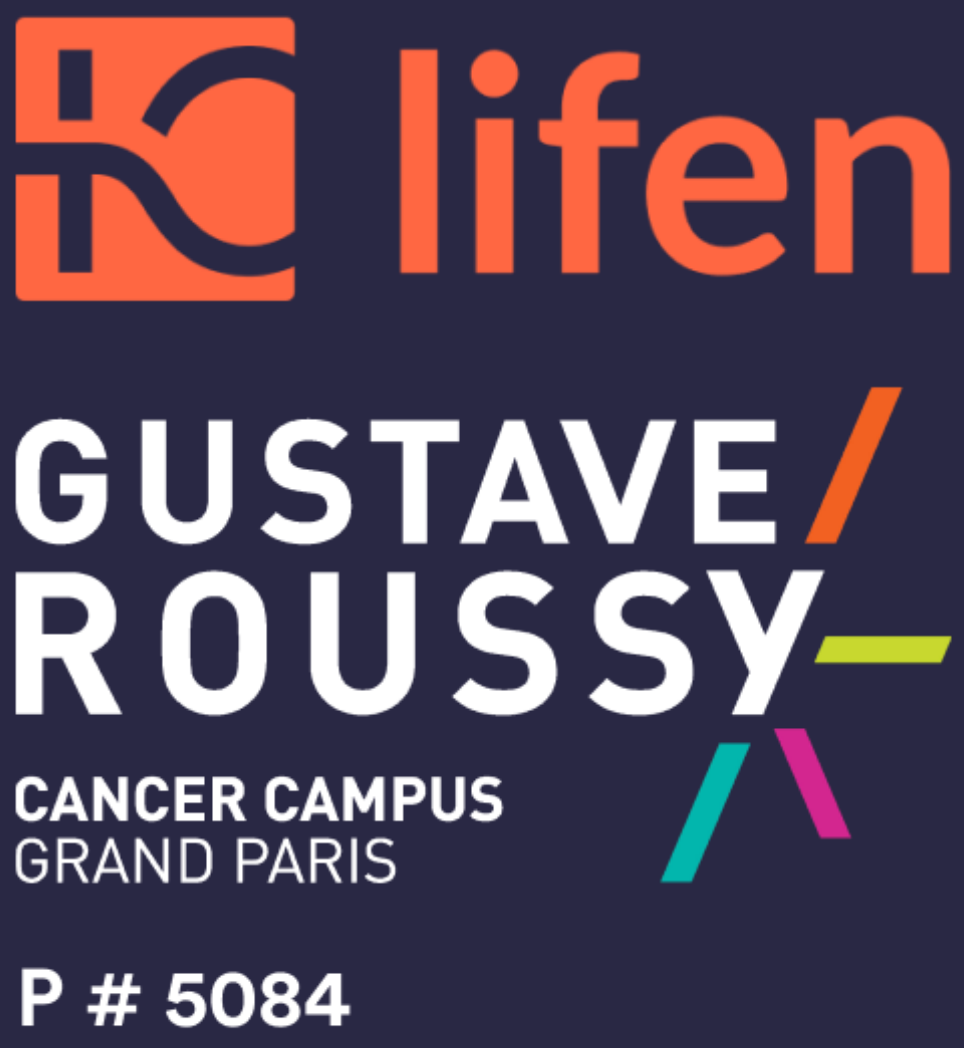


Using AI to Automatically Process Data from Unstructured Health Records of Patients with Lung Cancer

Mihaela Aldea¹, Pierre Rolland², Solenne Simon¹, Lodovica Zullo¹, Azeddine Djarallah², Aliette Poplu², Lisa Chuttoo², Muriel Wartelle³, Benjamin Vignal², Jean-Charles Louis², François Lion³, Arnaud Borie², David Planchard¹, Caroline Robert¹, Stefan Michiels⁴, Fabrice Barlesi¹, Franck Le Ouay², Benjamin Besse¹

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



BACKGROUND

- Structured databases created from electronic health records (EHR) are crucial for cancer research. Manual data entry into databases is both labor-intensive and error-prone.
- This study's objective was to create and validate an artificial intelligence (AI)-driven approach for automatically inputting lung cancer patient information from EHRs.

METHOD

POPULATION
Patients with thoracic cancer seen at Gustave Roussy between February 2021 and June 2023.

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 - January 2024.
- A schematic description of each variable.

METHOD

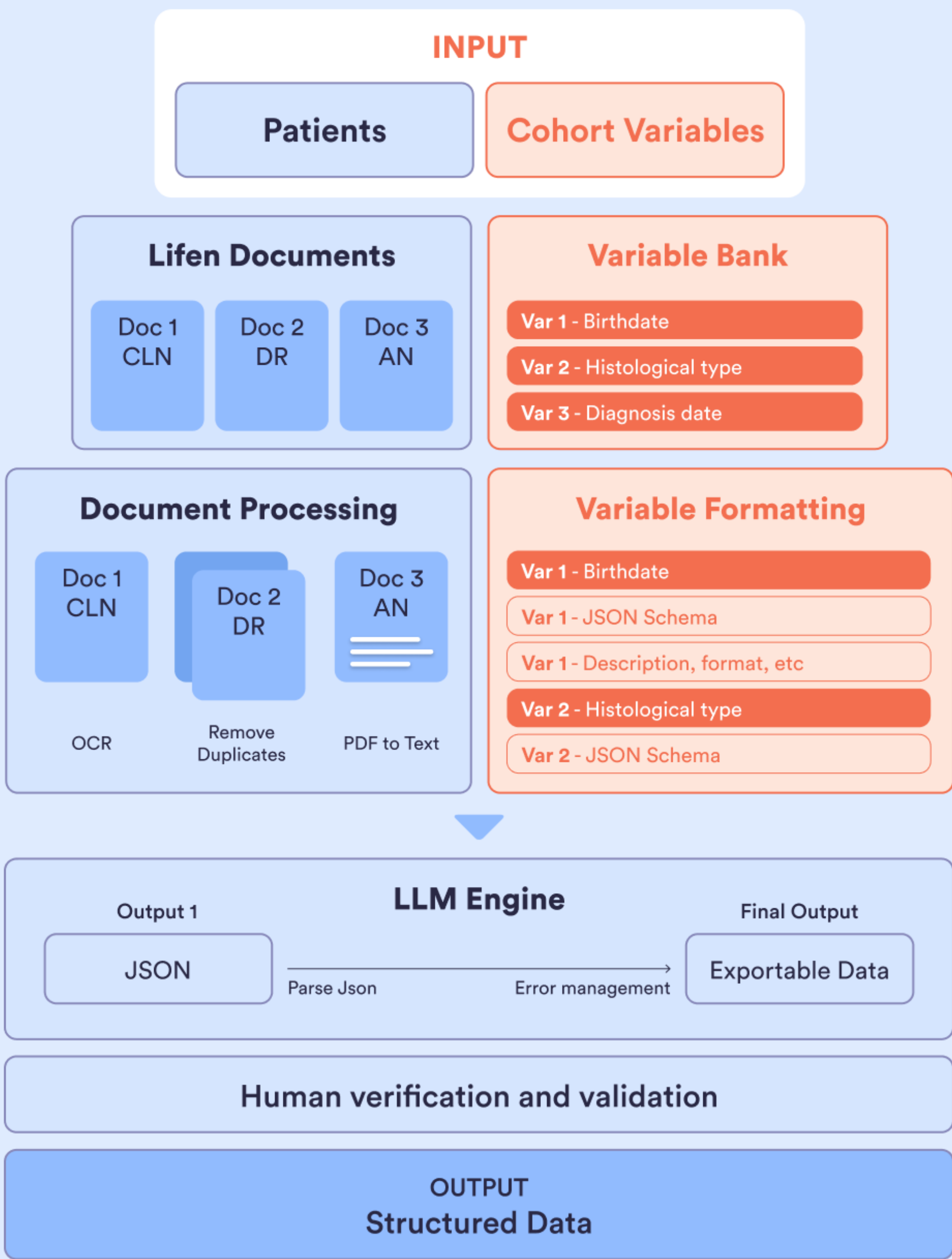
- Generative AI to find, quote and process variables into a structured form.
- Large language model (LLM) actions with prompt engineering and tailored few-shots examples.
- Mortality data were auto-extracted from the French public registry, INSEE.

OUTPUT
Demographics, risk factors, molecular profile, cancer history, treatment data, survival data.

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

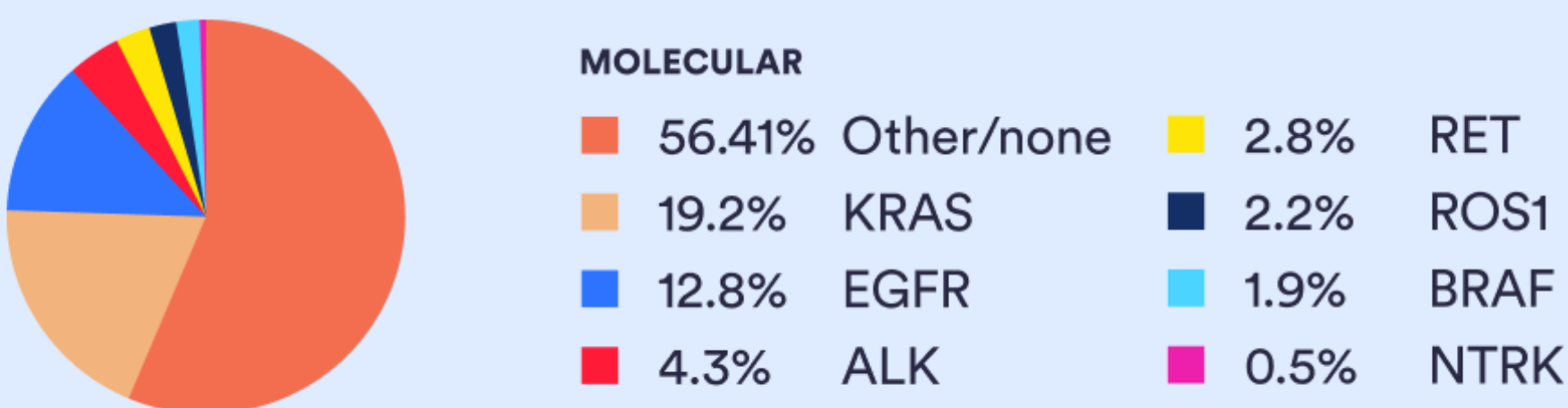
RESULTS

Algorithm



ADE Cohort

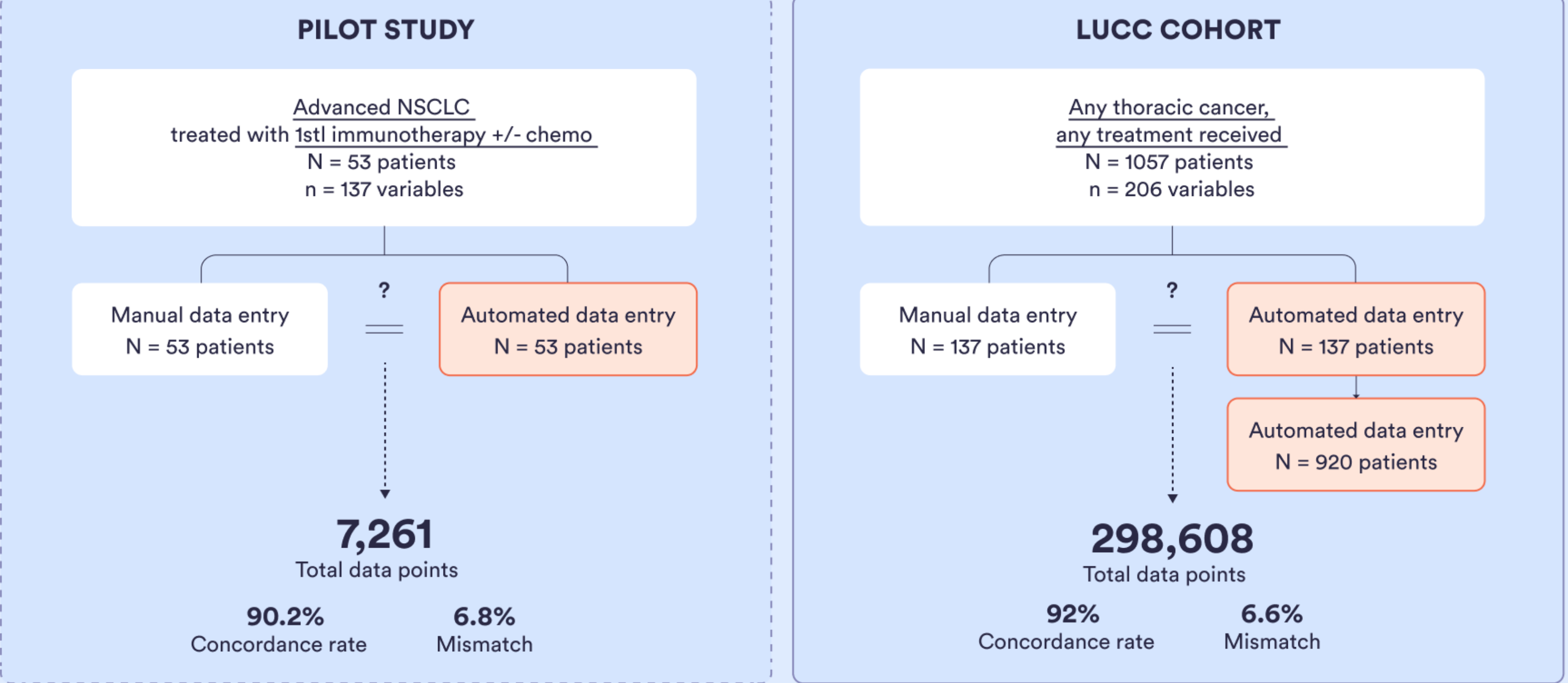
Genomic alterations in 861 NSCLC - automated data entry



Time per Patient



Concordance ADE vs MDE



CORRECTNESS 95% - 100%

Gender; Birthdate; Date of death
Asbestos exposure
COPD, myocardial infarction, autoimmune disease
Histology
Metastatic evolution anytime
Systemic treatment class, drugs, sequence

90% - 94%

Life status
Smoking (yes/no)
Dyslipidemia, diabetes
Family history of lung cancer
Metastatic from diagnosis
Molecular alterations
Sites of progression

80% - 89%

Current/former smoking status
Pack-years
Thromboembolic event
PDL1 score
Metastatic sites at time of each systemic treatment
Treatment discontinuation

70 - 79%

Cannabis consumption
Joint-year
Comorbidities: HTA
Date diagnosis and first metastasis
Stage cT, cN
Date of start treatment

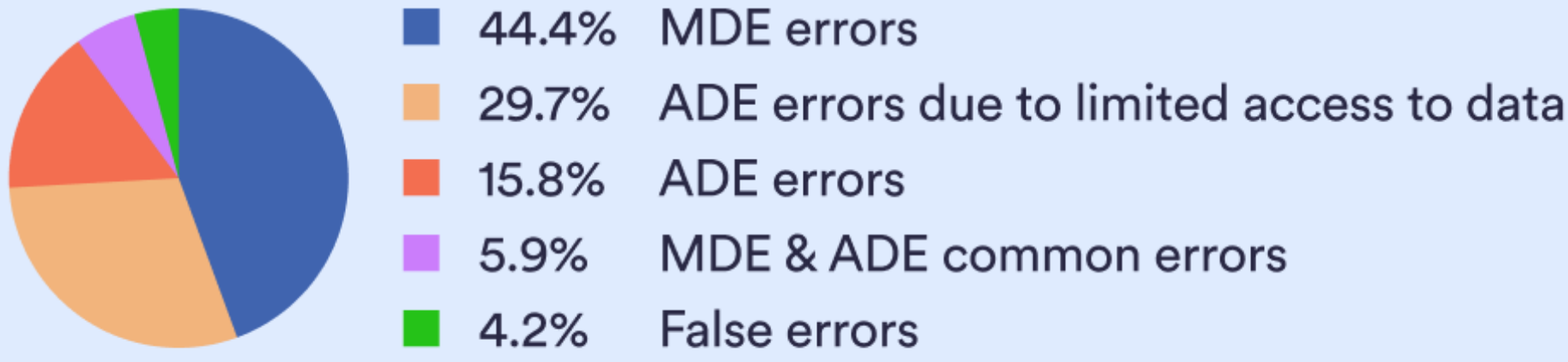
50 - 70%

Date first metastasis
Stage cM ;TMB value
Start date of each treatment I
Best objective response to each treatment I
Event of progression to each treatment I

< 50%

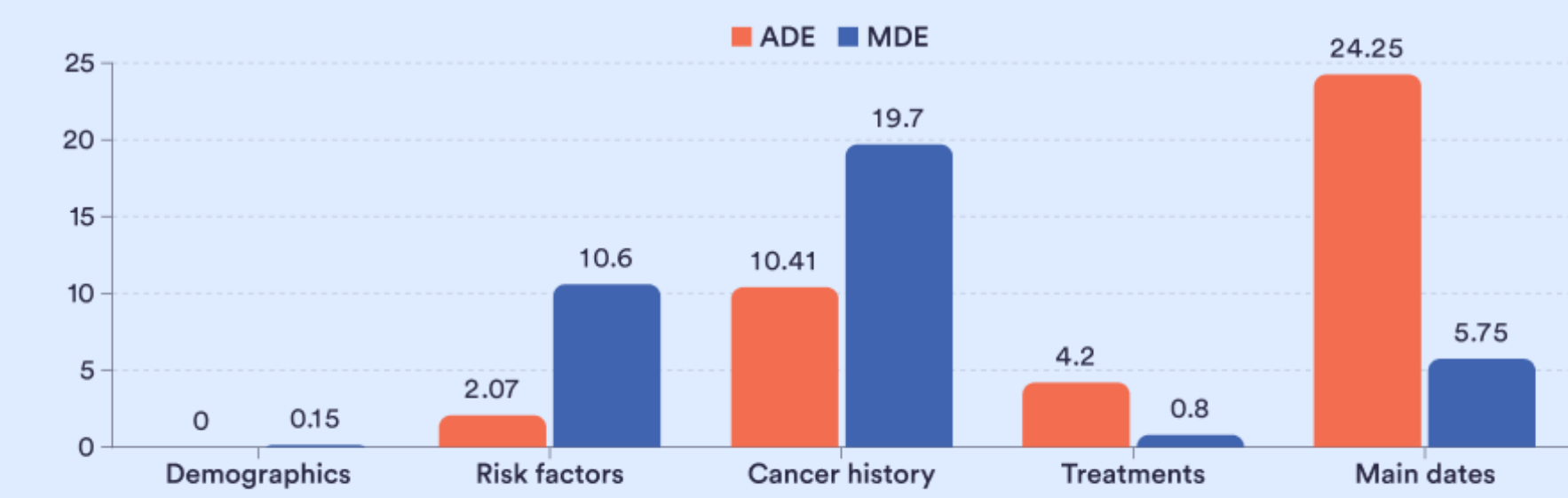
Date last follow-up for living patients
ECOG performance status at each treatment start
Date of last administration of treatment
Date of progression for each treatment

Mismatch



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

Missing Data



CONCLUSION

- Generative AI can identify and structure unstructured data from EHRs, with >90% concordance between ADE and MDE.
- High performance of ADE is seen with demographics, risk factors, comorbidities, histology, molecular profile and treatment types, while lower performance with dates (e.g. last follow-up, progression or last scan without progression).
- ADE cases of low correctness are often due to a lack of information in medical notes.
- ADE has the potential to enhance the efficiency, accuracy, and scalability of EHR-to-database conversions.