# Utilizing Al for Automated Data Entry and Analysis to Pre-screen Lung Cancer Clinical Trial Candidates

Mihaela Aldea<sup>1</sup>, Pierre Rolland<sup>2</sup>, Lodovica Zullo<sup>1</sup>, Solenne Simon<sup>1</sup>, Azeddine Djarallah<sup>2</sup>, Lisa Chuttoo<sup>2</sup>, Benjamin Vignal<sup>2</sup>, Jean-Charles Louis<sup>2</sup>, François Lion<sup>3</sup>, Arnaud Borie<sup>2</sup>, David

Planchard¹, Anas Gazzah⁴, Capucine Baldini⁴, Caroline Robert¹, Fabrice Andre¹, Fabrice Barlesi¹, Stefan Michiels⁵, Franck Le Ouay², Benjamin Besse¹

5 Biostatistics & Epidemiology, Gustave Roussy, Villejuif, France

1 Department of Medical Oncology, Gustave Roussy Cancer Center, Villejuif, France 2 Lifen, Paris, France 3 Informatic Team (DTNSI), Gustave Roussy, Villejuif, France 4 Informatic Team (DTNSI), Gustave Roussy, Villejuif, France 4 Informatic Team (DTNSI), Gustave Roussy, Villeju



- In cancer research, patient selection for clinical trials mainly depends on physicians manually identifying eligible patients. Automated approaches could increase the number of trials and reduce inclusion times.
- This study's objective was to evaluate the correctness and completeness of an artificial intelligence (AI)-driven approach for automatically inputting lung cancer patient information.

## **METHOD**

#### **POPULATION**

Patients with thoracic cancer seen at Gustave Roussy between Feb 2021 and June 2024.

#### MANUAL DATA ENTRY (MDE)

Manual retrospective collection of data in a secured RedCap database.

#### **AUTOMATED DATA ENTRY (ADE) – INPUT**

- Unstructured patient medical letters between February 2021 - July 2024.
- A schematic description of each variable.

#### **METHOD**

- Generative AI to find, quote and process variables into a structured form.
- Large language model (LLM) actions with prompt engineering and tailored few-shots examples.
- Mortality data were auto-extracted from the French public registry, INSEE.

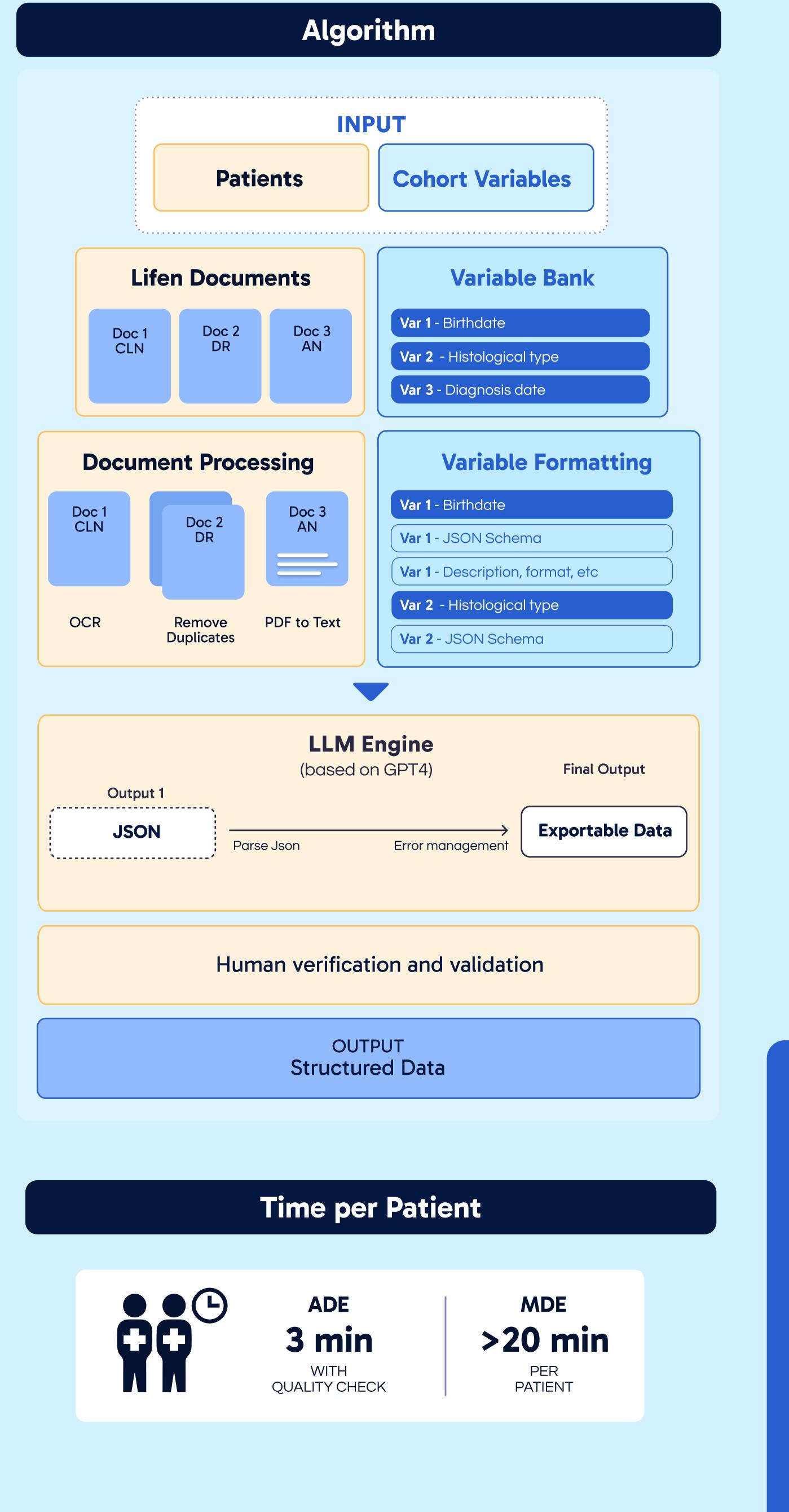
#### OUTPUT

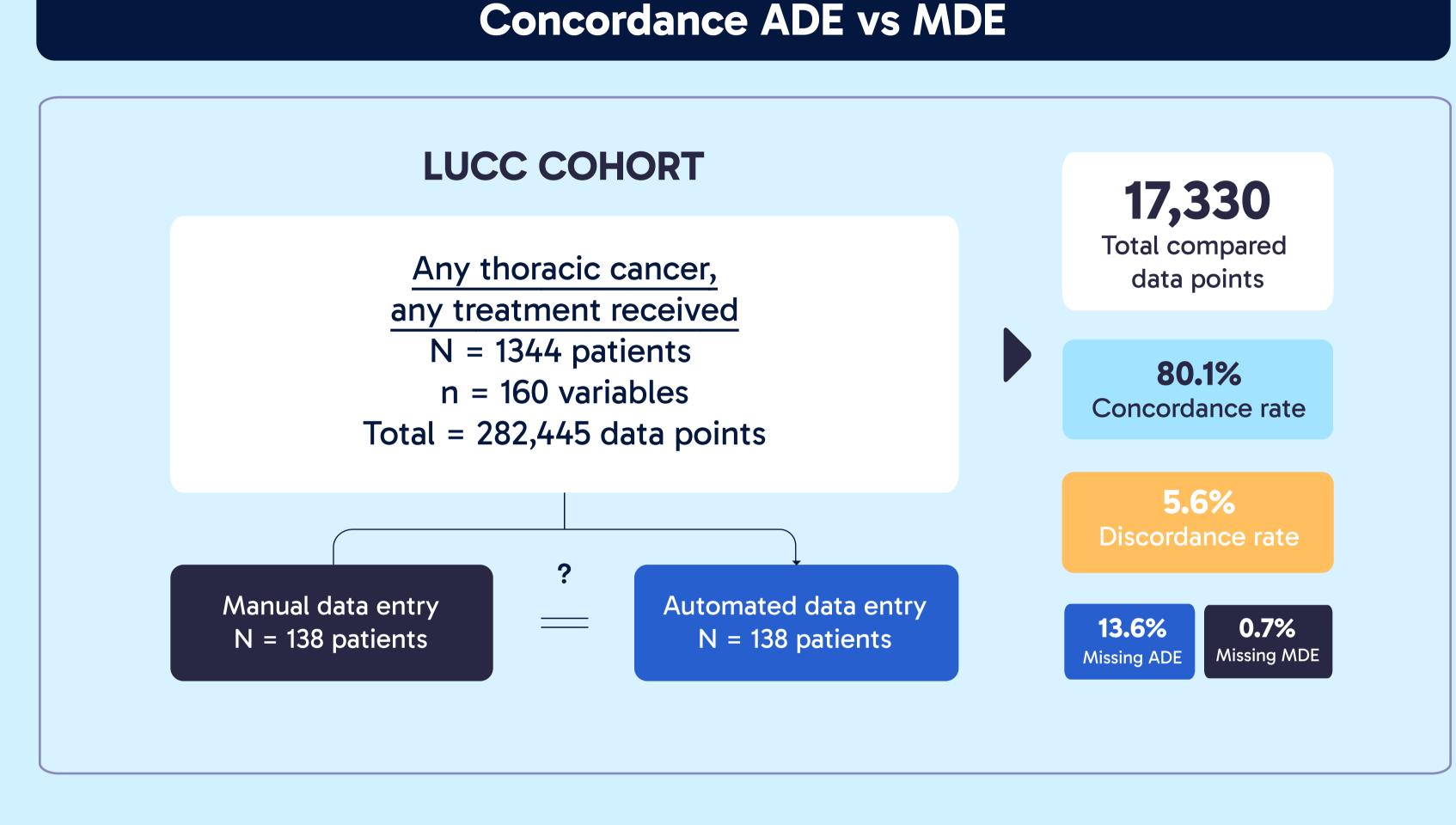
Demographics, disease characteristics, comorbidities, treatment history and life status.

### **METRICS**

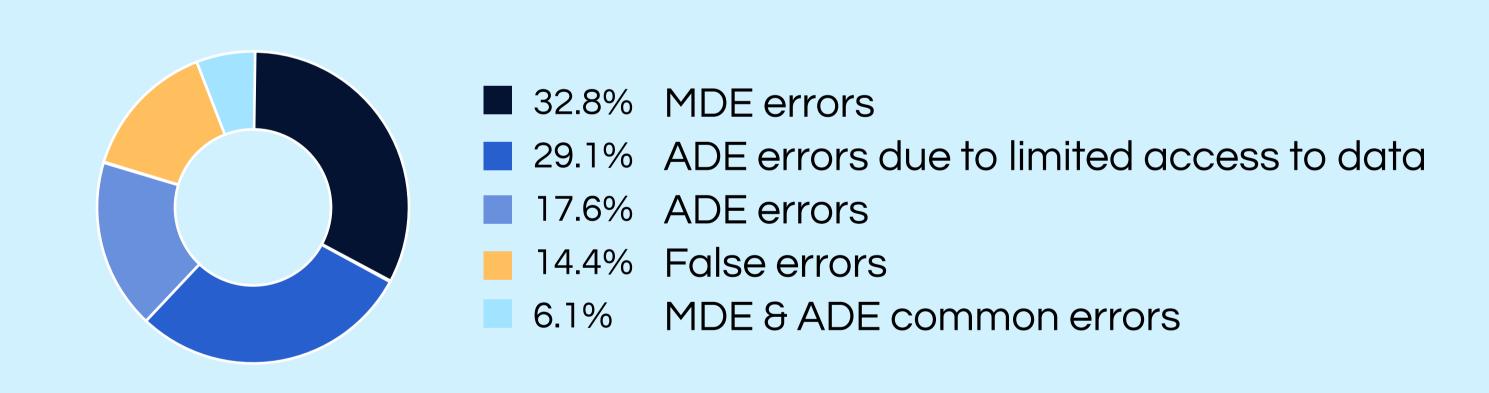
Concordance between comparable dates from MDE and ADE, secondary manual review for mismatches (senior physician), correctness (accuracy after checking), time per patient.

## **RESULTS**





#### Discordances check



ADE errors were mostly from data gaps in medical notes. Also, detailed information was often accessible to MDE in imaging or pathology reports, yet inaccessible to ADE.

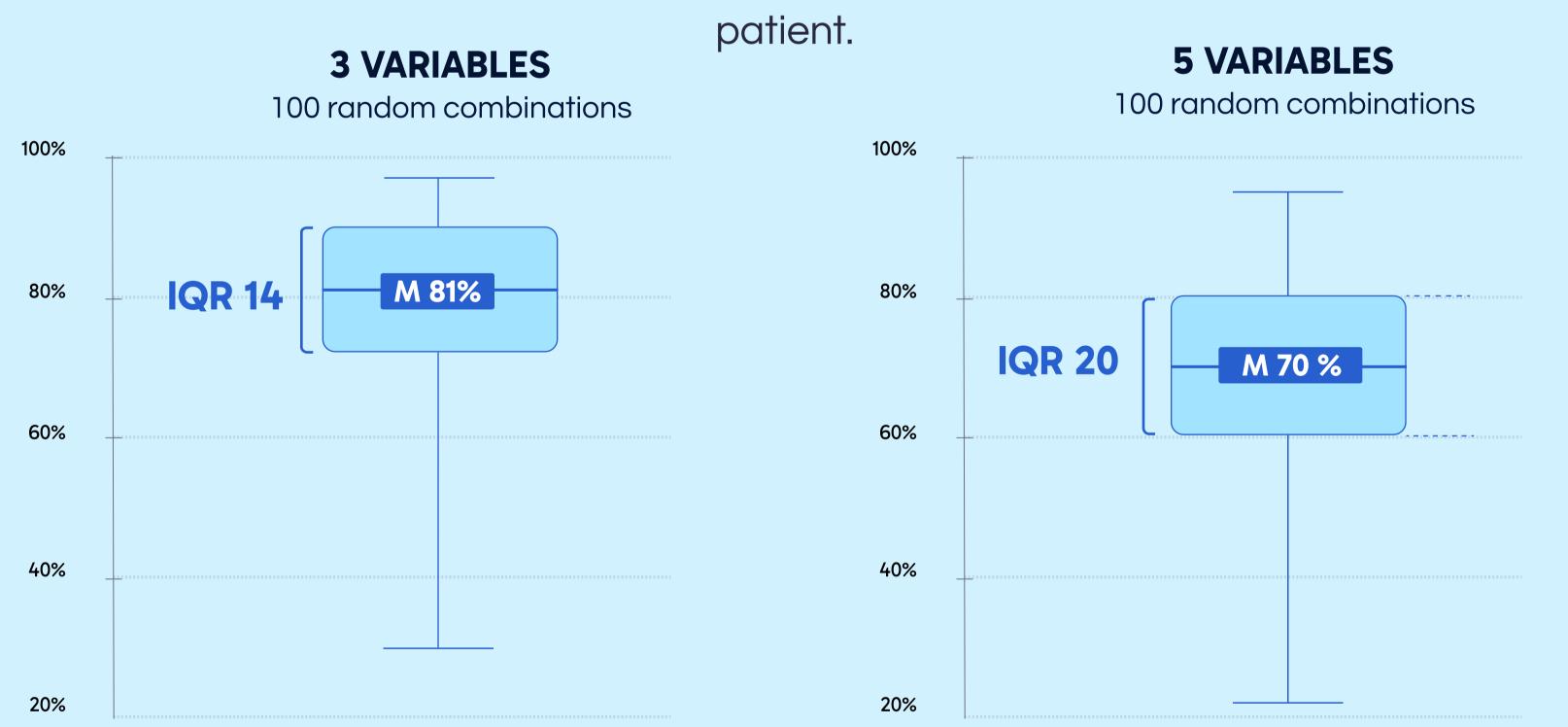
#### Correctness & Accuracy

Correctness was calculated after checking for discordances and excluding missing data on both sides.

#### 83.6% Overall correctness

CORRECTNESS 95% - 100%	85% - 94%	< 85%
Gender & Birthdate	Histology	Date first diagnosis/metastasis
Life status & date of death	Smoking status	Systemic treatment lines and
Comorbidities (Cardiac disease, autoimmune disease, HIV, hepatitis B and C, thromboembolic events etc.)	PDL1 expression	details
	Metastatic from diagnosis	Pack years  TAAD value
Molecular alterations (EGFR, BRAF, ROS1, RET, MET, HER2, PIK3CA, SMARCA4, KEAP1, NRG1, NTRK)	Metastatic anytime	TMB value  Brain and leptomeningeal  metastasis

Based on potential inclusion criteria of clinical trial, accuracy was calculated by how often multiple variables were correctly entered together at the same time for one



## CONCLUSION

- Generative AI can identify eligible clinical trial candidates with over 80% accuracy between ADE and MDE depending on selected criteria.
- High performance of ADE is seen with demographics, life status, histology, molecular alterations and comorbidities.
- ADE errors or missing data are often due to a lack of information in medical notes.
- ADE has the potential to enhance the efficiency, accuracy, and scalability of clinical trial pre-screening.