# Data extraction from full-text PDFs using Large Language Models for systematic reviews

Eitan Agai[1] , Alon Agai[1]

1 - PICO Portal, St. Petersburg, FL, USA

## Introduction

Rapid adoption of evidence-based decision-making in medicine, public health, and other fields has fueled a surge in evidence synthesis.

Traditional evidence synthesis methods are time-intensive and resource-demanding.

Advances in Artificial Intelligence (AI) and Large Language Models (LLMs) offer transformative potential to address these challenges by automating key tasks.

## Problem and Question

Issues with LLM hallucinations, lack of model transparency, and inconsistent outputs hinder reliability and scalability.

Can a multi-layered and AI-assisted approach to data extraction increase efficiency while maintaining accuracy during systematic review?

## Methods

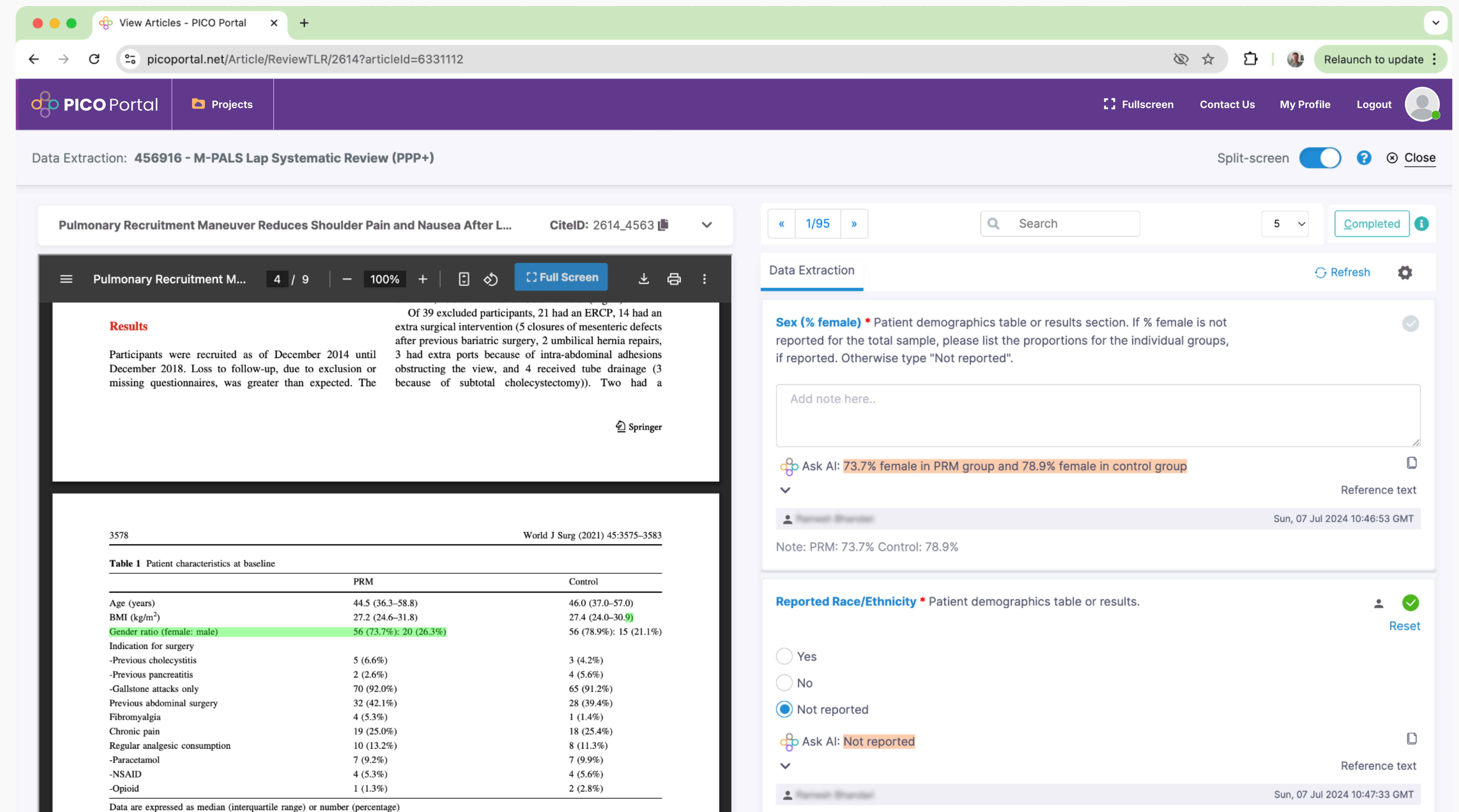We used an example systematic review with 193 included studies about digital health interventions.

Our workflow parsed each PDF to text, encoded the content, and applied task-specific prompt engineering with OpenAI GPT-4o. The model returned answers and highlighted supporting passages, allowing us to extract 40 data elements such as demographics/study characteristics, intervention details, and outcome data.

Human methodologists extracted the data assisted by the AI suggestions; we checked the accuracy of the AI by comparing the final human response to the AI recommended answer.

## Results

Overall, the AI reached 72.3% accuracy in its suggestions, with individual questions ranging from 52.9% to 99.3%. No custom prompt engineering was applied.

Productivity increased from an average of 80 extracted elements per hour for human-only extraction to 195 extracted elements per hour for AI-assisted extraction, an estimated time savings of 59%. This approach also appeared to reduce fatigue, allowing methodologists to work for an additional 1-2 hours without breaks.



**Figure**: AI-assisted data extraction (right) with interactive PDF (left).

## Conclusions

LLM technologies, combined with human oversight, have demonstrated the ability to reduce time, costs, and decision fatigue.

This approach represents a scalable solution to accelerate evidence syntheses across diverse fields including clinical research, public health policy making, and education.

REFERENCES [1] Khan et al. 2025 (DOI: 10.1101/2024.09.20.24314108); [2] Konet et al. 2023 (DOI: 10.1002/jrsm.1732); [3] Gartlehner et al. 2025 (DOI: 10.1101/2025.03.20.25324350); [4] Lieberum et al. 2025 (DOI:10.1016/j.jclinepi.2025.111746); [5] Gartlehner et al. 2024 (DOI: 10.1002/jrsm.1710); [6] Khraisha et al. 2024 (DOI: 10.1002/jrsm.1715); [7] Schmidt et al. 2024, ALTARS; [8] James et al. 2023 (DOI: 10.21203/rs.3.rs-3288515/v1); [9] Honghao et al. 2024 (DOI: 10.1001/jamanetworkopen.2024.12687); [10] Tyler et al. 2023 (DOI: 10.1101/2023.11.19.23298727); [11] Jacobsen et al. 2022 (DOI: 10.1186/s12874-022-01649-y); [12] Angelika et al. 2024 (DOI: 10.1101/2024.07.16.24310483).