Evidence creation, gathering, & synthesis with the tools of AI: potential & considerations

Chirag J Patel
NASEM: Exploring the types of evidence synthesis and communications in diet
and chronic disease relationships
July 10, 2025



chirag@hms.harvard.edu @chiragjp www.chiragjpgroup.org

1

Disclosures

- Research funding (had no role in the production of research or this presentation)
 - NIH (NIDDK, NIEHS, NIAID, NIA)
 - · ARPA-H
 - NSF
 - Vranos Foundation
 - Steven and Alexandra Cohen Foundation
- Compute credits: Google, Microsoft, Amazon, Oracle
- · Consulting: Sanofi, Janssens

Emerging types of AI and high level use cases in biomedical research and applications

Supervised Practice of medicine:

Faster, cheaper, more accurate diagnoses/ decisions?

Unsupervised **Precision medicine:**

Discovering new actionable disease states

Generative **Creating medicine:**

Proposing new drugs Chatting with a virtual doctor/expert

Reinforcement "Problem solving":

Thinking through problems to "reason"

JAMA | Original Investigation | INNOVATIONS IN HEALTH CARE DELIVERY

Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs

JAMA 2016 varun Gulshan, PhD: Lily Peng, MD, PhD: Marc Coram, PhD. Martin C. Stumpe, PhD; Derek Wu, BS; Arunachalam Narayanaswamy, PhD: Subhashini Venugopalan, MS; Kasumi Widner, MS; Tom Madams, MEng, Jorge Cuadros, OD, PhD, Ramasamy Kim, OD, ONB; Rajik Param, MS, ONB, Philip C. Nelson, SS; Jessica L. Meg, MD, MPH; Dale R. Vesbeter, PhD

IMPORTANCE Deep learning is a family of computational methods that allow an algorithm to program itself by learning from a large set of examples that demonstrate the desired behavior, removing the need to specify rules explicitly. Application of these methods to medical imaging requires further assessment and validation

 $\begin{tabular}{ll} \textbf{OBJECTIVE} & To apply deep learning to create an algorithm for automated detection of diabetic retinopathy and diabetic macular edema in retinal fundus photographs. \end{tabular}$

DESIGN AND SETTING A specific type of neural network optimized for image classification called a deep convolutional neural network was trained using a retrospective development called a deep convolutional neural network was trained using a retrospective development data set of 128 T5 retinal images, which were graded 3 to 7 times for diabetic retinopathy, diabetic macular edema, and image gradability by a panel of 54 US licensed ophthalmologists and ophthalmology senior residents between May and December 2015. The resultant algorithm was validated in January and February 2016 using 2 separate data sets, both graded by at least 7 US board-certified ophthalmologists with high intragrader consistency.

EXPOSURE Deep learning-trained algorithm.

IES AND MEASURES The sensitivity and specificity of the algorithm for detecting MAIN OUT COMES AND MASORES THE RESERVING AND SPECIALTLY OF THE algorithm for creecting referable diabetic retinopathy (RDR), defined as moderate and worse diabetic retinopathy, referable diabetic macular edema, or both, were generated based on the reference standard of the majority decision of the ophthalmologist panel. The algorithm was evaluated at 2 operating points selected from the development set, one selected for high specificity and another for high sensitivity.

RESULTS The EyePACS-1 data set consisted of 9963 images from 4997 patients (mean age, 54.4 years; 62.2% women; prevalence fDR, 683/8876 fully gradable images [7.8%); ht Messidor-2 data set had 1748 images from 874 patients (mean age, 57.6 years; 42.6% women; prevalence of RDR, Z54/1745 fully gradable images [14.6%). For detecting RDR, the algorithm dat an area under the receiver operating curve of 0.991 (9.5% cl. 0.988-0.993) for EyePACS-1 and 0.990 (95% cl. 0.986-0.993) for EyePACS-1 and 0.990 (95% cl. 0.9786-27%) and the specificity for EyePACS-1 the sensitivity was 98.5% (6.9% cl. 0.9786-27%) and the specificity was 98.5% (95% cl. 0.9786-997%) and the specificity was 98.5% (95% cl. 0.9786-997%) and the specificity was 98.5% (95% cl. 0.9786-997%). Using a second operating point with high sensitivity in the development set, for EyePACS-1 the sensitivity was 97.5% and specificity was 93.4% and for Messidor-2 the sensitivity was 96.1% and specificity was 93.9%. RESULTS The EyePACS-1 data set consisted of 9963 images from 4997 patients (mean age, 54.4

Foundation, Aravind Eye Care
System, Madurai, India (Kim); Shri

Question How does the performance of an automated deep learning algorithm compare with manual grading by ophthalmologists for identifying diabetic retinopathy in retinal fundus photographs?

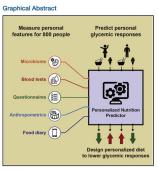
Finding In 2 validation sets of 9963 images and 1748 images, at the operating point selected for high specificity, the algorithm had 90.3% and 87.0% sensitivity and 98.1% and 98.5% specificity for detecting referable diabetic retinopathy, defined as moderate or worse diabetic retinopathy or referable macular edema by the majority decision of a panel of at least 7 US board-certified ophthalmologists. At the operating point selected for high sensitivity, the algorithm had 97.5% and 96.1% sensitivity and 93.4% and 93.9% specificity in the 2 validation sets.

Meaning Deep learning algorithms had high sensitivity and specificity for detecting diabetic retinopathy and macular edema in retinal fundus photographs.

Deep vision model (CNN) with high AUC Author Affiliations: Google Ind.
Mountain View. California (Galban,
Clear clinical bottleneck articulated
Peng, Ceam, Stumpe, Wil,
Narayanawamy, Venugopalan,
Viden, Madams, Meshon, Webster),
Department of Computer Science,
Unerusylor (Feas, Austin
(Venugopalan); EyerKCS Lt.C.
Comparison with human raters
Unequipolani; EyerKCS Lt.C.
Gadatae Grup, University of
Gada

Examples in the precision nutrition space: But need prospective data to learn of benefits

Cell **Personalized Nutrition by Prediction of Glycemic** Responses



Authors David Zeevi, Tal Korem, Niv Zmora, ..., Zamir Halpern, Eran Elinav, Eran Segal

eran.elinav@weizmann.ac.il (E.E.), eran.segal@weizmann.ac.il (E.S.)

People eating identical meals present high variability in post-meal blood glucose response. Personalized diets created with the help of an accurate created with the help of an accurate predictor of blood glucose response that integrates parameters such as dietary habits, physical activity, and gut microbiota may successfully lower post-meal blood glucose and its long-term metabolic consequences.

Cell 2015

- High interpersonal variability in post-meal glucose observed in an 800-person cohort
- Using personal and microbiome features enables accura
- Prediction is accurate and superior to common practi dependent cohort
- Short-term personalized dietary interventions success lower post-meal glucose

Novel subgroups of adult-onset diabetes and their association with outcomes: a data-driven cluster analysis of six variables

Emma Ahlqvist, Petter Storm, Annemari Käräjämäki*, Mats Martinell*, Mozhqan Dorkhan, Annelie Carlsson, Petter Vikman, Rashmi B Prasad, Dina Mansour Aly, Peter Almqren, Ylva Wessman, Nael Shaat, Peter Spéqel, Hindrik Mulder, Eero Lindholm, Olle Melander, Ola Hansson, Ulf Malmqvist, Åke Lernmark, Kaj Lahti, Tom Forsén, Tiinamaija Tuomi, Anders H Rosengren, Leif Groop

Lancet Diabetes and Endocrinology 2019

Are there more than 3 major types of diabetes?

... and do they matter clinically?

Unsupervised learning on diabetics: f(HbA1C, BMI, HOMA-IR, HOMA-B, age of diagnosis, glutamate carboxylate antibodies) highlight new phenotypes with potentially new risk factors

The Large Language Model ("ChatGPT") revolution: Vast feature set and getting better

- Summarization of text
- · Creation of text, images, and videos
- Programming and data processing
- "Thinking": solving problems and searching for abstract ideas
- Hallucination and references to false literature
- · Proliferation of false positives?

7

FDA NEWS RELEASE

FDA Launches Agency-Wide AI Tool to Optimize Performance for the American People

For Immediate Release: June 02, 2025



The U.S. Food and Drug Administration (FDA) today launched Elsa, a generative Artificial Intelligence (AI) tool designed to help employees—from scientific reviewers to investigators—work more efficiently. This innovative tool modernizes agency functions and leverages AI capabilities to better serve the American people.

"Following a very successful pilot program with FDA's scientific reviewers, I set an aggressive timeline to scale AI agency-wide by June 30," asid FDA Commissioner Marty Makary, M.D., M.P.H. "Today's rollout of Elsa is ahead of schedule and under budget, thanks to the collaboration of our in-house experts across the centers."

Built within a high-security GovCloud environment, Elsa offers a secure platform for FDA employees to access internal documents while ensuring all information remains within the agency. The models do not train on data submitted by regulated industry, safeguarding the sensitive research and data handled by FDA staff.

"Today marks the dawn of the AI era at the FDA with the release of Elsa, AI is no longer a distant promise but a dynamic force enhancing and optimizing the performance and potential of every employee," said FDA Chief AI Officer Jeremy Walsh. "As we learn how employees are using the tool, our development team will be able to add capabilities and grow with the needs of employees and the agency."

The agency is already using Elsa to accelerate clinical protocol reviews, shorten the time needed for scientific evaluations, and identify high-priority inspection targets.

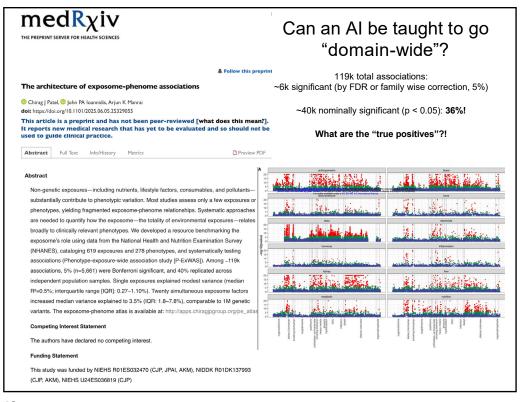
Can Large Language Models enhance and/or reduce the bottlenecks of the evidence synthesis process? Specific Questions:

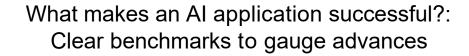
- What are the "causal machinery" behind LLMs? [unknown]
 - · What data were they trained on?
 - · How can these be refined (e.g., nutritional or interventional studies and papers)?
- · How good are they at information retrieval?: millions of papers to traverse
- · Can they distinguish study design artifacts from genuine biases? [to be evaluated]
 - · Can they reduce publication biases by uncovering file-drawer studies?
 - Do they know what is "high quality" vs. "low quality"? Randomized vs observational evidence?
 - · Can they infer across biological models (cell line vs. mouse vs human?)
- · Can they predict reproducibility? [to be evaluated]
 - .. or understand concepts such as triangulation or Bradford-Hill criteria?
 - · ... or inflate false positive findings? [yes]
- · Can they extract summary statistics? [to be evaluated; getting better]

9

Input annotated data (e.g., data dictionary) and produces a traceable paper NEJM Al Supply the goal (closed) or have the system create a hypothesis (autopilot or with some human intervention) Autonomous LLM-Driven Research - from Data to Literature search **Human-Verifiable Research Papers** Plan for testing hypothesis Table design Received: June 4, 2024; Revised: October 10, 2024; Accepted: October 17, 2024; Published: December 3, 2024 Literature search Human review JSIONS Our work demonstrates the potential for AI-driven acceleration of scien-overy in data-driven biomedical research and beyond, while enhancing, rather than Open Goal/"autopilot": use of public data (e.g., CDC BRFSS) 1 hour ~ hypothesis generation to code, run analyses, and generate the paper; run 5 times Fixed Goal: replication of existing paper with negative or positive findings; or findings with multiple steps (published after the GPT was trained) Provided the goal of the paper; ran 10 times Figure 3. Data-to-Paper Is Able to Autonomously Create Correct Papers for Simple Research Goals and Datasets While Human Copiloting Is Required to Ensure Accuracy in More Complex Settings.







Article

Highly accurate protein structure prediction with AlphaFold

https://doi.org/10.1038/s41586-021-038 Received: 11 May 2021 Accepted: 12 July 2021 Published online: 15 July 2021 Open access John Aumger⁴⁴⁴, Efchard Coura⁵⁴, Alexander Prizzi⁵⁴, Tim Green⁵⁴, Michael Figurnov⁵⁴, Old Finoneberger⁵⁴, Edwiny Turnyssunvisor⁵⁴, Beast Baste⁵⁴, Josephin Zelder⁵, Anna Potgopnich⁵⁴, Alex Bridghandr⁵⁴, Clemens Myeyer⁵⁴, Simon A. A. Kohl⁵⁴, Anna Potgopnich⁵⁴, Alex Bridghandr⁵⁴, Clemens Myeyer⁵⁴, Simon A. A. Kohl⁵⁴, Andere A. Balladr⁵⁴, Andere Courage ⁵⁴, Beather Michael ⁵⁴, Simolah Jahi⁵⁴, Simolah Jahi⁵⁴, Andere Media ⁵⁴, Simolah Jahi⁵⁴, Andere Media ⁵⁴, Simolah Jahi⁵⁴, Andere Myer ⁵⁴, Simolah Jahi⁵⁴, Andere Myer ⁵⁴, Simolah Jahi⁵⁴, Andere ⁵⁴, Simolah Jahi⁵⁴, Si

Proteins are essential to tille, and understanding their structure can facilitate an enchanticul understanding of their intection. Through encommons oper-primered better interestanticul understanding of their intection. Through encommons oper-primered better processes a small fraction of the billions of known protein sequences." Structural of their processes as the contractive of the billions of known protein sequences." Structural observations of their processes as the contractive of the contractive or plantasting effort excepted to determine a single procise instructure. Accurate computational approaches are needed external excepted and a single procise instructure. Accurate computational approaches are needed as sequence—the structure prediction component of the protein finding problems—in an appearance present problem from alloy though the observable only the same to calculate these an important open except problem from the soft provide instructure is proposed. "In a seminar observable and the seminar processes are necessary to be a seminar processes and the seminary opening problems are necessary to be a seminary opening to the contractive processes are necessary to the seminary opening a the Critical Assertion of our necessary to the contractive processes are necessary to the contractive processes and processes are necessary to processes approaches and processes and processes are necessary to processes are necessary to the contractive processes are necessar

The development of computational methods to predict the proceeded sing component problems the case of the problems of the language of the problems of the problems of the problems of the physical interactions on the evolutionary bilatory. The physical interaction of the physical interactions of the evolutionary bilatory in the physical interactions of the evolutionary bilatory. Although these time physical interactions thereof. Although these time of the time physical interaction is the physical interaction of the physical interaction of the physical interaction of the original interaction of the physical interaction of the original interaction of the physical physical interaction of the physical interaction of the physical physical interaction of the physical interaction of the physical physical interaction of the physical physical physical interaction of the physical interaction of the physical physical physical physical physical interaction of the physical phys

the steady growth of experimental protein structures deposit the Protein Data Bank (PUB)², the explosion of genomic seeper and the rapid development of deep learning techniques to inte and and experimental designation of the protein structures of the structure of the protein structures of the protein structures of the and evolutionary binkoy-based approved, producing real cases in variance and evolutionary binkoy-based approved superimentally and this intellecture of the protein structures of the protein structures of the bink study, we deseight of first, to our boundedge, computat the study where deserting the structures of the study of the securacy in a majority of cases. The neural network Alphair Gold the developed was entered into the CASP4 assessment (May 1–196).

Mind, London, LK. Yaboud of Biological Sciences, Seou Associal University, Seou, Supuk-News, "Authfald Intelligence Institute, Seoul National University, Seoul, South Korws, "Theology of Seoul National University, Seoul, South Korws, "Theology of Seouling, Seouling,

Challenge:

3D structure prediction from sequence is difficult and costly

Enabled by:

Critical Assessment of Protein Structure prediction

Database of "gold standard" labeled data and benchmarks for algorithms to exceed Experiments online since 1994-current day

https://predictioncenter.org/

What are the benchmarks for **evidence synthesis**?

13

Article Open access Published: 12 July 2023

Large language models encode clinical knowledge

Karan Singhal El, Shekoofeh Azizi El, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Alay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Abubakr Babiker, Nathaneal Schärli, Aakanksha Chowdhery, Phillip Mansfield, Dina Demmer-Fushman, Blaise Agüera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou ... Yukek Natarajan El + Show suthors

Nature 620, 172-180 (2023) | Cite this article

340k Accesses | 1763 Citations | 1222 Altmetric | Metrics



This article has been <u>updated</u>

Abstract

 $Large\ language\ models\ (LLMs)\ have\ demonstrated\ impressive\ capabilities,\ but\ the\ bar\ for the large language\ models\ (LLMs)\ have\ demonstrated\ impressive\ capabilities,\ but\ the\ bar\ for\ large\ l$ clinical applications is high. Attempts to assess the clinical knowledge of models typically rely on automated evaluations based on limited benchmarks. Here, to address these limitations, $we present \, Multi Med QA, a \, benchmark \, combining \, six \, existing \, medical \, question \, answering \, description \, descrip$ of medical questions searched online. HealthSearchOA, We propose a human evaluation framework for model answers along multiple axes including factuality, $comprehension, reasoning, possible \ harm \ and \ bias. \ In \ addition, we \ evaluate \ Pathways$ $L\underline{anguage\,Model^{I}\,(PaLM,a\,540\text{-}billion\,parameter\,LLM)\,and\,its\,instruction\text{-}tuned}\,variant, Flandard and the contraction of the contractio$ PaLM² on MultiMedQA. Using a combination of prompting strategies, Flan-PaLM achieves $state \hbox{-} of \hbox{-} the \hbox{-} art\ accuracy\ on\ every\ MultiMedQA\ multiple} \hbox{-} choice\ dataset\ (MedQA^3, A) \ and A) \ are the more properties of the more$ $MedMCQA^{\underline{4}}, PubMedQA^{\underline{5}} \ and \ Measuring \ Massive \ Multitask \ Language \ Understanding \ Massive \ Model \ Language \ Massive \ Model \ Massive \ Model \ Massive \ Model \ Massive \ Model \ Mo$ (MMLU) clinical topics6), including 67.6% accuracy on MedOA (US Medical Licensing Examstyle questions), surpassing the prior state of the art by more than 17%. However, human evaluation reveals key gaps. To resolve this, we introduce instruction prompt tuning, a parameter-efficient approach for aligning LLMs to new domains using a few exemplars. The resulting model, Med-PaLM, performs encouragingly, but remains inferior to clinicians. We show that comprehension, knowledge recall and reasoning improve with model scale and $instruction\ prompt\ tuning, suggesting\ the\ potential\ utility\ of\ LLMs\ in\ medicine.\ Our\ human\ prompt\ promp$ evaluations reveal limitations of today's models, reinforcing the importance of both evaluation frameworks and method development in creating safe, helpful LLMs for clinical

Benchmarking by medical QA: LLMs clearly capture medical textbook knowledge

But how do we assess how they do in practice with new data, emerging from RCTs and observational studies?

What is their false positive or hallucination rate?

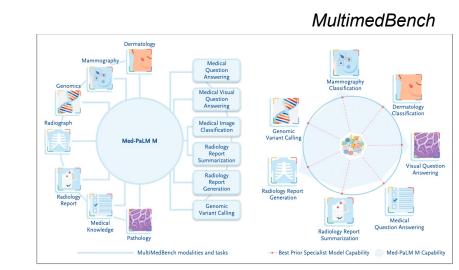


Figure 1. Med-PaLM M Overview.

A generalist biomedical artificial intelligence system should be able to handle a diverse range of biomedical data modalities and tasks, ideally using a single set of model weights to enable computationally efficient usage and elegant modeling of cross-modality interactions. To enable progress toward this overarching goal, we curated MultiMedBench, a benchmark spanning 14 diverse biomedical tasks, including question answering, visual question answering, image classification, radiology report generation and summarization, and genomic variant calling. Med-PaLM Multimodal (Med-PaLM M), our proof-of-concept for such a generalist biomedical artificial intelligence system (denoted by the shaded blue area), is competitive with or exceeds prior state-of-the-art results from specialist models (denoted by dotted red lines) on all tasks in MultiMedBench.

Pathways Learning Model (PaLM)

Sighal et al., Nature 2023 Tu et al., NEJM AI 2024

15

k Type	Modality	Dataset	Metric	SOTA†	PaLM-E (84B)	Med-PaLM M (84B)
Question answering	Text	Medical Question Answering (MedQA)	Accuracy	86.50% ¹⁶	28.83%	46.11%
		Medical Multiple-Choice Question Answering (MedMCQA)	Accuracy	72.30% ¹⁶	33.35%	47.60%
		PubMed Question Answering (PubMedQA)	Accuracy	81.80% ¹⁶	64.00%	71.40%
Report summarization	Radiology	Medical Information Mart for Intensive Care (MIMIC)-III	ROUGE-L	38.70% ¹⁷	3.30%	31.47%
			BLEU	16.20% ¹⁷	0.34%	15.36%
			F1-RadGraph	40.80%17	8.00%	33.96%
Visual question answering	Radiology	Visual Question Answering Radiology (VQA-RAD)	BLEU-1	71.03% ¹⁸	59.19%	69.38%
			F1	NA3	38.67%	59.90%
		Semantically-Labeled Knowledge- Enhanced Visual Question Answering (Slake-VQA)	BLEU-1	78.60% ¹⁹	52.65%	92.70%
			F1	78.10% ¹⁹	24.53%	89.28%
	Pathology	Pathology Visual Question Answering (Path-VQA)	BLEU-1	70.30% ¹⁹	54.92%	70.16%
			F1	58.40% ¹⁹	29.68%	59.51%
Report generation	Chest x-ray	MIMIC Chest X-ray (MIMIC-CXR)	Micro-F1-14	44.20% ²⁰	15.40%	53.56%
			Macro-F1-14	30.70% ²⁰	10.11%	39.83%
			Micro-F1-5	56.70% ²¹	5.51%	57.88%
			Macro-F1-5	NA3	4.85%	51.60%
			F1-RadGraph	24.40%14	11.66%	26.71%
			BLEU-1	39.48% ²⁰	19.86%	32.31%
			BLEU-4	13.30% ²¹	4.60%	11.31%
			ROUGE-L	29.60% ²²	16.53%	27.29%
			CIDEr-D	49.50% ²³	3.50%	26.17%
Image classification	Chest x-ray	MIMIC-CXR (5 conditions)	Macro-AUC	81.27% ²⁴	51.48%	78.35%
			Macro-F1	NA3	7.83%	36.83%
	Dermatology	PAD-UFES-20	Macro-AUC	NA2	63.37%	97.27%
			Macro-F1	NA2	1.38%	84.32%
	Mammography	VinDr-Mammo	Macro-AUC	64.50% ²⁸	51.49%	71.76%
			Macro-F1	NA3	16.06%	35.70%
		Curated Breast Imaging Subset of Digital Database for Screening Mammography (CBIS-DDSM) (mass)	Macro-AUC	NA	47.75%	73.09%
			Macro-F1	NA	7.77%	49.98%
		Curated Breast Imaging Subset of DDSMDigital Database for Screening Mammography (CBIS-DDSM) (calcification)	Macro-AUC	NA3	40.67%	82.22%
			Macro-F1	70.71% ²⁶	11.37%	63.81%
	Genomics (variant calling)	PrecisionFDA (Truth Challenge V2)	Indel-F1	99.40% ²⁷	53.01%	97.04%
			SNP-F1	99.70%27	52.84%	99.32%

Finding and assimilating evidence is resource intensive (1 SR ~ 62 weeks, Borah et al BMJ Open 2017), with implications on resources on future studies



Abstract Screening Full-text Screening Data Extraction Data extraction accuracy²

Institutional email Sign up → Read the preprint →

Can the paper screening and data extraction procedure be sped up or automated?

Elicit Covidence RevMan **DistillerSR**

17

medRxiv preprint doi: https://doi.org/10.1101/2025.06.13.25329541; this version posted June 19, 2025. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a CC-8F 4.0 International license.

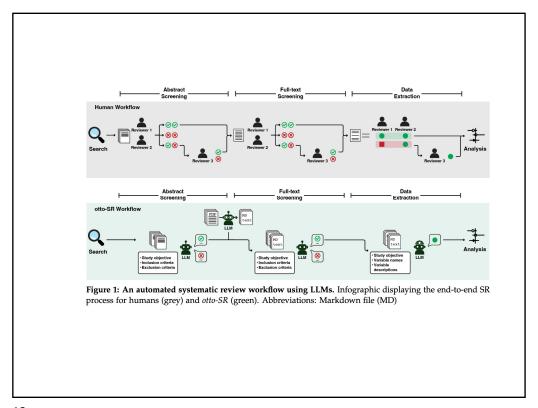
Automation of Systematic Reviews with Large Language Models

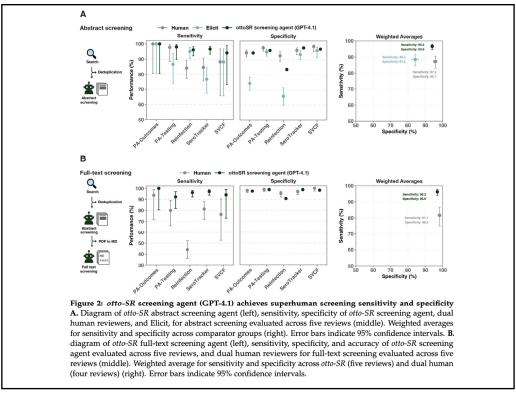
Christian Cao¹, Rohit Arora², Paul Cento³, Katherine Manta¹, Elina Farahani¹, Matthew Cecere¹, Anabel Christian (ao'; Rohit Arora'; Paul Cento'; Katherine Manta'; Elinia Farianni-, Matthew Cecere¹, Anabel Selemon¹, Jason Sang², Linig Ati Gong², Robert Kloosterman¹, Scott Jiang², Richard Seleh¹, Denis Margaliki¹, James Lin⁶, Jane Jomy¹, Jerry Xie², David Chen¹, Jaswanth Gorla¹, Sylvia Lee⁸, Kelvin Zhang², Harriet Ware², Mairead Whelan², Bijan Teja^{1,10}, Alexander A. Leung², Lina Ghosn^{11,12,13}, Rahul K. Arora², Allen S. Detsky¹, Michael Noetel¹⁴, David B. Emerson¹⁵, Isabelle Boutron^{11,12,13}, David Moher^{1,16,17}, George Church², Niklas Bobrovitz²

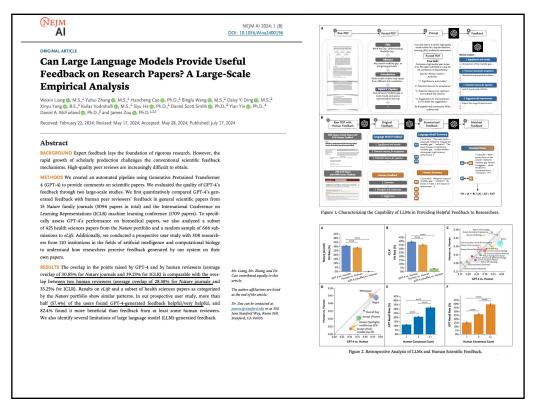
¹University of Toronto; ²Harvard Medical School; ³Independent Researcher; ⁴McGill University; ⁵University of British Columbia; ⁶Massachusetts Institute of Technology; ²University of Waterloo; ⁵Mount Sinai Hospital; ³University of Calgary; ³⁰St. Michael's Hospital; ³¹University for Sinai Cité; ¹³Université Sorbonne Paris Nord; ³²Cochrane France; ³²The University of Queensland; ³⁵Vector Institute; ³⁶Ottawa Hospital Research Institute; ³²University of Ottawa

Abstract

Systematic reviews (SRs) inform evidence-based decision making. Yet, they take over a year to complete, are prone to human error, and face challenges with reproducibility: limiting access to timely and reliable information. We developed otto-SR, an end-to-end agentic workflow using large language models (LLMs) to support and automate the SR workflow from initial search to analysis. We found that otto-SR outperformed traditional dual automate the Stw Owntow from initial search to analysis. We bound that arm 5-8 outperformed traditional automate human workflows in SB screening (atto-58: 96.7% sensitivity, 97.9% specificity) numan: 81.7% sensitivity, 98.1% specificity) and data extraction (atto-58: 93.1% accuracy; human: 79.7% accuracy). Using atto-58, we reproduced and updated an enlire issue of Cochrane reviews (n=12) in two days, representing approximately 12 work-years of traditional systematic review work. Across Cochrane reviews, atto-58 incorrectly excluded a median of 0 studies (IQR 0 to 0.25), and found a median of 2.0 (IQR 1 to 6.5) eligible studies likely missed by the original authors. Meta-analyses revealed that otto-SR generated newly statistically significant findings in 2 reviews and negated significance in 1 review. These findings demonstrate that LLMs can rapidly conduct and update systematic reviews with superhuman performance, laying the foundation for automated, scalable, and reliable evidence synthesis







21

Key observations: Al is getting better, but advances are resource and capital intensive

- Much of the reduction of practice by industry (e.g., access to infrastructure, energy, and human resources)
- Closed vs. open source LLMs: closed versions have dominated the market, but gap shrinking. What are the implications in evaluation?
- · "Human in the loop" approaches are proliferating
- · Automating procedures may lead to more false leads in the literature
- Success in AI deployment is related to availability of clear benchmarks (e.g., protein folding, disease screening)
- · How will we evaluate in research and practice?
- · How will consumers use and evaluate new technologies? (Not discussed here)