Leveraging Artificial Intelligence to Enhance the Quality and Efficiency of Real World Evidence Generation in HEOR

I. Delaroziere¹, J-E. Karcenty², T. Lacombe², A. Perez², L. Houvet², V. Susplugas², T. Vergnol¹, J-B. Angeloglou¹
Ospi, Lyon, France
Ospi, Paris, France

INTRODUCTION

Real World Data (RWD) and Real World Evidence (RWE) studies can be used to supplement randomized controlled trials (RCTs) and may help support clinical decision-making.

Recent RWE studies are shifting towards using secondary data sources generated through the delivery of hospital care to benefit from a higher granularity of clinical data. Electronic health records (EHRs), provide a variety of essential information such as socio-demographic characteristics, direct clinical data, medical diagnoses, prescribed and administered treatments, adverse events, results of biological and radiological examinations and much more that may not be reported in billing claim databases.

However, these data are mainly available in an unstructured manner (doctor's notes, discharge summaries,



RWD171

laboratory or imaging reports, consensus conference proceedings, etc.) and exploiting it requires a high-quality data structuring process to enable valid clinical assertions.

Treatments Hospital administrative data

Collected data related to patient's health status or delivery of care

METHOD

Objectives

We conducted an experiment in multiple French hospitals to assess how AI can translate Real World unstructured Data from EMRs into structured databases to generate Real World Evidence.

From Site Selection to Cohort Identification

Participating sites have been selected across French public and private hospitals already equipped with the Intelligence For Health (I4H) solution and actively willing to participate in this experiment. Data structuring and fine-tuning operations were performed in pseudonymized environments hosted within each site.

Members of the hospitals' clinical teams were onboarded to define variables with a high level of medical interest but difficult to obtain without human review of patients' EHRs. In parallel, the I4H solution enabled the medical staff to select a cohort





TRANSLATION OF PROTOCOL CRITERIA (ELIGILITY AND ENDPOINTS) IN VARIABLES

of patients of interest for their experiment based on broad criteria (age, sex, stay dates, or type of care received).

Automated Data Extraction Pipeline Using LLM to generate RWE

Our automated pipeline, powered by a large language model (LLM), enables efficient and consistent structuring of complex clinical data. Compared to eCRF, it offers significant time savings, eliminates the need for costly re-annotation when variables change, and ensures uniform annotations across sites, reducing discrepancies in multicentric trials.

Additionally, the LLM outperforms REGEX by better handling context, negations, and hypothetical scenarios. After validation by annotators, performance is evaluated automatically with metrics like precision, recall, and F1-score, achieving an average F1-score of 0.8 across key variables.

Once this IA pipeline was deemed of sufficient quality, it was replicated on other participating hospital's IT systems, leading to outputs that could be aggregated and analyzed as multicentric data.



Based on the selected variables and patient listing, a sample of hospital EHRs was extracted to a separate virtual environment where all personal information was removed. These texts were then analyzed through an iterative process using IA solutions to extract variables of interest in a structured manner.



RESULTS

LLMs allowed the use of Examining medical data summarization and analysis of healthcare data while still providing equivalent quality output compared to manual data entry.

In addition, **data entry time using our tool has been reduced from 120 minutes to 30 seconds per patient**, representing a substantial time saving for clinical staff.

Collaboration with hospital clinical teams has been instrumental in ensuring that the results produced are of the highest quality.

This has **proved useful for collecting key data** for health economics, such as observed efficacy of comparators, description of treatment sequences, list and frequency of adverse events for comparators, resources consumed, ...

CONCLUSION

Real World Data offers advantages over randomized controlled trials, such as timely and cost-effective data, large sample sizes for subpopulation analysis, and representation of real-world practices. However, challenges like biases, data quality issues, and the risk of misleading results hinder RWE's credibility.

We explore how **AI can improve the identification of clinically relevant data points** critical to the approval process through its ability to **better translate real-world unstructured data** from EHRs into structured databases, **increasing replicability and optimizing lead times and costs**.

The possibility of **easily iterating data collection once the automated pipeline is set up** offers significant advantages for recurring data collection and multicentric studies.