

Webinar Series 75

Kiat Praktis Menulis SYSTEMATIC LITERATURE REVIEW dari Planning hingga Submit

NARASUMBER
Muhammad Roil Bilad
Associate Professor dari Universitas Brunei Darussalam

H-Index Scopus	H-Index Scholar	I10-Index Scholar
54	62	265

MODERATOR
Sugeng Priyanto, S.S., M.IP
Pustakawan UPT Perpustakaan dan UNDIP Press

FASILITAS

- e-Sertifikat
- Materi Baru
- Ilmu yang Bermanfaat

Selasa, 30 Juni 2026
08.15- 11.30 WIB
Zoom Meeting
Meeting ID: 938 8731 0340
Passcode: UNDIP3026

TERBUKA UNTUK UMUM & GRATIS

Bersinergi • Literasi • Inovasi • Layanan Prima

Perpustakaan Undip • UNDIP Press • 082135876098 • penerbit.undip.ac.id • digilib.undip.ac.id

Materi

- Pengantar Gen-AI dan Workspace Chat-GPT
- Pengantar SLR
- Demo Prompting SLR
- Prompting SLR dengan CustomGPT + Rayyan

Pengenalan Gen-AI: Hype, Legalitas, Etik

AI **meniru** kecerdasan manusia dalam mesin.

Model bahasa AI **memahami** dan **menghasilkan** teks seperti manusia.



Keunggulan & Tantangan Terbesar

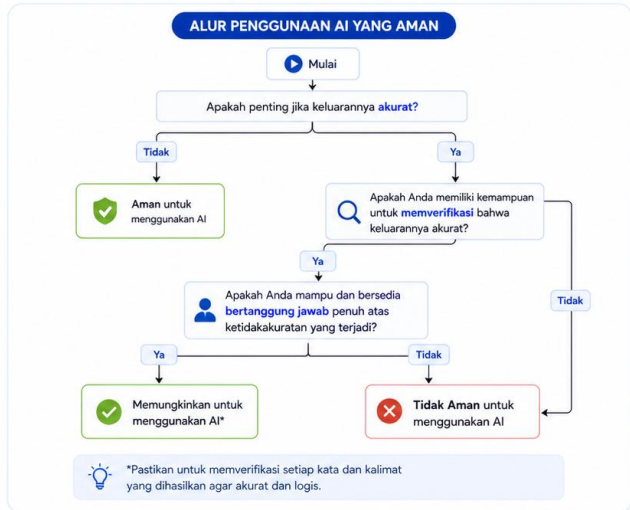
Text Prompt
Semua orang bisa melakukannya.

Text Understanding
Setiap orang bisa punya interpretasi berbeda.

Aspek Legal & Etis: Pengguna AI Bertanggungjawab

PANDUAN PENGGUNAAN
GENERATIVE ARTIFICIAL INTELLIGENCE (GenAI)
PADA PEMBELAJARAN DI PERGURUAN TINGGI

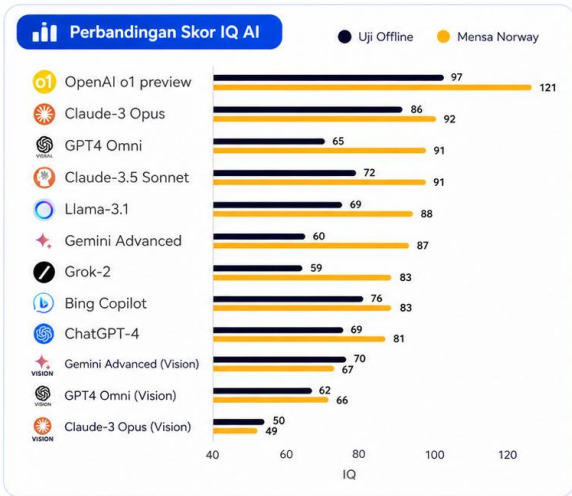
Aman • Etis • Akurat • Bertanggungjawab
Gunakan AI untuk belajar lebih baik dan berintegritas.



Panduan UNESCO terkait Alur Penggunaan ChatGPT yang Aman

- Verifikasi**: Selalu cek kebenaran informasi.
- Tanggung Jawab**: Anda bertanggung jawab atas hasil yang digunakan.
- Etika**: Gunakan untuk kebaikan dan sesuai aturan.
- Privasi**: Jaga data pribadi dan informasi sensitif.

Inteligensi AI Data: Oktober 2024



Model **OpenAI o1 preview** meraih skor **120** dalam uji IQ Mensa Norway.

RESEARCH ARTICLES

MACHINE LEARNING

Performance of a large language model on the reasoning tasks of a physician

Peter G. Bresnan¹, Thomas A. Buckley², Zohar Karpas³, Ethan Gao⁴, Evelyn Wu⁵, Pratik Shah⁶, Stephanie Cahani⁷, Rajeev Dasgupta⁸, Adam D. Heywood⁹, James A. Freed¹⁰, Andrew Chan¹¹, Daniel J. Mays¹², Jason Hwang¹³, Robert Gallo¹⁴, Lam G. McQuinn¹⁵, Heidi Mendenhall¹⁶, Christopher Lanza¹⁷, Misha Hershkov¹⁸, Matthew Dangelhoff¹⁹, Danielle Reibel²⁰, Daniel Restrepo²¹, Eric Horvitz²², Jonathan Chan²³, Ajay K. Menon²⁴, Adam Rodwin²⁵*

More than 65 years ago, complex clinical diagnostic reasoning cases were introduced as the gold standard for the evaluation of expert medical consulting systems, a standard that has held ever since. In this study, we report the results of a physician evaluation of a large language model (LLM) on challenging clinical cases across five experiments with a baseline of hundreds of physicians. The first part of a real-world study comparing human expert and artificial intelligence (AI) second opinions in randomly selected patients in the emergency room of a tertiary academic medical center. In all experiments, the LLM outperformed physician baselines and displayed continued improvement from prior generations of AI clinical decision support. Our study suggests that LLMs have reached most benchmarks of clinical reasoning, motivating the urgent need for prospective trials.

Artificial intelligence (AI) diagnostic support tools have been studied since the 1970s, with a landmark paper published to Lewis and Lusted (1), who advocated for case-based benchmarks as an evaluation standard, a standard that has held for more than the past half century (2-4). In particular, the *New England Journal of Medicine* (NEJM) clinicalopathological case conference series has been seen as an aspirational goal post, tested by every differential diagnosis generator spanning general internal medicine, pediatric, and surgical specialties, and natural language synthesis checkers. Recently, large language models (LLMs) have consistently outperformed older models on these challenging cases, motivating their relevance to professional training needs, mathematics education, software engineering, and engineering problems (5-22).

However, recent studies of LLMs in medicine have focused on narrow diagnostic tasks or on context and operational clinical expertise (23, 24). More importantly, most studies of LLMs for diagnosis and management to date have lacked human physician baselines. This was highlighted in previous generations of technology because of the overall poor performance of computerized models on benchmarks. GPT-4 outperformed in LLMs and increasing "benchmark saturation"

It is necessary to establish human baselines and study clinically grounded tasks. Here, we comprehensively evaluated diagnostic and management reasoning capabilities of an advanced LLM (OpenAI GPT-4) across several diagnostic and management reasoning tasks with human performance from hundreds of physicians. We further studied LLM second opinions in a simulated triage setting on expert physician baselines on randomly selected patients in a major academic tertiary care emergency department in Boston, Massachusetts.

Results
Quality of differential diagnosis and testing plans on NEJM clinicalopathological conferences
We first evaluated *o1*-preview using the clinicalopathological conference (CPC) published by the NEJM, a standard for the evaluation of differential diagnosis generated since the 1950s (1). There was substantial agreement between the two physicians evaluating the quality of *o1*-preview differential diagnosis (agreement on 120/148 cases (81%), inter-rater reliability $\kappa = 0.63$). *o1*-preview included the correct diagnosis in 101 of 148 cases (68%). The first diagnosis suggested was the correct diagnosis in 51% of cases (91% CL, 44 to 61%). When expanding to also include potentially helpful or very close diagnoses, *o1*-preview was accurate on 87% (91% CL, 81 to 94%) of cases (Fig. 1A). We did not find evidence of a significant difference in performance before and after the processing chain for *o1*-preview (CPC accuracy before, 75.5% accuracy after, $P = 0.59$, table S1). In a subset of 68 cases from a prior study (5), *o1*-preview outperformed a human physician baseline in both top and top-10 accuracy (table S1). On 70 cases used to evaluate GPT-4 in a prior study (7), *o1*-preview showed a response with the most or a very close diagnosis in 66.4% of cases, compared with 51.4% of cases (GPT-4 $P = 0.01$, Fig. 2B). Overall, *o1*-preview and GPT-4 performed identically on 66/70 (94.3%) of cases, *o1*-preview outperformed GPT-4 on 1/70 (1.4%) of cases, and GPT-4 outperformed *o1*-preview on 3/70 (4.3%) of cases (Fig. 3).

We next evaluated the ability of *o1*-preview to select the case diagnostic test to test the NEJM CPC for a subset of 104 cases. Two physicians scored the suggested test plus provided by *o1*-preview (agreement on 2/133/104 cases (97%), $\kappa = 0.26$) with respect to the actual management of the patient described in the CPC. The proportion of agreements was high, but the κ was low as a result of severe class imbalance. In 82.5% of cases, *o1*-preview selected the correct test to order. In another 12% of cases, the chosen testing plan was judged by the two physicians to be helpful and in 1.5% of cases, it would have been unhelpful (Fig. 2C). Reference case sets and examples of model outputs are shown in tables S2 to S5.

Presentation of reasoning in NEJM Healer diagnostic cases
We first evaluated reasoning over the NEJM Healer curriculum (25) that were also evaluated in a prior study using GPT-4 (16). NEJM Healer cases are virtual patient episodes designed for the assessment of clinical reasoning (25). There was substantial agreement of *o1*-preview (R-IDEA) scores, a validated scoring scale for evaluating the core domains of deconstructing clinical reasoning (27)—between the two physicians (agreement on 70/102 cases, $\kappa = 0.82$, Fig. 3B). The *o1*-preview and GPT-4 performed identically on 66/70 (94.3%) of cases, *o1*-preview outperformed GPT-4 on 1/70 (1.4%) of cases, and GPT-4 outperformed *o1*-preview on 3/70 (4.3%) of cases (Fig. 3).

¹Department of Internal Medicine, Beth Israel Deaconess Medical Center, Boston, MA, USA; ²Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA; ³Department of Computer Science, Harvard University, Cambridge, MA, USA; ⁴Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA; ⁵Department of Pathology, Harvard Medical School, Boston, MA, USA; ⁶Department of Internal Medicine, Brigham Young University, Provo, UT, USA; ⁷Department of Emergency Medicine, Massachusetts General Hospital, Boston, MA, USA; ⁸Department of Emergency Medicine, Massachusetts General Hospital, Boston, MA, USA; ⁹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁰Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹¹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹²Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹³Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁴Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁵Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁶Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁷Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁸Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²⁰Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²¹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²²Department of Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA; ²³Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²⁴Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²⁵Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA. *Corresponding author: Email: adam.rodwin@massgeneral.org

Science 30 April 2024



ARTICLE IN PRESS

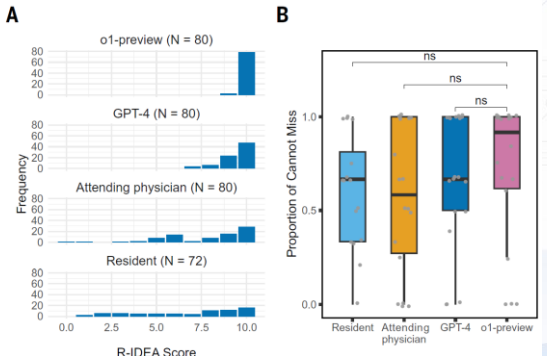


Fig. 3. Comparison of *o1*-preview, GPT-4, and physicians for clinical diagnostic reasoning. (A) Distribution of 312 R-IDEA scores stratified by respondents on 20 NEJM Healer cases. (B) Box plot of the proportion of cannot-miss diagnoses included in differential diagnosis for the initial triage presentation. The total sample size in this figure is 70, with 18 responses from attending physicians, GPT-4, and *o1*-preview, and 16 responses from residents. Two cases were excluded because the cannot-miss diagnoses could not be identified. ns, not significant.

RESEARCH ARTICLES

MACHINE LEARNING

Performance of a large language model on the reasoning tasks of a physician

Peter G. Bresnan¹, Thomas A. Buckley², Zohar Karpas³, Ethan Gao⁴, Evelyn Wu⁵, Pratik Shah⁶, Stephanie Cahani⁷, Rajeev Dasgupta⁸, Adam D. Heywood⁹, James A. Freed¹⁰, Andrew Chan¹¹, Daniel J. Mays¹², Jason Hwang¹³, Robert Gallo¹⁴, Lam G. McQuinn¹⁵, Heidi Mendenhall¹⁶, Christopher Lanza¹⁷, Misha Hershkov¹⁸, Matthew Dangelhoff¹⁹, Danielle Reibel²⁰, Daniel Restrepo²¹, Eric Horvitz²², Jonathan Chan²³, Ajay K. Menon²⁴, Adam Rodwin²⁵*

More than 65 years ago, complex clinical diagnostic reasoning cases were introduced as the gold standard for the evaluation of expert medical consulting systems, a standard that has held ever since. In this study, we report the results of a physician evaluation of a large language model (LLM) on challenging clinical cases across five experiments with a baseline of hundreds of physicians. The first part of a real-world study comparing human expert and artificial intelligence (AI) second opinions in randomly selected patients in the emergency room of a tertiary academic medical center. In all experiments, the LLM outperformed physician baselines and displayed continued improvement from prior generations of AI clinical decision support. Our study suggests that LLMs have reached most benchmarks of clinical reasoning, motivating the urgent need for prospective trials.

Artificial intelligence (AI) diagnostic support tools have been studied since the 1970s, with a landmark paper published to Lewis and Lusted (1), who advocated for case-based benchmarks as an evaluation standard, a standard that has held for more than the past half century (2-4). In particular, the *New England Journal of Medicine* (NEJM) clinicalopathological case conference series has been seen as an aspirational goal post, tested by every differential diagnosis generator spanning general internal medicine, pediatric, and surgical specialties, and natural language synthesis checkers. Recently, large language models (LLMs) have consistently outperformed older models on these challenging cases, motivating their relevance to professional training needs, mathematics education, software engineering, and engineering problems (5-22).

However, recent studies of LLMs in medicine have focused on narrow diagnostic tasks or on context and operational clinical expertise (23, 24). More importantly, most studies of LLMs for diagnosis and management to date have lacked human physician baselines. This was highlighted in previous generations of technology because of the overall poor performance of computerized models on benchmarks. GPT-4 outperformed in LLMs and increasing "benchmark saturation"

It is necessary to establish human baselines and study clinically grounded tasks. Here, we comprehensively evaluated diagnostic and management reasoning capabilities of an advanced LLM (OpenAI GPT-4) across several diagnostic and management reasoning tasks with human performance from hundreds of physicians. We further studied LLM second opinions in a simulated triage setting on expert physician baselines on randomly selected patients in a major academic tertiary care emergency department in Boston, Massachusetts.

Results
Quality of differential diagnosis and testing plans on NEJM clinicalopathological conferences
We first evaluated *o1*-preview using the clinicalopathological conference (CPC) published by the NEJM, a standard for the evaluation of differential diagnosis generated since the 1950s (1). There was substantial agreement between the two physicians evaluating the quality of *o1*-preview differential diagnosis (agreement on 120/148 cases (81%), inter-rater reliability $\kappa = 0.63$). *o1*-preview included the correct diagnosis in 101 of 148 cases (68%). The first diagnosis suggested was the correct diagnosis in 51% of cases (91% CL, 44 to 61%). When expanding to also include potentially helpful or very close diagnoses, *o1*-preview was accurate on 87% (91% CL, 81 to 94%) of cases (Fig. 1A). We did not find evidence of a significant difference in performance before and after the processing chain for *o1*-preview (CPC accuracy before, 75.5% accuracy after, $P = 0.59$, table S1). In a subset of 68 cases from a prior study (5), *o1*-preview outperformed a human physician baseline in both top and top-10 accuracy (table S1). On 70 cases used to evaluate GPT-4 in a prior study (7), *o1*-preview showed a response with the most or a very close diagnosis in 66.4% of cases, compared with 51.4% of cases (GPT-4 $P = 0.01$, Fig. 2B). Overall, *o1*-preview and GPT-4 performed identically on 66/70 (94.3%) of cases, *o1*-preview outperformed GPT-4 on 1/70 (1.4%) of cases, and GPT-4 outperformed *o1*-preview on 3/70 (4.3%) of cases (Fig. 3).

We next evaluated the ability of *o1*-preview to select the case diagnostic test to test the NEJM CPC for a subset of 104 cases. Two physicians scored the suggested test plus provided by *o1*-preview (agreement on 2/133/104 cases (97%), $\kappa = 0.26$) with respect to the actual management of the patient described in the CPC. The proportion of agreements was high, but the κ was low as a result of severe class imbalance. In 82.5% of cases, *o1*-preview selected the correct test to order. In another 12% of cases, the chosen testing plan was judged by the two physicians to be helpful and in 1.5% of cases, it would have been unhelpful (Fig. 2C). Reference case sets and examples of model outputs are shown in tables S2 to S5.

Presentation of reasoning in NEJM Healer diagnostic cases
We first evaluated reasoning over the NEJM Healer curriculum (25) that were also evaluated in a prior study using GPT-4 (16). NEJM Healer cases are virtual patient episodes designed for the assessment of clinical reasoning (25). There was substantial agreement of *o1*-preview (R-IDEA) scores, a validated scoring scale for evaluating the core domains of deconstructing clinical reasoning (27)—between the two physicians (agreement on 70/102 cases, $\kappa = 0.82$, Fig. 3B). The *o1*-preview and GPT-4 performed identically on 66/70 (94.3%) of cases, *o1*-preview outperformed GPT-4 on 1/70 (1.4%) of cases, and GPT-4 outperformed *o1*-preview on 3/70 (4.3%) of cases (Fig. 3).

¹Department of Internal Medicine, Beth Israel Deaconess Medical Center, Boston, MA, USA; ²Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA; ³Department of Computer Science, Harvard University, Cambridge, MA, USA; ⁴Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA; ⁵Department of Pathology, Harvard Medical School, Boston, MA, USA; ⁶Department of Internal Medicine, Brigham Young University, Provo, UT, USA; ⁷Department of Emergency Medicine, Massachusetts General Hospital, Boston, MA, USA; ⁸Department of Emergency Medicine, Massachusetts General Hospital, Boston, MA, USA; ⁹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁰Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹¹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹²Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹³Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁴Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁵Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁶Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁷Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁸Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ¹⁹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²⁰Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²¹Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²²Department of Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA; ²³Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²⁴Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA; ²⁵Department of Internal Medicine, Massachusetts General Hospital, Boston, MA, USA. *Corresponding author: Email: adam.rodwin@massgeneral.org

Science 30 April 2024



ARTICLE IN PRESS

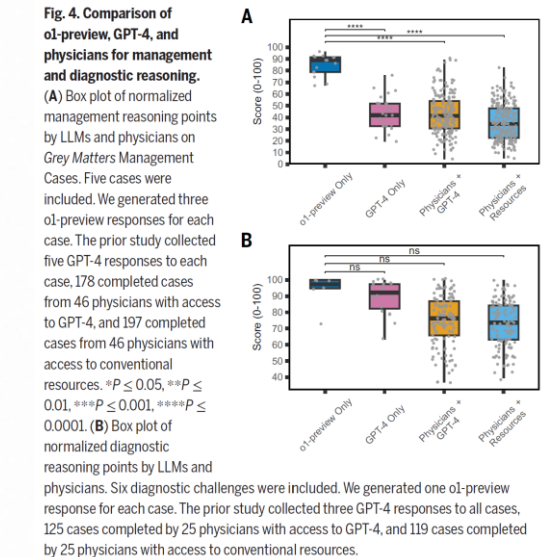


Fig. 4. Comparison of *o1*-preview, GPT-4, and physicians for management and diagnostic reasoning. (A) Box plot of normalized management reasoning points by LLMs and physicians on *Grey Matters* Management Cases. Five cases were included. We generated three *o1*-preview responses for each case. The prior study collected five GPT-4 responses to each case, 178 completed cases from 46 physicians with access to GPT-4, and 197 completed cases from 46 physicians with access to conventional resources. * $P \leq 0.05$, ** $P \leq 0.01$, *** $P \leq 0.001$, **** $P \leq 0.0001$. (B) Box plot of normalized diagnostic reasoning points by LLMs and physicians. Six diagnostic challenges were included. We generated one *o1*-preview response for each case. The prior study collected three GPT-4 responses to all cases, 125 cases completed by 25 physicians with access to GPT-4, and 119 cases completed by 25 physicians with access to conventional resources.



Pengenalan Chat GPT

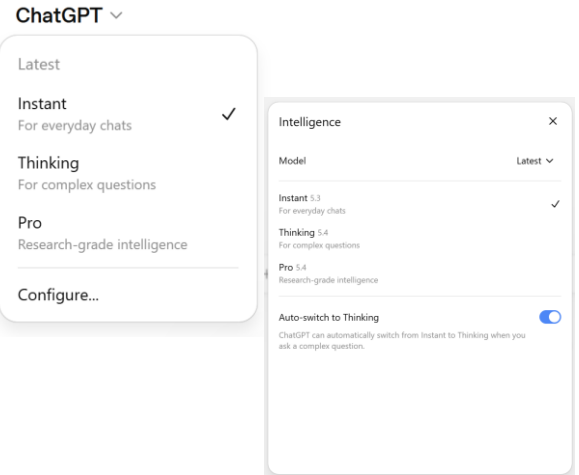
Model/Algoritma:

Jenis Akun dan Keamanan

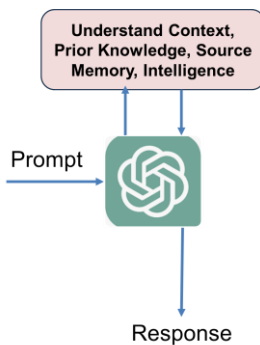
- **Data Tidak Aman: (Free, Plus)**
- Data Aman: Business, Enterprise, Pro

GPT Ecosystem

- Canvas: Kolaborasi dengan AI
- Project: management
- Membagi Chat (Sesama member GPT Teams)
- Custom-GPT



Cara kerja: Internal vs external knowledge (Context)



Attention Is All You Need

<p>Ashish Vaswani* Google Brain avaswani@google.com</p>	<p>Noam Shazeer* Google Brain noam@google.com</p>	<p>Niki Parmar* Google Research nikip@google.com</p>	<p>Jakob Uszkoreit* Google Research uszko@google.com</p>
<p>Llion Jones* Google Research llion@google.com</p>	<p>Aidan N. Gomez*¹ University of Toronto aidan@cs.toronto.edu</p>	<p>Lukas Kaiser* Google Brain lukasz.kaiser@google.com</p>	
<p>Illia Polosukhin*¹ illia.polosukhin@gmail.com</p>			

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

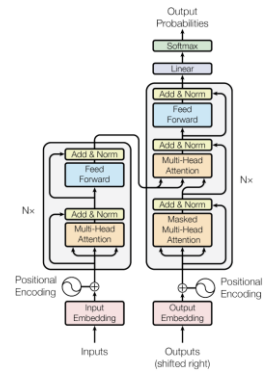


Figure 1: The Transformer - model architecture.

Pilih Model GPT yang Tepat

Gunakan model sesuai kompleksitas tugas, kedalaman analisis, dan kebutuhan waktu Anda.

Aspek	STANDARD	ADVANCED	DEEP RESEARCH
Model	Instant, Thinking-Standard	Thinking-Extended, Pro	Deep Research
Deskripsi	Tugas sehari-hari dan kebutuhan cepat	Analisis mendalam dan tugas kompleks	Riset komprehensif dan keputusan berbasis data
Kecepatan	Cepat	Sedang	Lebih lama
Kedalaman Analisis	Standar	Mendalam	Sangat mendalam
Akurasi	Baik	Tinggi	Sangat tinggi
Tugas Ideal	Aktivitas sehari-hari	Tugas kompleks & strategis	Riset & insight kritis
Penggunaan	Efisien	Profesional	Analisis tingkat lanjut

Keterbatasan Model Sederhana LLM (Model: Instant, Thinking-Standard)

- Tidak mendalam/spesifik
- Tanpa sumber/referensi
- Kesulitan pada konsep kompleks
- Sering keliru memilik fakta/fiksi
- Berhalusinasi



“Fully guided AI as writing assistance”

Strategic Prompting

Prompt Engineering

Prompting Framework

Memberi ruang bagi pengguna untuk berkreasi

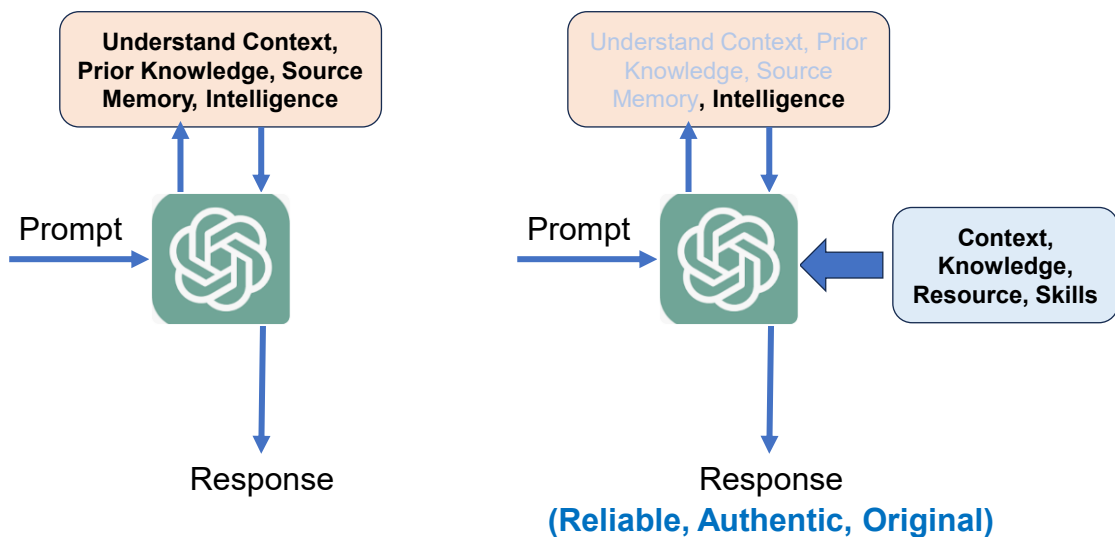
Keterbatasan Model Lanjut LLM (Pro, Thinking Extended)

- Tidak mendalam/spesifik
- Tanpa sumber/referensi
- Kesulitan pada konsep kompleks
- Sering keliru memilik fakta/fiksi
- Berhalusinasi
- Sedikit ruang untuk berkreasi/kreatifitas
- Kualitas Respond >>> Kapasitas Prompter
(AI lebih dominan)

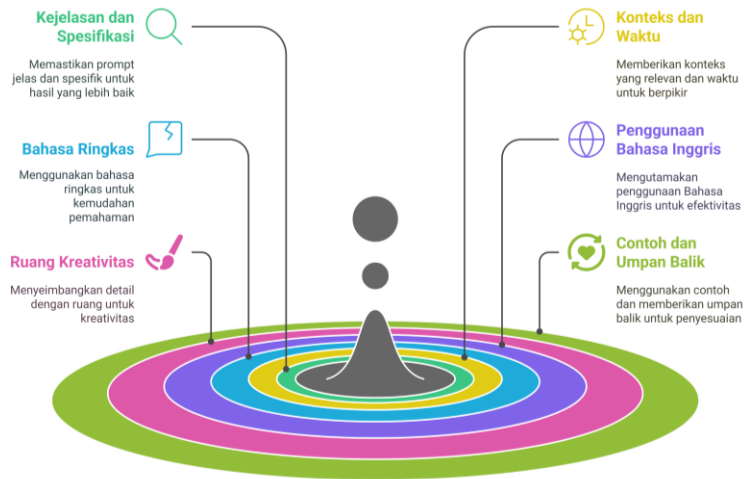


1-3 Prompting yang kompleks

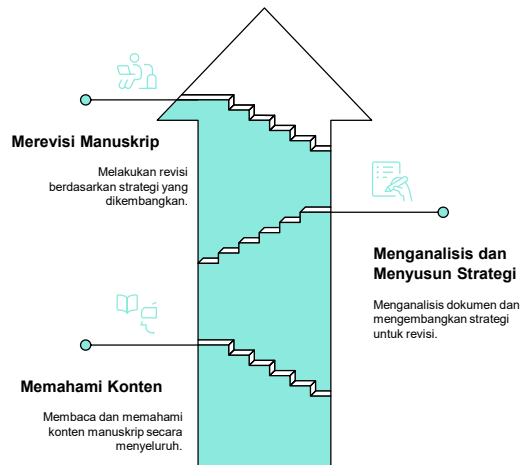
Teknik Prompting: Retrieval-Augmented Generation



Teknik Prompting: Instruction-Based Prompting AI → *Organisme Digital*



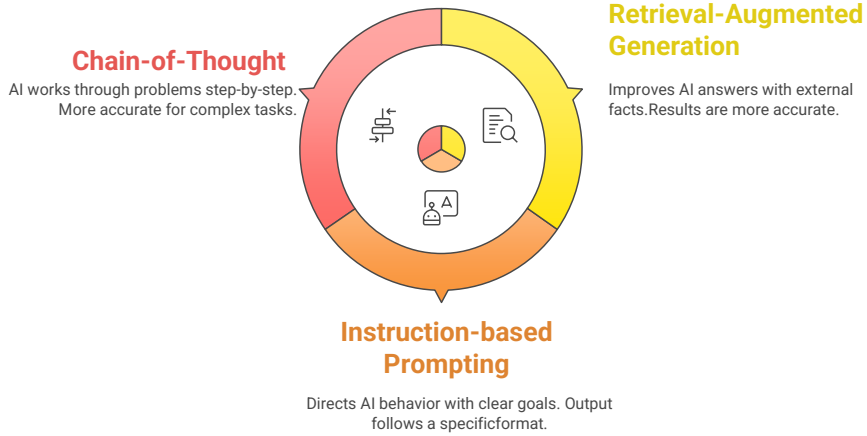
Teknik Prompting: *Chain-of-Thought*



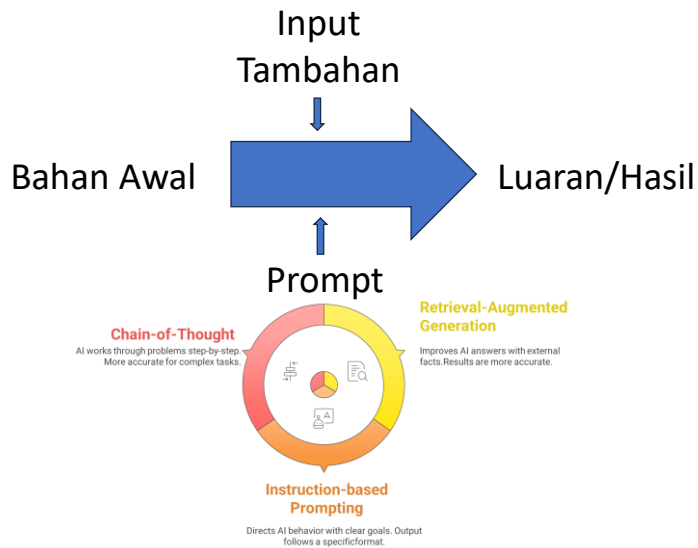
Teknik Prompting Lanjut

Pilihan Model

Tipe respon: Chat atau Canvas



Proses Prompting



3 Elemen Prompt

Instruksi: Tugas atau pertanyaan.

Input : Konteks, data, atau lampiran.

Model : Pilih model AI.



Rumus sederhana:

Instruksi jelas + input relevan + model tepat = hasil AI lebih optimal.

Pilihan Prompting

Pilih metode prompting sesuai kebutuhan dan tujuan kerja Anda.

Standard — Prompt Detail/Panjang

- **Panjang/Susah**: Prompting detail, harus **menguasai materi** dan **berpengalaman**
- **Liar**: Variasi respond tinggi, memerlukan revisi berkali-kali
- **Fleksibel**: Prompt bisa dirubah sesuai kemauan

Custom-GPT — Shortcut Prompt (Prompt detail di dalam sistem)

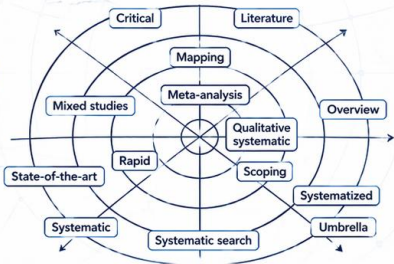
- **Singkat**: prompt pendek
- **Cepat**: hasil instan untuk kebutuhan khusus (i.e., Menulis Jurnal)
- **Terbatas**: perlu mendapatkan prompt dari developer
- **Rigid**: terbatas untuk fungsi tertentu

Dari Lautan Informasi Menuju Sintesis Pengetahuan

Tantangan utama penelitian: menavigasi beragam metodologi tinjauan.

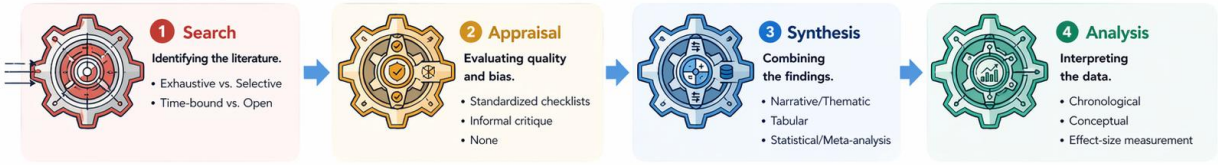
14 JENIS TINJAUAN

Peneliti harus menavigasi 14 jenis tinjauan yang berbeda.



“In 1753, James Lind conducted one of the first systematic evaluations of literature on scurvy, noting that before the subject could be understood, it was necessary to remove a great deal of rubbish.”

MESIN UNIVERSAL SINTESIS LITERATUR



Setiap jenis tinjauan adalah konfigurasi unik dari input dan batasan SALSAs.

10 Tahap Sistematis, Insight Bukan Kebetulan.

Satu langkah terlewat, seluruh hasil bisa bias.



Key Insight:
Systematic review bekerja seperti rantai. Kuat hanya sekuat mata rantai terlemahnya.

⚠ Tidak ada shortcut dalam ilmu. Ikuti semua tahap, hasil akan berbicara.

Perbedaan Empat Jenis Tinjauan

Narrative Review • Scoping Review • Systematic Review • Meta-Analysis

	1. NARRATIVE REVIEW	2. SCOPING REVIEW	3. SYSTEMATIC REVIEW	4. META-ANALYSIS
TUJUAN	Pemahaman luas topik melalui sintesis naratif.	Memetakan bukti: topik, konsep, jenis bukti, celah.	Menjawab pertanyaan spesifik dengan sintesis sistematis.	Menggabungkan hasil numerik dari studi serupa.
CAKUPAN	Luas, tanpa batasan ketat metode/desain.	Luas, berbagai jenis sumber & desain.	Spesifik, fokus pada pertanyaan jelas.	Spesifik, studi kuantitatif dengan kriteria ketat.
METODE	Fleksibel, tidak selalu prosedur baku.	Scoping (mis. PRISMA-ScR) iteratif.	Protokol ketat (mis. PRISMA) sistematis & replikatif.	Teknik statistik untuk menggabungkan effect size.
HASIL	Sintesis naratif, interpretasi, perspektif konseptual.	Peta bukti: tema, konsep, jenis intervensi, populasi, celah.	Sintesis tematik/kualitatif atau naratif untuk menjawab pertanyaan.	Estimasi efek gabungan (effect size) (mis. CI, p-value).
KEGUNAAN	Pemahaman awal, identifikasi konsep & arah penelitian.	Identifikasi cakupan, tren, dan celah penelitian.	Dasar keputusan berbasis bukti (praktik, kebijakan, penelitian).	Menentukan besaran efek intervensi/variabel secara lebih akurat.
KARAKTERISTIK KUNCI	Subjektif, berfokus pada interpretasi penulis.	Eksploratif, inklusif, memetakan lanskap bukti.	Sistematis, transparan, dapat direplikasi.	Kuantitatif, berbasis statistik, homogenitas studi diperlukan.



Claim of Novelty

Research gap – Novelty
Hypothesis – Objectives/Research questions

Data: Proof – Evidence

Confirmation of the claim
Answers to the research questions

Title

Abstract, Keyword

Introduction

Materials and Methods

Results Discussion

Conclusion

References

Dalam SLR, Literatur → Data

Data set 1

Data set 2

Data set 3

Data set 4

Title

Abstract

Introduction

Methods
(Structured, Systematic)

Results and Discussion

Data set 1

Data set 2

Data set 3

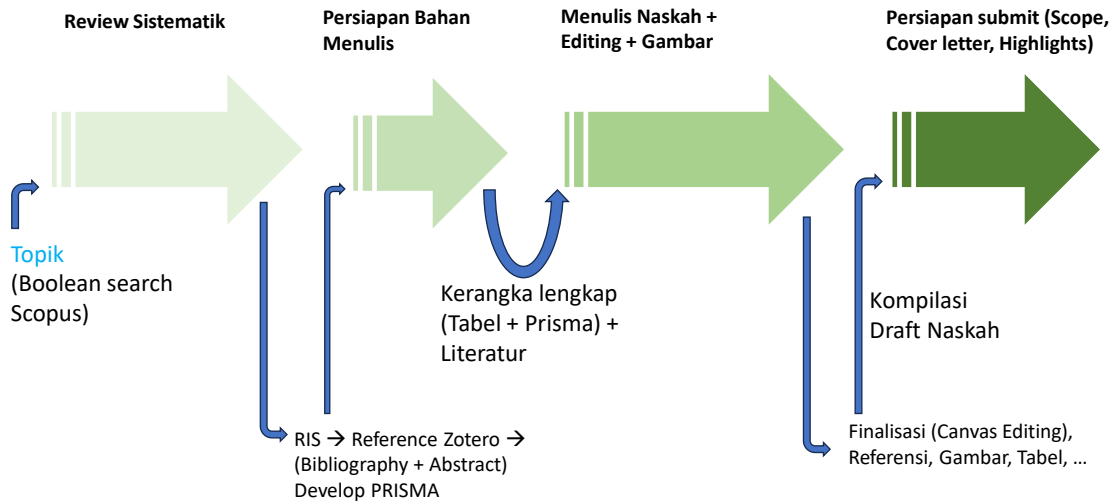
Data set 4

Conclusion

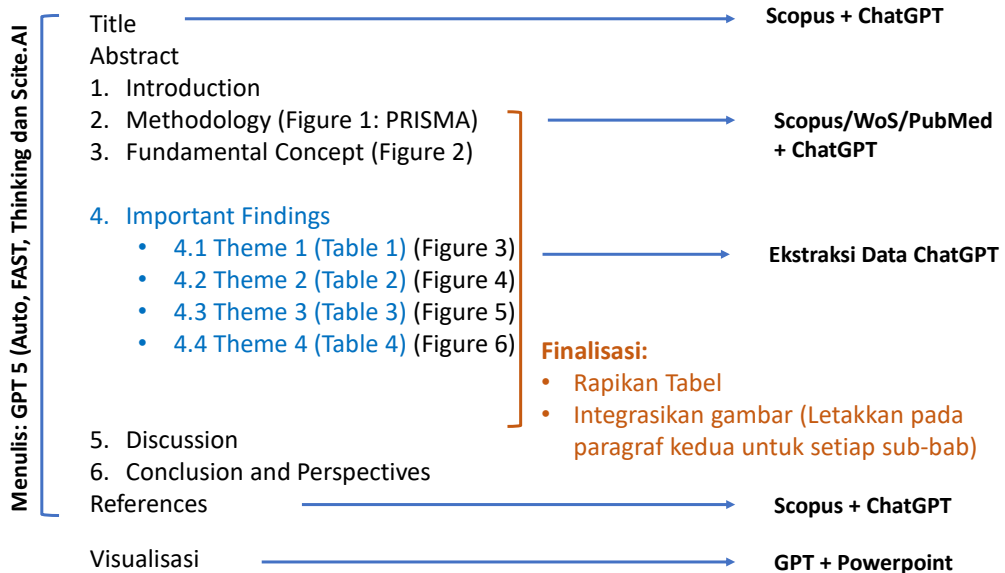
Progres Penelitian →

References

Workflow: SLR with GRA



Standar Kerangka SLR

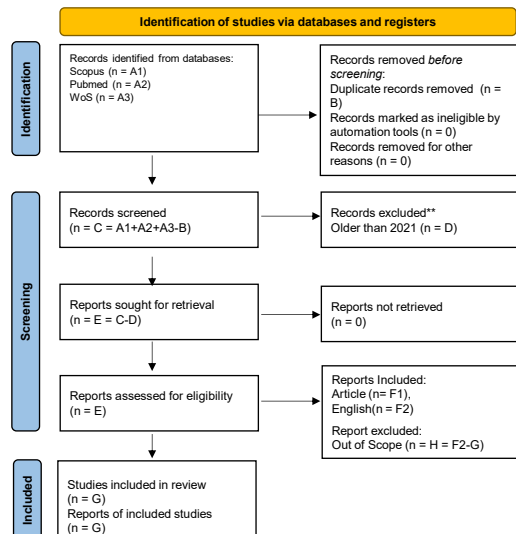


Prompting Framework SLR

Input/Model	Prompt	Output + action
Topik riset Instant	Prompt 0: Suggest Boolean Search for research topic! [Topic:]!	Keywords pencarian literatur Scopus (hanya satu database). Lakukan pencarian (Target 100-500), buat file References_Summary Mulai buat diagram Prisma!
References_Summary_SLR Thinking	Prompt 1: Suggest titles!	Tiga topik pilihan (pindahkan ke file MS Word , validasi, pilih sesuai minat)
Thinking	Prompt 2: Write the outline!	Canvas, kerangka keseluruhan (gabungkan ke file MS Word)
Thinking	Prompt 3: Extract data to construct Tables!	Canvas, Table 1-4, gabungkan ke file MS Word . Selesaikan Prisma, gabungkan ke file MS Word .
Thinking	Prompt 4: Generate references statements	Pernyataan literatur, gabungkan ke file MS Word . Boleh mengambil dari references summary saja (dengan prompt tambahan)
File MS Word Extended	Prompt 5: Write SLR manuscript	File Draft #1 SLR
Thinking	Prompt 6: Propose figures!	Deskripsi, judul dan penjelasan gambar, gabungkan ke file word
Instant	Prompt 7: Draw Figure XX !	Gambar (2-6) satu per satu, gabungkan ke file Draft #1 SLR
Instant	Prompt 8: Write captions and description!	Keterangan gambar XX (2-6) dan deskripsi, gabungkan ke file Draft #1 SLR

Catatan: Untuk menghitung jumlah referensi unik (included) di Table 1-4, gunakan prompt tambahan.

Diagram Prisma (Custom-GPT GRA)



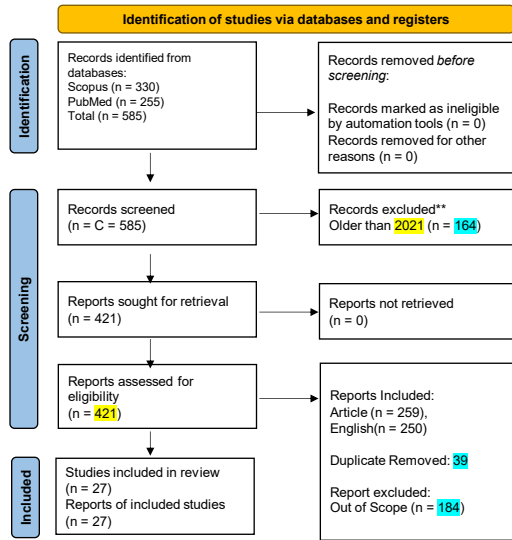
Kelebihan

- Cepat
- AI membantu dongkrak kualitas
- Cocok untuk bidang riset non-sensitive (control longgar)

Kekurangan

- Satu database (scopus)
- Exclusi dalam database suharusnya bukan bagian integral PRISMA
- Berbasis abstrak saja (tidak sesuai protocol)
- Tidak sesuai bidang riset
- Tidak terdaftar (i.e., Prospero)
- Tidak ada analisis Risiko bias (*Risk of Bias*)
- Tidak ada proses internal review (>1)

Diagram Prisma (Custom-GPT GRA) – Scopus + PubMed



(("cancer patient*" OR "oncology patient*" OR "patients with cancer" OR neoplasm*) AND ("malnutrition" OR "sarcopenia" OR "cancer cachexia") AND ("prediction model*" OR "risk prediction model*" OR "clinical prediction rule*" OR "risk score*" OR nomogram*))

Integrasi: Rayyan + Custom-GPT GRA)

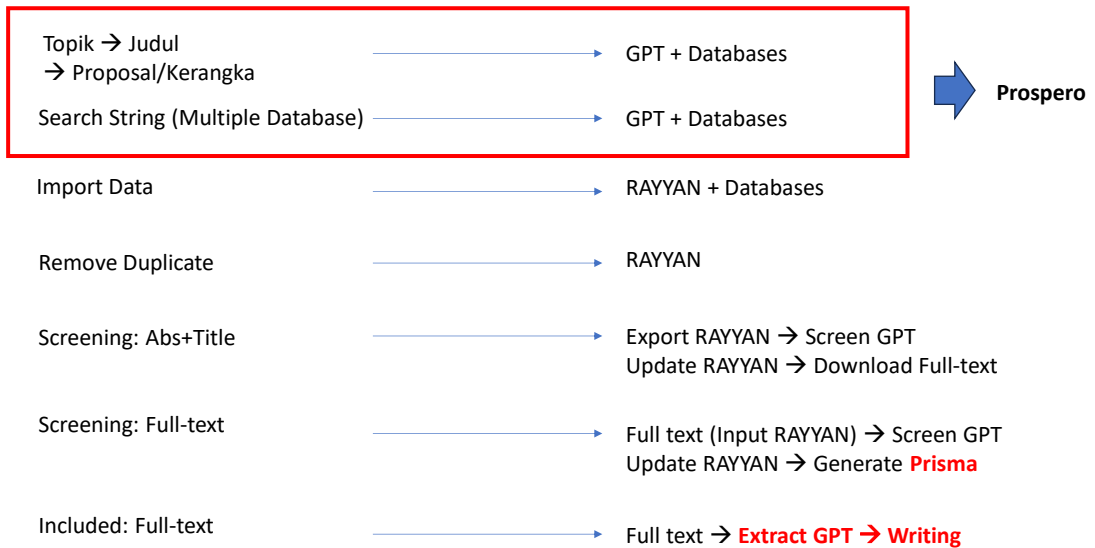
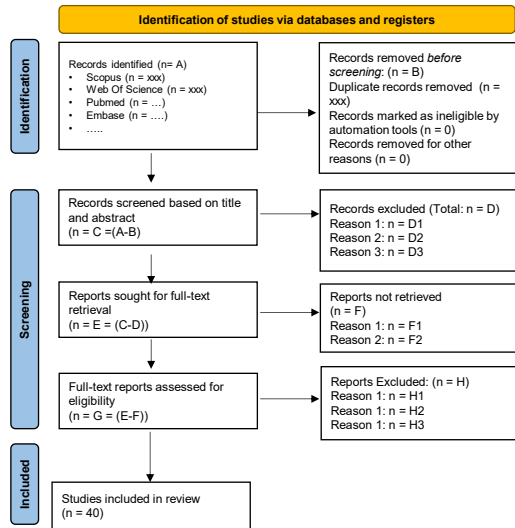


Diagram Prisma (RAYYAN)



Workflow RAYYAN

- Import Data (Multi-databases)
- Remove Duplicate (Otomatis)
- Screening Abs+Title: (Manual, >1 Peneliti)
- Screening Full-text : (Download → Upload → Pilih)
- Included: Full-text

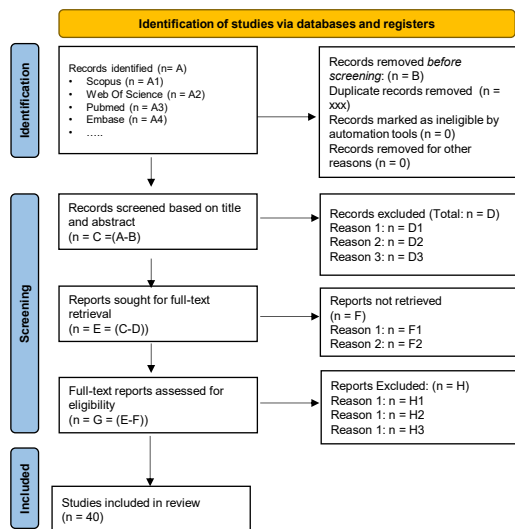
Kelebihan

- Proses Transparan
- Prudent (>1 peneliti)
- Record keeping
- Bagian metodologi → Supplementary

Kelemahan

- (sangat) Lama
- Hilang momentum dan focus
- Hasil terbatas keahlian peneliti

Integrasi: Rayyan + Custom-GPT GRA



Identifikasi study dengan kata kunci (**Prompt 0**)

- Database 1 (n=A1)
- Database 2 (n=A2),
- Database 3 (n=A3)
-
- (Total, n=A = (A1 + A2 + A3+))

Record screened: n = C = (A-B)

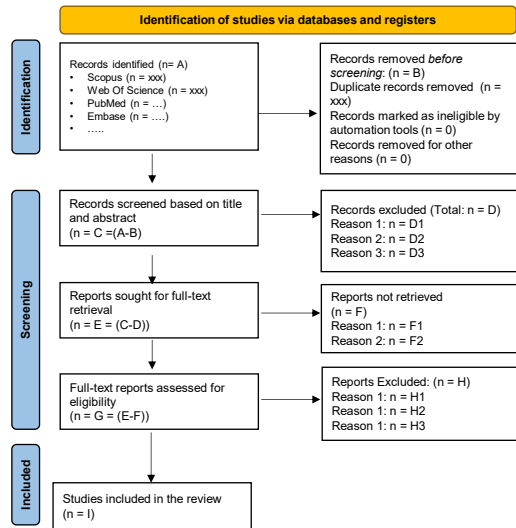
Import semua ke RAYYAN, deduplikasi (n = B)

Export Bibtext → Import Zotero
→ Buat References Summary

Prompt 1 → Pilih topik sesuai minat dan keahlian

Prompt 2 → Buat kerangka awal → Ajukan ke Prospero

Integrasi: Rayyan + Custom-GPT GRA



Record screened: $n = C = (A-B)$

Proses Screening → Export Bibtext → Import Zotero

→ Buat References Summary

Prompt 3 → (Tambahkan alasan eksklusi, kelompokkan jadi 3-4 kategori)

Record Excluded, Total: $n = D$

Reason 1: D1

Reason 2: D2

Reason 3: D3

Download Full text semua artikel terpilih included) (Total jumlah artikel = E)

Artikel yang tidak diperoleh (Total F), Alasan:

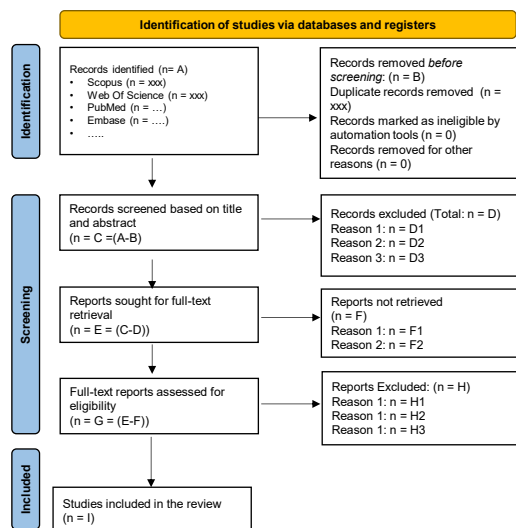
Reason 1: F1

Reason 2: F2

Reason 3: F3

Full-text reports assessed for eligibility, ($n = G = (E-F)$)

Integrasi: Rayyan + Custom-GPT GRA



[Lampirkan semua PDF n = G]

Prompt 3 (Ulang) → Refine Tables 1-4 based on full text. Remove studies out of scope (if any) found after assessing full-text! Add exclusion reasons (if any) grouped into 3-4 categories (if possible).

Reports Excluded: (n = H)

Reason 1: n = H1

Reason 1: n = H2

Reason 1: n = H3

Studies included in the review

(n = I) → Seluruh artikel full text pada Tabel 1-4.

Lanjutkan: **Prompt 4, Prompt 5, Prompt 6, Prompt 7, Prompt 8!**

Prompting Framework: Integrasi: Rayyan + Custom-GPT GRA

Model	Input	Prompt	Output + action
5.5 Instant	Topik riset	Prompt 0: Suggest Boolean Search for research topic [Topic:] in several suitable databases (provide the 3-5 suitable list!	Keywords pencarian literatur Scopus (gunakan multi-database). Lakukan pencarian (Target 100-500 untuk masing-masing database), buat file References_Summary Mulai buat diagram Prisma! (Gunakan template RAYYAN) Lakukan deduplikasi di RAYYAN Data Prisma: A (AB, A2, A3,), B (B1, B2, B3, ...), C = A-B
5.5 Thinking	References_Summary_SLR	Prompt 1: Suggest titles!	Tiga topik pilihan (pindahkan ke file MS Word , validasi, pilih sesuai minat)
5.5 Thinking		Title: Prompt 2: Write the outline! Prepare in an extensive form (3000-5000 words) suitable for Prospero reporting! Respond in Chat!	Canvas, kerangka keseluruhan (gabungkan ke file MS Word)
5.5 Thinking	Full teks semua artikel inklusi	Prompt 3: Extract data to construct Tables! Aim for 40-60 most relevant articles! (Tambahkan alasan eksklusi, kelompokkan alasan menjadi jadi 3-4 kategori)	Table 1-4 (Sementara). Download full-teks semua study inklusi Update Data Prisma: D (Alasan eksklusi: D1, D2, D3,), Data n = E = C-D Data F untuk yang tidak berhasil di download (F1, F2, F3,...) Data G = E-F

Prompting Framework: Integrasi: Rayyan + Custom-GPT GRA

Model	Input	Prompt	Output + action
	Semua full-text	Prompt 3: Refine Tables 1-4 based on full text. Remove studies out of scope (if any) found after assessing full-text! Add exclusion reasons (if any) grouped into 3-4 categories (if possible).	Table 1-4 (Final), gabungkan ke file MS Word . Update Data Prisma: F (Alasan eksklusi: H1, H2, H3,) Data I = G-F
5.5 Thinking		Prompt 4: Generate references statements	Pernyataan literatur, gabungkan ke file MS Word . Boleh mengambil dari references summary saja (dengan prompt tambahan) Selesaikan Prisma, gabungkan ke file MS Word .
5.5 Thinking Extended	File MS Word	Prompt 5: Write SLR manuscript	File Draft #1 SLR
5.5 Thinking		Prompt 6: Propose figures!	Deskripsi, judul dan penjelasan gambar, gabungkan ke file word
5.5 Thinking		Prompt 7: Draw all Figures!	Gambar, gabungkan ke file Draft #1 SLR
5.5 Thinking		Prompt 8: Write captions and description!	Keterangan gambar dan deskripsi, gabungkan ke file Draft #1 SLR