

1 **TITLE**

2 ChatGPT's performance in dentistry and allergy-immunology assessments: a
3 comparative study

4 **Accepted for publication: October 4, 2023**

5 **AUTHOR NAMES AND INSTITUTIONAL AFFILIATIONS**

6 Alexander Fuchs¹, Tina Trachsel², Roland Weiger¹, Florin Eggmann¹

7 ¹Department of Periodontology, Endodontology, and Cariology, University Center for
8 Dental Medicine Basel UZB, University of Basel, Basel, Switzerland

9 ²Division of Allergy, University Children's Hospital Basel, Basel, Switzerland

10 **CORRESPONDENCE**

11 Dr. med. dent. Florin Eggmann, Klinik für Parodontologie, Endodontologie und
12 Kariologie, Universitäres Zentrum für Zahnmedizin Basel UZB, Universität Basel,
13 Mattenstrasse 40, CH-4058 Basel, Schweiz. Tel. +41 61 267 26 80; E-Mail:
14 florin.eggmann@unibas.ch

15 **Acknowledgments**

16 Alexander Fuchs contributed to this work as part of the requirements for his Master of
17 Dental Medicine degree at the University of Basel.

18 **Conflicts of interest**

19 The authors declare no financial or non-financial conflicts of interest related to this
20 work.

21

22 TITLE

23 ChatGPT's performance in dentistry and allergy-immunology assessments: a
24 comparative study

25 ABSTRACT

26 Large language models (LLMs) such as ChatGPT have potential applications in
27 healthcare, including dentistry. Priming, the practice of providing LLMs with initial,
28 relevant information, is an approach to improve their output quality. This study aimed
29 to evaluate the performance of ChatGPT 3 and ChatGPT 4 on self-assessment
30 questions for dentistry, through the Swiss Federal Licensing Examination in Dental
31 Medicine (SFLEDM), and allergy and clinical immunology, through the European
32 Examination in Allergy and Clinical Immunology (EEAACI). The second objective was
33 to assess the impact of priming on ChatGPT's performance. The SFLEDM and
34 EEAACI multiple-choice questions from the University of Bern's Institute for Medical
35 Education platform were administered to both ChatGPT versions, with and without
36 priming. Performance was analyzed based on correct responses. The statistical
37 analysis included Wilcoxon rank sum tests ($\alpha=0.05$). The average accuracy rates in
38 the SFLEDM and EEAACI assessments were 63.3% and 79.3%, respectively. Both
39 ChatGPT versions performed better on EEAACI than SFLEDM, with ChatGPT 4
40 outperforming ChatGPT 3 across all tests. ChatGPT 3's performance exhibited a
41 significant improvement with priming for both EEAACI ($p=0.017$) and SFLEDM
42 ($p=0.024$) assessments. For ChatGPT 4, the priming effect was significant only in the
43 SFLEDM assessment ($p=0.038$). The performance disparity between SFLEDM and
44 EEAACI assessments underscores ChatGPT's varying proficiency across different
45 medical domains, likely tied to the nature and amount of training data available in each
46 field. Priming can be a tool for enhancing output, especially in earlier LLMs.
47 Advancements from ChatGPT 3 to 4 highlight the rapid developments in LLM

48 technology. Yet, their use in critical fields such as healthcare must remain cautious
49 owing to LLMs' inherent limitations and risks.

50

51 **KEY WORDS**

52 Allergology, artificial intelligence, dental education, clinical immunology, machine
53 learning, medical informatics applications

54 **Introduction**

55 Machine learning applications have brought about significant advancements in
56 medicine, including dentistry (DUCRET ET AL. 2022, SCHWENDICKE ET AL. 2022, HAUG &
57 DRAZEN 2023). Among these advancements, a notable development has been the
58 emergence of large language models (LLMs) with a conversational interface, such as
59 ChatGPT, Bard, Baidu's Ernie Bot, Claude 2, Llama 2, and the chatbot function of the
60 revamped Bing search engine.

61 These LLMs, underpinned by deep learning transformer architectures, are trained on
62 vast amounts of tokenized text data (VASWANI ET AL. 2017). This training allows them
63 to generate fluent, contextually pertinent responses based on the input they receive.
64 Their capabilities span a wide range of tasks, from answering questions, summarizing
65 texts, translating languages, to writing computer code.

66 Priming, the practice of providing LLMs with initial, contextually relevant information, is
67 a useful approach to enhance their output quality (RAFFEL ET AL. 2020). By initiating a
68 conversation with strategically chosen words, phrases, or longer text excerpts, users
69 can guide LLMs to produce more accurate and contextually congruous responses.

70 LLMs have many potential use cases in healthcare, including dentistry (EGGMANN ET
71 AL. 2023). For instance, healthcare professionals could soon leverage LLMs to
72 streamline routine administrative tasks and improve patient education. However, LLMs
73 come with a set of significant risks and some inherent limitations (MELLO & GUHA 2023).
74 Many LLMs operate with knowledge cutoffs, which means they lack up-to-date
75 information (DASHTI ET AL. 2023). Determining the reliability of their response sources
76 can be difficult, if not impossible (WALKER ET AL. 2023). Moreover, LLMs sometimes
77 produce answers that seem plausible but are incorrect, underscoring the importance
78 of human oversight (DASHTI ET AL. 2023). Given these limitations, there are serious

79 concerns regarding the utility and safety of LLMs, especially in high-stakes fields of
80 application such as healthcare (BEAM ET AL. 2023).

81 In light of the potential implications of LLMs for healthcare, rigorous evaluation of their
82 outputs is paramount. By assessing LLMs' performance against external
83 benchmarks—including reasoning, coding, and knowledge tests—one can discern
84 their strengths and weaknesses (KUNG ET AL. 2023). Such evaluations can then inform
85 strategies to enhance LLMs' performance and guard against incautious use.

86 A prime resource for such evaluations is the University of Bern's Institute for Medical
87 Education (IML). The IML hosts a digital platform offering a vast array of self-
88 assessment questions tailored for dental and medical students and healthcare
89 professionals (<https://self-assessment.measured.iml.unibe.ch/>). Among its offerings
90 are multiple-choice questions designed for dental students preparing for the Swiss
91 Federal Licensing Examination in Dental Medicine (SFLEDM) and allergists and
92 immunologists preparing for the European Examination in Allergy and Clinical
93 Immunology (EEAACI). These curated question banks present an ideal tool for
94 assessing and comparing the performance of ChatGPT across distinct medical fields.
95 Considering the importance of examining the output accuracy of LLMs, this study
96 pursues two objectives. First, it aims to compare the performance of ChatGPT 3 and
97 ChatGPT 4 in responding to the SFLEDM and EEAACI self-assessment questions.
98 Second, it seeks to evaluate the impact of priming on ChatGPT's performance in these
99 assessments.

100

101 **Materials and Methods**

102 **Input sources**

103 The SFLEDM and EEAACI self-assessment questions were obtained from the IML
104 platform on February 13, 2023. While SFLEDM questions were translated from
105 German to English, EEAACI questions were already available in English. Any
106 questions with images or illustrations were excluded. The questions were of two
107 multiple-choice formats:

- 108 • A-type questions: These comprised a stem (either a question or a case
109 scenario) followed by potential answers. The task was to identify the single most
110 appropriate answer. Within the SFLEDM and EEAACI self-assessments, these
111 questions had four and five options, respectively.
- 112 • Kprim-type questions: These also started with a stem, succeeded by four
113 related statements or answers. The task was to determine the correctness of
114 each of these statements or answers.

115 The study used 32 SFLEDM questions, comprising 22 A-type and 10 Kprim-type
116 questions. In total, 28 EEAACI questions were used, comprising 19 A-type and 9
117 Kprim-type questions. The terms of service of the IML platform restrict the
118 dissemination of these self-assessment questions, even though they are publicly
119 accessible at <https://self-assessment.measured.iml.unibe.ch/> (last accessed on
120 October 3, 2023). They are therefore not featured in this report.

121 **Priming**

122 The primers, utilized to provide context for the questions, encompassed details about
123 the respective test, main subject information with relevant keywords, as well as
124 information about the question format and response guidelines. They offered a
125 thorough overview of the examination, including insights into the organizing body,

126 exam purpose, and covered topics, while also instructing the use of scientific reasoning
127 and adherence to general guidelines of the respective field for answering questions.
128 Designed to be analogous in length, structure, style, and content, the primers for the
129 SFLEDM and EEAACI self-assessments underwent several optimization rounds using
130 ChatGPT 3, adhering to principles of effective prompt design. Each trial for the primed
131 groups consistently utilized the same primer.
132 Conversely, the non-primed groups received a prompt that only supplied basic
133 information about the question format and response guidelines, deliberately omitting
134 context about the examination or topics to maintain succinctness and avoid priming.
135 Supplementary Table S-I details the input texts used for both primed and non-primed
136 groups prior to administering the multiple-choice questions.

137 **Administering questions to ChatGPT**

138 The tests involving ChatGPT 3 and ChatGPT 4 took place on February 19, 2023, and
139 March 25, 2023, respectively. For each group, 20 trials were conducted. Before
140 initiating each trial, the entire chat history was cleared. A new chat window was then
141 opened to eliminate any potential context carryover. For the non-primed groups, the
142 input prompt contained brief instructions on answering the questions, followed by either
143 the A-type or Kprim-type questions. In contrast, for the primed groups, the primer was
144 introduced before presenting the questions.

145 **Performance assessment**

146 An unblinded investigator recorded ChatGPT's responses in a pilot-tested
147 spreadsheet. For A-type questions, a score of 1 point was given for correct answers
148 and 0 points for incorrect ones. For Kprim-type questions, correctly answering all four
149 related statements or answers earned 1 point. If three out of the four statements or
150 answers were evaluated correctly, 0.5 points were given. A score of 0 points was
151 assigned if fewer than three statements or answers were correctly evaluated.

152 **Statistical analysis**

153 For each trial, the attained points were presented as a percentage of the maximum
154 possible points. This percentage was chosen as a performance metric since the
155 maximum points varied between the SFLEDM and EEAACI self-assessments.
156 Descriptive statistics, including mean, standard deviation, median, and interquartile
157 range, were computed for each group. Analysis of the distribution within each group
158 revealed a non-normal distribution, verified using a graphical method (normal
159 probability plot). Performance was analyzed between primed and non-primed groups
160 for both the SFLEDM and EEAACI self-assessments. This comparison was made
161 within each ChatGPT version, as well as across the ChatGPT 3 and ChatGPT 4
162 subsets. The Wilcoxon rank sum test was used for this analysis.

163 To assess the impact of priming, the improvement due to priming was calculated for
164 both the SFLEDM and EEAACI self-assessments within the ChatGPT 3 and ChatGPT
165 4 subsets. To quantify the improvement, trials—both without and with priming—were
166 ranked within their respective groups based on the percentage of the maximum points
167 attained. These ranks were paired before subtracting the values of the non-primed
168 groups from the values of the primed group, producing 20 improvement values within
169 each group. Analysis utilizing a normal probability plot confirmed a non-normal
170 distribution of data. Consequently, the Wilcoxon rank sum test was used for group
171 comparisons. The level of significance was set at $\alpha=0.05$. The statistical analyses were
172 performed by an unblinded investigator using R software (version 4.2.2, R Core Team,
173 R Foundation for Statistical Computing, Vienna, Austria). The dataset generated and
174 analyzed in this study is available in an open repository (FUCHS ET AL. 2023).

175

176

177 **Results**

178 Table I and Figure 1 present the detailed results. Both ChatGPT 3 and ChatGPT 4
179 exhibited superior performance in the EEAACI compared with the SFLEDM
180 assessment ($p < 0.001$). Overall, ChatGPT 4 scored higher than ChatGPT 3 across all
181 groups ($p < 0.001$). The performance gap between ChatGPT 4 and ChatGPT 3 was
182 wider in the EEAACI assessment than in the SFLEDM assessment. In the SFLEDM
183 assessment, without and with priming, the average percentage point increases for
184 ChatGPT 4 over ChatGPT 3 were 5.1 and 4.1, respectively. In contrast, for the EEAACI
185 assessment, these increases were 18.2 (without priming) and 15.0 (with priming).

186 Priming significantly enhanced the performance of ChatGPT 3 in both the SFLEDM
187 ($p = 0.012$) and EEAACI ($p = 0.001$) assessments. For ChatGPT 4, while there was a
188 significant performance increase in the SFLEDM assessment due to priming ($p = 0.03$),
189 priming had no significant effect on the performance in the EEAACI assessment
190 ($p = 0.221$).

191 As shown in Table II, with ChatGPT 3, priming enhanced the performance for EAACI
192 more than for SFLEDM ($p = 0.037$). Conversely, when using ChatGPT 4, priming
193 improved the performance for SFLEDM performance more than for EEAACI ($p = 0.002$).

194

195 **Discussion**

196 This study compared ChatGPT 3's and ChatGPT 4's performance on SFLEDM and
197 EEAACI self-assessment questions. These multiple-choice questions served to
198 benchmark and contrast the LLMs' proficiency in the field of dentistry and allergy and
199 immunology. The results showed that both versions performed better on the EEAACI,
200 with ChatGPT 4 surpassing ChatGPT 3 in all tests. Priming notably improved ChatGPT
201 3's performance in both tests, but only impacted ChatGPT 4 in the SFLEDM
202 assessment.

203 The observed performance disparity between the EEAACI and SFLEDM assessments
204 suggests that ChatGPT's proficiency may vary across all medical specialties. One
205 plausible explanation for this disparity may lie in the nature of the data the LLM has
206 been trained on (PATCAS ET AL. 2022, BORNSTEIN 2023, WALKER ET AL. 2023). Most of
207 the medical literature, discussions, and queries available in open sources focus on
208 broader medical fields, with allergy and immunology being more extensively
209 represented than smaller branches of medicine such as dentistry. Furthermore, in
210 dentistry, diagnoses and treatments frequently rely heavily on physical examinations
211 and imaging, aspects that textual models such as ChatGPT are not adept at grasping.
212 By contrast, allergy and immunology, being more systemic and often reliant on patient
213 history and laboratory results, are better suited for textual analysis and understanding
214 by LLMs.

215 This study, while specifically pertaining to the SFLEDM and EEAACI assessments,
216 may offer broader implications for the application of LLMs in other medical domains.
217 The observed performance disparities and the impact of priming across different
218 assessments suggest that the effectiveness of LLMs can be significantly influenced by
219 the subject matter. Extending this to other domains, it becomes pivotal to consider the
220 availability and specificity of training data, as well as the inherent characteristics of the

221 medical field in question. For instance, medical specialties that heavily rely on textual
222 information and have abundant data available might observe better LLM performance,
223 akin to the results seen in the EEAACI assessment. Conversely, fields that depend
224 more on visual or practical elements may present additional challenges for LLMs, as
225 seen in the SFLEDM assessment. Further research is warranted to explore these
226 dynamics, identifying patterns and strategies to optimize LLMs' performance across
227 diverse medical specialties.

228 ChatGPT's performance has been studied across various medical knowledge
229 examinations, with the accuracy rates demonstrating considerable variation among
230 different tests and medical disciplines. A recent systematic review and meta-analysis
231 revealed that the performance range for ChatGPT 3.5 in these evaluations spanned
232 from 40% to 100%, with an average accuracy rate of 61.1% (LEVIN ET AL. 2023). The
233 mean performance of 63.3% in the SFLEDM assessment aligns with this average
234 across medical domains. In contrast, ChatGPT's performance in the EEAACI
235 assessment yielded a higher average score of 79.3%, placing it at the top range
236 compared with results from other studies (LEVIN ET AL. 2023). It is noteworthy that
237 ChatGPT 4, when primed, exceeded the commonly used passing threshold of 60% in
238 the SFLEDM assessment. This level of performance was observed in the EEAACI
239 assessment across the two examined ChatGPT iterations, regardless of priming.

240 In dental and medical education, LLM chatbots hold potential for supplementing
241 learning materials and providing interactive learning opportunities (ALI ET AL. 2023).
242 However, it is important to emphasize that while ChatGPT 3 and ChatGPT 4 showed
243 promise in answering self-assessment questions from the SFLEDM and EEAACI, they
244 should not be relied on for exam preparation. Despite their capabilities, these LLMs
245 frequently provide inaccurate or misleading information (MELLO & GUHA 2023). Relying
246 on ChatGPT for exam preparation, especially in critical fields like healthcare, could

247 lead to misconceptions and an incomplete understanding of the subject matter (LEVIN
248 ET AL. 2023, SAAD ET AL. 2023). Therefore, it is crucial to always approach their outputs
249 with caution and cross-reference with trusted educational resources. Today, LLMs
250 should serve merely as supplements to more traditional methods of information
251 seeking (MELLO & GUHA 2023).

252 The noticeable performance improvement from ChatGPT 3 to ChatGPT 4, available
253 exclusively to subscription fee payers, underscores the rapid advancements in LLM
254 development within a short period. Comparable enhancements in response quality
255 from ChatGPT 4 have been noted for queries related to dermatology and myopia
256 (LEWANDOWSKI ET AL. 2023, LIM ET AL. 2023). Moreover, while ChatGPT 3 operated
257 solely with text, ChatGPT 4 is multimodal, allowing it to accept and produce text and
258 image inputs and outputs. This shift to multimodality represents a substantial
259 enhancement in ChatGPT's functionality. The increasing adaptability of LLMs suggests
260 that they might soon serve as additional tools in specific use cases in healthcare (VAIRA
261 ET AL. 2023).

262 However, on the road to artificial general intelligence, LLMs underpinned by the next-
263 token-prediction paradigm are likely an off-ramp (MARCUS 2022). Their capabilities,
264 based on brute statistics, are impressive, but their genuine understanding remains
265 shallow (THIRUNAVUKARASU 2023). Medical professionals, including allergists,
266 immunologists, and dentists, are therefore not predicted to face major changes due to
267 the widespread adoption of LLM applications (THIRUNAVUKARASU 2023).

268 Priming and adept prompt design serve as strategic tools to guide LLMs towards
269 generating more contextually congruous responses (RAFFEL ET AL. 2020). The results
270 of this study are in line with this assertion, particularly with ChatGPT 3, where priming
271 significantly enhanced its performance in both the SFLEDM and EEAACI
272 assessments. However, whereas priming exhibited a significant impact on ChatGPT

273 4's performance in the SFLEDM assessment, its influence was negligible in the
274 EEAACI assessment. This difference underscores the evolving nature of LLMs and
275 suggests that as these models become more advanced, the relative impact of priming
276 may vary depending on the complexity and specificity of the task at hand.

277 This study has several limitations that warrant careful consideration. First, the
278 questions from the IML platform, specifically the SFLEDM and EEAACI self-
279 assessments, represented only a narrow spectrum of knowledge within dentistry,
280 allergy, and immunology. This limits the generalizability of the findings.

281 Second, tasks like answering board examination questions or retrieving information
282 from medical records have only a tangential connection to real-world care decisions
283 (MELLO & GUHA 2023). This means that assessments using such tasks as benchmarks
284 offer limited insight into a LLM's usefulness for clinical decision support (MELLO & GUHA
285 2023).

286 Third, the translation of SFLEDM questions from German to English introduced
287 potential biases, as nuances in language might affect the LLM's comprehension and
288 response accuracy.

289 Fourth, the exclusion of questions with images or illustrations omits a significant aspect
290 of medical assessments, which often rely on visual diagnostics and the interpretation
291 of data charts and graphs.

292 Fifth, an unblinded evaluator recorded and graded ChatGPT's responses to the
293 multiple-choice questions. Since the answer key for these questions was objective and
294 definitive, allowing no room for interpretation or discretion, calibration procedures,
295 evaluator blinding, and employment of multiple evaluators were foregone.
296 Nonetheless, to guard against potential biases inherent in unblinded assessments—
297 even when utilizing unequivocal answer keys—future investigations should consider
298 implementing evaluator calibration and blinding.

299

300 Sixth, by focusing solely on two versions of ChatGPT, the study did not capture the full
301 range of LLM capabilities across various models or iterations. These limitations
302 emphasize the critical need for additional research to thoroughly evaluate the
303 performance and potential impact of LLMs in medical disciplines.

304

305 **Conclusion**

306 Within the constraints of this study, the following conclusions were drawn:

- 307 • ChatGPT 3 and ChatGPT 4 both demonstrated stronger performance on the
308 EEAACI compared with the SFLEDM assessment. This performance disparity
309 highlights ChatGPT's varying proficiency across different medical domains,
310 likely influenced by the type and volume of training data available in each field.
- 311 • Priming improved ChatGPT 3's performance across both assessments. For
312 ChatGPT 4, while priming influenced results in the SFLEDM assessment, its
313 effect was negligible for the EEAACI. This underscores the nuanced influence
314 of priming as LLMs become more advanced.
- 315 • The progress from ChatGPT 3 to ChatGPT 4 reveals rapid advancements in
316 LLM development, including the shift to multimodality. Yet, their enhanced
317 capabilities notwithstanding, LLMs have major inherent limitations and risks,
318 emphasizing the need for cautious use in high-stakes fields such as healthcare.

319

320

321 **Table I** Results of the performance assessments

Assessment	LLM	Priming	N	Mean	SD	Median	IQR
SFLEDM	ChatGTP 3	None	20	59.3%	5.2%	59.4%	5.5%
		Yes	20	62.6%	3.3%	62.5%	3.1%
	ChatGTP 4	None	20	64.4%	2.8%	64.1%	3.5%
		Yes	20	66.7%	3.2%	66.4%	3.5%
EEAACI	ChatGTP 3	None	20	69.0%	3.7%	67.9%	5.4%
		Yes	20	72.9%	2.6%	73.2%	4.0%
	ChatGTP 4	None	20	87.2%	1.9%	87.5%	3.6%
		Yes	20	87.9%	1.9%	87.5%	2.2%

322 *EEAACI*, European Examination in Allergy and Clinical Immunology; *IQR*, interquartile range; *LLM*, large language323 model; *SD*, standard deviation; *SFLEDM*, Swiss Federal Licensing Examination in Dental Medicine

324

325 **Table II** Performance improvement through priming

Assessment	LLM	N	Mean	SD	Median	IQR
SFLEDM	ChatGTP 3	20	3.3%	2.2%	3.1%	3.1%
	ChatGTP 4	20	2.3%	0.9%	1.6%	1.6%
EEAACI	ChatGTP 3	20	3.9%	1.5%	3.6%	1.2%
	ChatGTP 4	20	0.7%	0.9%	0.0%	1.8%

326 *EEAACI*, European Examination in Allergy and Clinical Immunology; *IQR*, interquartile range; *LLM*, large language327 model; *SD*, standard deviation; *SFLEDM*, Swiss Federal Licensing Examination in Dental Medicine

328

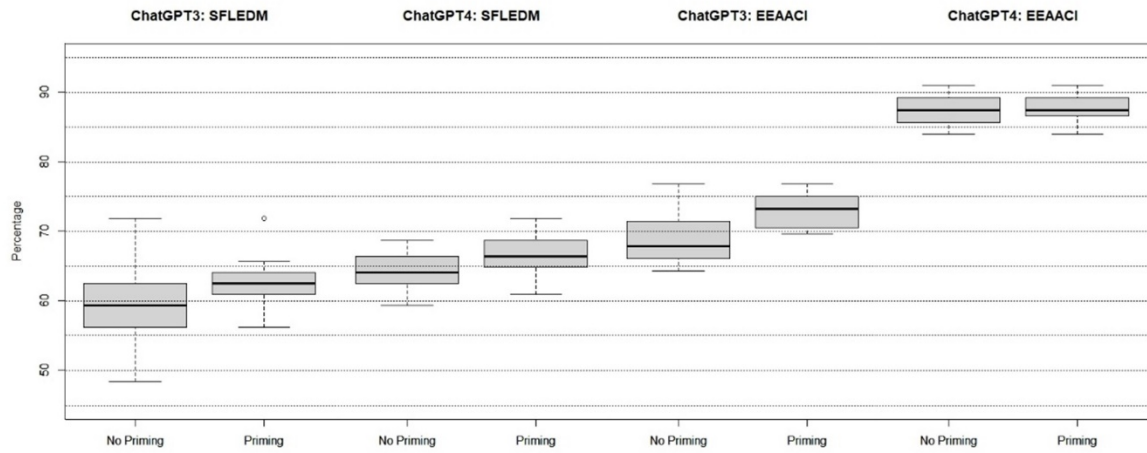
329

330

331

332

333



334

335 **Figure 1** Box plots depicting the distribution of assessment scores, represented by
 336 correct response rates, across the eight distinct groups.

337

338 **Supplementary Table S-I** Input texts used for the primed and non-primed groups
 339 before administering the multiple-choice questions to ChatGPT

Assessment	Multiple-choice format	Priming	Input text
EEAACI	A-type questions	None	Please answer the following single-choice questions. For each question, please clearly indicate which answer (A, B, C, D, or E) is correct. If uncertain, please clearly indicate the most probable answer. List the correct answers in a table. Number the questions as given.
		Yes	The European Academy of Allergy & Clinical Immunology (EAACI) has been conducting the European Examination in Allergology and Clinical Immunology annually since 2008. The exam is designed to test candidates' knowledge on a wide range of topics related to allergology, including allergens, dermatology, respiratory and pediatric allergy, anaphylaxis, venom hypersensitivity, drug and food hypersensitivity, as well as relevant issues such as pregnancy and allergology, occupational allergies, eosinophilic disorders, mastocytosis, and CI-INH deficiency. The exam also covers basic immunology and clinical immunology, including autoimmunity and immune deficiency. To aid allergists and immunologists in preparing for the exam, training questions are available for practice. Please answer the following single-choice questions. For each question, please clearly indicate which answer (A, B, C, D, or E) is correct. Use scientific reasoning and general guidelines for allergology and immunology to answer the questions correctly. If uncertain, please clearly indicate the most probable answer. List the correct answers in a table. Number the questions as given.
Kprim-type	questions	None	Please answer the following multiple-choice questions. For each question, please clearly indicate for each answer (A, B, C, and D) whether the answer is correct or incorrect. List the answers in a table. Number the questions as given.
		Yes	The European Academy of Allergy & Clinical Immunology (EAACI) has been conducting the European Examination in Allergology and Clinical Immunology annually since 2008. The exam is designed to test candidates' knowledge on a wide range of topics related to allergology, including allergens, dermatology, respiratory and pediatric allergy, anaphylaxis, venom hypersensitivity, drug and food hypersensitivity, as well as relevant issues such as pregnancy and allergology, occupational allergies, eosinophilic disorders, mastocytosis, and CI-INH deficiency. The exam also covers basic immunology and clinical immunology, including autoimmunity and immune deficiency. To aid allergists and immunologists in preparing for the exam, training questions are available for practice. Please answer the following multiple-choice questions. For each question, please clearly indicate for each answer (A, B, C, and D) whether the answer is correct or incorrect. Therefore, each question can have 0, 1, 2, 3, or 4 correct answers. Use scientific reasoning and general guidelines for allergology and immunology to answer the questions correctly. List the answers in a table. Number the questions as given.
SFLEDM	A-type questions	None	Please answer the following single-choice questions. For each question, please clearly indicate which answer (A, B, C, or D) is correct. If uncertain, please clearly indicate the most probable answer. List the correct answers in a table. Number the questions as given.
		Yes	You will be asked questions from the Swiss Federal Licensing Examination in dental medicine. The exam is designed to test candidates' knowledge on a wide range of topics related to dental medicine, including preventive dentistry, stomatology, oral health, cariology, restorative dentistry, endodontics, periodontics, dental implantology, prosthodontics, esthetic dentistry, pediatric dentistry, orthodontics, dental radiology, geriatric dentistry, special needs dentistry, and communication in the dentist-patient relationship. Additionally, the catalogue of learning objectives demands that candidates know the most common medical issues and corresponding treatment approaches in the following medical specialties: infectiology, internal medicine, oral and maxillofacial surgery, dermatology/allergy, psychiatry, geriatrics, and care for patients with special needs. To aid dental students in preparing for the Swiss Federal Licensing Examination in dental medicine, training questions are available for practice. Please answer the following single-choice questions. For each question, please clearly indicate which answer (A, B, C, or D) is correct. Use scientific reasoning and general guidelines for dentistry to answer the questions correctly. If uncertain,

please clearly indicate the most probable answer. List the correct answers in a table. Number the questions as given.

Kprim-type questions	None	Please answer the following multiple-choice questions. For each question, please clearly indicate for each answer (A, B, C, and D) whether the answer is correct or incorrect. List the answers in a table. Number the questions as given.
	Yes	You will be asked questions from the Swiss Federal Licensing Examination in dental medicine. The exam is designed to test candidates' knowledge on a wide range of topics related to dental medicine, including preventive dentistry, stomatology, oral health, cariology, restorative dentistry, endodontics, periodontics, dental implantology, prosthodontics, esthetic dentistry, pediatric dentistry, orthodontics, dental radiology, geriatric dentistry, special needs dentistry, and communication in the dentist-patient relationship. Additionally, the catalogue of learning objectives demands that candidates know the most common medical issues and corresponding treatment approaches in the following medical specialties: infectiology, internal medicine, oral and maxillofacial surgery, dermatology/allergy, psychiatry, geriatrics, and care for patients with special needs. To aid dental students in preparing for the Swiss Federal Licensing Examination in dental medicine, training questions are available for practice. Please answer the following multiple-choice questions. For each question, please clearly indicate for each answer (A, B, C, and D) whether the answer is correct or incorrect. Therefore, each question can have 0, 1, 2, 3, or 4 correct answers. Use scientific reasoning and general guidelines for dentistry to answer the questions correctly. List the answers in a table. Number the questions as given.

340 *EEAACI*, European Examination in Allergy and Clinical Immunology; *SFLEDM*, Swiss Federal Licensing
 341 Examination in Dental Medicine
 342
 343
 344
 345
 346
 347
 348
 349
 350

351 **Zusammenfassung**

352 **Einleitung**

353 Anwendungen der künstlichen Intelligenz (KI) können dem Gesundheitspersonal,
354 einschliesslich Zahnärzten, verschiedene Vorteile bieten. Grosse Sprachmodelle
355 (GSM) sind KI-Anwendungen, die mit grossen Mengen von Textdaten trainiert werden
356 und verschiedene sprachbezogene Aufgaben durchführen können. ChatGPT, ein
357 GSM mit einer Konversationsschnittstelle, wurde im November 2022 auf den Markt
358 gebracht und ist online verfügbar. Trotz seiner beeindruckenden Fähigkeiten hat
359 ChatGPT erhebliche Einschränkungen und Unzulänglichkeiten. Beispielsweise gibt
360 ChatGPT teilweise fehlerhafte Antworten oder stellt Fehlinformationen als Fakten dar.
361 Vor der Anwendung GSM in medizinischen Disziplinen ist es von grosser Bedeutung,
362 die Fähigkeiten und Grenzen von GSM zu verstehen. Ein interessanter Ansatz ist das
363 "Priming", bei einem GSM vorab relevante Informationen gegeben werden, um die
364 Qualität seiner Antworten zu verbessern. Diese Studie konzentriert sich auf die
365 Bewertung der Leistung von ChatGPT Versionen 3 und 4 in den medizinischen
366 Bereichen Zahnmedizin sowie Allergologie und klinische Immunologie, unter
367 besonderer Berücksichtigung des Priming-Effekts.

368 **Material und Methoden**

369 Zur umfassenden Evaluation von ChatGPT wurden Multiple-Choice-Fragen zur
370 Selbstbewertung in Zahnmedizin («Swiss Federal Licensing Examination in Dental
371 Medicine» [SFLEDM]) und Allergologie sowie klinischer Immunologie («European
372 Examination in Allergy and Clinical Immunology» [EEAACI]) vom Institut für
373 Medizinische Lehre der Universität Bern zusammengestellt. ChatGPT 3 und 4 wurden
374 unter zwei Bedingungen getestet: mit Priming und ohne Priming. Das Hauptkriterium
375 für die Leistungsbewertung war die Genauigkeitsrate, gemessen an der Anzahl korrekt

376 beantworteter Fragen. Die statistischen Analysen erfolgten mittels Wilcoxon-
377 Rangsummentests mit einem Signifikanzniveau von $\alpha = 0,05$.

378 **Resultate**

379 Im SFLEDM-Bereich betrug die durchschnittliche Genauigkeitsrate 63,3%. Im
380 Gegensatz dazu zeigte ChatGPT im EEAACI-Bereich mit einer durchschnittlichen
381 Genauigkeit von 79,3% eine überlegene Leistung. Beide ChatGPT-Modelle zeigten im
382 EEAACI-Bereich bessere Leistungen als im SFLEDM-Bereich. Bemerkenswert ist,
383 dass ChatGPT 4 durchgehend bessere Leistungen als ChatGPT 3 in beiden Bereichen
384 zeigte. In Bezug auf das Priming zeigte ChatGPT 3 sowohl bei den Fragen aus dem
385 EEAACI Bereich ($p=0,001$) als auch im SFLEDM Bereich ($p=0,012$) eine deutliche
386 Verbesserung bei Verwendung von Priming. Im Gegensatz dazu verbesserte sich die
387 Leistung durch Priming bei ChatGPT 4 nur im SFLEDM-Bereich signifikant ($p=0,03$).

388 **Diskussion**

389 Die unterschiedliche Leistung von ChatGPT in der Beantwortung von Multiple-Choice-
390 Fragen aus dem SFLEDM und EEAACI Bereich weist auf eine unterschiedliche
391 Kompetenz von GSM in verschiedenen medizinischen Bereichen hin. Diese
392 unterschiedliche Kompetenz könnte durch die Art und das Volumen der verfügbaren
393 Trainingsdaten für jeden Bereich beeinflusst werden. Priming erweist sich als
394 vorteilhafte Methode zur Leistungsverbesserung von GSM, besonders bei älteren
395 Versionen wie ChatGPT 3. Der signifikante Leistungszuwachs von ChatGPT 3 zu 4
396 unterstreicht die rasanten Entwicklungen in der GSM-Technologie. Dennoch ist beim
397 Einsatz von GSM im Gesundheitssektor, einschliesslich der Zahnmedizin, höchste
398 Sorgfalt und Umsicht angebracht, denn GSM weisen weiterhin zahlreiche Limitationen
399 und Risiken auf.

400

401

402 **Résumé**

403 **Introduction**

404 Les applications d'intelligence artificielle (IA) peuvent offrir divers avantages aux
405 professionnels de la santé, y compris aux dentistes. Les modèles de langage de
406 grande taille (abrégé LLM de l'anglais *large language model*) sont des applications d'IA
407 entraînées avec de grandes quantités de données textuelles et capables d'effectuer
408 différentes tâches liées à la langue. ChatGPT, un LLM doté d'une interface
409 conversationnelle, a été lancé en novembre 2022 et est disponible en ligne. Malgré
410 ses capacités impressionnantes, ChatGPT présente des limitations et des
411 insuffisances importantes. Par exemple, ChatGPT donne parfois des réponses
412 erronées ou présente des informations erronées comme des faits. En raison de la
413 nature critique des disciplines médicales, il est très important de comprendre les
414 capacités et les limites du LLM. Une approche intéressante est le "priming", qui
415 consiste à donner au LLM des informations pertinentes à l'avance afin d'améliorer la
416 qualité de ses réponses. Cette étude se concentre sur l'évaluation des performances
417 de ChatGPT versions 3 et 4 dans les domaines médicaux de la dentisterie ainsi que
418 de l'allergologie et de l'immunologie clinique, en accordant une attention particulière à
419 l'effet d'amorçage.

420 **Matériels et méthodes**

421 Pour une évaluation complète de ChatGPT, des questions à choix multiples d'auto-
422 évaluation en médecine dentaire («Swiss Federal Licensing Examination in Dental
423 Medicine» [SFLEDM]) et en allergologie et immunologie clinique («European
424 Examination in Allergy and Clinical Immunology» [EEAACI]) ont été compilées par
425 l'Institut d'enseignement médical de l'Université de Berne. ChatGPT 3 et 4 ont été
426 testés dans deux conditions : avec et sans amorçage. Le principal critère d'évaluation
427 des performances était le taux de précision, mesuré par le nombre de questions

428 auxquelles il a été répondu correctement. Les analyses statistiques ont été effectuées
429 à l'aide de tests de répartition des rangs de Wilcoxon avec un niveau de signification
430 de $\alpha = 0,05$.

431 **Résultats**

432 Dans le domaine SFLEDM, le taux de précision moyen était de 63,3%. En revanche,
433 ChatGPT a montré une performance supérieure dans le domaine EEAACI, avec une
434 précision moyenne de 79,3%. Les deux modèles ChatGPT ont montré de meilleures
435 performances dans le domaine EEAACI que dans le domaine SFLEDM. Il est à noter
436 que ChatGPT 4 a montré des performances systématiquement meilleures que
437 ChatGPT 3 dans les deux domaines. En ce qui concerne l'amorçage, ChatGPT 3 a
438 montré une nette amélioration lors de l'utilisation de l'amorçage, tant pour les questions
439 du domaine EEAACI ($p=0,001$) que pour le domaine SFLEDM ($p=0,012$). En revanche,
440 la performance de ChatGPT 4 ne s'est améliorée de manière significative par
441 l'amorçage que dans le domaine SFLEDM ($p=0,03$).

442 **Discussion**

443 La différence de performance de ChatGPT dans les réponses aux questions à choix
444 multiples des domaines SFLEDM et EEAACI pourrait indiquer une différence de
445 compétence des LLMs dans différents domaines médicaux. Cette différence de
446 compétence pourrait être influencée par le type et le volume des données
447 d'entraînement disponibles pour chaque domaine. L'amorçage s'avère être une
448 méthode avantageuse pour améliorer les performances du LLM, en particulier pour les
449 anciennes versions comme ChatGPT 3. L'augmentation significative des
450 performances de ChatGPT 3 à 4 souligne les développements rapides de la
451 technologie LLM. Toutefois, l'utilisation du LLM dans le secteur de la santé, y compris
452 la médecine dentaire, requiert la plus grande prudence et le plus grand soin. En effet,
453 les LLMs présentent encore de nombreuses limites et risques.

454 **References**

- 455 ALI K, BARHOM N, TAMIMI F, DUGGAL M: ChatGPT-A double-edged sword for
456 healthcare education? Implications for assessments of dental students. *Eur J Dent*
457 *Educ* (2023). doi: 10.1111/eje.12937.
- 458 BEAM A L, DRAZEN J M, KOHANE I S, LEONG T-Y, MANRAI A K, RUBIN E J: Artificial
459 intelligence in medicine. *N Engl J Med* 388: 1220–1221 (2023)
- 460 BORNSTEIN M M: Artificial intelligence and personalised dental medicine - just a hype
461 or true game changers? *Br Dent J* 234: 755 (2023)
- 462 DASHTI M, LONDONO J, GHASEMI S, MOGHADDASI N: How much can we rely on artificial
463 intelligence chatbots such as the ChatGPT software program to assist with scientific
464 writing? *J Prosthet Dent* (2023). doi: 10.1016/j.prosdent.2023.05.023.
- 465 DUCRET M, MÖRCH C-M, KARTEVA T, FISHER J, SCHWENDICKE F: Artificial intelligence
466 for sustainable oral healthcare. *J Dent* 127: 104344 (2022)
- 467 EGGMANN F, WEIGER R, ZITZMANN N U, BLATZ M B: Implications of large language
468 models such as ChatGPT for dental medicine. *J Esthet Restor Dent* (2023). doi:
469 10.1111/jerd.13046.
- 470 FUCHS A, TRACHSEL T, WEIGER R, EGGMANN F: ChatGPT's performance in dentistry
471 and allergy-immunology assessments: a comparative study (Version 1) [Data set].
472 Zenodo (2023). <https://doi.org/10.5281/zenodo.8331147>
- 473 HAUG C J, DRAZEN J M: Artificial intelligence and machine learning in clinical
474 medicine, 2023. *N Engl J Med* 388: 1201–1208 (2023)
- 475 KUNG T H, CHEATHAM M, MEDENILLA A, SILLOS C, DE LEON L, ELEPAÑO C, MADRIAGA M,
476 AGGABAO R, DIAZ-CANDIDO G, MANINGO J, TSENG V: Performance of ChatGPT on
477 USMLE: potential for AI-assisted medical education using large language models.
478 *PLOS Digit Health* 2: e0000198 (2023). doi: 10.1371/journal.pdig.0000198.

- 479 LEVIN G, HOESH N, BREZINOV Y, MEYER R: Performance of ChatGPT in medical
480 examinations: a systematic review and a meta-analysis. *BJOG* (2023). doi:
481 10.1111/1471-0528.17641.
- 482 LEWANDOWSKI M, ŁUKOWICZ P, ŚWIETLIK D, BARAŃSKA-RYBAK W: An original study of
483 ChatGPT-3.5 and ChatGPT-4 dermatological knowledge level based on the
484 dermatology specialty certificate examinations. *Clin Exp Dermatol* (2023). doi:
485 10.1093/ced/llad255.
- 486 LIM Z W, PUSHPANATHAN K, YEW S M E, LAI Y, SUN C-H, LAM J S H, CHEN D Z, GOH J H
487 L, TAN M C J, SHENG B, CHENG C-Y, KOH V T C, THAM Y-C: Benchmarking large
488 language models' performances for myopia care: a comparative analysis of
489 ChatGPT-3.5, ChatGPT-4.0, and Google Bard. *EBioMedicine* 95: 104770 (2023)
- 490 MARCUS G: Deep learning is hitting a wall. Accessed on October 3, 2023.
491 <https://nautil.us/deep-learning-is-hitting-a-wall-238440/>. (2022)
- 492 MELLO M M, GUHA N: ChatGPT and physicians' malpractice risk. *JAMA Health Forum*
493 4: e231938 (2023)
- 494 PATCAS R, BORNSTEIN M M, SCHÄTZLE M A, TIMOFTE R: Artificial intelligence in medico-
495 dental diagnostics of the face: a narrative review of opportunities and challenges.
496 *Clin Oral Investig* 26: 6871–6879 (2022)
- 497 RAFFEL C, SHAZEER N, ROBERTS A, LEE K, NARANG S, MATENA M, ZHOU Y, LI W, LIU P
498 J: Exploring the limits of transfer learning with a unified text-to-text transformer. *The*
499 *Journal of Machine Learning Research* 21: 5485–5551 (2020)
- 500 SAAD A, IYENGAR K P, KURISUNKAL V, BOTCHU R: Assessing ChatGPT's ability to pass
501 the FRCS orthopaedic part A exam: a critical analysis. *Surgeon* (2023). doi:
502 10.1016/j.surge.2023.07.001.

503 SCHWENDICKE F, CEJUDO GRANO DE ORO, J, GARCIA CANTU A, MEYER-LUECKEL H,
504 CHAURASIA A, KROIS J: Artificial intelligence for caries detection: value of data and
505 information. *J Dent Res* 101: 1350–1356 (2022)

506 THIRUNAVUKARASU A J: Large language models will not replace healthcare
507 professionals: curbing popular fears and hype. *J R Soc Med* 116: 181–182 (2023)

508 VAIRA L A, LECHIEN J R, ABBATE V, ALLEVI F, AUDINO G, BELTRAMINI G A, BERGONZANI
509 M, BOLZONI A, COMMITTERI U, CRIMI S, GABRIELE G, LONARDI F, MAGLITTO F,
510 PETROCELLI M, PUCCI R, SAPONARO G, TEL A, VELLONE V, CHIESA-ESTOMBA C M,
511 BOSCOLO-RIZZO P, SALZANO G, DE RIU G: Accuracy of ChatGPT-generated
512 information on head and neck and oromaxillofacial urgency: a multicenter
513 collaborative analysis. *Otolaryngol Head Neck Surg* (2023). doi: 10.1002/ohn.489.

514 VASWANI A, SHAZEER N, PARMAR N, USZKOREIT J, JONES L, GOMEZ A N, KAISER Ł,
515 POLOSUKHIN I: Attention is all you need. *Advances in Neural Information Processing*
516 *Systems* 30: 1–11 (2017)

517 WALKER H L, GHANI S, KUEMMERLI C, NEBIKER C A, MÜLLER B P, RAPTIS D A, STAUBLI S
518 M: Reliability of medical information provided by ChatGPT: assessment against
519 clinical guidelines and patient information quality instrument. *J Med Internet Res* 25:
520 e47479 (2023)

521