

Assessing the Performance of Artificial Intelligence (AI)-augmented Electronic Health Record (EHR) Data Abstraction for Clinical Trial Patient Screening

Study Protocol

NCT06561217

Version date: January 2025

Outline

1. Abstract
2. Overall objectives
3. Aims
 - 3.1 Primary outcome
 - 3.2 Secondary outcomes
4. Background
5. Study design
 - 5.1 Design
 - 5.2 Study duration
 - 5.3 Target population
 - 5.4 Accrual
 - 5.5 Key inclusion criteria
 - 5.6 Key exclusion criteria
6. Subject recruitment
7. Subject compensation
8. Study procedures
 - 8.1 Consent
 - 8.2 Procedures
9. Analysis plan
10. Investigators
11. Human research protection
 - 11.1 Data confidentiality
 - 11.2 Subject confidentiality
 - 11.3 Subject privacy
 - 11.4 Data disclosure
 - 11.5 Data safety and monitoring
 - 11.6 Risk/benefit
 - 11.6.1 Potential study risks
 - 11.6.2 Potential study benefits
 - 11.6.3 Risk/benefit assessment

1. Abstract

The identification of eligible patients for clinical trials is a key component - and a rate-limiting step - of the clinical research enterprise. Currently, the process of identifying eligible patients relies on manual chart review by clinical research staff, which can be time-consuming, labor intensive and prone to human error. As a result, many eligible patients may be overlooked and opportunities for trial participation may be missed: this is particularly true in the oncology space, where fewer than 5% of adult cancer patients enroll in trials. The integration of AI technology into the clinical trial patient screening process has potential to improve trial participation rates and the generalizability of research findings. This study aims to assess the performance (accuracy, efficiency) of AI-augmented patient identification and inform optimal integration of AI technology into clinical research screening processes.

2. Overall objectives

The objective of this study is to assess and compare the accuracy and efficiency of three different approaches to abstracting clinical data used to identify oncology patients who meet the inclusion criteria for participation in clinical trials. The three approaches under evaluation include: (1) an autonomous AI algorithm (Mendel AI; developed by artificial intelligence startup company Mendel) which analyzes patient medical records to extract relevant clinical facts (“AI-alone”); (2) a human researcher who manually reviews patient charts as per the current norm/practice (“Human-alone”); and (3) a human researcher utilizing AI augmentation (“Human+AI”), where Mendel AI serves as a supportive tool in the decision-making process by providing the researcher a list of elements abstracted by the AI algorithm and a rank-order list of patients most likely to meet inclusion criteria for a trial.

The study primarily aims to compare (1) the chart-level accuracy of the Human+AI collaboration relative to Human-alone given the relevance of this comparison for real-world clinical workflows, defined by the percentage of pre-identified chart elements classified correctly compared against a predetermined “gold standard”; and (2) the efficiency of the Human+AI vs. Human-alone arms, defined by the time per chart review in minutes, measured for each chart.

Our hypotheses are (1) the Human+AI arm will be non-inferior in accuracy when compared to the Human-alone arm, in relation to a predetermined “gold standard”, and (2) that a Human+AI arm will be superior in speed of abstraction when compared to Human-alone screening.

The identification of eligible patients for clinical trials is a critical component of clinical research, as it directly impacts patient recruitment, study enrollment, and the generalizability of research findings. Currently, the process of identifying eligible patients often relies on manual chart review by clinical research staff, which can be time-consuming, labor-intensive, and prone

to human error. Consequently, eligible patients may be overlooked, and opportunities for trial participation may be missed. The integration of AI technology into the patient identification process has the potential to enhance the accuracy and efficiency of this critical task, leading to improved clinical trial recruitment and outcomes.

This study holds important implications for the field of clinical research by evaluating the effectiveness of AI-augmented patient identification compared to traditional manual methods and autonomous AI algorithms. By examining the strengths and limitations of each approach, the study will provide valuable insights into the optimal integration of AI technology in clinical research processes. Furthermore, the results of this study have the potential to benefit patients by improving their access to clinical trials and increasing awareness of available treatment options. For clinical research institutions, enhancing the efficiency of patient identification can lead to more effective use of research resources and the potential for accelerated clinical trial timelines. Ultimately, the findings of this study may contribute to advancements in clinical research practices, promoting more equitable access to trials and facilitating the development of innovative treatments for patients with cancer.

3. Aims

3.1 Primary

Aim 1: Compare the **accuracy** of the Human+AI vs. Human-alone arms, in relation to a “Gold Standard” determined by multiple human abstractors with no time limitation.

3.2 Secondary

Aim 2: Compare the **efficiency** of the Human+AI vs. Human-alone arms

4. Background

Insufficient patient enrollment in clinical trials is a significant challenge faced by the clinical trials community, and it is often considered the rate-limiting step in completing clinical trials. The traditional approach to determining patient eligibility for clinical trials is manual eligibility screening (ES), which requires a labor-intensive review of patient records, leading to a high workload for clinical staff and a potential delay in trial enrollment. The problem of insufficient patient enrollment is symbolized by the fact that less than 5% of patients with cancer participate in clinical trials. [1] The manual eligibility screening process is resource-intensive, and reducing the burden of screening is a priority for improving trial participation rates and the efficiency of clinical research.

Recent advancements in artificial intelligence (AI) and natural language processing (NLP) offer potential solutions to this challenge. These technologies have the capability to automatically detect patients' eligibility for clinical trials by mapping coded clinical trial eligibility criteria to corresponding clinical information extracted from electronic health records (EHRs). This process involves extracting relevant eligibility criteria from clinical notes using machine learning-based NLP applications and comparing these extracted criteria with reference criteria from trial descriptions. Studies have shown that machine learning-based NLP applications provide high accuracy in extracting eligibility criteria from clinical notes and determining trial eligibility, with recall rates of up to 90.9% and precision rates of up to 89.7%. [2]

Automated ES algorithms have demonstrated success in significantly reducing the workload of clinical staff involved in trial-patient matching, thereby increasing the trial screening efficiency of oncologists. Additionally, the use of AI-based tools has the potential to enable participation of smaller practices, which are often left out from trial enrollment due to resource constraints. The integration of AI and NLP technologies in the clinical trial eligibility screening process can play a crucial role in increasing total patient enrollment, the speed of enrollment, as well as expanding access to trials, and enhancing the execution of clinical research.

[1] Unger et al. Role of Clinical Trial Participation in Cancer Research: Barriers, Evidence, and Strategies. *American Society of Clinical Oncology*. 2017. DOI: [10.14694%2FEDBK_156686](https://doi.org/10.14694%2FEDBK_156686)

[2] Meystre et al. Automatic Trial Eligibility Surveillance Based on Unstructured Clinical Data. *International Journal of Medical Informatics*. 2019. DOI: [10.1016/j.ijmedinf.2019.05.018](https://doi.org/10.1016/j.ijmedinf.2019.05.018)

5. Study design

5.1 Design

In this experimental study, we will employ a three-arm design including (1) AI-alone, (2) Human-alone, (3) Human+AI to assess the accuracy of abstracting clinical elements from the patient chart. The accuracy of the abstraction will be compared against a pre-made “gold standard”. As detailed elsewhere, the study is powered for a direct comparison of the human alone vs human + AI arms as a non-inferiority analysis. For Arm 1, an autonomous AI algorithm (Mendel AI) will analyze patient medical records to extract relevant clinical elements (“AI-alone”). For Arm 2, a human researcher will manually review patient charts as per the current norm/practice (“Human-alone”). For Arm 3, a human researcher will utilize Mendel AI augmentation such that Mendel AI serves as a supportive tool in the decision-making process by providing the researcher a list of elements abstracted by the AI algorithm and a rank-order list of patients most likely to meet inclusion criteria for a trial (“Human+AI”). **For this study, we will**

only be using information from deidentified patient charts in Mendel's database. No Penn data will be utilized in this study.

We will employ 2 clinical research coordinators (CRCs) who will review and abstract specific clinical elements while reviewing the patient's electronic health record (EHR) which will include both "structured" (machine-readable) and "unstructured" (traditionally non-machine-readable; Mendel AI is uniquely capable of processing) data. To ensure no inter-rater differences in accuracy the CRCs will be assessed on a practice set of three charts and must meet an 80% threshold of agreement on review elements as well as concordance in chart-level accuracy as measured by Cohen's Kappa (threshold 0.6).

The components of the EHR that may be reviewed to abstract the data elements of interest may include:

- Physician Progress Notes
- Radiology Reports
- Pathology Reports
- Laboratory Reports
- Genetic/molecular sequencing and other genetic testing results
- Discharge Summaries
- Infusion Therapy reports

To help assess accuracy, the "Gold Standard" (to which all study arms will be compared) will be a manually annotated version based on review of at least 2 human abstractors (plus a tie-breaker if there is a discrepancy). This process will be led by the Mendel team, with independent oversight by the Penn research team to ensure consistency.

First, the Mendel AI algorithm will review all patient charts, producing a machine-annotated version of each chart with AI-abstracted elements (AI-alone). Next, each CRC will review all patient charts. Each CRC will review half of their charts manually (Human-alone) and half of their charts using the machine-annotation and rank-order list (Human+AI). The order of review will be randomized such that manual reviews and AI-assisted reviews will be mixed, and such that each CRC will review half of the charts manually and half of the charts with AI assistance.

Secondary aim metrics (Efficiency) will be assessed by tracking the volume of completed prescreens per unit of time.

5.2 Study duration

We anticipate the "review" phase of this study where human researchers are abstracting data from patient charts to take roughly 6 months, followed by 6 months for analysis.

5.3 Target population

This study will include de-identified patient charts from Mendel's databases from patients in community oncology practices with a diagnosis of non-small cell lung cancer (NSCLC) or colorectal cancer (CRC) with a minimum of 5 patient documents and data within 5 years from the time of data extraction. No Penn data will be utilized in this study and we will use standard EHR datasets.

5.4 Accrual

Mendel AI upholds a database of clinical data for deidentified patients from participating community-based oncology practices. This data will be utilized for this study. As noted above, no Penn data will be utilized in this study.

5.5 Key inclusion criteria

A patient chart will be included for analysis based on the following three screening criteria:

- A diagnosis of NSCLC or colorectal cancer
- A minimum of 5 patient documents in the Mendel database
- Most recent document was within 5 years from the time of data extraction

6. Subject recruitment

No subjects will be involved with this work.

7. Subject compensation

There will be no compensation for the use of subjects' de-identified data collected as part of routine clinical care.

8. Study procedures

8.1 Consent

The data utilized for analysis will be obtained from existing electronic health records (EHR) that have already been collected as part of routine clinical care and collated in the Mendel AI datasets. As such, patients will not be directly involved in the research process, and the study will not require any new data collection or intervention from patients.

Most importantly, all patient data used in the study has been de-identified and anonymized to ensure the protection of patient privacy and confidentiality.

In light of these considerations, patients will not need to provide consent for their data to be included in the study. This approach aligns with the ethical standards for retrospective studies using de-identified data, and the study design has been developed to ensure compliance with relevant regulations and guidelines on patient privacy and data protection.

As noted before, no Penn data will be utilized in this study.

8.2 Procedures

Our study will follow the design outlined above (section 5.1).

There is no informed consent needed for the use of this data. There will be no participant recruitment or compensation. Potential risks and benefits of the study are reported elsewhere (section 11.6).

Two CRCs will be hired to analyze de-identified patient data. CRCs will access the data via a password-protected portal, using a UPHS desktop and/or laptop. De-identified patient chart data will be included for analysis based on having a tumor type of interest. They will then review patient charts and record clinical elements of interest in the Mendel abstraction tool, which will be secured in a HIPAA-compliant external enclave. Summary-level data will be exported into datasheets for analysis. These datasheets will only be used and analyzed on UPHS equipment, or shared via HIPAA-compliant means.

Data will be analyzed using statistical software including Excel, R code, and STATA. For details of the analysis plan, see section 9.

Data elements that will be abstracted across all arms of the study include the following (the below list may change slightly pending further discussion among the research team):

Cancer Type and Staging

- Cancer Type
- Stage - Tumor
- Stage - Node
- Stage - Mets

Performance

- ECOG Performance Status

Prior Systemic Treatment

- Any Prior Platinum Therapy (e.g., Cisplatin, Carboplatin, Oxaliplatin)
- Any Other Prior Chemotherapy (e.g., 5-fluorouracil, Leucovorin, Capecitabine, Irinotecan, Pemetrexed, Paclitaxel, Docetaxel, Gemcitabine)
- Any Prior Immunotherapy (e.g., Atezolizumab, Nivolumab, Pembrolizumab)

Biomarkers for CRC

- KRAS
- MSI-H
- MSS
- BRAF
- HER2
- NTRK
- DPYD

Biomarkers for NSCLC

- EGFR
- ALK
- ROS1
- KRAS
- BRAF
- NTRK
- cMET
- RET
- PDL1

Others

- Outcomes (remission, residual disease, recurrence, loco-regional recurrence, distant recurrence)
- Responses (complete response, stable disease, partial response, mixed response, disease progression)

9. Analysis plan

The primary study will be powered on a retrospective paired, non-inferiority design to evaluate abstraction accuracy of elements commonly used to help screen patients for clinical trials. The two arms are (1) human reviewers versus (2) human reviewers utilizing predictions from Mendel AI. The primary outcome, the percentage of gold standard items correctly abstracted, will be measured for each chart. Noninferiority can be established by comparing the difference in the

mean proportion correct between the human-augmented and Human-alone arms with a predefined non-inferiority margin, Δ .

The one-sided hypotheses to test noninferiority of average chart-level accuracy of a Human+AI abstraction approach, compared with a Human-alone approach are outlined below:

$$H_0: \mu_{h-ai} - \mu_h \leq -\Delta \text{ (Null hypothesis)}$$

$$H_a: \mu_{h-ai} - \mu_h > -\Delta \text{ (Alternative hypothesis)}$$

Test statistic:

$$T_n = \frac{(\mu_{h-ai} - \mu_h) + \Delta}{\sqrt{\frac{Var(D)}{n}}}$$

Sample size equation for a noninferiority, paired study design:

$$n = \frac{(z_{1-\beta} + z_{1-\alpha})^2 \sigma^2}{((\mu_{h-ai} - \mu_h) + \Delta)^2} = \frac{(z_{1-\beta} + z_{1-\alpha})^2 \sigma^2}{(D + \Delta)^2}$$

As a secondary analysis, we will test the superiority of the Human+AI arm for achieving greater average chart-level accuracy, compared with the Human-alone arm. Our hypotheses are:

$$H_0: \mu_{h-ai} \leq \mu_h \text{ (Null hypothesis)}$$

$$H_a: \mu_{h-ai} > \mu_h \text{ (Alternative hypothesis)}$$

Sample size equation for a superiority, paired study design:

$$n = \frac{(z_{1-\beta} + z_{1-\alpha})^2 Var(D)}{(\mu_{h-ai} - \mu_h)^2}$$

<u>Parameter</u>	<u>Description</u>
μ_{h-ai}	Mean proportion of criteria correctly abstracted by AI-augmented human reviewers
μ_h	Mean proportion of criteria correctly abstracted by Human-alone reviewers
$D = \mu_{h-ai} - \mu_h$	Mean difference in chart-level accuracy between arms
Δ	Non-inferiority margin, set as 0.05 for conservative sample size estimates.
$z_{1-\beta}$	Critical value at $1-\beta$, where $1-\beta$ represents the desired power and β represents Type II error

$Z_{1-\alpha}$	Critical value at $1-\alpha$, where α = Type I error
σ^2	$n^* \text{Var}(D)$, where n is the sample size and $\text{Var}(D)$ is the variance of the mean difference in chart-level accuracy

*Note: μ_{h-ai} , μ_h , $\Delta > 0$.

From our initial calculations, 400 individual patients will achieve an estimated statistical power greater than 90% for this non-inferiority study with a Type I error of 5%. This was determined from simulations for different realistic parameter scenarios for TNR, event rate, and concordance assumptions showing that the needed sample size rarely exceeded 400.

The goals of the secondary analyses are to compare the efficiency of screening among the three study arms. We will analyze efficiency, defined by comparing the average number of minutes per chart review, measured for each chart. A nonparametric, paired Wilcoxon Rank Sum test with continuity correction will compare overall median efficiency of chart abstraction between arms.

The following are post-hoc analyses to be conducted after the pre-specified main analyses. Superiority of chart-level accuracy for the Human+AI arm relative to Human-alone will be tested using a one-sided paired t-test with an unspecified inferiority margin. A sensitivity analysis will repeat the primary chart-level accuracy outcome analysis with a less strict definition of an accurate match to the gold standard set. To assess whether there is a difference in accuracy among the 12 trial eligibility criteria, criterion-level accuracy will be compared assessed between arms. To assess criterion-level accuracy, a two-sided, paired t-test will be conducted with Bonferroni correction for multiple hypotheses (i.e., $\alpha = 0.05 / n$, where $n=12$ was the number of criteria studied). For any criterion with a statistically significant difference in mean criterion-level accuracy, superiority of the Human+AI arm for criterion-level abstraction will be then assessed with an additional one-sided hypothesis with an unspecified superiority margin. Lastly, we repeat both the chart-level accuracy analysis and efficiency analyses while stratifying by chart characteristics that include chart complexity, length, order, and cancer type.

Descriptive statistics, including means, medians, proportions, and ranges, will be used to summarize each analyses. All statistical analyses will be conducted using appropriate statistical software, and will rely on a predetermined significance level (e.g., $\alpha = 0.05$). Effect sizes and confidence intervals will be reported where applicable to provide further context to the findings.

10. Investigators

Ezekiel J. Emanuel, MD, PhD (Principal Investigator) is Diane v.S. Levy and Robert M. Levy University Professor in Health Policy and Medical Ethics and Co-Director of the Healthcare Transformation Institute.

Ravi B. Parikh, MD, MPP (Sub-Investigator) is Assistant Professor of Health Policy and Medicine and Innovation Faculty at the Penn Center for Cancer Care Innovation (PC3I).

Other members of the Penn research team include Professor Jinbo Chen, PhD (Collaborator), William Ferrell, MPH (Project Manager, Perelman MEHP Department), Matthew Guido (Project Manager, Perelman MEHP Department), Liz Beothy (Project Manager, Clinical Trials Unit) Ryan O'Keefe, MD, MBA (Resident Physician, HUP) and Likhitha Kolla (MD-PhD Student).

Key personnel from Mendel.ai - including Karim Galil, MD (CEO) and Sailu Challapalli (Chief Product Officer) - will be involved as external collaborators.

11. Human research protection

11.1 Data confidentiality

Computer-based files will only be made available to personnel involved in the study through the use of access privileges and passwords. All identifiers will be removed from study-related information. Precautions will be put in place to ensure the data are secure by using passwords and HIPAA-compliant encryption. All patient data provided by Mendel.ai will be viewed by HIPAA-certified Research Coordinators trained in EHR screening within a secure enclave.

NOTE: No Penn data will be utilized in this study.

11.4 Data disclosure

The results of the study will be presented in aggregate form, with no disclosure of individual-level data. No information that identifies specific patients will be recorded as part of documenting and/or publishing the results of this research.

11.5 Data safety and monitoring

The investigators will provide oversight for the study evaluation of this project.

11.6 Risk/benefit

11.6.1 Potential study risks

The potential risks associated with this study are minimal given the research is focused on de-identified data from patient charts. Breach of data is a potential risk that will be mitigated by using HIPAA compliant and secure data platforms for analysis.

NOTE: No Penn data will be utilized in this study.

11.6.2 Potential study benefits

This is highly practical research because it assesses the effectiveness of AI-led and human-AI-assisted chart review for clinical trial eligibility screening. This has major implications for clinical research, as it could lead to significant time savings in screening patients for eligibility and ultimately lead to more cancer patients participating in clinical trials.

11.6.3 Risk/benefit assessment

The risk/benefit ratio is favorable given the potential benefit for significant practical knowledge that could be gained from understanding the potential for AI to review “unstructured” patient data and accelerate the process of identifying patient eligibility for clinical trials. Efforts have been put into place to minimize the risk of breach of data. If favorable outcomes are found, then there is a potential to leverage the insights to guide further research and future clinical practice.