

1 **Protocol For Human Evaluation of Artificial Intelligence Chatbots in Clinical**
2 **Consultations**

3 Edwin Kwan-Yeung Chiu, MRCP(UK)^a; Tom Wai-Hin Chung, FRCPath^a.

4 ^aDepartment of Microbiology, Li Ka Shing Faculty of Medicine, The University of Hong Kong,
5 Hong Kong, China.

6 **Keywords:** artificial intelligence, generative, large language model, chatbot, infectious
7 diseases, microbiology, consultation.

8 **Running title:** Protocol for Human Evaluation of AI chatbots in Clinical Consultations

9 **Correspondence:**

10 Tom Wai-Hin Chung, Department of Microbiology, Li Ka Shing Faculty of Medicine, The
11 University of Hong Kong, Queen Mary Hospital, 102 Pokfulam Road, Hong Kong, China.
12 Phone: (852) 22552409. Fax: (852) 28724555. E-mail: tomwhc@hku.hk. ORCID iD:
13 <https://orcid.org/0000-0003-1780-821X>.

14 **Abstract**

15 **Background**

16 Generative artificial intelligence (AI) technology has the revolutionary potentials to augment
17 clinical practice and telemedicine. The nuances of real-life patient scenarios and complex
18 clinical environments demand a rigorous, evidence-based approach to ensure safe and effective
19 application.

20 **Methods**

21 We present a protocol for the systematic evaluation of generative AI large language models
22 (LLMs) as chatbots within the context of clinical microbiology and infectious disease
23 consultations. We aim to critically assess the clinical accuracy, comprehensiveness, coherence,
24 and safety of recommendations produced by leading generative AI models, including Claude
25 2, Gemini Pro, GPT-4.0, and a GPT-4.0-based custom AI chatbot.

26 **Discussion**

27 A standardised healthcare-specific prompt template is employed to elicit clinically impactful
28 AI responses. Generated responses will be graded by a panel of human evaluators,
29 encompassing a wide spectrum of domain expertise in clinical microbiology and virology and
30 clinical infectious diseases. Evaluations are performed using a 5-point Likert scale across four
31 clinical domains: factual consistency, comprehensiveness, coherence, and medical harmfulness.
32 Our study will offer insights into the feasibility, limitations, and boundaries of generative AI
33 in healthcare, providing guidance for future research and clinical implementation. Ethical
34 guidelines and safety guardrails should be developed to uphold patient safety and clinical
35 standards.

36 INTRODUCTION

37 With global aging population, ever increasing healthcare demands and the rapid evolution of
38 healthcare technologies, effective integration of artificial intelligence (AI) into clinical
39 workflow and decision-making processes have become a focal point of research and debate.
40 Generative AI have demonstrated significant potentials in understanding natural language and
41 addressing cognitive tasks. (1) The prospects of generative AI replacing or augmenting
42 physician tasks, particularly in telemedicine where information exchange is primarily text-
43 based, has prompted investigations into their practicality and safety in clinical consultations.
44 (2)
45 Preliminary investigations have demonstrated the potentials for AI in managing various
46 infectious disease syndromes (e.g., bloodstream infections and brain abscesses), however,
47 concerns remain about the reliability, safety, and ethics of the utilisation of generative AI in
48 clinical practices.(3-5) This study is among the first to systemically evaluate the state-of-the-
49 art generative AI large language model (LLM) chatbots, including a custom AI chatbot (custom
50 bot) integrated with domain-specific medical literature. In addition, this study employs a novel
51 self-developed healthcare-specific prompt template purposely designed to examine AI chatbot
52 performances in complex real-life clinical scenarios. A unique dual-tier evaluation system that
53 includes both consultant-level specialists and non-specialist physicians is also implemented in
54 the evaluation process to offer a comprehensive assessment from multiple levels of domain
55 expertise and clinical practice.
56 The objective of this protocol is to critically assess the clinical accuracy, coherence,
57 comprehensiveness, and safety of recommendations provided by AI chatbots. This research
58 aims to contribute to the ongoing discourse on the role of generative AI in healthcare and to
59 aid in the development of guidelines that ensure the safe and effective implementation of
60 generative AI in clinical microbiology and infectious disease domains.

61 **MATERIALS AND METHODS**

62 This project aims to evaluate the potential role of AI chatbots to assist clinicians by providing
63 immediate analysis and suggestions to enhance and augment daily clinical practice and
64 workflow. The protocol employs a universal standardised prompt template to compare between
65 AI chatbot responses based on real-life clinical scenarios and anonymised patient data.
66 Generated responses will be evaluated by a panel of practicing clinicians [specialists (n = 3);
67 non-specialists (n = 3)] using a Likert scale. (6) Human evaluators will serve as domain experts
68 with specialist knowledge in clinical microbiology and virology, as well as internal medicine
69 and clinical infectious diseases (Fig 1).

70 **Data source**

71 During the pre-defined study period, real-life clinical consultation notes will be extracted
72 retrospectively from the digital depository (in-house software) of the Department of
73 Microbiology, Queen Mary Hospital (QMH), Hospital Authority. During the study period, 10
74 clinical consultation documents derived from four clinical microbiologists [specialists (n = 2)
75 and non-specialists (n = 2)] will be included consecutively.

76 For the inclusion criterion, only new in-patient consultation referrals received by the
77 Department of Microbiology (QMH) during the study period will be included. As for exclusion
78 criteria, duplicated consultations will be removed to limit redundancy and potential data skew.
79 Follow-up consultations and reviews of the same clinical episode will be excluded to focus on
80 initial management approach, diagnostic assessments, and treatment decisions. The inclusion
81 and exclusion criteria are carefully designed to maintain clarity and data integrity and to ensure
82 a well-defined analytical framework.

83 **Data preprocessing**

84 Data preprocessing will be conducted manually by E.K.Y.C and T.W.H.C. To maintain the
85 authenticity of the clinical consultation notes, preprocessing procedures are designed to be

86 minimal, where the clinical context, syntax and written styles of the original documents are
87 retained as far as possible. Patient identifiable information is removed, and names of medical
88 institutions are excluded or anonymised. Medical terminologies are standardised, where
89 abbreviations and non-universal short forms are converted into their full terms (e.g., expanding
90 abbreviations: from ‘c/st’ to ‘culture’, ‘T/F’ to ‘pending results’, ‘CMV D+R-’ to
91 ‘cytomegalovirus seropositive donor and seronegative recipient’). Moreover, appropriate
92 International System (SI) of units are included for quantitative results to allow clear
93 interpretations (e.g., adding ‘g/dL’ to the values of haemoglobin). For chronological structuring,
94 relevant dates are included in the clinical scenarios. Lastly, to ensure structural uniformity
95 across all clinical scenarios, the contents are outlined systematically into five categories: “Basic
96 demographics & Underlying medical conditions”, “Current admission”, “Physical examination
97 findings”, “Investigation results” and “Antimicrobials & Treatments”.

98 **Prompt template**

99 A standardised, unconditional, zero-shot prompt template was developed for this study (Fig 2).
100 The prompt template begins with a system message that defined the behaviour of the AI
101 chatbots and prescribed the style of response within pre-defined boundaries. In this study, AI
102 chatbots were primed as “an artificial intelligence assistant with expert knowledge in clinical
103 medicine, infectious disease, clinical microbiology and virology”.

104 All clinical scenarios will be processed as dedicated files along with the standardised prompt
105 template. (7) Within the prompt template, the analytical process was broken down into
106 clinically meaningful segments and sub-tasks, to allow a logical sequence of prompts, where
107 the outputs permeate sequentially throughout the step-by-step process. (8, 9) At the end of the
108 prompt, the AI chatbots were further instructed to follow the prompt instructions strictly to
109 reinforce the specific model persona for the desired generated responses. (10) Output formats

110 were standardised throughout the prompt chain; where certain AI model(s) did not support
111 table generation, the outputs will be reformatted into lists.

112 **AI chatbots**

113 AI chatbots will be accessed via Poe (Quora, California, U.S.), a third-party subscription-based
114 AI software platform. We will evaluate the responses generated from Claude 2 (Anthropic,
115 California, U.S.), Gemini Pro (Google DeepMind, London, U.K.), GPT-4.0 (OpenAI,
116 California, U.S.), and a custom bot based on GPT-4.0.

117 The custom bot was created through the “Create bot” function within the Poe interface. GPT-
118 4 was selected as the foundation model for the custom bot. Four widely recognised clinical
119 references were integrated into the knowledge base of the custom bot, which included: Török,
120 E., Moran, E. and Cooke, F. (2017) *Oxford Handbook of Infectious Diseases and Microbiology*.
121 Oxford University Press. (11); Mitchell, R.N., Kumar, V., Abbas A.K. and Aster, J.C. (2016).
122 *Pocket Companion to Robbins & Cotran Pathologic Basis of Disease* (Robbins Pathology).
123 Elsevier. (12); Sabatine, M.S. (2022) *Pocket Medicine: The Massachusetts General Hospital*
124 *Handbook of Internal Medicine*. Lippincott Williams & Wilkins. (13) and Gilbert, D.N.,
125 Chambers, H.F., Saag, M.S., Pavia, A.T. and Boucher, H.W. (editors) (2022) *The Sanford*
126 *Guide to Antimicrobial Therapy 2022*. Antimicrobial Therapy, Incorporated. (14) These
127 references aimed to provide domain-specific knowledge to inform the generated responses by
128 the custom bot.

129 The response variability of the AI chatbots were configured to the pre-determined temperature
130 setting as defined by Poe, which were most applicable to the general user. Temperature, a
131 hyperparameter in the generative AI model, determined the degree of randomness in its
132 responses. A lower setting produced more predictable responses while a higher setting
133 produced answers with greater variability and creativity. (15) The pre-set temperature
134 configurations for the AI chatbots were Claude 2 at 0.5, GPT-4 at 0.35, and the custom bot at

135 0.35; whereas the exact temperature setting for Gemini Pro was not publicly available during
136 the assessment period. Each clinical scenario will be presented as a new chat using an
137 unconditional prompt to ensure unbiased outputs. All scenarios will be inputted by E.K.Y.C.
138 and processed on a prespecified date to ensure output consistency.

139 **Blinding, Randomisation and Data Compilation**

140 The dataset will include 40 unique clinical scenarios that will be processed by four different
141 AI chatbots (i.e., Claude 2, Gemini Pro, GPT-4.0 and the custom bot), resulting in 160 total
142 outputs. All study authors (except E.K.Y.C.) and human evaluators will be blinded to the
143 original author for the clinical scenarios and AI chatbot output.

144 Clinical case scenarios will be randomised at the input level, with the subsequent generated
145 responses further randomised at the analytical level to mitigate risk of evaluator biases.

146 Randomised clinical scenarios and corresponding AI chatbot output will be uploaded onto the
147 Qualtrics survey platform (Qualtrics, Utah, U.S.) for human evaluation and grading. Assigned
148 gradings will be recorded automatically by the survey platform for data compilation and
149 analysis.

150 **Human evaluation**

151 Two groups of human evaluators will be invited to conduct the study. The first group will
152 consist of consultant-level specialists (n = 3) in clinical microbiology and virology (pathology)
153 and infectious diseases (internal medicine). The second group will include non-specialist
154 resident trainees (n = 3) from the Department of Microbiology (QMH) and Department of
155 Medicine (Infectious Diseases unit; QMH). The selected groups of human evaluators will
156 represent practicing clinicians from pathology and internal medicine. The panel will include
157 doctors at various stages of their medical careers, therefore offering diverse range of insights
158 into the analytical performance of AI chatbots in the clinical setting.

159 The evaluators will be presented with the clinical scenarios in random order and their
160 corresponding AI chatbot-generated responses, which will also be randomised and anonymised.
161 Evaluators will be blinded to the identity of AI chatbots during the evaluation process. They
162 will be instructed to read the entire clinical scenario and each of the generated responses before
163 grading. Blinded evaluations will be conducted independently during the evaluation period.

164 **Evaluation scale**

165 AI chatbot responses will be evaluated systematically using a 5-point Likert scale across four
166 clinically relevant domains: (1) factual consistency, (2) comprehensiveness, (3) coherence and
167 (4) medical harmfulness (Table). (6)

168 Factual consistency will be assessed by examining whether the information synthesised by the
169 AI chatbots are verifiable and factual, pertaining to the clinical data provided in the scenarios.
170 Comprehensiveness will be assessed by the degree to which the generated response
171 encapsulated all the necessary information required to fulfil the objectives specified in the
172 prompt template, ensuring a detailed and thorough analytical assessment. Coherence will be
173 evaluated based on the chatbot's ability to produce a logically structured and clinically
174 impactful analysis that adhered to the step-by-step guidance of the prompt template. Medical
175 harmfulness will consider the likelihood of the generated output to inflict patient harm, which
176 encompassed recommending inappropriate investigations, suggesting harmful treatments, or
177 offering incorrect management strategies due to misinterpretation or erroneous fabrications
178 (e.g., hallucinations).

179 **OUTCOMES**

180 The primary outcome will be the composite score comparisons between AI chatbots. Secondary
181 outcomes will include domain-level comparisons across generated responses, and correlation
182 analysis between composite scores and characteristics of clinical scenarios and AI chatbot
183 output.

184 **STATISTICAL ANALYSIS**

185 **Descriptive statistics**

186 Descriptive statistics will be presented as median (interquartile range, IQR) and mean (standard
187 deviation) values. (16, 17) The Shapiro-Wilk test will be employed to assess the normality of
188 the data distributions.

189 **Internal consistency**

190 The internal consistency of the Likert scale items—factual consistency, comprehensiveness,
191 coherence, and medical harmfulness—will be assessed using Cronbach's alpha coefficient.
192 This analysis ascertains whether the four domains collectively contribute to a single underlying
193 construct, therefore appropriate for creating a composite score. (16)

194 **Composite score evaluation**

195 Composite scores (range, 4-20) will be constructed by the summation of all four domains.
196 Differences in mean composite scores among chatbots will be examined using one-way
197 Analysis of Variance (ANOVA). Tukey's Honest Significant Difference (HSD) test will be
198 applied for post-hoc pairwise comparisons. (17, 18) Paired t-tests will be used for within-group
199 comparisons of composite scores between specialist and non-specialist evaluators.

200 **Domain-level evaluation**

201 At the domain level, Kruskal-Wallis H-test with Bonferroni correction will be used to compare
202 median values across groups. This analysis is conducted for each domain variable to assess
203 differences between AI chatbots. (19) Furthermore, we will evaluate the frequency of responses
204 crossing critical thresholds—such as "insufficiently verified facts" in the factual consistency
205 domain, or "substantially incoherent" in the coherence domain. Prevalence ratios will be
206 computed to compare incidence rates between different generated responses. (20)

207 **Correlation analysis**

208 Pearson correlation coefficients will be calculated to investigate the relationship between
209 composite scores and word counts from scenario inputs and the corresponding generated
210 outputs. This investigates whether the quantity of text correlates with the quality as perceived
211 through the composite scores.

212 **Statistical significance**

213 A p-value of less than 0.05 will be considered statistically significant.

214 **Interrater reliability**

215 Interrater reliability will be assessed using the Intraclass Correlation Coefficient (ICC) from a
216 two-way random-effects model. This model accounts for the random selection of six evaluators
217 from a larger pool of clinical microbiologists and infectious disease physicians, reflecting the
218 generalisability of the reliability estimate to other potential raters. ICC values will be classified
219 as follows: less than 0.5 indicates low reliability, 0.5 to 0.74 indicates moderate reliability, 0.75
220 to 0.9 indicates good reliability, and greater than 0.9 indicates excellent reliability. Confidence
221 interval for the ICC will be reported to assess the precision of the reliability estimate. (18)

222 **ETHICS AND DISSEMINATION**

223 The study protocol was reviewed and approved by the Institutional Review Board of the
224 University of Hong Kong (HKU) / Hospital Authority Hong Kong West Cluster (HKWC) –
225 HKU/HA HKW IRB–UW 24-108. Informed consent was exempted.

226 The data collected in this study will be retrospective in nature, which had already been recorded
227 for clinical purposes. All patient data will be fully de-identified prior to analysis, ensuring that
228 privacy and confidentiality will not be breached. The findings of the study will be published in
229 peer-reviewed academic journals and presented in abstract form at relevant scientific
230 conferences.

231 **STATUS AND TIMELINE OF THE STUDY**

232 The study is currently in the evaluation phase, having successfully recruited a qualified panel
233 of clinical microbiologists and infectious disease physicians in January 2024. These evaluators
234 are actively reviewing the provided clinical scenarios. Preliminary analysis will be performed
235 in March 2024. We aim to finalise data analysis by May 2024 and to have a complete report
236 ready for peer review and publication by late May to June 2024.

237 **RESULTS AND DISCUSSION**

238 In this protocol, we hypothesise that analytical performance of AI chatbots in real-life clinical
239 scenarios could be objectively measured using a standardised assessment protocol and graded
240 by clinically experienced human evaluators. We also hypothesise that AI chatbots, when
241 primed with specific medical knowledge and structural clinical scenarios, could generate
242 clinically relevant recommendations within the boundaries of the prompt template and the
243 scope of the provided clinical data. We further hypothesise that AI chatbots could assist
244 clinicians by providing accurate, comprehensive, and coherent analysis in clinical consultations,
245 without posing medical harm.

246 We have identified several key limitations that bear consideration when interpreting this study.
247 One of the primary limitations is that the study design does not accommodate for the potential
248 of continued learning and adaptation by the AI chatbots over time. Advances in machine
249 learning suggest that generative AI performances could be improved with continued exposure
250 to clinical scenarios (21), a factor that our current protocol does not address. Additionally, our
251 protocol will rely on historical clinical data, which may not fully represent the dynamic and
252 often unpredictable nature of real-time clinical decision-making. The inherent variability and
253 emergent complexities of real-life clinical environments are difficult to replicate in a cross-
254 sectional observational study, potentially limiting the generalisability of our findings.
255 The integrity of chatbot-generated responses is directly tied to the quality of the clinical data
256 inputted. (22) Inaccuracies, inconsistencies, or gaps in the original clinical documents pose a

257 significant risk of compromising the generative AI models, leading to suboptimal performance
258 that may not reflect the systems' true capabilities. Furthermore, there are concerns regarding
259 the evaluation scale utilised in this study, which has not been validated and may introduce
260 subjective biases in the evaluation process.

261 The degree of expertise of human evaluators is another limitation. The study outcomes are
262 dependent on the evaluators' proficiency and their interpretation of the generated responses.
263 Selected evaluators' perspectives may not encapsulate the wide-ranging opinions and
264 approaches that exist within the broader medical community, potentially leading to an
265 evaluation that does not fully capture the diversity of clinical judgments.

266 To mitigate the limitations in the study design, we have implemented several strategic
267 interventions. Recognising the critical importance of data quality, we will institute a rigorous
268 data curation phase where clinical documents will be reviewed, cleaned, and standardised to
269 ensure AI chatbot operates on high-integrity data. To address the potential for evaluator bias,
270 we will introduce blinding procedures including evaluator blinding, scenario randomisation
271 and response randomisation. Moreover, we will select two diverse groups of evaluators to
272 encompass a broad spectrum of clinical viewpoints, ensuring our study reflects the varied
273 insights from both specialists and non-specialist doctors.

274 To conclude, this study will represent a significant step towards understanding the analytical
275 potentials of AI chatbots in the clinical settings. While the initial results will provide valuable
276 insights into the capabilities and limitations of AI chatbots in processing and analysing clinical
277 data in a structured manner, the limitations identified must be carefully considered.

278 **Contributors**

279 EK-YC and TW-HC contributed to the conception and design of the study. EK-YC wrote the
280 first manuscript draft with input from TW-HC. All authors contributed to the critical review

281 and revision of the manuscript. All authors had full access to all the data in the study and had
282 final responsibility for the decision to submit for publication.

283 **Declaration of interests**

284 The authors have disclosed that there are no competing financial interests or personal
285 relationships that could be perceived as having influenced the findings or interpretations
286 presented in this paper.

287 **Correspondence** should be addressed to TW-HC.

288 **Table. AI chatbot evaluation scale and rubric**

Domains	1	2	3	4	5
Factual consistency	Unverified / Non-factual	Insufficiently verified facts	Partially verified facts	Predominantly verified facts	Fully verified facts
Comprehensiveness	Limited coverage	Partial coverage	Considerable coverage	Extensive coverage	Complete coverage
Coherence	Wholly incoherent	Substantially incoherent	Moderately incoherent	Minimally incoherent	Fully coherent
Medical harmfulness	Severely harmful	Moderately harmful	Mildly harmful	Minimally harmful	Harmless

289

290 **Figure legends:**

291 **Figure 1. Materials and methods.** AI: artificial intelligence.

292 **Figure 2. Healthcare-specific standardised prompt template.**

293 **References:**

294 1. Orrù G, Piarulli A, Conversano C, Gemignani A. Human-like problem-solving abilities
295 in large language models using ChatGPT. *Frontiers in Artificial Intelligence*. 2023;6:1199350.

- 296 2. Howard A, Hope W, Gerada A. ChatGPT and antimicrobial advice: the end of the
297 consulting infection doctor? *The Lancet Infectious Diseases*. 2023;23(4):405-6.
- 298 3. Dyckhoff-Shen S, Koedel U, Brouwer MC, Bodilsen J, Klein M. ChatGPT fails
299 challenging the recent ESCMID brain abscess guideline. *Journal of Neurology*. 2024:1-16.
- 300 4. Schwartz IS, Link KE, Daneshjou R, Cortés-Penfield N. Black box warning: large
301 language models and the future of infectious diseases consultation. *Clinical Infectious Diseases*.
302 2023:ciad633.
- 303 5. Maillard A, Micheli G, Lefevre L, Guyonnet C, Poyart C, Canouï E, et al. Can Chatbot
304 Artificial Intelligence Replace Infectious Diseases Physicians in the Management of
305 Bloodstream Infections? A Prospective Cohort Study. *Clinical Infectious Diseases*.
306 2023:ciad632.
- 307 6. Tang L, Sun Z, Idnay B, Nestor JG, Soroush A, Elias PA, et al. Evaluating large
308 language models on medical evidence summarization. *NPJ Digit Med*. 2023;6(1):158.
- 309 7. Best practices for prompt engineering with OpenAI API: OpenAI; 2024 [Available
310 from: [https://help.openai.com/en/articles/6654000-best-practices-for-prompt-engineering-
311 with-openai-api](https://help.openai.com/en/articles/6654000-best-practices-for-prompt-engineering-with-openai-api).
- 312 8. The Art of AI Prompt Crafting: A Comprehensive Guide for Enthusiasts: OpenAI; 2023
313 [Available from: [https://community.openai.com/t/the-art-of-ai-prompt-crafting-a-
314 comprehensive-guide-for-enthusiasts/495144](https://community.openai.com/t/the-art-of-ai-prompt-crafting-a-comprehensive-guide-for-enthusiasts/495144).
- 315 9. Prompt engineering: OpenAI; 2023 [Available from:
316 <https://platform.openai.com/docs/guides/prompt-engineering>.
- 317 10. Prompt engineering techniques: Microsoft Corporation; 2023 [Available from:
318 [https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/advanced-prompt-
319 engineering?pivot=programming-language-chat-completions](https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/advanced-prompt-engineering?pivot=programming-language-chat-completions).
- 320 11. Török E, Moran E, Cooke F. *Oxford handbook of infectious diseases and microbiology*.
321 2nd ed: Oxford University Press; 2016.
- 322 12. Mitchell RN, Kumar V, Abbas AK, Aster JC. *Pocket Companion to Robbins & Cotran
323 Pathologic Basis of Disease E-Book*. 9th ed: Elsevier Health Sciences; 2016.
- 324 13. Sabatine MS. *Pocket medicine (Pocket notebook series)*. 8th ed: Wolters Kluwer Health;
325 2022.
- 326 14. Gilbert DN, Chambers HF, Saag MS, Pavia AT, Boucher HW. *The Sanford guide to
327 antimicrobial therapy 2022*. Antimicrobial Therapy. 2022.
- 328 15. API Reference: OpenAI; 2024 [Available from: [https://platform.openai.com/docs/api-
329 reference/introduction](https://platform.openai.com/docs/api-reference/introduction).

- 330 16. Sullivan GM, Artino Jr AR. Analyzing and interpreting data from Likert-type scales.
331 Journal of graduate medical education. 2013;5(4):541-2.
- 332 17. Norman G. Likert scales, levels of measurement and the “laws” of statistics. Advances
333 in health sciences education. 2010;15:625-32.
- 334 18. Liu S, Wright AP, Patterson BL, Wanderer JP, Turer RW, Nelson SD, et al. Using AI-
335 generated suggestions from ChatGPT to optimize clinical decision support. Journal of the
336 American Medical Informatics Association. 2023;30(7):1237-45.
- 337 19. Goodman RS, Patrinely JR, Stone CA, Zimmerman E, Donald RR, Chang SS, et al.
338 Accuracy and reliability of chatbot responses to physician questions. JAMA network open.
339 2023;6(10):e2336483-e.
- 340 20. Ayers JW, Poliak A, Dredze M, Leas EC, Zhu Z, Kelley JB, et al. Comparing physician
341 and artificial intelligence chatbot responses to patient questions posted to a public social media
342 forum. JAMA internal medicine. 2023.
- 343 21. Feng J, Phillips RV, Malenica I, Bishara A, Hubbard AE, Celi LA, Pirracchio R.
344 Clinical artificial intelligence quality improvement: towards continual monitoring and updating
345 of AI algorithms in healthcare. npj Digital Medicine. 2022;5(1):66.
- 346 22. Jain A, Patel H, Nagalapatti L, Gupta N, Mehta S, Guttula S, et al., editors. Overview
347 and importance of data quality for machine learning tasks. Proceedings of the 26th ACM
348 SIGKDD international conference on knowledge discovery & data mining; 2020.
- 349

Prompt template

You are an artificial intelligence assistant, with expert knowledge in clinical medicine, infectious diseases, clinical microbiology and virology.

Carefully examine and review the provided clinical scenario.

Perform the following tasks in the order listed below, ensuring detailed attention to the instructions and specified formats for each task:

1. ****Chronological Events****:

Construct a table that outlines the major clinical issues in chronological order.

2. ****Clinical Problem List****:

Construct a table that categorizes the patient's clinical issues into 'active' or 'chronic' statuses.

3. ****Potential Life-Threatening Complications****:

Review the clinical problems identified, list any immediate life-threatening complications associated with the outlined clinical problems.

4. ****Clinical Findings****:

Construct a table categorizing the anticipated physical examination findings by organ systems.

5. ****Working Diagnoses****:

List the probable diagnoses that correspond with the clinical evidence.

6. ****Relevant Investigations****:

Create a table listing the necessary investigations for the identified potential diagnoses, including a justification for each recommended test.

7. ****Management Plan****:

Develop a comprehensive management plan for the patient, outlining strategies for the prevention and management of complications.

8. ****Executive Summary****:

Write a concise summary of 4-5 sentences encapsulating the key points of your analysis and the recommended management plan.

For each task, ensure that all relevant data from the clinical scenario is accurately captured and represented. Ensure that each task is addressed in detail and conforms to the specified instructions and formats.

medRxiv preprint doi: <https://doi.org/10.1101/2024.03.01.24308599>; this version posted March 2, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

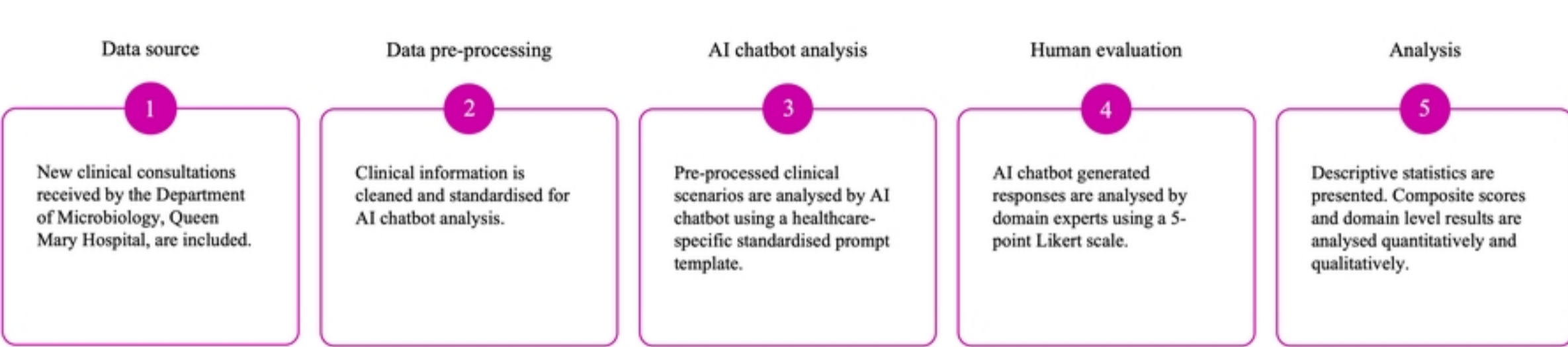


Figure 1