

---

## GPT-4 Identifies Cancer Phenotypes in Electronic Health Records



---

Extracting clinical phenotypes from electronic health records (EHRs) is an essential yet challenging task in healthcare informatics. This process involves identifying specific medical information such as disease stages, treatment modalities, and instances of recurrence. Accurate extraction of these phenotypes, particularly in complex conditions like non-small cell lung cancer (NSCLC), can significantly improve patient outcomes and advance medical research. [This study recently published in JAMIA Open](#) examines the application of OpenAI's GPT-4 model for extracting clinical phenotypes from EHRs, comparing its performance with other models, including GPT-3.5-turbo, Flan-T5, Llama-3-8B, and spaCy-based methods.

### The Challenge of Clinical Phenotype Extraction

Clinical phenotype extraction involves identifying medical concepts from unstructured text data found in EHRs. This task is crucial in understanding patient health, disease progression, and treatment efficacy, especially for diseases like NSCLC, which has diverse clinical presentations and treatments. Traditional methods for phenotype extraction, such as rule-based and machine learning-based approaches, often struggle with the variability and complexity of clinical language. Rule-based models like spaCy's medspaCy and scispaCy rely on predefined patterns and are limited by their inability to adapt to new or nuanced language uses.

Pre-trained transformer-based models offer a powerful alternative, particularly those using the GPT architecture. These models can understand and generate human-like text, making them particularly suited for tasks requiring deep linguistic and contextual understanding. Despite their potential, the application of these models in clinical settings, especially for detailed tasks like cancer phenotype extraction, has not been extensively explored.

### Methodology and Model Comparison

The study focused on a dataset of 13,646 clinical notes from 63 patients diagnosed with NSCLC, collected from Washington University in St. Louis. The primary objective was to identify critical phenotypes such as initial cancer stage, types of treatments administered, evidence of cancer recurrence, and the organs affected during recurrence. The performance of GPT-4 was compared against other models, including GPT-3.5-turbo, Flan-T5-xl, Flan-T5-xxl, Llama-3-8B, and spaCy-based models (medspaCy and scispaCy).

Each model's performance was assessed using precision, recall, and micro-F1 scores. These metrics provide a comprehensive measure of accuracy, balancing the identification of relevant instances (recall) with the precision of these identifications (precision). GPT-4 outperformed the other models across all metrics, demonstrating particularly high scores in precision and recall for the identification of cancer recurrence and treatment modalities. GPT-3.5-turbo, while slightly less accurate than GPT-4, still showed competitive performance, particularly regarding recall.

While effective in specific contexts, the spaCy-based models showed limitations in handling the complex and varied language found in clinical notes. These models struggled with identifying phenotypes expressed in diverse ways or embedded within complex sentence structures, highlighting the limitations of rule-based systems in capturing nuanced medical language.

### Discussion and Implications

The findings from this study underscore the significant advantages of using advanced transformer-based models like GPT-4 for clinical phenotype extraction. GPT-4's ability to understand and process complex medical language without the need for extensive preprocessing or rule-setting marks a substantial improvement over traditional methods. Its advanced pattern recognition capabilities enable it to accurately extract relevant information from unstructured text, making it particularly valuable in medical settings where precise and comprehensive data extraction is critical.

Moreover, the comparison with GPT-3.5-turbo reveals that while both models are effective, GPT-4 offers more consistent and contextually accurate outputs. This consistency is crucial in clinical applications, where even minor inaccuracies can lead to significant consequences in

patient care. The study also highlights the potential for GPT models to handle tasks beyond simple text generation, suggesting broader applications in clinical data processing and analysis.

The research also indicates that while rule-based models like medspaCy and scispaCy provide deterministic results and are useful for specific applications, they lack the flexibility and contextual understanding required for more complex tasks. The limitations observed in these models, particularly in handling contextual ambiguities and variations in language, underscore the need for more sophisticated tools like GPT-4 in clinical settings.

This study demonstrates the efficacy of GPT-4 in extracting clinical phenotypes from EHRs, offering a significant improvement over traditional methods. By leveraging advanced natural language processing capabilities, GPT-4 provides more accurate, comprehensive, and contextually relevant data extraction, which is critical for improving patient care and advancing medical research. Future research should focus on further optimising these models, exploring their applications in other medical domains, and addressing any remaining limitations related to data diversity and model variability.

**Source Credit:** [JAMIA OPEN](#)

**Image Credit:** [iStock](#)

Published on : Thu, 1 Aug 2024