

University of Groningen

Role of knowledge and reasoning processes as predictors of resident physicians' susceptibility to anchoring bias in diagnostic reasoning

Mamede, Sílvia; Zandbergen, Adrienne; de Carvalho-Filho, Marco Antonio; Choi, Goda; Goeijenbier, Marco; van Ginkel, Joost; Zwaan, Laura; Paas, Fred; Schmidt, Henk G

Published in:
BMJ Quality & Safety

DOI:
[10.1136/bmjqs-2023-016621](https://doi.org/10.1136/bmjqs-2023-016621)

IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

Document Version
Publisher's PDF, also known as Version of record

Publication date:
2024

[Link to publication in University of Groningen/UMCG research database](#)

Citation for published version (APA):

Mamede, S., Zandbergen, A., de Carvalho-Filho, M. A., Choi, G., Goeijenbier, M., van Ginkel, J., Zwaan, L., Paas, F., & Schmidt, H. G. (2024). Role of knowledge and reasoning processes as predictors of resident physicians' susceptibility to anchoring bias in diagnostic reasoning: a randomised controlled experiment. *BMJ Quality & Safety*, 33, 563-572. <https://doi.org/10.1136/bmjqs-2023-016621>

Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).








The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

Role of knowledge and reasoning processes as predictors of resident physicians' susceptibility to anchoring bias in diagnostic reasoning: a randomised controlled experiment

Sílvia Mamede ¹, Adrienne Zandbergen ²,
 Marco Antonio de Carvalho-Filho ³, Goda Choi,⁴ Marco Goeijenbier,^{5,6}
 Joost van Ginkel ⁷, Laura Zwaan ¹, Fred Paas ⁸,
 Henk G Schmidt ⁸

► Additional supplemental material is published online only. To view, please visit the journal online (<https://doi.org/10.1136/bmjqs-2023-016621>).

For numbered affiliations see end of article.

Correspondence to

Dr Sílvia Mamede;
s.mamede@erasmusmc.nl

Received 8 August 2023
 Accepted 26 January 2024
 Published Online First
 16 February 2024



► <http://dx.doi.org/10.1136/bmjqs-2024-017141>



© Author(s) (or their employer(s)) 2024. No commercial re-use. See rights and permissions. Published by BMJ.

To cite: Mamede S, Zandbergen A, de Carvalho-Filho MA, et al. *BMJ Qual Saf* 2024;**33**:563–572.

ABSTRACT

Background Diagnostic errors have been attributed to reasoning flaws caused by cognitive biases.

While experiments have shown bias to cause errors, physicians of similar expertise differed in susceptibility to bias. Resisting bias is often said to depend on engaging analytical reasoning, disregarding the influence of knowledge. We examined the role of knowledge and reasoning mode, indicated by diagnosis time and confidence, as predictors of susceptibility to anchoring bias. Anchoring bias occurs when physicians stick to an incorrect diagnosis triggered by early salient distracting features (SDF) despite subsequent conflicting information.

Methods Sixty-eight internal medicine residents from two Dutch university hospitals participated in a two-phase experiment. Phase 1: assessment of knowledge of discriminating features (ie, clinical findings that discriminate between lookalike diseases) for six diseases. Phase 2 (1 week later): diagnosis of six cases of these diseases. Each case had two versions differing exclusively in the presence/absence of SDF. Each participant diagnosed three cases with SDF (SDF+) and three without (SDF–). Participants were randomly allocated to case versions. Based on phase 1 assessment, participants were split into higher knowledge or lower knowledge groups. Main outcome measurements: frequency of diagnoses associated with SDF; time to diagnose; and confidence in diagnosis. **Results** While both knowledge groups performed similarly on SDF– cases, higher knowledge physicians succumbed to anchoring bias less frequently than their lower knowledge counterparts on SDF+ cases ($p=0.02$). Overall, physicians spent more time ($p<0.001$) and had lower confidence ($p=0.02$) on SDF+ than SDF– cases ($p<0.001$). However, when diagnosing SDF+ cases, the groups did not differ in time ($p=0.88$) nor in confidence ($p=0.96$).

Conclusions Physicians apparently adopted a more analytical reasoning approach when presented with distracting features, indicated by increased time and lower confidence, trying to combat bias. Yet, extended deliberation alone did not explain the observed

WHAT IS ALREADY KNOWN ON THIS TOPIC

⇒ Experimental research has shown cognitive biases to cause diagnostic error, but vulnerability to bias differs substantially even among physicians at similar level of training and clinical experience. Predictors of susceptibility to bias have been much debated.

WHAT THIS STUDY ADDS

⇒ Knowledge of features that discriminate between lookalike diseases was the major predictor of anchoring bias in diagnostic reasoning. Engaging in more analytical reasoning in itself was not sufficient to overcome bias, which instead depended on possessing the relevant diagnostic knowledge.

HOW THIS STUDY MIGHT AFFECT RESEARCH, PRACTICE OR POLICY

⇒ The quest for ways to prevent cognitive bias and diagnostic error, hitherto focused on improving physicians' reasoning processes, should be redirected for the development of effective strategies for learning the critical knowledge to avoid errors throughout undergraduate and graduate training.

performance differences between knowledge groups. Success in mitigating anchoring bias was primarily predicted by knowledge of discriminating features of diagnoses.

BACKGROUND

Five to fifteen per cent of all diagnoses are estimated to be wrong^{1 2} often with serious consequences to patients.^{3 4} Several strategies have been proposed to reduce diagnostic error but with limited success.^{5 6} Most diagnostic errors involve faults in physicians' reasoning,^{7 8} but the origin of these faults is unclear. They have often been attributed to the use of heuristics, shortcuts in reasoning used by physicians to make routine judgements.^{9 10} Heuristics are typically correct and efficient but can induce biases.^{3 11}

Some experimental evidence exists that bias can cause errors. For example, recent experiences with a disease led physicians to confuse a subsequent look-alike (but different) case with the disease seen before in a demonstration of 'availability bias'.^{12 13} Bias was also induced by 'salient distracting features' (hereafter, SDF) encountered early in a clinical case.^{14 15} SDFs are clinical findings that are irrelevant to the current problem but draw the physicians' attention, because they are strongly associated with a particular disease that seems a plausible diagnosis at first glance. SDF amplified diagnostic mistakes, as physicians often adhered to the (erroneous) diagnosis prompted by these features, failing to revise it when contradictory evidence emerged. This phenomenon, identified as 'anchoring bias' within the medical literature, has been associated with premature conclusion and is frequently cited as a major contributor to cognitive diagnostic errors.^{10 16 17}

Noteworthy, while most physicians indeed fell prey to bias in these experimental studies, a substantial fraction of them did not.^{12 13} Physicians at similar expertise level (as measured by training and years of clinical experience) differed in their ability to overcome the influence of the bias. Better understanding the sources of these differences may help develop strategies to counteract bias. A prominent view in the medical literature is that variation in resistance to bias is primarily determined by differences in diagnostic reasoning *processes*.⁹ Consequently, physicians are taught reasoning strategies to avoiding cognitive biases.^{9 18} Conversely, a different perspective assumes that what predicts susceptibility to bias is primarily disease knowledge, that is, knowledge of the associations between each disease and its signs and symptoms as stored in the physician's memory.¹⁹⁻²¹ Particularly important is probably knowledge of features that help discriminate between the disease and other usual plausible diagnoses for a patient with a particular clinical presentation. Strong knowledge of these *discriminating features* would make a physician less likely to overlook them when irrelevant cues become salient in bias-inducing circumstances.²¹

This hypothesis is not easily investigated because it requires measurements of *specific* rather than general knowledge as well as indicators of reasoning processes. Nevertheless, it has some preliminary

support. A recent study requested physicians to diagnose cases under conditions that induce availability bias and evaluated their knowledge of discriminating features.²² Lower knowledge physicians succumbed to the availability bias considerably more frequently than their higher knowledge counterparts. This observation cannot be explained by increased engagement in analytical reasoning (which would necessitate more time), given that both groups dedicated similar time-frames to diagnosis. The level of knowledge apparently predicted susceptibility to bias. These findings, however, are preliminary. The knowledge evaluation task, conducted immediately after the diagnostic task, may have been influenced by it, and time was the sole metric of reasoning processes. Moreover, the study investigated availability bias, and other types of biases are also deemed relevant.

The present study examined the role of knowledge of discriminating features in counteracting anchoring bias induced by SDF. Disease knowledge was evaluated 1 week *before* the physicians diagnosed clinical cases with and without SDF that have been shown to induce anchoring bias.^{14 15} Time spent in diagnosis and ratings of confidence in the diagnosis were obtained as indicators of reasoning mode. We hypothesised that knowledge of discriminating features would be the primary predictor of susceptibility to anchoring bias. Therefore, physicians who had more knowledge of discriminating features were expected to make fewer diagnostic errors related to anchoring bias without differences in the measures of reasoning mode.

METHODS

Study design

The study was a two-phase experiment consisting of: (1) evaluation of the physicians' knowledge about discriminating features of the diseases tested in the study, and (2) (1 week later) diagnosis of clinical cases of these diseases. The cases were manipulated so that SDFs were either present (SDF+ cases) or absent (SDF- cases). Two case sets were prepared alternating the manipulation across cases. Participants were randomly assigned to one of the sets in a counterbalanced within-subjects incomplete block design. We measured time required to diagnose and confidence in the diagnosis. In psychological research on reasoning, moving from intuitive towards more analytical reasoning is associated with increased response time²³⁻²⁵ and decreased confidence.^{26 27} Figure 1 presents the study design.

The study protocol was pre-registered on the Open Science Framework (online supplemental file 1). We followed the Consolidated Standards of Reporting Trials reporting guidelines as far as they applied to the present experiment.

Setting and participants

Participants were residents in internal medicine (In the country, residents in internal medicine are physicians

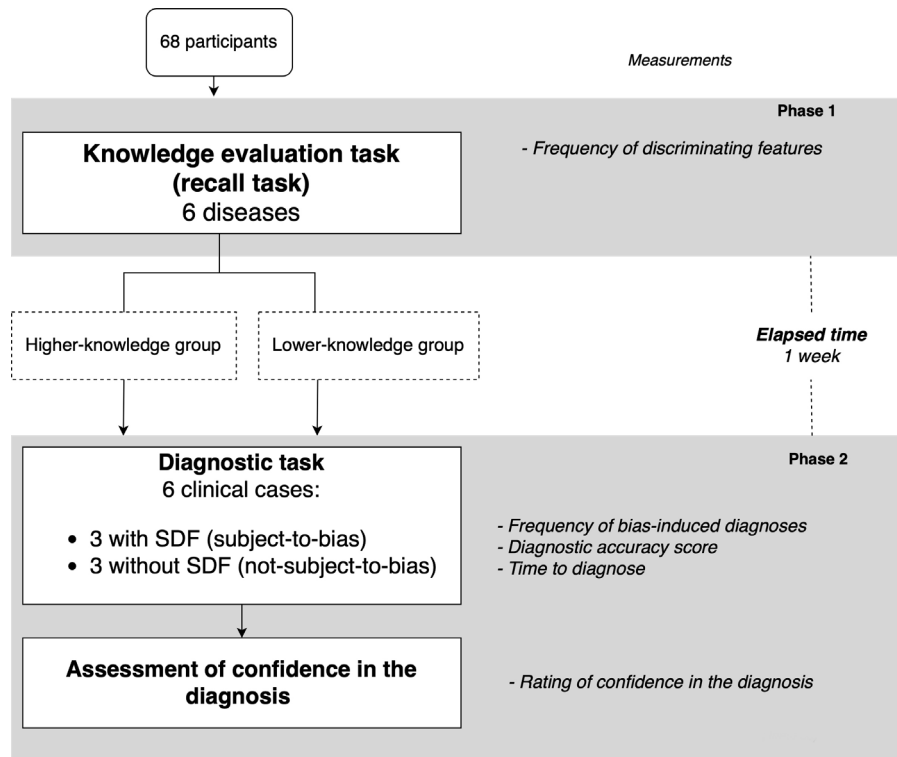


Figure 1 Diagram of the study design. SDF, salient distracting features.

in training to become internists after having graduated from medical schools which have a 6-year curriculum) at Erasmus Medical Center (Erasmus MC) and at the University Medical Centre Groningen in the Netherlands. Residents with at least 1 year of clinical experience were invited by the programme director to volunteer for the study. A €50/participant donation was made to the residents' association to organise extracurricular scientific activities. The study took place during regular educational days in 2021–2022, online at Erasmus MC and face to face in Groningen. (Because patterns of performance were similar regardless of session format, participants were aggregated for the analysis.) Written consent was obtained from all physicians.

Sample size estimation

A priori power analysis was performed, assuming a medium effect size (Cohen's $f=0.30$), as suggested by previous research,¹⁵ and the standard alpha level of 0.05. It indicated that a sample size of 62 participants would be sufficient to achieve a power of 0.80.

Materials and procedure

In phase 1, the physicians performed a recall task to assess knowledge of findings associated with each of the six diseases included in the diagnostic task (phase 2). A recall task assesses mental representations of disease knowledge by measuring how much about it

participants are able to recollect from memory under restricted time conditions.²⁸ The six tested diseases (see online supplemental file 2) were mixed with three unrelated diseases and presented one by one in random order on a computer screen. For each disease, the physician typed everything (s)he could remember about clinical findings (In a broad sense, including findings referring to medical history, complaints, results of physical examination and diagnostic tests) that were critical for the diagnosis of that particular disease and helped distinguish between the disease and other usual diagnoses. Based on a pilot with participants not involved in the study, a maximum time of 2.5 min was allocated for each disease.

In phase 2, six clinical cases were diagnosed (see online supplemental file 2). The cases were developed by board-certified internists and used in prior studies with internal medicine residents.^{12 15} None of the cases involved real patients' information. Each case was prepared in two versions which differed exclusively on the presence or absence of SDF (see table 1 for an example) but had the same most likely diagnosis. This manipulation has been used in previous studies.¹⁵

Two sets of cases were prepared, alternating the manipulation across the cases. The cases were presented one by one in random order on the Qualtrics research suite (Qualtrics, Provo, Utah), which automatically randomly allocated each participant to one of the case sets and registered responses and response time. The participant was asked to read the case and type the most likely diagnosis as fast as possible without

Table 1 Example of a clinical case (diagnosis: vitamin B₁₂ deficiency) used in the study with the version with and without salient distracting features (SDF)

Version without SDF	Version with SDF
<p>A 62-year-old man with uncomplicated diabetes in his otherwise unremarkable medical history was brought to the clinic by his son. The patient has been referred to your clinic because of increasing changes in mental state associated with emotional lability for 4 months. He also had problems with walking and reduced vision during this period. The patient reports sexual contacts only with his wife.</p> <p>Medication: NPH insulin.</p> <p>Health-related behaviours: social alcohol use, not drugs.</p> <p>Physical examination: Patient is in a wheelchair, cries often and does not understand much. BP 140/90 mm Hg; pulse 88/min; temperature 36.2°C. No signs of dehydration.</p> <p>Neurological examination: Patient walks slowly and with long strides; ataxia; difficulties keeping balance with open and closed eyes. Muscle power and reflexes are normal; Babinsky negative.</p> <p>Fundoscopy examination: no abnormalities.</p> <p>Further physical examination: no abnormalities.</p> <p>Diagnostic tests: Hb 5.6 mmol/L; Ht 32.9%; MCV 112 µm³; macrocytosis; leucocytes 5.2×10⁹/L with normal differentiation and hypersegmentation of the polymorphonuclear granulocytes; thrombocytes 154×10⁹/L; reticulocytes <1%; glucose 7.27 mmol/L; urea 9.3 mmol/L; creatinine 79.5 µmol/L; electrolytes, liver functions and blood gases: normal; TSH 5 mU/L.</p> <p>Antibodies to syphilis and HIV: negative.</p> <p>Cerebrospinal fluid analysis: no abnormalities.</p> <p>CT scan and MRI skull: no abnormalities.</p>	<p>A 62-year-old man with uncomplicated diabetes in his otherwise unremarkable medical history was brought to the clinic by his son who reports that 4 months ago, before the onset of the present complaints, his father <u>fell from his own height</u>, without loss of conscience. His son was very worried because <u>his grandfather died of dementia at the age of 67</u>. The patient has been referred to your clinic because of increasing changes in mental state associated with emotional lability for 4 months. He also had problems with walking and a reduced vision during this period. The patient reports sexual contacts only with his wife.</p> <p>Medication: NPH insulin.</p> <p>Health-related behaviours: social alcohol use, not drugs.</p> <p>Physical examination: Patient is in a wheelchair, cries often and does not understand much. BP 140/90 mm Hg; pulse 88/min; temperature 36.2°C. No signs of dehydration.</p> <p>Neurological examination: Patient walks slowly and with long strides; ataxia; difficulties keeping balance with open and closed eyes. Muscle power and reflexes are normal; Babinsky negative.</p> <p>Fundoscopy examination: no abnormalities.</p> <p>Further physical examination: no abnormalities.</p> <p>Diagnostic tests: Hb 5.6 mmol/L; Ht 32.9%; MCV 112 µm³; macrocytosis; leucocytes 5.2×10⁹/L with normal differentiation and hypersegmentation of the polymorphonuclear granulocytes; thrombocytes 154×10⁹/L; reticulocytes <1%; glucose 7.27 mmol/L; urea 9.3 mmol/L; creatinine 79.5 µmol/L; electrolytes, liver functions and blood gases: normal; TSH 5 mU/L.</p> <p>Antibodies to syphilis and HIV: negative.</p> <p>Cerebrospinal fluid analysis: no abnormalities.</p> <p>CT scan and MRI skull: no abnormalities.</p>
<p>SDFs are underlined.</p> <p>BP, blood pressure; Hb, haemoglobin; Ht, haematocrit; MCV, Mean Corpuscular Volume; NPH, Neutral Protamine Hagedorn; TSH, Thyroid Stimulating Hormone.</p>	

compromising accuracy.^{12 13} Each participant diagnosed three SDF+ cases and SDF− cases, but which case was diagnosed with and without SDF varied for each participant depending on the case set to which s/he had been assigned. After diagnosing all cases, each case was presented again, with only its initial sentence and the diagnosis given by the participant, who was asked to mark, on a 0–100% scale, how confident s/he was that the diagnosis was correct. Subsequently, the participants provided demographic information, answered probing questions about the study and received feedback on the correct diagnoses. (Online supplemental file 2 provides additional information.)

Outcome measurements

Knowledge was measured (phase 1) by the frequency of discriminating features mentioned by the participant in the recall task. Four board-certified internists (AZ, GC, MAdC-F, MG) worked independently to assign each clinical finding mentioned by the participants for each disease to one of three categories: ‘discriminating feature’, ‘correct finding but not critical for the diagnosis’ or ‘incorrect finding’. The internists formed triads, and the most frequent category attributed to each finding was used. For each participant, we summed up the number of discriminating features mentioned, and descriptive statistics were computed. Similar to what happens when students’ outcomes in tests are based on relative standard setting methods,²⁹ participants were assigned to either a lower knowledge or a higher knowledge group based on the median. (See online supplemental file 2 for further details.)

Our main outcomes of interest were: frequency of incorrectly giving the diagnosis associated with the SDF as the most likely diagnosis; diagnostic accuracy score; time to diagnose the case; and rating of confidence in the diagnosis. (Note that while the first outcome measures only errors *linked* to the SDF, the second outcome refers to any type of error.) These measurements were computed for SDF+ cases and for SDF− cases. First, the diagnoses given by the participants were categorised as either correct (accuracy score 1.0); partially correct (accuracy score 0.5); incorrect—associated with the SDF (accuracy score 0); and incorrect—not associated with the SDF (accuracy score 0). Triads of internists categorised the diagnoses, following a procedure similar to the one described for the knowledge evaluation task. The mean frequency of incorrect diagnoses associated with the SDF on SDF+ and SDF− cases was computed. Notice that comparing these two types of cases is critical to determine whether the incorrect diagnosis associated with the SDF can actually be attributed to bias. For example, dementia would be a plausible diagnosis in a patient with the clinical presentation displayed in [table 1](#) even in the absence of the SDF inserted in one of the versions. Therefore, only if the frequency of dementia among the diagnoses *increased* when the SDF were present, anchoring bias was considered to have occurred.

Missing data

Several variables in the data had missing values. Missing data were handled using multiple imputation,³⁰ creating 100 imputed datasets, by using the mice package in R.³¹ The dataset roughly contained three types of variables: numerical variables, numerical variables that were computed using other variables and categorical variables. Missing data on numerical variables were imputed using predictive mean matching,³² numerical variables computed from other variables were imputed using passive imputation³³ and categorical variables were imputed using a (multinomial) logistic regression model. Both outcomes and predictors were imputed. In total, the data contained 80 variables, all of which were used in the imputation process, either to be imputed, to serve as predictors for the missing values on other variables or both. Detailed information on the choices and procedures for the imputation are provided in online supplemental file 2.

Statistical analysis

Separate mixed analyses of variance with knowledge group (lower knowledge vs higher knowledge) as between-subjects factor and exposure to SDF (SDF+ and SDF-) as within-subjects factor were performed on the mean frequency of incorrect diagnoses associated with the SDF, mean diagnostic accuracy scores, mean time spent in diagnosis and mean confidence in the diagnosis. These analyses were performed both on the multiple imputation dataset and on the observed data (complete case analysis). Online supplemental file

2 provides additional information on the statistical analysis and its rationale.

RESULTS

Seventy-five physicians attended the two phases of the study. Seven were removed from the analysis (five denied consent to use their data; two resumed the tasks repeatedly due to interruptions), resulting in 68 participants.

The complete case analysis is presented in online supplemental file 2. The number of years in clinical practice did not differ between lower knowledge and higher knowledge groups (mean (SD), respectively: 2.71 (1.21); 2.70 (1.39); $p=0.98$). The complete case analysis yielded results similar to the analyses with imputation for missing data, which are presented below. (In all figures error bars represent 95% CIs.)

The frequency with which the diagnosis associated with the SDF was incorrectly given as the most likely diagnosis in SDF+ cases and SDF- cases is presented in figure 2. Overall, as expected, the error occurred more frequently when the SDFs were present in the case (ie, the case was subject to bias) than when they were absent ($p<0.001$), and physicians with higher knowledge made fewer errors than physicians with lower knowledge ($p=0.015$). However, a significant interaction effect ($p=0.020$) was found, with the difference in performance showing only on SDF+ cases. While the two groups performed similarly in the absence of SDF, higher knowledge physicians fell to anchoring bias less frequently when SDFs were present than their colleagues with lower knowledge.

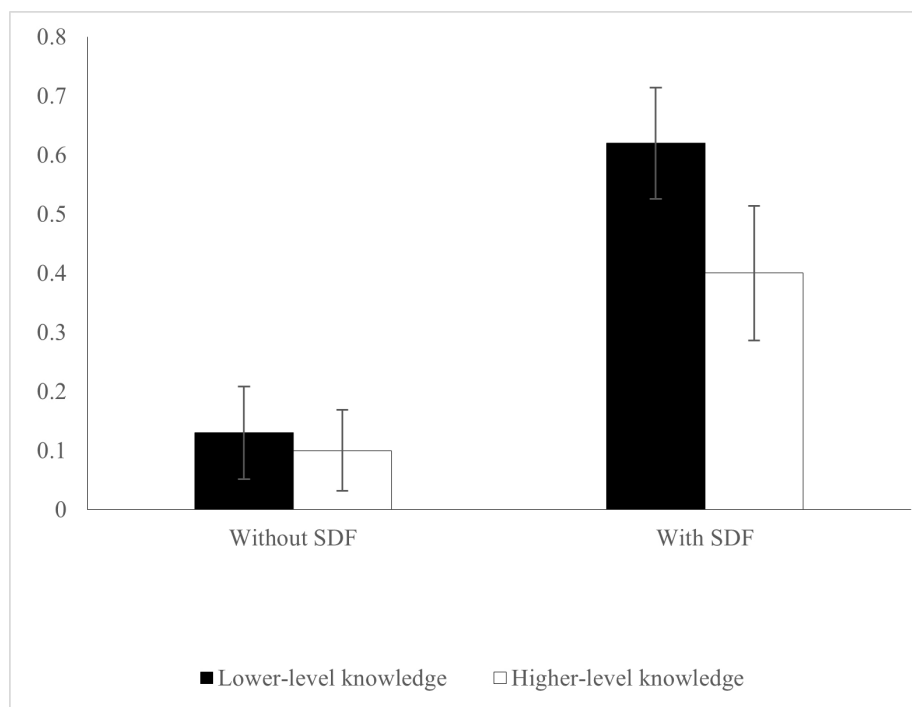


Figure 2 Incorrect diagnoses associated with salient distracting features (SDF) on cases with SDF and without SDF (mean frequency; range 0–1).

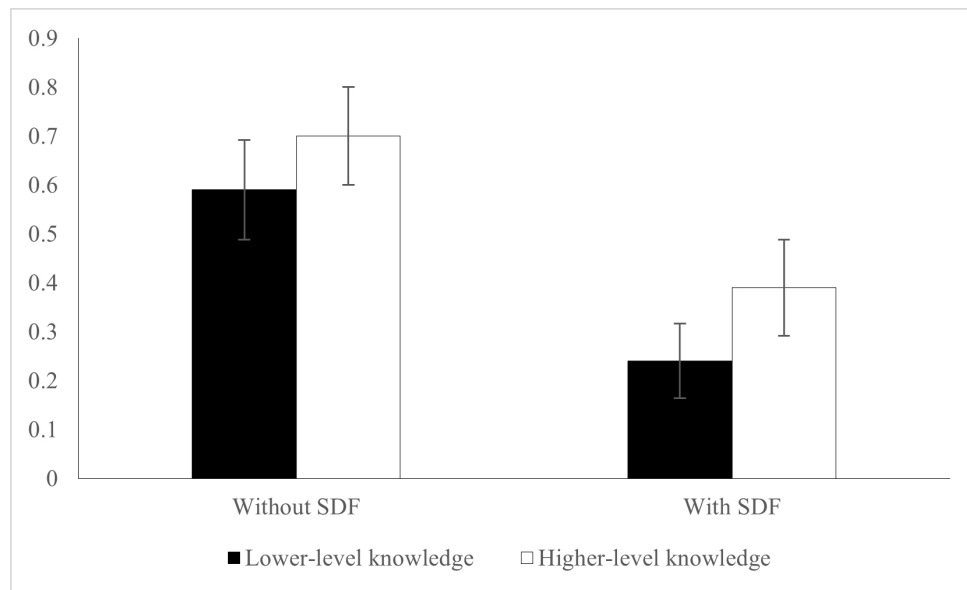


Figure 3 Diagnostic accuracy on cases with salient distracting features (SDF) and without SDF (mean diagnostic accuracy score; range 0–1).

Figure 3 presents the mean diagnostic accuracy scores when cases were subject to bias and not. Overall, higher knowledge physicians showed significantly higher accuracy relative to lower knowledge physicians ($p=0.014$), and the presence of SDF significantly decreased accuracy ($p<0.001$). Contrary to our expectation, the two groups did not significantly differ in how their diagnostic accuracy scores were hindered by the presence of SDF. Though the decrease in accuracy was lower in the higher knowledge than in the lower knowledge the interaction effect was not significant ($p=0.667$).

Figure 4 presents the indicators of reasoning mode—time spent in diagnosing and confidence in the diagnosis—in SDF+ and SDF– cases. Both knowledge groups invested more time to diagnose cases when SDFs were present than when they were not ($p<0.001$), but there was no overall significant difference between the two groups ($p=0.334$) and no interaction effect ($p=0.138$). Regarding rates of confidence, overall, the two knowledge groups did not significantly differ in

their confidence ratings ($p=0.087$), and the presence of SDF decreased confidence ($p=0.002$). However, a significant interaction effect was observed ($p=0.008$), because while the confidence of lower knowledge physicians remained basically the same when cases were subject to bias and not ($p=0.764$), higher knowledge physicians' confidence decreased in the presence of SDF ($p<0.001$).

DISCUSSION

In this experimental study, SDF encountered early in clinical cases induced anchoring bias, contributing to diagnostic errors. Among physicians from the same training programme and with similar years of clinical experience, susceptibility to bias was predicted by knowledge of critical diagnostic features that differentiated between alternative diagnoses for the clinical presentations displayed in the cases. When the cases contained SDF, thereby being subject to anchoring bias, physicians with less knowledge of discriminating features gave the bias-induced diagnosis considerably

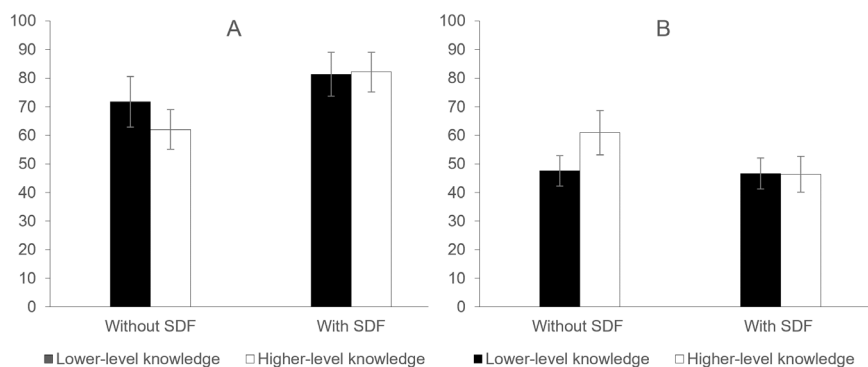


Figure 4 Indicators of reasoning mode: (A) time spent in diagnosis (mean time in seconds) and (B) confidence in the diagnosis (mean confidence; range 0–100) on cases with salient distracting features (SDF) and without SDF.

more frequently than their counterparts who had more comprehensive knowledge of these specific features. Overall, physicians needed more time and had lower confidence in cases with SDF compared with the same cases without SDF. However, the higher knowledge and lower knowledge groups did not differ in time or confidence when diagnosing the SDF+ cases. In psychological research on bias, more time to respond^{23–25} and lower confidence^{26–34} have been associated with more engagement in analytical reasoning. Therefore, our findings suggest that physicians from both groups moved towards more analytical reasoning when the cases were subject to bias. Because the two groups did not differ in these measures, reasoning mode per se cannot explain the difference in susceptibility bias. Instead, susceptibility was predicted by differences in knowledge of discriminating features.

Anchoring bias has been often pointed in the medical literature as a common threat to diagnostic reasoning.^{10–35} It appears frequently in retrospective investigations of actual diagnostic errors.^{7–8–36} Experienced clinicians often attribute their own diagnostic errors to sticking to an early hypothesis triggered by salient information even after disproving information becomes available.¹⁶ These studies suggest that repairing an incorrect initial hypothesis, though possible, is difficult. A hypothesis triggered by salient information comes to mind quickly and seems reasonable. After being generated, the hypothesis influences subsequent reasoning, potentially hindering gathering and/or correctly interpreting critical diagnostic information.^{37–38} In the present study, few modest SDFs were added to the cases, but even this little information was enough to distract the physicians.

An appealing (and widely discussed) idea is that clinicians should learn to detect when they are at risk of being biased, when they should stop and think further.^{39–40} Based on this idea, courses on clinical reasoning and cognitive debiasing strategies have been tried.¹⁸ Results so far do not seem promising, at least when error reduction was evaluated.^{20–41–42} An obstacle may be that detecting an error requires a sort of ‘diagnosis’ of one’s own diagnostic reasoning, which depends on cognitive processes that are as fallible as the ones involved in the clinical diagnosis. Another difficulty is that engaging in analytical reasoning per se does not guarantee that the bias-induced response is over-ridden.⁴³ There may always be physicians who fail to acknowledge that a problem requires further deliberation, proceeding instead with their intuitive (and potentially biased) response. However, in our study, the overall increase in diagnosis time on the SDF+ cases suggests that most physicians somehow recognised that the problem required further consideration. Indeed, psychological research shows that problems containing salient potentially bias-triggering information tend to lead people to think further before responding.^{43–45} Even those who end up giving

the biased response seem to engage in more analytical reasoning. Nevertheless, analytical reasoning can focus either on ‘justifying’ the initial response by searching for evidence that ends up strengthening it or on ‘overriding’ the initial response by reconsidering other findings and generating alternatives. The latter is more likely to arrive at the correct response.^{43–45} Physicians with stronger knowledge of discriminating features may be better equipped for the ‘overriding-format’ of analytical reasoning because they are more likely to recognise these features as relevant, which may raise alternative hypotheses and restructure initial reasoning. Less knowledgeable physicians who are not aware, for example, that hypersegmentation of the polymorphonuclear granulocytes would strongly support vitamin B₁₂ deficiency in the patient described in [table 1](#) could think about the case for a long time but still overlook these key features. Though apparently reasonable, these assumptions are conjectures demanding investigation.

Despite the seemingly obvious importance of knowledge to improve diagnosis, whether knowledge helps counteract bias is far from clear. In psychological research, studies comparing experts and novices’ susceptibility to bias have led to conflicting results.^{46–48} Furthermore, earlier experiments with physicians showed that difficulty in restructuring initial clinical reasoning increased with experience, making more experienced (most often older) physicians more vulnerable to anchoring bias.^{49–51} Noteworthy, these studies considered time in professional practice but did not assess physicians’ knowledge. However, even among physicians with similar training and years in clinical practice, differences in disease knowledge are unavoidable. If these differences are examined, as our study did, we may have a better picture of the role of knowledge in bias.

Educational strategies to reduce clinicians’ susceptibility to bias have hitherto focused on improving physicians’ reasoning *processes*, teaching about reasoning and biases. The role of knowledge in counteracting bias has been largely neglected.¹⁹ However, in a previous study we showed an intervention which enhanced relevant disease knowledge to reduce physicians’ vulnerability to bias in future cases.²¹ Consistently, a recent study associated deficiencies in diagnostic knowledge with adverse outcomes of primary care visits.⁵² Our findings reinforce the need to redirect our efforts to develop effective educational strategies for strengthening knowledge of discriminating features.

Intriguingly, the different susceptibility to bias in the two knowledge groups did not result in significantly different diagnostic accuracy. With regard to the frequency of the diagnosis associated with the SDF, the two knowledge groups did not significantly differ in how often (scarcely, as expected) they mentioned this diagnosis on SDF– cases. However, when the presence of SDF made the cases subject to bias, the

lower knowledge physicians incorrectly gave the bias-induced diagnosis considerably more frequently than the higher knowledge physicians. The results for diagnostic accuracy seem to show a similar pattern. The difference in diagnostic accuracy scores between the two knowledge groups was lower in SDF– cases than when the cases were subject to bias. Nevertheless, the interaction effect was not significant. The large variation within groups may have contributed to that. Furthermore, we only measured whether the physicians *possessed* the critical diagnostic knowledge but not if this knowledge was actually *used* during the diagnosis, which increases the noise. There may also be other reasons. The presence of SDF made the cases more ambiguous and, whereas the higher knowledge physicians could more easily exclude the bias-induced diagnosis than the lower knowledge physicians, higher knowledge physicians made other mistakes. Indeed, while the bias-induced diagnosis accounted for around 81% of the incorrect diagnoses given by the lower knowledge group, this proportion dropped to 65% in the higher knowledge group. Future research should examine this phenomenon.

The findings referring to confidence in the diagnosis are relevant to the current debate on diagnostic calibration. The two groups did not differ in confidence in SDF+ cases. Nonetheless, while the higher knowledge group reported lower confidence in SDF+ than in SDF– cases, the confidence reported by the lower knowledge group did not decrease when the cases were subject to bias. Knowledge seems to play a key role in susceptibility to bias and in diagnostic calibration.

This study has several limitations, some of them informing future research. We measured time and confidence, which though common in psychological research on bias^{24 26} are *indirect* indicators of reasoning. Research on the knowledge–bias interaction has just started, and these easy-to-obtain measurements sufficed as a first step. Future research should examine what takes place when participants think further about the case to determine features of analytical reasoning that help counteract bias. On a related issue (and limitation), we measured knowledge of discriminating features that participants *possessed* but not the extent to which they actually *used* this knowledge while diagnosing the cases. This requires research with other methodological tools. Another limitation is that knowledge level was classified based on performance aggregating all phase 1 diseases, which is not ideal considering content specificity. Nevertheless, case-level knowledge was highly correlated across all diseases (online supplemental file 2), suggesting that aggregated performance sufficiently captured knowledge differences. The use of written clinical cases leaves out important steps of the authentic clinical process. Noteworthy, an effect found with cases that provide physicians with all relevant information would possibly be larger rather than smaller when physicians

have to gather this information themselves. The influence of an incorrect initial hypothesis, for example, hinders gathering of critical information,³⁸ which would increase the risk of anchoring bias. Although accessing other resources is possible in real settings, physicians often fail to recognise this need.⁵³ Moreover, written cases allow for control and have proved an acceptable proxy for group differences in performance in practice.⁵⁴ The limited clinical experience of our participants makes it unclear whether the findings apply to more experienced physicians. Previous studies have shown that more experienced physicians tend to be more rather than less susceptible to bias; however, this susceptibility might also depend on specific clinical knowledge in this group.

In conclusion, in this study incorrect diagnoses induced by anchoring bias were more frequent among physicians with less knowledge of features that discriminate between lookalike diseases relative to more knowledgeable physicians. Regardless of their knowledge level, the physicians spent more time and reported lower confidence in the diagnosis of cases that contained SDF than in those without SDF. This suggests that the physicians realised, even if intuitively, that the case required further thought when it was subject to bias and moved towards more analytical reasoning. The lower susceptibility to bias cannot therefore be explained by differences in the degree of engagement in analytical reasoning. Differences in knowledge of discriminating features seem to have predicted it. These findings may help place refinement of diagnostic knowledge in the core of strategies to make clinicians less vulnerable to bias.

Author affiliations

¹Institute of Medical Education Research Rotterdam, Erasmus Medical Center, Rotterdam, The Netherlands

²Department of Internal Medicine, Erasmus Medical Center, Rotterdam, The Netherlands

³Wenckebach Institute (WIOO), University Medical Centre Groningen, Groningen, The Netherlands

⁴Department of Hematology, University Medical Centre Groningen, Groningen, The Netherlands

⁵Department of Intensive Care, Spaarne Gasthuis, Haarlem, The Netherlands

⁶Department of Intensive Care, Erasmus MC, Rotterdam, The Netherlands

⁷Department of Psychology, Methodology and Statistics, Leiden University, Leiden, The Netherlands

⁸Department of Psychology, Education and Child Studies, Erasmus Universiteit Rotterdam, Rotterdam, The Netherlands

Acknowledgements The authors are grateful to the residents who dedicated their scarce time to participate in the study and to Liona Ionescu for her help to format the manuscript.

Contributors All authors had full access to all the study data and take responsibility for the integrity of the data and the accuracy of the data analysis. Study conception and design: SM, HGS. Development of study materials: SM, AZ, MAdC-F. Acquisition of data: SM, AZ, MAdC-F, MG, GC, LZ. Statistical analysis: SM, HGS, JvG. Analysis or interpretation of data: all authors. Drafting of the manuscript: SM. Critical revision of the manuscript for important intellectual content: all authors. Administrative, technical or material support: SM, AZ, GC. Supervision: SM, HGS. Guarantor: SM.

Funding The authors have not declared a specific grant for this research from any funding agency in the public, commercial or not-for-profit sectors.

Competing interests None declared.

Patient consent for publication Not applicable.

Ethics approval This study involves human participants and was approved by the Ethics Review Committee of the Department of Psychology, Education and Child Studies (reference number: 20-026 ErC DPECS), Erasmus University Rotterdam. Participants gave informed consent to participate in the study before taking part.

Provenance and peer review Not commissioned; externally peer reviewed.

Data availability statement Data are available upon reasonable request. Data are available upon reasonable request and subject to institutional regulations.

Supplemental material This content has been supplied by the author(s). It has not been vetted by BMJ Publishing Group Limited (BMJ) and may not have been peer-reviewed. Any opinions or recommendations discussed are solely those of the author(s) and are not endorsed by BMJ. BMJ disclaims all liability and responsibility arising from any reliance placed on the content. Where the content includes any translated material, BMJ does not warrant the accuracy and reliability of the translations (including but not limited to local regulations, clinical guidelines, terminology, drug names and drug dosages), and is not responsible for any error and/or omissions arising from translation and adaptation or otherwise.

ORCID iDs

Sílvia Mamede <http://orcid.org/0000-0003-1187-2392>
 Adrienne Zandbergen <http://orcid.org/0000-0001-5056-5921>
 Marco Antonio de Carvalho-Filho <http://orcid.org/0000-0001-7008-4092>
 Joost van Ginkel <http://orcid.org/0000-0002-4137-0943>
 Laura Zwaan <http://orcid.org/0000-0003-3940-1699>
 Fred Paas <http://orcid.org/0000-0002-1647-5305>
 Henk G Schmidt <http://orcid.org/0000-0001-8706-0978>

REFERENCES

- Berner ES, Graber ML. Overconfidence as a cause of diagnostic error in medicine. *Am J Med* 2008;121(5 Suppl):S2–23.
- Singh H, Meyer AND, Thomas EJ. The frequency of diagnostic errors in outpatient care: estimations from three large observational studies involving US adult populations. *BMJ Qual Saf* 2014;23:727–31.
- Balogh EP, Miller BT, Ball JR, *et al.* Improving diagnosis in health care. In: *Committee on Diagnostic Error in Health Care. Improving diagnosis in health care*. Washington, D.C: The National Academies Press, 29 December 2015.
- Wallace E, Lowry J, Smith SM, *et al.* The epidemiology of malpractice claims in primary care: a systematic review. *BMJ Open* 2013;3:e002929.
- Dave N, Bui S, Morgan C, *et al.* Interventions targeted at reducing diagnostic error: systematic review. *BMJ Qual Saf* 2022;31:297–307.
- Prakash S, Sladek RM, Schuwirth L. Interventions to improve diagnostic decision making: a systematic review and meta-analysis on reflective strategies. *Medical Teacher* 2019;41:517–24.
- Graber ML, Franklin N, Gordon R. Diagnostic error in internal medicine. *Arch Intern Med* 2005;165:1493–9.
- Kachalia A, Gandhi TK, Puopolo AL, *et al.* Missed and delayed diagnoses in the emergency department: a study of closed malpractice claims from 4 liability insurers. *Ann Emerg Med* 2007;49:196–205.
- Croskerry P. From mindless to mindful practice — cognitive bias and clinical decision making. *N Engl J Med* 2013;368:2445–8.
- Redelmeier DA. Improving patient care. The cognitive psychology of missed diagnoses. *Ann Intern Med* 2005;142:115–20.
- Evans JS. The heuristic-analytic theory of reasoning: extension and evaluation. *Psychon Bull Rev* 2006;13:378–95.
- Mamede S, van Gog T, van den Berge K, *et al.* Effect of availability bias and reflective reasoning on diagnostic accuracy among internal medicine residents. *JAMA* 2010;304:1198–203.
- Schmidt HG, Mamede S, van den Berge K, *et al.* Exposure to media information about a disease can cause doctors to misdiagnose similar-looking clinical cases. *Acad Med* 2014;89:285–91.
- Mamede S, Splinter TAW, van Gog T, *et al.* Exploring the role of salient distracting clinical features in the emergence of diagnostic errors and the mechanisms through which reflection counteracts mistakes. *BMJ Qual Saf* 2012;21:295–300.
- Mamede S, van Gog T, van den Berge K, *et al.* Why do doctors make mistakes? A study of the role of salient distracting clinical features. *Academic Medicine* 2014;89:114–20.
- Balla J, Heneghan C, Goyder C, *et al.* Identifying early warning signs for diagnostic errors in primary care: a qualitative study. *BMJ Open* 2012;2:e001539.
- Restrepo D, Armstrong KA, Metlay JP. *Annals* clinical decision making: avoiding cognitive errors in clinical decision making. *Ann Intern Med* 2020;172:747–51.
- Reilly JB, Ogdie AR, Von Feldt JM, *et al.* Teaching about how doctors think: a longitudinal curriculum in cognitive bias and diagnostic error for residents. *BMJ Qual Saf* 2013;22:1044–50.
- Dhaliwal G. Premature closure? Not so fast. *BMJ Qual Saf* 2017;26:87–9.
- Norman GR, Monteiro SD, Sherbino J, *et al.* The causes of errors in clinical reasoning: cognitive biases, knowledge deficits, and dual process thinking. *Academic Medicine* 2017;92:23–30.
- Mamede S, de Carvalho-Filho MA, de Faria RMD, *et al.* 'Immunising' physicians against availability bias in diagnostic reasoning: a randomised controlled experiment. *BMJ Qual Saf* 2020;29:550–9.
- Mamede S, Goeijenbier M, Schuit SCE, *et al.* Specific disease knowledge as predictor of susceptibility to availability bias in diagnostic reasoning: a randomized controlled experiment. *J Gen Intern Med* 2021;36:640–6.
- Bonner C, Newell BR. In conflict with ourselves? An investigation of heuristic and analytic processes in decision making. *Memory & Cognition* 2010;38:186–96.
- De Neys W, Glumicic T. Conflict monitoring in dual process theories of thinking. *Cognition* 2008;106:1248–99.
- Rotgans JI, Schmidt HG, Rosby LV, *et al.* Evidence supporting dual-process theory of medical diagnosis: a functional near-infrared spectroscopy study. *Med Educ* 2019;53:143–52.
- Gangemi A, Bourgeois-Gironde S, Mancini F. Feelings of error in reasoning—in search of a phenomenon. *Thinking & Reasoning* 2015;21:383–96.
- Vartanian O, Beatty EL, Smith I, *et al.* The reflective mind: examining individual differences in susceptibility to base rate neglect with fMRI. *J Cogn Neurosci* 2018;30:1011–22.
- Westerman DL, Payne DG. Research methods in human memory. In: Davis SF, ed. *Handbook of research methods in experimental psychology*. Wiley-Blackwell, 2005: 346–65.

- 29 Norcini J, Guille R. Combining tests and setting standards. In: Van der Norman GR, Vleuten CPM, Newble DI, eds. *International handbook of research in medical education*. Dordrecht: Kluwer Academic Publishers, 2002.
- 30 Rubin DB. Multiple imputation for nonresponse in surveys. New York Wiley; 1987.
- 31 Van Buuren S, Groothuis-Oudshoorn CGM. mice: Multivariate imputation by chained equations in R. *J Stat Softw* 2011;45:1–67.
- 32 Rubin DB. Statistical matching using file concatenation with adjusted weights and multiple imputations. *J Bus Econ Stat* 1986;4:87–94.
- 33 van Buuren S. *Flexible imputation of missing data, 2nd ed.* Boca Raton, FL: Chapman & Hall/CRC Press, 2018.
- 34 Vartanian O, Lam TK, Maceda E, et al. Can a fast thinker be a good thinker? The neural correlates of base-rate neglect measured using a two-response paradigm. *Cogn Neuropsychol* 2021;38:365–86.
- 35 Croskerry P. The importance of cognitive errors in diagnosis and strategies to minimize them. *Acad Med* 2003;78:775–80.
- 36 Saposnik G, Redelmeier D, Ruff CC, et al. Cognitive biases associated with medical decisions: a systematic review. *BMC Med Inform Decis Mak* 2016;16:138.
- 37 Barrows HS, Norman GR, Neufeld VR, et al. The clinical reasoning of randomly selected physicians in general medical practice. *Clin Invest Med* 1982;5:49–55.
- 38 Kostopoulou O, Russo JE, Keenan G, et al. Information distortion in physicians' diagnostic judgments. *Med Decis Making* 2012;32:831–9.
- 39 Croskerry P. Cognitive forcing strategies in clinical decisionmaking. *Ann Emerg Med* 2003;41:110–20.
- 40 Croskerry P. When I say... cognitive debiasing. *Med Educ* 2015;49:656–7.
- 41 Sherbino J, Kulasegaram K, Howey E, et al. Ineffectiveness of cognitive forcing strategies to reduce biases in diagnostic reasoning: a controlled trial. *CJEM* 2014;16:34–40.
- 42 Staal J, Hooftman J, Gunput STG, et al. Effect on diagnostic accuracy of cognitive reasoning tools for the workplace setting: systematic review and meta-analysis. *BMJ Qual Saf* 2022;31:899–910.
- 43 Pennycook G, Fugelsang JA, Koehler DJ. What makes us think? A three-stage dual-process model of analytic engagement. *Cogn Psychol* 2015;80:34–72.
- 44 De Neys W. Bias and conflict: a case for logical intuitions. *Perspect Psychol Sci* 2012;7:28–38.
- 45 De Neys W, Bonnefon J-F. The “whys” and “whens” of individual differences in thinking biases. *Trends Cogn Sci* 2013;17:172–8.
- 46 Brewer NT, Chapman GB, Schwartz JA, et al. The influence of irrelevant anchors on the judgments and choices of doctors and patients. *Med Decis Making* 2007;27:203–11.
- 47 Mussweiler T, Englich B. Adapting to the Euro: evidence from bias reduction. *J Econ Psychol* 2003;24:285–92.
- 48 Schmittat SM, Englich B. If you judge, investigate! Responsibility reduces confirmatory information processing in legal experts. *Psychology, Public Policy, and Law* 2016;22:386–400.
- 49 Eva KW, Cunnington JPW. The difficulty with experience: does practice increase susceptibility to premature closure? *J Contin Educ Health Prof* 2006;26:192–8.
- 50 Eva KW, Link CL, Lutfey KE, et al. Swapping horses midstream: factors related to physicians' changing their minds about a diagnosis. *Academic Medicine* 2010;85:1112–7.
- 51 St-Onge C, Landry M, Xhignesse M, et al. Age-related decline and diagnostic performance of more and less prevalent clinical cases. *Adv in Health Sci Educ* 2016;21:561–70.
- 52 Vandergrift JL, Weng WF, Gray BM. The association between physician knowledge and inappropriate medications for older populations. *J Am Geriatr Soc* 2021;69:3584–94.
- 53 Meyer AND, Payne VL, Meeke DW, et al. Physicians' diagnostic accuracy, confidence, and resource requests. *JAMA Intern Med* 2013;173:1952.
- 54 Peabody JW, Luck J, Glassman P, et al. Comparison of vignettes, standardized patients, and chart abstraction: a prospective validation study of 3 methods for measuring quality. *JAMA* 2000;283:1715–22.

Study protocol

Title: The role of specific disease knowledge in counteracting bias in diagnostic reasoning

Scientific title: The influence of specific disease knowledge on physicians' susceptibility to bias in diagnostic reasoning: A randomized controlled trial

Authors: Sílvia Mamede (ErasmusMC, ESSB), Marco Goeijenbier (ErasmusMC), Adrienne Zandbergen (ErasmusMC), Marco de Carvalho-Filho (University of Minho and UMC Groningen), Laura Zwaan (ErasmusMC), Fred Paas (ESSB), Henk G. Schmidt (ErasmusMC, ESSB).

Background

Diagnostic errors have been attributed to flaws in physicians' reasoning associated with the use of heuristics,¹⁻⁴ experientially derived rules of thumb that serve as shortcuts in reasoning.⁵ For example, recent experiences with a particular disease led physicians to confuse subsequent cases that looked like it (but was in fact different) with the disease seen before.^{6,7} However, while most physicians indeed fell prey to bias in these experimental studies, a substantial fraction of them did not.^{6,7} Although these physicians appear to be at the same level of expertise (as measured by training level and years of clinical experience), some were therefore to be more susceptible to bias than others. The sources of these differences in susceptibility to bias are unclear. The present study aims to further examine the role of specific disease knowledge in counteracting bias in diagnostic reasoning. It will examine bias induced by "salient distracting features", which are findings in a case that tend to catch physicians' attention because though irrelevant to the case they are strongly associated with a disease that seems at first glance a plausible diagnosis.^{8,9} Physicians will diagnose the same set of clinical vignettes that have been manipulated to be either subject-to-bias (with salient distracting features) or not-subject-to-bias (without salient distracting features). Physicians knowledge about the diseases presented in the vignettes will be evaluated. By building upon the previous study and on psychological research, we have the following prior hypotheses:

H1: Overall, diagnostic accuracy will be lower on subjected-to-bias (StB) than on not-subjected-to-bias cases (Not-StB), but the difference in accuracy will be smaller (or absent) in the higher-knowledge group.

H2: Overall, the frequency of the wrong diagnoses induced by the salient distracting features will be higher on StB than Not-StB cases but lower in the higher-knowledge than in the lower-knowledge group.

H3: Overall, physicians will spend more time to diagnose StB than Not-StB cases without differences between the knowledge groups.

H4: Overall, the rating of confidence in the diagnosis will be lower on StB than Not-StB cases, but the difference will be lower in the higher-knowledge group.

H5: Overall, the diagnosis triggered by the salient distracting features will be mentioned more frequently when physician recall which diagnoses they had considered on StB than Not-StB cases, without differences between the two knowledge groups.

Methods

The study is an experiment consisting of three phases: (1) assessment of physicians' knowledge about the diseases used in the study; (2) diagnostic task; (3) assessment of case processing. Phase 1 will be carried out 2-4 hours before the subsequent phases, which will take place sequentially in a single session. Phase 1 and the subsequent phases will be presented to participants as different, independent studies. Phase 1 consists of a "recall task" which assesses the associations between the diseases and symptoms as stored in memory by measuring the amount and accuracy of findings that belong to a particular disease that participants are able to recollect from memory under restricted time conditions.¹⁰ In Phase 2, the physicians will diagnose the same set of cases, all with a confirmed diagnosis. Each participant diagnoses half of the cases with the salient distracting feature and half of the cases without the salient distracting features, in a balanced within-subjects incomplete block design. Finally, in Phase 3, physicians will indicate which alternative diagnoses they considered for each case solved in Phase 2. Despite being a post hoc recall, this request was sensitive enough to capture differences in reasoning in a previous study by our group on the effect of time restrictions on diagnostic accuracy.¹¹ Participants will also rate the confidence in the diagnosis given. Time required to give answers will be registered in all phases. Figure 1 presents the study design.

Participants

Physicians in training to become specialists in internal medicine or emergency medicine at the ErasmusMC and other teaching hospitals in the Netherlands with at least 1 year of clinical experience will be invited to voluntarily participate in the study. Those who volunteer will be recruited as participants. Each participant will receive a compensation of € 50,00. Written consent will be obtained from all participants.

Ethical approval

The study protocol has been approved by the Ethics Review Committee of the Department of Psychology, Education and Child Studies (DPECS), Erasmus University (Ethics Review Application 20-026).

Sample size estimation

A priori power analysis was performed by assuming a medium effect size (Cohen's $f = 0.30$), which has been detected in a previous study,⁹ as the to-be-detected effect size, the standard alpha level of 0.05, indicated that a sample size of 62 participants would be sufficient to achieve a power of 0.80.

Materials and procedure

The knowledge evaluation task will consist of a recall task of findings associated with each of the 6 diseases to be seen in the diagnostic task (Phase 2) plus 3 filler diseases. On each screen presenting the disease, the participant will be requested to type everything they can remember about clinical findings that are important for the diagnosis of the disease. The order of presentation of the diseases for the recall task will be randomized.

For the subsequent phases of the study, 9 written clinical cases will be used, (6 relevant cases and 3 fillers) (see Appendix 1). Each case will be prepared in two versions which differ exclusively on the presence of one-two salient distracting features (with or without). In both versions, the most likely diagnosis of the case will be the same. For each case, the participant will be asked to read the case and type the most likely diagnosis as fast as possible but without compromising accuracy. This stimulus to be fast aims at avoiding that participants engage in extensive reflection upon the cases, as reflective reasoning has shown to counteract the effect of bias.^{7,12} After diagnosing all cases, in Phase 3, each case is presented again, one by one, with its initial sentence that reminds the participants of the case and the diagnosis given by the participant in Phase 2. The participant is requested to recall the alternative diagnoses considered for the case and rate (in percentage) the confidence in the diagnosis given to the case.

Finally, the participants will be requested to provide demographic information, rate their experience with the diseases presented in the study by using a scale adopted in previous studies,¹³ and answer probing questions about their understanding of the study task, interruptions during the task and suspicion about the study purposes and manipulation. Feedback about the correct diagnoses of each case will then be offered by presenting each case again with its correct diagnosis to participants who opt to proceed to it.

The whole experiment will be conducted on an electronic environment by using the Qualtrics research suite, which automatically randomly allocates each participant to different conditions and registers the participant's responses and response time.

Outcome measurements and data analysis

The following measurements will be collected in each phase of the study:

Phase 1: Performance in the recall task.

Phase 2: Diagnostic accuracy; frequency of bias-induced diagnosis; time to diagnose.

Phase 3: Frequency of the diagnoses of the biasing-inducing cases (diagnosis triggered by the salient distracting features); ratings of confidence in the diagnosis.

Performance in the recall task (Phase 1) will be measured by the frequency of accurate clinical findings (counted as idea units) present in the recall protocols. To develop a standard for the task, three internists (A.Z.; M.G.; M.C.) will previously identify the critical diagnostic findings for each disease. A two-step procedure will be used for the scoring of the protocols. First, two researchers will independently count the units in a random sample of 10% of the protocols, and the interrater

agreement will be computed. If the interrater agreement is above 0.80, the count will proceed with a single evaluator. The frequencies will be analyzed by using descriptive statistics, and participants will be assigned to different knowledge-level groups based on their performance.

Separate mixed ANOVAs with knowledge group as between-subjects factor (lower knowledge vs higher knowledge) and exposure to bias (StB and NotStB) as within-subjects factors will be performed with the following dependent variables:

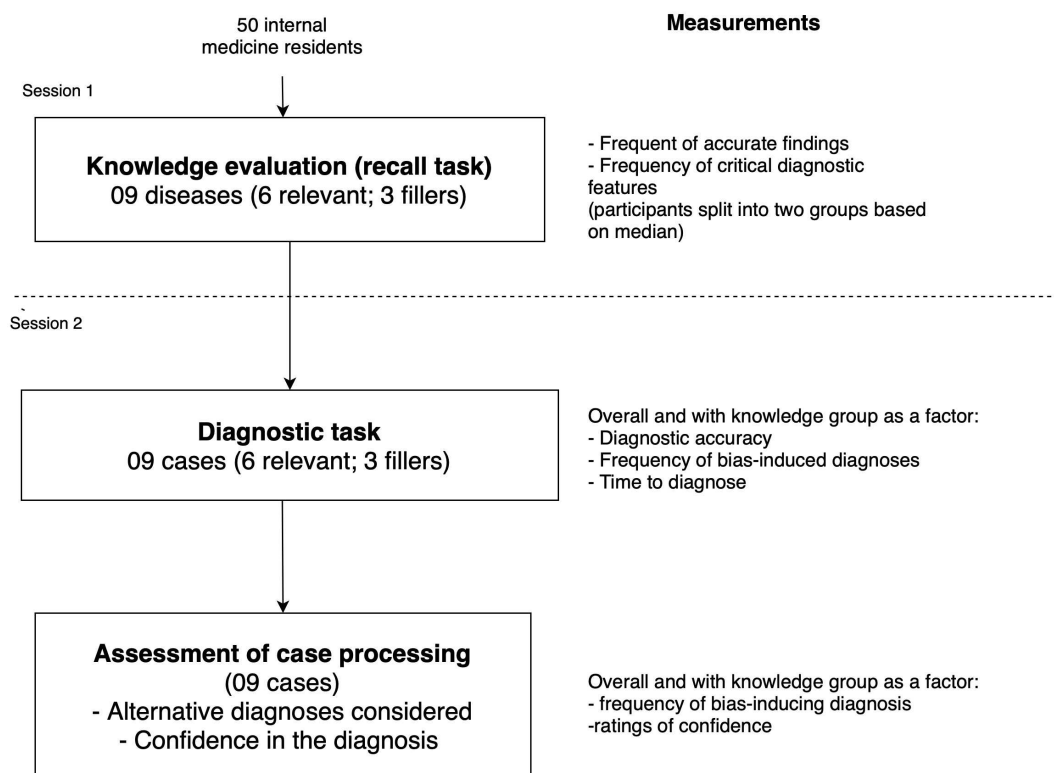
- mean diagnostic accuracy scores (Phase 2), aimed at testing H1.
- mean frequency of bias-induced diagnosis (diagnosis triggered by the salient distracting features) (Phase 2), aimed at testing H2.
- mean time spent in diagnosing the cases (Phase 2), aimed at testing H3.
- mean ratings of confidence in diagnosis (Phase 3), aimed at testing H4.
- mean frequency of the diagnoses of the bias-inducing phase mentioned in Phase 3, aimed at testing H5.

This analysis is based on the aggregated data for all cases and on the grouping of participants into two knowledge levels. However, depending on the variation in knowledge level across the diseases, additional knowledge groups will be formed and additional analysis at case-level will be performed. Eventual outliers (absolute values surpassing the threshold of 2 z-scores) on the main outcome measurements - time spent in the diagnosis (Phase 2), which indicates compliance with the treatment, and the main outcome variables related to diagnostic performance – will be excluded from the analysis.

References

1. Croskerry P. The importance of cognitive errors in diagnosis and strategies to minimize them. *Acad Med* 2003;78(8):775-80.
2. Croskerry P. From mindless to mindful practice--cognitive bias and clinical decision making. *N Engl J Med* 2013;368(26):2445-8.
3. Klein JG. Five pitfalls in decisions about diagnosis and prescribing. *BMJ* 2005;330(7494):781-3.
4. Redelmeier DA. Improving patient care. The cognitive psychology of missed diagnoses. *Ann Intern Med* 2005;142(2):115-20.
5. Tversky A, Kahneman D. Judgment under Uncertainty: Heuristics and Biases. *Science* 1974;185(4157):1124-31.
6. Mamede S, van Gog T, van den Berge K, et al. Effect of availability bias and reflective reasoning on diagnostic accuracy among internal medicine residents. *JAMA* 2010;304(11):1198-203.
7. Schmidt HG, Mamede S, van den Berge K, et al. Exposure to media information about a disease can cause doctors to misdiagnose similar-looking clinical cases. *Acad Med* 2014;89(2):285-91.
8. Mamede S, Splinter TA, van Gog T, et al. Exploring the role of salient distracting clinical features in the emergence of diagnostic errors and the mechanisms through which reflection counteracts mistakes. *BMJ Qual Saf* 2012;21(4):295-300.
9. Mamede S, van Gog T, van den Berge K, et al. Why do doctors make mistakes? A study of the role of salient distracting clinical features. *Acad Med* 2014;89(1):114-20.
10. Westerman DL, Payne DG. Research Methods in Human Memory. In: Davis SF, editor. *Handbook of Research Methods in Experimental Psychology*: Wiley-Blackwell; 2005. p. 346-65.
11. DA AL, Rotgans JI, Mamede S, et al. Factors underlying suboptimal diagnostic performance in physicians under time pressure. *Med Educ* 2018;52(12):1288-98.
12. Mamede S, van Gog T, van den Berge K, et al. Effect of Availability Bias and Reflective Reasoning on Diagnostic Accuracy Among Internal Medicine Residents. *Jama-Journal of the American Medical Association* 2010;304(11):1198-203.
13. Mamede S, Carvalho-Filho MA, de Faria RM, et al. "Immunising" physicians against availability bias in diagnostic reasoning: A randomised controlled experiment. *BMJ Q&S* in press.

Figure 1 – Diagram of the study design and outcome measurements



Appendix 1- Diseases and cases to be included in each phase of the study

Phase 1 Knowledge evaluation (Diseases)	Phase 3 Diagnostic task (Clinical cases)	Phase 3 Processing assessment (Clinical cases)
All participants	Half of the participants randomized to one version of the case half to the other*	All participants
1. Inflammatory bowel disease 2. Hyperthyroidism 3. Vitamin B12 deficiency 4. Addison's disease 5. Acute appendicitis 6. Acute bacterial endocarditis 3 fillers	1. Inflammatory bowel disease 2. Hyperthyroidism 3. Vitamin B12 deficiency 4. Addison's disease 5. Acute appendicitis 6. Acute bacterial endocarditis 3 fillers	1. Inflammatory bowel disease 2. Hyperthyroidism 3. Vitamin B12 deficiency 4. Addison's disease 5. Acute appendicitis 6. Acute bacterial endocarditis 3 fillers

(*) half of the cases with salient distracting features, half of the cases without salient distracting features

Supplementary file 2

Additional information on methods and results

Methods

Material and Procedure

Table 1 – Diseases evaluated in the recall task in phase 1 (knowledge evaluation task) and diagnoses of the clinical cases in phase 2 (diagnostic task)

Phase 1 – Knowledge evaluation	Phase 2 – Diagnostic task
<i>Test diseases</i>	<i>Test cases</i>
Addison's disease	Addison's disease
Hyperthyroidism	Hyperthyroidism
Acute bacterial endocarditis	Acute bacterial endocarditis
Inflammatory bowel disease	Inflammatory bowel disease
Vitamin B12 deficiency	Vitamin B12 deficiency
Appendicitis	Appendicitis
<i>Filler diseases</i>	<i>Filler cases</i>
Nephrotic syndrome	Alcoholic cirrhosis
Bacterial meningitis	Acute alcoholic pancreatitis
Acute viral hepatitis	Rheumatoid arthritis

(* Fillers are not relevant to the study and not analysed. They were mixed to the test diseases in phase 1 (knowledge evaluation task) and to the test cases in phase 2 (diagnostic task) to reduce the chance that participants easily recognize, when they start to diagnose the cases in phase 2, that they were being tested on the diseases saw in phase 1, which would facilitate the diagnosis.

In phase 2, the cases were presented one by one, and the participant was asked to read the case and type the most likely diagnosis as fast as possible without compromising accuracy. This stimulus to be fast aims at avoiding that participants engage in extensive reflection upon the cases, as reflective reasoning has shown to counteract the effect of bias.^(1, 2) Participants were informed that time required to diagnose a case ranged from 30 seconds to 1.5 minute in previous studies. The programme reminded the participant when 1.5 minute had passed but did not automatically move to the next case. Additional data was collected after the diagnosis of all cases but was not analysed for the present study due to the large number of missing values.

Missing data

The computer program did not set responses as compulsory for participants to proceed in the task. Thirty-seven out of the 68 participants skipped questions. Missing responses were observed in all variables, with 8% of all data missing. In total, the data contained 80 variables, all of which were used in the imputation process, either to be imputed, to serve as predictors for the missing values on other variables, or both. With a total sample size of $N = 68$, it is not possible for the missing values on a variable to be predicted by all other variables, as there would be more variables in the prediction model than there would be cases. Instead, a selection of predictors was made for each variable with missing data. For each numerical variable (not the ones computed using other variables), the number of predictors was based on the principle of at least 15 cases per predictor⁽³⁾ (thus, the more observed values, the more predictors). For categorical variables with j categories, the number of predictors was the number based on Steven's rule of thumb, divided by $j - 1$.

Once the number of predictors was determined for each variable with missing data, variables were selected to be predictors, where categorical predictors with k categories counted as $k - 1$ (dummy) predictors. The selection of these predictors was done as follows: first, all variables that were presumed to be related to the variable with missing data in the subsequent statistical analyses, were selected. If there were still predictors left to be selected after that, those variables were selected with either the highest squared correlation with the variable with missing data (in case both variables were numeric), the highest Cramér's V (in case both variables were categorical), or the highest η^2 (in case one variable was numeric and the other one was categorical). Predictors were added until the maximum number of predictors was reached.

It should be noted that an important assumption of multiple imputation is that the missing data are Missing at Random (e.g., Little & Rubin, 2002;(4) Rubin, 1976(5)). Unfortunately, this assumption cannot be checked on the available data. However, even when the MAR assumption does not hold, multiple imputation is still a better alternative to handle missing data than deleting incomplete cases from the analysis. Additionally, the MAR assumption becomes more plausible as more variables are included in the imputation procedure.(7) Note that we tried to include as many variables in the imputation procedure as possible, to make the MAR assumption most likely.

An additional clarification on the multiple imputation is that we opted for using the imputed values of the outcome variables in the analyses. Von Hippel (2008)(6) argued that imputing the outcome variable is justified and even recommendable because the imputed values on the outcome variable may help in predicting the missing values on the predictors. However, he went on arguing that it would be better not to use the imputed values on the outcome variables in the subsequent analysis because that may cause the standard errors to slightly increase. In other words, not imputing the outcome variable is not recommended but not

using the imputed values on the outcome variables in the analysis, is. The reason why we did use the imputed values of the outcome variable in the subsequent analyses, is that we wanted to keep all analyses comparable regarding sample size and imputed values used. Also, see van Ginkel et al. (2020, p. 305)(7) who used this as an argument for keeping the imputed values on the outcome variables in the analyses. Besides, using the imputed values of the outcome variables in the analysis will not lead to biased results, only in slightly less power, which we did not consider to be a serious problem.

Data analysis

Relative (norm-referenced) standard setting methods have been used in medical schools in the Netherlands as to define cut-off scores and rank students in achievement tests.(8) Similarly, we used a relative cut-off point to categorize participants as having either lower-knowledge or higher-knowledge of discriminating features based on the number of features reported in the recall task (phase 1). Employing median splits to create a categorical variable with “high” and “low” groups based on scores obtained in a continuous variable test is common practice for instance in research domains such as psychology, which typically uses analysis of variance to test influences of categorical predictors on an outcome variable. The categorization offers the advantages of simplifying data analysis and the presentation of results, with median splits offering acceptable results when the distribution of the original variable does not substantially deviate from normality.(9)

We opted for mixed ANOVAs for the statistical analysis, as previewed in the pre-registered protocol. However, it should be noted that multilevel modelling is an alternative way of dealing with missing data in repeated measures analysis, to multiple imputation followed by mixed ANOVAs. However, an advantage of multiple imputation over multilevel is that

firstly, it not only resolves the missing-data problem for the repeated measures analysis, but also for all other analyses applied to the same data. Secondly, unlike multilevel analysis, multiple imputation is capable of taking into account dependencies of the missing data on variables outside the analysis model, whereas multilevel can only pick up relations of the missing data with variables inside the multilevel model.

In a post-hoc analysis, we checked whether knowledge level as measured by the counting of findings recalled in the recall task in phase 1 correlated across all 6 diseases. We computed the correlations between disease-level measurements of knowledge. This helped clarify whether classifying knowledge in aggregate (i.e., based on the performance in the recall task in all 6 diseases taken together) actually provided an acceptable classification of knowledge level, considering that specific knowledge would be the key driver of reasoning success. All correlations were significant, with very high Pearson's correlation coefficients, ranging from $r = .64$ (between knowledge of inflammatory bowel disease and knowledge of acute bacterial endocarditis) and $r = .83$ (between knowledge of appendicitis and knowledge of acute bacterial endocarditis). These results indicate that the aggregated performance was a sufficient approximation of a participant's knowledge, thereby validating our analysis.

Results

Complete case analysis

Thirty-nine participants had no missing responses in the outcome variables. The two knowledge groups did not differ in number of years in clinical practice. Mean (standard deviation), lower-knowledge group, 2.71 (1.21), higher-knowledge group: 2.70 (1.39); $p = .98$.

The table below presents the mean (standard deviation into brackets) frequency of diagnoses associated with the SDF, diagnostic accuracy score, time spent in diagnosis and confidence in the diagnosis in cases with and without SDF (respectively, SDF+ and SDF-) for the two knowledge groups including only the participants with complete data.

Table 1. Participants' performance as a function of knowledge group and type of case

	Lower- knowledge (N = 19)	Higher- knowledge (N = 20)	Overall (N = 39)
Frequency of diagnoses associated with the SDF			
• SDF- cases	0.10 (0.19)	0.12 (0.19)	0.11 (0.19)
• SDF+ cases	0.61 (0.23)	0.40 (0.33)	0.50 (0.30)
Diagnostic accuracy score			
• SDF- cases	0.60 (0.28)	0.72 (0.24)	0.66 (0.26)
• SDF+ cases	0.27 (0.19)	0.42 (0.28)	0.35 (0.25)
Time spent in diagnosis			
• SDF- cases	66.51 (22.19)	62.24 (17.94)	64.32 (19.98)
• SDF+ cases	76.91 (27.25)	81.79 (21.75)	79.41 (24.38)
Confidence in the diagnosis			
• SDF- cases	50.12 (13.86)	62.49 (23.20)	56.31 (19.87)
• SDF+ cases	49.75 (14.64)	48.07 (20.58)	48.91 (17.64)

Below we present the results of the statistical tests for each outcome measurement. Effect size is measured by partial eta squared, with values of 0.01, 0.06, and 0.14 as benchmarks for small, medium, and large effect sizes, respectively.⁽¹⁰⁾

Frequency of diagnoses associated with the SDF

Main effect of SDF: $F(1, 37) = 74.61; p < .001; \eta^2 = 0.67$

Main effect of knowledge: $F(1, 37) = 2.49; p < .123; \eta^2 = 0.06$

Interaction effect: $F(1, 37) = 6.04; p < .019; \eta^2 = 0.14$

Diagnostic accuracy score

Main effect of SDF: $F(1, 37) = 33.85; p < .001; \eta^2 = 0.48$

Main effect of knowledge: $F(1, 37) = 4.98; p < .032; \eta^2 = 0.12$

Interaction effect: $F(1, 37) = 0.15; p = .70; \eta^2 = 0.004$

Time spent in diagnosis

Main effect of SDF: $F(1, 37) = 10.53; p = .002; \eta^2 = 0.22$

Main effect of knowledge: $F(1, 37) = 0.003; p < .956; \eta^2 = 0.00$

Interaction effect: $F(1, 37) = 0.98; p = .329; \eta^2 = 0.03$

Confidence in diagnosis

Main effect of SDF: $F(1, 37) = 4.51; p = .041; \eta^2 = 0.11$

Main effect of knowledge: $F(1, 37) = 1.19; p < .282; \eta^2 = 0.03$

Interaction effect: $F(1, 37) = 4.07; p = .051; \eta^2 = 0.10$

References

1. Mamede S, van Gog T, van den Berge K, Rikers RMJP, van Saase JLCM, van Guldener C, et al. Effect of Availability Bias and Reflective Reasoning on Diagnostic Accuracy Among Internal Medicine Residents. *Jama-Journal of the American Medical Association*. 2010;304(11):1198-203.
2. Schmidt HG, Mamede S, van den Berge K, van Gog T, van Saase JLCM, Rikers RMJP. Exposure to Media Information About a Disease Can Cause Doctors to Misdiagnose Similar-Looking Clinical Cases. *Academic Medicine*. 2014;89(2):285-91.
3. Stevens J. *Applied multivariate statistics for the social sciences*. 2nd ed. Hillsdale, NJ: Lawrence Erlbaum Associates; 1992.
4. Little RJA, Rubin DB. *Statistical analysis with missing data*. 2nd ed. New York: Wiley; 2002.
5. Rubin DB. Inference and missing data. *Biometrika*. 1976;63:581-92.
6. von Hippel PT. Regression with Missing Ys: An Improved Strategy for Analyzing Multiply Imputed Data. *Sociol Methodol*. 2007;37:83-117.
7. van Ginkel JR, Linting M, Rippe RCA, van der Voort A. Rebutting Existing Misconceptions About Multiple Imputation as a Method for Handling Missing Data. *J Pers Assess*. 2020;102(3):297-308.
8. Cohen-Schotanus J, van der Vleuten CPM. A standard setting method with the best performing students as point of reference: Practical and affordable. *Med Teach*. 2010;32(2):154-60.
9. DeCoster J, Gallucci M, Iselin AMR. Best Practices for using Median Splits, Artificial Categorization, and their Continuous Alternatives. *J Exp Psychopathology*. 2011;2(2):197-209.

10. Cohen J. Statistical power analysis for the behavioral sciences. Hillsdale, NJ: Lawrence Erlbaum Associates; 1988.