

# **Evaluation of ChatGPT-Generated Differential Diagnosis for Common Diseases with Atypical Presentation: Descriptive Research**

Kiyoshi Shikino, Taro Shimizu, Yuki Otsuka, Masaki Tago, Takahashi Hiromizu, Takashi Watari, Yosuke Sasaki, Gemmei Iizuka, Hiroki Tamura, Koichi Nakashima, Kotaro Kunitomo, Morika Suzuki, Sayaka Aoyama, Shintaro Kosaka, Teiko Kawahigashi, Tomohiro Matsumoto, Fumina Orihara, Toru Morikawa, Toshinori Nishizawa, Yoji Hoshina, Yu Yamamoto, Yuichiro Matsuo, Yuto Unoki, Hirofumi Kimura, Midori Tokushima, Satoshi Watanuki, Takuma Saito, Fumio Otsuka, Yasuharu Tokuda

Submitted to: JMIR Medical Education  
on: March 27, 2024

**Disclaimer:** © The authors. All rights reserved. This is a privileged document currently under peer-review/community review. Authors have provided JMIR Publications with an exclusive license to publish this preprint on its website for review purposes only. While the final peer-reviewed paper may be licensed under a CC BY license on publication, at this stage authors and publisher expressly prohibit redistribution of this draft paper other than for review purposes.

Table of Contents

Original Manuscript..... 5

Supplementary Files..... 22

    Figures ..... 23

        Figure 1..... 24

        Figure 2..... 25

        Figure 3..... 26

    Multimedia Appendixes ..... 27

        Multimedia Appendix 1..... 28

        Multimedia Appendix 2..... 28

# Evaluation of ChatGPT-Generated Differential Diagnosis for Common Diseases with Atypical Presentation: Descriptive Research

Kiyoshi Shikino<sup>1,2</sup> MD, MHPE, PhD; Taro Shimizu<sup>3</sup> MD, PhD, MSc, MPH, MBA; Yuki Otsuka<sup>4</sup> MD, PhD; Masaki Tago<sup>5</sup> MD, PhD; Takahashi Hiromizu<sup>6</sup> MD, PhD; Takashi Watari<sup>7</sup> MD, MHQS, PhD; Yosuke Sasaki<sup>8</sup> MD, PhD; Gemmei Iizuka<sup>9,10</sup> MD, PhD; Hiroki Tamura<sup>1</sup> MD, PhD; Koichi Nakashima<sup>11</sup> MD; Kotaro Kunitomo<sup>12</sup> MD; Morika Suzuki<sup>13</sup> MD; Sayaka Aoyama<sup>14</sup> MD; Shintaro Kosaka<sup>15</sup> MD; Teiko Kawahigashi<sup>16</sup> MD, PhD; Tomohiro Matsumoto<sup>17</sup> MD, DDS, PhD; Fumina Orihara<sup>17</sup> MD; Toru Morikawa<sup>18</sup> MD, PhD; Toshinori Nishizawa<sup>19</sup> MD; Yoji Hoshina<sup>20</sup> MD; Yu Yamamoto<sup>21</sup> MD; Yuichiro Matsuo<sup>22</sup> MD, MPH; Yuto Unoki<sup>23</sup> MD; Hirofumi Kimura<sup>23</sup> MD; Midori Tokushima<sup>24</sup> MD; Satoshi Watanuki<sup>25</sup> MD, MBA; Takuma Saito<sup>25</sup> MD; Fumio Otsuka<sup>4</sup> MD, PhD; Yasuharu Tokuda<sup>26,27</sup> MD, MPH, PhD

<sup>1</sup>Department of General Medicine Chiba University Hospital Chiba JP

<sup>2</sup>Department of Community-Oriented Medical Education Chiba University Graduate School of Medicine Chiba JP

<sup>3</sup>Department of Diagnostic and Generalist Medicine Dokkyo Medical University Tochigi JP

<sup>4</sup>Department of General Medicine Okayama University Graduate School of Medicine Dentistry and Pharmaceutical Sciences Okayama JP

<sup>5</sup>Department of General Medicine Saga University Hospital Saga JP

<sup>6</sup>Department of General Medicine Juntendo University Hospital Faculty of Medicine Tokyo JP

<sup>7</sup>Integrated Clinical Education Center Hospital Integrated Clinical Education Kyoto University Hospital Kyoto JP

<sup>8</sup>Department of General Medicine and Emergency Care Toho University School of Medicine Tokyo JP

<sup>9</sup>Center for Preventive Medical Sciences Chiba University Chiba JP

<sup>10</sup>Tama Family clinic Kanagawa JP

<sup>11</sup>Department of General Medicine Awa Regional Medical Center Chiba JP

<sup>12</sup>Department of General Medicine NHO Kumamoto Medical Center Kumamoto JP

<sup>13</sup>Department of General Internal Medicine National Hospital Organization Sendai Medical Center Miyagi JP

<sup>14</sup>Department of Internal Medicine Mito Kyodo General Hospital Ibaraki JP

<sup>15</sup>Hospital medicine Tokyo Metropolitan Hiroo Hospital Tokyo JP

<sup>16</sup>Department of Molecular and Human Genetics Baylor College of Medicine Houston US

<sup>17</sup>Division of General Medicine Nerima Hikarigaoka Hospital Tokyo JP

<sup>18</sup>Department of General Medicine Nara City Hospital Nara JP

<sup>19</sup>Department of General Internal Medicine St. Luke's International Hospital Tokyo JP

<sup>20</sup>Department of Neurology University of Utah Salt lake US

<sup>21</sup>Division of General Medicine Jichi Medical University Center for Community Medicine Tochigi JP

<sup>22</sup>Department of Clinical Epidemiology and Health Economics The University of Tokyo The Graduate School of Medicine Tokyo JP

<sup>23</sup>Department of General Internal Medicine Iizuka Hospital Fukuoka JP

<sup>24</sup>Saga Medical Career Support Center Saga University Hospital Saga JP

<sup>25</sup>Department of Emergency and General Medicine Tokyo Metropolitan Tama Medical Center Tokyo JP

<sup>26</sup>Muribushi Okinawa Center for Teaching Hospitals Okinawa JP

<sup>27</sup>Tokyo Foundation for Policy Research Tokyo JP

## Corresponding Author:

Kiyoshi Shikino MD, MHPE, PhD

Department of Community-Oriented Medical Education

Chiba University Graduate School of Medicine

1-8-1

Inohana

Chiba

JP

## Abstract

**Background:** Despite significant advancements in medical knowledge and medical diagnosis techniques, misdiagnosis remains a significant public health issue, contributing to mortality and morbidity worldwide. Artificial intelligence (AI), especially

models such as the Generative Pre-trained Transformer (GPT), has shown promise in enhancing diagnostic accuracy. However, the effectiveness of these AI models in diagnosing atypical presentations of common diseases has not been extensively explored.

**Objective:** This study aimed to assess the diagnostic accuracy of the AI model ChatGPT-4 in generating differential diagnoses for atypical presentations of common diseases, and to understand its reliance on patient history during the diagnostic process.

**Methods:** We utilized 25 clinical vignettes from the Journal of Generalist Medicine that presented atypical manifestations of common diseases. Two general medicine physicians categorized the cases based on atypicality. ChatGPT-4 was then employed to generate differential diagnoses, based on the clinical information provided. The concordance between AI-generated and final diagnoses was measured, with a focus on the top-ranked disease (top 1) and the top five differential diagnoses (top 5).

**Results:** ChatGPT-4's diagnostic accuracy decreased with an increase in atypical presentation. For Category 1 (C1) cases, the concordance rates were 17% for the top 1 and 67% for the top 5. Categories 3 (C3) and 4 (C4) showed a 0% concordance for top 1, and markedly lower rates for the top 5, indicating difficulties in handling highly atypical cases.

**Conclusions:** ChatGPT-4 demonstrates potential as an auxiliary tool for diagnosing typical and mildly atypical presentations of common diseases. However, its performance declines with greater atypicality. The findings of study underscores the need for AI systems to encompass a broader range of linguistic capabilities, cultural understanding, and diverse clinical scenarios to improve diagnostic utility in real-world settings. Clinical Trial: NA

(JMIR Preprints 27/03/2024:58758)

DOI: <https://doi.org/10.2196/preprints.58758>

## Preprint Settings

1) Would you like to publish your submitted manuscript as preprint?

Please make my preprint PDF available to anyone at any time (recommended).

Please make my preprint PDF available only to logged-in users; I understand that my title and abstract will remain visible to all users.

✓ **Only make the preprint title and abstract visible.**

No, I do not wish to publish my submitted manuscript as a preprint.

2) If accepted for publication in a JMIR journal, would you like the PDF to be visible to the public?

✓ **Yes, please make my accepted manuscript PDF available to anyone at any time (Recommended).**

Yes, but please make my accepted manuscript PDF available only to logged-in users; I understand that the title and abstract will remain visible to all users.

Yes, but only make the title and abstract visible (see Important note, above). I understand that if I later pay to participate in [a peer-reviewed journal](#).

## Original Manuscript

## Original Paper

# Evaluation of ChatGPT-Generated Differential Diagnosis for Common Diseases with Atypical Presentation: Descriptive Research

Kiyoshi Shikino,<sup>1,2\*</sup> MD, MHPE, PhD, Taro Shimizu,<sup>3</sup> MD, PhD, MSc, MPH, MBA, Yuki Otsuka,<sup>4</sup> MD, PhD, Masaki Tago,<sup>5</sup> MD, PhD, Takahashi Hiromizu,<sup>6</sup> MD, PhD, Takashi Watari,<sup>7</sup> MD, MHQS, PhD, Yosuke Sasaki,<sup>8</sup> MD, PhD, Gemmei Iizuka,<sup>9,10</sup> MD, PhD, Hiroki Tamura,<sup>2</sup> MD, PhD, Koichi Nakashima,<sup>11</sup> MD, Kotaro Kunitomo,<sup>12</sup> MD, Morika Suzuki,<sup>13</sup> MD, PhD, Sayaka Aoyama,<sup>14</sup> MD, Shintaro Kosaka,<sup>15</sup> MD, Takuma Saito, MD, Teiko Kawahigashi,<sup>16</sup> MD, PhD, Tomohiro Matsumoto,<sup>17</sup> MD, DDS, PhD, Fumina Orihara,<sup>17</sup> MD, Toru Morikawa,<sup>18</sup> MD, PhD, Toshinori Nishizawa,<sup>19</sup> MD, Yoji Hoshina,<sup>20</sup> MD, Yu Yamamoto,<sup>21</sup> MD, Yuichiro Matsuo,<sup>22</sup> MD, MPH, Yuto Unoki,<sup>23</sup> MD, Hirofumi Kimura,<sup>23</sup> MD, Midori Tokushima,<sup>24</sup> MD, Satoshi Watanuki,<sup>25</sup> MD, MBA, Takuma Saito,<sup>25</sup> MD, Fumio Otsuka,<sup>4</sup> MD, PhD, Yasuharu Tokuda,<sup>26,27</sup> MD, MPH, PhD

<sup>1</sup>Community-Oriented Medical Education, Chiba University Graduate School of Medicine, Chiba, Japan

<sup>2</sup>Department of General Medicine, Chiba University Hospital, Chiba, Japan

<sup>3</sup>Department of Diagnostic and Generalist Medicine, Dokkyo Medical University, Tochigi, Japan

<sup>4</sup>Department of General Medicine, Okayama University Graduate School of Medicine, Dentistry and Pharmaceutical Sciences, Okayama, Japan.

<sup>5</sup>Department of General Medicine, Saga University Hospital, Saga, Japan

<sup>6</sup>Juntendo University Hospital Faculty of Medicine, Tokyo, Japan

<sup>7</sup>Integrated Clinical Education Center Hospital Integrated Clinical Education, Kyoto University Hospital.

<sup>8</sup>Department of General Medicine and Emergency Care, Toho University School of Medicine

<sup>9</sup>Center for Preventive Medical Sciences, Chiba University, Chiba, Japan

<sup>10</sup>Tama Family clinic, Kanagawa, Japan

<sup>11</sup>General Medicine, Awa Regional Medical Center, Chiba, Japan

<sup>12</sup>General Medicine, NHO Kumamoto Medical Center, Kumamoto, Japan

<sup>13</sup>Department of General Internal Medicine, National Hospital Organization Sendai Medical Center, Miyagi, Japan

<sup>14</sup>Department of Internal Medicine, Mito Kyodo General Hospital

<sup>15</sup>Hospital medicine, Tokyo Metropolitan Hiroo Hospital, Tokyo, Japan

<sup>16</sup>Department of Molecular and human Genetics, Baylor College of Medicine, Texas, USA

<sup>17</sup>Division of General Medicine, Nerima Hikarigaoka Hospital, Tokyo, Japan

<sup>18</sup>Department of General Medicine, Nara City Hospital

<sup>19</sup>Department of General Internal Medicine, St. Luke's International Hospital, Tokyo, Japan

<sup>20</sup>Department of Neurology, University of Utah, Utah, U.S.

<sup>21</sup>Division of General Medicine, Center for Community Medicine, Jichi Medical University, Tochigi, Japan

<sup>22</sup>Department of Clinical Epidemiology and Health Economics, The Graduate School of Medicine, The University of Tokyo, Tokyo, Japan

<sup>23</sup>Department of General Internal Medicine, Iizuka Hospital, Fukuoka, Japan

<sup>24</sup>Saga Medical Career Support Center, Saga University Hospital, Saga, Japan

<sup>25</sup>Department of Emergency and General Medicine, Tokyo Metropolitan Tama Medical Center, Tokyo, Japan

<sup>26</sup>*Muribushi Okinawa Center for Teaching Hospitals, Okinawa, Japan*

<sup>27</sup>*Tokyo Foundation for Policy Research, Tokyo, Japan*

### Corresponding Author:

Kiyoshi	Shikino,	MD,	MHPE,	PhD
Department of Community-Oriented Medical Education Chiba University Graduate School of Medicine				
1-8-1,				Inohana
Chiba,				2608677
Japan				
Phone:	81	43	222	7171
Email:	<a href="mailto:kshikino@gmail.com">kshikino@gmail.com</a>			

## ABSTRACT

**Background:** The persistence of diagnostic errors, despite advances in medical knowledge and diagnostics, highlights the importance of understanding atypical disease presentations and their contribution to mortality and morbidity. Artificial intelligence (AI), particularly Generative Pre-trained Transformers like ChatGPT-4, holds promise for improving diagnostic accuracy, but requires further exploration in handling atypical presentations.

**Objective:** This study aimed to assess the diagnostic accuracy of the AI model ChatGPT-4 in generating differential diagnoses for atypical presentations of common diseases, with a focus on the model's patient history reliance during the diagnostic process.

**Methods:** We utilized 25 clinical vignettes from the Journal of Generalist Medicine, characterizing atypical manifestations of common diseases. Two general medicine physicians categorized the cases based on atypicality. ChatGPT-4 was then employed to generate differential diagnoses, based on the clinical information provided. The concordance between AI-generated and final diagnoses was measured, with a focus on the top-ranked disease (top 1) and the top five differential diagnoses (top 5).

**Results:** ChatGPT-4's diagnostic accuracy decreased with an increase in atypical presentation. For Category 1 (C1) cases, the concordance rates were 17% for the top 1 and 67% for the top 5. Categories 3 (C3) and 4 (C4) showed a 0% concordance for top 1, and markedly lower rates for the top 5, indicating difficulties in handling highly atypical cases. Chi-squared tests revealed no significant difference in the top 1 differential diagnosis accuracy between less atypical (C1+C2) and more atypical (C3+C4) groups ( $\chi^2(1, N=25) = 2.07, p = .131$ ). However, a significant difference was found in the top 5 analyses, with less atypical cases showing higher accuracy ( $\chi^2(1, N=25) = 4.01, p = .048$ ).

**Conclusions:** ChatGPT-4 demonstrates potential as an auxiliary tool for diagnosing typical and mildly atypical presentations of common diseases. However, its performance declines with greater atypicality. The study findings underscore the need for AI systems to encompass a broader range of linguistic capabilities, cultural understanding, and diverse clinical scenarios to improve diagnostic utility in real-world settings.

**Keywords:** atypical presentation; ChatGPT; common disease; diagnostic accuracy; diagnosis; patient safety



## Introduction

For the past decade, medical knowledge and diagnostic techniques have expanded globally, becoming more accessible with remarkable advancements in clinical testing and useful reference systems [1]. Despite these advancements, misdiagnosis significantly contributes to mortality, making it a noteworthy public health issue [2,3]. Studies have revealed discrepancies between clinical and postmortem autopsy diagnoses in at least 25% of cases, with diagnostic errors contributing to approximately 10% of deaths, and to 6–17% of hospital adverse events [4–8]. The significance of atypical presentations as a contributor to diagnostic errors is especially notable, with recent findings suggesting that such presentations are prevalent in a substantial portion of outpatient consultations, and are associated with a higher risk of diagnostic inaccuracies [9]. This underscores the persistent challenge in diagnosing patients correctly due to the variability in disease presentation and reliance on medical history, which comprises approximately 80% of the medical diagnosis [10,11].

The advent of artificial intelligence (AI) in healthcare, particularly through natural language processing (NLP) models such as the Generative Pre-trained Transformer (GPT), has opened new avenues in medical diagnosis [12]. Recent studies on AI medical diagnosis across various specialties—including neurology [13], dermatology [14], radiology [15], and pediatrics [16]—have shown promising results, improving diagnostic accuracy, efficiency, and safety. Among these developments, ChatGPT-4, a state-of-the-art AI model developed by OpenAI, has demonstrated remarkable capabilities in understanding and processing medical language, significantly outperforming its predecessors in medical knowledge assessments, and potentially transforming medical education and clinical decision support systems [12,17].

Notably, one study found that ChatGPT could pass the United States Medical Licensing Examination (USMLE), highlighting its potential in medical education and medical diagnosis [18,19]. Moreover, in controlled settings, ChatGPT-4 has shown over 90% accuracy in diagnosing common diseases with typical presentations, based on chief complaints and patient history [20]. However, while research has examined the diagnostic accuracy of AI chatbots, including ChatGPT models, in generating differential diagnoses for complex clinical vignettes derived from general internal medicine (GIM) department case reports, their diagnostic accuracy in handling atypical presentations of common diseases remains less explored [21,22]. There has been a notable study aimed at evaluating the accuracy of the differential diagnosis lists generated by both third- and fourth-generation ChatGPT models using case vignettes from case reports published by the Department of GIM of Dokkyo Medical University Hospital, Japan. ChatGPT-4 was found to achieve a correct diagnosis rate within the top 10 differential diagnosis lists, top five lists, and top diagnoses of 83%, 81%, and 60%, respectively — rates comparable to those by physicians. Although the study highlights the potential of ChatGPT-4 as a supplementary tool for physicians, particularly in the context of GIM, it also underlines the importance of further investigation into the diagnostic accuracy of ChatGPT with atypical disease presentations (Figure 1). Given the crucial role of patient history in diagnosis and the inherent variability in disease presentation, our study expands upon this foundation to assess the accuracy of ChatGPT-4 in diagnosing common diseases with atypical presentations [23].

### **[Figure 1 here]**

More specifically, this study aims to evaluate the hypothesis that the diagnostic accuracy of AI, exemplified by ChatGPT-4, declines when dealing with atypical presentations of common diseases. We hypothesize that despite the known capabilities of AI in recognizing typical disease patterns, its performance will be significantly challenged when presented with clinical cases that deviate from these patterns, leading to reduced diagnostic precision. Consequently,

this study seeks to systematically assess this hypothesis and explore its implications for the integration of AI in clinical practice. By exploring the contribution of AI-assisted medical diagnoses to common diseases with atypical presentation and patient history, the study assesses the accuracy of ChatGPT in reaching a clinical diagnosis based on the medical information provided. By reevaluating the significance of medical information, our study contributes to the ongoing discourse on optimizing diagnostic processes — both conventional and AI-assisted.

## Methods

### Study Design, Settings, and Participants

This study utilized a series of 25 clinical vignettes from a special issue of Generalist Medicine (International Standard Serial Number 2188-8051, Japanese), published on March 5, 2024. These vignettes, which exemplify atypical presentations of common diseases, were selected for their alignment with our research aim to explore the impact of atypical disease presentations in AI-assisted diagnosis. The clinical vignettes were derived from real patient cases and curated by an editorial team specializing in general internal medicine, with final edits by KS. Each case included comprehensive details such as age, gender, chief complaints, medical history, medication history, current illness, and physical examination findings, along with the ultimate and initial misdiagnoses.

An expert panel comprising two general medicine and medical education physicians, TS and YO, initially reviewed these cases. After deliberation, they selected all 25 cases that exemplified atypical presentations of common diseases. Subsequently, TS and YO evaluated their degree of atypicality and categorized them into four distinct levels, using the following definition as a guide: *'Atypical presentations have a shortage of prototypical features. These can be defined as features that are most frequently encountered in patients with the disease, features encountered in advanced presentations of the disease, or simply features of the disease commonly listed in medical textbooks. Atypical presentations may also have features with unexpected values [24].'* Category 1 was assigned to cases that were closest to the typical presentations of common diseases, whereas Category 4 was designated to those that were markedly atypical. In instances where TS and YO did not reach consensus, a third expert, KS, was consulted. Through collaborative discussions, the panel reached a consensus on the final category for each case, ensuring a systematic and comprehensive evaluation of the atypical presentation of common diseases (Figure 2).

#### [Figure 2 here]

Our analysis was conducted on March 12, 2024, utilizing ChatGPT-4's proficiency in Japanese. The language processing was enabled by the standard capabilities of the ChatGPT-4 model, requiring no additional adaptations or programming by our team. We exclusively used text-based input for the generative AI, excluding tables or images to maintain a focus on linguistic data. This approach is consistent with the typical constraints of language-based AI diagnostic tools. Inputs to ChatGPT-4 consisted of direct transcriptions of the original case reports in Japanese, ensuring the authenticity of the medical information was preserved. We measured the concordance between AI-generated differential diagnoses and the vignettes' final diagnoses, as well as the initial misdiagnoses. Our investigation entailed inputting clinical information—including medical history, physical examination, and laboratory data—into ChatGPT, followed by posing the question 'List of differential diagnoses in order of likelihood, based on the provided vignettes' information,' labeled as 'GAI differential diagnoses.'

## Data Collection and Measurements

We assigned the correct diagnosis for each of these 25 cases as “Final diagnosis.” We then used ChatGPT to generate differential diagnoses (GAI differential diagnoses). For each case, ChatGPT was prompted to create a list of differential diagnoses. Patient information was provided in full each time, without incremental inputs. The concordance rate between “Final diagnosis,” “Misdiagnosis,” and “GAI differential diagnoses” was then assessed. To extract a list of diagnoses from ChatGPT, we concluded each input session with the phrase “List of differential diagnoses in order of likelihood, based on the provided vignettes’ information.” We measured the percentage in which the final diagnosis or misdiagnosis was included in the top-ranked disease (top 1) and within the top five differential diagnoses (top 5) generated by ChatGPT (Figure 3).

[Figure 3 here]

## Data Analysis

Two board-certified physicians working in the medical diagnostic department of our facility judged the concordance between the AI-proposed diagnoses and the final diagnosis. The two physicians are GIM board-certified. The number of postgraduate years of the physicians was 7 and 17, respectively. A diagnosis was considered to match if the two physicians agreed to the concordance. We measured the interrater reliability with the K coefficient (0.8-1.0 = almost perfect, 0.6-0.8 = substantial, 0.4-0.6 = moderate, and 0.2-0.4 = fair) [25]. To further analyze the accuracy of the top 1 and top 5 diagnoses, we used Chi-square or Fisher's exact test, as appropriate. Statistical analyses were conducted using IBM SPSS Statistics for Windows 26.0 (IBM Corp. Armonk, NY), with the level of significance set at  $P < .05$ .

## Disclosure of Generative Artificial Intelligence Usage

In this study, generative artificial intelligence was utilized to create differential diagnoses for cases published in medical journals. However, it was not used in actual clinical practice. Similarly, no generative artificial intelligence was used in our manuscript writing.

## Ethics Approval

Our research did not involve humans, medical records, patient information, observations of public behaviors, or secondary data analyses; thus, it was exempt from ethical approval, informed consent requirements, and institutional review board approval. Additionally, as no identifying information was included, the data did not need to be anonymized or de-identified. We did not offer any compensation because there were no human participants in the study.

## Results

The 25 clinical vignettes comprised 11 male and 14 female patients, with ages ranging from 21 to 92 years. All individuals were older than 20 years, and eight were older than 65 years. Table 1, Supplementary file 1, and Supplementary file 2 present these results. The correct final diagnosis listed in the Journal of Generalist Medicine clinical vignette as common disease presenting atypical symptoms (labeled as “Final diagnosis”) showed that “GAI differential

diagnoses” and “Final diagnosis” coincided 12% (3/12) within the first list differential diagnosis, while “GAI differential diagnoses” and “Final diagnosis” had a concordance rate of 44% (11/25) within five differential diagnoses. The interrater reliability was substantial (Cohen’s kappa = 0.84).

**Table 1.** List of answers and diagnoses provided by ChatGPT  
Category 1 (C1) being closest to typical, and Category 4 (C4) being most atypical.

Case	Age	Gender	Final diagnosis	Category	GAI diagnosis rank
1	34	F	Caffeine intoxication	1	0
2	40	F	Asthma	1	1
3	55	F	Obsessive-compulsive disorder	1	3
4	58	M	Drug-induced enteritis	1	3
5	38	F	Cytomegalovirus infection	1	3
6	29	M	Acute HIV infection	1	5
7	62	M	Cardiogenic cerebral embolism	2	1
8	70	M	Cervical epidural hematoma	2	0
9	70	F	Herpes zoster	2	0
10	86	F	Hemorrhagic gastric ulcer	2	0
11	77	M	Septic arthritis	2	3
12	78	F	Compression fracture	2	0
13	45	M	Infective endocarditis	2	0
14	21	F	Ectopic pregnancy	2	1
15	55	F	Non-ST elevation myocardial infarction	2	2
16	54	F	Hypoglycemia	3	0
17	77	F	Giant cell arteritis	3	0
18	60	M	Adrenal insufficiency	3	4
19	38	F	Generalized anxiety disorder	3	0
20	24	F	Graves' disease	4	4
21	31	M	Acute myeloblastic leukemia	4	0
22	76	F	Elderly onset rheumatoid arthritis	4	0
23	45	M	Appendicitis	4	0
24	92	M	Rectal cancer	4	0
25	60	M	Acute aortic dissection	4	0

GAI diagnosis rank: The high priority differential diagnosis rank generated by ChatGPT.

Final diagnosis: Final correct diagnosis listed in the Journal of Generalist Medicine clinical vignette as common disease presenting atypical symptoms.

The analysis of the concordance rates between the "GAI differential diagnoses" generated by ChatGPT and the "Final diagnosis" from the Journal of Generalist Medicine

revealed distinct patterns across the four categories of atypical presentations (Table 2). For the top 1 differential diagnosis, Category 1 (C1) cases, which were closest to typical presentations, showed a concordance rate of 17%, whereas Category 2 (C2) cases exhibited a slightly higher rate of 22%. Remarkably, Categories 3 (C3) and 4 (C4), which represent more atypical cases, demonstrated no concordance (0%) in the top 1 differential diagnosis.

When the analysis was expanded to the top five differential diagnoses, the concordance rates varied across categories. C1 cases showed a significant increase in concordance to 67%, indicating a better performance of the “GAI differential diagnoses” when considering a broader range of possibilities. C2 cases had a concordance rate of 44%, followed by C3 cases at 25% and C4 cases at 17%.

Table 2. Concordance rates of AI-generated differential diagnoses by atypicality category Category 1 (C1) being closest to typical, and Category 4 (C4) being most atypical.

Category	Rank 1 (n)	Rank 2 (n)	Rank 3 (n)	Rank 4 (n)	Rank 5 (n)	Misdiagnosis (n)	Top 1 (%)	Top 5 (%)
C1	1	0	3	0	0	2	17	67
C2	2	1	1	0	0	5	22	44
C3	0	0	0	1	0	3	0	25
C4	0	0	0	1	0	5	0	17

To assess the diagnostic accuracy of ChatGPT across varying levels of atypical presentations, we employed chi-squared tests. Specifically, we compared the frequency of correct diagnoses in the top 1 and top 5 differential diagnoses provided by ChatGPT for cases categorized as C1+C2 (less atypical) versus C3+C4 (more atypical). For the top 1 differential diagnosis, there was no statistically significant difference in the number of correct diagnoses between the less atypical (C1+C2) and more atypical (C3+C4) groups ( $\chi^2(1, N=25) = 2.07, p = .131$ ). However, when expanding the analysis to the top 5 differential diagnoses, we found a statistically significant difference, with the less atypical group (C1+C2) demonstrating a higher number of correct diagnoses compared to the more atypical group (C3+C4) ( $\chi^2(1, N=25) = 4.01, p = .048$ ).

## DISCUSSION

This study provides insightful data on the performance of ChatGPT-4 in diagnosing common diseases with atypical presentations. Our findings offer a nuanced view of the capacity of AI-driven differential diagnoses across varying levels of atypicality. In the analysis of the concordance rates between “GAI differential diagnoses” and “Final diagnosis,” we observed a decrease in diagnostic accuracy as the degree of atypical presentation increased.

The performance of ChatGPT-4 in Category 1 (C1) cases, which are the closest to typical presentations, was moderately successful, with a concordance rate of 17% for the top 1 diagnosis and 67% within the top five. This suggests that when the disease presentation closely aligns with the typical characteristics known to the model, ChatGPT-4 is relatively reliable at

identifying a differential diagnosis list that coincides with the final diagnosis. However, the utility of ChatGPT-4 appears to decrease as atypicality increases, as evidenced by the lower concordance rates in Category 2 (C2), and notably more so in Categories 3 (C3) and 4 (C4), where the concordance rates for the top 1 diagnosis fell to 0%. Similar challenges were observed in another 2024 study [26], where the diagnostic accuracy of ChatGPT varied depending on the disease etiology, particularly in differentiating between CNS and non-CNS tumors.

It is particularly revealing that in the more atypical presentations of common diseases (C3 and C4), the AI struggled to provide a correct diagnosis, even within the top five differential diagnoses, with concordance rates of 25% and 17%, respectively. These categories highlight the current limitations of AI in medical diagnosis when faced with cases that deviate significantly from the established patterns within its training data [27].

By leveraging the comprehensive understanding and diagnostic capabilities of ChatGPT-4, this study aims to re-evaluate the significance of patient history in AI-assisted medical diagnosis, and contribute to optimizing diagnostic processes [28]. Our exploration of ChatGPT-4's performance in processing atypical disease presentations not only advances our understanding of AI's potential in medical diagnosis [23], but also underscores the importance of integrating advanced AI technologies with traditional diagnostic methodologies to enhance patient care and reduce diagnostic errors.

The contrast in performance between the C1 and C4 cases can be seen as indicative of the challenges AI systems currently face with complex clinical reasoning requiring pattern recognition. Atypical presentations can include uncommon symptoms, rare complications, or unexpected demographic characteristics, which may not be well-represented in the datasets used to train AI systems [29]. Furthermore, these findings can inform the development of future versions of AI medical diagnosis systems, and guide training curricula to include a broader spectrum of atypical presentations.

This study underscores the importance of the continued refinement of AI medical diagnosis systems, as highlighted by the recent advances in AI technologies and their applications in medicine. Studies published in 2024 [30-32] provide evidence of the rapidly increasing capabilities of large language models (LLMs) like GPT-4 in various medical domains, including oncology, where AI is expected to significantly impact precision medicine [30]. The convergence of text and image processing, as seen in multimodal AI models, suggests a qualitative leap in AI's ability to process complex medical information, which is particularly relevant for our findings on AI-assisted medical diagnostics [30]. These developments reinforce the potential of AI tools like ChatGPT-4 in bridging the knowledge gap between machine learning developers and practitioners, and their role in simplifying complex data analyses in medical research and practice [31]. However, as these systems evolve, it is crucial to remain aware of their limitations and the need for rigorous verification processes to mitigate the risk of errors, which can have significant implications in clinical settings [32]. This aligns with our observation of decreased diagnostic accuracy in atypical presentations and the necessity for cautious integration of AI into clinical practice. It also points to the potential benefits of combining AI with human expertise to compensate for current AI limitations, and enhance diagnostic accuracy [33].

Our research suggests that while AI, particularly ChatGPT-4, shows promise as a supplementary tool for medical diagnosis, reliance on this technology should be balanced with expert clinical judgment, especially in complex and atypical cases [28,29]. The observed concordance rate of 67% for C1 cases indicates that even when not dealing with extremely atypical presentations, cases with potential pitfalls may result in AI medical diagnosis accuracy lower than the 80–90% estimated by existing studies [10,11]. This revelation highlights the

need for cautious integration of AI in clinical settings, acknowledging that its diagnostic capabilities, while robust, may still fall short in certain scenarios [34,35].

## Limitations

Despite the strengths of our research, the study has certain limitations that must be noted when contextualizing our findings. First, the external validity of the results may be limited as our dataset comprises only 25 clinical vignettes, sourced from a special issue of the *Journal of Generalist Medicine*. While these vignettes were chosen for their relevance to the study's hypothesis on atypical presentations of common diseases, the size of the dataset and its origin from mock scenarios rather than real patient data may limit the generalizability of our findings. This sample size may not adequately capture the variability and complexities typically encountered in broader clinical practice, and thus, might not be sufficient to firmly establish statistical generalizations. This limitation is compounded by the exclusion of pediatric vignettes, which narrows the demographic range of our findings and potentially reduces their applicability across diverse age groups.

Second, ChatGPT's current linguistic capabilities predominantly cater to English, presenting significant barriers to patient-provider interactions that may occur in other languages. This raises concerns about the potential for miscommunication and subsequent misdiagnosis in non-English medical consultations. This underscores the essential need for future AI models to exhibit a multilingual capacity that can grasp the subtleties inherent in various languages and dialects, as well as the cultural contexts within which they are used.

Finally, the diagnostic prioritization process of ChatGPT did not always align with clinical probabilities, potentially skewing the perceived effectiveness of the AI model. Additionally, it must be acknowledged that our research utilized ChatGPT-4, which is not a publicly available model. Consequently, the results obtained using ChatGPT-4 may not be directly generalizable to other large language models, especially open-source models like Llama3, which might have different underlying architectures and training datasets. Moreover, since our study relied on clinical vignettes which are mock scenarios, the potential for bias based on the cases is significant. The lack of real demographic diversity in these vignettes means that the findings may not accurately reflect the social or regional nuances, such as ethnicity, prevalence of disease, or cultural practices, that could influence diagnostic outcomes. This limitation suggests a need for careful consideration when applying these AI tools across different geographic and demographic contexts to ensure the findings are appropriately adapted to local populations. This emphasizes the necessity for AI systems to be evaluated in diverse real-world settings to understand their effectiveness comprehensively and mitigate any bias. This distinction is important to consider when extrapolating our study's findings to other AI systems. Future studies should not only refine AI's diagnostic reasoning, but also explore the interpretability of its decision-making process, especially when errors occur. ChatGPT should be considered as a supplementary tool in medical diagnosis, rather than a standalone solution. This reinforces the necessity for combined expertise, where AI supports—but does not replace—human clinical judgment. Further research should expand these findings to a wider range of conditions, especially prevalent diseases with significant public health impacts, to thoroughly assess the practical utility and limitations of AI in medical diagnosis.



## Conclusions

Our study contributes valuable evidence for the ongoing discourse on the role of AI in medical diagnosis. This study provides a foundation for future research to explore the extent to which AI can be trained to recognize increasingly complex and atypical presentations, which is critical for its successful integration into clinical practice.

## Acknowledgements

The authors thank the members of Igaku-Shoin, Tokyo, Japan, for permission to use the clinical vignettes. Igaku-Shoin did not participate in designing and conducting the study; data analysis and interpretation; preparation, review, or approval of the paper; or the decision to submit the paper for publication. The authors thank Dr. Mai Hongo, Saka General Hospital, for providing a clinical vignette. The authors also thank Editage for the English language review.

## Authors' Contributions

KS, TW, TS, YO, MT, HT, YS, and YT designed the study. TS and YO checked the atypical case categories. MT and HT confirmed the diagnoses. KS wrote the first draft and analyzed the research data. All of the authors created atypical common clinical vignettes and published them in the Journal of General Medicine. KS, TS, and HT critically revised the manuscript. All of the authors have checked the final version of the manuscript.

## Data Availability

The datasets generated and analyzed in this study are available from the corresponding author upon reasonable request.

## Conflicts of Interest

None declared.

## Abbreviations

AI: artificial intelligence

Chat GPT: Chat Generative Pre-trained Transformer

## REFERENCES

1. Brown MP, Lai-Goldman M, Billings PR. Translating Innovation in Diagnostics: Challenges and Opportunities. *Genomic and Personalized Medicine*. 2009:367–377.
2. Omron R, Kotwal S, Garibaldi BT, Newman-Toker DE. The diagnostic performance feedback “calibration gap”: Why clinical experience alone is not enough to prevent serious diagnostic errors. *AEM Educ Train* 2018;2:339-342.
3. Balogh EP, Miller BT, Ball JR, editors. Improving diagnosis in health care. Washington, DC: National Academies Press; 2015.

4. Friberg N, Ljungberg O, Berglund E, Berglund D, Ljungberg R, Alafuzoff I, et al. Cause of death and significant disease found at autopsy. *Virchows Arch* 2019;475:781-788.
5. Shojania KG, Burton EC, McDonald KM, Goldman L. Changes in rates of autopsy-detected diagnostic errors over time: a systematic review. *JAMA* 2003;289:2849.
6. Schmitt BP, Kushner MS, Wiener SL. The diagnostic usefulness of the history of the patient with dyspnea. *J Gen Intern Med* 1986;1:386-393.
7. Kuijpers CCHJ, Fronczek J, Van De Goot FRW, Niessen HWM, Van Diest PJ, Jiwa M. The value of autopsies in the era of high-tech medicine: discrepant findings persist. *J Clin Pathol* 2014;67:512-519.
8. Ball JR, Balogh E. Improving diagnosis in health care: highlights of a report from the national academies of sciences, engineering, and medicine. *Ann Intern Med* 2016;164:59.
9. Harada Y, Otaka Y, Katsukura S, Shimizu T. Prevalence of atypical presentations among outpatients and associations with diagnostic error. *Diagnosis (Berl)*. 2023;11(1):40-48.
10. Hampton JR, Harrison MJ, Mitchell JR, Prichard JS, Seymour C. Relative contributions of history-taking, physical examination, and laboratory investigation to diagnosis and management of medical outpatients. *BMJ* 1975;2:486-489.
11. Peterson MC, Holbrook JM, Hales DV, Smith NL, Staker LV. Contributions of the history, physical examination, and laboratory investigation in making medical diagnoses: *Obstet Gynecol Surv* 1992;47:711-712.
12. Alowais SA, Alghamdi SS, Alsuhbeyany N, Alqahtani T, Alshaya AI, Almohareb SN, Aldairem A, Alrashed M, Bin Saleh K, Badreldin HA, Al Yami MS, Al Harbi S, Albekairy AM. Revolutionizing healthcare: The role of artificial intelligence in clinical practice. *BMC Med Educ*. 2023;23(1):689.
13. Giannos P. Evaluating the limits of AI in medical specialisation: ChatGPT's performance on the UK Neurology Specialty Certificate Examination. *BMJ Neurol Open*. 2023;5(1):e000451.
14. Passby L, Jenko N, Wernham A. Performance of ChatGPT on dermatology Specialty Certificate Examination multiple choice questions. *Clin Exp Dermatol*. 2023;llad197.
15. Srivastav S, Chandrakar R, Gupta S, Babhulkar V, Agrawal S, Jaiswal A, Prasad R, Wanjari MB. ChatGPT in radiology: The advantages and limitations of Artificial Intelligence for medical imaging diagnosis. *Cureus*. 2023;15(7):e41435.
16. Andykarayalar R, Surapaneni KM. ChatGPT in Pediatrics: Unraveling its significance as a clinical decision support tool. *Indian Pediatr*. 2024:S097475591600610. Online ahead of print.
17. Mugahed A. Al-Antari. Artificial Intelligence for Medical Diagnostics—Existing and Future AI Technology! *Diagnostics (Basel)*. 2023;13(4): 688.
18. Mihalache A, Huang RS, Popovic MM, Muni RH. ChatGPT-4: An assessment of an upgraded artificial intelligence chatbot in the United States Medical Licensing Examination. *Med Teach*. 2024;46(3):366-372.
19. Kung TH, Cheatham M, Medenilla A, Sillos C, De Leon L, Elepano C, Madriaga M, Aggabao R, Diaz-Candido G, Maningo J, et al. Performance of ChatGPT on USMLE: potential for AI-assisted medical education using large language models. *PLOS Digit Health*. 2023;2(2):e0000198.
20. Fukuzawa F, Yanagita Y, Yokokawa D, Uchida S, Yamashita S, Li Y, Shikino K, Tsukamoto T, Noda K, Uehara T, Ikusaka M. Importance of patient history in Artificial Intelligence-assisted medical diagnosis: Comparison study. *JMIR Med Educ*. 2024;10:e52674.
21. Rao A, Pang M, Kim J, Kamineni M, Lie W, Prasad AK, Landman A, Dreyer K, Succi MD. Assessing the utility of ChatGPT throughout the entire clinical workflow: Development and usability study. *J Med Internet Res*. 2023;25:e48659.
22. Hirosawa T, Kawamura R, Harada Y, Mizuta K, Tokumasu K, Kaji Y, Suzuki T, Shimizu T.

- ChatGPT-generated differential diagnosis lists for complex case-derived clinical vignettes: Diagnostic accuracy evaluation. *JMIR Med Inform.* 2023;11:e48808.
23. Suthar PP, Kounsai A, Chhetri L, Saini D, Dua SG. Artificial Intelligence (AI) in radiology: A deep dive into ChatGPT 4.0's accuracy with the American Journal of Neuroradiology's (AJNR) "Case of the Month". *Cureus.* 2023;15(8):e43958.
  24. Kostopoulou O, Delaney BC, Munro CW. Diagnostic difficulty and error in primary care: a systematic review. *Fam Pract.* 2008;25(6):400-13.
  25. Landis JR, Koch GG. An application of hierarchical kappa-type statistics in the assessment of majority agreement among multiple observers. *Biometrics.* 1977;33 (2):363-374.
  26. Horiuchi D, Tatekawa H, Shimono T, Walston SL, Takita H, Matsushita S, Oura T, Mitsuyama Y, Miki Y, Ueda D. Accuracy of ChatGPT generated diagnosis from patient's medical history and imaging findings in neuroradiology cases. *Neuroradiology.* 2024;66(1):73-79
  27. Umapathy VR, Rajinikanth BS, Samuel Raj RD, Yadav S, Munavarah SA, Anandapandian PA, Mary AV, Padmavathy KRA. Perspective of Artificial Intelligence in disease diagnosis: A review of current and future endeavours in the medical field. *Cureus.* 2023;15(9):e45684.
  28. Mizuta K, Hirokawa T, Harada Y, Shimizu T. Can ChatGPT-4 evaluate whether a differential diagnosis list contains the correct diagnosis as accurately as a physician? *Diagnosis (Berl).* 2024. Online ahead of print.
  29. Ueda D, Walston SL, Matsumoto T, Deguchi R, Tatekawa H, Miki Y. Evaluating GPT-4-based ChatGPT's clinical potential on the NEJM quiz. *BMC Digital Health.* 2024;2:4.
  30. Truhn D, Eckardt JN, Ferber D, Kather JN. Large language models and multimodal foundation models for precision oncology. *NPJ Precis Oncol.* 2024;8(1):72.
  31. Tayebi Arasteh S, Han T, Lotfinia M, Kuhl C, Kather JN, Truhn D, Nebelung S. Large language models streamline automated machine learning for clinical studies. *Nat Commun.* 2024;15(1):1603.
  32. Lee P, Bubeck S, Petro J. Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine. *N Engl J Med.* 2023;388(13):1233-1239.
  33. Harada T, Shimizu T, Kaji Y, Suyama Y, Matsumoto T, Kosaka C, Shimizu H, Nei T, Watanuki S. A perspective from a case conference on comparing the diagnostic process: Human diagnostic thinking vs. Artificial Intelligence (AI) decision support tools. *Int J Environ Res Public Health.* 2020;17(17):6110.
  34. Voelker R. The Promise and Pitfalls of AI in the Complex World of Diagnosis, Treatment, and Disease Management. *JAMA.* 2023;330(15):1416-1419.
  35. Takagi S, Watari T, Erabi A, Sakaguchi K. Performance of GPT-3.5 and GPT-4 on the Japanese Medical Licensing Examination: Comparison Study. *JMIR Med Educ.* 2023;9:e48002.

Figure Legends

Figure 1. Study motivation.

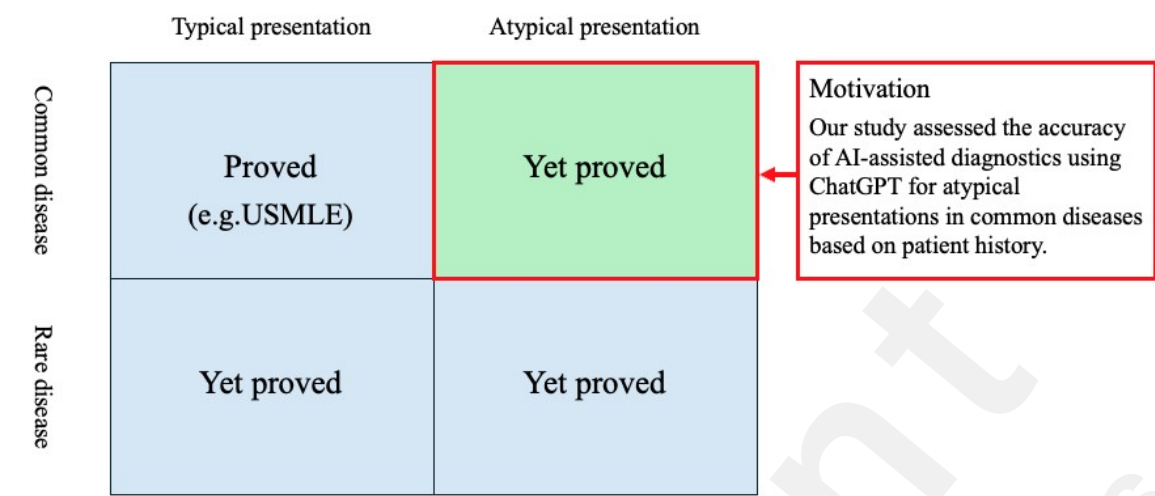


Figure 2. Categories of common disease with atypical presentation.

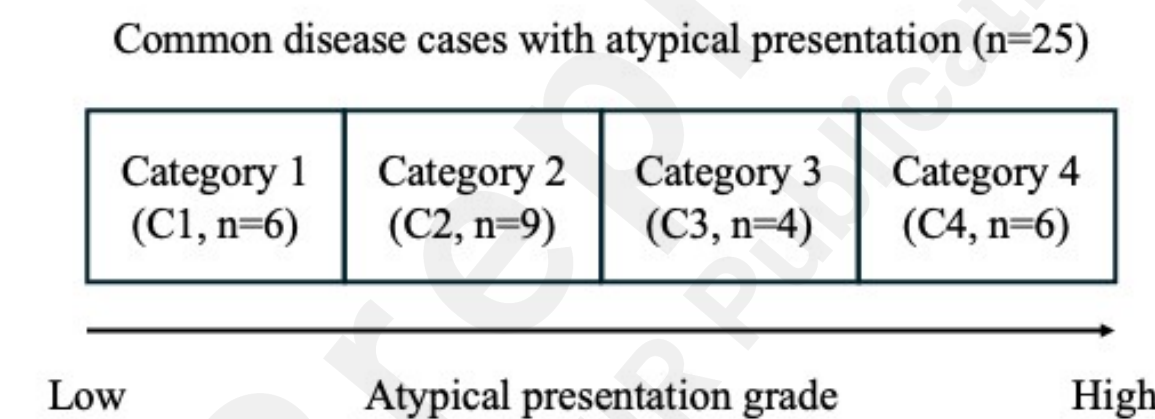
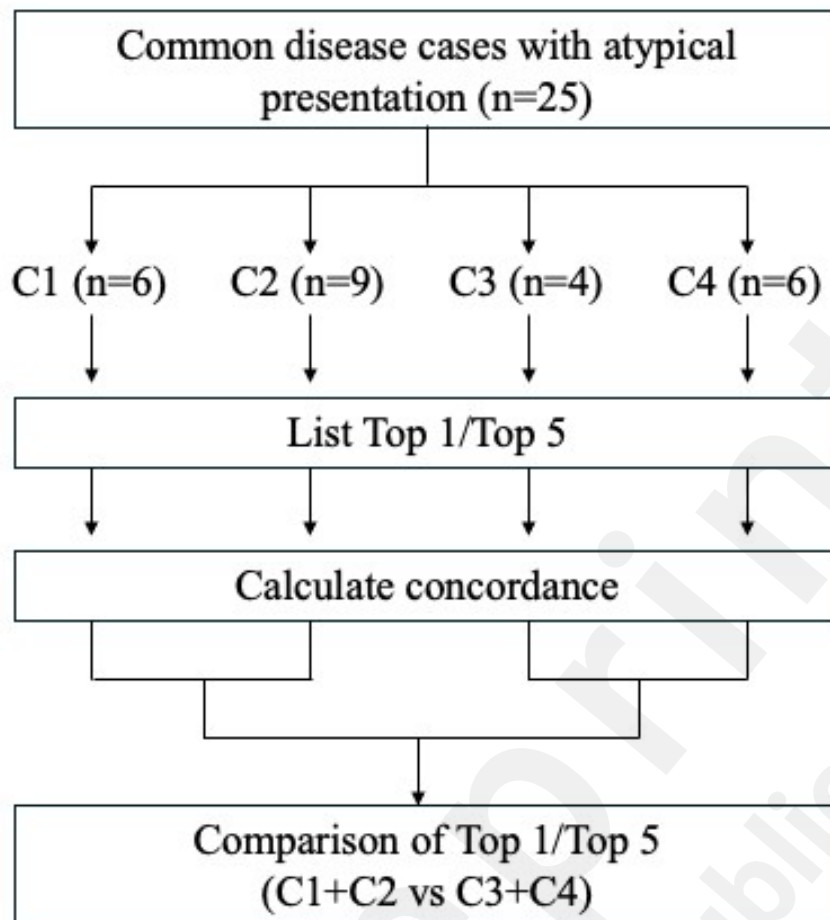


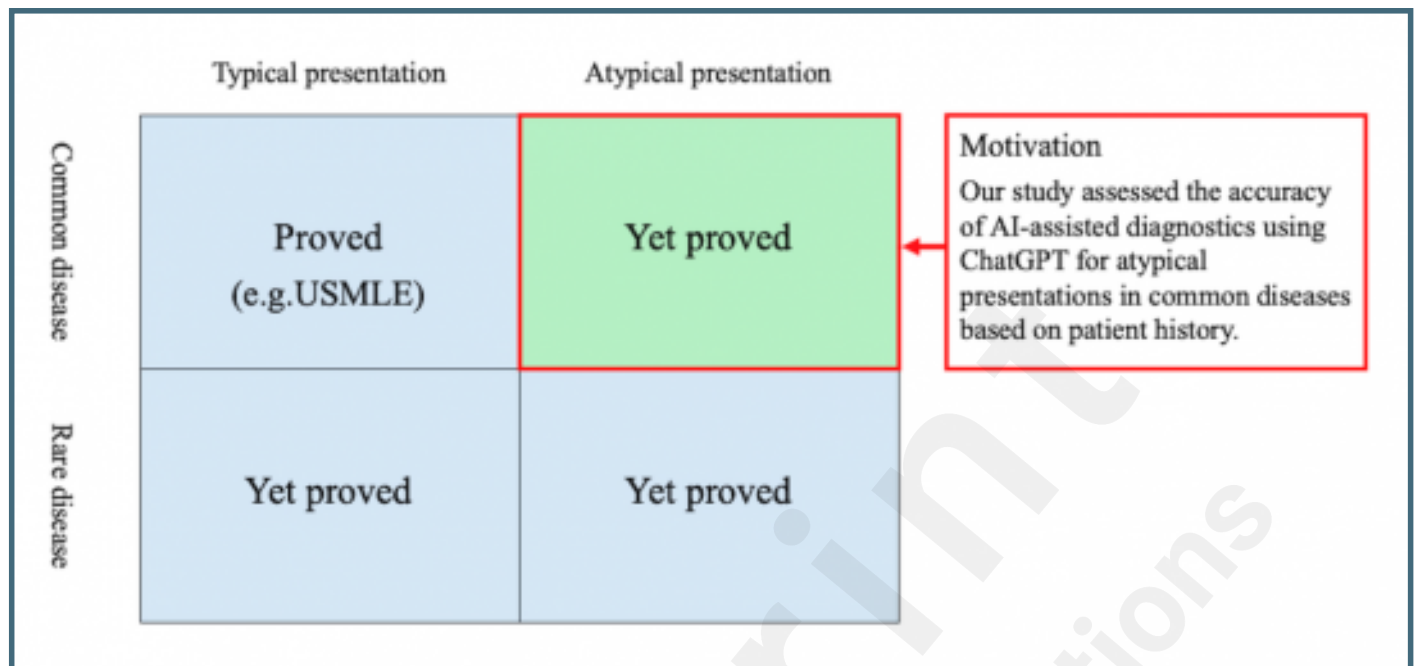
Figure 3. Study flow.



## Supplementary Files

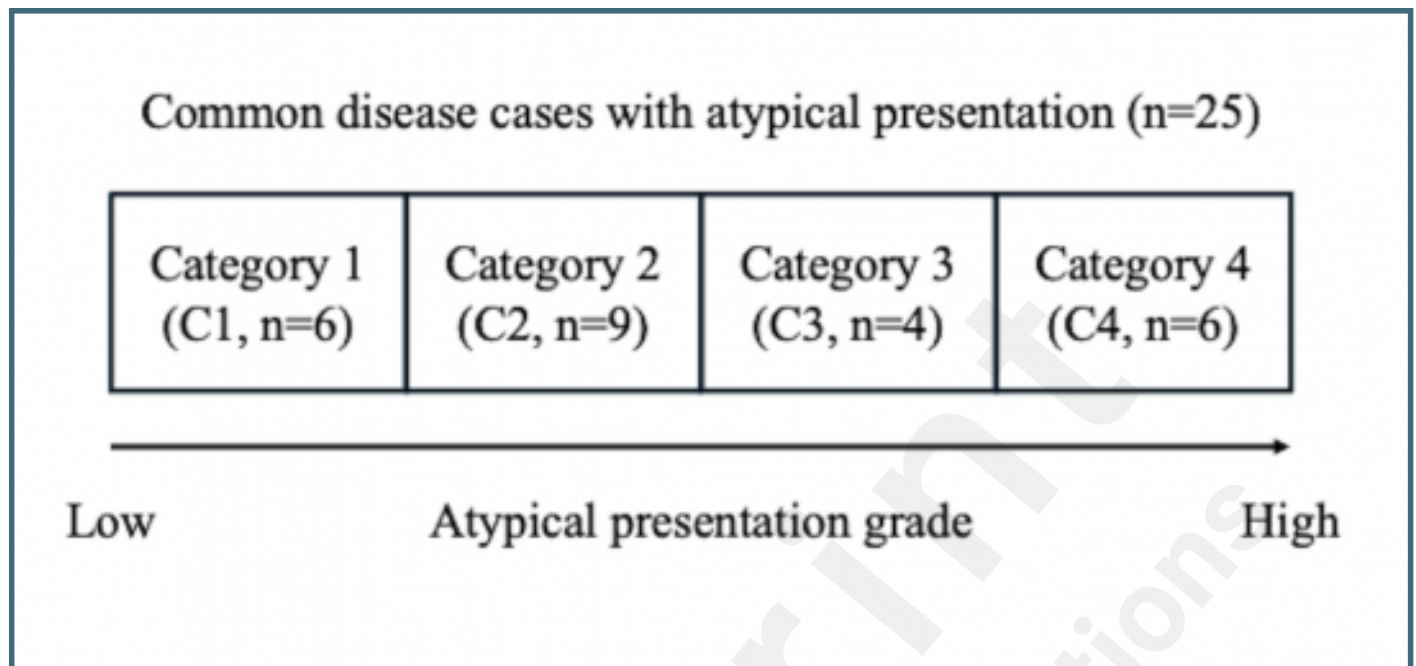
## Figures

Study motivation.

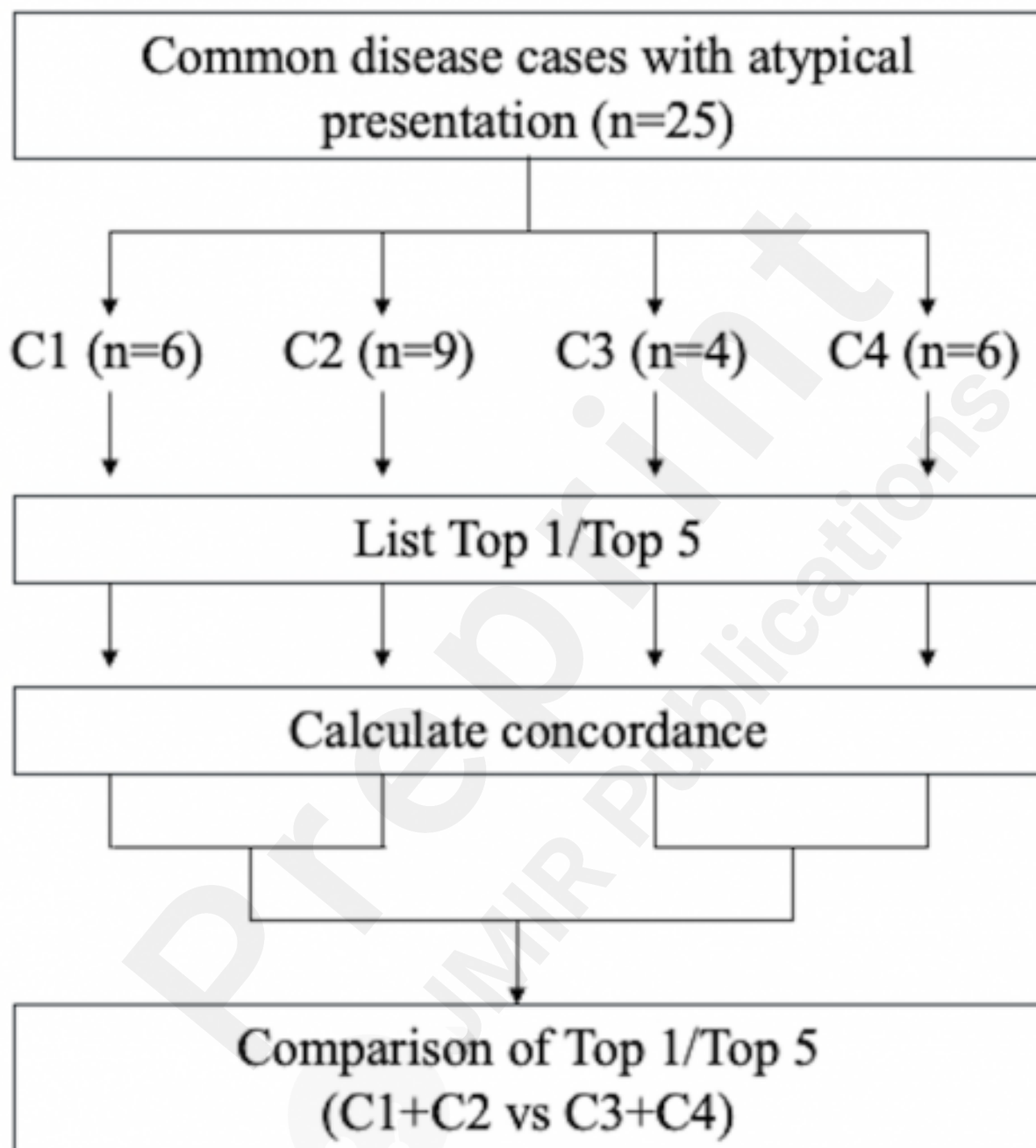




Categories of common disease with atypical presentation.



Study flow.



## **Multimedia Appendixes**

Supplementary file 1.

URL: <http://asset.jmir.pub/assets/4a54c66a1698e2d01483a52a1c976f09.docx>

Transcript of the conversation with ChatGPT and the answers to all the questions.

URL: <http://asset.jmir.pub/assets/a8b68e5794f1672c9b0c6982e1cebd38.docx>

