

## Background

With the rapid development of big data and deep learning techniques in recent years, natural language processing (NLP) systems have been widely used in the healthcare domain. However, evaluating an NLP system for a specific clinical use case is challenging due to the complexity of creating the gold standard corpus (GSC).

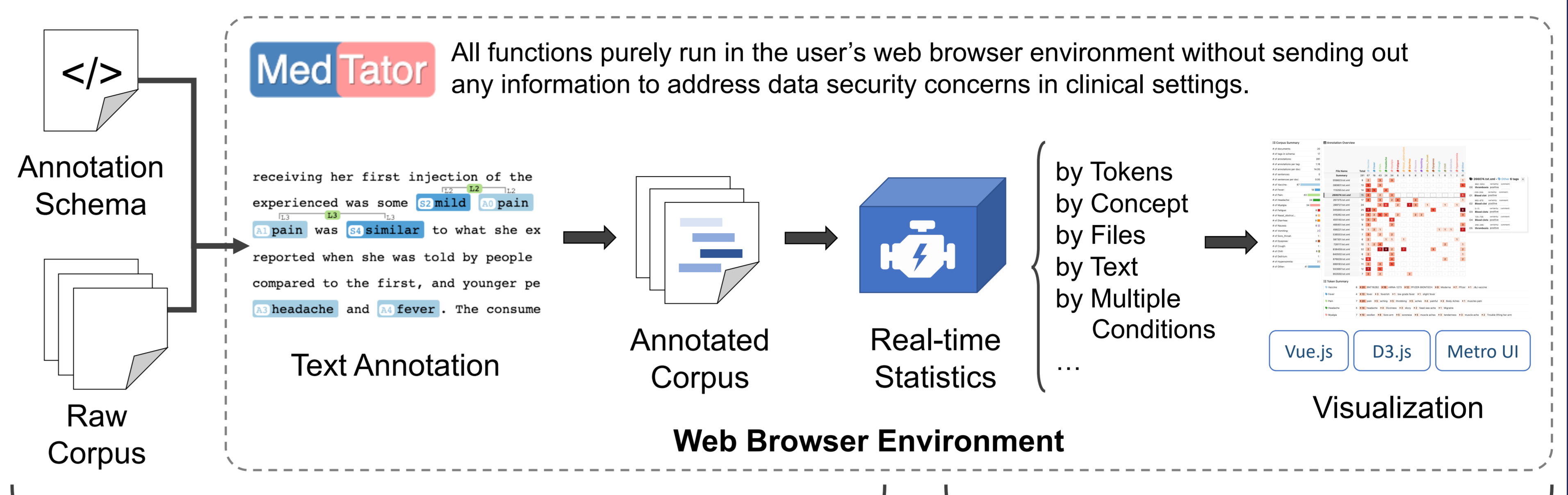
Existing text annotation tools can provide powerful features to cover various needs for corpus development, but few tools support various needs during the annotation process.

To address these needs, we developed multiple corpus visualization features in **MedTator**, a **serverless** annotation tool to create GSC for building NLP systems.

You can check our online demo at: <https://ohnlp.github.io/MedTator/>



## System Overview



Users can annotate the corpus by using the Annotation Tab in MedTator, and the text annotations are saved in memory.

The annotations are analyzed in real-time and shown in the Statistics Tab.

## Visualization and Interactivity Designs

### Statistics Tab

MedTator provides a separate Statistics Tab to visualize the annotated corpus

The Statistics Tab provides a comprehensive overview of the annotated corpus. It includes a **Corpus Summary** with metrics like the number of documents, tags in schema, and annotations per tag. A **Token Summary** lists various concepts such as Vaccine, Fever, Pain, Headache, Myalgia, Fatigue, Nasal\_obstruction, Diarrhea, Nausea, Vomiting, Sore\_throat, Dyspnea, Cough, Chill, Delirium, and Hypersomnia. A **Heatmap** visualizes the distribution of these concepts across different files. A **Color Legend** explains the color encodings used for different concepts. A **Download** button allows users to export the data as a CSV file.

The annotation overview and other statistical results can be exported as a multi-tab Excel file for sharing and further usage. Users can also export the basic statistics to a CSV file.

The annotations from two annotators can be compared to calculate the inter-annotator agreement for evaluation

### Adjudication Tab

adjudicating annotation results of two annotators

The Adjudication Tab is used for comparing annotations from two different annotators. It displays **Inter-annotator Agreement (IAA)** metrics such as Overall F1 score and Agreement percentage. A **Tag Name** list shows the agreement for various concepts. A **Multi-level Agreement** section displays bar charts for Overall, Concept-level, and File-level agreements. A **Text Matched Files** section lists files with their respective agreement scores. A **Download** button is available for exporting the results.

The multi-level agreements are displayed in corresponding bar charts:

- Overall agreement
- Concept-level agreement
- File-level agreement

By clicking each bar, the user can check detailed tags for adjudication.

### Annotation Tab

MedTator provides all core functions for text annotation tasks for multiple documents

The Annotation Tab is the primary interface for adding and managing annotations. It features a **Schema File** (YAML/JSON) and an **Annotation File** (XML). The main text area shows a document with highlighted entities and their corresponding tags. A **Tag List** on the right side provides a detailed view of all tags, including their ID, spans, text, and attributes. A **Entity Tags** dropdown menu allows users to select from a list of predefined concepts for annotation.

When new tags are added or new annotated files are added, the real-time statistics will be updated within the user's web browser.

Users can switch to the Statistics Tab to check the annotation progress or analyze the imported corpus at any time.

## Future Work

We are still working on improving the performance and usability of MedTator and its corpus visualization module, but we have been using MedTator in our internal NLP projects for corpus development, adjudication, and building rule-based NLP systems.

We welcome you to use MedTator in your projects and leave any comments on the MedTator's repo issues: <https://github.com/OHNLPMedTator/issues>