

Evaluating the Limitations of Large Language Models in Therapeutic Decision-making for patients with Aortic Stenosis

AUTHORS

*Tobias Roeschl, MD^{1,2,3,4}, *Marie Hoffmann, PhD^{1,2,4}, Djawid Hashemi, MD^{2,3,4,5}, Felix Rarreck^{1,2}, Nils Hinrichs, MSc^{1,2,4}, Tobias D. Trippel, MD^{2,4,5}, Axel Unbehaun, MD^{1,2}, Christoph Klein, MD^{2,5}, Jörg Kempfert, MD^{1,2}, Henryk Dreger, MD^{2,5}, Benjamin O'Brien, MD^{2,6,7}, Gerhard Hindricks, MD^{2,5}, Felix Balzer, MD, PhD^{2,8}, Volkmar Falk, MD^{1,2,4,9}, Alexander Meyer, MD^{1,2,4,10}

¹Deutsches Herzzentrum der Charité (DHZC), Department of Cardiothoracic and Vascular Surgery, Augustenburger Platz 1, 13353 Berlin, Germany

²Charité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt Universität zu Berlin, and Berlin Institute of Health, Augustenburger Platz 1, 13353 Berlin, Deutschland/Germany

³Berlin Institute of Health at Charité – Universitätsmedizin Berlin, BIH Biomedical Innovation Academy, BIH Charité Digital Clinician Scientist Program, Charitéplatz 1, 10117 Berlin, Germany

⁴DZHK (German Centre for Cardiovascular Research), partner site Berlin, Berlin, Germany

⁵Deutsches Herzzentrum der Charité (DHZC), Department of Cardiology, Angiology and Intensive Care Medicine, Augustenburger Platz 1, 13353 Berlin, Germany

⁶Deutsches Herzzentrum der Charité (DHZC), Department of Cardiac Anesthesiology and Intensive Care Medicine, Augustenburger Platz 1, 13353 Berlin, Germany

⁷St Bartholomew's Hospital and Barts Heart Centre, Department of Perioperative Medicine, London EC1A 7BE, UK

⁸Charité – Universitätsmedizin Berlin, Institute of Medical Informatics, Charitéplatz 1, 10117 Berlin, Germany

⁹Translational Cardiovascular Technologies, Institute of Translational Medicine, Department of Health Sciences and Technology, Swiss Federal Institute of Technology, 8092 Zürich, Switzerland.

¹⁰Berlin Institute for the Foundations of Learning and Data – TU Berlin, Berlin, Germany

* both authors contributed equally

Corresponding Author:

Prof. Dr. med. Alexander Meyer
Deutsches Herzzentrum der Charité
Campus Virchow-Klinikum
Augustenburger Platz 1
13353 Berlin, Germany
Tel.: +49 (30) 4593 – 2001
Email: alexander.meyer@dhzc-charite.de

Keywords

- Large language models
- Treatment decision-making
- Aortic stenosis
- Artificial intelligence
- Clinical practice guidelines
- Medical data processing

Abstract

Aims

Large language models (LLMs) have shown promise in therapeutic decision-making comparable to medical experts, but these studies have used specially prepared patient data. The aim of this study was to determine whether LLMs can make guideline-adherent treatment decisions based on real-world patient data.

Methods and Results

We conducted a retrospective study of 80 patients with severe aortic stenosis who were scheduled for either surgical (SAVR, n=24) or transcatheter aortic valve replacement (TAVR, n=56) by our institutional heart team in 2022. Various LLMs (BioGPT, GPT-3.5, GPT-4 and GPT-4 Turbo, Llama-2, Mistral, and PaLM-2) were queried using either deidentified original medical reports or manually generated case summaries to determine the most guideline-adherent treatment. Agreement with the Heart Team was measured using Cohen's kappa coefficients, reliability using intraclass correlation coefficients (ICCs), and fairness using frequency bias indices (FBIs) with FBIs >1 indicating bias towards TAVR. When presented with original medical reports, LLMs showed poor performance (kappa: -0.47-0.09, ICC: 0.0-0.91, FBI: 0.95-1.53). The LLMs' performance improved substantially when case summaries were used as input and additional guideline knowledge was added to the prompt (kappa: -0.02-0.62, ICC: 0.01-0.97, FBI: 0.46-1.24). Qualitative analysis revealed instances of hallucinations in all LLMs tested.

Conclusion

Our findings suggest that even advanced LLMs currently make informed treatment decisions only with extensively pre-processed data, not with original patient data. Unreliable responses, bias and hallucinations pose significant health risks and highlight the need for caution in applying LLMs to real-world clinical decision-making.

Introduction

Large Language Models (LLMs) have recently demonstrated their impressive capabilities in medicine, exemplified by passing medical board exams¹, making correct diagnoses in complex clinical cases² and excelling in physician-patient communication.³

Most recently, the use of LLMs in therapeutic decision making has been trialed. Several studies have shown that LLMs can make treatment decisions for patients with oncological and cardiovascular diseases that are in substantial agreement with the respective treatment decisions made by clinical experts in tumor boards⁴⁻⁷ and Heart Teams (HTs)⁸. However, a common feature of these studies was that the LLMs did not make treatment decisions based on real-world patient data in its original format (e.g., discharge letters, imaging reports, etc.) but based on pre-processed data.

In clinical practice, relevant patient data such as patient characteristics, comorbidities, tumor stages and imaging results are typically available in free-text format, either as medical text reports or as text entries in the electronic health record – a format that is likely to persist in the near future. In the studies, however, decision-relevant patient data were extracted from the original medical reports by the investigators in a pre-processing step before being provided to the LLMs as input in a concise and high-quality form. However, it is still unknown to what extent LLMs can make treatment decisions based on the original medical data, a scenario that could lead to a significant reduction in physician workload and potentially increase guideline adherence and thus improve patient care.

In this study, we investigate the impact of data representation, i.e. using original medical reports versus case summaries either manually written by physicians or generated by LLMs, on the performance of LLMs in therapeutic decision-making.

As our study population, we selected patients with severe aortic stenosis (AS). This cohort was chosen because the parameters relevant to decision-making are readily quantifiable, the potential for resource optimization is substantial, and the prevalence of the condition is increasing. If left untreated, AS is associated with high morbidity and mortality.⁹ Treatment modalities for severe AS include surgical aortic valve replacement (SAVR), transcatheter aortic valve replacement (TAVR) and, to a lesser extent, medical therapy. The choice of the optimal treatment modality depends on several clinical variables, including patient age, estimated surgical risk, comorbidities and anatomical factors, as specified in the 2021 guidelines of the European Society of Cardiology (ESC) and the European Association for Cardio-Thoracic Surgery (EACTS) guidelines for the management of valvular heart disease.¹⁰ The 2021 ESC/EACTS Guidelines strongly endorse an active, collaborative consultation with the multidisciplinary Heart Team (HT). HTs are comprised of cardiologists, cardiac surgeons,

cardiac imaging specialists and cardiac anesthesiologists. In HT meetings, these experts review a patient's condition based on patient data laboriously extracted from medical reports before arriving at a treatment decision using a guideline-based approach.

Methods

We presented patient data to an LLM to obtain a treatment decision of either SAVR or TAVR. We assessed the degree of concordance between the treatment decisions provided by the LLM and the treatment decisions provided by the HT. Furthermore, we assessed decidability, reliability and fairness of the LLMs. Finally, we compared the performance of seven state-of-the-art LLMs to the performance of a simple non-LLM reference model. In an ablative manner, we studied the effect of using case summaries instead of the original medical reports and adding guideline knowledge to the prompt, respectively, resulting in four distinct experiments (Figure 1). In addition, we conducted a fifth experiment to determine whether the best-performing LLM could make sound treatment decisions based on case summaries created by the LLM itself in an intermediate step.

Study population

This study included patients treated at a heart center. We screened all patients with severe degenerative AS who were scheduled for a HT meeting in our hospital information system at one campus of our center in 2022, we identified 80 patients with sufficiently digitized documentation. As part of a quaternary care center, our institutional HT receives preselected patients scheduled for invasive AS treatment. Therefore, the number of patients recommended for conservative treatment at our institution is negligible. As a result, we decided to limit the possible therapeutic options for this study to SAVR and TAVR, excluding conservative therapy. This study was approved by the research ethics committee of Charité – Universitätsmedizin Berlin (EA1/146/23).

Data collection

Medical reports were available as Portable Document Format (PDF) files in our hospital information system. For each patient, we included the following pre-procedural reports: the two most recent discharge letters (including letters from external clinics), the invasive coronary angiography report, the echocardiography report, the CT scan report, and the HT protocol. We manually anonymized these reports. HT meeting protocols are standardized documents that contain decision-relevant patient characteristics, such as comorbidities, surgical risk scores, as well as the final treatment decision of the HT (Supplementary material online, Figure S1). A detailed description of our institutional HT is provided in the Supplementary.

Large Language Models

The study employed several state-of-the-art LLMs, namely GPT-3.5¹¹, GPT-4¹² and GPT-4 Turbo by OpenAI, and PaLM 2 by Google¹³. In addition, we used the open-source models Mistral-7B¹⁴, Llama 2 by Meta¹⁵ and BioGPT¹⁶. These LLMs had either demonstrated proficiency on similar tasks or had undergone pre-training on medical literature. Model details are provided in Supplementary Table S1. We set model hyperparameters to the default values, except for the temperature τ (creativity), which we set to zero in line with previous studies in the medical domain. The model hyperparameters were

set to the default values, with the exception of the temperature, which was set to zero in accordance with previous studies in the medical domain.¹⁷ The temperature is a hyperparameter in LLMs that controls the randomness of the model's output. A low value results in a more deterministic and focused model output, thereby reducing variability and creativity. A detailed description of how we accessed the LLMs and handled input size constraints is given in the Supplementary material online (Section Supplementary methods).

Reference model

The reference model represented an algorithmic emulation of the 2021 ESC/EACTS Guidelines for the management of valvular heart disease.¹⁰ More specifically, the reference model assigned patients to either SAVR or TAVR according to a flowchart (Supplementary material online, Figure S2) and relevant clinical variables (Table S4, Table S5).¹⁰ Model details are provided in the Supplementary.

Experiments

Five experiments were conducted to investigate the effects of data pre-processing on LLM performance:

RAW: In the RAW experiment, anonymized medical reports (i.e., the two most recent discharge letters, the invasive coronary angiography report, the echocardiography report, and the CT scan report) were concatenated, and stored in a unified text file. This file was programmatically inserted into a prompt template. Each prompt included an introductory or continuation phrase and concluded with a request for a treatment decision (Supplementary material online, Table S6).

RAW+: As it is unknown whether the LLMs we used had sufficient guideline knowledge, we compiled a resume of relevant CPG content from the ESC/EACTS guidelines.¹⁰ We added this resume to the prompt along with the unified text reports.

SUM: To study the effect of content compression, we replaced the original medical reports used in RAW with concise case summaries. These case summaries were created based on patient data extracted from the HT protocols.

SUM+: Case summaries were used as input and enriched with the CPG resume (Figure 1).

SUM_{LLM+}: In a first query, case summaries were generated from the original medical reports by GPT-4 Turbo before these LLM-generated case summaries were used as input enriched with the CPG resume in a second query.

Prompt templates, the CPG resume and an exemplary case summary are shown in the Supplementary material online (Table S6).

The LLMs' responses were manually reviewed and categorized as either "TAVR", "SAVR" or "indeterminate". Indeterminate responses occur when the model output does not match the available answer choices or when the model determines that there is insufficient information to support a

decision (Supplementary material online, Table S7). To assess reliability and obtain robust estimates of performance metrics, the LLMs were presented with the same prompt input 10 times in succession for each experiment and patient (hereafter referred to as ‘runs’) to obtain a treatment decision. To prevent memory bias, a new chat session was initiated for each run.

Performance metrics

We quantified concordance by measuring accuracy and interrater agreement. Accuracy was calculated as the proportion of treatment decisions that agreed with those made by the HT. Interrater agreement was estimated using Cohen’s kappa coefficients. Decidability was quantified as the proportion of determinate treatment decisions. Bias was quantified using frequency bias indices (FBI), defined as the ratio of predicted to observed treatment decisions for TAVR.

Due to the limitations of individual metrics, we used three different metrics to quantify reliability: intraclass correlation coefficients (ICCs), normalized Shannon entropy, and the proportion of unanimously accurate treatment decisions. A detailed description of the performance metrics, including strategies for handling indeterminate responses, is provided in the Supplementary material online (Table S8).

Statistical analysis

Patient characteristics for patients who received SAVR vs. TAVR were compared using Student’s t-test for normally distributed continuous variables and Mann-Whitney U test for non-normally distributed continuous variables. The Shapiro-Wilk test was used to assess normality. The chi-squared test was used for binary variables and Fisher’s exact test for sparse binary data.

Accuracy and Cohen’s kappa were computed with Python’s `sklearn.metrics` package (version 1.2.2). ICCs were calculated based on a one-way random effects, absolute agreement, single-rater model¹⁸ using Python’s `pingouin` package (version 0.5.3).

Results

Patient characteristics

A total of 80 patients with severe AS who were discussed at our institutional HT in 2022 were included. Of these patients, 24 (30 %) underwent SAVR while 56 (70 %) underwent TAVR. Patient characteristics are presented in Table 1.

Qualitative analysis

The LLMs' output ranged from nonsensical treatment recommendations (e.g., heart transplant) and purely fabricated content, to correctly assessing the patient's status, choosing the correct treatment option, and supporting the treatment decision with additional anatomical insights (Table 2).

Qualitative analysis revealed that smaller models (e.g., BioGPT) tended to provide conflicting treatment recommendations for the same patient. In contrast, the frontier models (e.g., GPT-4, PaLM 2) consistently provided the same treatment recommendation when presented with the same patient data repeatedly over 10 runs.

In each experiment, all LLMs produced hallucinations of varying severity and frequency. These included instructional, contextual and factual inconsistencies (Table 2).

Quantitative analysis

Figure 2 and Table S9 (Supplementary material online) present the performance metrics. In the RAW experiment, LLMs achieved accuracies of up to 71 % and Cohen's kappa coefficients of up to 0.09. Some LLMs gave indeterminate treatment recommendations in up to 54 % of cases (e.g., GPT-3.5) and showed low reliability as evidenced by low ICCs and high entropy values (e.g., Mistral). FBIs were substantially higher than 1.0 for all LLMs except BioGPT. The reference model generally outperformed the LLMs in the RAW experiment regarding the metrics we assessed.

In the RAW+ experiments, performance metrics did not change substantially. However, the performance metrics improved in the SUM experiment and peaked in the SUM+ experiment, where some LLMs (e.g., GPT-4 Turbo) drew level with the reference model.

A general trend towards more concordant treatment decisions, fewer indeterminate responses, increased reliability, and less bias towards TAVR was observed with increasing data pre-processing and information enrichment efforts from RAW to SUM+ (Figure 2, Figure 3).

When instructing GPT-4 Turbo to generate case summaries as an intermediate step, we observed an accuracy of 74 %, a Cohen's kappa coefficient of 0.17, and an FBI of 1.2.

Discussion

To our knowledge, this is the first study to evaluate the impact of input data representation, including real-world medical data, on the ability of LLMs to make guideline-adherent treatment decisions.

Current LLMs make incorrect decisions based on original clinical data

Our analysis reveals that LLMs struggled to process original medical reports effectively, often outputting ‘TAVR’ or providing indeterminate responses. The LLMs showed low agreement with the HT, exhibited undecidability and unreliability, and displayed a strong bias towards TAVR. The high accuracies observed with some LLMs in the RAW experiment can be largely attributed to the class imbalance within our patient cohort, where 70 % of patients received TAVR.

LLMs require extensive data pre-processing to make sound therapeutic decisions

Performance improved substantially when physician-made case summaries were used as input, and when guideline knowledge was added to the prompts. GPT-4 and GPT-4 Turbo stood out as the most capable LLMs in our experiments. When given case summaries and a CPG resume, these two models showed substantial agreement with HT treatment decisions and drew level with the reference model in terms of accuracy, interrater agreement, decidability, and bias.

When GPT-4 Turbo was instructed to first generate a case summary before making a treatment decision, its performance was found to be intermediate between the results observed with original patient data and physician-generated case summaries.

Data representation affects LLM performance

An explanation for the underperformance of LLMs in the RAW experiment is not immediately apparent due to their opaque nature and the lack of established tools that allow the direct examination of input-output correlations. However, the underperformance cannot be attributed to a lack of or incorrectly applied guideline knowledge since the performance in RAW+ was similar to RAW and since LLMs can supposedly apply clinical knowledge to clinical cases as shown in their ability to pass medical board exams.^{1,19} This, together with the significant improvement in almost all performance measures observed when providing case summaries instead of original medical reports, suggests that the representation of input data is the most critical factor in LLM performance.

This finding is consistent with the fact that virtually all studies in which LLMs have been shown to make sound treatment decisions, used pre-processed clinical data as input.⁴⁻⁸ Of particular note is the study by Salihu et al.⁸ In this study, data from patients with severe AS were provided to GPT-4 to obtain a treatment decision for either TAVR, SAVR or conservative management. Patient data were provided in the form of a standardized multiple-choice questionnaire with 14 key clinical variables as

input, similar to our SUM experiments. GPT-4 treatment decisions were in substantial agreement with HT treatment decisions, a finding that we were able to reproduce in our experiments.

Similarly, in studies on tasks beyond therapeutic decision-making, such as answering board exam questions¹ and diagnosing complex clinical cases^{2,20}, LLMs performed particularly well when the input data were concise and information-dense.

Basic research has indicated that LLMs struggle with lengthy texts²¹ spanning over multiple prompts, potentially leading to memory loss²² and texts with a low signal-to-noise ratio.²³ A study by Levy et al²⁴ demonstrated that LLM reasoning performance declined notably with increasing input length. Specifically, the authors observed a 26 % drop in LLM performance when the input length was artificially increased from 250 to 3,000 tokens – i.e. a range of input lengths comparable to our study (see Supplementary material online Table S3).

Recently, Hager et al.²⁵ investigated the ability of LLMs to correctly diagnose patients presenting to the emergency department with abdominal pain. In this study, it was shown that deliberately withholding relevant clinical information from the LLMs paradoxically improved their diagnostic accuracy. Overall, this implies that LLMs are sensitive to both the signal-to-noise ratio and the sheer quantity of information provided.

LLMs may not be ideally suited for clinical decision-making

The results obtained with pre-processed patient data in our study and in previous studies demonstrate the potential of LLMs in medicine. However, the use of curated and pre-processed data does not reflect the clinical situation: To this day, the communication of clinical data within hospitals is largely based on unstructured free text.

Healthcare professionals have high expectations of AI to reduce their workload. This is not the case when physicians must manually extract and prepare key patient data for LLMs, as data extraction, not the actual decision-making task, is usually the most labor-intensive step. Interestingly, the observation that the performance with physician-generated summaries (i.e., SUM+) was substantially better than with summaries created by GPT-4 Turbo (i.e., SUM_{LLM+}), suggests that current LLMs are not yet adequately capable of pre-processing clinical data at the level of physicians. Once key patient data has been extracted and made available as input to a decision support model, the question arises as to why, of all machine learning (ML) models, LLMs should be used for therapeutic decision-making.

Selecting the best treatment modality for a patient is essentially a classification task. Some traditional ML models, such as tree-based models, are specifically developed for this purpose contrasting with LLMs, which are designed to generate text. Very simple reference models performed similarly to GPT-4 in our study and Salihu et al.'s⁸ study, suggesting that more sophisticated non-LLM models might generally be better decision-makers than LLMs if trained accordingly. In addition, non-

generative models do not exhibit undesirable behaviors such as hallucinations and unreliability^{26,27}, and provide explainability and established measures of uncertainty quantification, two hallmarks of reasonable AI²⁸ that are currently not adequately implemented for LLMs^{29,30}.

Another hallmark of reasonable AI is to address algorithmic bias. It is conceivable that the bias we observed in virtually all LLMs in our study could be due to LLMs being exposed to an abundance of TAVR-related internet literature during training compared to SAVR, subsequently influencing their treatment decisions.

Overall, we propose using LLMs to extract clinical data³¹ and generate input for downstream non-LLMs, which then perform the decision-making. While this strategy should ideally exploit the strengths of LLMs and well-established ML classifiers, its effectiveness remains to be proven in future studies.

Limitations

Our study is subject to certain limitations. These include a small patient cohort from a single center and the retrospective nature of our investigation. Nevertheless, the size of our study cohort (n=80), was comparable to previous key publications^{2,32} studying the performance of LLMs in medicine, and we assume that our patient cohort was sufficiently large given the clear trends we observed.

The intentionally vague instruction to make decisions according to “the guidelines” left it unclear which specific guidelines had to be followed. For instance, the 2021 ESC guidelines differ from the 2020 ACC/AHA guidelines³³. This ambiguity could contribute to the poor performance in the RAW experiment, as the 2021 ESC guidelines were used as benchmark. However, in the RAW+ experiments, the content of the 2021 ESC guidelines was included in the prompt, yet the performance did not improve substantially. Therefore, it must be assumed that this ambiguous instruction was not the driving factor for the poor performance in the RAW experiment.

The HT decisions against which we compared the LLMs’ treatment decisions could themselves be non-objective and deviate from the CPGs. We manually reviewed the HT treatment decisions and found no substantial deviations from the CPGs. Since human decision-makers (i.e. physicians) ultimately make the treatment decisions, the ground truth in experiments such as ours is inherently susceptible to some degree of subjectivity.

Given the limited cohort size and the considerable length of the medical reports, few-shot prompting or fine-tuning was not a viable option. We did not employ more sophisticated prompting techniques, such as Chain-of-Thought³⁴, and confined hyperparameter tuning to the temperature parameter.

Conclusions

Our experiments have been among the most challenging tasks LLMs have been asked to perform in the medical sciences. Overall, we conclude that LLMs are currently not suitable as decision makers

for the treatment of patients with severe AS, as our results suggest that a) LLMs require elaborate pre-processing of patient data to make informed treatment decisions, and b) LLMs are currently not able to pre-process original patient data on par with medical experts. Thus, we do not share the medical community's concern that staff will be replaced by artificial intelligence³⁵ in clinical decision-making in the near future.

Our findings suggest that LLMs should be used cautiously, particularly by medical laypersons seeking medical advice, such as second opinions. Users without extensive domain knowledge may receive treatment recommendations at a level similar to our RAW experiments. This is because medical laypersons may not be able to support prompts with guideline knowledge or create case summaries of sufficient quality but will only be able to use original medical reports. The study by Hager et al.²⁵ suggests that LLMs perform poorly when collecting additional patient data sequentially, as physicians would during a patient-physician dialogue. This suggests that the alternative to our approach - not providing all clinical data to the LLM at once, but having medical laypersons provide essential information incrementally during a chat session - is also likely to lead to suboptimal therapeutic recommendations.

Finally, medical laypersons may not be able to recognize hallucinations as effectively as medical professionals. This, combined with the eloquent and persuasive linguistic style of most LLMs, has the potential to mislead users by creating an illusion of greater certainty than warranted, aggravating the hazardous effects of incorrect treatment recommendations.

Figures and Tables

Figures

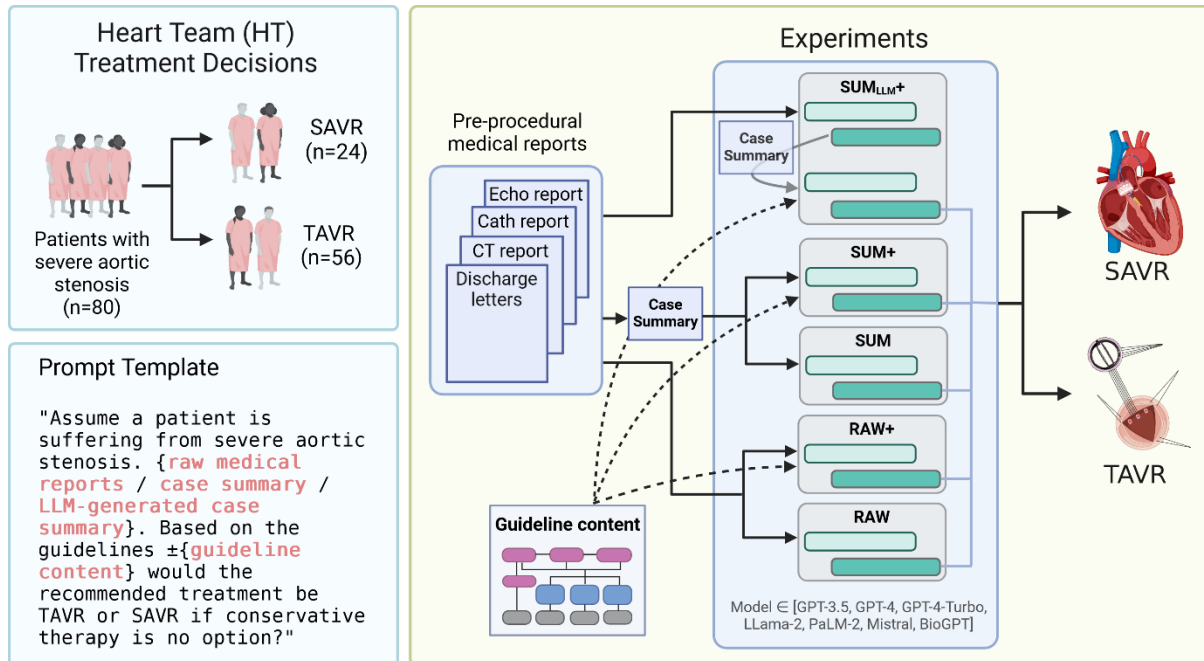


Figure 1: Experimental design

We presented the clinical data of 80 patients suffering from severe aortic stenosis (AS) to a Large Language Model (LLM) to receive a treatment decision for either surgical aortic valve replacement (SAVR) or transcatheter aortic valve replacement (TAVR). To investigate whether injecting guideline knowledge (RAW+) into the prompt and/or using case summaries (SUM, SUM+, SUM_{LLM+}) instead of the original medical reports (RAW) improves LLM performance, we conducted a total of five experiments. Case summaries included only decision-relevant patient data and were either manually created by physicians or created by an LLM (SUM_{LLM+}).

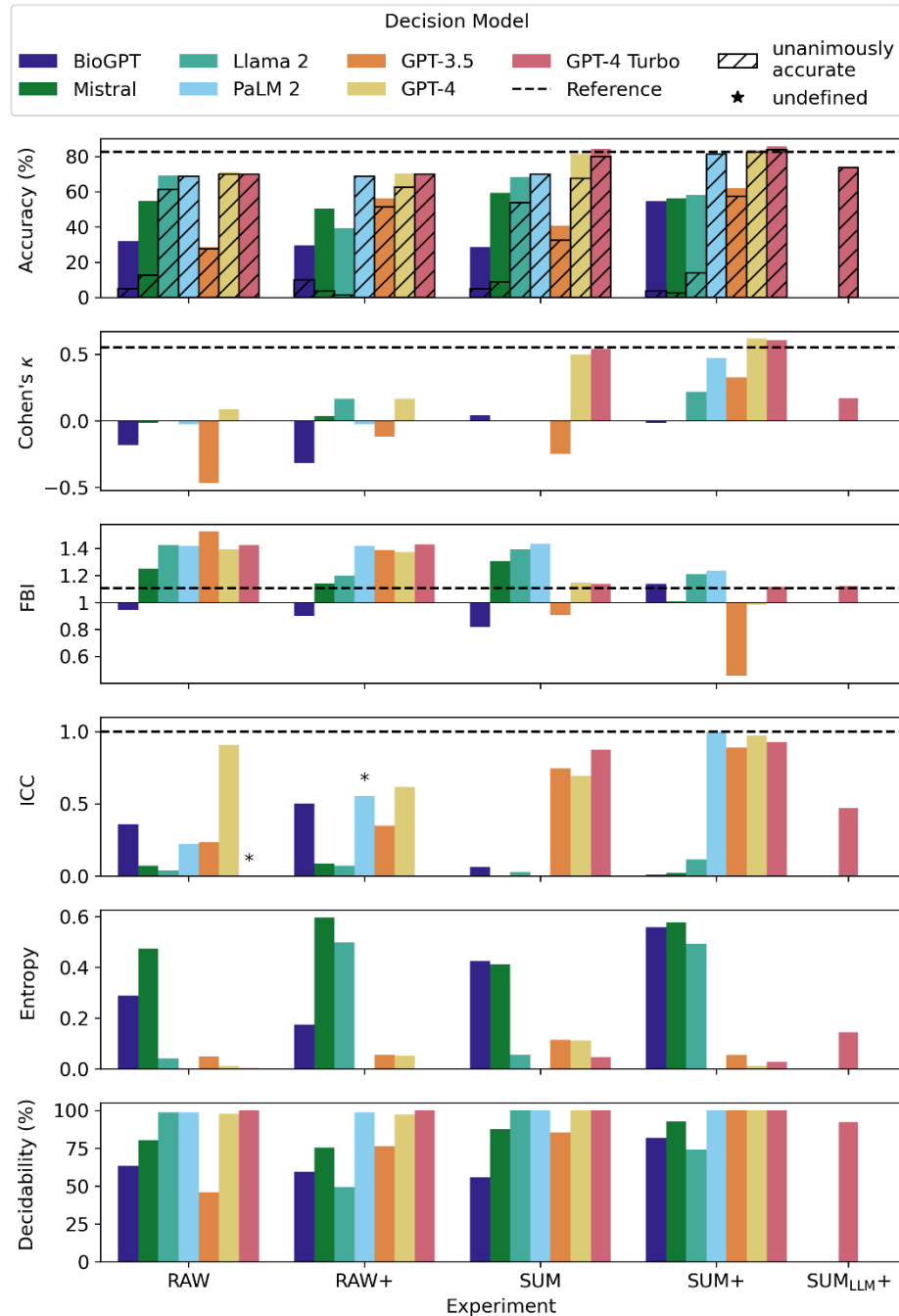


Figure 2: Performance metrics

Performance metrics of the Large Language Models (LLMs) are shown for the five experiments we conducted. The dashed line represents the reference model. Cohen's kappa coefficients ≤ 0 indicate no agreement, 0.01-0.20 slight, 0.21-0.40 fair, 0.41-0.60 moderate, 0.61-0.80 substantial, and 0.81-1.0 almost perfect agreement³⁶ with the Heart Team's treatment decisions. Frequency bias indices (FBIs) > 1 indicate bias towards TAVR and < 1 bias towards SAVR. Intraclass correlation coefficients (ICCs) < 0.5 indicate poor, 0.50-0.75 moderate, 0.75-0.90 good, > 0.90 excellent test-retest reliability.¹⁸ Instances where ICCs were undefined are marked by asterisks. Entropy values close to zero indicate low and entropy values close to one indicate high output variation, respectively. Decidability was defined as the proportion of non-indeterminate treatment decisions. The exact numerical values for the performance metrics are displayed in the Supplementary material online (Table S9). Abbreviations as in Figure 1.

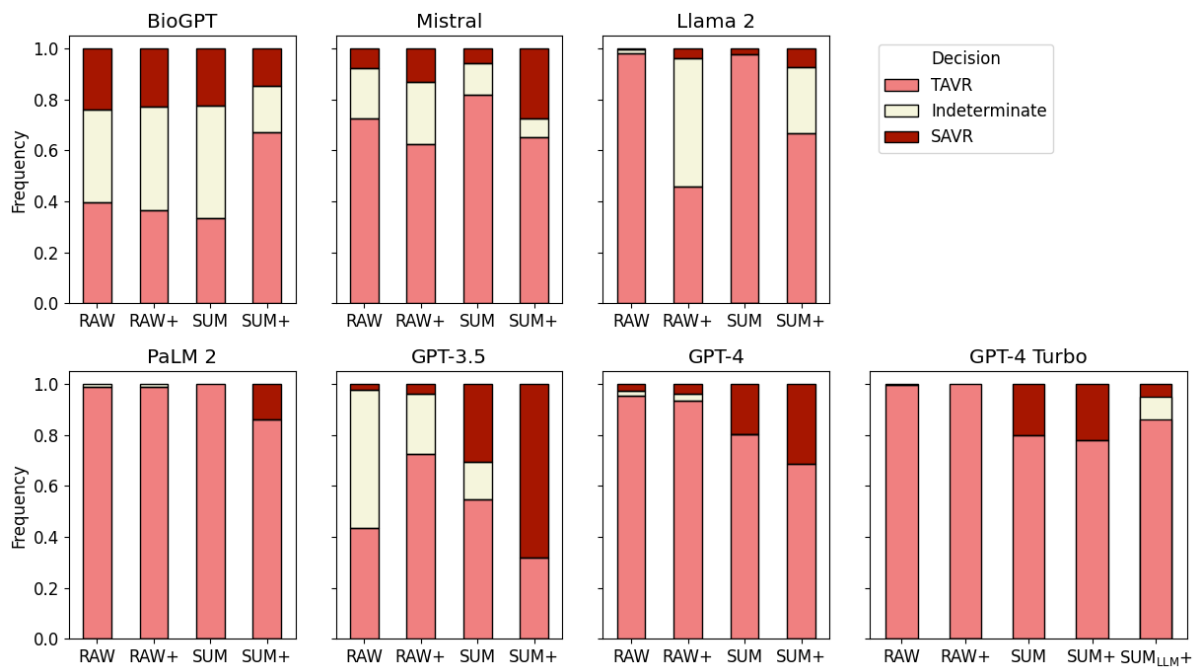


Figure 3: Frequencies of treatment decisions

Frequency distributions of the Large Language Models' treatment decisions are shown. A general trend towards increasing decidability and an increasing proportion of treatment decisions favoring SAVR could be observed between the RAW and the SUM+ experiment. Abbreviations as in Figure 2.

Tables

Variable	With Data (%)	Overall	SAVR (n = 24)	TAVR (n = 56)	P-Value
Age (years)	100.0	77.74 ± 7.5	70.71 ± 6.1	80.75 ± 5.8	<0.001
Female sex	100.0	36 (45.0)	8 (33.3)	28 (50.0)	0.26
Height (cm)	100.0	168.1 ± 11.0	172.5 ± 11.0	166.3 ± 10.6	<0.05
Body mass (kg)	100.0	76.3 ± 17.0	79.0 ± 16.0	75.1 ± 17.4	0.35
Body mass index (kg/m ²)	100	26.0 [23.0, 29.7]	25.9 [23.2, 29.0]	26.2 [23.0, 29.8]	0.66
Logistic EuroSCORE	31.2	6.8 [4.5, 13.0]	4.5 [2.2, 6.8]	8.4 [5.0, 16.0]	0.20
EuroSCORE II	98.8	2.6 [1.6, 4.5]	1.8 [1.1, 3.1]	2.9 [1.8, 4.9]	<0.05
STS score	76.2	2.8 [1.6, 4.5]	1.4 [1.1, 3.0]	3.3 [2.1, 4.5]	0.12
Left-ventricular ejection fraction (%)	100	60.0 [54.3, 61.3]	60.0 [48.8, 62.0]	60.0 [55.0, 60.0]	0.28
Aortic valve opening area (cm ²)	100	0.70 [0.60, 0.80]	0.80 [0.68, 0.80]	0.70 [0.60, 0.80]	0.18
Arterial hypertension	100	59 (73.8)	18 (75.0)	41 (73.2)	1.0
Diabetes mellitus	100	22 (27.5)	6 (25.0)	16 (28.6)	0.96
Hyperlipidemia	100	51 (63.7)	13 (54.2)	38 (67.9)	0.36
Previous cardiac surgery	100	1 (1.2)	0 (0.0)	1 (1.8)	1.0
Frailty	100	7 (8.8)	0 (0.0)	7 (12.5)	0.17
Sequelae of chest radiation	100	0 (0.0)	0 (0.0)	0 (0.0)	1.0
Porcelain aorta	100	0 (0.0)	0 (0.0)	0 (0.0)	1.0
Expected patient-prosthesis mismatch	100	1 (1.2)	0 (0.0)	1 (1.8)	1.0
Severe chest deformation or scoliosis	100	7 (8.8)	1 (4.2)	6 (10.7)	0.60
Severe coronary artery disease requiring revascularization	100	6 (7.5)	5 (20.8)	1 (1.8)	<0.05
Left-ventricular ejection fraction ≤ 40 %	100	6 (7.5)	3 (12.5)	3 (5.4)	0.52

Active neoplasia	100	7 (8.8)	2 (8.3)	5 (8.9)	1.0
Liver cirrhosis	100	1 (1.2)	0 (0.0)	1 (1.8)	1.0
Chronic obstructive pulmonary disease (GOLD stage ≥ 3)	100	5 (6.2)	1 (4.2)	4 (7.1)	1.0
Pulmonary arterial hypertension	100	8 (10.0)	3 (12.5)	5 (8.9)	0.94
Under immunosuppressive therapy	100	10 (12.5)	2 (8.3)	8 (14.3)	0.71

Table 1: Patient characteristics

Values are mean \pm standard deviation for continuous, normally distributed data, median and [interquartile range] for continuous, non-normally distributed data, and n (%) for binary data.

EuroSCORE: European system for cardiac operative risk evaluation, SAVR: Surgical Aortic Valve Replacement, STS: Society of Thoracic Surgeons, TAVR: Transcatheter Aortic Valve Replacement

Model	Experiment	Patient characteristics	LLM response	HT treatment decision	Interpretation
PaLM2	RAW	male, 55-59 years* old, EuroSCORE II: 0,55 %, no comorbidities except diffuse, mild coronary atherosclerosis and arterial hypertension, no relevant anatomical aspects to consider.	R: “Based on the provided information, the recommended treatment for the patient would be transcatheter aortic valve replacement (TAVR). The patient has severe aortic stenosis, which is causing symptoms and has been confirmed by echocardiography and cardiac catheterization. The patient is also at high risk for surgical aortic valve replacement (SAVR) due to their age, comorbidities, and the presence of diffuse coronary artery disease. TAVR is a less invasive procedure than SAVR and is associated with a lower risk of complications. Therefore, TAVR would be the preferred treatment option for this patient.”	SAVR	The LLM incorrectly recommends TAVR based on falsely extracted and/or interpreted patient characteristics (factual inconsistency). According to the literature, it has not yet been conclusively determined whether TAVR is associated with fewer complications than SAVR ³⁶ (factual inconsistency/bias).
BioGPT	RAW	male, 65-69 years* old, EuroSCORE- II: 7,2 %, post-cardiac arrest syndrome, frailty, long-term mechanical ventilation, liver cirrhosis, no relevant anatomical aspects to consider.	R: “Surgical replacement of the aortic valve. The patient was referred to the operating room for AVR. The operation was performed under general anesthesia and extracorporeal circulation ... The operation was performed with a median sternotomy. The ascending aorta was dissected free and opened...and the coronary arteries were visualized...”	TAVR	The LLM falsely states a surgery report from a procedure this patient never received and that was not mentioned in the prompt (context inconsistency).
BioGPT	RAW	female, 75-79 years * old, EuroSCORE- II: 2.4 %, STS score: 2.9%, COPD (GOLD Stage III), pulmonary hypertension, frailty, no relevant anatomical aspects to consider.	R: “Heart transplantation. - The patient's family decided against surgery. - The patient died due to heart failure.”	TAVR	The LLM did not follow the instruction to choose between SAVR and TAVR but opted for heart transplant (instruction inconsistency). The LLM fabricated a clinical course that was not mentioned in the prompt (context inconsistency).
GPT35	SUM	female, 70-74 years * old, EuroSCORE- II: 1.6 %, STS score: 1.1 %, no relevant comorbidities, no relevant anatomical aspects to consider.	R₁: “Based on the given information, the recommended treatment for a person with severe aortic valve stenosis and no other significant comorbidities would be SAVR (Surgical Aortic Valve Replacement) if conservative therapy is not an option.” R₂: “Based on the provided information, the recommended treatment for a person with severe aortic valve stenosis, who does not have any contraindications or	SAVR	The LLM gave opposite answers to the exact same question (unreliability). The LLM incorrectly recommends TAVR based on the absence of comorbidities

			<i>relevant comorbidities, would be TAVI (Transcatheter Aortic Valve Implantation) rather than SAVR (Surgical Aortic Valve Replacement)."</i>		(factual inconsistency).
PaLM2	RAW+	male, 55-59 years * old, EuroSCORE- II: 0.55 %, STS score: 0.7 %, no comorbidities except arterial hypertension, no relevant anatomical aspects to consider.	<i>R: "Based on the guidelines, the recommended treatment would be transcatheter aortic valve replacement (TAVR). The patient is 55-59 years old and at high surgical risk (STS-PROM score >8 % or EuroSCORE II > 8 %) due to severe frailty and previous cardiac surgery. TAVR is the preferred treatment for patients who are inoperable or at high surgical risk."</i>	SAVR	The LLM falsely claims that the patient is at high surgical risk due to frailty and previous cardiac surgery – which are both entirely fabricated (factual inconsistency).
GPT35	SUM+	female, 80-84 years * old, logistic EuroSCORE: 8.44 %, EuroSCORE II: 1.82 %, STS score: 4.33 %, stage 3A chronic kidney disease, no relevant anatomical aspects to consider.	<i>R: "Based on the guidelines provided, the recommended treatment for the patient would be SAVR (surgical aortic valve replacement) if conservative therapy is not an option. This is because the patient is 80-84 years old and falls into the older age category, and their logistic EuroSCORE is 8.44 %, which is above the threshold for TAVR (transcatheter aortic valve replacement) recommendation."</i>	TAVR	The LLM incorrectly recites the association between the logistic EuroSCORE and surgical risk prohibiting SAVR (context inconsistency). The LLM correctly assigns the patient to the older age category but then incorrectly recommends SAVR as the optimal therapy (logical inconsistency).
GPT4	RAW	female, 65-69 years * old, EuroSCORE- II: 2.5 %, STS score: 1.4 %, no relevant comorbidities, ascending aortic aneurysm (48 mm) mentioned in the CT scan report.	<i>R: "Based on the guidelines, the recommended treatment would be Surgical Aortic Valve Replacement (SAVR) as the patient has a severe aortic stenosis and an ascending aortic aneurysm. Transcatheter Aortic Valve Implantation (TAVI) is typically recommended for patients who are at high or prohibitive surgical risk, which does not seem to be the case here."</i>	SAVR	Correct treatment decision plus additional extraction of decision-relevant clinical data.
LLaMA2	SUM+	male, 65-69 years * old, EuroSCORE- II: 1.29 %, STS score: 3.04 %, COPD, no relevant anatomical aspects to consider.	<i>R₁: "Based on the guidelines provided, the recommended treatment for this patient would be TAVI. The patient's age (65-69 years old) and low surgical risk (EuroSCORE II 1.29 %, STS score 3.037 %) make them a suitable candidate for TAVI. Additionally, the patient does not have any contraindications for TAVI, such as severe frailty, previous cardiac surgery, or significant thorax deformity."</i>	SAVR	The LLM correctly identifies low surgical risk but then incorrectly rejects SAVR as therapeutic option (logical inconsistency). The LLM incorrectly recites factors favoring TAVR from the prompt (context inconsistency).

			<p>R₂: "Based on the information provided, the patient's EuroSCORE II is 1.29 % and STS score is 3.037 %, which indicates that they are at low surgical risk. The patient is also younger than 75 years old. Therefore, according to the guidelines, SAVR is the recommended treatment option for this patient."</p>		<p>Although the LLM correctly decided for SAVR this time, it gave the opposite response to the exact same question (unreliability).</p>
--	--	--	--	--	--

Table 2: Representative responses from the Large Language Models (LLMs)

The LLMs' treatment responses included well-informed decisions but also hallucinations ranging from obvious misinformation to absurd treatment recommendations and logical errors. We largely adhered to the taxonomy for the description of hallucinations established by Huang et al.²⁶ R:

Response of the LLM with subscripts indicating responses to the same question (obtained during ten runs). Green: Correct or useful, red: incorrect or harmful. Abbreviations as in Table 1. * Exact patient ages are replaced by age ranges for this table to ensure anonymization. In the experiments the exact patient age in years was used.

Acknowledgements

Dr. Roeschl and Dr. Hashemi are participants in the BIH Charité Digital Clinician Scientist Program funded by the Charité – Universitätsmedizin Berlin, and the Berlin Institute of Health at Charité (BIH).

We thank Michael Gudo (MORPHISTO GmbH) for providing access to GPT-4 and Hadi El Ali (B.Sc.), University of Bayreuth, for contributing to the illustration of Figure 1.

Funding

This work was supported by the German Centre for Cardiovascular Research (DZHK), funded by the German Federal Ministry of Education and Research, and the Charité – Universitätsmedizin Berlin.

D.H. received two grants from the DZHK (Grant Number: 81X3100214 and Grant Number: 81X3100220).

T.R. and D.H. are participants in the BIH Charité Digital Clinician Scientist Program funded by the Charité – Universitätsmedizin Berlin, and the Berlin Institute of Health at Charité (BIH).

Conflicts of Interests

Djawid Hashemi reports financial engagements beyond the scope of the presented work. These activities include consultation services and speaking engagements for companies including AstraZeneca, Bayer Vital, Boehringer Ingelheim, Coliquio, and Novartis.

Tobias D. Trippel reports on the potential conflict of interest by holding shares of Microsoft, Amazon and Palantir Technologies.

Axel Unbehaun serves as physician proctor to Boston Scientific, Edwards Lifesciences, and Medtronic.

Jörg Kempfert reports personal fees from Edwards, personal fees from LSI, outside the submitted work.

Benjamin O'Brien declares Research funding from the British Heart Foundation and the National Institute for Health Science Research and relevant financial activities outside the submitted work with following commercial entities: Teleflex, Abiomed in relation to consultancy fees.

Felix Balzer reports funding from Medtronic and grants from the German Federal Ministry of Education and Research, grants from the German Federal Ministry of Health, grants from the Berlin Institute of Health, personal fees from Elsevier Publishing, grants from Hans Böckler Foundation, other from Robert Koch Institute, grants from Einstein Foundation, and grants from Berlin University Alliance outside the submitted work.

Volkmar Falk declares relevant financial activities outside the submitted work with following commercial entities: Medtronic GmbH; Biotronik SE & Co.; Abbott GmbH & Co. KG; Boston Scientific; Edwards Lifesciences; Berlin Heart; Novartis Pharma GmbH; JOTEC GmbH; Zurich Heart. In relation to: Educational Grants (including travel support); Fees for lectures and speeches; Fees for professional consultation; Research and study funds.

Alexander Meyer declares the receipt of consulting and lecturing fees from Medtronic, lecturing fees from Bayer, and consulting fees from Pfizer. Alexander Meyer is founder and managing director of x-cardiac GmbH.

The remaining authors have no conflicts of interest to disclose.

Data sharing statement

The (anonymized) data underlying this article will be shared on reasonable request to the corresponding author.

Author Contributions

Conception and design of the study and literature review: TR, MH, DH, AM. Data collection: DH, FR. Analysis and interpretation of the data: MH, TR, AM, NH. Drafting of the manuscript: TR, MH, DH, AM, NH. All authors: revising and editing the manuscript.

References

1. Kung TH, Cheatham M, Medenilla A, et al. Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models. *PLOS Digital Health*. 2023;2(2):e0000198. doi:10.1371/journal.pdig.0000198
2. Kanjee Z, Crowe B, Rodman A. Accuracy of a Generative Artificial Intelligence Model in a Complex Diagnostic Challenge. *Jama*. Jul 3 2023;330(1):78-80. doi:10.1001/jama.2023.8288
3. Tu T, Palepu A, Schaekermann M, et al. Towards Conversational Diagnostic AI. *arXiv preprint arXiv:240105654*. 2024;
4. Sorin V, Klang E, Sklair-Levy M, et al. Large language model (ChatGPT) as a support tool for breast tumor board. *NPJ Breast Cancer*. May 30 2023;9(1):44. doi:10.1038/s41523-023-00557-8
5. Aghamaliyev U, Karimbayli J, Giessen-Jung C, et al. ChatGPT's Gastrointestinal Tumor Board Tango: A limping dance partner? *Eur J Cancer*. May 7 2024;205:114100. doi:10.1016/j.ejca.2024.114100
6. Kozel G, Gurses ME, Gecici NN, et al. Chat-GPT on brain tumors: An examination of Artificial Intelligence/Machine Learning's ability to provide diagnoses and treatment plans for example neuro-oncology cases. *Clin Neurol Neurosurg*. Apr 2024;239:108238. doi:10.1016/j.clineuro.2024.108238
7. Lukac S, Dayan D, Fink V, et al. Evaluating ChatGPT as an adjunct for the multidisciplinary tumor board decision-making in primary breast cancer cases. *Arch Gynecol Obstet*. Dec 2023;308(6):1831-1844. doi:10.1007/s00404-023-07130-5
8. Salihu A, Meier D, Noirclerc N, et al. A study of ChatGPT in facilitating Heart Team decisions on severe aortic stenosis. *EuroIntervention*. 2024;20(8):e496-e503. doi:10.4244/eij-d-23-00643
9. Roth GA, Mensah GA, Johnson CO, et al. Global Burden of Cardiovascular Diseases and Risk Factors, 1990–2019: Update From the GBD 2019 Study. *Journal of the American College of Cardiology*. 2020/12/22/ 2020;76(25):2982-3021. doi:10.1016/j.jacc.2020.11.010
10. Vahanian A, Beyersdorf F, Praz F, et al. 2021 ESC/EACTS Guidelines for the management of valvular heart disease: Developed by the Task Force for the management of valvular heart disease of the European Society of Cardiology (ESC) and the European Association for Cardio-Thoracic Surgery (EACTS). *European Heart Journal*. 2021;43(7):561-632. doi:10.1093/eurheartj/ehab395
11. Ye J, Chen X, Xu N, et al. A comprehensive capability analysis of gpt-3 and gpt-3.5 series models. *arXiv preprint arXiv:230310420*. 2023;
12. Achiam J, Adler S, Agarwal S, et al. Gpt-4 technical report. *arXiv preprint arXiv:230308774*. 2023;
13. Anil R, Dai AM, Firat O, et al. Palm 2 technical report. *arXiv preprint arXiv:230510403*. 2023;
14. Jiang AQ, Sablayrolles A, Mensch A, et al. Mistral 7B. *arXiv preprint arXiv:231006825*. 2023;
15. Touvron H, Martin L, Stone K, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:230709288*. 2023;
16. Luo R, Sun L, Xia Y, et al. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings in Bioinformatics*. 2022;23(6):bbac409.
17. Liévin V, Hother CE, Winther O. Can large language models reason about medical questions? *arXiv preprint arXiv:220708143*. 2022;
18. Koo TK, Li MY. A Guideline of Selecting and Reporting Intraclass Correlation Coefficients for Reliability Research. *J Chiropr Med*. Jun 2016;15(2):155-63. doi:10.1016/j.jcm.2016.02.012
19. Cai Y, Wang L, Wang Y, et al. MedBench: A Large-Scale Chinese Benchmark for Evaluating Medical Large Language Models. *arXiv preprint arXiv:231212806*. 2023;
20. Eriksen AV, Möller S, Ryg J. Use of GPT-4 to Diagnose Complex Clinical Cases. *NEJM AI*. 2024;1(1):Alp2300031. doi:doi:10.1056/Alp2300031
21. Liu NF, Lin K, Hewitt J, et al. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:230703172*. 2023;
22. Sejnowski TJ. Large Language Models and the Reverse Turing Test. *Neural Comput*. Feb 17 2023;35(3):309-342. doi:10.1162/neco_a_01563

23. Wang B, Wei C, Liu Z, Lin G, Chen NF. Resilience of Large Language Models for Noisy Instructions. *arXiv preprint arXiv:240409754*. 2024;
24. Levy M, Jacoby A, Goldberg Y. Same Task, More Tokens: the Impact of Input Length on the Reasoning Performance of Large Language Models. *arXiv preprint arXiv:240214848*. 2024;
25. Hager P, Jungmann F, Holland R, et al. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature Medicine*. 2024/07/04 2024;doi:10.1038/s41591-024-03097-1
26. Huang L, Yu W, Ma W, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:231105232*. 2023;
27. Singhal K, Azizi S, Tu T, et al. Large language models encode clinical knowledge. *Nature*. 2023/07/12 2023;doi:10.1038/s41586-023-06291-2
28. Sivarajah U, Wang Y, Olya H, Mathew S. Responsible Artificial Intelligence (AI) for Digital Health and Medical Analytics. *Inf Syst Front*. Jun 5 2023:1-6. doi:10.1007/s10796-023-10412-7
29. Luo H, Specia L. From Understanding to Utilization: A Survey on Explainability for Large Language Models. *arXiv preprint arXiv:240112874*. 2024;
30. Liu L, Pan Y, Li X, Chen G. Uncertainty Estimation and Quantification for LLMs: A Simple Supervised Approach. *arXiv preprint arXiv:240415993*. 2024;
31. Dagdelen J, Dunn A, Lee S, et al. Structured information extraction from scientific text with large language models. *Nature Communications*. 2024/02/15 2024;15(1):1418. doi:10.1038/s41467-024-45563-x
32. Huang AS, Hirabayashi K, Barna L, Parikh D, Pasquale LR. Assessment of a Large Language Model's Responses to Questions and Cases About Glaucoma and Retina Management. *JAMA Ophthalmology*. 2024;doi:10.1001/jamaophthalmol.2023.6917
33. Otto CM, Nishimura RA, Bonow RO, et al. 2020 ACC/AHA Guideline for the Management of Patients With Valvular Heart Disease: A Report of the American College of Cardiology/American Heart Association Joint Committee on Clinical Practice Guidelines. *Circulation*. 2021;143(5):e72-e227. doi:doi:10.1161/CIR.0000000000000923
34. Wei J, Wang X, Schuurmans D, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*. 2022;35:24824-24837.
35. Fogo AB, Kronbichler A, Bajema IM. AI's Threat to the Medical Profession. *JAMA*. 2024;doi:10.1001/jama.2024.0018
36. Virtanen MPO, Eskola M, Jalava MP, et al. Comparison of Outcomes After Transcatheter Aortic Valve Replacement vs Surgical Aortic Valve Replacement Among Patients With Aortic Stenosis at Low Operative Risk. *JAMA Network Open*. 2019;2(6):e195742-e195742. doi:10.1001/jamanetworkopen.2019.5742