

# BMJ Open Measuring differential attainment: a longitudinal analysis of assessment results for 1512 medical students at four Scottish medical schools

David Hope ,<sup>1</sup> Avril Dewar,<sup>1</sup> Eleanor J Hothersall,<sup>2</sup> John Paul Leach,<sup>3</sup> Isobel Cameron,<sup>4</sup> Alan Jaap<sup>1</sup>

**To cite:** Hope D, Dewar A, Hothersall EJ, *et al*. Measuring differential attainment: a longitudinal analysis of assessment results for 1512 medical students at four Scottish medical schools. *BMJ Open* 2021;**11**:e046056. doi:10.1136/bmjopen-2020-046056

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2020-046056>).

Received 22 October 2020  
Accepted 29 July 2021



© Author(s) (or their employer(s)) 2021. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

<sup>1</sup>Medical Education Unit, The University of Edinburgh College of Medicine and Veterinary Medicine, Edinburgh, UK

<sup>2</sup>School of Medicine, University of Dundee, Dundee, UK

<sup>3</sup>School of Medicine, Dentistry, and Nursing, University of Glasgow, Glasgow, UK

<sup>4</sup>School of Medicine and Dentistry, University of Aberdeen, Aberdeen, UK

## Correspondence to

Dr David Hope;  
[david.hope@ed.ac.uk](mailto:david.hope@ed.ac.uk)

## ABSTRACT

**Objective** To measure Differential Attainment (DA) among Scottish medical students and to explore whether attainment gaps increase or decrease during medical school.

**Design** A retrospective analysis of undergraduate medical student performance on written assessment, measured at the start and end of medical school.

**Setting** Four Scottish medical schools (universities of Aberdeen, Dundee, Edinburgh and Glasgow).

**Participants** 1512 medical students who attempted (but did not necessarily pass) final written assessment.

**Main outcome measures** The study modelled the change in attainment gap during medical school for four student demographical categories (white/non-white, international/Scottish domiciled, male/female and with/without a known disability) to test whether the attainment gap grew, shrank or remained stable during medical school. Separately, the study modelled the expected versus actual frequency of different demographical groups in the top and bottom decile of the cohort.

**Results** The attainment gap grew significantly for white versus non-white students ( $t(449.39)=7.37$ ,  $p=0.001$ ,  $d=0.49$  and 95% CI 0.34 to 0.58), for internationally domiciled versus Scottish-domiciled students ( $t(205.8)=-7$ ,  $p=0.01$ ,  $d=-0.61$  and 95% CI  $-0.75$  to  $-0.42$ ) and for male versus female students ( $t(1336.68)=3.54$ ,  $p=0.01$ ,  $d=0.19$  and 95% CI 0.08 to 0.27). International, non-white and male students received higher marks than their comparison group at the start of medical school but lower marks by final assessment. No significant differences were observed for disability status. Students with a known disability, Scottish students and non-white students were over-represented in the bottom decile and under-represented in the top decile.

**Conclusions** The tendency for attainment gaps to grow during undergraduate medical education suggests that educational factors at medical schools may—however inadvertently—contribute to DA. It is of critical importance that medical schools investigate attainment gaps within their cohorts and explore potential underlying causes.

## INTRODUCTION

Promoting fairness in assessment is a key priority. Success in medicine should be

## Strengths and limitations of this study

- This the largest study to date investigating longitudinal attainment gaps within undergraduate medical education.
- By evaluating differential attainment longitudinally, the study tests whether attainment gaps are due to pre-existing differences or emerge during medical school.
- The study has sufficient power to detect small/medium effects by pooling data from multiple cohorts and institutions.
- All contributing schools were based in Scotland, and care should be taken when generalising to other contexts.
- The study methodology cannot fully explain the mechanisms behind such attainment gaps

determined by ability rather than background characteristics like ethnicity, sex or socioeconomic status (SES).<sup>1</sup> There is an increasing emphasis on educational processes being ‘fair’ to candidates of diverse backgrounds: besides the legal and regulatory requirements,<sup>2</sup> there is growing acceptance that evaluating fairness should be a routine part of test construction and assessment.<sup>3</sup>

Despite this, candidates continue to experience different outcomes in medical education and training because they have characteristics that lead to them being treated differently by staff, students and patients. The tendency for outcomes to vary in this fashion is usually termed differential attainment (DA). It influences every stage of medical education and is a global phenomenon with similar problems manifesting in a range of contexts.<sup>4,5</sup> The varying treatment of some groups influences the likelihood of candidates completing medical school and affects selection methods.<sup>6–8</sup> Performance on measures of success at or just beyond



graduation shows a similar pattern,<sup>9 10</sup> and for example, ethnically white UK graduates are given higher marks than non-white UK graduates in postgraduate examinations with typically moderate ( $d=0.22$ ) effects.<sup>11</sup> After graduation, ethnically non-white and female doctors experience barriers to success on a range of professional and educational outcomes.<sup>12–14</sup> Students from under-represented backgrounds are substantially less likely to be awarded high ratings from their clerkship directors, less likely to be given honours and less likely to be given honour society membership.<sup>15</sup>

Such compelling evidence has led to calls to establish the mechanisms of DA, but this is challenging. Many historical assumptions—such as the idea that examiners are biased against some candidate groups—remain commonly cited despite evidence to the contrary.<sup>16 17</sup> Examiner bias does not appear to explain DA in postgraduate clinical examinations<sup>18</sup> or written assessment.<sup>19</sup> Qualitative research has emphasised a range of possible factors that can contribute to DA, including trust between trainers and trainees and the process by which those in difficulty are identified and referred to support networks.<sup>20–22</sup> Other research has suggested that unconscious biases may alter training pathways or assessment in the workplace.<sup>4 13 23 24</sup> Some authors now recommend a programmatic approach whereby each component of training is separately reviewed.<sup>25</sup>

As a result, evidence for the existence of DA is very strong, but we have so far only a limited understanding of the mechanisms by which it operates or even whether DA increases or decreases with time spent in medical education. Compounding this, while a great deal of research has been carried out on access to medical school and postgraduate assessment, relatively little work has evaluated DA on assessment *during* medical school. In a large meta-analysis, eleven of fourteen published studies examining undergraduate medical education used a single site, and two of the remaining studies used only two sites.<sup>11</sup> Combined with the tendency to monitor attainment at only a single time point (typically finals), we know little of whether DA is of similar magnitude for different medical schools or remains stable during medical school.

This is an obvious limitation given the role of medical schools in providing the foundation of medical education and training. Due to the diversity of intakes, assessment choices, curriculum design and performance on postgraduate assessment,<sup>26 27</sup> investigating DA at medical schools may help in several ways. By comparing different institutions, the effect of different recruitment strategies, curriculum types and policies on fairness in medical education can be explored. If the magnitude of DA is highly variable across institutions, it argues for a relatively larger role in medical school policy in creating DA. If DA remains consistent despite varying institutional contexts, it argues either that DA is explained by factors outside of medical school control or that no current approaches are identifiably superior or inferior. By examining the data longitudinally, it becomes possible to explore whether DA

increases or decreases over time. If DA is present from the earliest part of medical education, this suggests different mechanisms than if DA is minimally present at the beginning but then grows with time. Such work can therefore significantly improve medical education and support a fairer experience for doctors.

In this study, we used data from four Scottish medical schools operating within a common regulatory framework. Our aim was to evaluate longitudinal DA across undergraduate medical education in 1512 medical students, exploring disability status, domicile, ethnicity and gender. Here, we report on the longitudinal effects of DA for these groups and the impact of DA on student rank.

## METHODS

### Participants

Participants were undergraduate medical students who had attempted (but not necessarily passed) a major written (multiple choice question) assessment near the end of medical school. All institutions operated under the UK medical education system,<sup>2</sup> and new graduates typically embarked on a 2-year foundation training programme as a doctor.

In total, 1512 medical students were eligible for inclusion in the study. To be eligible, a student had to (a) have attempted (but not necessarily passed) the final written assessment, (b) have made the attempt by the end of data collection and (c) have provided demographical information.

The 1512 students represented 74% of all available participants within the period of this study. Excluded subjects were typically those who had exited medical school before final assessment, experienced an interruption of study or intercalated close to the end of the study period and so had not yet sat finals. Due to the complexity of discontinuation, it is theoretically possible for a student to graduate up to 9 years after starting a 5-year programme, which makes confirmation of discontinuation challenging. Candidates who did not attempt final assessment prior to the end of the period of data collection are not included in any analyses presented here.

Table 1 summarises the partner schools, total sample sizes and assessments used. All schools offered 5-year MBChBs (Bachelor of Medicine and Surgery). The first 2 years of each programme involved an introduction to the fundamentals of medicine, anatomy, social issues around healthcare and working with peers. Each programme offered an opportunity to intercalate, whereby candidates spent an additional year studying a topic in greater depth before returning to the core programme. In the later years, candidates rotated through a series of clinical placements to develop the skills and knowledge necessary to work as a junior doctor.

In each school, candidates sat a written assessment at the end of their first year. These featured multiple choice questions (MCQs) and, for two schools, short answer

**Table 1** Participants, data ranges and assessments used

School name	Sample size	Data range	First year assessment	Final assessment
University of Aberdeen	104	2014/2017	MCQ and SAQ	MCQ and SAQ
University of Dundee	202	2013/2016 and 2014/2017	MCQ	MCQ
University of Edinburgh	871	2009/2013, 2010/2014, 2011/2015, 2012/2016 and 2013/2017	MCQ and SAQ	MCQ
University of Glasgow	335	2014/2018 and 2015/2019	MCQ	MCQ

Note: Data range described the first/final year of assessment data for each cohort. ‘Multiple Choice Questions’ (MCQs) require students to select the correct answer from a series of options. ‘Short Answer Questions’ (SAQs) require students to type or write a short answer. All assessments were written rather than clinical.

questions (SAQs). For each question, candidates were presented with a scenario and question. For MCQs, candidates selected the correct answer from a list, whereas for SAQs, candidates provided a short, written answer. The assessment was blueprinted based on programme learning outcomes and standard set by experts familiar with the curriculum.

Near the end of medical school, candidates sat another written assessment. Three schools delivered this in the final year, while one (the University of Aberdeen) delivered it at the very end of the prefinal year. The blueprinting and standard setting process was the same as in the early assessment.

In each case, the assessments acted as a progression barrier: candidates needed to achieve a satisfactory mark to progress to either second year or graduation. A review by the authors identified that although there were some variations in curricula and teaching methods, there were no significant differences in content and structure of assessments between programmes that would impact cross-school comparisons of DA.

**Table 2** describes the participants according to important demographical characteristics. We report whether the candidate did or did not have a known disability, where they were domiciled before starting medical school, their ethnicity and their gender. All recorded data were self-reported. For ethnicity and domicile, we aggregate data across many subcategories into broad groups such as ‘Scottish domicile’ or ‘white.’ While a more detailed breakdown would be helpful, the small numbers in many groups prohibit this. The demographical characteristics selected for study are based partly on the concept of a ‘protected characteristic’ for which there is a legal obligation to promote equality within the UK,<sup>28</sup> partly on demographical characteristics known to be important from past research and partly on availability of data. To give two examples of data availability, marital status and sexual orientation had levels of missingness that were too high to achieve necessary levels of power. The four categories described here (known/no known disability, international, non-EU/Scottish domicile, non-white/white and female/male) represent all those selected for full analysis, and all analyses have sufficient power to detect medium effects. We selected Scottish (as opposed to the

whole UK) domicile due to Scottish-domiciled candidates having already experienced the Scottish legislative and educational framework and having selected a medical school relatively close to home. Furthermore, differences in the funding approach in Scotland compared with the rest of the UK made merging the two groups less defensible. Non-Scottish-domiciled UK students were included in the other comparisons, and so for example, an English-domiciled student who provided valid information on gender would have been reported for that analysis.

SES was recorded in the dataset in two forms. First, candidates had the opportunity to list parental occupation. Over 90% of candidates did not fill this in. A second proxy for SES was candidate postcode, which can be converted into an index of multiple deprivation.<sup>29</sup> However, it was not possible to effectively compare Scottish, non-Scottish UK and international measures of SES within a single dataset. As such we did not explore this covariate further in the present study.

### Data protection and ethics

This project represented a considerable challenge under data protection legislation and required a careful and thorough evaluation of ethical issues. To ensure data protection, a designated team member undertook an honorary contract with each partner and worked in tandem with a data custodian at that school. This meant individualised data were never transferred outside of the school servers, and a thorough anonymisation protocol was used to verify that no ‘unique’ combinations could identify candidates from their data patterns. Ethical approval was granted by the ethics committee for the College of Medicine and Veterinary Medicine at the University of Edinburgh (reference: 2018/7) and then separately approved by an ethics board and a data protection officer at each of the other schools. All participants gave informed consent. Prior to data analysis, all partners agreed to disseminate the results in public and to representatives of the study population: in this case, medical student organisations.

When describing inequities, researchers must ensure individuals are described fairly and appropriately, without discriminatory language. Throughout this paper, we have used language that shows that group membership itself does not *cause* an attainment gap and is never a direct

**Table 2** Demographical characteristics of the study sample

Demographical characteristic	Category	Institution	N	Total n	
Disability	Known disability	Aberdeen	13	102	1512
		Dundee	13		
		Edinburgh	74		
		Glasgow	2		
	No known disability	Aberdeen	91	1410	
		Dundee	189		
		Edinburgh	797		
		Glasgow	333		
Domicile	EU (non-UK)	Aberdeen	2	44	1512
		Dundee	17		
		Edinburgh	14		
		Glasgow	11		
	International	Aberdeen	9	146	
		Dundee	12		
		Edinburgh	88		
		Glasgow	37		
	Rest of the UK	Aberdeen	24	500	
		Dundee	40		
		Edinburgh	354		
		Glasgow	82		
	Scotland	Aberdeen	69	822	
		Dundee	133		
		Edinburgh	415		
		Glasgow	205		
Ethnicity	Non-white	Aberdeen	27	298	1512
		Dundee	21		
		Edinburgh	157		
		Glasgow	93		
	White	Aberdeen	77	1143	
		Dundee	165		
		Edinburgh	665		
	Unknown	Glasgow	236		
		Dundee	16	71	
		Edinburgh	49		
Gender	Female	Aberdeen	67	877	1512
		Dundee	129		
		Edinburgh	480		
		Glasgow	201		
	Male	Aberdeen	37	635	
		Dundee	73		
		Edinburgh	391		
		Glasgow	134		

Candidates of 'unknown' ethnicity, 'EU (non-UK)' and 'Rest of the UK' domicile students are not included in any analyses described in the present study. All demographical characteristics relied on self-report data.

determinant of performance and instead likely reflects systemic societal issues. We have provided some additional references that may be helpful in exploring language

choice when describing historically under-represented groups.<sup>4 20</sup>

### Patient and public involvement

The study was carried out exclusively on medical students and did not involve patients in any way. As such, there was no patient or public involvement.

### Statistical analyses

Each medical school has a locally designed curriculum and assessment environment. We investigate written assessment as the most comparable form of assessment, as the available clinical examinations vary considerably across the schools in both timing and format. To allow like-for-like comparisons across different written assessments, we converted each cohort of data to z-scores.<sup>30</sup>

A z-score is a standardised measurement, where a score of zero indicates the candidate has received exactly the mean mark on the assessment and a score of +/-1 indicates they have received a mark one SD above or below the mean, respectively. This is analytically helpful because it allows for comparisons where relative (rather than absolute) differences are important. If a candidate from one medical school receives a mark of 75 and a candidate from another medical school receives a mark of 70 on two different assessments, it is difficult to know who is more capable. But if the z-score for each candidate is zero, this indicates they are of the same level of ability *relative to their peers* and that they are both average.

We used the Shapiro-Wilk test to model residual values to test for normality.<sup>31</sup> Although the normality parameters were violated ( $W=0.99$  and  $p<0.001$ ), further investigation suggested that parametric testing would still be more appropriate as parametric tests are more effective at minimising the risk of false positives where the group sample sizes and SD vary across groups.<sup>32</sup> Sample sizes were sufficient to detect small effects at 80% power for ethnicity, gender and domicile, whereas for disability status, the unequal group sizes and small numbers of students self-reporting a disability allowed for only medium effects at 80% power.<sup>33</sup> Due to the low sample sizes *within* each medical school, it was not feasible to compare intermedical school variability with sufficient power. Likewise, it was not possible to compare intersectional DA (eg, ethnicity *and* gender). We used Welch's t-test for significance testing as a more robust alternative to other t-tests.<sup>34</sup> All analyses were carried out using R.<sup>35</sup>

### Design choices

We made several design choices that influence the final dataset. Most importantly, by only including candidates who reach final assessment, we exclude the majority of those who experienced major difficulties early in their studies. However, the only alternative is to either measure graduation rates, which prevents granular analyses as the overwhelming majority of students pass medical school,<sup>36</sup> or attempt some form of imputation to estimate final performance of candidates who never



reached that stage of education, with significant uncertainty over the accuracy of such estimates. We opt for a simple approach of reporting data only where fully available. One consequence of this is that variability is higher in final assessment than in first year, with more candidates performing poorly, so most z-score change values were negative. For example, it would be possible for a candidate to receive an A in the first year and an F in the final year and participate in our study, but it would not be possible for the reverse to be true—unless the student successfully resat assessment *and* then completed within the specified timeframe. This can be considered a form of ‘survival bias’, and approaches to the problem always require trade-offs.<sup>37</sup>

To investigate survival bias, we compared the ratios of those who did to those who did not provide final year assessment results for each group. For example, we compared the ratio of non-white/white completers to non-white/white non-completers. No differences in the ratios were detected for any studied group. This likely reflects the fact that non-completion (by the end of the present study) was due to a variety of factors and did not in itself indicate academic difficulty.

Following this, we carried out a number of comparisons. First, we calculated the z-score for each student in their first year and then the final assessment. We explored the equivalence of school. We compared z-score *change* between groups to see whether attainment gaps were growing or shrinking during medical school. Finally, we ranked all candidates to see who would appear in either the top or bottom decile for the final assessment.

### RESULTS

We first tested whether the performance profiles of each school were sufficiently similar to pool data into a single sample. We compared the shapes of the distributions, frequencies of outliers and overall variability of each cohort. After confirming the equivalence of the cohorts, we pooled all data into a combined sample of 1512 students.

Table 3 provides a summary of (a) the z-score for each demographical characteristic per assessment, (b) the relative change in z-score over time and (c) whether the z-score change *between* groups is significant. For the present study, we are not interested in the attainment gap at either the start or end of medical school—but whether the magnitude of the gap changes over time. We found that the gap grew significantly for white versus non-white students ( $t(449.39)=7.37$ ,  $p=0.001$ ,  $d=0.49$  and 95% CI 0.34 to 0.58), for internationally domiciled versus Scottish-domiciled students ( $t(205.8) = -7$ ,  $p=0.01$ ,  $d=0.61$  and 95% CI  $-0.75$  to  $-0.42$ ) and for male versus female students ( $t(1336.68)=3.54$ ,  $p=0.01$ ,  $d=0.19$  and 95% CI 0.08 to 0.27). No significant differences were observed for candidates with versus without a known disability.

**Table 3** Z-score change during medical school study

Demographical characteristic	Category	First assessment (mean)	First assessment (SD)	Final assessment (mean)	Final assessment (SD)	Change (mean)	Change (SD)	Significance/CI
Disability	Known disability	-0.15	0.94	-0.38	0.73	-0.18	0.93	
	No known disability	0.09	0.89	-0.05	0.93	-0.1	0.95	
Domicile	International	0.46	0.83	-0.4	0.92	-0.57	0.92	$(-0.75$ to $-0.42)^*$
	Scotland	-0.08	0.91	-0.05	0.9	0.01	0.97	
Ethnicity	Non-white	0.15	0.93	-0.34	1.06	-0.45	0.96	$(0.34$ to $0.58)^*$
	White	0.04	0.89	-0.02	0.88	0	0.92	
Gender	Female	0.03	0.89	-0.01	0.88	-0.03	0.93	$(0.08$ to $0.27)^*$
	Male	0.14	0.89	-0.2	0.98	-0.2	0.97	

\*Indicates statistical significance at  $p=0.001$ . 95% CIs are given for significant results. For model values, see text. Statistical significance indicates the relative attainment gap between categories changed significantly during the course of the study.

**Table 4** Rankings of top and bottom decile by demographical characteristic

Demographical characteristic	N	Category	N category	Percentage	Decile	N in decile	Expected percentage	Actual percentage
Disability	1512	Known disability	102	6.75	1	5	0.68	0.33
					10	14	0.93	
		No known disability	1410	93.25	1	145	9.32	9.59
					10	136	8.99	
Domicile	968	International	146	9.66	1	21	0.97	1.39
					10	19	1.26	
		Scotland	822	54.37	1	78	5.44	5.16
					10	85	5.62	
Ethnicity	1441	Non-white	298	19.71	1	24	1.97	1.59
					10	54	3.57	
		White	1143	75.6	1	115	7.56	7.61
					10	92	6.08	
Gender	1512	Female	877	58	1	85	5.8	5.62
					10	81	5.36	
		Male	635	42	1	65	4.2	4.3
					10	69	4.56	

N indicates the total sample size for that characteristic, while N category indicates the sample size for the individual category. Percentage indicates the proportion of students from that category in the overall sample. Decile 1 is the highest (ie, best performing) decile; decile 10 is the lowest (ie, worst performing) decile. N in decile gives the number of candidates who actually appeared in that decile, and the difference between the expected and actual percentage shows whether the category is over-represented or under-represented.

For the three significant analyses, non-white, internationally domiciled and male candidates were awarded a relatively higher score at the start of medical school. By the end of medical school, they were respectively awarded a lower score than white, Scottish-domiciled and female students. The effect size was medium when testing ethnicity and domicile and small for testing gender. In summary, non-white, internationally domiciled and male students experienced a relative decline in their achieved marks at medical school, which cannot be explained by low attainment before or in the first year of medical school.

Finally, we estimated how often medical students of different demographics would appear in the top and bottom decile based on their z-scores versus their expected frequencies based purely on how many existed in each category. Table 4 summarises the details.

Decile 1 is the highest-scoring decile, and decile 10 is the lowest-scoring decile. Students with a known disability, Scottish students and non-white students are over-represented in the bottom decile and under-represented in the top decile. Students with no known disability and white students are over-represented in the top decile and under-represented in the bottom decile. International students and male students are over-represented in *both* the top and bottom decile. Female students are under-represented in the top and bottom decile.

This analysis shows that many groups exhibit differences not just in mean performance but also in variability, with some candidates being under-represented and over-represented at the extremes of the distribution.

## DISCUSSION

### Statement of principal findings

DA exists within Scottish medical schools, with small to medium effects. The analysis described here demonstrates both the considerable difficulty in organising datasets to longitudinally investigate DA and the ongoing importance of such work. Even among successful medical students—and the overwhelming majority of those described in the present dataset have become doctors—DA exists. The fact that many attainment gaps grow during medical school suggests educational factors within medical schools may promote DA.

### Strengths and weaknesses of the study

It is important not to overstate the findings. Small to medium effect sizes are consequential and impact student education, but there remains considerable variance between students of all groups. In this dataset, candidates across the attainment continuum were present in every group. In addition, the core purpose of medical education—graduating a safe doctor—has been met for almost all participants in the dataset. The gaps observed here must be placed in this context. Finally, as until

now we have operated in an environment with almost no published data, there is a risk that organisations that attempt to directly engage with the problem of DA are criticised for the differences they reveal, which may in turn drive reluctance to explore the issue in depth. It is important that stakeholders support the exploration of DA across the sector.

This study represents a novel attempt to understand DA not as a fixed factor, but as a changing influence on student performance and behaviour. The sample size and range suggest we can be confident the findings are potentially generalisable to other UK medical schools. By opting for a straightforward methodology, we believe the findings are robust and can inform future policy.

Despite this, there are limitations. The challenges of organising a longitudinal study using data from a range of institutions with varying outcome measures should not be understated. We have made design choices—such as excluding those who failed before reaching finals—which may influence the pattern of results. Due to the relatively small sample sizes of some groups, it was not possible to explore ‘intersectional’ DA for, for example, candidates who were non-white and female.<sup>38</sup> Due to the nature of the available data on SES, we were not able to include SES as a covariate in the present study. All candidate demographics were self-reported, and so, some information could theoretically be inaccurate. While we consider the curricula and assessment of the institutions to be sufficiently similar to allow for a combined analysis, it is possible that local factors may have created some unidentified sources of variance.

The lack of a shared, standardised assessment across schools required the use of z-scores (or an equivalent method), and the presence of a standardised assessment, such as the forthcoming UK Medical Licensing Assessment, would have greatly simplified the analysis.<sup>39</sup>

Data collection was challenging, and it was clear that there was no expectation during data creation that assessment-level data would be required 5 or 10 years after the assessment was sat. Medical education data should be thought of as ‘perishable’—it is possible that even relatively recent data are being lost, overwritten or rendered inaccessible. If medical educators wish to investigate DA across time, it is critical that better data collection practices are implemented, and historic data sources should be secured and documented in national-level databases.<sup>40</sup> The alternative is that we may establish excellent prospective analyses for which we will have no useful data for up to a decade.

### Comparison with other studies and unanswered questions

DA exists across medical education systems across the world and should always be considered when designing teaching and assessment.<sup>4 5</sup> Our findings support and extend past work exploring DA in postgraduate medical education<sup>9 12 13 21</sup> and at medical school.<sup>15 24</sup> Importantly, our study also confirms that we remain unclear, as a sector, on the mechanisms behind DA.<sup>18 19</sup> All organisations

involved in medical education must proactively consider how they approach fairness in medical education and evaluate the impact of DA.

The limitations described above are logical opportunities for future work. Exploring the impact of SES, analysing intersectional characteristics and studying those who do not graduate may offer insights into both the scope and mechanisms of DA. Exploring candidate domicile in a more granular fashion (such as measuring the distance between home and their selected medical school) may be helpful, especially alongside measurements of SES. Importantly, the design challenges highlighted here will persist until institutions develop rigorous frameworks to investigate long-term changes in student performance.

### IMPLICATIONS AND CONCLUSIONS

The present study demonstrates DA changes in magnitude during undergraduate medical education. Combined with evidence that candidates of some groups are less likely to be given awards<sup>15</sup> and more likely to experience prejudice,<sup>24</sup> it is very plausible that some of the mechanisms of DA are located in, or caused by, aspects of medical education within medical schools. As such, institutions must consider the possibility that their actions contribute to DA and develop appropriate policies for investigation and correction.<sup>14</sup>

**Twitter** Eleanor J Hothersall @e\_hothersall

**Contributors** Dr IC, Dr EJH and Professor JPL were each responsible for sourcing data, describing the context and exploring the results in their institutions. AD was responsible for sourcing data at her institution and then collating all the data and running the initial analyses. Dr DH organised the project, designed the analyses, was primarily responsible for writing the paper and is the guarantor for the content. Dr AJ acted as supervisor for all the project work and reviewed the analyses. All authors have separately reviewed the manuscript and provided input in developing the final analyses and paper.

**Funding** The Scottish Medical Education Research Consortium (SMERC) provided funding to allow the research project to take place. The funding was used to pay for administrator and researcher time to collate and analyse the data. The funder had no direct input into the analyses chosen or the reporting of the results. The researchers were independent from the funder, and all researchers had access to the data and can take responsibility for the integrity of the data and the accuracy of the data analysis.

**Competing interests** All authors have completed the ICMJE Uniform Disclosure Form at [www.icmje.org/doi\\_disclosure.pdf](http://www.icmje.org/doi_disclosure.pdf) and declare that all authors had financial support from the Scottish Medical Education Research Consortium (SMERC) for the submitted work, no financial relationships with any organisations that might have an interest in the submitted work in the previous three years and no other relationships or activities that could appear to have influenced the submitted work.

**Patient and public involvement** Patients and/or the public were not involved in the design, conduct, reporting or dissemination plans of this research.

**Patient consent for publication** Not required.

**Ethics approval** Ethical approval was granted by the ethics committee for the College of Medicine and Veterinary Medicine at the University of Edinburgh (reference: July 2018) and then separately approved by an ethics board and a data protection officer at each of the other schools. All participants gave informed consent. Prior to data analysis, all partners agreed to disseminate the results in public and to representatives of the study population: in this case, medical student organisations. This information is reproduced in the main text.

**Provenance and peer review** Not commissioned; externally peer reviewed.

**Data availability statement** No data are available. Due to the sensitivity of the dataset—including confidential information on student demographics and assessment scores—we are unable to share raw data.

**Open access** This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

#### ORCID iD

David Hope <http://orcid.org/0000-0001-6623-2857>

## REFERENCES

- Tsouroufli M, Malcolm I, Equality MI. Equality, diversity and fairness in medical education: international perspectives. *Med Educ* 2015;49:4–6.
- General Medical Council. *Outcomes for graduates*. Manchester: General Medical Council, 2015.
- American Educational Research Association, American Psychological Association, National Council on Measurement in Education. *Standards for educational and psychological testing*. Washington, DC: AERA, 2014.
- Murphy M. Teaching and learning about sexual diversity within medical education: the promises and pitfalls of the informal curriculum. *Sex Res Soc Policy* 2019;16:84–99.
- Prideaux D, Roberts C, Eva K, et al. Assessment for selection for the health care professions and specialty training: consensus statement and recommendations from the Ottawa 2010 conference. *Med Teach* 2011;33:215–23.
- O'Neill L, Hartvigsen J, Wallstedt B, et al. Medical school dropout-testing at admission versus selection by highest grades as predictors. *Med Educ* 2011;45:1111–20.
- Patterson F, Knight A, Dowell J, et al. How effective are selection methods in medical education? A systematic review. *Med Educ* 2016;50:36–60.
- Cliffordson C. Selection effects on applications and admissions to medical education with regular and Step-Wise admission procedures. *Scandinavian Journal of Educational Research* 2006;50:463–82.
- MacKenzie RK, Cleland JA, Ayansina D, et al. Does the UKCAT predict performance on exit from medical school? a national cohort study. *BMJ Open* 2016;6:e011313.
- Pershing S, Co JPT, Katznelson L. The new USMLE step 1 paradigm: an opportunity to Cultivate diversity of excellence. *Acad Med* 2020;95:1325–8.
- Woolf K, Potts HWW, McManus IC. Ethnicity and academic performance in UK trained doctors and medical students: systematic review and meta-analysis. *BMJ* 2011;342:d901.
- Linton S. Taking the difference out of attainment. *BMJ* 2020;368:m438.
- Klein R, Julian KA, Snyder ED, et al. Gender bias in resident assessment in graduate medical education: review of the literature. *J Gen Intern Med* 2019;34:712–9.
- Ufomata E, Merriam S, Puri A, et al. A policy statement of the Society of general internal medicine on tackling racism in medical education: reflections on the past and a call to action for the future. *J Gen Intern Med* 2021;36:1077–81.
- Teherani A, Hauer KE, Fernandez A, et al. How small differences in assessed clinical performance amplify to large differences in grades and awards: a cascade with serious consequences for students underrepresented in medicine. *Acad Med* 2018;93:1286–92.
- Woolf K. Differential attainment in medical education and training. *BMJ* 2020;368:m339.
- Yeates P, Woolf K, Benbow E, et al. A randomised trial of the influence of racial stereotype bias on examiners' scores, feedback and recollections in undergraduate clinical exams. *BMC Med* 2017;15:1–11.
- MR JUSTICE MITTING. The Queen on the application of Bapio Action Ltd [Ciamant] v Royal College of General Practitioners [First Defendant] and General Medical Council [Second Defendant], in the High Court of Justice, Queen's Bench Division, The Administrative Court. 10th April 2014. EWHC 1416 (Admin) 2014, 2014. Available: <http://www.rcgp.org.uk/news/2014/may/~media/Files/News/Judicial-Review-Judgment-14-April-2014.ashx>
- Hope D, Adamson K, McManus IC, et al. Using differential item functioning to evaluate potential bias in a high stakes postgraduate knowledge based assessment. *BMC Med Educ* 2018;18:64.
- Cleland J, Fahey Palma T, Palma TF. "Aspirations of people who come from state education are different": how language reflects social exclusion in medical education. *Adv Health Sci Educ Theory Pract* 2018;23:513–31.
- Woolf K, Rich A, Viney R, et al. Perceived causes of differential attainment in UK postgraduate medical training: a national qualitative study. *BMJ Open* 2016;6:e013429.
- Atewologun D, Kline R, Ochieng M. Fair to refer? reducing disproportionality in fitness to practise concerns reported to the GMC, 2019. Available: <https://www.gmc-uk.org/about/what-we-do-and-why/data-and-research/research-and-insight-archive/fair-to-refer>
- Kristoffersson E, Diderichsen S, Verdonk P, et al. To select or be selected - gendered experiences in clinical training affect medical students' specialty preferences. *BMC Med Educ* 2018;18:268.
- Cheng L-F, Yang H-C. Learning about gender on campus: an analysis of the hidden curriculum for medical students. *Med Educ* 2015;49:321–31.
- Karani R, Varpio L, May W, et al. Commentary: racism and bias in health professions education: how educators, faculty developers, and researchers can make a difference. *Acad Med* 2017;92:S1.
- McManus IC, Elder AT, de Champlain A, et al. Graduates of different UK medical schools show substantial differences in performance on MRCP(UK) part 1, part 2 and PACES examinations. *BMC Med* 2008;6:5.
- Devine OP, Harborne AC, McManus IC. Assessment at UK medical schools varies substantially in volume, type and intensity and correlates with postgraduate attainment. *BMC Med Educ* 2015;15:146.
- Davies C, Ferreira N, Morris A, et al. The equality act 2010: five years on. *International Journal of Discrimination and the Law* 2016;16:61–5.
- Noble M, Wright G, Smith G, et al. Measuring multiple deprivation at the Small-Area level. *Environ Plan A* 2006;38:169–85.
- Abdi H. Z-scores. In: *encyclopedia of measurement and statistics*. Thousand Oaks (CA: Sage, 2007: 1055–8.
- Shapiro SS, Wilk MB. An analysis of variance test for normality (complete samples). *Biometrika* 1965;52:591–611.
- Zimmerman DW. Invalidation of parametric and nonparametric statistical tests by concurrent violation of two assumptions. *The Journal of Experimental Education* 1998;67:55–68.
- Cohen J. A power primer. *Psychol Bull* 1992;112:155–9.
- Delacre M, Lakens D, Leys C. Why psychologists should by default use Welch's t-test instead of student's t-test. *International Review of Social Psychology* 2017;30:92–101.
- Ihaka R, Gentleman R. R: a language for data analysis and graphics. *Journal of Computational and Graphical Statistics* 1996;5:299–314.
- Arulampalam W, Naylor RA, Smith JP. A hazard model of the probability of medical school drop-out in the UK. *Journal of the Royal Statistical Society: Series A* 2004;167:157–78.
- Zhou Z, Rahme E, Abrahamowicz M, et al. Survival bias associated with time-to-treatment initiation in drug effectiveness evaluation: a comparison of methods. *Am J Epidemiol* 2005;162:1016–23.
- Morrison N, Chimkupete P. Double jeopardy: black and female in medicine. *Clin Teach* 2020;17:566–8.
- Archer J, Lynn N, Coombes L, et al. The medical licensing examination debate. *Regul Gov* 2017;11:315–22.
- Dowell J, Cleland J, Fitzpatrick S, et al. The UK medical education database (UKMED) what is it? why and how might you use it? *BMC Med Educ* 2018;18:1–8.