# Comparative Analysis of Information Quality in Pediatric Otorhinolaryngology: Clinicians, Residents, and Large Language Models

Otolaryngology—
Head and Neck Surgery
2025, Vol. 173(1) 228–236
© 2025 The Author(s).
Otolaryngology—Head and Neck
Surgery published by Wiley
Periodicals LLC on behalf of
American Academy of
Otolaryngology—Head and Neck
Surgery Foundation.
DOI: 10.1002/ohn.1225
http://otojournal.org

Eleonora M. C. Trecca, MD, PhD<sup>1,2</sup> 

Note 

Note

#### **Abstract**

Objective. Pediatric otorhinolaryngology (ORL) addresses complex conditions in children, requiring a tailored approach for patients and families. With artificial intelligence (Al) gaining traction in medical applications, this study evaluates the quality of information provided by large language models (LLMs) in comparison to clinicians, identifying strengths and limitations in the field of pediatric ORL.

Study Design. Comparative blinded study.

Setting. Controlled research environment using LLMs.

Methods. Fifty-four items of increasing difficulty, namely 18 theoretical questions, 18 clinical scenarios, and 18 patient

questions, were posed to ChatGPT-3.5, -4.0, -40, Claude-3, Gemini, Perplexity, Copilot, a second-year resident, and an expert in the field of pediatric ORL. The Quality Analysis of Medical Artificial Intelligence (QAMAI) tool was used for blinded evaluation of the quality of medical information by a panel of expert members from the Young Otolaryngologists Group of the Italian Society of ORL and the International Federation of ORL Societies.

Results. LLMs performed comparably to specialist in theoretical and standardized clinical scenarios, with Bing Copilot achieving the highest QAMAI scores. However, AI responses lacked transparency in citing reliable sources and were less effective in addressing patient-centered questions. Poor

#### **Corresponding Author:**

Eleonora M. C. Trecca, MD, PhD, Department of Otorhinolaryngology and Maxillofacial Surgery, IRCCS Research Hospital Casa Sollievo della Sofferenza, Viale Cappuccini, I, 71013 San Giovanni Rotondo, Foggia, Italy. Email: eleonoramc.trecca@gmail.com; e.trecca@operapadrepio.it EleonoraTrecca@Twitter.com

<sup>&</sup>lt;sup>1</sup>GOS, Young Otolaryngologists Group of the Italian Society of Otorhinolaryngology—Head and Neck Surgery, Rome, Italy

<sup>&</sup>lt;sup>2</sup>Department of Otorhinolaryngology and Maxillofacial Surgery, IRCCS Research Hospital Casa Sollievo della Sofferenza, San Giovanni Rotondo, Foggia, Italy

<sup>&</sup>lt;sup>3</sup>Department of Clinical and Experimental Medicine, University of Foggia, Foggia, Italy

<sup>&</sup>lt;sup>4</sup>Department of Otolaryngology–Head and Neck Surgery, ASST Lariana, Ospedale Sant'Anna, University of Insubria, Como, Italy

<sup>&</sup>lt;sup>5</sup>Department of Biotechnology and LifeSciences, University of Insubria, Varese, Italy

<sup>&</sup>lt;sup>6</sup>Otorhinolaryngology–Head and Neck Surgery Unit, AORN "San Pio", Benevento, Italy

<sup>&</sup>lt;sup>7</sup>CEINGE-Advanced Biotechnology, Naples, Italy

<sup>&</sup>lt;sup>8</sup>Department of Otolaryngology and Head and Neck Surgery, Fondazione IRCCS Ca' Granda Ospedale Maggiore Policlinico of Milan, Milan, Italy

Department of Clinical Sciences and Community Health, University of Milan, Milan, Italy

<sup>&</sup>lt;sup>10</sup>Department of Otolaryngology-Head and Neck Surgery, Division of Laryngology and Bronchoesophagology, Baudour, Belgium

<sup>&</sup>lt;sup>11</sup>Department of Surgery, UMONS Research Institute for Health Sciences and Technology, University of Mons (UMons), Mons, Belgium

<sup>&</sup>lt;sup>12</sup>Department of Otolaryngology-Head Neck Surgery, Foch Hospital, School of Medicine, UFR Simone Veil, Université Versailles Saint-Quentinen-Yvelines (Paris Saclay University), Paris, France

<sup>&</sup>lt;sup>13</sup>Department of Otolaryngology, Kore University, Enna, Italy

<sup>&</sup>lt;sup>14</sup>Department of Otolaryngology-Head and Neck Surgery, IRCCS Azienda Ospedaliero-Universitaria of Bologna, Bologna, Italy

<sup>&</sup>lt;sup>15</sup>Alma Mater Studiorum - University of Bologna, Bologna, Italy

<sup>&</sup>lt;sup>16</sup>Otorhinolaryngology Unity, Head and Neck Department, Meyer Children's Hospital IRCCS, Florence, Italy

<sup>&</sup>lt;sup>17</sup>ENT Unit, Morgagni Pierantoni Hospital, Azienda USL della Romagna, Forlì, Forlì-Cesena, Italy

<sup>&</sup>lt;sup>18</sup>Department of Otolaryngology Head and Neck Surgery, Santa Maria delle Croci Hospital, AUSL della Romagna, Ravenna, Italy

Trecca et al. 229

interrater agreement among reviewers highlighted challenges in distinguishing human-generated from Al-generated responses. Rhinology topics received the highest scores, whereas laryngology and patient-centered questions showed lower agreement and performance.

Conclusion. LLMs show promise as supportive resources in pediatric ORL, particularly in theoretical learning and standardized cases. However, significant limitations remain, including source transparency and contextual communication in patient interactions. Human oversight is essential to mitigate risks. Future developments should focus on refining Al capabilities for evidence-based and empathetic communication to support both clinicians and families.

# **Keywords**

artificial intelligence, ChatGPT, digital health, eHealth, large language models, otolaryngology, pediatric otorhinolaryngology, pediatrics

Received December 18, 2024; accepted February 27, 2025.

Pediatric otorhinolaryngology (ORL) focuses on diagnosing and treating a wide range of diseases in pediatric patients. This field requires a tailored approach that addresses both the unique physiological characteristics of children and the concerns of their families. Unlike adult ORL, pediatric ORL often involves treating patients from birth through adolescence. This continuum of care presents challenges, as many conditions are congenital, anatomical spaces are small and variable, and the quality of life must be carefully considered in treatment decisions. In this context, parents play a crucial role as primary decision-makers, often having to choose the best treatment path for their children, which may include complex surgical interventions.<sup>1</sup>

In recent years, artificial intelligence (AI) has emerged as a transformative tool in the medical field, with applications ranging from clinical decision support to patient selfmanagement, real-world drug research, and assistance in research studies.<sup>2</sup> In ORL, AI has shown promise for automating classification tasks, analyzing clinical data, and simulating preoperative outcomes, which can aid physicians in providing precise, personalized care.<sup>3</sup> Among AI advancements, language-based models like ChatGPT (OpenAI, Microsoft) have gained popularity due to their ability to generate human-like responses. These models can make medical information more accessible to patients and clinicians; however, concerns persist regarding the accuracy and reliability of the medical advice they provide, as errors or misunderstandings in responses could pose risks to patient safety.4-6

Given these opportunities and challenges, this study seeks to evaluate the potential use of AI in pediatric ORL, specifically in comparison to medical practitioners, to assess the strengths and limitations of AI as a supportive tool in this complex subspecialty.

The aim of this study is to assess the quality of medical information provided by large language models (LLMs) in pediatric ORL. We intend to evaluate the accuracy and reliability of AI responses across three key areas: theoretical knowledge, clinical decision-making, and patient-centered advice. Comparisons will be made with responses from human practitioners of varying expertize, including residents and experienced clinicians. This study also aims to identify both the advantages and pitfalls of using AI in pediatric ORL clinical practice, with particular emphasis on aspects that impact physicians and families.

#### Materials and Methods

This study evaluated responses to a set of 54 questions, designed to represent a range of complexities in pediatric ORL (Supplemental Appendices 1-3, available online). These questions were divided into three categories: 18 theoretical questions, 18 clinical scenarios, and 18 patientcentered questions. The three types of questions were grouped and pertained to the fields of pediatric otology, rhinology, and laryngology for a total of 162 items. For this study, three experts—a pediatric otologist (M.R.), a rhinologist (A.M.d.L.), and a laryngologist (I.C.V.) generated fictitious results from hypothetical patients inspired by their clinical activities. For these reasons, patient safety was guaranteed, and formal ethical approval was waived. The items were designed to range from easy to more complex and were graded on a Likert scale (1-5) as easy (1-2), medium (3-4), and difficult (5). Agreement on the difficulty levels of the items was reached by the first and last authors (E.M.C.T., V.D.).

Responses were obtained from a range of AI models, including ChatGPT (OpenAI, Microsoft, versions 3.5, 4.0, and 40), Claude-3 (Anthropic), Gemini (Google), Perplexity AI (Inc.), and Copilot (Microsoft). For comparison, responses were also collected from a second-year resident (G.M.) and an experienced pediatric ORL specialist (M.T.-Z.). All answers were extracted from the AI tool during the week of September 23 to 29, 2024, whereas resident and specialist were given 1 month to answer the questions (September 2024).

A panel of experts (M.G., A.M., and J.R.L.) from the Young Otolaryngologists Group of the Italian Society of ORL and the Young Otolaryngologists of IFOS (International Federation of ORL Societies) evaluated each response, scoring them on the accuracy and relevance of the medical information provided. The panel, namely reviewer 1 (Rev1), reviewer 2 (Rev2), and reviewer 3 (Rev3), was selected based on their proven record of a high number of publications in the field of pediatric ORL.

The Quality Analysis of Medical AI (QAMAI)<sup>7</sup> tool was used for a blinded assessment of response quality; the

scale analyzes six parameters of medical information: accuracy, clarity, relevance, completeness, sources, and usefulness. Each parameter is evaluated using a Likert scale from 1 (strongly agree) to 5 (strongly disagree). This methodology allowed for an objective comparison between AI-generated responses and those provided by human practitioners at different levels of expertize, enabling us to assess the feasibility of using AI as a supportive tool in the clinical practice of pediatric ORL.

We also asked the panel (Rev1, Rev2, and Rev3) to indicate for each item whether the author of the response was a human or if the answer was AI-generated. This step was included to ensure a blinded evaluation and to assess the ability of the AI tool to appear human-like in its responses. Furthermore, this approach was implemented to mitigate any risk of bias, as the QAMAI is a score created to evaluate the quality of medical information generated by AI. In this article, however, it was also used to assess the quality of medical information provided by humans.

# Statistical Analysis

The statistical analysis aimed to evaluate the overall quality of responses provided by AI tools compared to clinicians in the field of pediatric ORL. QAMAI scores were summarized as total and mean scores, based also on subscores (ie, accuracy, clarity, relevance, etc.). Results were presented using descriptive statistics across categories (ie, responder type, question difficulty, and topics). Scores were averaged and tested for normality, using the Shapiro-Wilk test to guide the choice of parametric or nonparametric tests. Analysis of variance (ANOVA) or Kruskal-Wallis tests were employed to find differences in QAMAI scores among the question variables (ie, topic, difficulty, etc.) and responder type (ie, resident, Claude-3, others AI, etc.). Post hoc pairwise Wilcoxon test with Bonferroni correction was applied for multiple pairwise comparisons. The adjusted P-value is considered statistically significant if it is below .05. SPSS<sup>®</sup> Statistics applies the correction by maintaining the threshold at .05 and multiplying the retrieved P-value by the number of comparisons.

Interrater reliability was measured using Fleiss' kappa for multirater agreement across reviewers.

The IBM® SPSS® Statistics version 25 and R statistics were used to perform statistical analysis.

# **Results**

## Rater Agreement Among Reviewers

Multirater agreement among reviewers, evaluated using Fleiss' kappa, resulted in .05, indicating poor agreement. The three reviewers provided the same answers in only 35.8% of the cases (n = 174). Rev3 showed a stronger tendency to classify responses as AI-generated, whereas Rev1 and Rev2 exhibited a more balanced classification pattern. Specifically, Rev1 and Rev2 classified, respectively,

292 (60.1%) and 290 (59.7%) answers as AI, against 376 (77.4%) of Rev3. Agreement varied according to the difficulty of the questions, decreasing from 40.1% for easy questions to 36.4% for medium-difficulty questions and 30.9% for hard questions (P = .003). In topic-based analysis, the highest agreement was recorded in rhinology (50%), followed by otology (33.3%), and laryngeal topics (24.1%) (P = .001). When classifying patient-generated questions as either human- or AI-generated, agreement was the lowest (29.6%), whereas similar values were observed for clinical scenario and theoretical questions (37% and 40.7%, respectively) (P = .005). Data with specific subgroup topic, category, and difficulty are in **Table I** with complete agreement.

When assessing the accuracy of the reviewers in correctly identifying AI- versus human-generated responses, Rev3 achieved the highest accuracy, correctly classifying 84% of responses. This result is likely due to the previously described tendency to assign responses as AI-generated. Rev1 and Rev2 had similar classification accuracy rates, correctly identifying 58.4% and 59.3% of responses, respectively (P < .001).

# Overall Quality of Medical Information Provided by LLMs

Looking at the quality of the answers provided, the average QAMAI total score was 22.9 (SD = 3.1), with a median of 23 (range: 13.7-30). The average score across single subfields was consistent, with a mean of 4 (SD = 0.6) for usefulness, 4 (SD = 0.5) for completeness, 4.1 (SD = 0.6) for relevance and accuracy, and 4.2 (SD = 0.5) for clarity. The lowest score was recorded for the provided sources, with a mean value of 2.5 (SD = 1).

### Comparisons Based on Topic, Category, and Difficulty

Average total OAMAI score, SD, median, minimum, and maximum values combined for topic, category, and difficulty are collected in **Table 2**. The Kruskal-Wallis test found statistically significant differences among topics and categories (P = .007 and P = .02; difficulty resulted in a nonstatistically significant difference, P = .381). Based on actual responders, grouped as AI, specialist, and resident, no statistically significant difference was found from the Kruskal-Wallis test (P = .155) (mean<sub>AI</sub> = 23, SD = 2.9; mean<sub>resident</sub> = 22.1, SD = 4; mean<sub>specialist</sub> = 23.6, SD = 3.4). Two-way ANOVA models were built to investigate differences in mean total QAMAI scores between responder groups (AI, resident, and specialist) and, respectively, the topic, category, and difficulty. A single main effects analysis showed that there was a statistically significant higher score for the specialist only when compared to the resident (adj. P = .03), whereas no differences were found when comparing their answers to those provided by the AI. A single main effects analysis found similar results when considering the topic alone,

and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons

Trecca et al. 231

Table 1. Combined Group Based on Topic, Category, and Difficulty With Respective Full Agreement Among Reviewers

| Question variables                 | Cases of full agreement | Total | Cases of full agreement, % |
|------------------------------------|-------------------------|-------|----------------------------|
| Larynx: clinical scenario, easy    | 6                       | 18    | 33.3                       |
| Larynx: clinical scenario, hard    | 4                       | 18    | 22.2                       |
| Larynx: clinical scenario, medium  | 4                       | 18    | 22.2                       |
| Larynx: patient's question, easy   | 5                       | 18    | 27.8                       |
| Larynx: patient's question, hard   | 0                       | 18    | 0.0                        |
| Larynx: patient's question, medium | 4                       | 18    | 22.2                       |
| Larynx: theoretical, easy          | 4                       | 18    | 22.2                       |
| Larynx: theoretical, hard          | 7                       | 18    | 38.9                       |
| Larynx: theoretical, medium        | 5                       | 18    | 27.8                       |
| Oto: clinical scenario, easy       | 6                       | 18    | 33.3                       |
| Oto: clinical scenario, hard       | 6                       | 18    | 33.3                       |
| Oto: clinical scenario, medium     | 5                       | 18    | 27.8                       |
| Oto: patient's question, easy      | 5                       | 18    | 27.8                       |
| Oto: patient's question, hard      | 7                       | 18    | 38.9                       |
| Oto: patient's question, medium    | 5                       | 18    | 27.8                       |
| Oto: theoretical, easy             | 8                       | 18    | 44.4                       |
| Oto: theoretical, hard             | 5                       | 18    | 27.8                       |
| Oto: theoretical, medium           | 7                       | 18    | 38.9                       |
| Rhino: clinical scenario, easy     | 9                       | 18    | 50.0                       |
| Rhino: clinical scenario, hard     | 9                       | 18    | 50.0                       |
| Rhino: clinical scenario, medium   | H                       | 18    | 61.1                       |
| Rhino: patient's question, easy    | 9                       | 18    | 50.0                       |
| Rhino: patient's question, hard    | 3                       | 18    | 16.7                       |
| Rhino: patient's question, medium  | 10                      | 18    | 55.6                       |
| Rhino: theoretical, easy           | 13                      | 18    | 72.2                       |
| Rhino: theoretical, hard           | 9                       | 18    | 50.0                       |
| Rhino: theoretical, medium         | 8                       | 18    | 44.4                       |

with higher scores achieved in the rhinology group only when compared to otology (adj. P = .001). The interaction between responder type and topic was not statistically significant. Two-way ANOVA was also conducted to examine the effect of responders and category with respect to the total mean QAMAI score. There was a statistically significant interaction between responders and question category (P = .02). The specialist achieved higher scores in theoretical scenarios compared to the resident answering clinical and theoretical questions (respectively, adj. P = .02, and = .01) and AI answering patients' questions (adj. P = .04). When difficulty was considered in the model, its simple main effects and interaction with the responder group were not statistically significant different (P = .48 and = .99, respectively).

# LLMs Comparisons

When comparing different AI models, a statistically significant difference was found (P = .023), with Bing Copilot reporting the highest mean score (24.1, SD = 3). This score was significantly higher than that of Google Gemini (adj. P = .032) and ChatGPT-40 (adj. P = .020). Notably, Bing Copilot also achieved a significantly

higher mean total QAMAI score compared to residents (adj. P = .026), whereas no significant difference was found between specialists and any of the AI models (**Table 3**).

# Analysis of QAMAI Subfields

In the last instance, single QAMAI subfields were analyzed to determine which variables had the greatest impact on the overall total score. Means and SDs are depicted in **Table 4**. Overall, most subfields received high score, while the "resources provided" subfield had the lowest scores, particularly in the rhinology topic. The Kruskal-Wallis test showed no statistically significant differences in the mean scores of QAMAI subfields based on the difficulty of the questions or category, except for the "source provided" subfield (P = .005). However, all subfields showed statistically significant differences based on topics (P < .001). The highest scores were awarded in the rhinology topic, whereas otology and larynx topics had similar mean scores. When analyzing the responders' QAMAI subfield scores, aspects such as accuracy, clarity, relevance, and usefulness did not show significant differences. However, variable scores were observed for

10976817, 2025, 1, Downloaded from https://aao-hnsfjournals.onlinelibrary.wiley.com/doi/10.1002/ohn.1225 by Universite De Mons (Umons), Wiley Online Library on [06/10/2025]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms

**Table 2.** Average Total Quality Analysis of Medical Artificial Intelligence (QAMAI) Score, SD, Median, Minimum (min), and Maximum (max) Recorded Values for Topic, Category, and Difficulty

| Торіс  | Category           | Difficulty | Mean for total QAMAI score | SD   | Median | Min  | Max  |
|--------|--------------------|------------|----------------------------|------|--------|------|------|
| Larynx | Clinical scenario  | Easy       | 23.3                       | 3.24 | 22.3   | 18.7 | 29.7 |
| Larynx | Clinical scenario  | Hard       | 22.9                       | 2.73 | 23.3   | 19.3 | 28   |
| Larynx | Clinical scenario  | Medium     | 23.5                       | 2.86 | 23.5   | 18.3 | 28.3 |
| Larynx | Patient's question | Easy       | 22.3                       | 3.41 | 22.3   | 15.7 | 27.3 |
| Larynx | Patient's question | Hard       | 20.7                       | 2.35 | 20     | 17.7 | 25.3 |
| Larynx | Patient's question | Medium     | 20.4                       | 3.13 | 19.7   | 15.7 | 26   |
| Larynx | Theoretical        | Easy       | 24.2                       | 2.58 | 24.2   | 20.7 | 29   |
| Larynx | Theoretical        | Hard       | 23.9                       | 2.95 | 23     | 17.7 | 28.3 |
| Larynx | Theoretical        | Medium     | 24.2                       | 2.52 | 23.8   | 20   | 30   |
| Oto    | Clinical scenario  | Easy       | 22.6                       | 3.63 | 22.3   | 13.7 | 27.3 |
| Oto    | Clinical scenario  | Hard       | 22.4                       | 3.6  | 22.3   | 13.7 | 27   |
| Oto    | Clinical scenario  | Medium     | 21.8                       | 3.85 | 22     | 15   | 28   |
| Oto    | Patient's question | Easy       | 21.9                       | 3.88 | 22.2   | 13.7 | 28   |
| Oto    | Patient's question | Hard       | 22.3                       | 3.45 | 22.8   | 16   | 27   |
| Oto    | Patient's question | Medium     | 22.5                       | 3.7  | 22.7   | 13.7 | 27.3 |
| Oto    | Theoretical        | Easy       | 22.6                       | 4.19 | 22.8   | 13.7 | 29.7 |
| Oto    | Theoretical        | Hard       | 22.8                       | 3.63 | 23.3   | 14.3 | 27   |
| Oto    | Theoretical        | Medium     | 22.7                       | 3.87 | 23     | 15   | 28.3 |
| Rhino  | Clinical scenario  | Easy       | 22.7                       | 1.76 | 22.8   | 17.7 | 25.3 |
| Rhino  | Clinical scenario  | Hard       | 23.2                       | 1.56 | 23.5   | 19.7 | 25.3 |
| Rhino  | Clinical scenario  | Medium     | 22.6                       | 2.21 | 22.7   | 16.7 | 25.7 |
| Rhino  | Patient's question | Easy       | 24.5                       | 2.47 | 24.8   | 18.7 | 27.3 |
| Rhino  | Patient's question | Hard       | 23.7                       | 2.17 | 23.5   | 20.3 | 28   |
| Rhino  | Patient's question | Medium     | 24.8                       | 1.95 | 24.7   | 21   | 27.7 |
| Rhino  | Theoretical        | Easy       | 23.9                       | 3.07 | 24     | 17   | 30   |
| Rhino  | Theoretical        | ,<br>Hard  | 22.4                       | 2.07 | 22.7   | 18.3 | 25.3 |
| Rhino  | Theoretical        | Medium     | 24.6                       | 1.88 | 24.2   | 21.7 | 28   |

**Table 3.** Total Mean Quality Analysis of Medical Artificial Intelligence Scores Specific to Resident, Specialist, and Single Artificial Intelligence (AI) Models

| Real responder | Mean | Std. deviation | N   |
|----------------|------|----------------|-----|
| Resident       | 22.1 | 4.0            | 54  |
| Specialist     | 23.6 | 3.4            | 54  |
| ChatGPT-3.5    | 23.2 | 2.6            | 54  |
| ChatGPT-4.0    | 23.1 | 2.7            | 54  |
| Perplexity Al  | 22.9 | 2.9            | 54  |
| Google Gemini  | 22.3 | 2.7            | 54  |
| Bing Copilot   | 24.1 | 3.0            | 54  |
| Claude-3       | 22.9 | 2.7            | 54  |
| ChatGPT-40     | 22.3 | 3.5            | 54  |
| Total          | 22.9 | 3.1            | 486 |

the "sources provided" and completeness (P < .001 and P = .048, respectively). Specifically, the expert received higher scores compared to the resident, whereas no statistically significant differences were found among the other comparisons.

# **Discussion**

To the best of our knowledge, this is the first study to evaluate the quality of information provided by seven LLMs in the field of pediatric ORL using blinded evaluation with human expertize. By evaluating responses from various LLMs, alongside those of a resident and an experienced specialist, we were able to assess the potential of AI as a supportive tool in clinical practice.

Our analysis revealed that AI models can achieve comparable quality scores to human respondents, particularly specialists, in theoretical and clinical scenarios. This indicates that AI tools have reached a level of sophistication where they can contribute meaningfully to decision-making in pediatric ORL. Among the AI models, Bing Copilot emerged as the top performer, underscoring the variability in performance across different platforms. This is consistent with the current literature with Bing Copilot demonstrating superior performance across several key metrics, including readability and comprehensiveness, and superior accuracy compared to other LLMs in various areas of research. In a related study, interestingly, Copilot demonstrated remarkable performance, achieving the second-highest

Table 4. Mean Quality Analysis of Medical Artificial Intelligence (QAMAI) Subfield Scores by Group, Topic, and Category

|          |                    |            | OAMAI total      | Accuracy        | Clarity         | Relevance       | Completeness     | Source          | Usefulness      |
|----------|--------------------|------------|------------------|-----------------|-----------------|-----------------|------------------|-----------------|-----------------|
| Topic    | Category           | Responder  | (mean±SD)        | (mean ± SD)     | (mean±SD)       | (mean±SD)       | ,<br>(mean ± SD) | (mean±SD)       | (mean±SD)       |
| Larynx   | Clinical scenario  | ₹          | 23.26 ± 2.77     | $3.97 \pm 0.51$ | 4.17±0.45       | 3.98 ± 0.49     | 3.99 ± 0.45      | 3.18 ± 0.62     | 3.96 ± 0.50     |
| Larynx   | Clinical scenario  | Resident   | $22.33 \pm 2.80$ | $3.89 \pm 0.45$ | 4.11±0.40       | $3.72 \pm 0.64$ | $3.72 \pm 0.39$  | $3.05 \pm 0.61$ | $3.83 \pm 0.46$ |
| Larynx   | Clinical scenario  | Specialist | $24.05 \pm 4.14$ | $4.17 \pm 0.70$ | $4.17 \pm 0.55$ | $4.17 \pm 0.7$  | $4.17 \pm 0.69$  | $3.22 \pm 0.93$ | $4.17 \pm 0.69$ |
| Larynx   | Patient's question | ¥          | $20.96 \pm 2.68$ | $3.59 \pm 0.45$ | $3.87 \pm 0.40$ | $3.61 \pm 0.44$ | $3.59 \pm 0.42$  | $2.71 \pm 0.66$ | $3.59 \pm 0.44$ |
| Larynx   | Patient's question | Resident   | $21.83 \pm 4.42$ | $3.89 \pm 0.69$ | $3.89 \pm 0.72$ | $3.89 \pm 0.69$ | 3.61 ± 0.77      | $2.72 \pm 0.83$ | $3.83 \pm 0.78$ |
| Larynx   | Patient's question | Specialist | $21.39 \pm 4.47$ | $3.72 \pm 0.71$ | $3.72 \pm 0.74$ | $3.78 \pm 0.65$ | $3.72 \pm 0.71$  | $2.72 \pm 1.00$ | $3.72 \pm 0.71$ |
| Larynx   | Theoretical        | ¥          | $24.20 \pm 2.55$ | $4.17 \pm 0.44$ | $4.28 \pm 0.43$ | $4.23 \pm 0.45$ | $4.17 \pm 0.49$  | $3.23 \pm 0.64$ | 4.11 ± 0.46     |
| Larynx   | Theoretical        | Resident   | $22.44 \pm 2.75$ | $3.83 \pm 0.40$ | $4.22 \pm 0.50$ | $3.94 \pm 0.38$ | $3.77 \pm 0.45$  | $2.88 \pm 0.62$ | $3.78 \pm 0.45$ |
| Larynx   | Theoretical        | Specialist | $25.17 \pm 2.89$ | $4.39 \pm 0.49$ | $4.33 \pm 0.52$ | $4.38 \pm 0.49$ | $4.28 \pm 0.44$  | $3.61 \pm 0.83$ | $4.17 \pm 0.55$ |
| Oto<br>O | Clinical scenario  | ₹          | $22.53 \pm 3.38$ | $3.96 \pm 0.57$ | $4.09 \pm 0.53$ | $3.94 \pm 0.57$ | $3.95 \pm 0.52$  | $2.62 \pm 0.84$ | $3.97 \pm 0.58$ |
| Oto<br>O | Clinical scenario  | Resident   | $20.67 \pm 5.78$ | $3.55 \pm 1.00$ | $3.94 \pm 0.77$ | $3.50 \pm 0.98$ | $3.66 \pm 0.84$  | $2.50 \pm 1.22$ | 3.50 ± 1.00     |
| Oto      | Clinical scenario  | Specialist | $21.89 \pm 3.09$ | $3.89 \pm 0.54$ | $4.05 \pm 0.39$ | $3.89 \pm 0.54$ | $3.94 \pm 0.53$  | $2.22 \pm 0.62$ | $3.89 \pm 0.54$ |
| ot<br>O  | Patient's question | ₹          | $22.15 \pm 3.52$ | $3.86 \pm 0.61$ | $4.10 \pm 0.53$ | $3.87 \pm 0.62$ | $3.88 \pm 0.54$  | $2.55 \pm 0.82$ | $3.88 \pm 0.62$ |
| Oto<br>O | Patient's question | Resident   | $23.67 \pm 5.08$ | $4.05 \pm 0.88$ | $4.28 \pm 0.64$ | $4.00 \pm 0.90$ | $4.17 \pm 0.75$  | $3.00 \pm 1.05$ | $4.17 \pm 0.91$ |
| oto<br>O | Patient's question | Specialist | $21.17 \pm 2.70$ | $3.72 \pm 0.44$ | $4.00 \pm 0.47$ | $3.72 \pm 0.44$ | $3.78 \pm 0.50$  | $2.17 \pm 0.41$ | $3.78 \pm 0.50$ |
| of<br>O  | Theoretical        | ₹          | $22.57 \pm 3.43$ | $3.91 \pm 0.57$ | $4.11 \pm 0.52$ | $3.92 \pm 0.58$ | $3.96 \pm 0.54$  | $2.70 \pm 0.82$ | $3.96 \pm 0.57$ |
| of<br>O  | Theoretical        | Resident   | $21.89 \pm 5.80$ | $3.89 \pm 1.00$ | $4.05 \pm 0.74$ | 3.89 ± 1.00     | $3.67 \pm 0.81$  | $2.61 \pm 1.27$ | $3.78 \pm 1.02$ |
| ot<br>O  | Theoretical        | Specialist | $24.44 \pm 4.55$ | $4.28 \pm 0.77$ | $4.39 \pm 0.53$ | $4.28 \pm 0.77$ | $4.28 \pm 0.68$  | $2.94 \pm 1.18$ | $4.28 \pm 0.77$ |
| Rhino    | Clinical scenario  | ₹          | 22.81 ± 1.72     | $4.20 \pm 0.35$ | $4.32 \pm 0.37$ | $4.48 \pm 0.30$ | $4.17 \pm 0.40$  | $1.52 \pm 0.64$ | $4.13 \pm 0.37$ |
| Rhino    | Clinical scenario  | Resident   | $21.39 \pm 2.06$ | $4.11 \pm 0.58$ | $4.22 \pm 0.40$ | $4.17 \pm 0.28$ | $3.83 \pm 0.55$  | $1.11 \pm 0.27$ | $3.95 \pm 0.39$ |
| Rhino    | Clinical scenario  | Specialist | $24.39 \pm 1.60$ | $4.61 \pm 0.49$ | $4.67 \pm 0.21$ | $4.72 \pm 0.25$ | $4.50 \pm 0.28$  | $1.33 \pm 0.52$ | $4.55 \pm 0.34$ |
| Rhino    | Patient's question | ₹          | $24.27 \pm 2.14$ | $4.52 \pm 0.40$ | $4.50 \pm 0.35$ | $4.67 \pm 0.34$ | $4.40 \pm 0.40$  | $1.73 \pm 0.90$ | $4.43 \pm 0.40$ |
| Rhino    | Patient's question | Resident   | $24.67 \pm 3.04$ | $4.55 \pm 0.40$ | $4.67 \pm 0.42$ | $4.72 \pm 0.33$ | $4.50 \pm 0.59$  | 1.78 ± 0.96     | $4.44 \pm 0.54$ |
| Rhino    | Patient's question | Specialist | $24.67 \pm 2.18$ | $4.72 \pm 0.32$ | $4.61 \pm 0.53$ | $4.83 \pm 0.28$ | $4.50 \pm 0.50$  | $1.50 \pm 0.66$ | $4.50 \pm 0.50$ |
| Rhino    | Theoretical        | ₹          | $23.91 \pm 2.04$ | $4.38 \pm 0.35$ | $4.45 \pm 0.29$ | $4.40 \pm 0.29$ | $4.21 \pm 0.38$  | $2.17 \pm 0.97$ | $4.30 \pm 0.33$ |
| Rhino    | Theoretical        | Resident   | $19.83 \pm 2.68$ | $3.83 \pm 0.55$ | $3.83 \pm 0.62$ | $3.89 \pm 0.45$ | $3.39 \pm 0.49$  | $1.39 \pm 0.13$ | $3.50 \pm 0.62$ |
| Rhino    | Theoretical        | Specialist | $25.28 \pm 2.27$ | $4.67 \pm 0.42$ | $4.72 \pm 0.32$ | $4.67 \pm 0.30$ | $4.50 \pm 0.41$  | $2.11 \pm 0.93$ | $4.61 \pm 0.39$ |
|          |                    |            |                  |                 |                 |                 |                  |                 |                 |

Abbreviation: Al, artificial intelligence.

score among 108 otolaryngologists who participated in a public health care system ORL job competition examination. <sup>10</sup>

However, AI models were less effective when addressing patient-centered questions, suggesting limitations in contextual understanding and empathy. In fact, although LLMs are capable of identifying and addressing emotions, they lack the ability to perform deep reflective analysis of emotional experiences and the motivational aspects of emotions. This limitation is particularly evident in the use of generic empathic phrases that lack genuine contextual adaptation and a tendency toward overly verbose or unnecessarily elaborate replies. These findings underscore the need for significant improvements in AI performance to generate nuanced and contextually appropriate empathic communication. 11,12 Furthermore, they highlight the necessity for more robust and standardized evaluation strategies to effectively measure such soft skills in complex clinical contexts, such as those involving pediatric patients and their families, who often pose a wide range of emotionally charged and intricate questions.

Our results evidenced that AI models demonstrated strong performance in clarity, relevance, and usefulness of their responses. However, they consistently underperformed in providing transparent and reliable sources, which is a key limitation for integrating AI into evidence-based clinical practice. As highlighted by other authors, <sup>13,14</sup> ChatGPT has a tendency to generate references that may be accurate, erroneous, or entirely inexistent, depending on their alignment with existing records in scientific databases. This shortcoming highlights the need for improving AI models' ability to cite reliable medical literature and enhance the trustworthiness of their outputs to avoid risks of fabrication and plagiarism. <sup>15</sup>

Reviewer agreement decreased as question difficulty increased, with laryngology and patient-centered questions showing the lowest agreement rates. Interestingly, AI tools and specialist achieved higher scores in rhinology-related questions, suggesting that this topic might be more suited for AI-driven analysis due to a relatively higher degree of standardization in the field. <sup>16,17</sup>

The potential for AI in pediatric ORL is significant, especially for supporting clinicians in theoretical learning and standardized clinical scenarios. However, its current limitations—particularly in providing source transparency and effectively addressing nuanced patient-centered queries—must be addressed before widespread adoption. Additionally, the integration of AI into clinical workflows should be approached cautiously, with human oversight to mitigate risks associated with erroneous or incomplete AI-generated advice.

A major strength of this study is that, to the best of our knowledge, it evaluates the highest number of LLMs to date in this context. This broad inclusion allows for a comprehensive understanding of AI performance across different models, providing valuable insights into their potential role in pediatric ORL practice. Although the

results provide meaningful insights, several limitations of the study should be acknowledged. First, the clinical scenarios were generated based on fictional patient cases, which, although inspired by real practice, may not fully capture the complexity and variability of actual clinical encounters. Also, responses were evaluated in a controlled setting rather than in real-world clinical environments where the dynamic nature of interactions with patients, caregivers, and health care providers could influence outcomes. Second, the study focused on a specific set of AI tools available at a particular time. Current advancements in AI may yield different results, limiting the generalizability of our findings to future models. Finally, although reviewers were blinded, inherent biases may have influenced their judgments. Notably, the QAMAI is a validated scoring system specifically designed to evaluate the quality of medical information generated by AI. Extending its use to assess responses provided by specialist and resident may introduce potential biases. However, reviewers often encountered objective difficulties in distinguishing whether responses were human- or AI-generated, a challenge particularly evident in Rev3's tendency to classify responses as AI-generated. Therefore, these observations support the suitability of QAMAI for research purposes, even when applied to the evaluation of human responses.

Our study underscores the importance of continued development and training of AI systems. Enhancing AI models' capabilities in sourcing reliable evidence and improving their understanding of patient-centered contexts could bridge existing gaps. Future research should focus on the longitudinal assessment of AI implementation in clinical practice to evaluate its impact on patient outcomes, clinician efficiency, and educational training.<sup>18</sup> In fact, especially following the COVID-19 pandemic, ORL resident education has increasingly incorporated asynchronous learning methods, including web-based platforms.<sup>19</sup> Although traditional textbook- and didactic-based education remains the gold standard, there is a growing need to consider a paradigm shift within academic ORL education. This shift should include the integration of high-quality electronic resources, where, after thorough evaluation and refinement, LLMs could play a central role.<sup>20</sup>

### **Conclusions**

This study demonstrates that LLMs can perform comparably to human respondents in many areas of pediatric ORL, with specialists only marginally outperforming the highest-performing AI tools. However, LLMs still faces significant challenges, particularly in patient communication and source transparency. These findings highlight the dual nature of AI in clinical practice: as a promising tool with considerable potential and as one requiring cautious integration and refinement. AI can serve as a valuable complementary resource in clinical practice and medical education, particularly for residents and early-career clinicians. However, its limitations underscore the need

Trecca et al. 235

for continued development and human oversight to ensure patient safety and optimal care outcomes. With iterative advancements and careful integration, AI could become a transformative tool in pediatric ORL, supporting clinicians and enhancing the quality of care for young patients and their families.

#### **Author Contributions**

All authors (Eleonora M. C. Trecca, Vito Carlo Alberto Caponio, Mario Turri-Zanoni, Antonella Miriam di Lullo, Michele Gaffuri, Jérôme R. Lechien, Antonino Maniaci, Giuseppe Maruccio, Marella Reale, Irene Claudia Visconti, and Virginia Dallari) meet the ICMJE criteria: (1) Substantial contributions to the conception or design of the work; or the acquisition, analysis, or interpretation of data for the work. (2) Drafting the work or revising it critically for important intellectual content. (3) Final approval of the version to be published. (4) Agreement to be accountable for all aspects of the work in ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved.

#### **Disclosures**

Competing interests: The authors declare that there is no conflicts of interest.

Funding source: None.

# **Data Availability Statement**

Supplementary material is available upon request from the corresponding author. For further details, Appendices 1-3 are available.

#### Supplemental Material

Additional supporting information is available in the online version of the article.

#### **ORCID iD**

Eleonora M. C. Trecca http://orcid.org/0000-0001-6490-1746 Vito Carlo Alberto Caponio http://orcid.org/0000-0001-5080-5921

Mario Turri-Zanoni http://orcid.org/0000-0002-3678-9088 Antonella Miriam di Lullo https://orcid.org/0000-0002-4482-6782

Michele Gaffuri http://orcid.org/0000-0002-5435-685X
Jérôme R. Lechien http://orcid.org/0000-0002-0845-0845
Antonino Maniaci http://orcid.org/0000-0002-1251-0185
Giuseppe Maruccio https://orcid.org/0009-0008-6474-2122
Marella Reale http://orcid.org/0000-0002-9716-8605
Irene Claudia Visconti http://orcid.org/0000-0002-5831-1042
Virginia Dallari http://orcid.org/0000-0001-9260-1003

# X (formerly known as Twitter)

Eleonora M. C. Trecca X@EleonoraTrecca

#### References

1. Trecca EMC, Gaffuri M, Molinari G, et al. Impact of the COVID-19 pandemic on paediatric otolaryngology: a nationwide study. *Acta Otorhinolaryngol Ital*. 2023;43(5): 352-359. doi:10.14639/0392-100X-N2452

 Ariyaratne S, Iyengar KP, Nischal N, Chitti Babu N, Botchu R. A comparison of ChatGPT-generated articles with human-written articles. *Skeletal Radiol*. 2023;52(9): 1755-1758. doi:10.1007/s00256-023-04340-5

- 3. Seifen C, Huppertz T, Gouveris H, et al. Chasing sleep physicians: ChatGPT-40 on the interpretation of polysomnographic results. *Eur Arch Otrhinolaryngol*. Published online October 20, 2024. doi:10.1007/s00405-024-08985-3
- Májovský M, Černý M, Kasal M, Komarc M, Netuka D. Artificial intelligence can generate fraudulent but authenticlooking scientific medical articles: Pandora's box has been opened. J Med Internet Res. 2023;25:e46924. doi:10.2196/46924
- 5. Vaishya R, Misra A, Vaish A. ChatGPT: is this version good for healthcare and research? *Diabetes Metab Syndr Clin Res Rev.* 2023;17(4):102744. doi:10.1016/j.dsx.2023.102744
- Alkaissi H, McFarlane SI. Artificial hallucinations in ChatGPT: implications in scientific writing. *Cureus*. Published online February 19, 2023. doi:10.7759/cureus.35179
- 7. Vaira LA, Lechien JR, Abbate V, et al. Validation of the Quality Analysis of Medical Artificial Intelligence (QAMAI) tool: a new tool to assess the quality of health information provided by AI platforms. *Eur Arch Otrhinolaryngol*. 2024;281(11):6123-6131. doi:10.1007/s00405-024-08710-0
- 8. Lim B, Lirios G, Sakalkale A, Satheakeerthy S, Hayes D, Yeung JMC. Assessing the efficacy of artificial intelligence to provide peri-operative information for patients with a stoma. *ANZ J Surg.* Published online December 2, 2024. doi:10.1111/ans.19337
- Kaftan AN, Hussain MK, Naser FH. Response accuracy of ChatGPT 3.5 Copilot and Gemini in interpreting biochemical laboratory data a pilot study. *Sci Rep.* 2024;14(1):8233. doi:10.1038/s41598-024-58964-1
- Mayo-Yáñez M, Lechien JR, Maria-Saibene A, Vaira LA, Maniaci A, Chiesa-Estomba CM. Examining the performance of ChatGPT 3.5 and Microsoft Copilot in otolaryngology: a comparative study with otolaryngologists' evaluation. *Indian J Otolaryngol Head Neck Surg.* 2024;76(4):3465-3469. doi:10. 1007/s12070-024-04729-1
- 11. Vzorin GD, Bukinich AM, Sedykh AV, Vetrova II, Sergienko EA. The emotional intelligence of the GPT-4 large language model. *Psychol Russ State Art*. 2024;17(2): 85-99. doi:10.11621/pir.2024.0206
- 12. Sorin V, Brin D, Barash Y, et al. Large language models and empathy: systematic review. *J Med Internet Res.* 2024;26: e52597. doi:10.2196/52597
- Lechien JR, Briganti G, Vaira LA. Accuracy of ChatGPT-3.5 and -4 in providing scientific references in otolaryngology–head and neck surgery. *Eur Arch Otrhinolaryngol*. 2024;281(4): 2159-2165. doi:10.1007/s00405-023-08441-8
- Frosolini A, Franz L, Benedetti S, et al. Assessing the accuracy of ChatGPT references in head and neck and ENT disciplines. Eur Arch Otrhinolaryngol. 2023;280(11): 5129-5133. doi:10.1007/s00405-023-08205-4
- 15. Elali FR, Rachid LN. AI-generated research paper fabrication and plagiarism in the scientific community. *Patterns*. 2023;4(3):100706. doi:10.1016/j.patter.2023.100706

- Ghosh Moulic A, Gaurkar SS, Deshmukh PT. Artificial intelligence in otology, rhinology, and laryngology: a narrative review of its current and evolving picture. *Cureus*. Published online August 2, 2024. doi:10.7759/cureus.66036
- Radulesco T, Saibene AM, Michel J, Vaira LA, Lechien JR. ChatGPT-4 performance in rhinology: a clinical case series. *Int Forum Allergy Rhinol*. 2024;14(6):1123-1130. doi:10. 1002/alr.23323
- 18. Dallari V, Sacchetto A, Saetti R, Calabrese L, Vittadello F, Gazzini L. Is artificial intelligence ready to replace specialist doctors entirely? ENT specialists vs ChatGPT: 1-0, ball at

- the center. *Eur Arch Otrhinolaryngol*. 2024;281(2):995-1023. doi:10.1007/s00405-023-08321-1
- Malka RE, Marinelli JP, Newberry TR, Carlson ML, Bowe SN. Asynchronous learning among otolaryngology residents in the United States. *Am J Otolaryngol*. 2022;43(5):103575. doi:10.1016/j.amjoto.2022.103575
- Dallari V, Liberale C, De Cecco F, et al. The role of artificial intelligence in training ENT residents: a survey on ChatGPT, a new method of investigation. *Acta Otorhinolaryngol Ital*. 2024;44(3):161-168. doi:10.14639/0392-100X-N2806