



1 [DATA ARTICLE TEMPLATE V.19 (DECEMBER 2024)]

2 ARTICLE INFORMATION

3 Article title

4 RAGCare-QA: A Benchmark Dataset for Evaluating Retrieval-Augmented Generation Pipelines in
5 Theoretical Medical Knowledge

6 Authors

7 Jovana Dobрева^{a,*}, Ivana Karasmanakis^b, Filip Ivanisevic^b, Tadej Horvat^b, Dimitar Kitanovski^a, Matjaz
8 Gams^b, Kostadin Mishev^a, Monika Simjanoska Misheva^a

9 Affiliations

10 ^aFaculty of Computer Science and Engineering, Ss. Cyril and Methodius University, Skopje, North
11 Macedonia

12 ^bDepartment of Intelligent Systems, Jozef Stefan Institute, Ljubljana, Slovenia

13 Corresponding author's email address and Twitter handle

14 jovana.dobрева@finki.ukim.mk

15 Keywords

16 Medical Education, Knowledge Assessment; Retrieval-Augmented Generation; Multiple Choice
17 Questions; Medical Knowledge Base; Healthcare AI; Theoretical Medicine

18

19 Abstract

20 The paper introduces RAGCare-QA, an extensive dataset of 420 theoretical medical knowledge
21 questions for assessing Retrieval-Augmented Generation (RAG) pipelines in medical education and
22 evaluation settings. The dataset includes one-choice-only questions from six medical specialties
23 (Cardiology, Endocrinology, Gastroenterology, Family Medicine, Oncology, and Neurology) with three
24 levels of complexity (Basic, Intermediate, and Advanced). Each question is accompanied by the best
25 fit of RAG implementation complexity level, such as Basic RAG (315 questions, 75.0%), Multi-vector
26 RAG (82 questions, 19.5%), and Graph-enhanced RAG (23 questions, 5.5%). The questions emphasize
27 theoretical medical knowledge on fundamental concepts, pathophysiology, diagnostic criteria, and
28 treatment principles important in medical education. The dataset is a useful tool for the assessment
29 of RAG-based medical education systems, allowing researchers to fine-tune retrieval methods for
30 various categories of theoretical medical knowledge questions.

31 SPECIFICATIONS TABLE

32

Subject	Health Sciences, Medical Sciences & Pharmacology
---------	--

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.



Specific subject area	Theoretical medical knowledge assessment with RAG pipelines in educational contexts
Type of data	Data Types: JSON Data Format: One-choice-only questions with annotations
Data collection	The RAGCare-QA dataset was compiled from medical text-books, educational materials, and clinical guidelines primarily from European medical literature. Questions were formulated to test theoretical medical knowledge across different complexity levels. Each question was analyzed by medical education experts to determine optimal RAG pipeline suitability.
Data source location	Institution: Faculty of Computer Science and Engineering, Ss. Cyril and Methodius University City/Town/Region: Skopje, Balkan Region Country: North Macedonia Primary Data Source: European Medical Textbooks and Educational Literature
Data accessibility	Repository name: RAGCare-QA Data identification number: DOI to be assigned Direct URL to data: https://huggingface.co/datasets/ChatMED-Project/RAGCare-QA
Related research article	<i>None</i>

33

34

35 VALUE OF THE DATA

- 36 - RAGCare-QA dataset is designed to benchmark state-of-the-art RAG architectures
- 37 recommendations for theoretical medical knowledge through 420 human annotated single-
- 38 choice questions, well-distributed in 6 different medical specialties.
- 39 - Researchers can leverage this resource to build more effective educational tools that adapt
- 40 their retrieval strategies based on question complexity and medical specialty.
- 41 - The dataset fills a gap in medical AI by providing a standardized benchmark that supports the
- 42 development of AI-based adaptive educational tools.
- 43 - The dataset classifies each question by the most suitable RAG architecture, Basic, Multi-
- 44 vector, or Graph-enhanced, needed for context retrieval, enabling precise performance
- 45 comparisons across retrieval strategies.



- 46 - The dataset can serve as a foundation for development of specialized retrieval strategies to
47 enhance learning outcomes in medical education.

48 BACKGROUND

49 The integration of artificial intelligence (AI) in medical education has gained significant momentum,
50 with RAG systems showing particular promise for knowledge assessment and educational content
51 delivery [1, 2]. While large language models (LLMs) demonstrate substantial medical knowledge,
52 their performance in educational contexts is significantly enhanced when combined with specialized
53 retrieval systems that access curated medical educational content [3, 4].

54 Theoretical medical knowledge assessment has particular requirements distinguishing it from clinical
55 problem-solving exercise. Didactic questions require precise recall of bottom-line ideas,
56 pathophysiologic processes, and traditional medical principles from structured sets of knowledge [5,
57 6]. The depth of medical theoretical knowledge, from straightforward definitions to intricate
58 pathophysiological correlations, requires sophisticated retrieval mechanisms that can be calibrated
59 for differing levels of cognition.

60 Existing medical QA datasets, including PubMedQA [5], MedMCQA [6], and specialized disease-
61 focused collections [7], primarily emphasizing clinical decision-making rather than systematic
62 evaluation of retrieval pipelines for educational content. These datasets evaluate model
63 performance against established medical knowledge, however, do not address how different RAG
64 pipelines influence learning and assessment outcomes in educational settings. The landscape of RAG
65 pipelines offers multiple approaches, each with distinct advantages for educational applications.

66 Basic RAG implementations provide straightforward document retrieval suitable for direct factual
67 queries common in foundational medical education [8]. Multi-vector RAG models show promise in
68 managing various types of educational content, ranging from dictionary-style definitions to elaborate
69 explanations, making them a good fit for inclusive medical education applications [9, 10]. Graph-
70 augmented RAG systems excel at representing hierarchical medical knowledge structures and
71 complex concept relationships essential for advanced theoretical understanding [11, 12]. These
72 pipeline designs have shown to be useful in educational settings where information must be
73 accessed from several angles and sources of knowledge.

74 There is a notable gap in the medical AI community regarding systematic approaches to choosing
75 suitable RAG pipelines for different categories of theoretical medical knowledge questions, which
76 often leads to suboptimal design of educational systems. Although recent advances, such as
77 HuatuoGPT

78 [13] and other education-oriented AI models, have contributed to medical education, they have
79 largely prioritized architectural improvements over the optimization of retrieval strategies for
80 instructional content.

81 This dataset bridges this important gap by offering a systematically annotated set of theoretical
82 medical questions with specific RAG pipeline recommendations, allowing evidence-based retrieval
83 strategy selection for medical education use cases.

84



85 DATA DESCRIPTION

86 Dataset Structure and Composition

87 The RAGCare-QA dataset comprises 420 theoretical medical knowledge questions systematically
88 distributed across medical specialties, complexity levels, and RAG implementation categories. Table 1
89 provides a detailed breakdown of the dataset structure.

90 Table 1: Dataset Breakdown by Medical Specialty and Complexity Levels.

Medical Specialty	Basic	Intermediate	Advanced	Total
Cardiology	32	36	19	87
Endocrinology	28	32	14	74
Family Medicine	26	41	14	81
Gastroenterology	20	22	12	54
Neurology	20	22	15	57
Oncology	24	28	15	67
Total	150	181	89	420

91

92 Complexity Level Distribution:

- 93 • Basic (150 questions): Fundamental medical concepts, definitions, and straightforward factual
94 knowledge.
- 95 • Intermediate (181 questions): Moderate complexity questions in-
96 volving pathophysiology, diagnostic criteria, and treatment principles.
- 97 • Advanced (89 questions): Complex theoretical scenarios requiring deep understanding of medical
98 mechanisms and advanced concepts.

99 RAG Implementation Complexity Distribution:

- 100 • Basic RAG: 315 questions (75.0%) - Direct factual queries with explicit terminology and
101 straightforward retrieval requirements.
- 102 • Multi-vector RAG: 82 questions (19.5%) - Questions requiring diverse knowledge sources and
103 multiple representation approaches.
- 104 • Graph-enhanced RAG: 23 questions (5.5%) - Complex relationship-
105 based queries requiring structured knowledge representation.

106 Each entry in the dataset follows a structured format with the following fields:

- 107 • Type: Medical specialty classification.
- 108 • Question: Multiple-choice question with options (a-e).
- 109 • Answer: Correct answer designation (a, b, c, d, or e).
- 110 • Text Answer: Correct answer in textual format



- 111 • Reference: Citation of the medical literature source.
- 112 • Page: Specific page reference within the source.
- 113 • Context: Relevant text from source that supports the answer.
- 114 • Label: Optimal RAG pipeline classification.
- 115 • Complexity: Difficulty level based on the depth of medical knowledge and clinical reasoning
- 116 required to answer the question

117 Dataset Statistics and Characteristics

118 The RAGCare-QA dataset exhibits diverse features across question complexity and content structure.

119 Table 2 provides detailed statistical analysis of the dataset composition.

120 Table 2: RAGCare-QA Statistical Features.

Feature	Minimum	Average	Maximum
Question Length (words)	12	47	112
Answer Options	5	5	5
Context Length (words)	25	89	245
Publication Years	1985	2015	2024
Reference Types: Medical Textbooks (63%), Peer-reviewed Journals (30%), Other Sources (7%)			

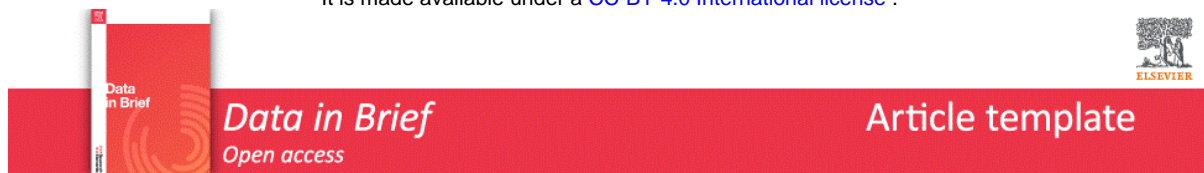
121

122 **Answer Distribution Analysis:** The correct answers are well- distributed across all five options,
123 ensuring balanced assessment: Option A (22%), Option B (20%), Option C (19%), Option D (21%),
124 Option E (18%). This distribution prevents systematic bias toward specific answer positions. Source
125 Reference Analysis: Reference Source Distribution: The dataset incorporates a balanced mix of
126 authoritative medical sources: medical textbooks (63%), peer-reviewed journal articles (30%), and
127 other specialized medical resources (7%). The primary sources include "Interna medicina" (6th
128 edition, 2022) [14] and "Harrison's Principles of Internal Medicine" [15] for foundational knowledge,
129 supplemented by high-impact journal publications. Publication years span from 1985 to 2024,
130 ensuring coverage of both foundational medical knowledge and current research developments.

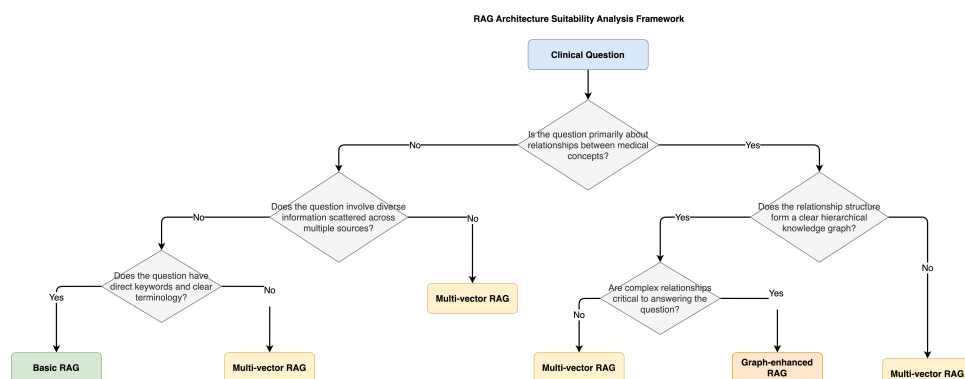
131 Complexity Progression: Questions demonstrate clear complexity escalation from Basic (average 32
132 words) to Intermediate (average 48 words) to Advanced (average 65 words), reflecting increasing
133 cognitive load and knowledge integration requirements

134 Question Examples by RAG Pipeline Type

135 The dataset encompasses various types of theoretical medical knowledge questions, each designed
136 to test specific aspects of medical understanding, with optimal RAG pipeline assignments based on
137 the framework shown in Figure 1.



138



139

140 Figure 1: RAG Pipeline Classification Decision Tree. The figure presents the systematic decision-making
141 framework used to classify theoretical medical questions into their optimal RAG pipeline categories. The
142 decision tree evaluates key characteristics such as relationship complexity between medical concepts,
143 information distribution patterns, query structure complexity, and reasoning requirements. Starting from the
144 initial question analysis, the framework guides classification into Basic RAG (for direct factual queries), Multi-
145 vector RAG (for questions requiring diverse knowledge sources), or Graph-enhanced RAG (for complex
146 relationship-based queries). This evidence-based classification ensures optimal matching between question
147 types and retrieval pipelines for medical education applications.

148 Complexity Level Characteristics

149 The questions in the dataset are categorized into three complexity levels: Basic, Intermediate, and
150 Advanced based on the depth of medical knowledge required, the cognitive effort involved, and the
151 nature of information integration needed for accurate resolution:

- 152 • Basic Level Questions focus on fundamental definitions, basic pathophysiology, and direct factual
153 knowledge. These questions typically involve single-concept retrieval and straightforward medical
154 terminology recognition.
- 155 • Intermediate Level Questions involve moderate complexity scenarios requiring understanding of
156 disease mechanisms, diagnostic approaches, and treatment principles. These questions often require
157 integration of multiple medical concepts and represent the largest portion of the dataset.
- 158 • Advanced Level Questions present complex theoretical scenarios demanding deep understanding
159 of pathophysiological mechanisms, differential diagnosis considerations, and advanced medical
160 principles. These questions often involve sophisticated medical reasoning and comprehensive
161 knowledge integration.

162 RAG Pipeline Classification Framework

163 Each question underwent systematic analysis to determine its optimal RAG pipeline suitability using
164 a structured evaluation framework. The decision-making process follows a systematic classification
165 tree as illustrated in Figure 1, while the resulting pipeline types are compared in Figure 2. The
166 classification framework evaluates several key factors:

167 Basic RAG Suitability (78.75% of questions): Questions with explicit terminology, direct factual
168 content, and straightforward retrieval requirements. These typically involve definition-based queries



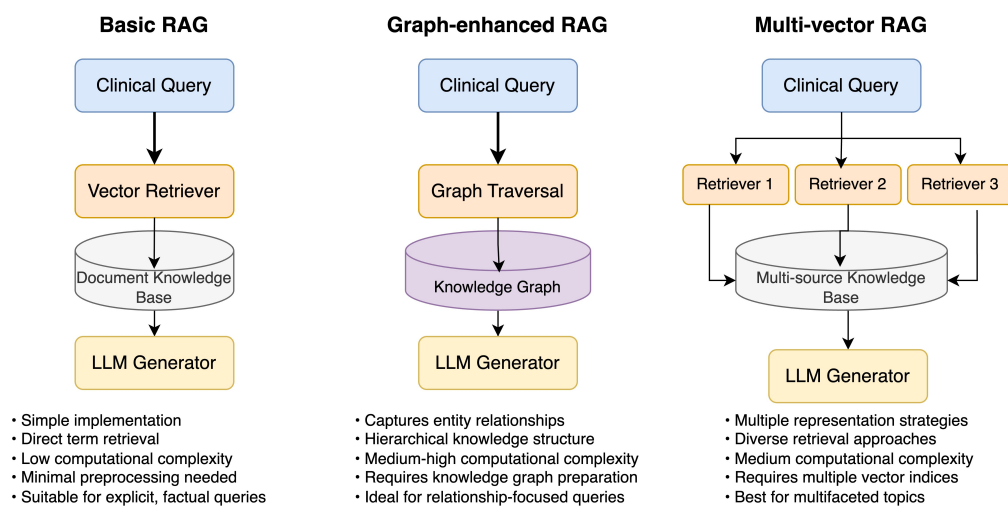
169 or simple factual recall that can be effectively answered through standard document retrieval. Multi-
 170 vector RAG Suitability (20.5% of questions): Questions requiring diverse knowledge sources,
 171 multiple perspectives on medical concepts, or integration of information from various medical
 172 domains. These questions benefit from multiple representation strategies and comprehensive
 173 knowledge coverage.

174 Graph-enhanced RAG Suitability (5.5% of questions): Questions involving complex medical
 175 relationships, hierarchical knowledge structures, and entity interconnections that benefit from
 176 graph-based knowledge representation. These represent sophisticated theoretical queries requiring
 177 advanced reasoning capabilities.

178 By systematically considering information complexity, retrieval needs, domain breadth, and cognitive
 179 demands, the classification process ensures each question is aligned with the most suitable RAG
 180 pipeline to maximize performance in educational applications.

Comparison of RAG Pipeline Types

Pipeline variants optimized for different clinical question types



Key Characteristics Comparison:

Characteristic	Basic RAG	Graph-enhanced RAG	Multi-vector RAG
Knowledge Representation	Flat document chunks	Structured entity relationships	Multiple representations
Optimal for Question Types	Direct factual queries	Relationship-focused queries	Multifaceted topics
Implementation Complexity	Low	Medium-High	Medium

The above comparison illustrates how different RAG pipelines process information flows and handle various clinical query types.

181
 182
 183 Figure 2: Comparison of RAG Pipeline Types. The figure illustrates the key differences between Basic
 184 RAG (simple document retrieval), Graph-enhanced RAG (structured knowledge representation),
 185 and Multi-vector RAG (diverse representation approaches). For each pipeline type, the diagram
 186 shows information flow patterns, typical use cases, and relative computational complexity.
 187 The comparison table highlights how different medical question types benefit from specific
 188 design approaches optimized for theoretical medical knowledge retrieval.

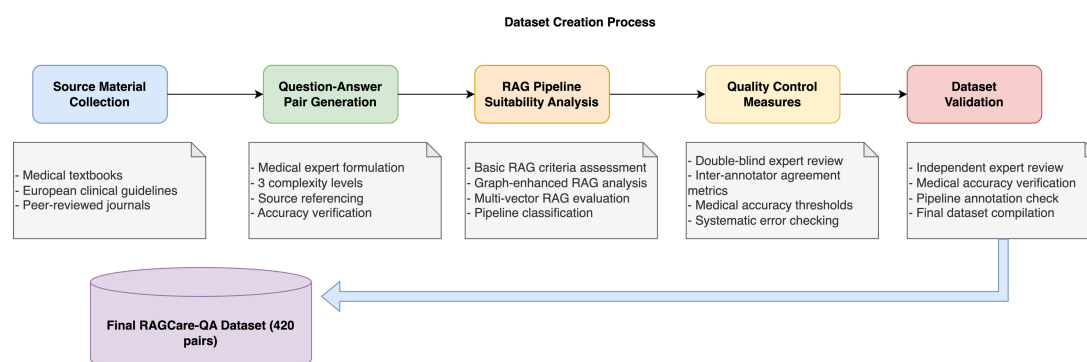


189 EXPERIMENTAL DESIGN, MATERIALS AND METHODS

190 Dataset Creation Process

191 The RAGCare-QA dataset was developed through a systematic multi-stage process specifically
192 designed for theoretical medical knowledge assessment, as illustrated in Figure 3. The creation
193 process began with comprehensive collection of source materials from authoritative European
194 medical textbooks, educational curricula, and clinical guidelines commonly used in medical
195 education programs.

196 The source materials were systematically collected from authoritative medical references across
197 multiple licensing categories to ensure comprehensive coverage and copyright compliance. The
198 primary source was "Interna medicina" (6th edition, 2022), editors Mitja Košnik, Dušan Štajer, and
199 colleagues, published by Medicinska fakulteta (University of Ljubljana Medical Faculty), Slovensko
200 zdravniško društvo (Slovenian Medical Association), and Buča [14], contributing 180 questions (42.9%
201 of the dataset). Additional sources included "Onkologija: Učbenik za študente medicine" (1st edition,
202 2018) by Strojan & Hočevár, published by Onkološki inštitut Ljubljana [16], contributing 44 questions
203 (10.5%) and licensed under Creative Commons Attribution-NonCommercial-NoDerivatives 4.0
204 International. "Harrison's Principles of Internal Medicine" (20th edition) [15] contributed 19
205 questions (4.5%) through transformed contextual synthesis. Open access educational resources from
206 StatPearls and NCBI Bookshelf [17] provided 56 questions (13.3%), also licensed under Creative
207 Commons Attribution-NonCommercial-NoDerivatives 4.0 International. Peer-reviewed journal
208 articles from high-impact publications contributed 39 questions (9.3%), and the remaining 82
209 questions (19.5%) were derived from additional educational materials including medical curricula
210 and clinical practice guidelines. This multi-tiered approach ensures comprehensive coverage of
211 established medical knowledge across six targeted specialties while maintaining rigorous copyright
212 compliance standards.



213

214 Figure 3: Dataset Creation Process. The figure illustrates the systematic five-stage approach
215 used to develop the RAGCare-QA dataset. The process begins with source material collection
216 from European medical literature, followed by question-answer pair generation by medical
217 experts across six specialties and three complexity levels. Each pair then undergoes RAG
218 pipeline suitability analysis using specific criteria for Basic RAG, Graph-enhanced RAG, and Multi-
219 vector RAG approaches. Quality control measures ensure medical accuracy and appropriate
220 complexity classification. The process concludes with dataset validation by independent
221 experts. The final dataset contains 420 theoretical medical knowledge questions with
222 comprehensive annotations and source references.



223 Medical education experts and subject matter specialists formulated questions across six medical
224 specialties, ensuring coverage of fundamental theoretical concepts essential for medical training.
225 Questions were designed as multiple-choice items with five options (a-e), following standard medical
226 education assessment formats. Each question targets specific learning objectives and cognitive levels
227 appropriate for medical students and healthcare professionals.

228 The question development process prioritized theoretical knowledge over clinical case-based
229 scenarios, focusing on pathophysiology, disease mechanisms, diagnostic criteria, pharmacological
230 principles, and anatomical knowledge. This approach ensures the dataset serves as a comprehensive
231 resource for foundational medical education rather than clinical decision-making training.

232 For each question, detailed contextual information was extracted from source materials, including
233 specific page references and relevant text passages that support the correct answer. This
234 contextualization enables effective training and evaluation of RAG pipelines by providing rich source
235 material for retrieval processes.

236 RAG Pipeline Analysis and Annotation

237 Each question underwent systematic analysis to determine its optimal RAG pipeline using a
238 structured evaluation framework. The analysis considered multiple factors including information
239 structure, retrieval complexity, knowledge domain requirements, and cognitive processing needs.

240 **Basic RAG Classification:** Questions were classified as Basic RAG suitable when they involved direct
241 factual retrieval, explicit terminology, clear question structure, and straightforward answer pathways.
242 These questions typically require simple document-based retrieval without complex relation- ship
243 processing and represent the majority of the dataset (78.75%).

244 **Multi-vector RAG Classification:** Questions requiring diverse information sources, multiple conceptual
245 perspectives, cross-domain knowledge integration, or comprehensive coverage of medical topics
246 were classified as Multi-vector RAG suitable (20.5%). These questions benefit from multiple retrieval
247 strategies and diverse representation approaches.

248 **Graph-enhanced RAG Classification:** A meaningful subset of questions (5.5%) requiring complex
249 relationship modeling, hierarchical knowledge representation, and sophisticated entity
250 interconnection analysis were classified as Graph-enhanced RAG suitable. These questions involve
251 the most complex theoretical reasoning scenarios requiring advanced graph-based retrieval
252 approaches.

253 The annotation process involved medical education experts working in conjunction with AI specialists
254 to ensure both medical accuracy and appropriate technical classification using the decision
255 framework shown in Figure 1. Inter-rater reliability was maintained through systematic review
256 processes and consensus-building approaches.

257 Quality Assurance and Validation

258 The dataset underwent comprehensive quality assurance to ensure medical accuracy, appropriate
259 difficulty progression, and correct RAG pipeline classification. Medical content was validated against
260 authoritative sources, with particular attention to European medical practice standards and edu-
261 cational requirements.



262 Each question's difficulty level was validated through expert review, ensuring appropriate
263 classification into Basic, Intermediate, and Advanced categories. The progression from basic factual
264 knowledge to complex theoretical understanding reflects authentic medical education pathways.

265 RAG pipeline annotations were validated through cross-verification processes, where three experts
266 independently assessed question suitability for different retrieval approaches. Disagreements were
267 resolved through consensus meetings and detailed discussion of classification criteria. The final
268 dataset underwent validation by independent medical and AI experts to ensure accuracy of medical
269 content and appropriate RAG pipeline annotations, as shown in the final stage of Figure 3.

270 LIMITATIONS

271 The RAGCare-QA dataset has several limitations that researchers should consider. The content
272 primarily reflects European medical education standards and may not fully represent global medical
273 education approaches or regional variations in medical practice. The dataset focuses on six major
274 medical specialties, potentially limiting applicability to other medical do- mains such as surgery,
275 pediatrics, or psychiatry.

276 The multiple-choice format, while standard in medical education, may not capture all forms of
277 theoretical medical knowledge assessment used in educational settings. The dataset's theoretical
278 focus excludes practical clinical skills, procedural knowledge, and patient interaction scenarios that
279 form important components of comprehensive medical education.

280 The complexity level classifications, while expert-validated, may not align perfectly with all
281 educational frameworks or institutional standards. The RAG implementation complexity shows a
282 reasonable balance among Ba- sic RAG (75.0%) for straightforward queries, Multi-vector RAG (19.5%)
283 for complex knowledge integration, and Graph-enhanced RAG (5.5%) for relationship-based
284 reasoning, which may provide sufficient diversity for com- prehensive retrieval architecture
285 evaluation.

286 Language limitations exist as the dataset is primarily in English with some source materials in
287 Slovenian, potentially affecting applicability in multilingual educational contexts. The static nature of
288 the knowledge base may require periodic updates to maintain currency with evolving medical under-
289 standing and educational standards.

290 ETHICS STATEMENT

291 This research complies with ethical publication guidelines and institutional review requirements. The
292 dataset construction involved no human subjects research, animal experimentation, or collection of
293 sensitive personal data. All source materials were derived from publicly available medical literature,
294 educational resources, and established clinical guidelines.

295 **Source Material Usage and Copyright Compliance:** The dataset creation process employed
296 systematic content transformation methodologies to ensure copyright compliance while maintaining
297 educational research integrity. Question development involved creating original multiple-choice
298 assessments based on established medical knowledge, with all contextual information paraphrased
299 and synthesized from source materials rather than reproduced verbatim.



300 For materials licensed under Creative Commons Attribution-NonCommercial-NoDerivatives 4.0
301 International, including "Onkologija: Učbenik za študente medicine" [16] (44 questions, 10.5%) and
302 StatPearls/NCBI Bookshelf resources [17] (56 questions, 13.3%), educational use is explicitly
303 permitted with proper attribution. For "Interna medicina" [14] (180 questions, 42.9%), published by
304 University of Ljubljana Medical Faculty in collaboration with Slovenian Medical Association, question
305 development followed academic fair use principles for educational research involving university-
306 published materials. For "Harrison's Principles of Internal Medicine" [15] (19 questions, 4.5%),
307 contextual information was substantially paraphrased and synthesized to present medical concepts
308 in original formulations, preserving only bibliographic references and page citations for academic
309 attribution purposes.

310 The systematic content paraphrasing process involved medical experts reformulating all medical
311 concepts into novel assessment formats without reproducing any substantial portions of original
312 textual content from any source. Only essential bibliographic information (references and page
313 numbers) was preserved to maintain academic integrity and enable verification of medical accuracy.
314 This comprehensive paraphrasing methodology represents original scholarly work that enhances
315 medical education research while respecting intellectual property rights through complete textual
316 transformation of all source materials.

317 No proprietary or confidential medical information was included in the dataset. The questions and
318 answers represent established medical knowledge reformulated into original assessment items and
319 do not include experimental or unvalidated medical information.

320 CRediT AUTHOR STATEMENT

321 Jovana Dobrova: Conceptualization, Data curation, Investigation, Methodology, Writing - Original
322 draft, Writing - Review & editing. Ivana Karasmanakis: Medical validation, Resources. Filip Ivanisevic:
323 Medical validation, Resources. Tadej Horvat: Data curation, Formal analysis, Software, Validation. &
324 editing. Dimitar Kitanovski: Investigation, Resources, Validation. & editing. Matjaz Gams: Supervision.
325 Kostadin Mishev: Supervision, Formal analysis, Methodology, Writing - Review & editing. Monika
326 Simjanoska Misheva: Conceptualization, Project administration, Supervision, Writing - Review &
327 editing.

328 ACKNOWLEDGEMENTS

329 Views and opinions expressed are, however, those of the author(s) only and do not necessarily
330 reflect those of the European Union or the European Research Executive Agency. Neither the
331 European Union nor the granting authority can be held responsible for them.

332 Funded by the European Union under Horizon Europe project ChatMED - Bridging Research
333 Institutions to Catalyze Generative AI Adoption by the Health Sector in the Widening Countries (grant
334 agreement ID: 101159214).

335 DECLARATION OF COMPETING INTERESTS

336 The authors declare no competing financial interests or personal relationships that could influence
337 the work reported in this paper.



338 REFERENCES

- 339 [1] K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-
340 Lewis, S. Pfohl, et al., Large language models encode clinical knowledge, *Nature* 620 (7972) (2023)
341 172–180. doi:10.1038/s41586-023-06291-2.
- 342 [2] K. Singhal, T. Tu, J. Gottweis, R. Sayres, E. Wulczyn, L. Hou, K. Clark, S. Pfohl, H. Cole-Lewis, D.
343 Neal, et al., Towards expert-level medical question answering with large language models, *arXiv*
344 preprint arXiv:2305.09617 (2023). doi:10.48550/arXiv.2305.09617.
- 345 [3] C. Zakka, A. Chaurasia, R. Shad, A. Karimpour, J. Kim, S. Kashyap, et al., Retrieval-augmented
346 generation for generative artificial intelligence in health care, *npj Health Systems* 1 (1) (2024) 4.
347 doi:10.1038/s44401-024-00004-1.
- 348 [4] A. B. Abacha, S. Gayen, J. J. Lau, S. Thomas, D. Demner-Fushman, Iryonlp at mediq-corr
349 2024: Tackling the medical error detection & correction task on the shoulders of medical agents, in:
350 Proceedings of the 5th Clinical Natural Language Processing Workshop, 2023, pp. 353–
351 361. doi:10.18653/v1/2023.clinicalnlp-1.36.
- 352 [5] Q. Jin, B. Dhingra, Z. Liu, W. Cohen, X. Lu, Pubmedqa: A dataset for biomedical research
353 question answering, in: Proceedings of the 2019 Conference on Empirical Methods in Natural
354 Language Processing and the 9th International Joint Conference on Natural Language Processing
355 (EMNLP-IJCNLP), 2019, pp. 2567–2577. doi:10.18653/v1/D19-1259.
- 356 [6] A. Pal, L. K. Umapathi, M. Sankarasubbu, Medmcqa: A large-scale multi-subject multi-choice
357 dataset for medical domain question answering, in: Proceedings of the Conference on Health,
358 Inference, and Learning, PMLR, 2022, pp. 248–260. doi:10.48550/arXiv.2203.14371.
- 359 [7] D. Jin, E. Pan, N. Oufattole, W.-H. Weng, H. Fang, P. Szolovits, What disease does this patient
360 have? a large-scale open domain question answering dataset from medical exams, *Applied Sciences*
361 11 (14) (2021) 6421. doi:10.48550/arXiv.2009.13081.
- 362 [8] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Kuttler, M. Lewis, W.-t. Yih, T.
363 Rocktaschel, et al., Retrieval-augmented generation for knowledge-intensive nlp tasks, in: *Advances*
364 *in Neural Information Processing Systems*, Vol. 33, 2020, pp. 9459–9474.
365 doi:10.48550/arXiv.2005.11401.
- 366 [9] R. Wu, S. Chen, X. Su, Y. Zhu, Y. Liao, J. Wu, A multisource retrieval question answering
367 framework based on rag, *arXiv preprint arXiv:2405.19207* (2024). doi:10.48550/arXiv.2405.19207.
- 368 [10] H. Wang, W. Huang, Y. Deng, R. Wang, Z. Wang, Y. Wang, F. Mi, J. Z. Pan, K.-F. Wong, Unims-
369 rag: A unified multi-source retrieval-augmented generation for personalized dialogue systems, *arXiv*
370 preprint arXiv:2401.13256 (2024). doi:10.48550/arXiv.2401.13256.
- 371 [11] Y. Gao, R. Li, E. Croxford, J. Caskey, B. W. Patterson, M. Churpek, T. Miller, D. Dligach, M.
372 Afshar, Leveraging medical knowledge graphs into large language models for diagnosis prediction:
373 Design and application study, *JMIR AI* 4 (2025) e58670. doi:10.2196/58670.



- 374 [12] R. Yang, H. Liu, E. Marrese-Taylor, Q. Zeng, Y. Ke, W. Li, L. Cheng, Q. Chen, J. Caverlee, Y.
375 Matsuo, et al., Kg-rank: Enhancing large language models for medical qa with knowledge graphs and
376 ranking techniques, arXiv preprint arXiv:2403.05881 (2024). doi:10.48550/arXiv.2403.05881.
- 377 [13] H. Zhang, J. Chen, F. Jiang, F. Yu, Z. Chen, J. Li, G. Chen, X. Wu, Z. Zhang, Q. Xiao, et al.,
378 Huatuogpt, towards taming language model to be a doctor, arXiv preprint arXiv:2305.15075 (2023).
379 doi:10.48550/arXiv.2305.15075.
- 380 [14] Košnik, M., Štajer, D., Blinc, A., Jug, B., Kocjan, T., & Koželj, M. (Eds.). (2022). Interna medicina
381 (6th ed.). Medicinska fakulteta, Slovensko zdravniško društvo, Buča.
- 382 [15] D. Kasper, A. Fauci, S. Hauser, D. Longo, J. Jameson, J. Loscalzo, Harrison's principles of
383 internal medicine, 19e, Vol. 1, Mcgraw-hill New York, NY, USA:, 2015.
- 384 [16] Strojan, P., & Hočevár, M. (Eds.). (2018). Onkologija: Učbenik za študente medicine
385 [Oncology: Textbook for medical students] (1st ed.). Onkološki inštitut Ljubljana. Available under
386 Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.
- 387 [17] StatPearls Publishing. (2025). StatPearls [Internet]. Treasure Island (FL): StatPearls Publishing;
388 2025 Jan-. Available from: <https://www.ncbi.nlm.nih.gov/books/> Licensed under Creative Commons
389 Attribution-NonCommercial-NoDerivatives 4.0 International.