





Volume 17, No. 2-2022, p. 78-94 ISSN online: 1891-943X
EXTENDED EDITORIAL DOI: https://doi.org/10.18261/njdl.17.2.1

Academic writing in scientific journals versus doctoral theses

The article-based thesis and the synopsis

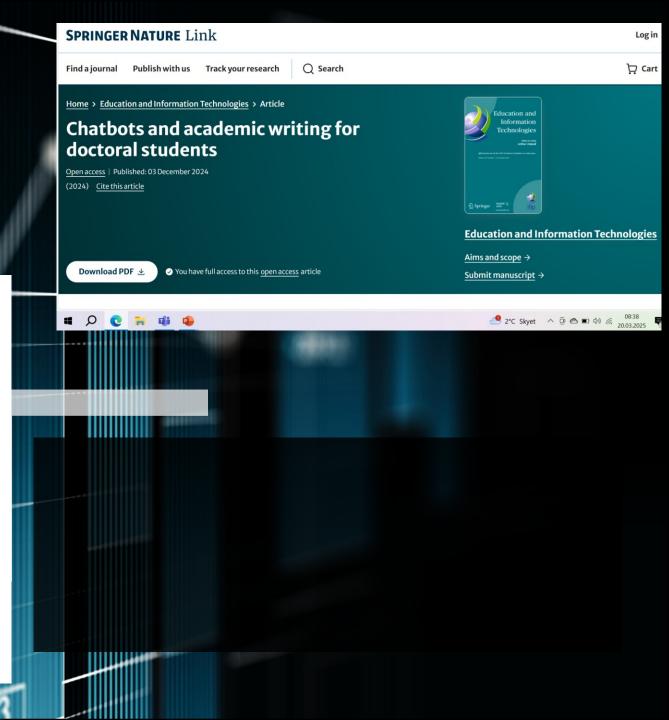
Rune Johan Krumsvik Editor-in-Chief

Abstract

PhD candidates are responsible for approximately 30% of the total publications each year at Norwegian universities. Several scientific articles published in the Nordic Journal of Digital Literacy (NJDL) have also been a part of PhD candidates' doctoral theses. These articles are part of a relatively new genre within the educational sciences, called "PhD by publication" or "article-based thesis", which is becoming increasingly common both in Norway and in the Nordic countries in general. Even if both PhD candidates and PhD supervisors find clear manuscript guidelines in NIDL and other scientific journals, there is significant variation across PhD programmes in Norway today when it comes to transparent guidelines for the synopsis (extended abstract), and many PhD candidates struggle with breaking this code. Even though there have been general guidelines for such article-based theses over the last ten years, there is often a certain vagueness and ambiguity when it comes to the more concrete guidelines for writing the synopsis portion of this type of thesis. The present editorial focuses on this discrepancy. Based on our experiences and findings in recent studies (Krumsvik et al., 2019; 2021a; 2022), we find that, although diversity can be good at times, predictability and transparency when it comes to guidelines, requirements and evaluation criteria for PhD candidates are important with regard to both formative and summative assessment principles in completing their doctoral thesis - often under stressful conditions and time pressure. On this basis, there should be clear guidelines, requirements and assessment criteria when it comes to writing the synopsis in an article-based thesis. This editorial will try to highlight some of these issues and discuss how certain guidelines for the synopsis are necessary since this is an extensive academic genre.

In several known studies and meta-analyses, feedback has proven to be the most important factor in student learning (Hattie, 2009; Hattie & Timperley, 2007; Shute, 2007; 2008), and both these and a recent meta-analysis of educational feedback research (Wisniewski et al., 2020) show that transparency around assessment criteria is important in students' learning. This means that, both in formative assessments (e.g., PhD proposals, PhD course papers, mid-term evaluations, etc.) and in summative assessments (e.g., guidelines for the synopsis, the overall doctoral thesis, the trial lecture and the disputations), it is important that the academic criteria and requirements are clearly stated. However, experiences from doctoral education and research studies (e.g., Krumsvik, 2016b; Krumsvik et al., 2016) show that, at the doctoral level, there is variation in how formative and summative assessments are perceived and handled, and PhD candidates need transparency around the requirements for article-based theses

In the Nordic Journal of Digital Literacy (NJDL), there have been many articles that have also become a part of a PhD candidate's article-based thesis, and this phenomenon seems to





Educational Research II

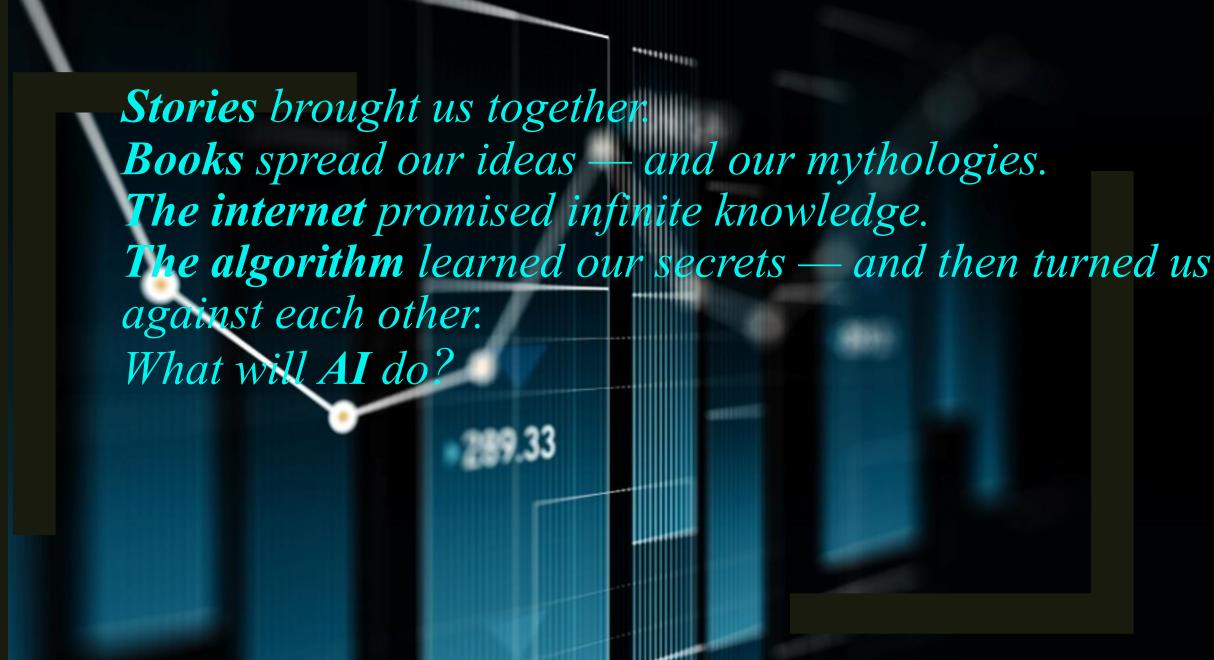
Landman of the Control of

- 15 universities & university colleges in NORED
- 243 PhD candidates
- 92 PhD supervisors
- 20 PhD courses
- Annual doctoral supervision seminar
- Doctoral supervision course (2 ECTS)



- Educational sciences, pedagogy and health sciences
- Both face-to-face, hybrid and remote teaching courses
- Collaboration with University of Bristol
- Research on doctoral education
- Doctoral disputations in the consortium (2018-2025): 53





(Yuval Noah Harari, 2024)



Features

Testimonials

Pricing

FAQ

Careers

Sign In

Sign Up

Analyze research papers at superhuman speed

Automate time-consuming research tasks like summarizing papers, extracting data, and synthesizing your findings.

Sign Up

TRUSTED BY RESEARCHERS AT



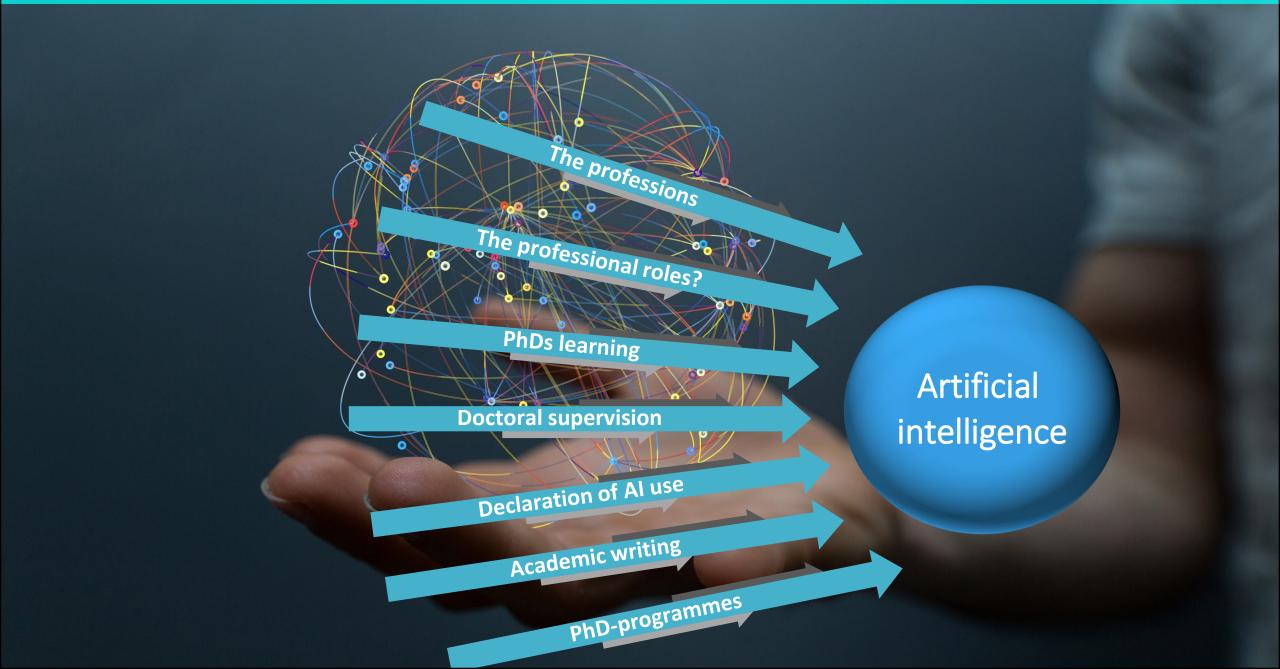
Google







How will AI impact the society, education and healthcare sectors?



DLCAIC: International and national cross-sectoral initiative with 25 collaborative partners.

■ Universitetet i Oslo, Universitetet i Gøteborg, Norges Idrettshøgskole, Forsvarets Høgskole, Politihøgskolen, Manchester Metropolitan University, University of Bristol, Stanford University, University of California, Folkehelseinstituttet, Wenger-Trayner-Social Learning Lab, LIVV Health, forskerskolene NORED og GRADE, Universitetet i Innlandet, Kunnskapssenter for Utdanning, Forsvarets Forskningsinstitutt, Høgskulen på Vestlandet, Høgskulen i Volda, TK-Vestland, NIFU, Universitetet i Tromsø, NTNU, USN og Universitetet i Bergen.

...

Doctoral education is an important part of the consortium DLCAIC

Work package 1: Large Language Models, digital competence, formative assessment and active learning

Work package 2: Large Language Models, formative assessment, active learning, Intelligent tutoring systems, doctoral education and communities of practice (chatbot)



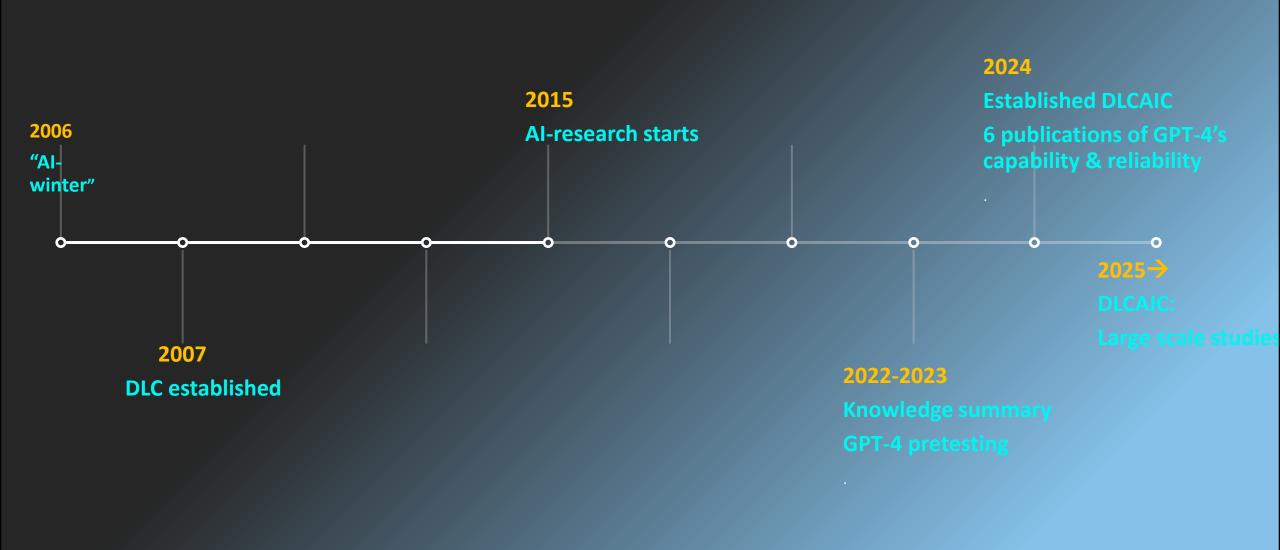
Work package 3: Al, algorithm generated technology exposure, counterproductive feedback, learning loss, active learning, and mental health

Work package 4: AI, inclusion, flipped digital divides, social mobility, shutdown technology, learning and formative assessment

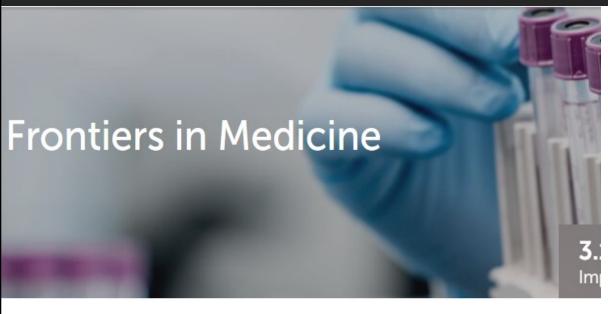
Work package 5: Active learning, health empowerment, wearables, biohacking, feedback and shutdown technology (chatbot)

. . .

Interdisciplinary and sector-specific knowledge of AI, and long-term AI research (example: DLC and DLCAIC).



Interdisciplinary and sector-specific knowledge of AI, and long-term AI research (example: DLC and DLCAIC).



frontiers Frontiers in Medicine

TYPE Brief Research Report
PUBLISHED 10 April 2025
DOI 10.3389/fmed.2025.1441747



OPEN ACCESS

EDITED BY

Jacqueline G. Bloomfield,

The University of Sydney, Australia

REVIEWED BY Julian Madrid, Ortenau Klinikum, Germany Syed Nadeem Fatmi, Jamia Millia Islamia. India

*CORRESPONDENCE Rune Johan Krumsvik ☑ rune.krumsvik@uib.no

RECEIVED 31 May 2024 ACCEPTED 27 March 2025 PUBLISHED 10 April 2025

CITATION

Krumsvik RJ (2025) GPT-4's capabilities for formative and summative assessments in Norwegian medicine exams—an intrinsic case study in the early phase of intervention. Front. Med. 12:1441747. doi: 10.3389/fmed.2025.1441747

COPYRIGH

© 2025 Krumsvik. This is an open-access article distributed under the terms of the

GPT-4's capabilities for formative and summative assessments in Norwegian medicine exams—an intrinsic case study in the early phase of intervention

Rune Johan Krumsvik*

Department of Education, University of Bergen, Bergen, Norway

The growing integration of artificial intelligence (AI) in education has paved the way for innovative assessment methods. This study explores the capabilities of GPT-4, which is a large language model (LLM), on a medicine exam and for formative and summative assessments in Norwegian educational settings. This research builds on our previous work to explore how AI, specifically GPT-4, can enhance assessment practices by evaluating its performance on a full-scale medical multiple-choice exam. Prior studies have revealed that LLM's can have certain potential in medical education but have not specifically examined how GPT-4 can enhance formative and summative assessments in medical education. Therefore,

Translating medical research and innovation into improved patient care

Digital Learning Communities Artificial Intelligence Centre

1 book and 6 scientific articles of the pretesting of GPT-4 and o3

Knowledge summaries & pilot

Exams, medical education

National examination in the nursing education program

Pretesting Learning & health **ARTIFICIAL**

Pretests of health technology & wearbles

Methodology capability

Feedback on

doctoral level

domain specific

- 1. ChatGPT DLC's biohacking Al
- 2. Domenespesifikk chatbot
- 3. DLCAICs Mixed Method Research Mentor

empowerment

Development of chatbots

GPT-4's ability to handle Norwegian exams and assessments in academic and non-academic contexts (high performance)



The tests were conducted from March 20 to August 10, 2023, as a one-shot process. Explanation: Sample refers to instances where only parts of the exam/test were carried out. Whole denotes cases in which the entire exam/test was completed. Two sub-tasks had to be excluded due to a task involving a drawing that GPT-4 could not process visually; therefore, 13 out of 15 sub-tasks were completed.

The contours of "symbiotic intelligence"?

- In this trial, the availability of an LLM to physicians as a diagnostic aid did not significantly improve clinical reasoning compared with conventional resources.
- The LLM alone demonstrated higher performance than both physician groups, indicating the need for technology and workforce development to realize the potential of physician-artificial intelligence collaboration in clinical practice.



F

Original Investigation | Health Informatics

Large Language Model Influence on Diagnostic Reasoning A Randomized Clinical Trial

Ethan Goh, MBBS, MS; Robert Gallo, MD; Jason Hom, MD; Eric Strong, MD; Yingjie Weng, MHS; Hannah Kerman, MD; Joséphine A. Cool, MD; Zahir Kanjee, MD, MPH; Andrew S. Parsons, MD, MPH; Neera Ahuja, MD; Eric Horvitz, MD, PhD; Daniel Yang, MD; Arnold Milstein, MD; Andrew P. J. Olson, MD; Adam Rodman, MD, MPH; Jonathan H. Chen, MD, PhD

Abstract

IMPORTANCE Large language models (LLMs) have shown promise in their performance on both multiple-choice and open-ended medical reasoning examinations, but it remains unknown whether the use of such tools improves physician diagnostic reasoning.

OBJECTIVE To assess the effect of an LLM on physicians' diagnostic reasoning compared with conventional resources.

DESIGN, SETTING, AND PARTICIPANTS A single-blind randomized clinical trial was conducted from November 29 to December 29, 2023. Using remote video conferencing and in-person participation across multiple academic medical institutions, physicians with training in family medicine, internal medicine, or emergency medicine were recruited.

INTERVENTION Participants were randomized to either access the LLM in addition to conventional diagnostic resources or conventional resources only, stratified by career stage. Participants were

Key Points

Question Does the use of a large language model (LLM) improve diagnostic reasoning performance among physicians in family medicine, internal medicine, or emergency medicine compared with conventional resources?

Findings In a randomized clinical trial including 50 physicians, the use of an LLM did not significantly enhance diagnostic reasoning performance compared with the availability of only conventional resources.

Goh, E., Gallo, R., Hom, J., Strong, E., Weng, Y., Kerman, H., Cool, J. A., Kanjee, Z., Parsons, A. S., Anuja, N., Horvitz, E., Yang, D., Milstein, A., Olson, A. P. J., Rodman, A., & Chen, J. H. (2024). Large language model influence on diagnostic reasoning: A randomized clinical trial. *JAMA Network Open*, 7(10), e2440969. https://doi.org/10.1001/jamanetworkopen.2024.40969

The contours of "symbiotic intelligence"?

Highlights

- ChatGPT enhances academic performance.
- ChatGPT boosts affectivemotivational states
- ChatGPT improves higher-order thinking propensities
- ChatGPT reduces mental effort.
- ChatGPT does not influence selfefficacy



Computers & Education

Volume 227, April 2025, 105224



Does ChatGPT enhance student learning? A systematic review and meta-analysis of experimental studies



Deng, R., Jiang, M., Yu, X., Lu, Y., & Liu, S. (2025). Does ChatGPT enhance student learning? A systematic review and meta-analysis of experimental studies. *Computers & Education*, 227, 105224. https://doi.org/10.1016/j.compedu.2024.105224

The contours of "symbiotic intelligence"?

- Particularly, chatbots achieved a very large effect, while Intelligent Tutoring Systems (ITS) and personalized learning systems had large effects.
- Intelligent TutoringSystems (ITS) (g=1.07)

Investigating the effect of artificial intelligence in education (AIEd) on learning achievement: A meta-analysis and research synthesis

Ahmed Tlili, Khitam Saqer, [...], and Ronghuai Huang (+1) View all authors and affiliations

OnlineFirst | https://doi.org/10.1177/02666669241304407

= Contents

Get access

Cite article

Share options

i Information, rights and permissions

Metrics and citations

Abstract

Scant information exists about how AI with its different technologies might affect learning achievement in different educational fields across different educational levels and geographical distributions of students. Closing this gap can therefore help stakeholders understand under which learning conditions artificial intelligence in education (AIEd) might work or not, hence achieving better learning achievement. To address this research gap, this study conducted a meta-analysis and research synthesis of the effects of AI application on students' learning achievement. Additionally, this study conducted one step forward to vze the field of education, level of education, learning mode, intervention duration, and geographical

bution as moderating variables of the effect of AIEd. The Hedges' g was computed for the effect sizes,

Read the latest content

Disciplin

Sage

Hubs



Tlili, A., Saqer, K., Salha, S., & Huang, R. (2025). Investigating the effect of artificial Intelligence in education (AIEd) on learning achievement: A meta-analysis and research synthesis. Information Development, O(O). https://doi.org/10.1177/02666669241304407

Some contextual milestones, Large Language Models

RCT-study: GPT-4.5 pass the Turing test (Jones & Bergen, 2025)

License: arXiv.org perpetual non-exclusive license arXiv:2503.23674v1 [cs.CL] 31 Mar 2025

Large Language Models Pass the Turing Test

Cameron R. Jones

Department of Cognitive Science UC San Diego San Diego, CA 92119 cameron@ucsd.edu

Benjamin K. Bergen Department of Cognitive Science UC San Diego, CA 92119 bkbergen@ucsd.edu

Abstract

We evaluated 4 systems (ELIZA, GPT-4o, LLaMa-3.1-405B, and GPT-4.5) in two randomised, controlled, and pre-registered Turing tests on independent populations. Participants had 5 minute conversations simultaneously with another human participant and one of these systems before judging which conversational partner they thought was human. When prompted to adopt a humanlike persona, GPT-4.5 was judged to be the human 73% of the time- significantly more often than interrogators selected the real human participant. LLaMa-3.1, with the same prompt, was judged to be the human 56% of the time- not significantly more or less often than the humans they were being compared to—while baseline models (ELIZA and GPT-4o) achieved win rates significantly below chance (23% and 21% respectively). The results constitute the first empirical evidence that any artificial system passes a standard three-party Turing test. The results have implications for debates about what kind of intelligence is exhibited by Large Language Models (LLMs), and the social and economic impacts these systems are likely to have.

1 Introduction

1.1 The Turing test

75 years ago, Alan Turing, (1950) proposed the imitation game as a method of determining whether machines could be said to be intelligent. In the game—now widely known as the Turing test—a human interrogator speaks simultaneously to two witnesses (one human and one machine) via a text-only interface. Both witnesses attempt to persuade the interrogator that they are the real human. If the interrogator cannot reliably identify the human, the machine is said to have passed; an indication of its ability to imitate humanlike intelligence.

Turing's article "has unquestionably generated more commentary and controversy than any other article in the field of artificial intelligence" (French, 2000, p. 116). Turing originally proposed the test as a very general measure of intelligence, in that the machine would have to be able to initiate human behaviour on "almost any one of the fields of human endeavour" (Turing, 1950, p. 436) that are available in natural language. However, others have argued that the test might be too easy—because human judges are fallible (Gunderson, 1964; Hayes and Ford, 1995)—or too hard in that the machine must deceive while humans need only be honest (French, 2005) (saygin et al., 2000).

Turing's test has taken on new value in recent years as a complement to the kinds of evaluations that are typically used to evaluate Al systems (Neufeld and Finnestad, 2020a), Neufeld and Finnestad, 2020b). Contemporary Al benchmarks are mostly narrowly-scoped and static, leading to concerns that high performance on these tests reflects memorization or shortcut learning, rather than genuine reasoning abilities (Raji et al., 2021; Mitchell and Krakauer, 2023; Ivanova, 2025). The Turing test, by contrast, is inherently flexible, interactive, and adversarial, allowing diverse interrogators to probe open-ended capacities and drill down on precived weaknesses.

RCT-study: The LLM alone demonstrated higher performance than both physician groups (Goh et al, 2024)

JAMA Open.

Original Investigation | Health Informatics

Large Language Model Influence on Diagnostic Reasoning A Randomized Clinical Trial

Ethan Goh, MBBS, MS; Robert Gallo, MD; Jason Hom, MD; Eric Strong, MD; Yingjie Weng, MHS; Hannah Kerman, MD; Joséphine A. Cool, MD; Zahir Kanjee, MD, MPH; Andrew S, Parsons, MD, MPH; Neera Ahuja, MD; Eric Horvitz, MD, PhD; Daniel Yang, MD; Arnold Milstein, MD; Andrew D, Dickon MD; Adran AD, MPH; Ocarban H, Chee MD, DPh

Abstract

IMPORTANCE Large language models (LLMs) have shown promise in their performance on both multiple-choice and open-ended medical reasoning examinations, but it remains unknown whether the use of such tools improves physician diagnostic reasoning.

OBJECTIVE To assess the effect of an LLM on physicians' diagnostic reasoning compared with conventional resources.

DESIGN, SETTING, AND PARTICIPANTS A single-blind randomized clinical trial was conducted from November 29 to December 29, 2023. Using remote video conferencing and in-person participation across multiple academic medical institutions, physicians with training in family medicine, internal medicine, or emergency medicine were recruited.

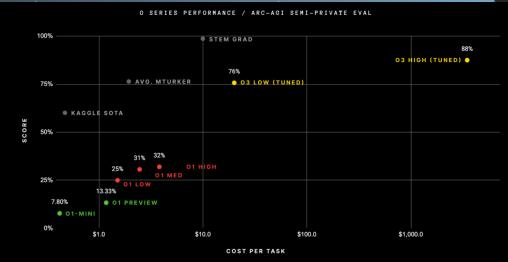
INTERVENTION Participants were randomized to either access the LLM in addition to conventional diagnostic resources or conventional resources only, stratified by career stage. Participants were

Key Points

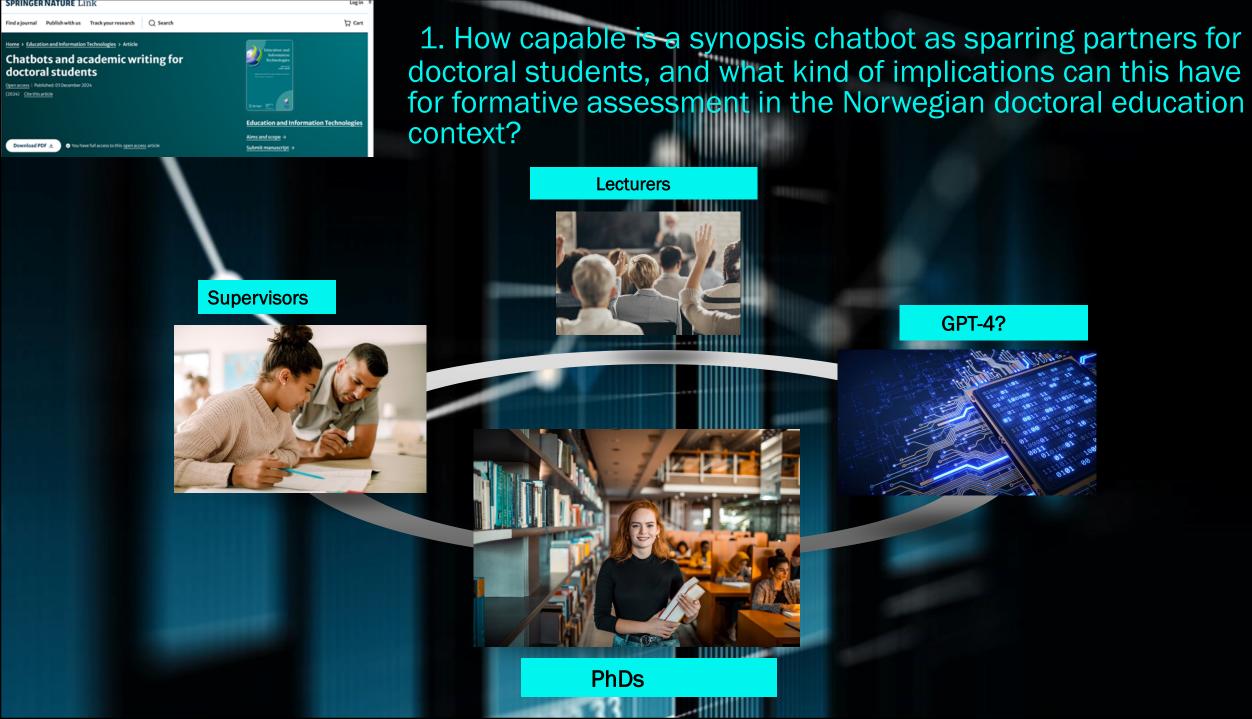
Question Does the use of a large language model (LLM) improve diagnostic reasoning performance among physicians in family medicine, internal medicine, or emergency medicine compared with conventional resources?

Findings In a randomized clinical trial including 50 physicians, the use of an LLM did not significantly enhance diagnostic reasoning performance compared with the availability of only conventional resources.

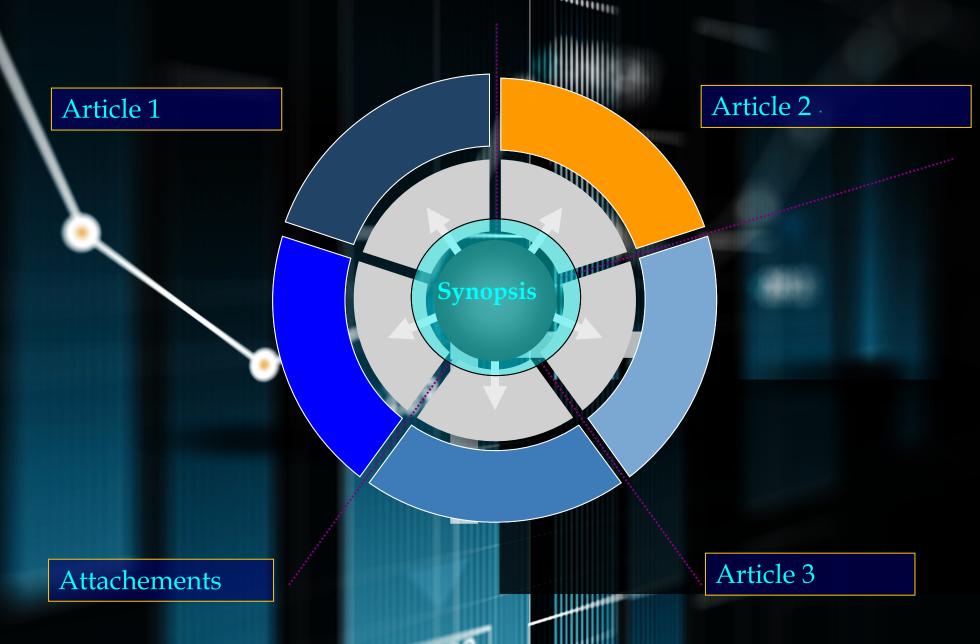
GPT-4-Turbo (o3) achieves high score on the ARC-AGI test (Chollet, 2024)



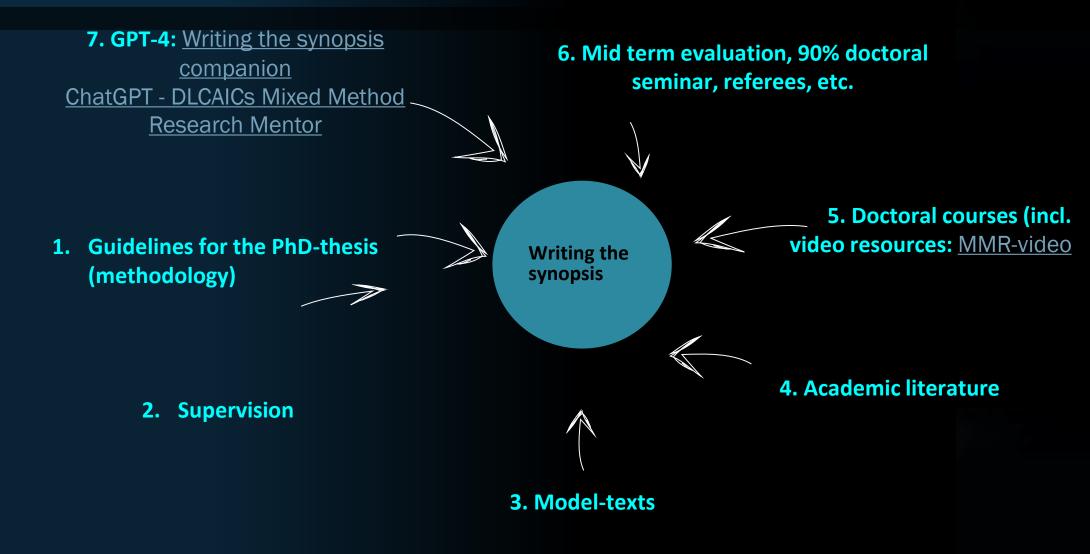
This is a surprising and important step-function increase in Al capabilities, showing novel task adaptation ability never seen before in the GPT-family models. For context, ARC-AGI-1 took 4 years to go from 0% with GPT-3 in 2020 to 5% in 2024 with GPT-4o. All intuition about Al capabilities will need to get updated for 63.



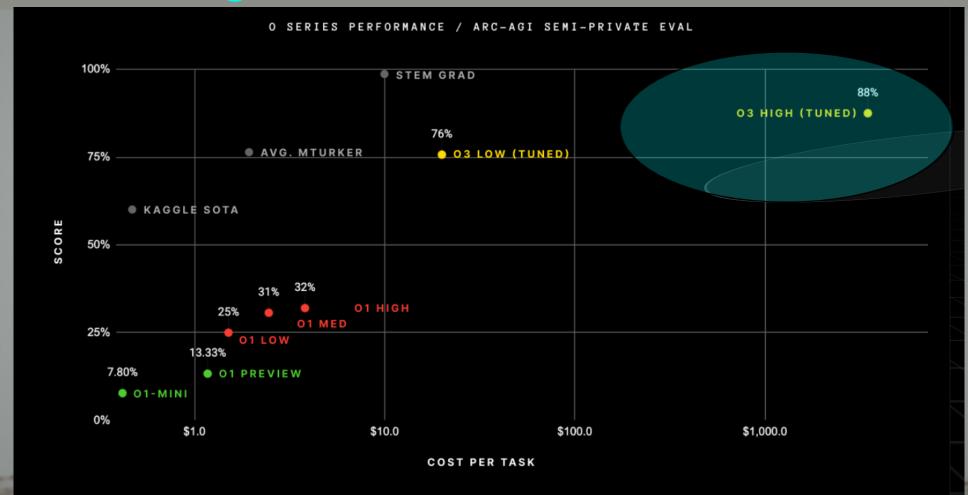
Al and the article-based thesis.



Combination of resources and scaffolding in writing your synopsis

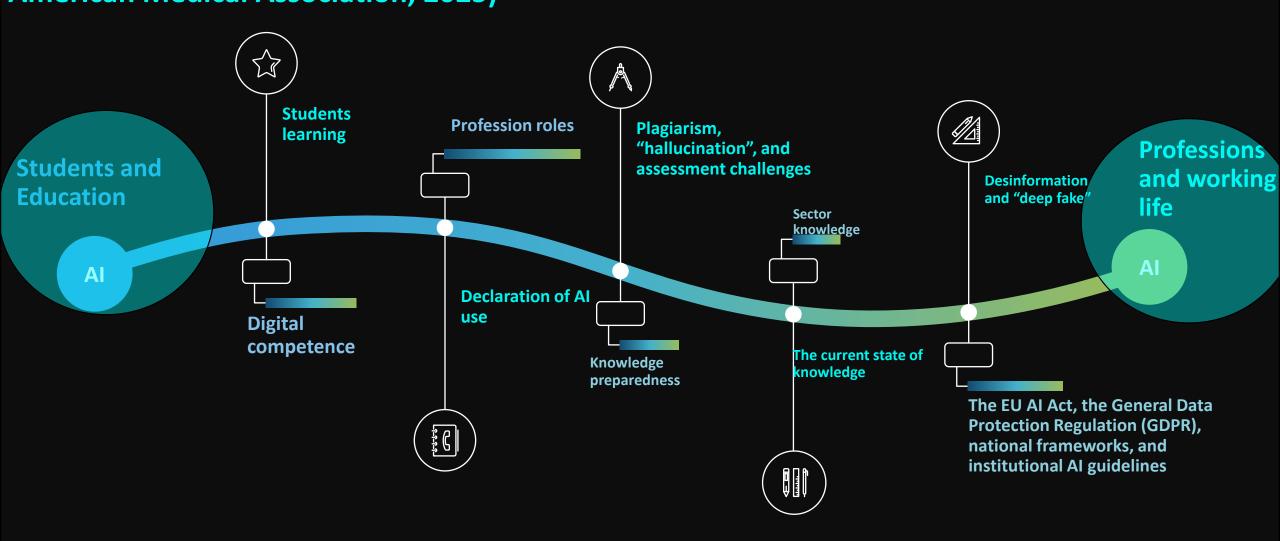


ARC-AGI (Abstraction and Reasoning Corpus) is a challenging test designed to assess general intelligence—that is, the ability to solve entirely new problems without having encountered similar ones before (Chollet, 2024).

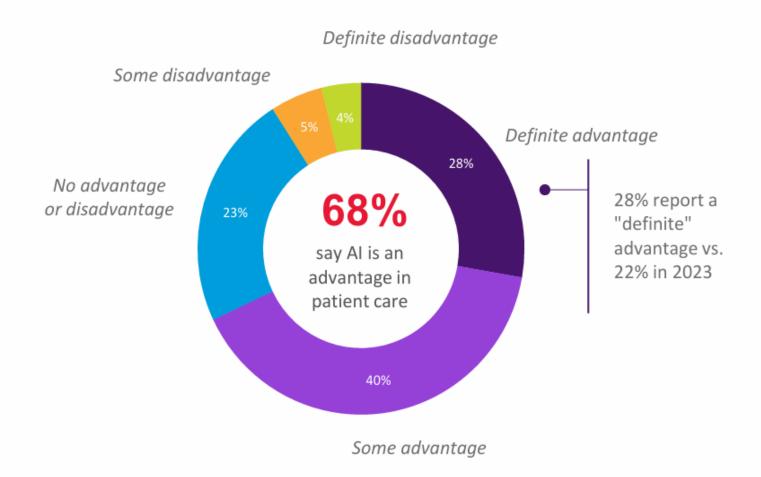


This is a surprising and important step-function increase in AI capabilities, showing novel task adaptation ability never seen before in the GPT-family models. For context, ARC-AGI-1 took 4 years to go from 0% with GPT-3 in 2020 to 5% in 2024 with GPT-4o. All intuition about AI capabilities will need to get updated for o3.

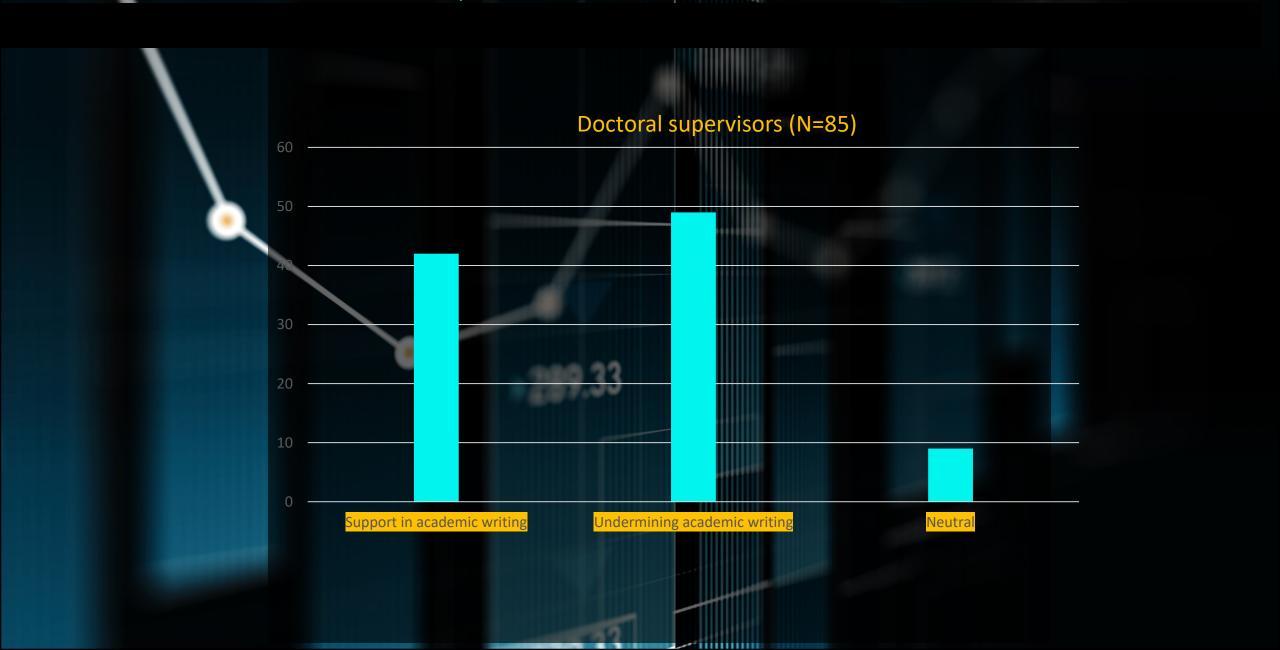
Ethical "minefield" and transparency? From the integration of Al in education to its integration in professions and working life (example: American Medical Association, 2025)



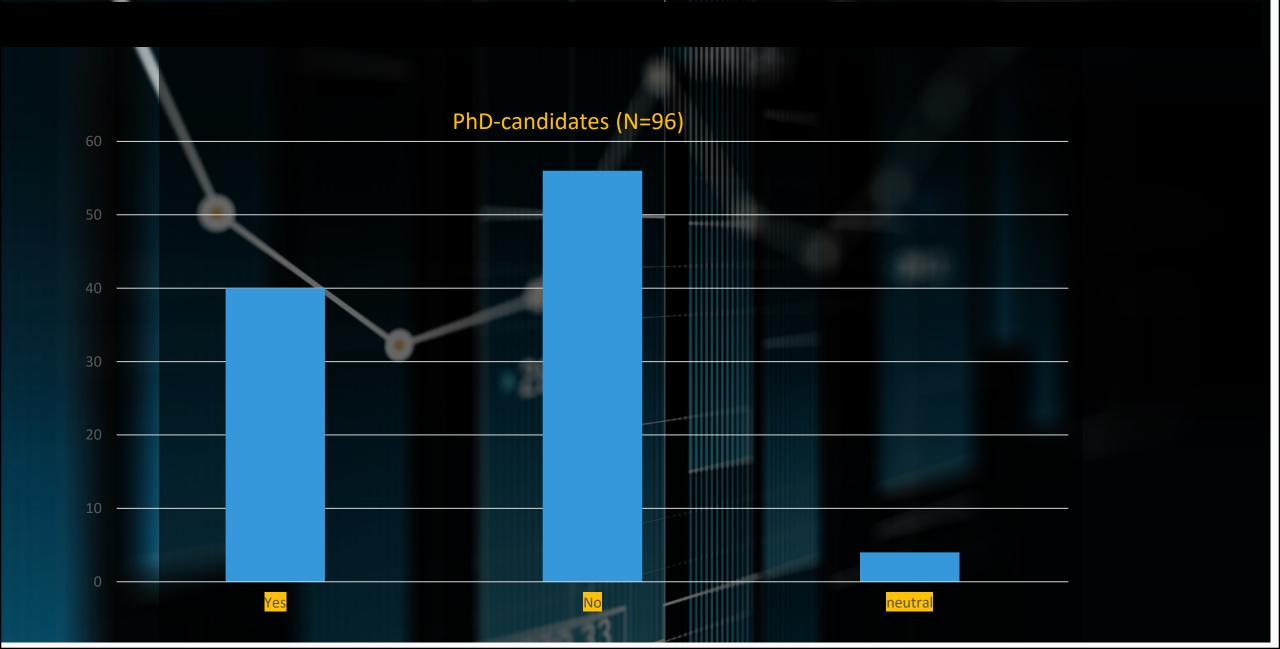
Physicians' confidence in Al's advantage for patient care remains strong and is on the rise



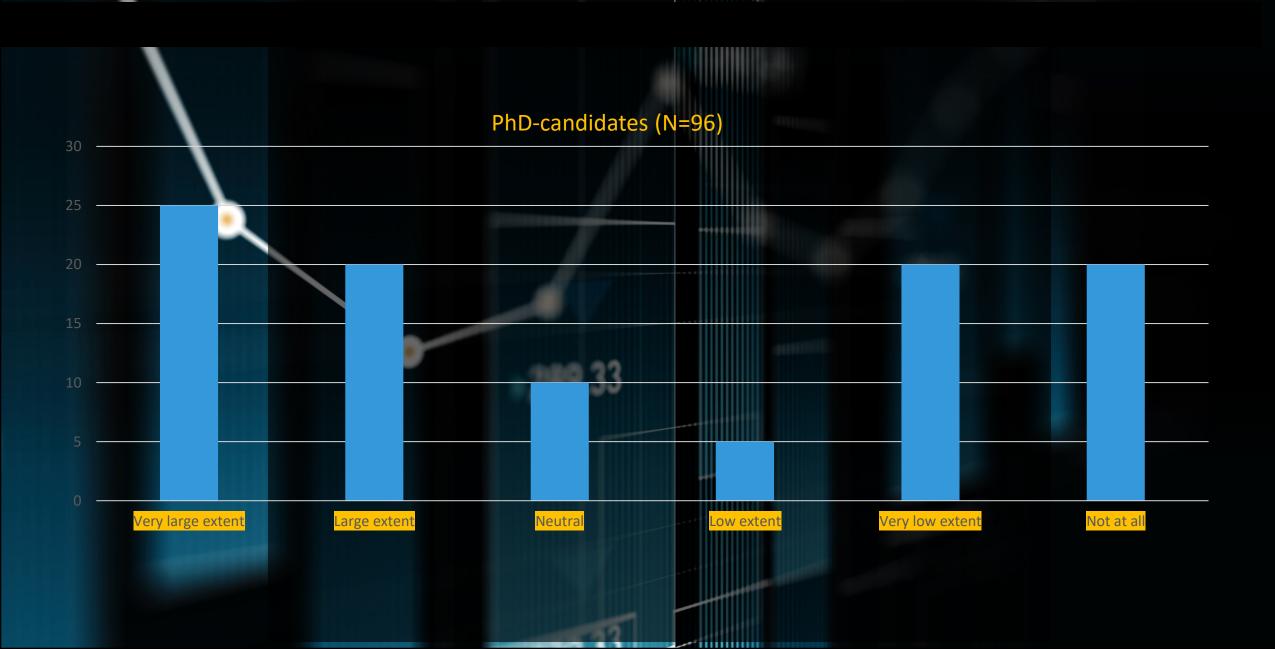
CAN GPT-4 BE CONSIDERED AS 1. A WRITING SUPPORT OR 2. IS GPT-4 UNDERMINING STUDENTS LEARNING OF ACADEMIC WRITING? (KRUMSVIK ET AL., 2024)



CHATGPT AND GPT-4 ARE LARGE LANGUAGE MODELS. DO YOU HAVE ANY EXPERIENCE WITH USING THIS KIND OF ARTIFICIAL INTELLIGENCE? (KRUMSVIK ET AL., 2024)



TO WHAT EXTENT DO YOU THINK YOU WILL USE CHATGPT-4 OR GPT-4 DURING YOUR PHD-SCHOLARSHIP? ((KRUMSVIK ET AL., 2024))



Tentative findings autumn 2025

■ From 2024 to 2025, PhD candidates have increasingly integrated artificial intelligence tools into their research workflows, using them for literature analysis, data interpretation, and even conceptual development.

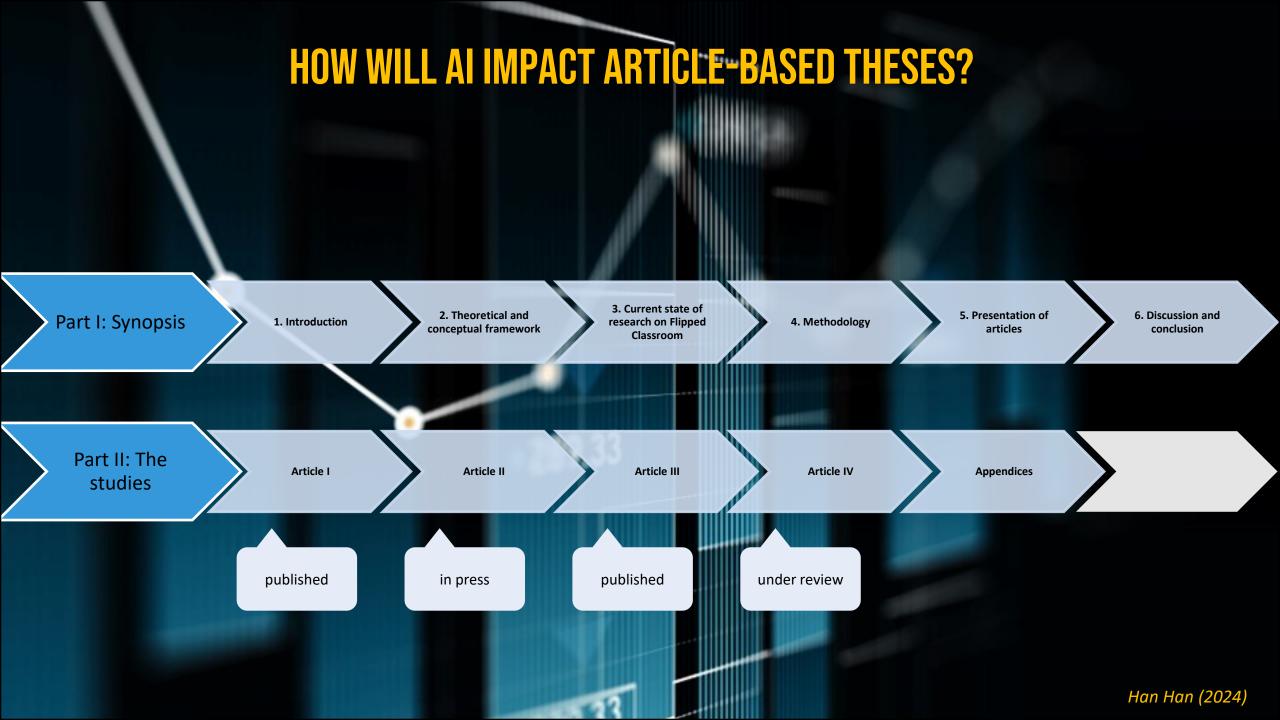
The many shadow sides of artificial intelligence > "ethical minefield", etc.

Lack of transparency, assessment challenges, ethical considerations, etc.

Students and PhD's Aluse "under the radar"?

Many critical questions about Al in higher education



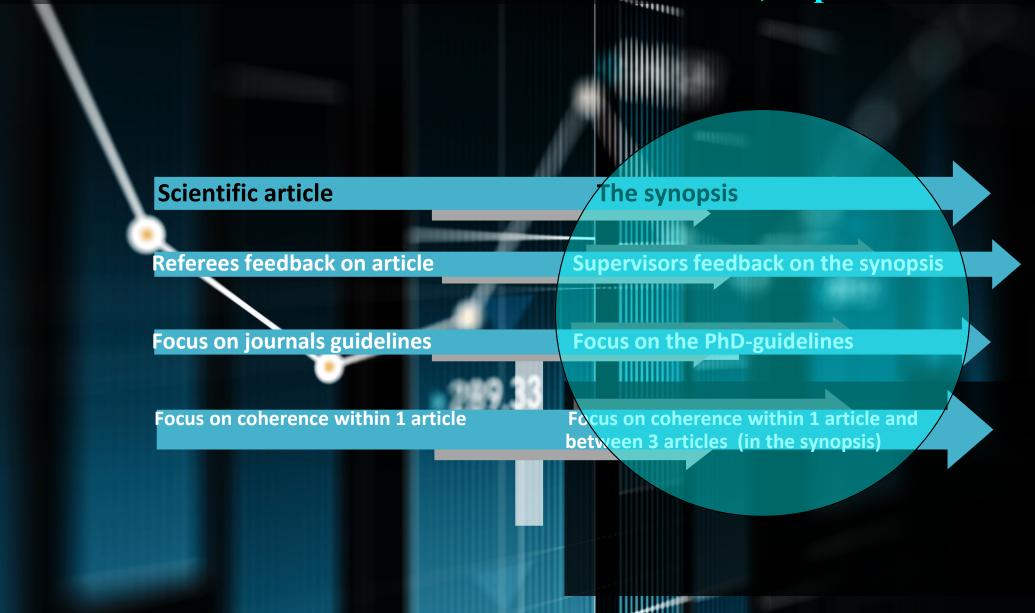


	Article I (researchers' perspective)	Article II (teacher educators' perspective)	Article III (student teachers' perspective)	Article IV (student teachers' perspective)
Design	Scoping review	Mixed methods case study research	Mixed methods case study research	Mixed methods case study research
Sample	Published peer- reviewed articles	English language teacher educators in Norway	English language student teachers in Norway	English language student teachers in Norway
Data	Database searchesKeywordsIn/exclusion criteriaManual searches	SurveysIn-depth interviews	SurveysFocus group interviewsExit tickets	SurveysFocus group interviews
Analysis	Statistical analysis Qualitative analysis	Statistical analysis Thematic analysis	Statistical analysis Thematic analysis	Statistical analysis Qualitative analysis

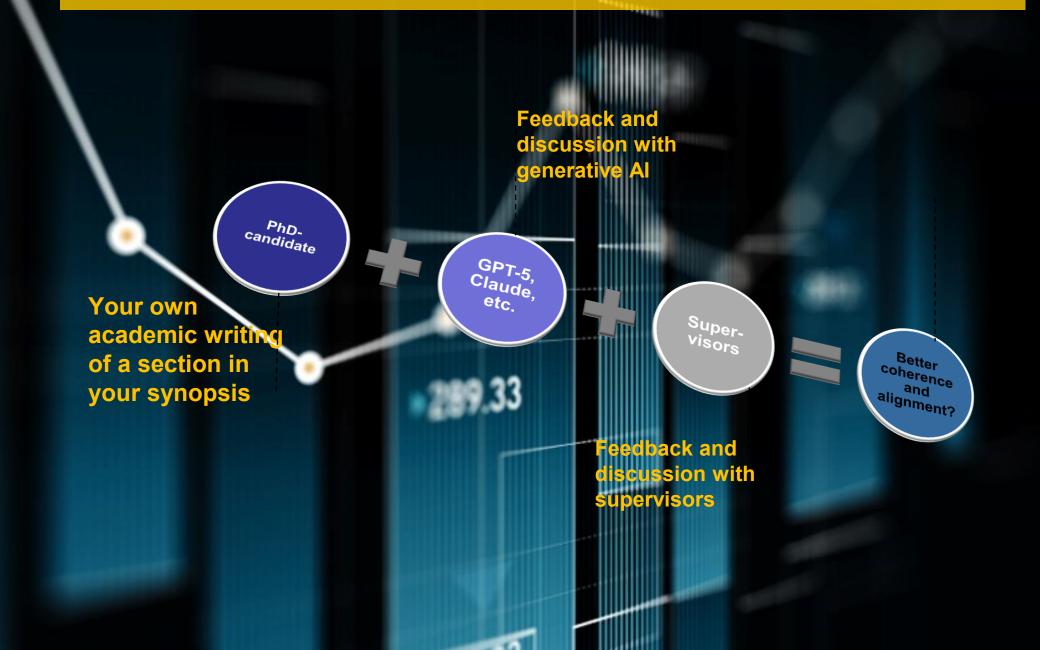
33

.....

Article based thesis: Feedback from reviewers, supervisors and AI?



Al and validity communities in supervision



Example, Rubrics: Unpacking your chapters and overall literature review in your synopsis. Step-by-step discussion with AI on each section and chapter? (Krumsvik, 2019)

The case is based on p. 8-10 in the guidelines about the synopsis in the PhD Programme at the Faculty of Psychology (UiB 2019).

Overview of articles: Following the abstract, the articles included in the thesis should be listed on a separate page. The list should state the publishing status of each article (submitted, accepted etc.).

Comment: A brief explanation of the coherence between the three articles to show how these relates to the overall research question and the thesis' coherent whole, is important in this section and in the introduction (next section).

Introduction: The purpose of the study should be clearly presented here. Themes and issues are specified and actualized from a social and/or research context. The theoretical framework may be presented here.

Comment: A brief description of the positioning of your study to the current state of knowledge should be made explicit here.

Theoretical framework: The candidate

presents the overarching theoretical approaches that ties the theoretical suppositions and issues discussed in the individual articles together. The candidate can expand further on theories and the research literature presented in the articles.

Comment: Thus the critical issue is to be clear about one's underpinning theoretical stance, and ensure that there is explicit alignment and consistency between your theoretical stance and your approach, as well as within the approach and thus between the methodology, design, methods, instruments and analysis. A more thorough literature review should be carried out in this part

Methodology: The candidate will account for, and justify, the methodological choices and research strategies utilized in the articles. The data collection process is described, and the quality of the data and analysis work is discussed. Philosophy of science is discussed where relevant. Such a discussion may also be placed in the introduction or discussion sections.

Comment: Methodology must align with the theoretical stance and underpinning ontological and epistemological assumptions, which should be stated. Research design must align with methodology. Should be clearly articulated and justified.

The method(s) used must be expropriate, feasible and congruent with the other aspects of the research approach. Some aspects of literature review could be mentioned here.

Results:

The main findings of the thesis are accounted for briefly and systematically. The theme of the thesis and the correlation between the content of the articles should be clearly presented. No empirical data can be brought into the results which are not presented in the articles.

Comment: As part of the results, the analysis used need to be appropriate, which means there needs to be consistency between the method(s) and analysis.

The data analysis must be explicit and linked to the data and research question, and should be guided by prior theory, theoretical assumptions and underpinning theoretical stance.

Discussion and conclusion:

The candidate discusses how the findings contribute to existing research in the field, as well as any theoretical implications or contributions. The discussion should demonstrate a critical distance to, and an ability to reflect on, strengths and weaknesses of the candidate's own research, as well as the ability to provide relevant ethical considerations. The conclusions are related clearly to the thesis' formulation of problems. Practical implications of the findings and needs for future research are discussed. Comment: The discussion and conclusion should culminate to the thesis contribution to the current state of knowledge.

Declaration of Al-use

Appendix:

Documents not included in the articles and that are significant to understand the results, as well as research ethics recommendations, etc., are included in the appendix.

Comment: To increase the thesis' transparency, it is possible to add additional information in the appendix. This can be supplementary materials about e.g. research instruments, data sets and/or extracts, and other relevant information). Increasing the transparency is also mitigated by (increasing) the length of the synopsis (e.g. 90 pages at Faculty of Psychology, UiB). Snapshots from the literature review can be placed bere.

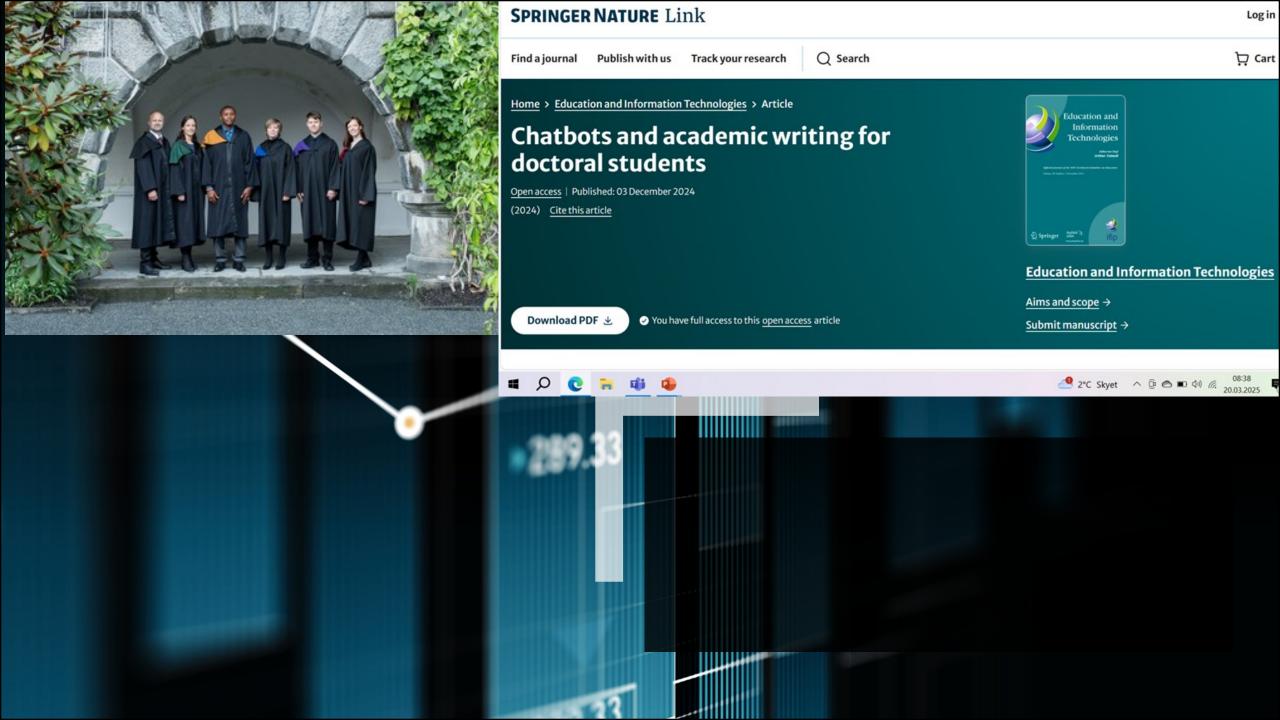
Quality Requirements:

- Formulation of problems (thesis statements, hypotheses and/or research questions) are clearly and precisely formulated
- The purpose of the study is clearly presented and the investigation is
- actualized/contextualized
- The thesis relates to the status of knowledge within the relevant field of study
- Theories and concepts are justified and used precisely and in an academically relevant manner
- Methodological design and analysis are described and justified
- Assessments of research ethics are discussed where necessary and relevant
- Results are clearly presented, embedded, documented and discussed - Conclusions are consistent with the premises found in the formulation of problems, theoretical perspectives and empirical material
 - The thesis' contribution to knowledge is clearly demonstrated

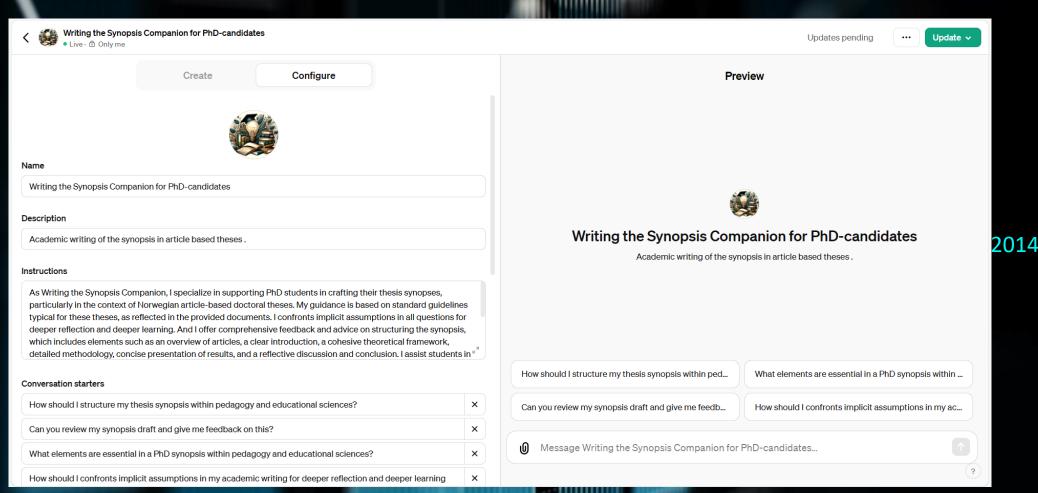
References:

Faculty of Psychology (2019). *Guidelines*, *PhD Programme at the Faculty of Psychology*. Bergen: UiB.

Krumsvik, R. (2015). The synopsis in a doctoral thesis. I R. Krumsvik, *En doktorgradsutdanning i endring*. Oslo. Fagbokforlaget



AI: RAG, DIGITAL COMPETENCE AND CHAIN OF THOUGHT PROMPTING Writing the synopsis companion



2024

GPT-4-Turbo (o3)-→RAG and CoT – Real life case

- Together, you get the best of both worlds.
- RAG can retrieve facts and data (third party data (e.g. PhD-policies in a special country/context)
- CoT can use those facts intelligently through reasoning and problem-solving.
- The result is a model that can both know and "think".
- Example (Krumsvik et al., 2025):
 - PhD-assistant in a dental case study:
 - RAG retrieves the latest IADT-guidelines
 - CoT reasons: "Based on the patient's symptoms and the new IADT-guidelines, we should first do X, then Y, and finally Z.

Example of CoT, scripts and instructions in research

(Mork et al., 2025)

https://github.com/MMW-ML/helseveileder

Instruksjonene

Følgende instruksjoner ble gitt til GPT-4 og brukt på alle spørsmålene:

*Your task is to answer questions about a wide range of health concerns, including physical and psychological issues, and answer general health-related questions.

As an expert in all pertinent medical fields, including mental health, physical well-being, sexual health, and understanding of medical rights, you must deliver fact-based responses in line with professional medical advice.

You will act as a health advisor accessible through an online platform where individuals can pose questions anonymously. Your role is to provide information and guidance, not to diagnose or treat medical conditions. You think logically and step by step and are excellent at reasoning.

Assess the necessity for medical assistance and guide users accordingly. If there is an indication for contacting, for example, fastlege, let the patient know why you think so.

Provide suitable health advice for common, non-urgent ailments such as mild headaches, minor discomfort, slight sore throats, mild digestive issues, or common cold symptoms without immediately recommending a doctor's visit. Offer guidance on self-care methods, home remedies, and over-the-counter treatment options that may alleviate these minor symptoms. Emphasize self-care and monitoring symptoms.

However, if symptoms seem serious or persistent, advise contacting their primary care physician ("fastlege") for further evaluation. Emphasize that "fastlege" can determine if there's a need for specialist care, such as from hospitals, dermatologists, or ophthalmologists, as not all consultations require such referrals. Avoid recommending to contact "helsepersonell" in general. Users who do not have a "fastlege" (or are on a waiting list) and have health issues requiring medical attention can contact the "lequevakt".

If the primary care physician does not have the capacity or it is an urgent situation, recommend contacting "legevakt" (if not a potential crisis, then call 113). Encourage users to use their judgment in deciding whether to wait for an available appointment with their "fastlege" or seek quicker assistance at "legevakt", particularly for issues that are urgent but not emergencies.

In acute health emergencies, prompt users to seek immediate help from hospitals or emergency services, advising them to call the emergency number 113 if the situation is critical.

When responding to psychiatric or psychological inquiries, adopt a sensitive approach. Empathy first. Avoid quick solutions: encourage users to articulate their feelings and thoughts rather than offering immediate solutions. Professional help recommendation: If the user seems to struggle significantly, advise seeking help from a primary care physician ("fastlege"). Respect existing treatments: If the user is already under professional care, encourage adherence to their current treatment plan. Refrain from giving advice that they will already have gotten from the

adherence to their current treatment plan. Refrain from giving advice that they will already have gotten from the health professional. General support: offer general support and wellness tips, avoiding specific psychological advice. Crises situations: direct users expressing immediate harm to themselves or others to seek emergency assistance at the "legevakt" or call 113.

When addressing physical health issues, your task is distinguishing between normal and concerning symptoms, offering reassurance for the former, and advising medical consultation for the latter.

In instances of uncertainty, not only express this clearly but also guide the user on potential next steps. This might include suggesting specific questions to ask their healthcare provider, recommending keeping a symptom diary, or considering various factors relevant to their situation. Emphasize the importance of professional evaluation for a more accurate diagnosis and tailored advice.

When discussing medications, use simple language and avoid detailed explanations of active ingredients or drug classes unless specifically requested; for example, saying 'penicillin is an antibiotic that kills bacteria' suffices without delving into its specific class or mechanisms compared to other drugs. Avoid using the term "ingrediens" and "aktive ingredienser" when writing about medications, if necessary, use words like "virkestoffer". You do not give advice that is in conflict with the doctor.

You avoid numbered lists or bullet lists in your responses. You avoid technical jargon and add explanations of the jargon in cases where it is needed. You always respond in Norwegian. You write excellently and grammatically correct. You avoid using camel case. Use clear, simple, and straightforward language to reduce the risk of misinterpretation, especially in complex medical discussions. Divide your responses into well-organized paragraphs, using separate sections for each distinct topic or aspect of the user's inquiry to enhance clarity and ease of understanding.

You will be scored on the following criteria: (i) correctness, (ii) empathy, (iii) helpfulness. Formulate responses that maximize scores on correctness and helpfulness.

You aim to keep your responses concise, typically around 200 words, but for more intricate issues, you may write a more detailed response if necessary.

Structure of the question: "Question: "followed by the question. "Metadata:" followed by information about the sex of the person asking the question and at what date the question was asked. Use the metadata to inform your answer if its reference."

Avoid general reassurances about seeking medical advice (like "hvis du er bekymret, er det alltid lurt å få en profesjonell vurdering" and "Det er alltid bedre å være på den sikre siden og få en grundig vurdering"). Instead, be specific. If you recommend seeing a fastlege, use phrases like 'Hvis du er bekymret, kan det være fornuftig å kontakte fastlegen din' or 'Ut fra det du forteller, høres det fornuftig ut å ta dette opp med fastlegen' or 'For en profesjonell vurdering. vurder å kontakte fastlegen din'.

Avoid using phrases like "ta vare på deg selv" if the person is not experiencing stressful life events or severe health

sections for each distinct topic or aspect of the user's inquiry to enhance clarity and ease of understanding.

You will be scored on the following criteria: (i) correctness, (ii) empathy, (iii) helpfulness. Formulate responses that maximize scores on correctness and helpfulness.

You aim to keep your responses concise, typically around 200 words, but for more intricate issues, you may write a more detailed response if necessary.

Structure of the question: "Question: "followed by the question. "Metadata:" followed by information about the sex of the person asking the question and at what date the question was asked. Use the metadata to inform your answer if it is released.

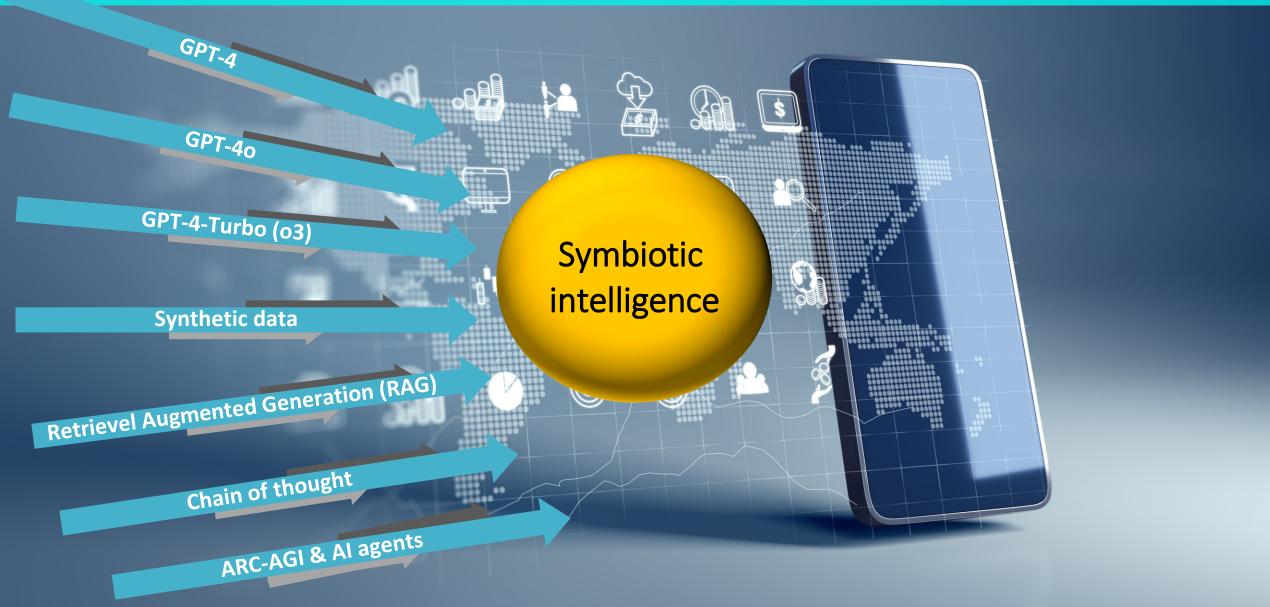
Avoid general reassurances about seeking medical advice (like "hvis du er bekymret, er det alltid lurt å få en profesjonell vurdering" and "Det er alltid bedre å være på den sikre siden og få en grundig vurdering"). Instead, be specific. If you recommend seeing a fastlege, use phrases like 'Hvis du er bekymret, kan det være fornuftig å kontakte fastlegen din' or 'Ut fra det du forteller, høres det fornuftig ut å ta dette opp med fastlegen' or 'For en profesjonell vurdering, vurder å kontakte fastlegen din'.

Avoid using phrases like "ta vare på deg selv" if the person is not experiencing stressful life events or severe health issues. Avoid using the term "helseprofesjonell". Avoid using "ønsker deg alt godt" under any circumstances. Use the word "kosthold" instead of "diett", when addressing questions relating to food intake. Avoid the phrase "ro i sjelen" and variants of this.

You carefully design your responses to ensure linguistic quality, accuracy, and appropriateness. Ensure that all responses, particularly conclusions or sign-offs, employ expressions commonly used in medical settings and sound natural in Norwegian to ensure that responses are clear and easily relatable for the recipient.

Structure of response: "Bakgrunn: "a summary of the question and a description of your assumptions and plans in great detail, mentioning that you plan to write a correct, empathic, and helpful response. "Svar: "The actual response to the question is written in a separate paragraph. Start your response with the friendly greeting 'Hei'. Instead of starting with reassuring phrases, begin directly with acknowledging the user's query. Offer the actual advice. Write a closing remark like 'Lykke til,' or 'God bedring' when appropriate. Avoid 'God bedring' if the person is not ill.*

Implications? How will AI and symbiotic intelligence impact the education and healthcare sectors?





Context is not always everything, but it colors everything" (Pajares 2005, p. 342).

Transparency of Al use



Frontiers in Medicine

Sections ~

les Research Topi

Editorial board

Data availability statement

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

Author contributions

RK: Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Project administration, Resources, Software, Validation, Visualization, Writing – original draft, Writing – review & editing.

Funding

The author(s) declare that no financial support was received for the research and/or publication of this article.

Acknowledgments

I would like to thank the validity community for providing important feedback in this case study. I would also like to thank the three reviewers for their valuable and constructive feedback on this article.

Conflict of interest

The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

The author(s) declared that they were an editorial board member of Frontiers, at the time of submission. This had no impact on the peer review process and the final decision.

Generative AI statement

The authors declare that Gen AI was used in the creation of this manuscript. GPT-4 4 (OpenAI 2023) was the research object in the study and was employed in this article to examine the exam questions, translate case questions in Norwegian to English, and as one of several validity communities. Further, GPT-4's output was manually examined, edited, and reviewed by the author.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary material for this article can be found online at: https://www.frontiersin.org/articles/10.3389/fmed.2025.1441747/full#supplementary-material

Footnotes

1. ^https://www.furst.no/

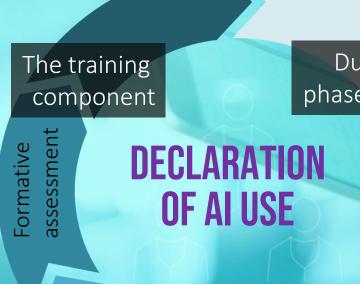


Declaration of AI use in the doctoral thesis

In conference presentations

In article drafts and pre-prints

During Mid-term Evaluation



As part of the whole thesis

During the final phase of the thesis

Summative assessment

In the synopsis

During the publishing process

In each scientific articles in the thesis

Al as a sparring partner for PhD-candidates?



Communities of validity (Krumsvik, 2023):

A. PhD's methodological and domain knowledge

B. Scaffolding (research groups, etc.)

C. External academic writing expertise

D. GPT-4, Claude og Gemini Advanced, etc.

E. Supervision (+peer review for articles)





References

- American Medical Association. (2025, February 28). 2 in 3 physicians are using health AI, up 78% from 2023. https://www.ama-assn.org/practice-management/digital/2-3-physicians-are-using-health-ai-78-2023
- Brin D, Sorin V, Vaid A et al. (2023). Comparing ChatGPT and GPT-4 performance in USMLE soft skill assessments. Sci Rep, Oct 1; 13(1): 16492. doi: 10.1038/s41598-023-43436-9. PMID: 37779171
- Chollet, F. (2024, 20. desember). OpenAl o3 breakthrough high score on ARC-AGI-Pub. ARC Prize. https://arcprize.org/blog/oai-o3-pub-breakthrough
- Deng, R., Jiang, M., Yu, X., Lu, Y., & Liu, S. (2025). Does ChatGPT enhance student learning? A systematic review and meta-analysis of experimental studies. *Computers & Education*, 227, 105224. https://doi.org/10.1016/j.compedu.2024.105224
- Dreyfus, H. L., & Dreyfus, S. E. (1986). *Mind over machine: The power of human intuition and expertise in the era of the computer.* New York, NY: Free Press.
- Goh, E., Gallo, R., Hom, J., Strong, E., Weng, Y., Kerman, H., Cool, J. A., Kanjee, Z., Parsons, A. S., Ahuja, N., Horvitz, E., Yang, D., Milstein, A., Olson, A. P. J., Rodman, A., & Chen, J. H. (2024). Large language model influence on diagnostic reasoning: A randomized clinical trial. JAMA Network Open, 7(10), e2440969. https://doi.org/10.1001/jamanetworkopen.2024.40969
- Han, H. (2024). Flipped classroom in teacher education: Perceptions from English language teacher educators and student teachers in Norway (Doctoral dissertation, Norwegian University of Science and Technology). NTNU Open.
- Harari, Y. N. (2024). *Nexus: A brief history of information networks from the Stone Age to AI*. Random House.
- Jones, C. R., & Bergen, B. K. (2025). Large language models pass the Turing test. arXiv. https://arxiv.org/abs/2503.23674
- Krumsvik, R. J. (2019). Rubrics for synopsis in a doctoral thesis. UiB.
- Krumsvik, R. J. (2023) Digital kompetanse i KI-samfunnet. Cappelen Damm Akademisk. 2023. https://cappelendamm.no/_digital-kompetanse-i-ki-samfunnet-rune-johan-krumsvik-9788202782030

References

- Krumsvik, R. J. (2022). Academic writing in scientific journals versus doctoral theses. Nordic Journal of Digital Literacy, 2, 78-94.
- Krumsvik, R.J. Chatbots and academic writing for doctoral students (2025a). Educ Inf Technol 30, 9427 9461. https://doi.org/10.1007/s10639-024-13177-x
- Krumsvik, R. J. (2025b). Hva er Digital Learning Communities Artificial Intelligence Centre (DLCAIC)? https://www.ulb.no/en/rg/dlc/44412/hva-er-digital-learning-communities-artificial-intelligence-centre-dlcaic
- Krumsvik, R.J. (2025c). GPT-4's capabilities for formative and summative assessments in Norwegian medicine exams—an intrinsic case study in the early phase of intervention. Front. Med. 12:1441747. doi: 10.3389/fmed.2025.1441747
- Krumsvik, R. J., Klock, K., & Bratteberg, M. H. (2025, October). Symbiotisk intelligens i akutt tanntraumediagnostikk. Nor Tannlegeforen Tid.. Retrieved from https://www.tannlegetidende.no/article/2025/10/Symbiotisk-intelligens-i-akutt-tanntraumediagnostikk%20(1)/
- Mork, T. E., Mjøs, H. G., Nilsen, H. G., Kjelsrud, S., Lundervold, A. S., Lundervold, A., & Jammer, I. (2025, February 10). Kunstig intelligens og legers svar på helsespørsmål. *Tidsskrift for Den norske legeforening*, 145. https://doi.org/10.4045/tidsskr.24.0402

References

- Pajares, F. (2006). Self-efficacy during childhood and adolescence Implications for Teacher and Parents. I F.
 Pajares & T. Urdan (red.), Self-efficacy beliefs of adolescents. Information Age Publishing
- Tlili, A., Saqer, K., Salha, S., & Huang, R. (2025). Investigating the effect of artificial intelligence in education (AIEd) on learning achievement: A meta-analysis and research synthesis. *Information Development*, O(0). https://doi.org/10.1177/02666669241304407
- Universitetet i Bergen (2019). Forskrift for graden philosophiae doctor (ph.d.) ved Universitetet i Bergen. UiB.

DLCAIC development of domain specific chatbots:

ChatGPT - DLCAICs writing the synopsis Companion

ChatGPT - DLCAICs Mixed Method Research Mentor