRT-c, eight were removed for failing to complete enough trials, and one for achieving accuracy scores on both tasks that surpassed three SDs from the mean score. No participant responded quicker than 150ms on any trial, and no participant was excluded for giving too many abnormally slow responses. The final sample consisted of 100 participants aged 18–30 years (49 male; M age = 24.8 years, SD = 3.5). The same screening procedures were applied to the participants who only completed the inverted version of the BFRT-r, resulting in the exclusion of one individual. A final sample of 55 participants aged 18–30 years (27 male; M age = 24.9 years, SD = 3.6) therefore proceeded to the analysis phase.

## Results

Mean accuracy performance on the upright version of the BFRT-r was 78.24% (SD = 9.20). Given different numbers of responses were required in different sections of each test, we followed the precedent of Rossion and Michel (2018) in analysing overall task completion times, rather than the average response time per trial. The mean overall task completion time for the BFRT-r was 251.22 seconds (SD = 128.08).

Mean accuracy performance on the BFRT-c was 80.44% (SD = 8.52), and mean task completion time was 165.51 seconds (SD = 69.58 seconds). The performance of this sample on the BFRT-c is therefore comparable to the norms presented by Rossion and Michel (2018), who reported a mean accuracy of 82.98% (SD = 6.37), and a mean task completion time of 180.85 seconds (SD = 59.86).

Accuracy performance on the BFRT-r strongly correlated with the BFRT-c ( $r = .657$ ,  $p < .001$ ). To fully compare the two tasks, a mixed 2 (test: BFRT-c, BFRT-r) x 2 (gender: male, female) ANOVA was carried out. There was a main effect of test,  $F(1,98) = 8.871$ ,  $p = .004$ ,  $\eta^2 = .083$ , in that individuals performed significantly better on the BFRT-c (M = 80.44%, SD = 8.52) than the BFRT-r (M = 78.24%, SD = 9.20; see Table 1). There was also a significant main effect of gender over the two tests,  $F(1,98) = 4.818$ ,  $p = .031$ ,  $\eta^2 = .047$  (see Figure 2), where females (M = 81.05%, SD = 8.31) outperformed males (M = 77.57%, SD = 9.11). However, there was no significant interaction between test and gender,  $F(1,98) = 2.198$ ,  $p = .141$ .

A 2 (test: BFRT-c, BFRT-r) x 2 (gender) ANOVA on task completion times revealed a main effect of test,  $F(1,98) = 101.714$ ,  $p < .001$ ,  $\eta^2 = .509$ , with participants taking longer to complete the BFRT-r (M = 251.22s, SD = 128.08) than the BFRT-c (M = 165.51s, SD = 69.58; see Table 1). There was no main effect of gender,  $F(1,98) = 0.424$ ,  $p = .517$  (see Figure 3), nor an interaction between test and gender,  $F(1,98) = 0.028$ ,  $p = .868$ .

Following the precedent of Rossion and Michel (2018), the task's internal reliability was assessed by correlating performance on even versus odd items, considering only the second part of the test in which three responses are made per trial. The inter-item correlation was significant for accuracy rates (mean score for the eight even items = 18.59/24, SD = 2.45; mean score odd items = 20.30/24, SD = 2.43;  $r_{SB}$  [Spearman–Brown] = .775,  $p < .001$ ). The interitem correlation was even higher for trial completion times (mean trial completion times for the eight even items = 96.70s, SD = 55.58s; mean trial completion times for the eight odd items = 88.99s, SD = 47.15s;  $r_{SB} = .967$ ,  $p < .001$ ).

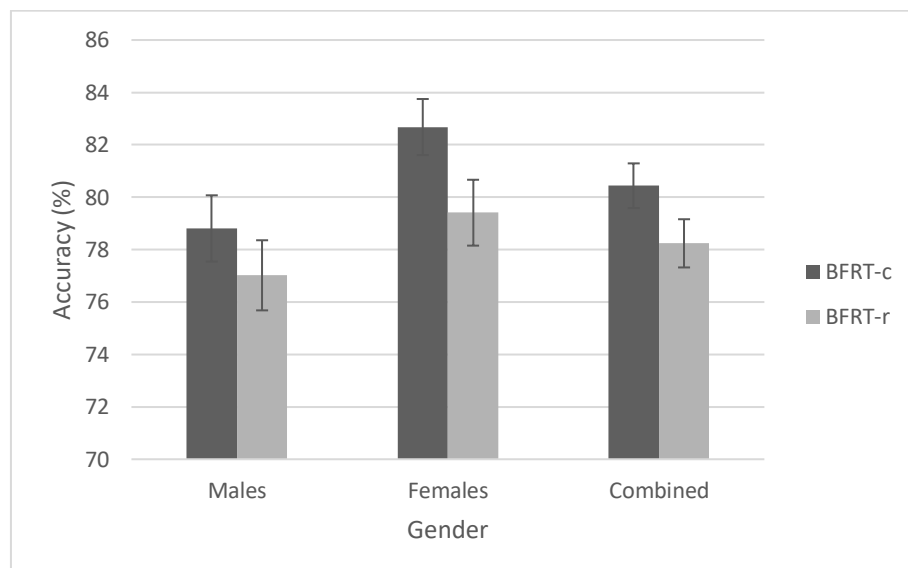


Figure 2: Mean accuracy on the BFRT-r and BFRT-c for males and females separately, and the overall sample. Bars represent standard error.

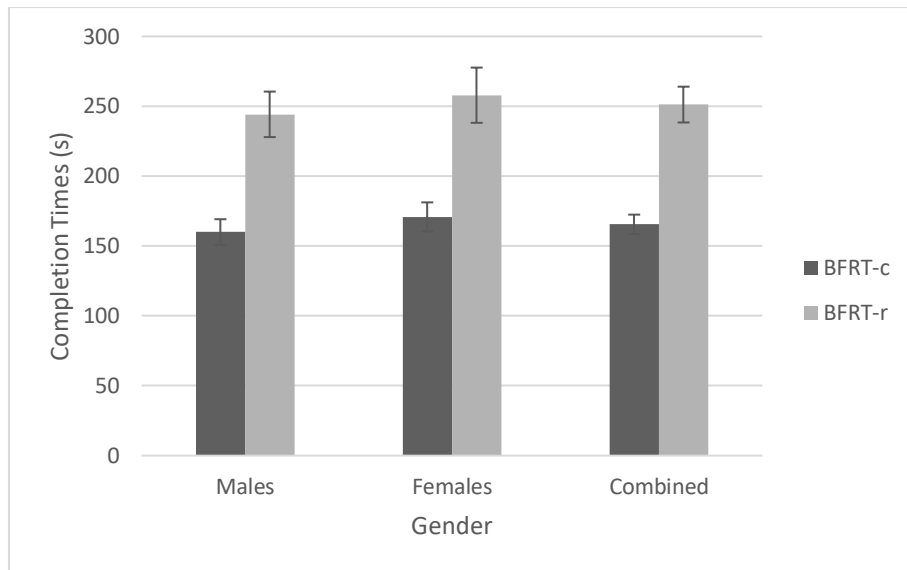


Figure 3: Mean task completion times on the BFRT-r and BFRT-c for males and females separately, and overall sample. Bars represent standard error.

Table 1. Descriptive data (means and standard deviations) for the upright versions of the BFRT-c and BFRT-r for younger controls (Experiment 1; M age = 24.8), and older controls and DPs (Experiment 2; M age = 48.2 and 51.3 respectively). Accuracy is presented as a percentage, and completion times in seconds.

	<b>BFRT-c accuracy</b>	<b>BFRT-c completion times</b>	<b>BFRT-r accuracy</b>	<b>BFRT-r completion times</b>
<b>Younger controls</b> (N = 100)	80.44 (8.52)	165.51 (69.58)	78.24 (9.20)	251.22 (128.08)
<b>Older controls</b> (N = 129)	82.23 (9.88)	241.07 (93.92)	77.59 (11.19)	341.93 (159.63)
<b>DPs (N = 31)</b>	74.91 (7.97)	346.69 (139.28)	67.68 (7.96)	495.95 (202.45)

Comparison between overall accuracy scores on the upright ( $M = 78.24\%$ ,  $SD = 9.20$ ) and inverted ( $M = 56.67\%$ ,  $SD = 10.64$ ) versions of the BFRT-r revealed a substantial inversion effect,  $t(153) = 13.198$ ,  $p < .001$ ,  $d = 2.17$ . However, there was no significant difference between upright BFRT-r ( $M = 251.22s$ ,  $SD = 128.08$ ) task completion times and inverted task completion times ( $M = 249.90$  seconds,  $SD = 146.88$ ),  $t(153) = 0.058$ ,  $p = .954$ .

### *Summary*

Here, we present the BFRT-r: a new test of face perception that adopts the same paradigm as the original BFRT (as per the BFRT-c), but uses more naturalistic images to accommodate within-person variation in facial images. As the BFRT-r follows the procedure of the BFRT-c, the test continues to be simple and quick to administer, with an approximate completion time of four minutes in typical young adults. Initial analyses reveal that the BFRT-r has good internal reliability with strong inter-item correlations. It also has a strong inversion effect according to accuracy (although not completion times), suggesting that it taps face- rather than image-processing mechanisms.

Comparison of performance on the two tasks indicates that the BFRT-r is more difficult than the original version. Importantly, typical participants are able to score well above chance on the BFRT-r (the lowest score was 55.56%; chance performance is 46.30%). Further, the norming data reported here ( $M = 78.24\%$ ,  $SD = 9.20$ ) would enable clinical participants to score two SDs below the mean without performing at chance level. Indeed, those with DP often show impaired face perception skills that are not completely abolished, but these skills are not completely abolished to the point that they are scoring at chance level. Thus, a suitable task requires this availability of scores (e.g. Bate et al., 2019; Biotti et al., 2019; Righart & Gelder, 2007).



## Experiment 2

Having explored the validity of the BFRT-r in younger participants, we next sought to examine the diagnostic utility of the updated version in individuals with DP. In particular, we examine (a) the additional benefit of evaluating response times as well as accuracy in atypical participants, and (b) whether there is a case for administration of multiple versions of the same task when screening for face perception deficits.

### Method

#### *Participants*

Thirty-one participants with a prior diagnosis of DP took part in this study. They had previously taken part in an objective screening session and scored atypically on at least two of three diagnostic tests: the CFMT (Duchaine & Nakayama, 2006), the CFPT (Duchaine et al., 2007), and a famous faces test (e.g. Bate, Bennetts, Gregory et al., 2019; Bennetts et al., 2015; Murray & Bate, 2019), following existing diagnostic protocols (Dalrymple & Palermo, 2015; Bate & Tree, 2017: see supplementary material for their diagnostic results). Eight were male, and they were aged between 40 and 59 years ( $M = 51.3$  years,  $SD = 5.6$ ).

Because our DP sample were older than the younger adults reported in Experiment 1, a new set of 138 older control participants ( $M$  age = 48.0 years,  $SD = 5.7$ ; 68 females) were recruited for age-matched comparison. These individuals were again recruited via the *Prolific* online recruitment platform, in exchange for a small financial incentive. All DP and control participants were Caucasian, and reported no history of socio-emotional, neurological or psychiatric disorder (including mild cognitive impairment), and had normal or corrected-to-normal vision.

Following the same data-processing strategies as in Experiment 1, data from nine control participants were removed: eight failed to elicit the required number of responses on more than 33% of the trials on the BFRT-r, and an additional participant took an abnormally long time to complete both the BFRT-r and BFRT-c. This resulted in a final sample of 129 (64 female) control participants, aged between 40 and 60 years ( $M = 48.2$  years,  $SD = 5.7$ ). The same exclusion criteria were applied to the DP data as for the control participants; no DP data were removed from the analysis.

### *Materials and procedure*

All participants completed the BFRT-r and BFRT-c in that order, online, via the testing platform *Testable*.

## Results

### *Age and gender*

For the new older controls, high correlations were observed between performance on the two versions of the Benton on both the accuracy ( $r = .753$ ,  $p = .001$ ) and task completion time ( $r = .866$ ,  $p = .001$ ) measures. Further, BFRT-r accuracy performance did not differ between the new set of older control participants ( $M = 77.59\%$ ,  $SD = 11.19$ ) and the younger sample reported in Experiment 1 ( $M = 78.24\%$ ,  $SD = 9.20$ ),  $t(227) = 0.470$ ,  $p = .148$ . However, overall task completion times were slower in older ( $M = 341.93$  seconds,  $SD = 159.63$ ) compared to younger ( $M = 251.22$  seconds,  $SD = 128.08$ ) controls,  $t(227) = 4.641$ ,  $p = .022$ ,  $d = .63$  (see Figure 4). The same pattern emerged for the BFRT-c: younger ( $M = 80.44\%$ ,  $SD = 8.52$ ) and older ( $M = 82.23\%$ ,  $SD = 9.88$ ) controls performed similarly in terms of accuracy,  $t(227) = 1.437$ ,  $p = .116$ , but younger controls ( $M = 165.51$  seconds,  $SD = 69.58$ ) completed the test significantly faster than older controls ( $M = 241.07$  seconds,  $SD = 93.92$ ),  $t(227) = 6.738$ ,  $p =$

.012,  $d = .91$  (see Figure 5). Thus, subsequent analyses only compared the performance of DPs to the older control group. No gender effects were found on either the BFRT-r or BFRT-c in this age group ( $ps > .05$ ).

#### *DP performance: Group analyses*

A mixed 2 (test: BFRT-c, BFRT-r) x 2 (group: DP, older controls) ANOVA was conducted to explore overall group differences in accuracy scores (see Table 1). There was a significant main effect of group, whereby DP participants scored significantly poorer ( $M = 71.29\%$ ,  $SD = 7.96$ ) than control participants ( $M = 79.91\%$ ,  $SD = 10.54$ ),  $F(1,158) = 21.03$ ,  $p < .001$ ,  $\eta^2 = .088$  (see Figure 4). There was also a main effect of test,  $F(1,158) = 61.94$ ,  $p < .001$ ,  $\eta^2 = .282$ : scores on the BFRT-c ( $M = 80.81\%$ ,  $SD = 9.95$ ) were higher than those on the BFRT-r ( $M = 75.67\%$ ,  $SD = 11.33$ ). There was no significant interaction between test and group,  $F(1,158) = 2.96$ ,  $p = .088$ .

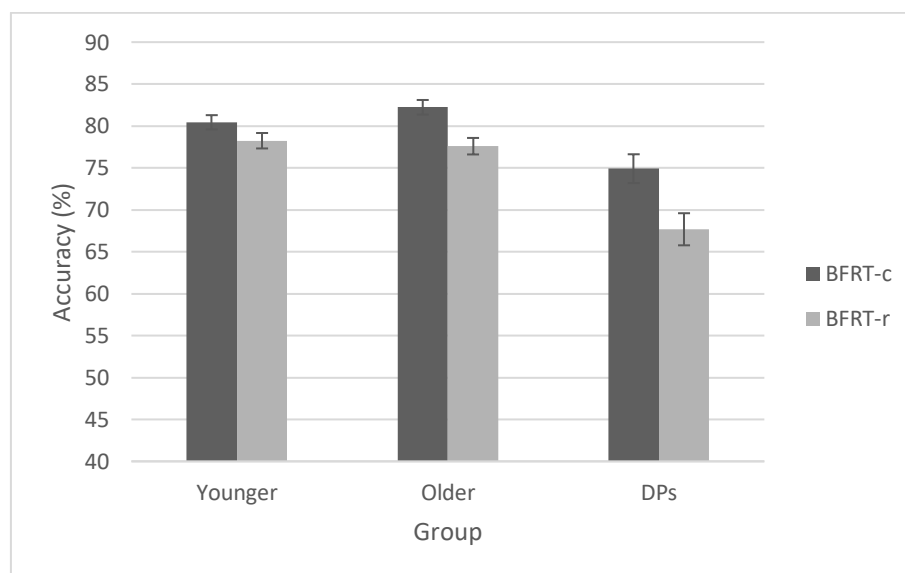


Figure 4: Mean accuracy on the BFRT-c and BFRT-r for younger control participants, older control participants, and the DP group. Bars represent standard error.

To investigate any differences in task completion times, a 2 (test: BFRT-c, BFRT-r) x 2 (group: DP, controls) ANOVA was conducted (see Table 1). A significant main effect of group indicated that DPs took longer to complete the tests ( $M = 421.32$  seconds,  $SD = 170.86$ ) than controls ( $M = 291.50$  seconds,  $SD = 126.78$ ),  $F(1,158) = 24.51$ ,  $p < .001$ ,  $\eta^2 = .134$  (see Figure 5). A significant main effect of test indicated that participants took longer to complete the BFRT-r ( $M = 371.77$  seconds,  $SD = 178.82$ ) than the BFRT-c ( $M = 261.53$  seconds,  $SD = 111.87$ ),  $F(1, 158) = 159.89$ ,  $p < .001$ ,  $\eta^2 = .503$ . There was also a significant interaction between test and group,  $F(1, 158) = 5.99$ ,  $p = .016$ ,  $\eta^2 = .037$ . Pairwise comparisons revealed that DPs took significantly longer than controls to complete both the BFRT-c ( $M = 346.69$ ,  $SD = 139.28$  and  $M = 241.07$ ,  $SD = 93.92$  respectively:  $p < .001$ ) and BFRT-r ( $M = 495.95$ ,  $SD = 202.45$  and  $M = 341.93$ ,  $SD = 178.82$  respectively:  $p < .001$ ).

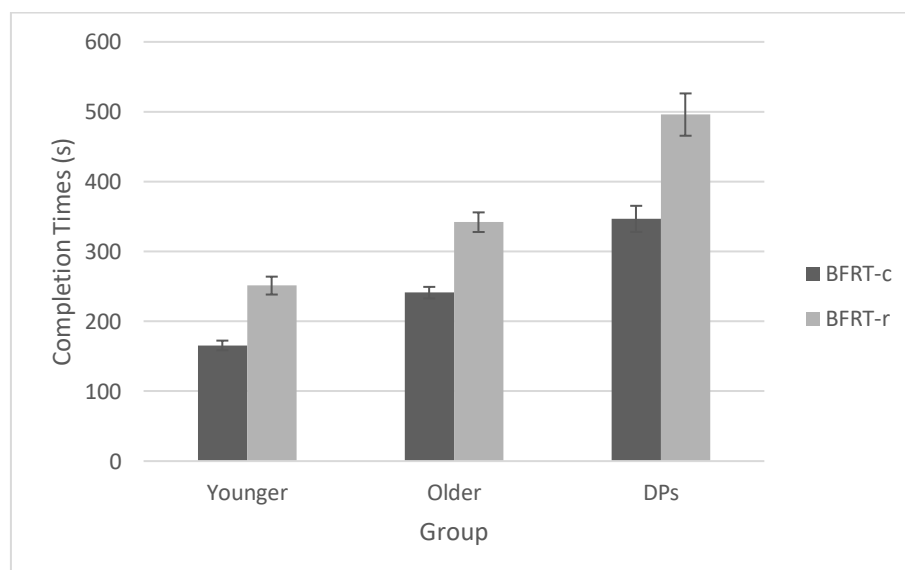


Figure 5: Mean task completion times on the BFRT-c and BFRT-r for younger control participants, older control participants, and the DP group. Bars represent standard error.

Because we also held CFPT upright accuracy scores for all our DPs as part of their background diagnostic profiles (see supplementary material), we were able to investigate

whether accuracy or completion time on both versions of the BFRT were associated with this indicator. However, no significant correlations were observed between CFPT performance and any of the four BFRT measures (all  $ps > .490$ ). Given the vast differences in paradigm, this is perhaps unsurprising. On the other hand, accuracy ( $r = 0.538, p = .002$ ) and completion time ( $r = 0.788, p = .001$ ; sequential Bonferroni correction for multiple correlations applied) on the two versions of the BFRT were highly correlated in DP participants, as was found for control participants in both experiments. Because of the lack of association between CFPT and BFRT scores, we did not proceed to use CFPT scores to further interpret individual patterns of performance on either version of the BFRT (see below). Indeed, it is not possible to infer from the current methodology whether either the CFPT or BFRT offers a “true” indicator of face perception, and we instead focus on consistency of individual performance across the two versions of the BFRT.

### *Single-case Analyses*

To examine the importance of assessing both accuracy and response times on face matching tests, each DP’s performance on the BFRT-r was evaluated on both parameters on a case-by-case basis (see Table 2). As all participants were over the age of 40, their scores and completion times were compared to that of the older control group (see Table 1). The  $z$ -score used as a cut-off for typicality varies within the DP literature, with some authors using two SDs from the mean (Biotti et al., 2019; Bowles et al., 2009; Bate, Bennetts, Tree et al., 2019) and others 1.7 SDs (DeGutis et al., 2012, 2014; Palermo et al., 2017; White et al., 2016). Here, to allow for recording error, and to err on the conservative side when determining face perception is intact (given it is currently assumed that the process is impaired in most DPs: Bate, Bennetts, Gregory et al., 2019; Biotti et al., 2019), we present the findings in terms of a 1.7 SD cut-off.

Fifteen of the 31 DPs (48.39%) performed within the typical range on both the BFRT-r and BFRT-c according to both accuracy and completion time measures (see Table 2). Notably few borderline scores were detected: the closest score to a cut-off was a  $z$ -score of -1.64 on the BFRT-r (DPM06), with the vast majority of other scores occurring within 1.25 SDs of the control mean. Of the 16 DPs who showed at least some impairment, eight exceeded cut-off on both tasks, according to at least one measure. An additional five participants were only impaired on the BFRT-r, and three participants were only impaired on the BFRT-c. Strikingly, only three DPs displayed impairments on accuracy alone, whereas 12 DPs only showed impairments on the completion time measure. Thus, task completion time was the primary indicator of impairment on both versions of the test.

Table 2. Normalised accuracy scores and task completion times for the 31 DP participants on the BFRT-r and BFRT-c. Note that negative  $z$ -scores represent poorer performance for accuracy, and positive scores indicate slower completion times.

<b>Participant ID</b>	<b>BFRT-r Accuracy</b>	<b>BFRT-r Completion Time</b>	<b>BFRT-c Accuracy</b>	<b>BFRT-c Completion Time</b>
DPM01	-1.14	0.34	-1.39	0.87
DPM02	0.02	-0.81	0.30	-0.96
DPM03	-0.31	-0.12	-0.45	1.07
DPM04	-1.97*	-0.06	-0.64	-0.24
DPM05	-0.31	-0.45	-1.01	0.79
DPM06	-1.64	2.11*	-0.64	1.54
DPM07	-0.81	-0.01	-1.39	-0.52
DPM08	-1.64	2.55*	-0.64	4.31 *

DPF01	-0.98	0.54	-1.39	1.02
DPF02	-0.65	0.82	0.11	1.44
DPF03	-0.48	0.86	-0.45	1.64
DPF04	-0.65	2.59 *	0.30	1.99 *
DPF05	-0.65	2.70 *	-0.26	4.66 *
DPF06	-0.98	1.42	-0.45	3.55 *
DPF07	-1.97 *	-1.41	-3.45 *	-1.56
DPF08	-1.14	2.74 *	-1.39	2.04 *
DPF09	-0.31	1.94 *	0.08	2.25 *
DPF10	-2.13 *	2.07 *	-1.20	1.53
DPF11	-0.31	1.11	-0.82	1.79 *
DPF12	-0.65	-0.15	-1.39	-0.08
DPF13	0.51	2.22 *	0.30	1.30
DPF14	0.35	3.06 *	-0.45	1.14
DPF15	-0.98	1.34	-1.39	0.36
DPF16	-1.14	0.21	0.11	-0.56
DPF17	-2.30 *	0.24	-2.14 *	-0.48
DPF18	-1.31	1.31	-0.64	2.10 *
DPF19	-1.14	3.24 *	-0.83	3.29 *
DPF20	0.18	-0.09	-0.64	0.15
DPF21	-0.98	-0.04	0.49	-0.09
DPF22	-1.47	-0.11	-1.01	0.23
DPF23	-0.48	-0.29	-0.45	0.27

\* denotes an atypical z-score (+/- 1.7)

### *Summary*

This study examined performance on the BFRT-r and BFRT-c in DPs and matched older adult control participants. Comparison of the control data to the younger participants tested in Experiment 1 indicates that accuracy performance is consistent across the two age groups, but older participants took longer to complete both tasks. Thus, performance of atypical participants needs to be compared to an age-matched control group. While a very small gender effect was found for younger controls in Experiment 1, this did not emerge for the older controls in this Experiment, and we therefore recommend for simplicity that gender-specific norms are not required for any age group in either version of the test. Consistent with the findings of Experiment 1, all participants found the BFRT-r more challenging than the original version, in terms of both accuracy and completion time.

As a group, the DPs performed more poorly than controls on both the BFRT-r and BFRT-c according to both accuracy and completion time measures (see Table 1). However, akin to previous work (Bate, Bennetts, Gregory et al., 2019; Burns et al., 2017; Le Grand et al., 2006; Minnebusch et al., 2007), case-by-case analyses indicated considerable heterogeneity in DPs' face perception performance, with approximately half the sample displaying intact face perception skills on both tests. Consistent deficits in face perception were noted across both versions of the test in eight DPs. However, only the BFRT-r detected impairment for five DPs, and only the BFRT-c for the remaining three, suggesting some cases may be missed by administration of a single face perception task (see also Murray & Bate, 2020). Strikingly, the importance of monitoring completion time was observed in 12 individuals who achieved typical accuracy scores in both tests, but slow completion times on at least one version of the task.



## General Discussion

In this paper, we introduce a new version of the BFRT (the BFRT-r), with updated stimuli that address recent theoretical progress in the face recognition literature. We sought to examine the validity of the BFRT-r in typical participants, and provide norming data for younger (aged 18–35 years) and older (aged 40–60 years) adult populations. Although our control samples did not include individuals aged 36–39 years, patterns observed in previous work (e.g. Bate, Bennetts, Gregory et al., 2019; Bowles et al., 2009) suggest clinical cases within this age range should be compared to the younger control group. That is, individuals aged 35–39 years perform similarly to those aged 35 and younger on face processing tasks. Finally, we investigated the test’s utility in identifying face perception difficulties in DP.

First, we replicated several known advantages of the BFRT. As the BFRT-r procedure is identical to the computerised version of the original BFRT (the BFRT-c: Rossion & Michel, 2018), the test is known to follow a simple procedure and is quick to administer. Here, we found that the BFRT-r takes approximately four to six minutes for typical participants to complete (with longer completion times in older adults), and approximately eight minutes for older adults with DP. Further, we noted the particular importance of monitoring task completion times for accurate diagnosis of face perception impairments. This is clearly facilitated by use of a computerised rather than pencil-and-paper format (see Rossion & Michel, 2018), and also lends the task to online administration – a particularly important concern in very recent times. Notably, online administration was used here in both experiments, and the resulting strong internal reliability and inter-item correlations directly support this mode of implementation. Moreover, the BFRT-r elicits a strong inversion effect, evidencing content-validity akin to other tests of face processing (Busigny & Rossion, 2010; Duchaine et al., 2007; Duchaine & Nakayama, 2006; McKone et al., 2011).

It should be noted that no measure of BFRT performance was found to be associated with CFPT accuracy scores in the DP sample. This finding warrants further investigation given both paradigms fall under the umbrella of “face perception tasks”, yet it seems likely that they tap rather different aspects of face perception performance. This could potentially have important implications for DP screening protocols, and, pending further investigation, it may be prudent for researchers to administer all three tests to their DP participants. This protocol would allow for face perception impairments to be screened across different paradigms, while the two versions of the BFRT will allow consistency of performance to be more closely monitored.

The main adaptation of the new BFRT-r concerns the new images, both in terms of visual quality, and in addressing important theoretical concerns within the field. While the original BFRT images were highly constrained and were presumably captured in the same setting on the same day and using the same camera, our images embraced the natural variability which typically occurs when viewing the same person on different occasions in everyday life. The photographs were captured over different days (sometimes months apart), using a variety of cameras, showing the person from varying viewpoints and distances from the camera, and in different lighting conditions. By moving away from the tightly controlled conditions that prevail in existing tests of face perception, it is likely that we move closer towards the circumstances of everyday face perception, providing a more ecologically valid diagnostic test (Burton, 2013). In addition, the use of more ambient facial images also overcomes previous concerns that extra-facial or distinguishing features could be used by clinical participants to achieve typical scores on the BFRT (Duchaine & Weidenfeld, 2003; however, it is noted that there needs to be further research to explore whether participants use the (minimal) extra-facial information and/or eyebrows when completing the BFRT-r). Given both controls and DPs found the new version of the BFRT to be more difficult than the original version, the new

stimuli have likely gone some way towards addressing this issue. While it could be argued that participants performed better on the BFRT-c than the BFRT-r due to practice effects with the paradigm (the BFRT-r was always completed first), this seems unlikely as our BFRT-c norms are strikingly similar to those reported by Rossion and Michel (2018).

Importantly, our data also indicated that face perception skills are heterogeneous in DP – a factor that has been highlighted in previous work (Bate, Bennetts, Tree et al., 2019; Burns et al., 2017; Le Grand et al., 2006; Minnebusch et al., 2007). Here, we found that approximately half of our DP sample presented with no impairments on either version of the BFRT. While eight of the remaining 16 DPs consistently displayed deficits on both versions of the BFRT, five were only detected on the BFRT-r and three by the BFRT- c (note that we did not attempt to further clarify these patterns using CFPT scores given the lack of association between the two paradigms). Together, these patterns of performance highlight the importance of administering more than one task when screening for face perception deficits. This is particularly true for a condition such as DP, where face recognition difficulties appear to mostly be lifelong and do not accompany any other form of dysfunction. This allows many people with DP to develop elaborate compensatory strategies that may help them with particular facial stimuli or task paradigms, allowing them to obscure their difficulties (Adams et al., 2019). The case for repeat-testing aligns with our recent demonstration of the importance of repeat-screening for face memory deficits in DP, given the possibility that “typical” scores can be achieved by chance or due to low task reliability (Murray & Bate, 2020).

One further way to address this issue, particularly in tasks of face perception, is to place more emphasis on completion times, given accurate scores may be obtained by spending a long time on a task. Consistent with existing work (e.g. Bukach, et al., 2006; Busigny & Rossion, 2010; Delvenne et al., 2004; Jansari et al., 2015; Rossion & Michel, 2018), the finding reported here that 12 DPs were only impaired on completion time (but not accuracy) on either, or both

tests, highlights the importance of assessing both measures. Indeed, longer completion times may reflect the use of laboured face processing strategies and methods which ultimately lead to a correct response. However, it is important to note that our older adult controls took longer to complete the task than younger adults, although the same effect did not emerge for accuracy. Thus, we strongly suggest that age-matched norms are used for identifying impaired performance on this task. Additionally, with this finding in mind, the BFRT-r likely offers itself to be a suitable task for examining age-related changes in face processing within the typical population.

In conclusion, this paper presents an updated version of the BFRT with new theoretically-motivated stimuli. This task is suitable for rapid online administration, can detect face perception deficits in DP, and offers an opportunity for repeat-screening for consistency of performance when coupled with the BFRT-c. The task can be shared with other researchers on request.

#### Open Practices Statement

None of the experiments were pre-registered. The data are available as supplementary material. The BFRT-r stimuli and dataset are available via the Open Science Framework, and can be accessed here:

[https://osf.io/vza3m/?view\\_only=404f6d1971924759b126d46cba1d25b7](https://osf.io/vza3m/?view_only=404f6d1971924759b126d46cba1d25b7)

#### References

Adams, A., Hills, P., Bennetts, R., & Bate, S. (2019). Coping strategies for developmental prosopagnosia. *Neuropsychological Rehabilitation*, 1-20.

<https://doi.org/10.1080/09602011.2019.1623824>

- Annaz, D., Karmiloff-Smith, A., Johnson, M., & Thomas, M. (2009). A cross-syndrome study of the development of holistic face recognition in children with autism, Down syndrome, and Williams syndrome. *Journal of Experimental Child Psychology, 102*(4), 456-486. <https://doi.org/10.1016/j.jecp.2008.11.005>
- Barton, J. (2008). Structure and function in acquired prosopagnosia: Lessons from a series of 10 patients with brain damage. *Journal of Neuropsychology, 2*(1), 197-225. <https://doi.org/10.1348/174866407x214172>
- Barton, J., & Corrow, S. (2016). The problem of being bad at faces. *Neuropsychologia, 89*, 119-124. <https://doi.org/10.1016/j.neuropsychologia.2016.06.008>
- Bate, S., & Bennetts, R. (2015). The independence of expression and identity in face-processing: evidence from neuropsychological case studies. *Frontiers in Psychology, 6*. <https://doi.org/10.3389/fpsyg.2015.00770>
- Bate, S., Bennetts, R., Gregory, N., Tree, J., Murray, E., & Adams, A., Bobak, A.K., Penton, T., Yang, T., & Bannisy, M.J. (2019). Objective Patterns of Face Recognition Deficits in 165 Adults with Self-Reported Developmental Prosopagnosia. *Brain Sciences, 9*(6), 133. <https://doi.org/10.3390/brainsci9060133>
- Bate, S., Bennetts, R., Tree, J., Adams, A., & Murray, E. (2019). The domain-specificity of face matching impairments in 40 cases of developmental prosopagnosia. *Cognition, 192*, 104031. <https://doi.org/10.1016/j.cognition.2019.104031>
- Bate, S., Frowd, C., Bennetts, R., Hasshim, N., Murray, E., & Bobak, A.K., Wills, H., & Richards, S. (2018). Applied screening tests for the detection of superior face

recognition. *Cognitive Research: Principles and Implications*, 3(1).

<https://doi.org/10.1186/s41235-018-0116-5>

Bate, S., Frowd, C., Bennetts, R., Hasshim, N., Portch, E., Murray, E., & Dudfield, G. (2019).

The consistency of superior face recognition skills in police officers. *Applied*

*Cognitive Psychology*, 33(5), 828-842. <https://doi.org/10.1002/acp.3525>

Bate, S., Haslam, C., Tree, J., & Hodgson, T. (2008). Evidence of an eye movement-based

memory effect in congenital prosopagnosia. *Cortex*, 44(7), 806-819.

<https://doi.org/10.1016/j.cortex.2007.02.004>

Bate, S., & Tree, J. (2017). The Definition and Diagnosis of Developmental

Prosopagnosia. *Quarterly Journal of Experimental Psychology*, 70(2), 193-200.

<https://doi.org/10.1080/17470218.2016.1195414>

Behrmann, M., Avidan, G., Marotta, J., & Kimchi, R. (2005). Detailed Exploration of Face-

related Processing in Congenital Prosopagnosia: 1. Behavioral Findings. *Journal of Cognitive Neuroscience*, 17(7), 1130-1149.

<https://doi.org/10.1162/0898929054475154>

Bennetts, R., Butcher, N., Lander, K., Udale, R., & Bate, S. (2015). Movement cues aid face

recognition in developmental prosopagnosia. *Neuropsychology*, 29(6), 855-860.

<https://doi.org/10.1037/neu0000187>

Benton, A. L., Sivan, A. B., Hamsher, K. D. S., Varney, N. R., & Spreen, O. (1983). Facial

recognition: Stimulus and multiple choice pictures. In A. L. Benton, A. B. Sivan, K. D. S. Hamsher, N. R. Varney, & O. Spreen (Eds.), *Contribution to*

*neuropsychological assessment* (pp. 30–40). Oxford University Press.

- Benton, A. L., & Van Allen, M. W. (1968). Impairment in facial recognition in patients with cerebral disease. *Transactions of the American Neurological Association*, *93*, 38–42.
- Biotti, F., Gray, K., & Cook, R. (2019). Is developmental prosopagnosia best characterised as an apperceptive or mnemonic condition?. *Neuropsychologia*, *124*, 285-298.  
<https://doi.org/10.1016/j.neuropsychologia.2018.11.014>
- Biotti, F., Wu, E., Yang, H., Jiahui, G., Duchaine, B., & Cook, R. (2017). Normal composite face effects in developmental prosopagnosia. *Cortex*, *95*, 63-76.  
<https://doi.org/10.1016/j.cortex.2017.07.018>
- Bowles, D., McKone, E., Dawel, A., Duchaine, B., Palermo, R., & Schmalzl, L., Rivolta, D., Wilson, E., & Yovel, G. (2009). Diagnosing prosopagnosia: Effects of ageing, sex, and participant–stimulus ethnic match on the Cambridge Face Memory Test and Cambridge Face Perception Test. *Cognitive Neuropsychology*, *26*(5), 423-455.  
<https://doi.org/10.1080/02643290903343149>
- Bukach, C. M., Bub, D. N., Gauthier, I., & Tarr, M. J. (2006). Perceptual expertise effects are not all or none: Spatially limited perceptual expertise for faces in a case of prosopagnosia. *Journal of Cognitive Neuroscience*, *18*, 48–63.  
<https://doi.org/10.1162/089892906775250094>
- Burns, E., Martin, J., Chan, A., & Xu, H. (2017). Impaired processing of facial happiness, with or without awareness, in developmental prosopagnosia. *Neuropsychologia*, *102*, 217-228. <https://doi.org/10.1016/j.neuropsychologia.2017.06.020>
- Burton, M.A. (2013). Why has research in face recognition progressed so slowly? The importance of variability. *Quarterly Journal of Experimental Psychology*, *66*(8), 1467-1485. <https://doi.org/10.1080/17470218.2013.800125>

- Burton, M.A., White, D., & McNeill, A. (2010). The Glasgow Face Matching Test. *Behavior Research Methods*, 42, 286–291. <https://doi.org/10.3758/BRM.42.1.286>
- Busigny, T., & Rossion, B. (2010). Acquired prosopagnosia abolishes the face inversion effect. *Cortex*, 46, 965–981. <https://doi.org/10.1016/j.cortex.2009.07.004>
- Dalrymple, K., & Palermo, R. (2015). Guidelines for studying developmental prosopagnosia in adults and children. *Wiley Interdisciplinary Reviews: Cognitive Science*, 7(1), 73-87. <https://doi.org/10.1002/wcs.1374>
- De Luca, M., Pizzamiglio, M., Di Vita, A., Palermo, L., Tanzilli, A., Dacquino, C., & Piccardi, L. (2019). First the nose, last the eyes in congenital prosopagnosia: Look like your father looks. *Neuropsychology*, 33(6), 855-861. <https://doi.org/10.1037/neu0000556>
- De Renzi, E., Faglioni, P., Grossi, D., & Nichelli, P. (1991). Apperceptive and Associative Forms of Prosopagnosia. *Cortex*, 27(2), 213-221. [https://doi.org/10.1016/s0010-9452\(13\)80125-6](https://doi.org/10.1016/s0010-9452(13)80125-6)
- DeGutis, J., Chatterjee, G., Mercado, R., & Nakayama, K. (2012). Face gender recognition in developmental prosopagnosia: Evidence for holistic processing and use of configural information. *Visual Cognition*, 20(10), 1242-1253. <https://doi.org/10.1080/13506285.2012.744788>
- DeGutis, J., Cohan, S., & Nakayama, K. (2014). Holistic face training enhances face processing in developmental prosopagnosia. *Brain*, 137(6), 1781-1798. <https://doi.org/10.1093/brain/awu062>



- Delvenne, J.F., Seron, X., Coyette, F., & Rossion, B. (2004). Evidence for perceptual deficits in associative visual (prosop)agnosia: A single-case study. *Neuropsychologia*, *42*, 597–612. <https://doi.org/10.1016/j.neuropsychologia.2003.10.008>
- Duchaine, B., Germine, L., & Nakayama, K. (2007). Family resemblance: Ten family members with prosopagnosia and within-class object agnosia. *Cognitive Neuropsychology*, *24*(4), 419-430. <https://doi.org/10.1080/02643290701380491>
- Duchaine, B., & Nakayama, K. (2004). Developmental prosopagnosia and the Benton Facial Recognition Test. *Neurology*, *62*(7), 1219-1220. <https://doi.org/10.1212/01.wnl.0000118297.03161.b3>
- Duchaine, B., & Nakayama, K. (2006). The Cambridge Face Memory Test: Results for neurologically intact individuals and an investigation of its validity using inverted face stimuli and prosopagnosic participants. *Neuropsychologia*, *44*(4), 576-585. <https://doi.org/10.1016/j.neuropsychologia.2005.07.001>
- Duchaine, B., & Weidenfeld, A. (2003). An evaluation of two commonly used tests of unfamiliar face recognition. *Neuropsychologia*, *41*(6), 713-720. [https://doi.org/10.1016/s0028-3932\(02\)00222-1](https://doi.org/10.1016/s0028-3932(02)00222-1)
- Esins, J., Schultz, J., Stemper, C., Kennerknecht, I., & Bühlhoff, I. (2016). Face Perception and Test Reliabilities in Congenital Prosopagnosia in Seven Tests. *I-Perception*, *7*(1), 204166951562579. <https://doi.org/10.1177/2041669515625797>
- Geskin, J., & Behrmann, M. (2017). Congenital prosopagnosia without object agnosia? A literature review. *Cognitive Neuropsychology*, *35*(1-2), 4-54. <https://doi.org/10.1080/02643294.2017.1392295>

- Hasson, U., Avidan, G., Deouell, L., Bentin, S., & Malach, R. (2003). Face-selective Activation in a Congenital Prosopagnosic Subject. *Journal of Cognitive Neuroscience*, *15*(3), 419-431. <https://doi.org/10.1162/089892903321593135>
- Herlitz, A., & Lovén, J. (2013). Sex differences and the own-gender bias in face recognition: A meta-analytic review. *Visual Cognition*, *21*(9-10), 1306-1336. <https://doi.org/10.1080/13506285.2013.823140>
- Farah, M. J. (1990). Visual Agnosia: Disorders of object recognition and what they tell us about normal vision. MIT Press.
- Jansari, A., Miller, S., Pearce, L., Cobb, S., Sagiv, N., Williams, A. L., Tree, J., & Hanley, J. R. (2015). The man who mistook his neuropsychologist for a popstar: When configural processing fails in acquired prosopagnosia. *Frontiers in Human Neuroscience*, *9*, 390. <https://doi.org/10.3389/fnhum.2015.00390>
- Kennerknecht, I., Plümpe, N., Edwards, S., & Raman, R. (2006). Hereditary prosopagnosia (HPA): the first report outside the Caucasian population. *Journal of Human Genetics*, *52*(3), 230-236. <https://doi.org/10.1007/s10038-006-0101-6>
- Le Grand, R., Cooper, P., Mondloch, C., Lewis, T., Sagiv, N., de Gelder, B., & Maurer, D. (2006). What aspects of face processing are impaired in developmental prosopagnosia?. *Brain and Cognition*, *61*(2), 139-158. <https://doi.org/10.1016/j.bandc.2005.11.005>
- Lovén, J., Herlitz, A., & Rehnman, J. (2011). Women's own-gender bias in face recognition memory. The role of attention at encoding. *Experimental Psychology*, *58*, 333-340. <https://doi.org/10.1027/1618-3169/a000100>

- McKone, E., Hall, A., Pidcock, M., Palermo, R., Wilkinson, R., Rivolta, D., Yovel, G., David, J.M., & O'Connor, K.B. (2011). Face ethnicity and measurement reliability affect face recognition performance in developmental prosopagnosia: Evidence from the Cambridge Face Memory Test–Australian. *Cognitive Neuropsychology*, *28*(2), 109-146. <https://doi.org/10.1080/02643294.2011.616880>
- Minnebusch, D., Suchan, B., Ramon, M., & Daum, I. (2007). Event-related potentials reflect heterogeneity of developmental prosopagnosia. *European Journal of Neuroscience*, *25*(7), 2234-2247. <https://doi.org/10.1111/j.1460-9568.2007.05451.x>
- Murray, E., & Bate, S. (2019). Self-ratings of face recognition ability are influenced by gender but not prosopagnosia severity. *Psychological Assessment*, *31*(6), 828-832. <https://doi.org/10.1037/pas0000707>
- Murray, E., & Bate, S. (2020). Diagnosing developmental prosopagnosia: repeat assessment using the Cambridge Face Memory Test. *Royal Society Open Science*, *7*(9), 200884. <https://doi.org/10.1098/rsos.200884>
- Palermo, R., Willis, M., Rivolta, D., McKone, E., Wilson, C., & Calder, A. (2011). Impaired holistic coding of facial expression and facial identity in congenital prosopagnosia. *Neuropsychologia*, *49*(5), 1226-1235. <https://doi.org/10.1016/j.neuropsychologia.2011.02.021>
- Palermo, R., Rossion, B., Rhodes, G., Laguesse, R., Tez, T., & Hall, B., Albonico, A., Malaspina, M., Daini, R., Irons, J., Al-Janabi, S., Taylor, L.C., Rivolta, D., & McKone, E. (2017). Do People Have Insight into their Face Recognition Abilities?. *Quarterly Journal of Experimental Psychology*, *70*(2), 218-233. <https://doi.org/10.1080/17470218.2016.1161058>

- Rabin, L., Barr, W., & Burton, L. (2005). Assessment practices of clinical neuropsychologists in the United States and Canada: A survey of INS, NAN, and APA Division 40 members. *Archives of Clinical Neuropsychology*, *20*(1), 33-65.  
<https://doi.org/10.1016/j.acn.2004.02.005>
- Rezlescu, C., Danaila, I., Miron, A., & Amariei, C. (2020). More time for science: Using Testable to create and share behavioral experiments faster, recruit better participants, and engage students in hands-on research. *Progress in Brain Research*, *253*, 243-262.  
<https://doi.org/10.1016/bs.pbr.2020.06.005>
- Righart, R., & de Gelder, B. (2007). Impaired face and body perception in developmental prosopagnosia. *Proceedings of The National Academy of Sciences*, *104*(43), 17234-17238. <https://doi.org/10.1073/pnas.0707753104>
- Robertson, D., Noyes, E., Dowsett, A., Jenkins, R., & Burton, A. (2016). Face Recognition by Metropolitan Police Super-Recognisers. *PLOS ONE*, *11*(2), e0150036.  
<https://doi.org/10.1371/journal.pone.0150036>
- Rossion, B., & Michel, C. (2018). Normative accuracy and response time data for the computerized Benton Facial Recognition Test (BFRT-c). *Behavior Research Methods*, *50*(6), 2442-2460. <https://doi.org/10.3758/s13428-018-1023-x>
- Sachse, M., Schlitt, S., Hainz, D., Ciaramidaro, A., Walter, H., & Poustka, F., Bolte, S., & Freitag, C.M. (2014). Facial emotion recognition in paranoid schizophrenia and autism spectrum disorder. *Schizophrenia Research*, *159*(2-3), 509-514.  
<https://doi.org/10.1016/j.schres.2014.08.030>

- Shah, P., Sowden, S., Gaule, A., Catmur, C., & Bird, G. (2015). The 20 item prosopagnosia index (PI20): relationship with the Glasgow face-matching test. *Royal Society Open Science*, 2(11), 150305. <https://doi.org/10.1098/rsos.150305>
- Van Belle, G., Busigny, T., Lefèvre, P., Joubert, S., Felician, O., Gentile, F., & Rossion, B. (2011). Impairment of holistic face perception following right occipito-temporal damage in prosopagnosia: Converging evidence from gaze-contingency. *Neuropsychologia*, 49(11), 3145-3150. <https://doi.org/10.1016/j.neuropsychologia.2011.07.010>
- White, D., Rivolta, D.A., Burton, M., Al-Janabi, S., & Palermo, R. (2017). Face Matching Impairment in Developmental Prosopagnosia. *Quarterly Journal of Experimental Psychology*, 70(2), 287-297. <https://doi.org/10.1080/17470218.2016.1173076>
- Young, A., & Burton, A. (2017). Recognizing Faces. *Current Directions in Psychological Science*, 26(3), 212-217. <https://doi.org/10.1177/0963721416688114>
- Young, A. W., Newcombe, F., de Haan, E. H., Small, M., & Hay, D. C. (1993). Face perception after brain injury. Selective impairments affecting identity and expression. *Brain*, 116, 941–959. <https://doi.org/10.1093/brain/116.4.941>