

Title: Large language models approach expert-level clinical knowledge and reasoning in ophthalmology: A head-to-head cross-sectional study

Authors: Arun James Thirunavukarasu^{1,2,*}, Shathar Mahmood¹, Andrew Malem³, William Paul Foster^{1,4}, Rohan Sanghera¹, Refaat Hassan¹, Sean Zhou⁵, Shiao Wei Wong⁶, Yee Ling Wong⁶, Yu Jeat Chong⁷, Abdullah Shakeel¹, Yin-Hsi Chang⁸, Benjamin Kye Jyn Tan⁹, Nikhil Jain¹⁰, Ting Fang Tan¹¹, Saaeha Rauz^{7,12}, Daniel Shu Wei Ting^{11,13,14}, Darren Shu Jeng Ting^{7,12,15*}

Affiliations:

1. University of Cambridge School of Clinical Medicine, Cambridge, United Kingdom
2. Oxford University Academic Graduate School, University of Oxford, Oxford, United Kingdom
3. Eye Institute, Cleveland Clinic Abu Dhabi, Abu Dhabi Emirate, United Arab Emirates
4. Department of Physiology, Development and Neuroscience, University of Cambridge, Cambridge, United Kingdom
5. West Suffolk NHS Foundation Trust, Bury St Edmunds, United Kingdom
6. Manchester Royal Eye Hospital, Manchester University NHS Foundation Trust, Manchester, United Kingdom
7. Birmingham and Midland Eye Centre, Sandwell and West Birmingham NHS Foundation Trust, Birmingham, United Kingdom
8. Department of Ophthalmology, Chang Gung Memorial Hospital, Linkou Medical Center, Taoyuan, Taiwan

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

9. Yong Loo Lin School of Medicine, National University of Singapore,
Singapore
10. Bedfordshire Hospitals NHS Foundation Trust, Luton and Dunstable, United
Kingdom
11. Singapore Eye Research Institute, Singapore National Eye Centre, Singapore,
Singapore
12. Academic Unit of Ophthalmology, Institute of Inflammation and Ageing,
University of Birmingham, Birmingham, United Kingdom
13. Duke-NUS Medical School, Singapore, Singapore
14. Byers Eye Institute, Stanford University, Palo Alto, California, United States of
America
15. Academic Ophthalmology, School of Medicine, University of Nottingham,
Nottingham, United Kingdom

***Correspondence details:**

Dr Darren Ting: Academic Ophthalmology, School of Medicine, University of
Nottingham, Nottingham, United Kingdom; ting.darren@gmail.com

Dr Arun Thirunavukarasu: University of Cambridge School of Clinical Medicine,
Cambridge, United Kingdom; ajt205@cantab.ac.uk

Funding: DSWT is supported by the National Medical Research Council, Singapore (NMCR/HSRG/0087/2018; MOH-000655-00; MOH-001014-00), Duke-NUS Medical School (Duke-NUS/RSF/2021/0018; 05/FY2020/EX/15-A58), and Agency for Science, Technology and Research (A20H4g2141; H20C6a0032). DSJT is supported by a Medical Research Council / Fight for Sight Clinical Research

Fellowship (MR/T001674/1). These funders were not involved in the conception, execution, or reporting of this review.

Ethical approval: Ethical approval was not required for this study as human patients were not involved in any experiments.

Competing interests: AM is a member of the Panel of Examiners of the Royal College of Ophthalmologists and performs unpaid work as an FRCOphth examiner. DSWT holds a patent on a deep learning system to detect retinal disease. DSJT authored the book used in the study and receives royalty from its sales. The other authors have no competing interests to declare.

Word count: 3548

ABSTRACT

Objective

To evaluate the clinical potential of large language models (LLMs) in ophthalmology using a more robust benchmark than raw examination scores.

Materials and methods

GPT-3.5 and GPT-4 were trialled on 347 questions before GPT-3.5, GPT-4, PaLM 2, LLaMA, expert ophthalmologists, and doctors in training were trialled on a mock examination of 87 questions. Performance was analysed with respect to question subject and type (first order recall and higher order reasoning). Masked ophthalmologists graded the accuracy, relevance, and overall preference of GPT-3.5 and GPT-4 responses to the same questions.

Results

The performance of GPT-4 (69%) was superior to GPT-3.5 (48%), LLaMA (32%), and PaLM 2 (56%). GPT-4 compared favourably with expert ophthalmologists (median 76%, range 64-90%), ophthalmology trainees (median 59%, range 57-63%), and unspecialised junior doctors (median 43%, range 41-44%). Low agreement between LLMs and doctors reflected idiosyncratic differences in knowledge and reasoning with overall consistency across subjects and types ($p>0.05$). All ophthalmologists preferred GPT-4 responses over GPT-3.5 and rated the accuracy and relevance of GPT-4 as higher ($p<0.05$).

Discussion

In view of the comparable or superior performance to trainee-grade ophthalmologists and unspecialised junior doctors, state-of-the-art LLMs such as GPT-4 may provide useful medical advice and assistance where access to expert ophthalmologists is limited. Clinical benchmarks provide useful assays of LLM capabilities in healthcare before clinical trials can be designed and conducted.

Conclusion

LLMs are approaching expert-level knowledge and reasoning skills in ophthalmology. Further research is required to develop and validate clinical applications to improve eye health outcomes.

Keywords

ChatGPT; large language model; natural language processing; decision support techniques; artificial intelligence; AI; deep learning; ophthalmology; examination; eye care; chatbot

INTRODUCTION

Generative Pre-trained Transformer 3.5 (GPT-3.5) and 4 (GPT-4) are large language models (LLMs) trained on datasets containing hundreds of billions of words from articles, books, and other internet sources.^{1,2} ChatGPT is an online chatbot which uses GPT-3.5 or GPT-4 to provide bespoke responses to human users' queries.³ LLMs have revolutionised the field of natural language processing, and in recent months, ChatGPT has attracted significant attention in medicine for attaining passing level performance in medical school examinations and providing more accurate and empathetic messages than human doctors in response to patient queries on a social media platform.³⁻⁶ While GPT-3.5 performance in more specialised examinations has been inadequate, GPT-4 is thought to represent a significant advancement in terms of medical knowledge and reasoning.^{3,7,8} Other LLMs in wide use include Pathways Language Model 2 (PaLM 2) and Large Language Model Meta AI 2 (LLaMA 2).^{3,9,10}

Applications and trials of LLMs in ophthalmological settings has been limited despite ChatGPT's performance in questions relating to 'eyes and vision' being superior to other subjects in an examination for general practitioners.^{7,11} ChatGPT has been trialled on the North American Ophthalmology Knowledge Assessment Program (OKAP), and Fellowship of the Royal College of Ophthalmologists (FRCOphth) Part 1 and Part 2 examinations. In both cases, relatively poor results have been reported for GPT-3.5, with significant improvement exhibited by GPT-4.¹²⁻¹⁶ However, previous studies are afflicted by two important issues which may affect their validity and interpretability. First, so-called 'contamination', where test material features in the pretraining data used to develop LLMs, may result in inflated

performance as models recall previously seen text rather than using clinical reasoning to provide an answer. Second, examination performance in and of itself provides little information regarding the potential of models to contribute to clinical practice as a medical-assistance tool.³ Clinical benchmarks are required to understanding the meaning and implications of scores in ophthalmological examinations attained by LLMs and are a necessary precursor to clinical trials of LLM-based interventions.

Here, we used FRCOphth Part 2 examination questions to gauge the ophthalmological knowledge base and reasoning capability of LLMs using fully qualified and currently training ophthalmologists as clinical benchmarks. These questions were not freely available online, minimising the risk of contamination. The FRCOphth Part 2 Written Examination tests the clinical knowledge and skills of ophthalmologists in training using multiple choice questions with no negative marking and must be passed to fully qualify as a specialist eye doctor in the United Kingdom.

METHODS

Question extraction

FRCOphth Part 2 questions were sourced from a textbook for doctors preparing to take the examination.¹⁷ This textbook is not freely available on the internet, making the possibility of its content being included in LLMs' training datasets unlikely.¹ All 360 questions from the textbook's six chapters were extracted. Two researchers matched the subject categories of the practice papers' questions to those defined in the Royal College of Ophthalmologists' documentation concerning the FRCOphth Part 2 written examination. Similarly, two researchers categorised each question as

first order recall or higher order reasoning, corresponding to ‘remembering’ and ‘applying’ or ‘analysing’ in Bloom’s taxonomy, respectively.¹⁸ Disagreement between classification decisions was resolved by a third researcher casting a deciding vote. Questions containing non-plain text elements such as images were excluded as these could not be inputted to the LLM applications.

Trialling large language models

Every eligible question was inputted into ChatGPT (GPT-3.5 and GPT-4 versions; OpenAI, San Francisco, California, United States of America) between April 29 and May 10, 2023. The answers provided by GPT-3.5 and GPT-4 were recorded and their whole reply to each question was recorded for further analysis. If ChatGPT failed to provide a definitive answer, the question was re-trialled up to three times, after which ChatGPT’s answer was recorded as ‘null’ if no answer was provided. Correct answers (‘ground truth’) were defined as the answers provided by the textbook and were recorded for every eligible question to facilitate calculation of performance. Upon their release, Bard (Google LLC, Mountain View, California, USA) and HuggingChat (Hugging Face, Inc., New York City, USA) were used to trial PaLM 2 (Google LLC) and LLaMA (Meta, Menlo Park, California, USA) respectively on the portion of the textbook corresponding to a 90-question examination, adhering to the same procedures between June 20 and July 2, 2023.

Clinical benchmarks

To gauge the performance, accuracy, and relevance of LLM outputs, five expert ophthalmologists who had all passed the FRCOphth Part 2 (E1-E5), three trainees (residents) currently in ophthalmology training programmes (T1-T3), and two unspecialised (*i.e.* not in ophthalmology training) junior doctors (J1-J2) first

answered the 90-question mock examination independently, without reference to textbooks, the internet, or LLMs' recorded answers. As with the LLMs, doctors' performance was calculated with reference to the correct answers provided by the textbook. After completing the examination, ophthalmologists graded the whole output of GPT-3.5 and GPT-4 on a Likert scale from 1-5 (very bad, bad, neutral, good, very good) to qualitatively appraise accuracy of information provided and relevance of outputs to the question used as an input prompt. For these appraisals, ophthalmologists were blind to the LLM source (which was presented in a randomised order) and to their previous answers to the same questions, but they could refer to the question text and correct answer and explanation provided by the textbook. Procedures are comprehensively described in the protocol issued to the ophthalmologists (Supplementary Material 2).

As our null hypothesis was that LLMs and doctors would exhibit similar performance, prospective power analysis was conducted which indicated that 63 questions were required to identify a 10% superior performance of an LLM to human performance at a 5% significance level (type 1 error rate) with 80% power (20% type 2 error rate). This indicated that the 90-question examination in our experiments was more than sufficient to detect ~10% differences in overall performance. The whole 90-question mock examination was used to avoid over- or under-sampling certain question types with respect to actual FRCOphth papers. To verify that the mock examination was representative of the FRCOphth Part 2 examination, expert ophthalmologists were asked to rate the difficulty of questions used here in comparison to official examinations on a 5-point Likert scale ("much

easier”, “somewhat easier”, “similar”, “somewhat more difficult”, “much more difficult”).

Statistical analysis

Performance of doctors and LLMs were compared using chi-squared (χ^2) tests. Agreement between answers provided by doctors and LLMs was quantified through calculation of Kappa statistics, interpreted in accordance with McHugh’s recommendations.¹⁹ To further explore the strengths and weaknesses of the answer providers, performance was stratified by question type (first order fact recall or higher order reasoning) and subject using a chi-squared or Fisher’s exact test where appropriate. Likert scale data corresponding to the accuracy and relevance of GPT-3.5 and GPT-4 responses to the same questions were analysed with paired *t*-tests with the Bonferroni correction applied to mitigate the risk of false positive results due to multiple-testing—parametric testing was justified by a sufficient sample size.²⁰ Statistical significance was concluded where $p < 0.05$. For additional contextualisation, examination statistics corresponding to FRCOphth Part 2 written examinations taken between July 2017 and December 2022 were collected from Royal College of Ophthalmologists examiners’ reports.²¹ These statistics facilitated comparisons between human and LLM performance in the mock examination with the performance of actual candidates in recent examinations.

Statistical analysis was conducted in R (version 4.1.2; R Foundation for Statistical Computing, Vienna, Austria), and figures were produced in Affinity Designer (version 1.10.6; Serif Ltd, West Bridgford, Nottinghamshire, United Kingdom).

RESULTS

Questions sources

Of 360 questions in the textbook, 347 questions (including 87 of the 90 questions from the mock examination chapter) were included.¹⁷ Exclusions were all due to non-text elements such as images and tables which could not be inputted into LLM chatbot interfaces. The distribution of question types and subjects within the whole set and mock examination set of questions is summarised in Table 1 and Table S1 alongside performance.

GPT-4 represents a significant advance on GPT-3.5 in ophthalmological knowledge and reasoning

Overall performance over 347 questions was significantly higher for GPT-4 (61.7%) than GPT-3.5 (48.41%; $\chi^2=12.32$, $p<0.01$), with results detailed in Figure S1 and Table S1. ChatGPT performance was consistent across question types and subjects (Table S1). For GPT-4, no significant variation was observed with respect to first order and higher order questions ($\chi^2 = 0.22$, $p=0.64$), or subjects defined by the Royal College of Ophthalmologists (Fisher's exact test over 2000 iterations, $p = 0.23$). Similar results were observed for GPT-3.5 with respect to first and second order questions ($\chi^2 = 0.08$, $p = 0.77$), and subjects (Fisher's exact test over 2000 iterations, $p = 0.28$). Performance and variation within the 87-question mock examination was very similar to the overall performance over 347 questions, and subsequent experiments were therefore restricted to that representative set of questions.

GPT-4 compares well with other LLMs, junior and trainee doctors and ophthalmology experts

Performance in the mock examination is summarised in Figure 1 – GPT-4 (69%) was

the top-scoring model, performing to a significantly higher standard than GPT-3.5 (48%; $\chi^2 = 7.33$, $p < 0.01$) and LLaMA (32%; $\chi^2 = 22.77$, $p < 0.01$), but statistically similarly to PaLM 2 (56%) despite a superior score ($\chi^2 = 2.81$, $p = 0.09$). LLaMA exhibited the lowest examination score, significantly weaker than GPT-3.5 ($\chi^2 = 4.58$, $p = 0.03$) and PaLM-2 ($\chi^2 = 10.01$, $p < 0.01$) as well as GPT-4.

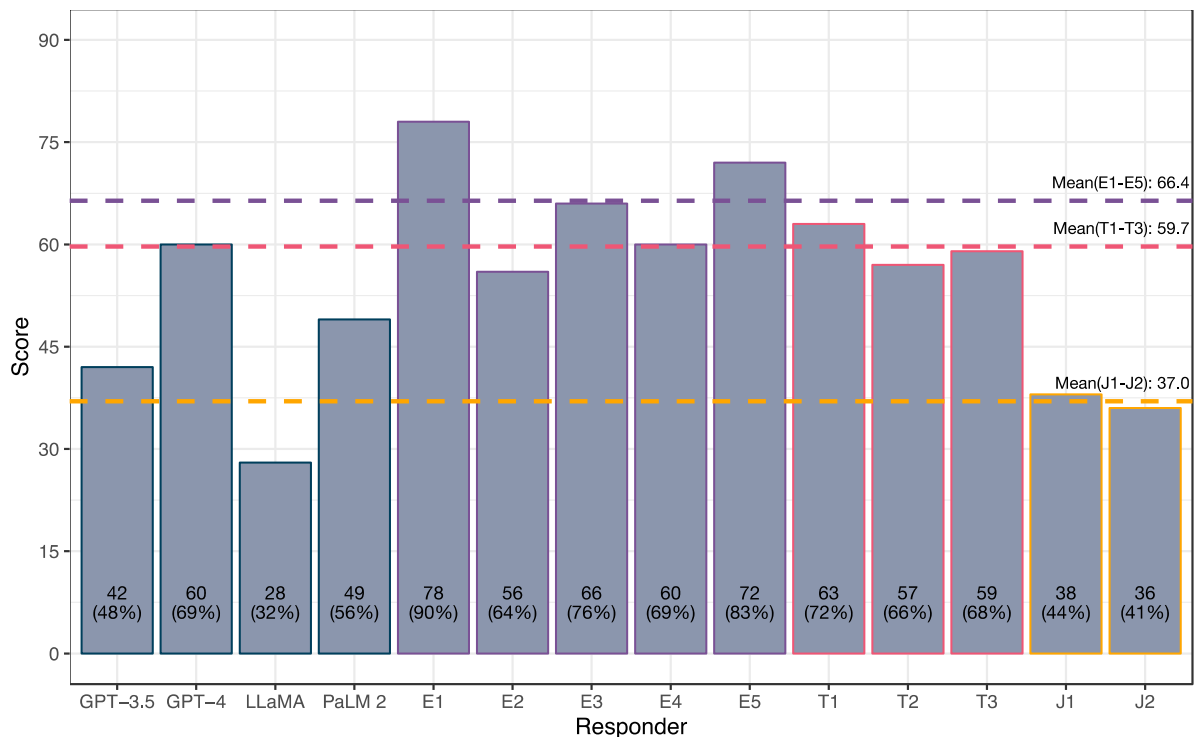


Figure 1: FRCOphth Part 2 performance of LLMs and doctors of variable expertise. Examination performance in the 87-question mock examination used to trial LLMs (GPT-3.5, GPT-4, LLaMA, and PaLM 2), expert ophthalmologists (E1-E5), ophthalmology trainees (T1-T3), and unspecialised junior doctors (J1-J2). Dotted lines depict the mean performance of expert ophthalmologists (66/87; 76%), ophthalmology trainees (60/87; 69%), and unspecialised junior doctors (37/87; 43%). The performance of GPT-4 lay within the range of expert ophthalmologists and ophthalmology trainees.

The performance of GPT-4 was statistically similar to the mean score attained by expert ophthalmologists (Figure 1; $\chi^2 = 1.18$, $p = 0.28$). Moreover, GPT-4's performance exceeded the mean mark attained across FRCOphth Part 2 written examination candidates between 2017-2022 (66.06%), mean pass mark according to standard setting (61.31%), and the mean official mark required to pass the examination after adjustment (63.75%), as detailed in Table S2. In individual comparisons with expert ophthalmologists, GPT-4 was equivalent in 3 cases (χ^2 tests, $p > 0.05$; Table S3), and inferior in 2 cases (χ^2 tests, $p < 0.05$; Table 2). In comparisons with ophthalmology trainees, GPT-4 was equivalent to all three ophthalmology trainees (χ^2 tests, $p > 0.05$; Table 2). GPT-4 was significantly superior to both unspecialised trainee doctors (χ^2 tests, $p < 0.05$; Table 2). Doctors were anonymised in analysis, but their ophthalmological experience is summarised in Table S3. Unsurprisingly, junior doctors (J1-J2) attained lower scores than expert ophthalmologists (E1-E5; $t = 7.18$, $p < 0.01$), and ophthalmology trainees (T1-T3; $t = 11.18$, $p < 0.01$), illustrated in Figure 1. Ophthalmology trainees approached expert-level scores with no significant difference between the groups ($t = 1.55$, $p = 0.18$). None of the other LLMs matched any of the expert ophthalmologists, mean mark of real examination candidates, or FRCOphth Part 2 pass mark.

Expert ophthalmologists agreed that the mock examination was a faithful representation of actual FRCOphth Part 2 Written Examination papers with a mean and median score of 3/5 (range 2-4/5).

LLM strengths and weaknesses are similar to doctors

Agreement between answers given by LLMs, expert ophthalmologists, and trainee doctors was generally absent ($0 \leq \kappa < 0.2$), minimal ($0.2 \leq \kappa < 0.4$), or weak ($0.4 \leq \kappa$

< 0.6), with moderate agreement only recorded for one pairing between the two highest performing ophthalmologists (Figure 2; $\kappa = 0.64$).¹⁹ Disagreement was primarily the result of general differences in knowledge and reasoning ability, illustrated by strong negative correlation between Kappa statistic (quantifying agreement) and difference in examination performance (Pearson's $r = -0.63$, $p < 0.01$). Answer providers with more similar scores exhibited greater agreement overall irrespective of their category (LLM, expert ophthalmologist, ophthalmology trainee, or junior doctor).

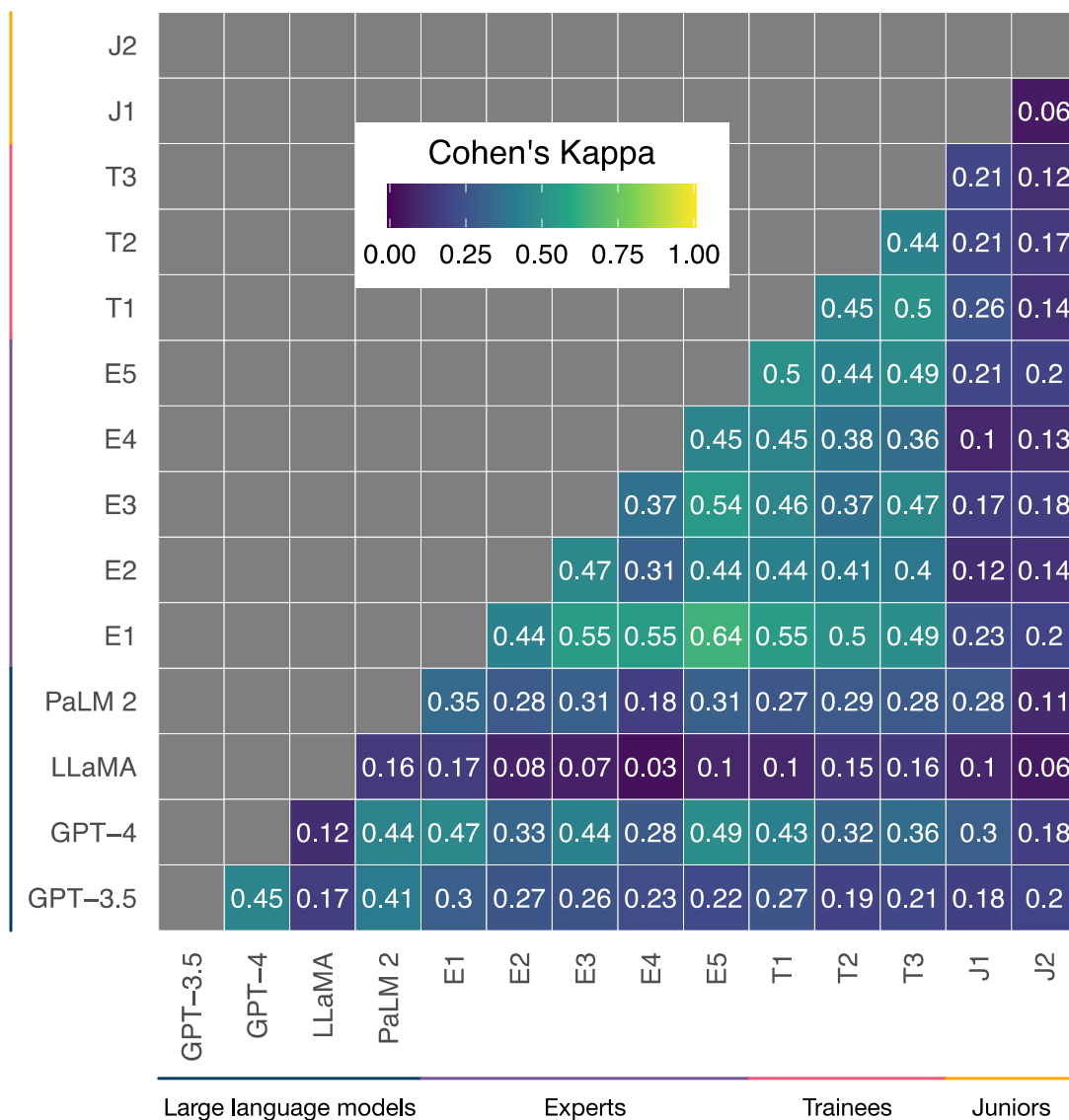


Figure 2: Heat map of Kappa statistics quantifying agreement between answers

given by LLMs, expert ophthalmologists, and trainee doctors. Agreement correlates strongly with overall performance and stratification analysis found no particular question type or subject was associated with better performance of LLMs or doctors, indicating that LLM knowledge and reasoning ability is general across ophthalmology rather than restricted to particular subspecialties or question types.

Stratification analysis was undertaken to identify any specific strengths and weaknesses of LLMs with respect to expert ophthalmologists and trainee doctors (Table 1, Table S4). No significant difference between performance in first order fact recall and higher order reasoning questions was observed among any of the LLMs, expert ophthalmologists, ophthalmology trainees, or unspecialised junior doctors (Table S4; χ^2 tests, $p > 0.05$). Similarly, only J1 (junior doctor yet to commence ophthalmology training) exhibited statistically significant variation in performance between subjects (Table S4; Fisher's exact tests over 2000 iterations, $p = 0.02$); all other doctors and LLMs exhibited no significant variation (Fisher's exact tests over 2000 iterations, $p > 0.05$). To explore whether consistency was due to an insufficient sample size, similar analyses were run for GPT-3.5 and GPT-4 performance over the larger set of 347 questions (Table S1; Table S4). As with the mock examination, no significant differences in performance across question types (Table S4; χ^2 tests, $p > 0.05$) or subjects (Table S4; Fisher's exact tests over 2000 iterations, $p > 0.05$) were observed.

LLM examination performance translates to subjective preference indicated by expert ophthalmologists

Ophthalmologists' appraisal of GPT-4 and GPT-3.5 outputs indicated a marked preference for the former over the latter, mirroring objective performance in the

mock examination and over the whole textbook. GPT-4 exhibited significantly (t -test with Bonferroni correction, $p < 0.05$) higher accuracy and relevance than GPT-3.5 according to all five ophthalmologists' grading (Table 3). Differences were visually obvious, with GPT-4 exhibiting much higher rates of attaining the highest scores for accuracy and relevance than GPT-3.5 (Figure 3). This superiority was reflected in ophthalmologists' qualitative preference indications: GPT-4 responses were preferred to GPT-3.5 responses by every ophthalmologist with statistically significant skew in favour of GPT-4 (χ^2 test, $p < 0.05$; Table 3).

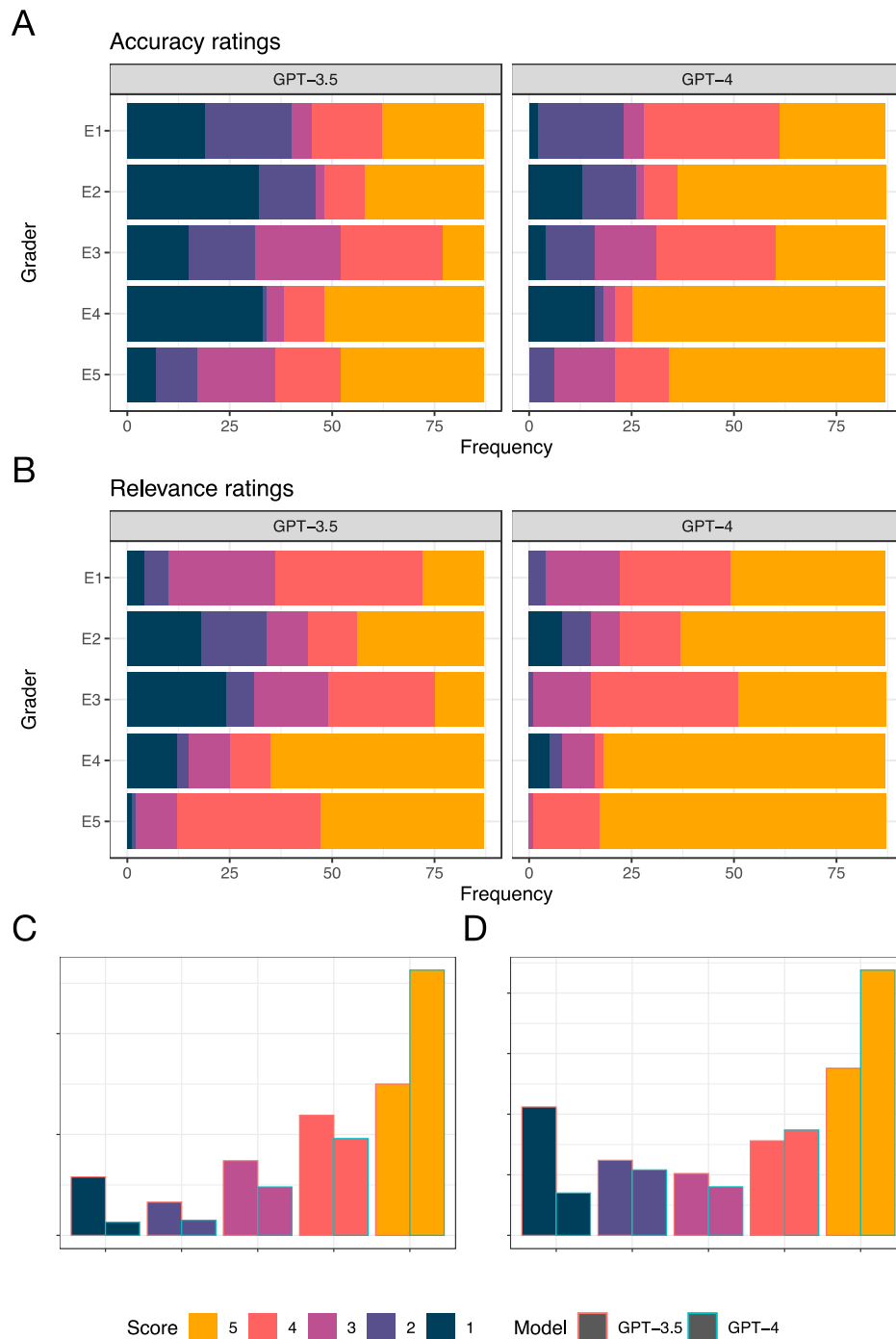


Figure 3: Accuracy and relevance of GPT-3.5 and GPT-4 in response to ophthalmological questions. Accuracy (A) and relevance (B) ratings were provided by five expert ophthalmologists for ChatGPT (powered by GPT-3.5 and GPT-4) responses to 87 FRCOphth Part 2 mock examination questions. In every case, the accuracy and relevance of GPT-4 is significantly superior to GPT-3.5 (t-test with

Bonferroni correct applied, $p < 0.05$). Pooled scores for accuracy (C) and relevance (D) from all five raters are presented in the bottom two plots, with GPT-3.5 (left bars) compared directly with GPT-4 (right bars).

DISCUSSION

Here, we present a clinical benchmark to gauge the ophthalmological performance of LLMs, using a source of questions with very low risk of contamination as the utilised textbook is not freely available online.¹⁷ Previous studies have suggested that ChatGPT can provide useful responses to ophthalmological queries, but often use online question sources which may have featured in LLMs' pretraining datasets.^{7,12,15,22} In addition, our employment of multiple LLMs as well as fully qualified and training doctors provides novel insight into the potential and limitations of state-of-the-art LLMs through head-to-head comparisons which provide clinical context and quantitative benchmarks of competence in ophthalmology. Subsequent research may leverage our questions and results to gauge the performance of new LLMs and applications as they emerge.

We make three primary observations. First, performance of GPT-4 compares well to expert ophthalmologists and ophthalmology trainees, and exhibits pass-worthy performance in an FRCOphth Part 2 mock examination. PaLM 2 did not attain pass-worthy performance or match expert ophthalmologists' scores but was within the spread of trainee doctors' performance. LLMs are approaching human expert-level knowledge and reasoning in ophthalmology, and significantly exceed the ability of non-specialist clinicians (represented here by unspecialised junior doctors) to answer ophthalmology questions. Second, clinician grading of model outputs suggests that GPT-4 exhibits improved accuracy and relevance when compared

with GPT-3.5. Development is producing models which generate better outputs to ophthalmological queries in the opinion of expert human clinicians, which suggests that models are becoming more capable of providing useful assistance in clinical settings. Third, LLM performance was consistent across question subjects and types, distributed similarly to human performance, and exhibited comparable agreement between other LLMs and doctors when corrected for differences in overall performance. Together, this indicates that the ophthalmological knowledge and reasoning capability of LLMs is general rather than limited to certain subspecialties or tasks. LLM-driven natural language processing seems to facilitate similar—although idiosyncratic—clinical knowledge and reasoning to human clinicians, with no obvious blind spots precluding clinical use.

Similarly dramatic improvements in the performance of GPT-4 relative to GPT-3.5 have been reported in the context of the North American Ophthalmology Knowledge Assessment Program (OKAP).^{13,15} State-of-the-art models exhibit far more clinical promise than their predecessors, and expectations and development should be tailored accordingly. Results from the OKAP also suggest that improvement in performance is due to GPT-4 being more well-rounded than GPT-3.5.¹³ This increases the scope for potential applications of LLMs in ophthalmology, as development is eliminating weaknesses rather than optimising in narrow domains. This study shows that well-rounded LLM performance compares well with expert ophthalmologists, providing clinically relevant evidence that LLMs may be used to provide medical advice and assistance. Further improvement is expected as multimodal foundation models, perhaps based on LLMs such as GPT-4, emerge and facilitate compatibility with image-rich ophthalmological data.^{3,23,24}

Limitations

This study was limited by three factors. First, examination performance is an unvalidated indicator of clinical aptitude. We sought to ameliorate this limitation by employing expert ophthalmologists, ophthalmology trainees, and unspecialised junior doctors answering the same questions as clinical benchmarks; and compared LLM performance to real cohorts of candidates in recent FRCOphth examinations. However, it remains an issue that comparable performance to clinical experts in an examination does not necessarily demonstrate that an LLM can communicate with patients and practitioners or contribute to clinical decision making accurately and safely. Early trials of LLM chatbots have suggested that LLM responses may be equivalent or even superior to human doctors in terms of accuracy and empathy, and experiments using complicated case studies suggest that LLMs operate well even outside typical presentations and more common medical conditions.^{4,25,26} In ophthalmology, GPT-3.5 and GPT-4 have been shown to be capable of providing precise and suitable triage decisions when queried with eye-related symptoms.^{22,27} Further work is now warranted in conventional clinical settings.

Second, while the study was sufficiently powered to detect a less than 10% difference in overall performance, the relatively small number of questions in certain categories used for stratification analysis may mask significant differences in performance. Testing LLMs and clinicians with more questions may help establish where LLMs exhibit greater or lesser ability in ophthalmology. Furthermore, researchers using different ways to categorise questions may be able to identify specific strengths and weaknesses of LLMs and doctors which could help guide design of clinical LLM interventions.

Finally, experimental tasks were ‘zero-shot’ in that LLMs were not provided with any examples of correctly answered questions before it was queried with FRCOphth questions from the textbook. This mode of interrogation entails the maximal level of difficulty for LLMs, so it is conceivable that the ophthalmological knowledge and reasoning encoded within these models is actually even greater than indicated by results here.¹ Future research may seek to fine-tune LLMs by using more domain-specific text during pretraining and fine-tuning, or by providing examples of successfully completed tasks to further improve performance in that clinical task.³

Future directions

Autonomous deployment of LLMs is currently precluded by inaccuracy and fact fabrication. Our study found that despite meeting expert standards, state-of-the-art LLMs such as GPT-4 do not match top-performing ophthalmologists.²⁸ Moreover, there remain controversial ethical questions about what roles should and should not be assigned to inanimate AI models, and to what extent human clinicians must remain responsible for their patients.³ However, the remarkable performance of GPT-4 in ophthalmology examination questions suggests that LLMs may be able to provide useful input in clinical contexts, either to assist clinicians in their day-to-day work or with their education or preparation for examinations.^{3,13,14,27} GPT-4 may prove especially useful where access to ophthalmologists is limited: provision of advice, diagnosis, and management suggestions by a model with FRCOphth Part 2-level knowledge and reasoning ability is likely to be superior to non-specialist doctors and allied healthcare professionals working without support, as their exposure to and knowledge of eye care is limited.^{27,29}

However, close monitoring is essential to avoid mistakes caused by inaccuracy or fact fabrication.³⁰ Clinical applications would also benefit from an uncertainty indicator reducing the risk of erroneous decisions.⁷ As LLM performance often correlates with the frequency of query terms' representation in the model's training dataset, a simple indicator of 'familiarity' could be engineered by calculating the relative frequency of query term representation in the training data.^{7,31} Users could appraise familiarity to temper their confidence in answers provided by the LLM, perhaps reducing error. Moreover, ophthalmological applications require extensive validation, preferably with high quality randomised controlled trials to conclusively demonstrate benefit (or lack thereof) conferred to patients by LLM interventions. Trials should be pragmatic so as not to inflate effect sizes beyond what may generalise to patients once interventions are implemented at scale.^{3,32,33} In addition to patient outcomes, practitioner-related variables should also be considered: interventions aiming to improve efficiency should be specifically tested to ensure that they reduce rather than increase clinicians' workload.³

Conclusion

According to comparisons with expert and trainee doctors, state-of-the-art LLMs are approaching expert-level performance in advanced ophthalmology questions. GPT-4 attains pass-worthy performance in FRCOphth Part 2 questions and exceeds the scores of some expert ophthalmologists. As top-performing doctors exhibit superior scores, LLMs do not appear capable of replacing ophthalmologists, but state-of-the-art models could provide useful advice and assistance to non-specialists or patients where access to eye care professionals is limited.^{27,28} Further research is required to design LLM-based interventions which may improve eye

health outcomes, validate interventions in clinical trials, and engineer governance structures to regulate LLM applications as they begin to be deployed in clinical settings.³⁴

Abbreviations

AI: artificial intelligence

FRCOphth: Fellowship of the Royal College of Ophthalmologists

GPT: Generative Pretrained Transformer

LLaMA: Large Language Model Meta AI

LLM: large language model

OKAP: Ophthalmology Knowledge Assessment Program

PaLM: Pathways Language Model

Funding

There were no funders which supported this project.

Ethical approval

Ethical approval was not required for this study, as human patients were not involved.

Author contributions

AJT and DSJT conceived the study and coordinated the team. AJT, DSJT, and YLW contributed to study design. AJT, WF, RH, SM, and RS undertook LLM data collection. YLW, SWW, AM, YJC, SZ, AS, YHC, BKJT, NJ, and TFT answered questions independently to benchmark LLM performance. YLW, SWW, AM, YJC, and SZ appraised LLM outputs. AJT conducted data analysis and visualisation. AJT and DSJT drafted the manuscript. All authors participated in manuscript redrafting.

DSJT provided academic supervision. All authors approved the submitted version of the manuscript.

Acknowledgements

The authors extend their thanks to Mr Arunachalam Thirunavukarasu (Betsi Cadwaladr University Health Board) for his advice and assistance with recruitment.

References

1. Brown, T. *et al.* Language Models are Few-Shot Learners. in *Advances in Neural Information Processing Systems* vol. 33 1877–1901 (Curran Associates, Inc., 2020).
2. OpenAI. GPT-4 Technical Report. Preprint at <https://doi.org/10.48550/arXiv.2303.08774> (2023).
3. Thirunavukarasu, A. J. *et al.* Large language models in medicine. *Nat Med* 1–11 (2023) doi:10.1038/s41591-023-02448-8.
4. Ayers, J. W. *et al.* Comparing Physician and Artificial Intelligence Chatbot Responses to Patient Questions Posted to a Public Social Media Forum. *JAMA Internal Medicine* (2023) doi:10.1001/jamainternmed.2023.1838.
5. Gilson, A. *et al.* How Does ChatGPT Perform on the United States Medical Licensing Examination? The Implications of Large Language Models for Medical Education and Knowledge Assessment. *JMIR Med Educ* 9, e45312 (2023).
6. Kung, T. H. *et al.* Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models. *PLOS Digital Health* 2, e0000198 (2023).
7. Thirunavukarasu, A. J. *et al.* Trialling a Large Language Model (ChatGPT) in General Practice With the Applied Knowledge Test: Observational Study

- Demonstrating Opportunities and Limitations in Primary Care. *JMIR Medical Education* **9**, e46599 (2023).
8. Nori, H., King, N., McKinney, S. M., Carignan, D. & Horvitz, E. Capabilities of GPT-4 on Medical Challenge Problems. Preprint at <http://arxiv.org/abs/2303.13375> (2023).
 9. Google. PaLM 2 Technical Report. Preprint at <https://ai.google/static/documents/palm2techreport.pdf> (2023).
 10. Touvron, H., Martin, L. & Stone, K. Llama 2: Open Foundation and Fine-Tuned Chat Models. Preprint at <https://ai.meta.com/research/publications/llama-2-open-foundation-and-fine-tuned-chat-models/> (2023).
 11. Ting, D. S. J., Tan, T. F. & Ting, D. S. W. ChatGPT in ophthalmology: the dawn of a new era? *Eye (Lond)* (2023) doi:10.1038/s41433-023-02619-4.
 12. Antaki, F., Touma, S., Milad, D., El-Khoury, J. & Duval, R. Evaluating the Performance of ChatGPT in Ophthalmology: An Analysis of its Successes and Shortcomings. *Ophthalmology Science* **0**, (2023).
 13. Teebagy, S., Colwell, L., Wood, E., Yaghy, A. & Faustina, M. Improved Performance of ChatGPT-4 on the OKAP Exam: A Comparative Study with ChatGPT-3.5. 2023.04.03.23287957 Preprint at <https://doi.org/10.1101/2023.04.03.23287957> (2023).
 14. Raimondi, R., Tzoumas, N., Salisbury, T., Di Simplicio, S. & Romano, M. R. Comparative analysis of large language models in the Royal College of Ophthalmologists fellowship exams. *Eye* 1–4 (2023) doi:10.1038/s41433-023-02563-3.

15. Mihalache, A., Popovic, M. M. & Muni, R. H. Performance of an Artificial Intelligence Chatbot in Ophthalmic Knowledge Assessment. *JAMA Ophthalmology* **141**, 589–597 (2023).
16. Thirunavukarasu, A. J. ChatGPT cannot pass FRCOphth examinations: implications for ophthalmology and large language model artificial intelligence. *Eye News* vol. 30 (2023).
17. Ting, D. S. J. & Steel, D. *MCQs for FRCOphth Part 2*. (Oxford University Press, 2020).
18. Adams, N. E. Bloom’s taxonomy of cognitive learning objectives. *J Med Libr Assoc* **103**, 152–153 (2015).
19. McHugh, M. L. Interrater reliability: the kappa statistic. *Biochemia Medica* **22**, 276–282 (2012).
20. Sullivan, G. M. & Artino, A. R. Analyzing and Interpreting Data From Likert-Type Scales. *J Grad Med Educ* **5**, 541–542 (2013).
21. Part 2 Written FRCOphth Exam. *The Royal College of Ophthalmologists* <https://www.rcophth.ac.uk/examinations/rcophth-exams/part-2-written-frcophth-exam/>.
22. Tsui, J. C. *et al.* Appropriateness of ophthalmic symptoms triage by a popular online artificial intelligence chatbot. *Eye* 1–2 (2023) doi:10.1038/s41433-023-02556-2.
23. Nath, S., Marie, A., Ellershaw, S., Korot, E. & Keane, P. A. New meaning for NLP: the trials and tribulations of natural language processing with GPT-3 in ophthalmology. *British Journal of Ophthalmology* **106**, 889–892 (2022).

24. Kline, A. *et al.* Multimodal machine learning in precision health: A scoping review. *npj Digit. Med.* **5**, 1–14 (2022).
25. Kulkarni, P. A. & Singh, H. Artificial Intelligence in Clinical Diagnosis: Opportunities, Challenges, and Hype. *JAMA* (2023)
doi:10.1001/jama.2023.11440.
26. Singhal, K. *et al.* Large language models encode clinical knowledge. *Nature* 1–9 (2023) doi:10.1038/s41586-023-06291-2.
27. Waisberg, E. *et al.* GPT-4 for triaging ophthalmic symptoms. *Eye* 1–2 (2023)
doi:10.1038/s41433-023-02595-9.
28. Thirunavukarasu, A. J. Large language models will not replace healthcare professionals: curbing popular fears and hype. *J R Soc Med* **116**, 181–182 (2023).
29. Alsaedi, M. G., Alhujaili, H. O., Fairaq, G. S., Alwdaan, S. A. & Alwadan, R. A. Emergent Ophthalmic Disease Knowledge among Non-Ophthalmologist Healthcare Professionals in the Western Region of Saudi Arabia: Cross-Sectional Study. *The Open Ophthalmology Journal* **16**, (2022).
30. Bakken, S. AI in health: keeping the human in the loop. *Journal of the American Medical Informatics Association* **30**, 1225–1226 (2023).
31. Biderman, S. *et al.* Pythia: A Suite for Analyzing Large Language Models Across Training and Scaling. Preprint at <https://doi.org/10.48550/arXiv.2304.01373> (2023).
32. Tossaint-Schoenmakers, R., Versluis, A., Chavannes, N., Talboom-Kamp, E. & Kasteleyn, M. The Challenge of Integrating eHealth Into Health Care: Systematic

Literature Review of the Donabedian Model of Structure, Process, and Outcome. *J Med Internet Res* **23**, e27180 (2021).

33. Thirunavukarasu, A. J., Hassan, R., Limonard, A. & Savant, S. V. Accuracy and reliability of self-administered visual acuity tests: Systematic review of pragmatic trials. *PLOS ONE* **18**, e0281847 (2023).
34. Meskó, B. & Topol, E. J. The imperative for regulatory oversight of large language models (or generative AI) in healthcare. *npj Digit. Med.* **6**, 1–6 (2023).

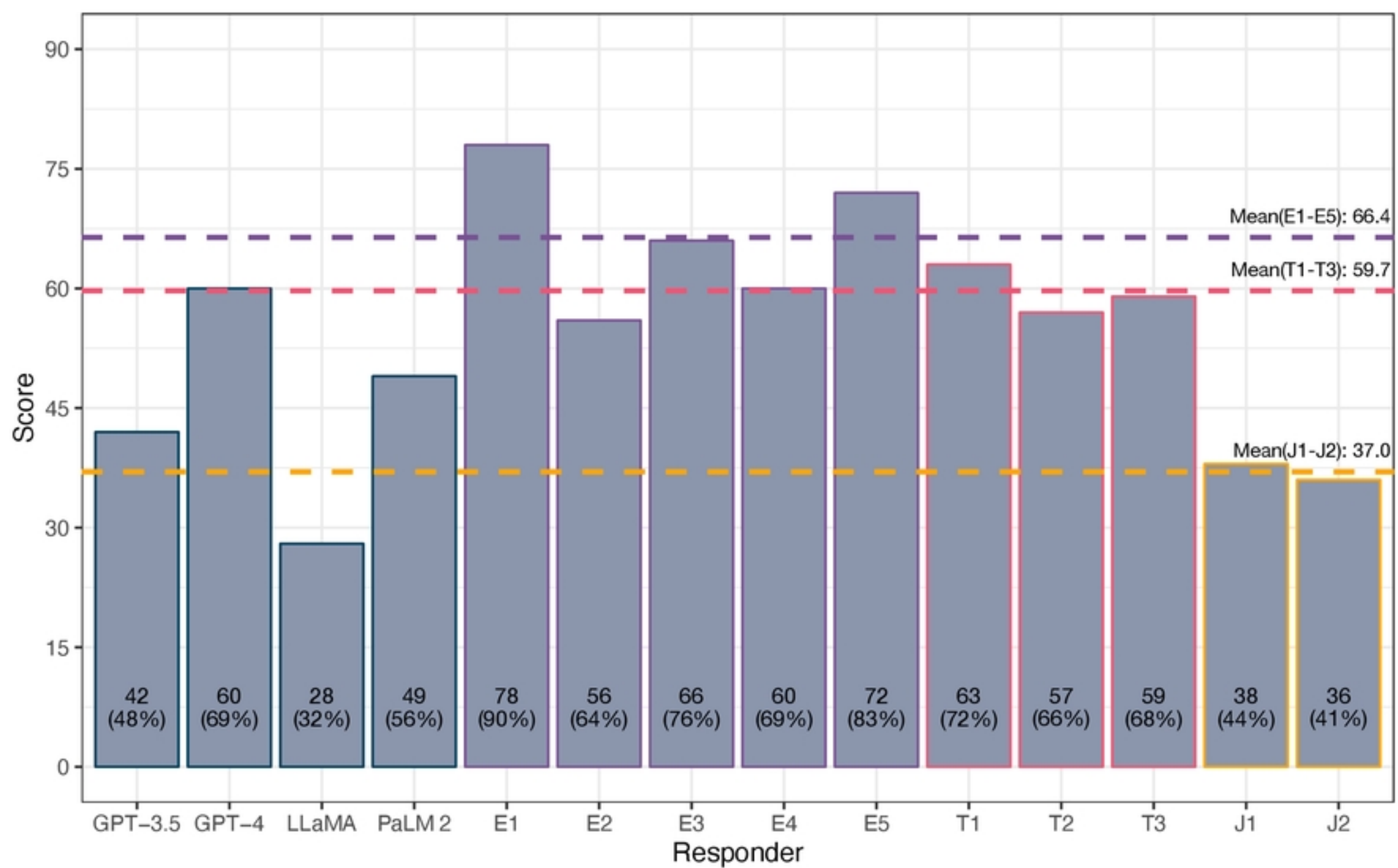


Figure 1

medRxiv preprint doi: <https://doi.org/10.1101/2023.07.31.23293474>; this version posted August 6, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

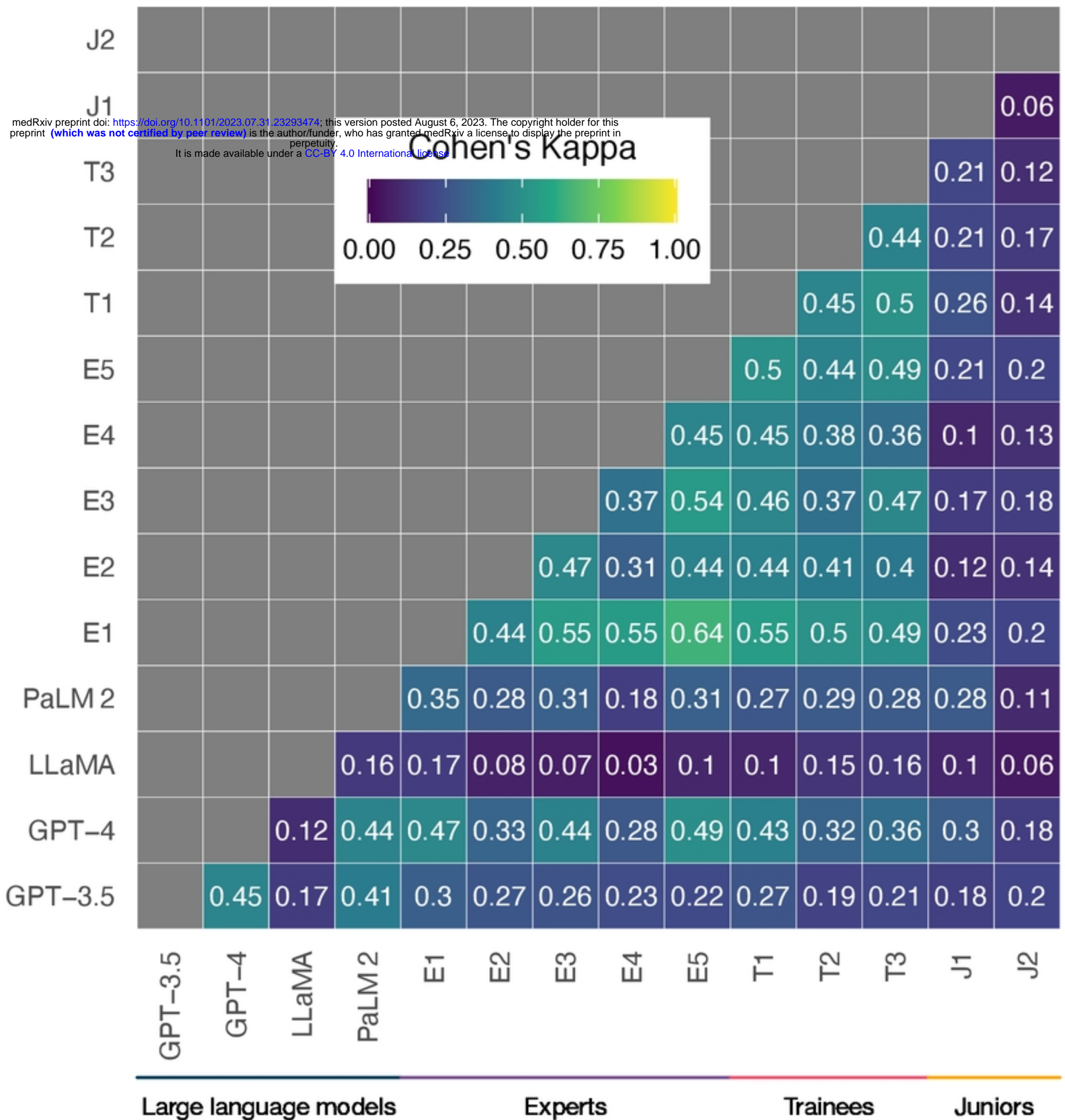
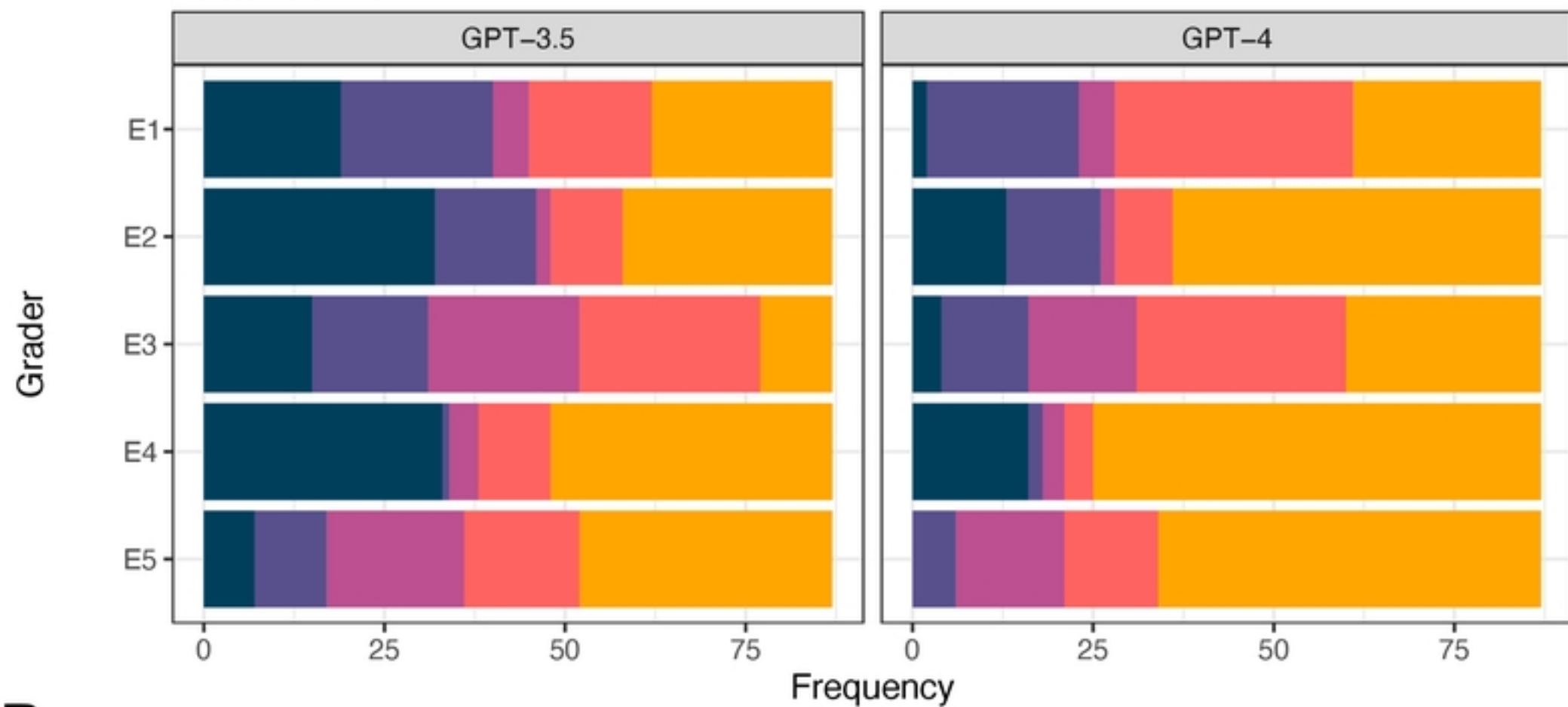


Figure 2

A

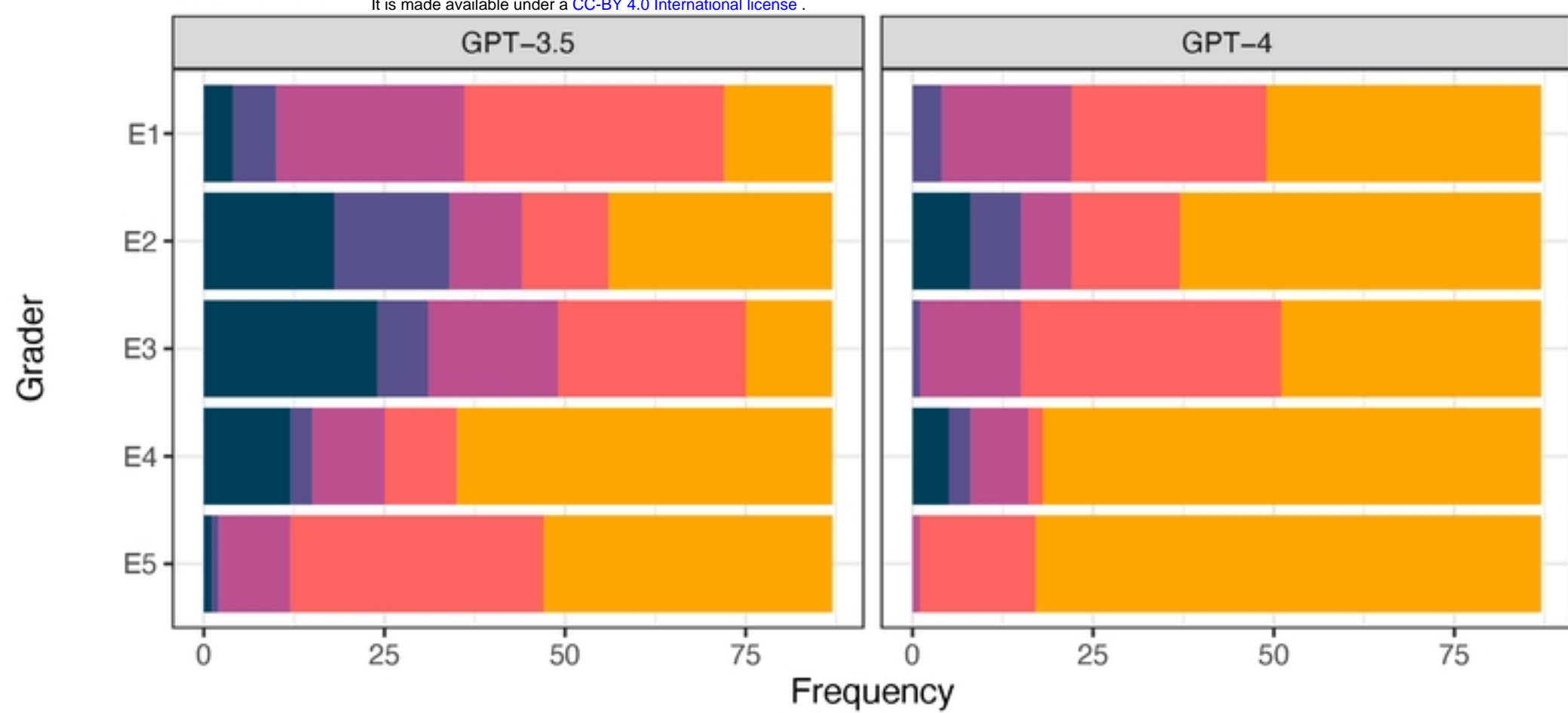
Accuracy ratings



B

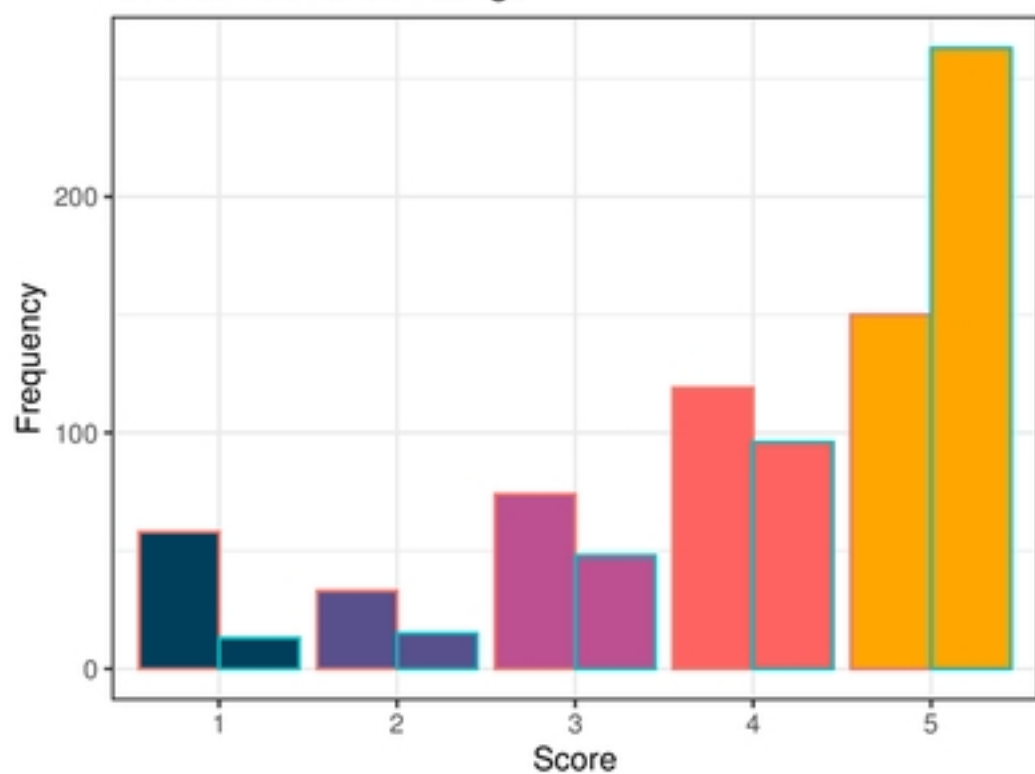
medRxiv preprint doi: <https://doi.org/10.1101/2023.07.31.23293474>; this version posted August 6, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Relevance ratings



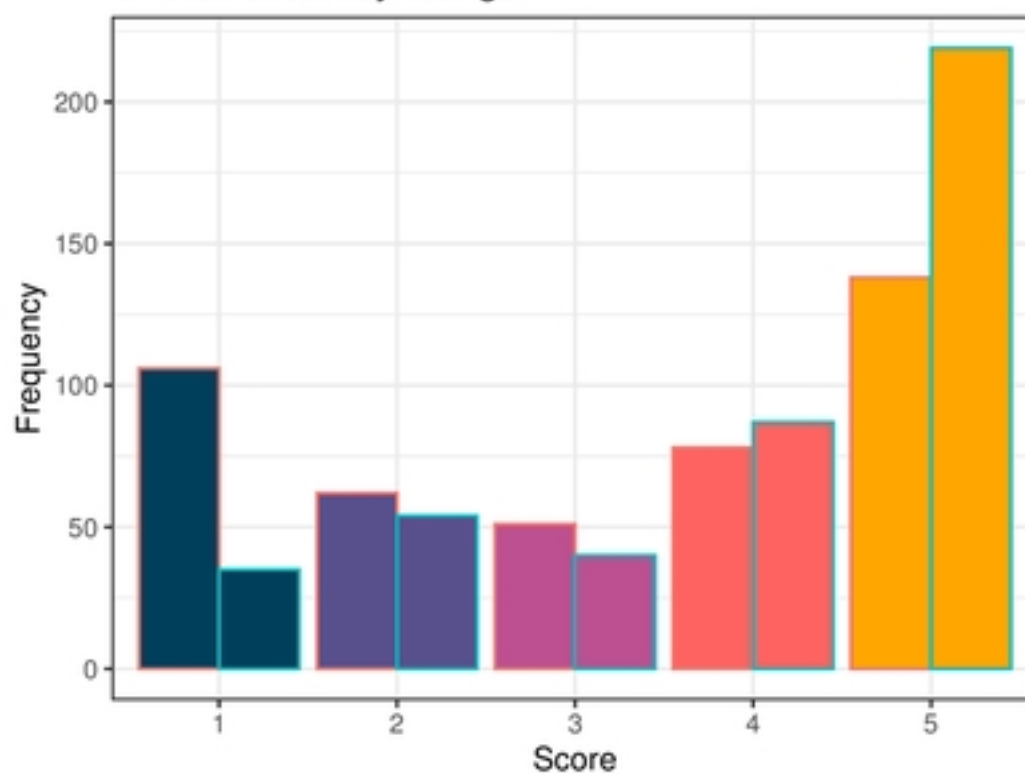
C

Overall relevance ratings



D

Overall accuracy ratings



Score 5 4 3 2 1

Model GPT-3.5 GPT-4

Figure 3