

Background

A gold standard annotated corpus is usually indispensable for building high-quality language models for clinical natural language processing.

Existing text annotation tools can provide powerful features to cover various needs for corpus development, but few tools provide real-time visualization to support the needs of annotation analysis during the annotation process.

To address this need, we developed a corpus visualization module in **MedTator**, a **serverless** annotation tool.

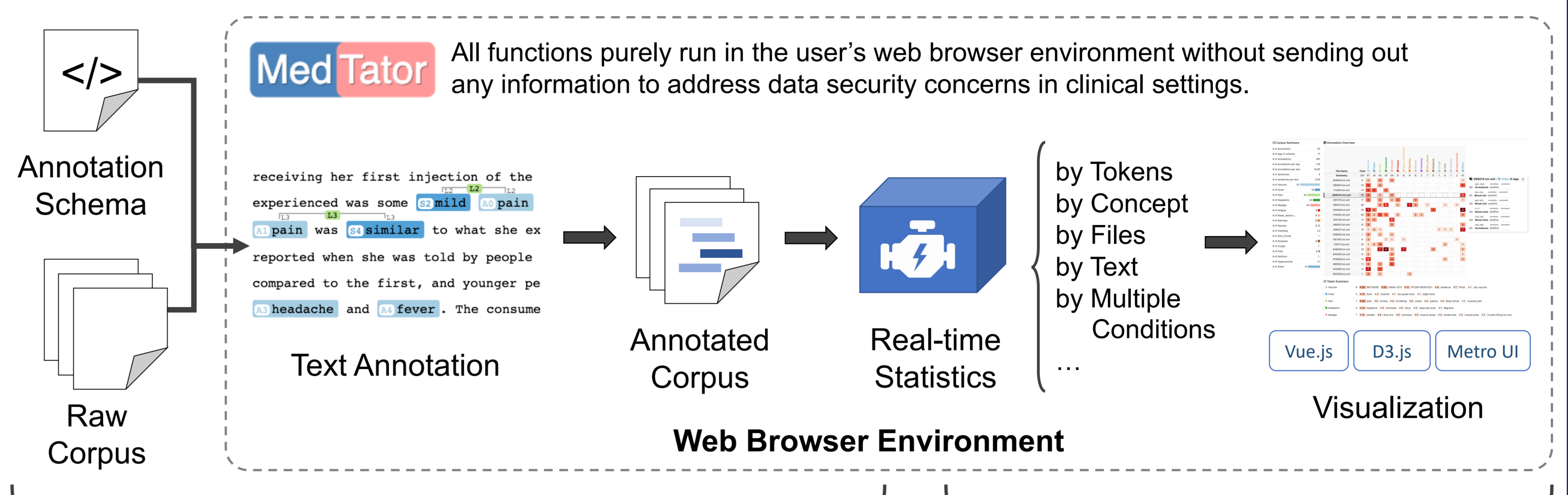
MedTator doesn't require any runtime installation! You can check our online demo at:

<https://ohnlp.github.io/MedTator/>

Source codes are available at: <https://github.com/OHNLPMedTator>



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

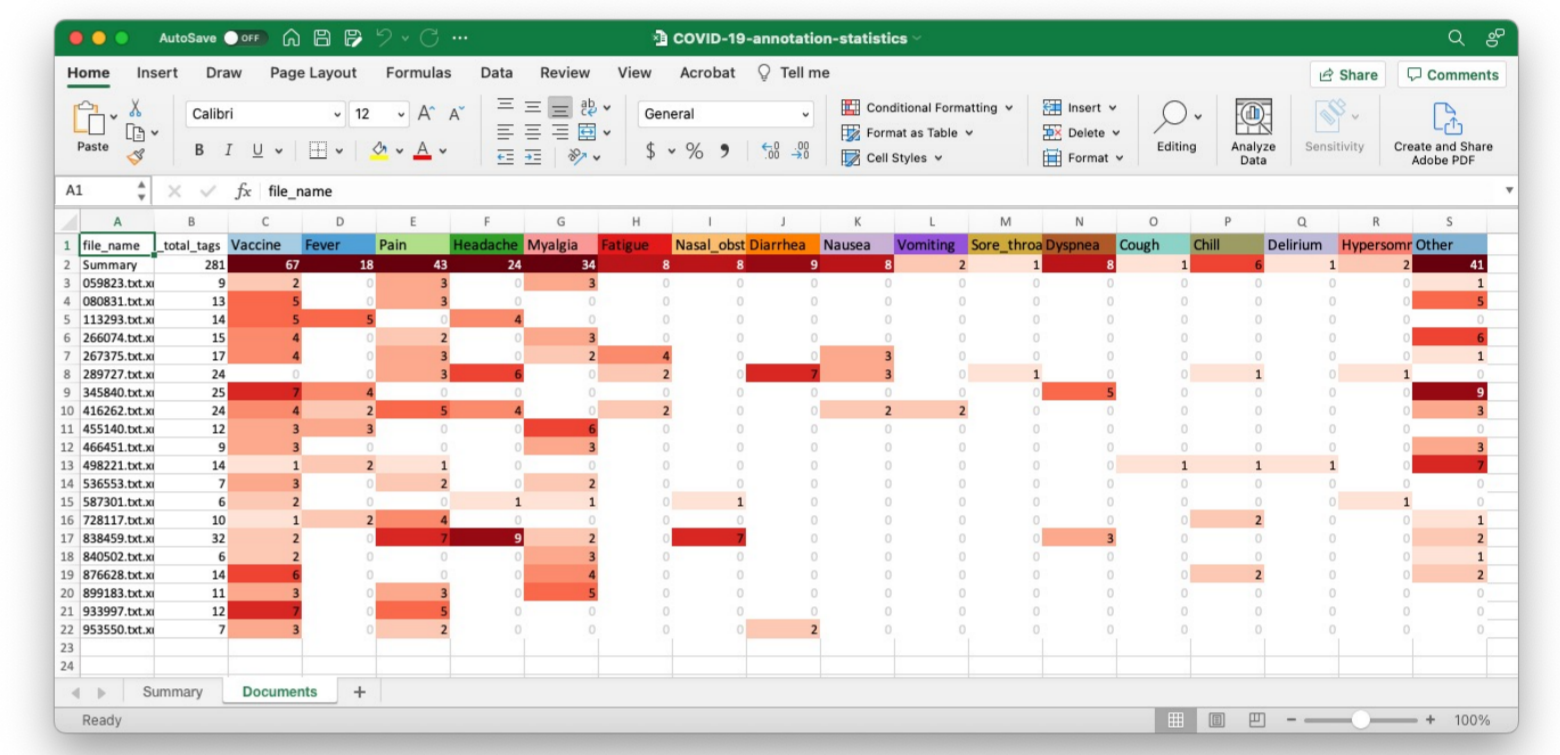
MedTator provides a separate Statistics Tab to visualize the annotated corpus

The Statistics Tab interface includes a navigation bar with options like Annotation, Statistics, Export, Adjudication, and Converter. It features a **Corpus Summary** section with various metrics (e.g., # of documents, # of tags in schema). The **Annotation Overview** section displays a heatmap where rows represent files and columns represent different annotation concepts. A **Token Summary** section lists various concepts and their occurrences. A **Export Tab** section allows users to export data in various formats (Text, Tag & Sentence, IOB2/BIO, etc.).

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 real-time statistics are shared by the Export Tab for exporting the results to other formats.

The source files of each token are listed, and the user can click the file name to look back to the selected annotation file in the Annotation Tab for revision.



Export Tab

the annotations and statistics can be exported to multiple formats

The Export Tab interface shows options for exporting data in various formats: Text, Tag & Sentence, IOB2/BIO (tsv), BioC (xml), MedTagger (zip), spaCy (jsonl), and How to use. An **Export Preview** section shows a sample of the exported data, including text snippets and their corresponding annotations.

IOB2/BIO format for training NLP models (e.g., fine-tuning BERT NER models)

MedTagger format for building rule-based NLP systems or spaCy based rulesets

Annotation Tab

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

The Annotation Tab interface includes a navigation bar with options like Annotation, Statistics, Export, Adjudication, and Converter. It features a **Annotation File** section where users can load and save files. The **Annotation** section shows a text document with various annotations (e.g., "Blood clots", "Sore arm", "Pain"). A **Entity Tags** dropdown menu is visible on the right. The **All Tags** section at the bottom lists all annotated tags with their IDs, spans, and text.

When new tags are added or new annotated files are added, the real-time statistics will be updated in the backend.

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>