

1  
2  
3 **The influence of candidates' physical attributes on assessors' ratings in clinical practice.**

4  
5  
6 **An experimental video-based simulation study.**

7  
8  
9  
10  
11  
12 **Sam AH<sup>1</sup>, Reid MD<sup>1</sup>, Thakerar V<sup>1</sup>, Gurnell M<sup>2</sup>, Westacott R<sup>3</sup>, Yeates P<sup>4</sup>, Reed MWR<sup>5</sup>,**  
13  
14 **Brown CA<sup>6</sup>,**

15  
16  
17  
18  
19  
20  
21 <sup>1</sup> Imperial College School of Medicine, Imperial College London, UK

22  
23 <sup>2</sup> Wellcome Trust-MRC Institute of Metabolic Science, University of Cambridge and NIHR  
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  
Cambridge Biomedical Research Centre, Addenbrooke's Hospital, Cambridge, UK

<sup>3</sup> Birmingham Medical School, University of Birmingham, Edgbaston, Birmingham, UK

<sup>4</sup> School of Medicine, Keele University, Keele, Staffordshire, UK and Fairfield General  
Hospital, Pennine Acute Hospitals NHS Trust, Bury, Lancashire, UK.

<sup>5</sup> Brighton and Sussex Medical School, University of Sussex, Brighton, UK

<sup>6</sup> Division of Health Sciences, Warwick Medical School, University of Warwick, UK

**Corresponding author:**

Dr Celia Brown

Warwick Medical School,

The University of Warwick,

Coventry, UK

[Celia.Brown@warwick.ac.uk](mailto:Celia.Brown@warwick.ac.uk)

## 1 **ABSTRACT**

### 2 **Background**

3 Assessments of physician competence in the work-place are common and often contribute to  
4 high-stakes assessments. Previous research suggests that assessors' judgements can be  
5 influenced by candidates' physical attributes. We investigated whether simulated candidates'  
6 scores were influenced by assessor bias based on tattoos, hair colour, and a regional accent.

### 7 **Methods**

8 We used an experimental, video-based, single-blinded, randomised, internet-based design.  
9 We created videos of simulated medical intern performances of a clinical examination at four  
10 different standards of competence. Four videos were also created of simulated candidates  
11 performing at a 'clear pass' standard, with either no stereotypical attribute (CPX), purple hair  
12 (CPH), tattoos (CPT) or a Liverpool English accent (CPA). Assessors were randomly  
13 assigned to watch five videos including the "clear pass" candidate without an attribute and  
14 one of the "clear pass" candidates with an attribute and asked to give an overall global grade  
15 for each candidate. We compared the global grades for the clear pass candidates with and  
16 without attributes.

### 17 **Results**

18 **Ninety-eight** assessors were included in the analysis. The total scores for the candidates with  
19 stereotyped attributes were not significantly lower than the candidate with no attribute.  
20 Assessors showed moderate levels of agreement between the global grades awarded for all  
21 the candidates. The global grades awarded to candidate with a stereotypical attribute were not  
22 significantly lower than for those without.

### 23 **Conclusions**

1  
2  
3 1 The presence of tattoos, purple hair, or a regional accent did not systematically negatively  
4  
5 2 influence the grade or score awarded by assessors to candidates in observed clinical  
6  
7  
8 3 examination scenarios.  
9

#### 10 4 **KEYWORDS**

11  
12  
13  
14 5 Assessment, bias, medicine, clinical  
15

#### 16 6 **PRACTICE POINTS**

- 17  
18  
19  
20 7 - Assessments of competence based on observations of practice are common in  
21  
22 8 healthcare settings.  
23  
24 9 - Individual assessor bias based on candidate characteristics has been previously  
25  
26 10 documented.  
27  
28  
29 11 - Systematic bias based on hair colour, tattoos, and UK regional accent does not seem  
30  
31 12 to negatively impact the scores or grades awarded by assessors when rating competent  
32  
33 13 candidates.  
34  
35

#### 36 14 37 38 39 15 **INTRODUCTION**

40  
41  
42 16 Ratings based on observations of a physician's competence in practice by senior colleagues  
43  
44 17 occur frequently and have traditionally contributed to learning in the workplace as part of an  
45  
46 18 apprenticeship model (Swanwick 2005). Workplace-based assessments of competence, such  
47  
48 19 as the mini-clinical evaluation exercise (mini-CEX), have been generally supported (Hatala et  
49  
50 20 al. 2006; Norcini & Burch 2007) and are increasingly being integrated into postgraduate  
51  
52 21 curricula across the world (Miller & Archer 2010). However, concerns have been raised  
53  
54 22 about the validity and reliability of such methods and their use as part of high-stakes  
55  
56 23 assessments (Hawkins et al. 2010). It is well established that assessors are prone to variability  
57  
58  
59  
60

1  
2  
3 1 due to cognitive biases such as leniency, inconsistency, and the halo effect (McManus et al.  
4  
5 2 2006; Iramaneerat & Yudkowsky 2007; Harasym et al. 2008). Individual examiners have also  
6  
7 3 been shown to rely on value-based judgements which are prone to stereotype bias (Williams  
8  
9 4 et al. 2003). Attempts to reduce the impact of these sources of assessor variability have  
10  
11 5 shown limited effect (Cook et al. 2009) such that they may ultimately threaten the validity  
12  
13 6 and objectivity of the assessment format (Hawkins et al. 2010). This paper contributes to the  
14  
15 7 developing understanding of sources of assessor variability due to bias.  
16  
17  
18  
19  
20  
21  
22  
23

24 9 The role of assessor inferences about candidate attributes such as body language, accent and  
25  
26 10 appearance has been shown to contribute to ratings (Kogan et al. 2011), but have not been  
27  
28 11 explored in a large-scale study. Further work regarding the origins of assessor variability in  
29  
30 12 direct observation assessments has resulted in a proposed model of ‘information integration’  
31  
32 13 by assessors which describes the formation of a general impression of a candidate first,  
33  
34 14 followed by the generation of domain scores second, rather than the reverse process which is  
35  
36 15 the intended method of such systems (Yeates et al. 2013). Previous studies have shown that  
37  
38 16 some physical attributes such as an individual’s ethnicity have an impact on their attainment  
39  
40 17 in both undergraduate and postgraduate medical examinations (Woolf et al. 2011). This effect  
41  
42 18 may be partly attributed to bias on behalf of the assessors but its overall origins are not clear  
43  
44 19 (Yeates et al. 2013). It is also apparent that amongst physicians in clinical practice bias based  
45  
46 20 on ethnicity persists and contributes to healthcare disparity for patients (Stone & Moskowitz  
47  
48 21 2011; Dovidio & Fiske 2012; Moskowitz et al. 2012).  
49  
50  
51  
52  
53

54 23 Stereotypes amongst the general population about those with tattoos (Wohlrab et al. 2007),  
55  
56 24 extremes of hair colour (Beddow 2011) and accents (Gluszek & Dovidio 2010) are  
57  
58 25 widespread. In particular, Liverpool English accents have been shown to be perceived as less  
59  
60

1 trustworthy than Standard Southern British English (SSBE) (Torre et al. 2018) and lower in  
 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

1 trustworthy than Standard Southern British English (SSBE) (Torre et al. 2018) and lower in  
 2 prestige and social attractiveness (Bishop et al. 2005). Activation of these stereotypes has  
 3 been shown to have an impact on real-world outcomes such as success in job interviews,  
 4 average salary and perceived professionalism (Johnston 2010; Deprez-Sims & Morris 2010;  
 5 Ruetzler et al. 2012), but there has been no work done to explore their role in assessment  
 6 within healthcare professionals' education. Despite best efforts, physicians remain prone to  
 7 the same implicit biases as the general population which may unconsciously impact decision  
 8 making (Chapman et al. 2013). In some cases bias based on these stereotypes is more overt,  
 9 such that physicians have been shown to openly express a preference for their colleagues to  
 10 be dressed according to established norms and where individuals deviate from this standard  
 11 peers may perceive this as a professionalism concern (Gjerdingen et al. 1987). Stereotypes  
 12 are more likely to be activated and result in bias when judgements are mentally demanding  
 13 (Macrae et al. 1994), for example during medical exams (Tavares & Eva 2014). Physical  
 14 attributes therefore present a potential source of bias that may influence assessor ratings and  
 15 challenge the validity of workplace-based assessments of competence in clinical practice.  
 16 This study therefore sought to establish whether the presence of a variety of physical  
 17 attributes amongst candidates performing at a standardised level had any effect on assessors'  
 18 ratings.

## 21 **METHODS**

### 22 *Study Design*

23 We used an experimental, video-based, single-blinded, randomised, internet-based design.

### 25 *Procedure*

1  
2  
3 1 Seven 10-minute videos were created of simulated candidates completing a clinical  
4  
5 2 examination typical of those observed in practice (a cranial nerve examination). Volunteer  
6  
7 3 Clinical Teaching Fellows affiliated with Imperial College London were recruited for this  
8  
9 4 role. All simulated candidates were female, of white ethnicity, and a similar age to avoid  
10  
11 5 potential confounding based on these factors. Four of the videos demonstrated the simulated  
12  
13 6 candidates performing the examination at one of four overall performance levels: 'clear fail'  
14  
15 7 (CF), 'borderline' (BD), 'clear pass' (CPX) or 'good' (GD). The other three videos showed a  
16  
17 8 candidate performing at a 'clear pass' level but with either purple hair (CPH), tattoos (CPT),  
18  
19 9 or a Liverpool English accent (CPA). The simulated candidates in all videos except CPA  
20  
21 10 performed with a SSBE accent. Each candidate followed a script created by a panel of  
22  
23 11 experienced examiners to ensure they were performing at the appropriate level and to  
24  
25 12 standardise those performing at the 'clear pass' level. Twelve sets of five videos were then  
26  
27 13 created; with every set including a video of a candidate performing at each of the overall  
28  
29 14 performance levels as well as one video of a candidate with a physical attribute performing at  
30  
31 15 a 'clear pass' level (*Appendix 1*). The ordering of the five videos differed across the 12 sets to  
32  
33 16 mitigate any bias associated with ordering effects. Each participant was randomly allocated to  
34  
35 17 one of the 12 video sets.  
36  
37  
38  
39  
40  
41  
42  
43  
44

### 19 ***Recruitment and Consent***

46  
47 20 The study was approved by the Medical Education Ethics Committee at Imperial College  
48  
49 21 London (MEEC1718-105). Each medical school in the UK was contacted via the Medical  
50  
51 22 Schools Council and invited to take part in the study. Heads of assessment at each medical  
52  
53 23 school were encouraged to invite a representative sample of assessors to participate in the  
54  
55 24 study via the study website. Participants were informed that they were taking part in a study  
56  
57 25 exploring inter-rater reliability amongst assessors but were not informed that the study aimed  
58  
59  
60

1  
2  
3 1 to evaluate the impact of physical attributes on scores and performance levels. No identifiable  
4  
5 2 information was collected about the participants. Participants were required to be clinicians  
6  
7 3 with at least one prior experience of formally assessing medical students in clinical  
8  
9  
10 4 examinations. Participants were informed that completion of the marksheets for all five  
11  
12 5 videos and submission of the post-completion questionnaire was evidence of consent.  
13  
14 6 Participants were able to withdraw from the process by closing the web browser at any time  
15  
16 7 prior to completion of the study but due to the lack of collection of identifiable data, were not  
17  
18 8 able to withdraw after submitting their results. Any incomplete data, where participants did  
19  
20 9 not view and score all five videos, were not used in the analysis.  
21  
22  
23  
24 10

### 25 26 11 ***Measures***

27  
28 12 Participants were asked to assess the candidates at the level expected of a foundation year 1  
29  
30 13 doctor (equivalent to a medical intern). Participants viewed the five videos and were provided  
31  
32 14 with a blank mark sheet to complete alongside each video (*Figure 1*). Participants marked  
33  
34 15 each candidate in four domains; 'Physical examination', 'Identify physical signs and the most  
35  
36 16 likely diagnosis', 'Clinical management skills', and 'Interpersonal skills'. Each domain was  
37  
38 17 scored between 0 - 4, with a maximum possible total score of 16. Participants were also asked  
39  
40 18 to assign each candidate a global grade of either 'clear fail', 'borderline', 'clear pass' or  
41  
42 19 'good'. Participants were able to return to mark sheets for previous candidates but were not  
43  
44 20 able to pause, rewind or replay the videos, to reflect the contemporaneous nature of rating a  
45  
46 21 competency in practice. Following completion of the mark sheets for all five videos,  
47  
48 22 participants were asked to confirm their assessment experience, job role, gender, ethnicity  
49  
50 23 and the geographical region where they worked.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

25 [FIGURE 1]

### *Statistical analysis*

Data management and analysis were conducted using Stata V16. Total scores and global grades for each candidate with an attribute were compared with those of the clear pass candidate without an attribute. Individual Wilcoxon matched-pairs signed-ranks tests were used to compare the total scores. **Weighted kappa analysis was used to compare the global grades, followed by a Wilcoxon analysis to measure the direction of any disagreement.** A p-value of less than 0.0167 was required for statistical significance to account for multiple comparisons within each type of score.

## RESULTS

### *Participants*

One-hundred and twenty assessors participated in the study, **of whom 98 were included in the analysis** (five assessors were removed due to a self-reported lack of experience; seventeen participants did not complete viewing and rating all of the candidates). Table 1 shows the demographic details of all participants included in the analysis. Participants included in the analysis came from ten distinct regions across the UK. Participants were varied in their level of experience and job role. The number of participants who viewed and rated each of the 12 sets of videos was comparable.

[TABLE 1]



1

2

3 **Total Score**

4 Total scores for all the clear pass candidates ranged from 8 to 16/16 (median 14, interquartile  
5 range [IQR] 12 to 15. *Figure 2*). **The modal scores for each candidate were as follows;**  
6 **CPX = 14, CPH = 16, CPT = 12, CPA = 14.** Individual Wilcoxon matched-pairs signed-  
7 ranks tests were performed on the total scores for the clear pass candidate with no attribute  
8 when compared to each clear pass candidate with an attribute. For the candidate with purple  
9 hair (CPH) this indicated that scores were statistically significantly higher (median paired  
10 difference 1, range -4 to 8,  $Z=2.42$ ,  $p=0.01$ ). There was no significant difference between the  
11 total scores for CPX and CPT (median paired difference 1, range -4 to 3,  $Z=1.68$ ,  $p=0.09$ ) or  
12 between CPX and CPA (median paired difference -1, range -3 to 2,  $Z=1.26$ ,  $p=0.22$ )

13  
14 [FIGURE 2]15  
16 **Global Grade**

17 Global Grades for all candidates varied from borderline to good. **A weighted kappa analysis**  
18 **using linear weights showed individual assessors had moderate agreement between the**  
19 **global grades awarded to CPX and CPT ( $K=0.412$ ,  $p=0.007$ ) and to CPA ( $K=0.446$ ,**  
20  **$p=0.004$ ).** There was no significant agreement between the global grade awarded to  
21 **CPH when compared to CPX ( $K=0.158$ ,  $p=0.129$ ).** A Wilcoxon matched-pairs signed-  
22 ranks analysis was performed to measure the direction of this difference by applying  
23 numerical values to each global grade, where fail=1, borderline=2, clear pass=3 and  
24 good=4. This showed no statistically significant difference ( $Z=2.13$ ,  $p=0.06$ ) but  
25 confirmed the median paired difference for CPH was 1 grade higher than for CPX.

1  
2  
3 1 *Figure 3* shows the number of assessors giving each configuration of global grades to each  
4  
5 2 candidate.  
6  
7  
8 3

9  
10 4 [FIGURE 3]  
11  
12 5  
13  
14

## 15 6 **DISCUSSION**

16  
17 7 For the first time we have compared the influence of hair colour, tattoos and accent on the  
18  
19 8 ratings clinicians give to simulated performances of clinical examinations by candidates.  
20

21 9 There was no negative impact on the global grades awarded by assessors despite the presence  
22  
23 10 of stereotyped physical attributes. Similarly, the total scores for clear pass candidates with  
24  
25 11 physical attributes were not significantly lower than for the candidate without these  
26  
27 12 characteristics. Interestingly, assessors gave higher total scores and global grades to the  
28  
29 13 candidate with purple hair than to the candidate performing at the same level without a  
30  
31 14 physical attribute. These findings are largely reassuring and suggest that any assessor bias  
32  
33 15 based on the presence of tattoos, hair colour and accent does not negatively influence their  
34  
35 16 judgement. This finding is in keeping with previous studies that suggest examiner bias is not  
36  
37 17 responsible for the differential attainment amongst minority ethnic medical students (Yeates  
38  
39 18 et al. 2017). The higher scores and global grades awarded to the candidate with purple hair  
40  
41 19 may represent a positive contrast effect based on the presence of a notable characteristic  
42  
43 20 which lead the candidate to stand out when compared to others (Yeates et al. 2015).  
44  
45  
46  
47  
48

49 21 However, any explanation for the difference in total scores is speculative at this stage and is  
50  
51 22 likely to require further research.  
52  
53  
54  
55

56 24 The study used a randomised, single-blinded, controlled methodology to explore the  
57  
58 25 influence of candidates' physical attributes on assessor ratings. However, the study does have  
59  
60

1  
2  
3 1 some limitations. The study used video recordings of simulated performances and it is  
4  
5 2 therefore possible that in real life assessors may be more or less vulnerable to bias than they  
6  
7 3 were in this study. **Further work should continue to explore the impact of bias in real-life**  
8  
9 4 **assessments.** We necessarily used different actors for each performance and whilst every  
10  
11 5 attempt was made to control for other sources of variability in the performances between  
12  
13 6 candidates by standardising for age, gender and ethnicity and using a script, it is possible that  
14  
15 7 minor variations between candidates persisted. All participants were volunteers and therefore  
16  
17 8 it is possible that they are not a representative sample of the population of assessors as a  
18  
19 9 whole. The study only explored the impact of physical attributes amongst white, female  
20  
21 10 candidates performing at a clear pass standard and it is important to note that these findings  
22  
23 11 may not be generalisable to candidates of other demographics, or to those performing at  
24  
25 12 different levels. The study explored the impact of physical attributes in the context of an  
26  
27 13 observed performance of a cranial nerve examination and we cannot exclude that different  
28  
29 14 effects may occur in other types of assessment, particularly when they are more cognitively  
30  
31 15 demanding for assessors. **We also recognise the impact mark schemes and global grading**  
32  
33 16 **systems may have on outcomes and our results may therefore not be generalisable if**  
34  
35 17 **significantly different scoring rubrics are used. Further work is still needed to explore if**  
36  
37 18 **other physical attributes such as choice of attire may have an impact on assessor**  
38  
39 19 **ratings. It is also worth noting that any systematic effect of bias based on stereotype**  
40  
41 20 **activation may vary over time as societal attitudes towards individual attributes also**  
42  
43 21 **change.**  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

## 23 CONCLUSION

24 **Within the context of an online simulated assessment there does not appear to be any**  
25 **systematic effect of negative stereotype bias from assessors when rating competent**

1  
2  
3 1 **candidates with tattoos, purple hair or a Liverpudlian accent when compared to a**  
4  
5 2 **candidate without these characteristics.**  
6  
7  
8 3  
9  
10 4  
11  
12 5  
13  
14

## 15 6 **ACKNOWLEDGEMENTS**

17  
18 7 The authors are grateful to all UK medical school assessment leads for their help in recruiting  
19  
20 8 assessors. The authors are also grateful to the Medical Schools Council for administrative  
21  
22 9 support with the study.  
23  
24  
25  
26 10  
27  
28

## 29 11 **DECLARATION OF INTERESTS**

30  
31  
32 12 MG is supported by the National Institute for Health Research (NIHR) Cambridge  
33  
34 13 Biomedical Research Centre. CAB is supported by the NIHR Applied Research  
35  
36 14 Collaboration (ARC) West Midlands. PY is funded by the NIHR Clinician Scientist Award.  
37  
38 15 The views expressed are those of the author(s) and not necessarily those of the NHS, the  
39  
40 16 NIHR or the Department of Health and Social Care.  
41  
42  
43  
44 17  
45  
46

## 47 18 **FUNDING**

48  
49  
50 19 The Medical Schools Council funded the recruitment of the simulated candidates, simulated  
51  
52 20 patient and sourcing of the recording equipment for this study.  
53  
54  
55  
56 21  
57  
58 22  
59  
60

1  
2  
3 **1 REFERENCES**  
4

- 5  
6 2 Beddow M. 2011. Hair color stereotypes and their associated perceptions in relationships and  
7  
8 3 the workplace. [place unknown]; [accessed 2020 Feb 10].  
9  
10 4 <https://pdfs.semanticscholar.org/57fd/6d85010f6db3475ee9acb9b651a9c6dcca27.pdf>  
11  
12  
13 5 Bishop H, Coupland N, Garrett P. 2005. Conceptual accent evaluation: Thirty years of accent  
14  
15 6 prejudice in the UK. *Acta Linguist Hafniensia*. 37(1):131–154.  
16  
17  
18 7 Chapman EN, Kaatz A, Carnes M. 2013. Physicians and implicit bias: How doctors may  
19  
20 8 unwittingly perpetuate health care disparities. *J Gen Intern Med*. 28(11):1504–1510.  
21  
22  
23 9 Cook DA, Dupras DM, Beckman TJ, Thomas KG, Pankratz VS. 2009. Effect of rater training  
24  
25 10 on reliability and accuracy of mini-CEX scores: A randomized, controlled trial. *J Gen Intern*  
26  
27 11 *Med*. 24(1):74–79.  
28  
29  
30 12 Deprez-Sims AS, Morris SB. 2010. Accents in the workplace: Their effects during a job  
31  
32 13 interview. *Int J Psychol [Internet]*. [accessed 2020 Feb 10] 45(6):417–426.  
33  
34 14 <http://www.ncbi.nlm.nih.gov/pubmed/22044081>  
35  
36  
37 15 Dovidio JF, Fiske ST. 2012. Under the radar: How unexamined biases in decision-making  
38  
39 16 processes in clinical interactions can contribute to health care disparities. *Am J Public Health*.  
40  
41 17 102(5):945–952.  
42  
43  
44 18 Gjerdingen DK, Simpson DE, Titus SL. 1987. Patients' and physicians' attitudes regarding  
45  
46 19 the physician's professional appearance. *Arch Intern Med [Internet]*. [accessed 2020 Feb 10]  
47  
48 20 147(7):1209–12. <http://www.ncbi.nlm.nih.gov/pubmed/3606278>  
49  
50  
51 21 Gluszek A, Dovidio JF. 2010. The way they speak: a social psychological perspective on the  
52  
53 22 stigma of nonnative accents in communication. *Pers Soc Psychol Rev [Internet]*. [accessed  
54  
55 23 2020 Feb 10] 14(2):214–37. <http://www.ncbi.nlm.nih.gov/pubmed/20220208>  
56  
57  
58  
59  
60

- 1  
2  
3 1 Harasym PH, Woloschuk W, Cunning L. 2008. Undesired variance due to examiner  
4  
5 2 stringency/leniency effect in communication skill scores assessed in OSCEs. *Adv Heal Sci*  
6  
7 3 *Educ.* 13(5):617–632.  
8  
9  
10 4 Hatala R, Ainslie M, Kassen BO, Mackie I, Roberts JM. 2006. Assessing the mini-clinical  
11  
12 5 evaluation exercise in comparison to a national specialty examination. *Med Educ.*  
13  
14 6 40(10):950–956.  
15  
16  
17 7 Hawkins RE, Margolis MJ, Durning SJ, Norcini JJ. 2010. Constructing a validity argument  
18  
19 8 for the mini-Clinical Evaluation Exercise: a review of the research. *Acad Med* [Internet].  
20  
21 9 [accessed 2020 Feb 10] 85(9):1453–61. <http://www.ncbi.nlm.nih.gov/pubmed/20736673>  
22  
23  
24  
25 10 Iramaneerat C, Yudkowsky R. 2007. Rater errors in a clinical skills assessment of medical  
26  
27 11 students. *Eval Heal Prof* [Internet]. [accessed 2020 Feb 10] 30(3):266–283.  
28  
29 12 <http://www.ncbi.nlm.nih.gov/pubmed/17693619>  
30  
31  
32  
33 13 Johnston DW. 2010. Physical appearance and wages: Do blondes have more fun? *Econ Lett.*  
34  
35 14 108(1):10–12.  
36  
37  
38 15 Kogan JR, Conforti L, Bernabeo E, Iobst W, Holmboe E. 2011. Opening the black box of  
39  
40 16 clinical skills assessment via observation: a conceptual model. *Med Educ* [Internet].  
41  
42 17 [accessed 2020 Feb 10] 45(10):1048–60. <http://www.ncbi.nlm.nih.gov/pubmed/21916943>  
43  
44  
45  
46 18 Macrae CN, Milne AB, Bodenhausen G V. 1994. Stereotypes as energy-saving devices: A  
47  
48 19 peek inside the cognitive toolbox. *J Pers Soc Psychol* [Internet]. [accessed 2020 Feb 10]  
49  
50 20 66(1):37–47. <http://doi.apa.org/getdoi.cfm?doi=10.1037/0022-3514.66.1.37>  
51  
52  
53 21 McManus I, Thompson M, Mollon J. 2006. Assessment of examiner leniency and stringency  
54  
55 22 ('hawk-dove effect') in the MRCP(UK) clinical examination (PACES) using multi-facet  
56  
57 23 Rasch modelling. *BMC Med Educ* [Internet]. [accessed 2020 Feb 11] 6(1):42.  
58  
59  
60

- 1  
2  
3 1 <https://bmcmmededuc.biomedcentral.com/articles/10.1186/1472-6920-6-42>  
4  
5  
6 2 Miller A, Archer J. 2010. Impact of workplace based assessment on doctors' education and  
7  
8 3 performance: A systematic review. *BMJ*. 341(7775):710.  
9  
10  
11 4 Moskowitz GB, Stone J, Childs A. 2012. Implicit stereotyping and medical decisions:  
12  
13 5 Unconscious stereotype activation in practitioners' thoughts about African Americans. *Am J*  
14  
15 6 *Public Health*. 102(5):996–1001.  
16  
17  
18 7 Norcini J, Burch V. 2007. Workplace-based assessment as an educational tool: AMEE Guide  
19  
20 8 No. 31. *Med Teach*. 29(9–10):855–871.  
21  
22  
23 9 Ruetzler T, Taylor J, Reynolds D, Baker W, Killen C. 2012. What is professional attire  
24  
25 10 today? A conjoint analysis of personal presentation attributes. *Int J Hosp Manag*. 31(3):937–  
26  
27 11 943.  
28  
29  
30 12 Stone J, Moskowitz GB. 2011. Non-conscious bias in medical decision making: What can be  
31  
32 13 done to reduce it? *Med Educ* [Internet]. [accessed 2020 Feb 10] 45(8):768–776.  
33  
34 14 <http://www.ncbi.nlm.nih.gov/pubmed/21752073>  
35  
36  
37 15 Swanwick T. 2005. Informal learning in postgraduate medical education: from cognitivism to  
38  
39 16 “culturism.” *Med Educ* [Internet]. [accessed 2020 Feb 19] 39(8):859–865.  
40  
41 17 <http://doi.wiley.com/10.1111/j.1365-2929.2005.02224.x>  
42  
43  
44 18 Tavares W, Eva KW. 2014. Impact of rating demands on rater-based assessments of clinical  
45  
46 19 competence. *Educ Prim Care*. 25(6):308–318.  
47  
48  
49 20 Torre I, Goslin J, White L, Zanatto D. 2018. Trust in artificial voices: A “congruency effect”  
50  
51 21 of first impressions and behavioural experience. In: *ACM Int Conf Proceeding Ser* [Internet].  
52  
53 22 New York, New York, USA: Association for Computing Machinery; [accessed 2020 Nov  
54  
55 23 21]; p. 1–6. <http://dl.acm.org/citation.cfm?doid=3183654.3183691>  
56  
57  
58  
59  
60

- 1  
2  
3 1 Williams RG, Klamen DA, McGaghie WC. 2003. Cognitive, Social and Environmental  
4 Sources of Bias in Clinical Performance Ratings. Teach Learn Med [Internet]. [accessed 2020  
5 Feb 10] 15(4):270–292. <http://www.ncbi.nlm.nih.gov/pubmed/14612262>  
6  
7  
8 2 Wohlrab S, Stahl J, Kappeler PM. 2007. Modifying the body: Motivations for getting  
9 tattooed and pierced. Body Image [Internet]. [accessed 2020 Feb 10] 4(1):87–95.  
10  
11 <http://www.ncbi.nlm.nih.gov/pubmed/18089255>  
12  
13 3 Woolf K, Potts HWW, McManus IC. 2011. Ethnicity and academic performance in UK  
14 trained doctors and medical students: systematic review and meta-analysis. BMJ [Internet].  
15  
16 [accessed 2020 Feb 10] 342(7797):d901. <http://www.ncbi.nlm.nih.gov/pubmed/21385802>  
17  
18 4 Yeates P, Cardell J, Byrne G, Eva KW. 2015. Relatively speaking: contrast effects influence  
19 assessors' scores and narrative feedback. Med Educ [Internet]. [accessed 2020 Mar 30]  
20  
21 49(9):909–919. <http://doi.wiley.com/10.1111/medu.12777>  
22  
23 5 Yeates P, O'Neill P, Mann K, Eva K. 2013. Seeing the same thing differently: mechanisms  
24 that contribute to assessor differences in directly-observed performance assessments. Adv  
25  
26 Health Sci Educ Theory Pract [Internet]. [accessed 2020 Feb 10] 18(3):325–41.  
27  
28 <http://www.ncbi.nlm.nih.gov/pubmed/22581567>  
29  
30 6 Yeates P, Woolf K, Benbow E, Davies B, Boohan M, Eva K. 2017. A randomised trial of the  
31 influence of racial stereotype bias on examiners' scores, feedback and recollections in  
32  
33 undergraduate clinical exams. BMC Med [Internet]. [accessed 2020 Apr 21] 15(1):179.  
34  
35 <http://bmcmmedicine.biomedcentral.com/articles/10.1186/s12916-017-0943-0>  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51

## 52 **APPENDIX 1**

### 53 *Video Ordering*

54  
55  
56  
57  
58  
59  
60



Version	Video 1	Video 2	Video 3	Video 4	Video 5
1	CPX	CPH	BL	CF	GD
2	CPX	CF	CPH	GD	BL
3	CPX	BL	GD	CPH	CF
4	CPX	GD	CF	BL	CPH
5	CPX	CPT	BL	CF	GD
6	CPX	CF	CPT	GD	BL
7	CPX	BL	GD	CPT	CF
8	CPX	GD	CF	BL	CPT
9	CPX	CPA	BL	CF	GD
10	CPX	CF	CPA	GD	BL
11	CPX	BL	GD	CPA	CF
12	CPX	GD	CF	BL	CPA

1 Key: CF – clear fail, BL – borderline, CPX – clear pass, no discernible attribute, CPH – clear  
 2 pass, purple hair, CPT – clear pass, tattoo on both forearms, CPA – clear pass, regional  
 3 accent, GD – good.

4

All participants N=98											
Demographics	CPH n=32			CPT n=34			CPA n=32				
	n	(%)	n	(%)	n	(%)	n	(%)	n	(%)	
<b>Experience</b>											
None	0	0.00	0	0.00	0	0.00	0	0.00	0	0.00	
1-2 Exams	6	6.12	2	6.25	1	2.94	3	9.38			
3-4 Exams	13	13.27	5	15.63	6	17.65	2	6.25			
5+ Exams	79	80.61	25	78.13	27	79.41	27	84.38			
<b>Job Role</b>											
Consultant	40	40.82	9	28.13	17	50.00	14	43.75			
Primary Care Physician	33	33.67	0	0.00	0	0.00	0	0.00			
Specialty Training years 3+	1	1.02	1	3.13	0	0.00	0	0.00			
Core Training or Specialty Training years 1-2	1	1.02	0	0.00	0	0.00	1	3.13			
Other/please specify role & grade if appropriate	22	22.45	6	18.75	9	26.47	7	21.88			
Prefer not to say	1	1.02	1	3.13	0	0.00	0	0.00			
<b>Gender</b>											
Male	44	44.90	13	40.63	14	41.18	17	53.13			
Female	54	55.10	19	59.38	20	58.82	15	46.88			
<b>Ethnicity</b>											
Asian	16	16.33	6	18.75	7	20.59	3	9.38			
Black African/Caribbean	3	3.06	0	0.00	1	2.94	2	6.25			
White	75	76.53	24	75.00	25	73.53	26	81.25			
Mixed/multiple	1	1.02	1	3.13	0	0.00	0	0.00			
Other/please specify	2	2.04	1	3.13	0	0.00	1	3.13			
Prefer not to say	1	1.02	0	0.00	1	2.94	0	0.00			
<b>Region</b>											
East Anglia	16	16.33	3	9.38	9	26.47	4	12.50			
East Midlands	5	5.10	3	9.38	1	2.94	1	3.13			
London	15	15.31	4	12.50	6	17.65	5	15.63			
North West	6	6.12	1	3.13	0	0.00	5	15.63			
Scotland	19	19.39	8	25.00	6	17.65	5	15.63			
South East	7	7.14	3	9.38	2	5.88	2	6.25			
South West	8	8.16	3	9.38	2	5.88	3	9.38			
Wales	2	2.04	0	0.00	0	0.00	2	6.25			
West Midlands	4	4.08	1	3.13	2	5.88	1	3.13			
Yorkshire and the Humber	16	16.33	6	18.75	6	17.65	4	12.50			

1  
2  
3  
4  
5  
6  
7  
8 *Table 1: Participant descriptives for all participants, and for the participants who rated the*  
9 *performance of the candidates performing at a 'clear pass' level who also had the presence*  
10 *of a physical attribute. CPH = Clear Pass, Purple Hair. CPT = Clear Pass, Tattoo. CPA =*  
11 *Clear Pass, Accent.*  
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 Peer-Review Only

**Mark Sheet: Cranial Nerve Examination****Domain 1. Physical examination****Task:** Examines the cranial nerves (I-XII)

Excellent	(4)
Good	(3)
Adequate	(2)
Fail	(1)
Severe fail	(0)

**Domain 2. Identifying physical signs and the most likely diagnosis****Task:** Reports abnormal findings and offers the most likely diagnosis

Excellent	(4)
Good	(3)
Adequate	(2)
Fail	(1)
Severe fail	(0)

**Domain 3. Clinical management skills****Task:** Explains management of patient

Excellent	(4)
Good	(3)
Adequate	(2)
Fail	(1)
Severe fail	(0)

**Domain 4. Interpersonal skills****Task:** Communicates appropriately with the patient and examiner

Excellent	(4)
Good	(3)
Adequate	(2)
Fail	(1)
Severe fail	(0)

**Global Grade**

Good  
Clear Pass  
Borderline  
Fail

*Figure 1: Sample mark sheet*

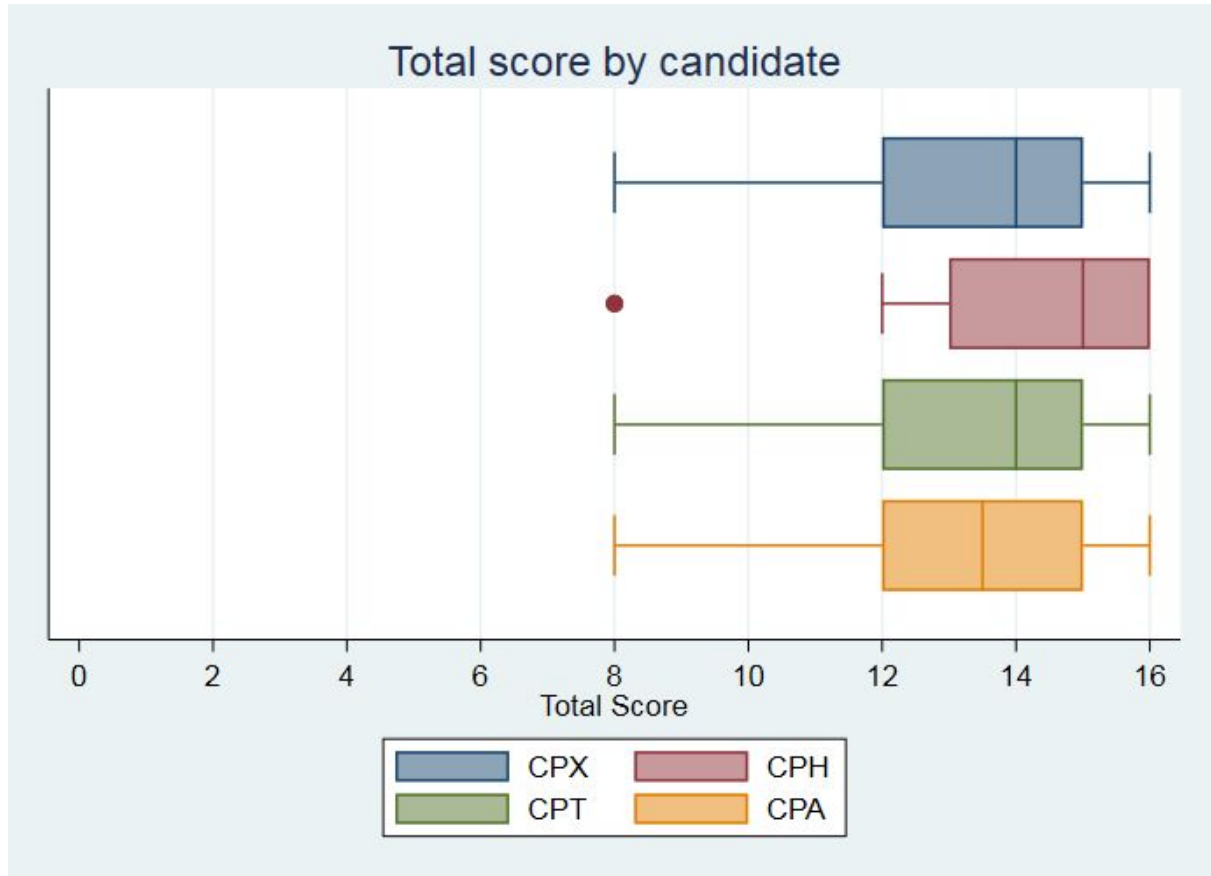


Figure 2: Total scores by candidate; CPX = Clear pass, no attribute, CPH = Clear pass, purple hair, CPT = Clear pass, tattoo, CPA = Clear pass, accent

		CPX			
		Fail	Borderline	Clear Pass	Good
CPH	Fail	0 (0.0)	0 (0.0)	0 (0.0)	0 (0.0)
	Borderline	0 (0.0)	0 (0.0)	0 (0.0)	0 (0.0)
	Clear Pass	0 (0.0)	0 (0.0)	3 (9.4)	2 (6.3)
	Good	0 (0.0)	1 (3.1)	8 (25)	18 (56.3)

		CPX			
		Fail	Borderline	Clear Pass	Good
CPT	Fail	0 (0.0)	0 (0.0)	0 (0.0)	0 (0.0)
	Borderline	0 (0.0)	0 (0.0)	0 (0.0)	0 (0.0)
	Clear Pass	0 (0.0)	0 (0.0)	10 (29.4)	3 (8.8)
	Good	0 (0.0)	0 (0.0)	7 (20.6)	14 (41.2)

		CPX			
		Fail	Borderline	Clear Pass	Good
CPA	Fail	0 (0.0)	0 (0.0)	0 (0.0)	0 (0.0)
	Borderline	0 (0.0)	0 (0.0)	0 (0.0)	0 (0.0)
	Clear Pass	0 (0.0)	0 (0.0)	10 (31.3)	7 (21.9)
	Good	0 (0.0)	0 (0.0)	2 (6.3)	13 (40.6)

Figure 3: Number of assessors giving each combination of global grades to the candidate with no clear attribute (CPX) and the candidates with purple hair (CPH), tattoos (CPT), and an accent (CPA).

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

Amir H. Sam , BSc, MBBS, PhD, FRCP, SFHEA, is head of Imperial College School of Medicine and consultant physician and endocrinologist at Imperial College Healthcare NHS Trust.

Michael D. Reid , MA, MB, BChir, MACadMed, MRCP, was a Clinical Education & Research Fellow at Imperial College London and is now a trainee in Geriatric Medicine at Kingston Hospital.

Viral Thakerar , MBBS, MRCP, MRCGP, is the lead for year 1 and 2 clinical placements at Imperial College School of Medicine and a practising general practitioner.

Mark Gurnell , PhD, MA(MEd), FHEA, FAcadMed, FRCP, is Clinical SubDean at the University of Cambridge School of Clinical Medicine and Professor of Clinical Endocrinology at Institute of Metabolic Science & Department of Medicine.

Rachel Westacott , MB, ChB, FRCP, is a Senior Lecturer in Medical Education at Birmingham Medical School and an acute medicine consultant at University Hospitals of Leicester NHS Trust (CCT in Nephrology).

Peter Yeates , MRCP, PhD, is a senior lecturer in medical education research and a consultant in acute and respiratory medicine. His interests focus on assessor cognition and technology-enhanced assessment.

Malcolm W. R. Reed , MD, BMedSci, MBChB, FRCS, is a breast cancer surgeon who has been Dean of Brighton and Sussex Medical School since 2014 having moved from Sheffield University Medical School where he was head of Undergraduate Assessment for medicine. He is currently Co-Chair of Medical Schools Council and Chair of the education subcommittee.

Celia A. Brown , PhD, SFHEA, is an Associate Professor in Quantitative Methods at Warwick Medical School. She has research interests in selection and assessment and teaches quantitative methods at all levels in Higher Education.

ew Only