

An Open Natural Language Processing (NLP) Framework for EHR-based Clinical Research: A Case Demonstration Using the National COVID Cohort Collaborative (N3C)

Sijia Liu^{1,*}, Andrew Wen^{1,*}, Liwei Wang^{1,*}, Huan He^{1,*}, Sunyang Fu^{1,*}, Robert Miller², Andrew Williams², Daniel Harris³, Ramakanth Kavuluru³, Mei Liu⁴, Noor Abu-el-rub⁴, Dalton Schutte⁵, Rui Zhang⁵, Masoud Rouhizadeh⁶, John D. Osborne⁷, Yongqun He⁸, Umit Topaloglu⁹, Stephanie S Hong¹⁰, Joel H Saltz¹¹, Thomas Schaffter¹², Emily Pfaff¹³, Christopher G. Chute¹⁰, Tim Duong¹⁴, Melissa A. Haendel¹⁵, Rafael Fuentes¹⁶, Peter Szolovits¹⁷, Hua Xu¹⁸, Hongfang Liu¹, National COVID Cohort Collaborative (N3C) Natural Language Processing (NLP) Subgroup, National COVID Cohort Collaborative (N3C)

1. Department of Artificial Intelligence and Informatics, Mayo Clinic.
2. Tufts Clinical and Translational Science Institute, Tufts Medical Center.
3. Department of Internal Medicine, University of Kentucky.
4. Department of Internal Medicine, University of Kansas Medical Center.
5. Department of Pharmaceutical Care & Health Systems, University of Minnesota at Twin Cities.
6. Department of Pharmaceutical Outcomes & Policy, University of Florida
7. Department of Computer Science, University of Alabama at Birmingham.
8. University of Michigan Medical School.
9. Wake Forest Baptist Medical Center.
10. Department of Medicine, Johns Hopkins University.
11. Department of Biomedical Informatics, Stony Brook University.
12. Sage Bionetwork.

13. Department of Medicine, University of North Carolina Chapel Hill.
14. Albert Einstein College of Medicine.
15. Center for Health AI, University of Colorado Anschutz Medical Campus.
16. Alex Informatics.
17. Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology.
18. School of Biomedical Informatics, University of Texas Health Science Center at Houston.

* These authors contributed equally to this study

Abstract

Despite recent methodology advancements in clinical natural language processing (NLP), adoption of clinical NLP models within the clinical and translational research community remains hindered by issues with ETL process heterogeneity and human factor variations. In this study, we proposed an open NLP development framework with the aim of addressing these issues. The viability of such a platform was evaluated on a COVID-19 use case through sites participating in the National COVID Cohort Collaborative (N3C). As part of our assessment of the impact of single vs. multi-site NLP algorithm development, we evaluated the performance of both an NLP ruleset developed solely using a single site's clinical narratives as well as one further refined using a synthetic derived dataset sourced from three sites (Mayo, UKen, and UMN). The single-site ruleset resulted in performances of 0.876, 0.706, and 0.694 in F-scores for Mayo, Minnesota, and Kentucky test datasets, respectively, while the multi-site NLP ruleset improved performances to 0.884, 0.769 and 0.806. The results of our use case test run inform us of the importance of a multi-site federated development, evaluation, and implementation framework. As such, we aim to meet

this need with our framework by providing the tools necessary to conduct NLP development in a collaborative manner through consensus building, process coordination, and best practice sharing.

Introduction

Over the past decade, Electronic Health Record (EHR) systems have been increasingly implemented at US healthcare institutions. Large amounts of detailed longitudinal patient information, including lab tests, medications, disease status, and treatment outcomes, have consequently been accumulated and made electronically available. These large clinical databases are valuable data sources for clinical and translational research. As a result, major initiatives have been established to exploit this crucial resource, including the Clinical and Translational Science Awards (CTSA) Program's National Center for Data to Health (CD2H)/National COVID Cohort Collaborative (N3C)^{1,2}, the Electronic Medical Records and Genomics (eMERGE) Network³, the Patient-Centered Outcomes Research Institute's (PCORI) Clinical Research Networks (CRNs)⁴, the NIH All of Us Research Program⁵, and the Observational Health Data Science and Informatics (OHDSI) Consortia with demonstrated successes^{6,7,8,9}.

One common challenge faced by those initiatives is, however, the prevalence of clinical information embedded in unstructured text¹⁰. Compared to structured data entry, text is a more conventional way in the healthcare environment to document impressions, clinical findings, assessments, and care plans. Even with the advent of sophisticated EHR systems, studies have shown that capturing health information fully in structured format through data entry is unlikely to happen and a blended model where physicians use templates when and where possible and dictate the details of a patient visit in text¹¹.

Natural language processing (NLP) has been promoted as having a great potential to extract information from text¹². NLP algorithms can generally be categorized into using either symbolic

or statistical methods¹³. Since the turn of the century, machine learning algorithms (i.e., statistical NLP) have gained increased prominence in clinical NLP research¹⁴. Nevertheless, a substantial portion of clinical NLP use cases leverages symbolic techniques given that dictionary or rule-based methodologies suffice to meet the information needs of many clinical applications under specific use cases. In the context of EHR-based clinical research, NLP has been leveraged to assist information extraction and knowledge conversion at different stages of research including feasibility assessment, eligibility criteria screening, data elements extraction, and text data analytics. As a result, an increasing number of clinical research benefits from state-of-the-art NLP solutions and have been reported ranging from disease study areas^{15, 16, 17, 18} to drug-related studies^{19, 20}. A majority of existing clinical NLP studies are, however, done within a mono-institutional environment¹³, which may suffer from limited external validity and research inclusiveness. Compared with single-site research, multisite research potentially offers larger sample size, more adequate representation of participant demographics (e.g., age, gender, race, ethnicity, and social-economic status), and more diverse investigator expertise, which may ultimately yield a higher level of research evidence^{21, 22, 23, 24}.

Despite a plethora of recent advances in adopting NLP for clinical research, there have been barriers towards adoption of NLP solutions in clinical and translation research, especially in multi-site settings. The root causes of these barriers can be categorized into two major reasons: 1) heterogeneity of ETL (extract, transform, load) processes between differing sites with their own disparate EHR environments, and 2) human factor variation in gold standard corpus development processes.

ETL Process Heterogeneity. The challenges faced by NLP development and evaluation to facilitate the secondary use of EHR data originate from the complex, voluminous, and dynamic nature of the data being documented and stored within a heterogeneous set of disparate, institution specific, EHR implementations. Variations in EHR system vendors, data infrastructure (e.g., unified, ontology driven, and de-centralized), and institutions' modes of operation can lead to

idiosyncratic ways of clinical documentation, transformed, and representation²⁵. Collecting these data would require a significant expenditure of effort to locate, retrieve, and link EHR data into a specific format²⁶. This variability in ETL processes required to support a high level of data heterogeneity brings additional challenges in the adoption of NLP for clinical and translational research, which substantially limits both the cross-institutional interoperability of developed NLP solutions and the reproducibility of the associated evaluations.

Human factor variation in gold standard corpus development process. The process of developing, evaluating, and deploying NLP solutions in both mono- and multi-site environments can be task-specific, iterative, and complex, often involving a multitude of stakeholders with diverse backgrounds^{13,26}. A key step prior to model development is corpus annotation, the process of developing a gold standard by marking the occurrence of both task-defined sets of clinical information as well as their associated interpretative linguistic features (e.g., certainty, status) within text documents. Due to the complexity of clinical language, creating such gold standard corpora requires significant expenditure of domain expertise and time as clinical experts regularly make decisions directly affecting study cohort, annotation guideline, and task definitions. Studies have discovered potential biases in clinical decision making and interpretation of clinical guidelines²⁷, in coding of clinical terminologies²⁸, and in interpretation of imaging findings²⁹. This issue can be further exacerbated when conducting multi-site collaborations due to inter-site variations in care practice^{30,31}, ultimately affecting the validity and reliability of the resulting gold standard corpus. A coordinated, transparent, and collaborative platform is therefore needed to promote open team science collaboration in NLP algorithm development and evaluation through consensus building, process coordination, and best practice sharing.

Built upon our previous work^{32,33}, here, we proposed an open NLP development framework to address the aforementioned issues through the following components: 1) an interoperable NLP infrastructure for incorporation of different NLP engines utilizing a clinical common data model for data source interfacing and representation with the aim of reducing the impact of the

heterogeneity of ETL processes; 2) a transparent multi-site participation workflow on corpus development and evaluation with the aim of addressing the variation in data abstraction and annotation processes between sites; and 3) a user-centric crowdsourcing interface for collaborative ruleset development that enables effectively and efficiently gathering, synthesizing, and fusing site-specific knowledge and findings. To demonstrate the viability of the framework, we conducted a case study where we developed, evaluated, and implemented an NLP algorithm for extracting^{34, 35, 36} COVID-19 signs and symptoms to support the National COVID Cohort Collaborative (N3C).

Results

Framework Description

The framework itself consists of a data ingestion layer, a processing layer, and a data persistence layer. The architecture of the proposed framework is illustrated in Figure 3. The **data ingestion layer** works as the data collector with the ability to read text from a configurable variety of data sources such as relational databases or file systems including and load them into the NOTE table of OMOP CDM. The **processing layer** serves as the NLP engine where information extraction from raw texts happens given a set of heuristic rules created by various NLP engines. By default, as an example implementation, the MedTagger³⁷ NLP engine is provided, although alternative NLP engines can be substituted by wrapping their respective NLP pipelines to conform to a provided API specification. After the term modifiers added by contextual rules from ConText Algorithm³⁸ around the extracted condition mentions, these conditions will compose clinical events with temporal information. The reason we opt for a symbolic solution is due to its simplicity, transparency, and interpretability as the outcomes are fully deterministic based on the definition of the rules. When the baseline rulesets and dictionaries are made available to the public, they can therefore be easily refined by different users from different sites. The **data persistence**

layer stores resulting extracted NLP artifacts in the OMOP CDM NOTE_NLP table as the events are extracted from NLP systems.

The framework is distributed as open-source software under the Apache 2.0 license via Github in three parts: 1) ETL Backbone (<https://github.com/OHNLP/Backbone>) with an example NLP engine (<https://github.com/OHNLP/MedTagger>), 2) process documentation (<https://github.com/OHNLP/N3C-NLP-Documentation/wiki>), 3) open-source collaborative platform for developing NLP rulesets (<https://github.com/OHNLP/OHNLP TK>). The demo homepage (Figure 2(a) - <https://ohnlp4covid-dev.n3c.ncats.io/>) demonstrates the N3C NLP engine outputs on annotating clinical text using the baseline rulesets and dictionary. The annotations are from components of Sign/Symptom extractor, temporal information extractor and dictionary lookup extractor. To further customize each model, the users can visit “Rule Editor” (https://ohnlp4covid-dev.n3c.ncats.io/ie_editor) and the “Dictionary Builder” (https://ohnlp4covid-dev.n3c.ncats.io/dict_builder) page (Figure 2(b)). Figure 2(c) provides an example of the rules editing interface with the baseline COVID-19 ruleset. The rulesets can be tested in real time by clicking the “Upload and test” button, where the rulesets will be uploaded, and the NLP engine will be generated for testing and debugging purposes. As a use case study, we also provide an example NLP project for extracting signs/symptoms related to COVID-19 that was developed as an example use case for this framework. The elements with original texts such as text snippets and concept mentions are truncated before submission.

N3C Case Study

NLP Algorithm Development and Evaluation: Table 1 shows the annotation corpora statistics. A COVID-19 sign/symptom ruleset was produced consisting of 17 concepts. The IAA of the annotated corpus was 0.686 F1-score for Mayo, 0.516 for UMN and 0.211 for UKen. Two NLP algorithms were evaluated in this study. One was developed based solely on the narratives sourced from a single site (Mayo Clinic). The other used the resulting NLP algorithm from the single site

and fine-tuned based on the annotated training data from an additional two sites (UMN and UKen). Table 2 shows the performance of the single-site NLP algorithm and Table 3 shows the performance of the multi-site NLP algorithm. The single-site ruleset resulted in performances of 0.876, 0.706, and 0.694 in F-scores for Mayo, Minnesota, and Kentucky test datasets, respectively, while the multi-site NLP ruleset improved performances to 0.884, 0.769 and 0.806. The performance of the multi-site NLP algorithm was better than that of the single-site NLP algorithm, but both showed a degrading trend from Mayo site to other sites.

Tables 4, 5 and 6 show the results of error analysis for the three sites. For FP, major discrepancies between the NLP algorithm and the gold standard were due to the NLP algorithm extracting mentions that are not COVID signs/symptoms but for instruction/patient education, adverse events/indication of treatment, clinical goal/precaution, template, etc. It should be noted that gold standards were not always correct, and in some notes, it was hard to judge if the mentions are COVID signs/symptoms when symptoms are not appearing with COVID or de-identified dates are inconsistent. For FN, reasons include NLP algorithm not complete, tokenization error due to de-identification process, template, and annotation errors.

Discussion

In this study, we proposed an open NLP development framework with the following properties: an interoperable NLP infrastructure, a transparent multi-site participation workflow, and a user-centric crowdsourcing interface. The key goal of this framework is to facilitate multi-site collaborative development, evaluation, and implementation of NLP algorithms. The framework has been implemented to support efforts conducted by the National COVID Cohort Collaborative (N3C) to enable the utilization of unstructured text in high throughput.

Here, we have presented our results from running our framework using a centralized annotation process on texts sourced from multiple sites after de-identification, with the aim of assessing the impact on NLP algorithm development (single-site algorithm vs multi-site algorithm). Several pragmatic implementation challenges were discovered that may impact the intermediate and final NLP results. We observed that IAA varied greatly between the three sites despite the fact that annotators had been trained using de-identified Mayo notes (0.686 F1-score for Mayo, 0.516 for UMN, 0.211 for UKen). Firstly, utilizing a centralized annotation approach, the process of text data collection took a very long time because each site needs to complete de-identification before sharing data. Secondly, it was a challenge for annotators to work on annotation tasks that spanned a long period of time. Thirdly, the shared data sets were usually small, and as such, annotators had no chance to do annotation training using these outside notes, and it was hard for them to get familiar with the disparate variety of document structures from other sites.

Both multi-site and single-site NLP algorithms showed a degrading trend in performance from Mayo site to other sites, albeit this issue being less prominent in the multi-site NLP algorithm as compared to the single-site NLP algorithm. The data sharing issues also impacted NLP algorithm performance. First, training sets from outside institutions were very small due to the small number of shared notes, causing difficulties in developing comprehensive rules as features, patterns, and contextual information that could appear in third party narratives could not be fully represented in such a small sample. Second, de-identification processes could cause text span issues that may impact the input text format and thus NLP algorithm performance. The algorithm performance for algorithms developed through a centralized mode was therefore not ideal for immediate use at multiple sites, as additional local fine-tuning is still needed before final implementation and application.

Our experiment results showed that a centralized approach towards multi-site NLP algorithm development is suboptimal for advancing the adoption of NLP techniques in the clinical and translational research community, this further support our proposed federated method. The

experiment also demonstrates that deployment of NLP algorithms for multi-site studies needs to be done in each local site. To ensure the scientific rigor of the data generated, each site need to perform annotation and evaluation on their own while collectively contributing to NLP algorithm development and refinement. Since the NLP models are to be shared in rule-based systems, the models can be shared without the concerns typically associated with language resources involving the Protected Health Information (PHI) issue.

In the proposed workflow, each site will evaluate the NLP algorithms for concept extraction by creating a gold standard corpus based on the common annotation guidelines. The federated evaluation can be deployed leveraging cloud computing through a centralized controller where NLP algorithms can be distributed to each institution. NLP Sandbox¹ is an example of such an evaluation framework, which uses Docker³⁹ containers to encapsulate algorithm implementations. By adopting this process, the evaluation only happens behind each institution's firewall, and only the summary statistics on NLP algorithm performance (i.e., no raw data containing PHI) is transferred out of the firewall. Performance statistics, such as the precision, recall, and F1-score, as defined depending on the experimental setting, can be obtained in near real-time and can thus be used as part of continuous development workflows.

This federated process offers several benefits. For instance, when conducting error analysis, we discovered that contexts played an important role in this case study. Error analyses showed it was not a trivial task to extract COVID signs/symptoms, as their occurrence is not necessarily isolated only to occurring due to COVID, and as they could appear as adverse events/indication of treatment, or in instruction/patient education, or clinical goal/precaution, etc. This posed a challenge not only for annotation, but also for the NLP algorithm development. One benefit of the

¹ NLP Sandbox: <https://github.com/nlpsandbox>

federated annotation and development process is that these contexts can be systematically incorporated by local expertise in the annotation process.

Deployment of a federated development framework requires the participation of multiple sites. Adoption can, however, be hindered by the fact that the process of translating NLP algorithms into implementation is complex, much like the “bench to bedside” process that translates laboratory discoveries into patient care. To facilitate participation in our federated method, we have developed a further suite of tools such as MedTator⁴⁰ and best practice guidelines⁴¹. MedTator, a serverless annotation tool, aims to provide an intuitive and interactive user interface for high-quality annotation corpus generation. The best practice guideline contains detailed instructions for facilitating multisite annotation practice with the following key activities: task formulation, cohort screening, annotation guideline development, annotation training, annotation production, and adjudication.

Simply having the toolsets be available, is, however, insufficient. Pragmatically, we have seen that there is a hyper focus on novel methods in academia with competing as opposed to collaborative priorities in NLP algorithm development. Our experience suggests that a collaborative development process for NLP algorithms is needed for truly implementable and useful multi-site NLP solutions. This is one of the key goals we seek to achieve with the Open Health Natural Language Processing (OHNLP) Collaboratory and have thus positioned our framework’s workflow to facilitate this task. Additionally, we recognize that it is not simply a software problem, a local workforce is also needed at each institution. As a consequence of conducting coordinated development of NLP algorithms deployed using our framework as a solution for consortia-specific tasks such as with the N3C, we simultaneously build the human workforce locally at institutions necessary to conduct the federated development, evaluation, and implementation of NLP algorithms using our framework.

Methods

Design Principles

Incorporating standards and interoperability. A common barrier to the widespread adoption of NLP in clinical research is the need to transform input and outputs to conform to part of an overall pipeline. While seemingly straightforward, such a task is difficult without prior significant investment in associated infrastructure and dedicated software development. It is therefore desirable to leverage existing infrastructure where possible and incorporate such an effort into the distributed NLP pipeline to reduce technical burden on the end user.

There is, however, significant variation in terms of available infrastructure and data availability amongst different institutions. Creating a solution that is immediately suitable for all these environments out of the box would be immensely challenging. For that reason, we sought to leverage existing data modeling efforts that are likely to be already adopted by academic medical institutions to standardize the data ingestion and output process. In our implementation, we chose the Observational Health Data Sciences and Informatics' Observational Medical Outcomes Partnership common data model (OHDSI/OMOP CDM) to handle input of clinical narratives via the NOTE table and output via the NOTE_NLP table. This brings the advantage that input/output is now standardized: so long as institutions have already transformed their clinical data into the OMOP CDM, and/or their downstream NLP-reliant applications read from the OMOP CDM database, no additional technical development burden is needed.

It is important to note that standardization as a default only serves to simplify adoption for those who already have a solution complying with the standard and cannot be a comprehensive solution. A purely OMOP CDM reliant solution is not ideal, as not all institutions will have their own OMOP CDM instance and standing up such an instance to just use a pipeline may produce undue burden. For that reason, input/output in our infrastructure is modularized, and can be substituted at will: the default OMOP CDM I/O utilizes a variant of SQL-based data extractors/writers, and the

specific query and connection strings used can be substituted via plaintext configuration changes. Additionally, SQL-based I/O is not the only supported setting, a variety of other data sources including Elasticsearch, google cloud storage, amazon s3, and plaintext are included as well as configuration-swappable options.

Crowdsourcing algorithm development. To promote collaboration and sharing efforts between participants in the algorithm development process, we built a crowdsourcing platform for domain experts to upload, customize, and examine their NLP algorithms in an interactive web application. Users can create keyword-based and rule-based algorithms and test the performance in the online environment instantly. The crowdsourcing platform consists of three modules based on our NLP system to support expert collaboration, including dictionary builder, regular expression rule set editor, and detection result visualization.

The dictionary builder can extend the keyword collection used by the algorithm. Users can customize particular terms from the ontology database such as CIDO ⁴² and MONDO ⁴³. The regular expression rule set editor provides an integrated interface to help users customize their own regular expression rule set (on top of an existing dictionary, if desired), to support use cases such as extraction of new symptoms, treatments, or outcomes. The detection result visualization is designed based on Brat annotation tool ⁴⁴ to check the results generated by different methods.

Case study

The National COVID Cohort Collaborative³⁶ (N3C) is a novel partnership that includes the Clinical and Translational Science Awards (CTSA) Program hubs, the National Center for Advancing Translational Science (NCATS), the Center for Data to Health (CD2H) and the community, focusing on collaborative sharing of structured EHR data. Access to unstructured data is limited due to protection of PHI and clinical care decision logics, that were further contributing to NLP infrastructure lacking within the consortia. However, structured data does not show the whole picture from the EHR perspective, greatly restricting research activities. In this case study,

extraction of COVID-19 signs and symptoms was used as a case study to investigate the viability of the proposed framework among sites participating in the N3C.

Centralizing gold standard corpus development. Due to resources and time constraints at each of the N3C sites, we opted to conduct the gold standard corpus development process in a centralized manner. A collection of de-identified and synthesized clinical documents was gathered from participating sites through an existing de-identification effort led by the NCATS Clinical Data to Health (CD2H). The N3C deidentification and synthetic text generation workflow is illustrated in Figure 1. Specifically, clinical notes from patients with positive COVID-19 test results from three institutions, Mayo Clinic, the University of Kentucky (UKen), and the University of Minnesota at Twin Cities (UMN) were initially collected. Notes that were not office visit notes (e.g., nurse calls, etc.), notes that had fewer than 1000 characters, and notes that were authored more than 14 days prior to the date of the patient's earliest positive COVID-19 test result were further filtered out. A total of 369 clinical notes from these sites that met these criteria were randomly selected, de-identified using the de-identification program developed by the Medical College of Wisconsin followed by manual review. The removed PHI identifiers are replaced by the programmatically added synthetic texts. We collected 20 signs and symptoms of COVID-19 as a basic COVID-19 concept set according to the recommendations from the CDC and Mayo Clinic. Five out of the 20 concepts are emergency warning signs including dyspnea, chest pain, delirium, hypersomnia and cyanosis. We then gathered formal definitions of each clinical concept from the Coronavirus Infectious Disease Ontology (CIDO) ⁴². Based on the Open Biological and Biomedical Ontology (OBO) Foundry library, CIDO concepts were imported from 45 ontologies, and it uses Human Phenotype Ontology (HPO) ⁴⁵ for phenotypes. Some representative phenotypes shown in COVID-19 have been imported to the CIDO. However, if the chosen COVID-19 clinical concepts were not collected by the CIDO, we re-pulled them from the HPO to the CIDO. We also gathered cross-reference concept codes from the CIDO including UMLS ⁴⁶, SNOMED-CT ⁴⁷, MeSH ⁴⁸, HPO, MeDDRA ⁴⁹.

We selected available clinical notes from both inpatients and outpatients in the two-week window preceding the order date of the first positive COVID-19 result as the annotation cohort. After the text data was collected from participating sites, the same annotation process was completed by the annotator team from Mayo Clinic to generate the gold standard annotations on COVID-19 signs and symptoms. There are 313 clinical notes from Mayo Clinic, 20 notes from UKen and 36 notes from UMN. Annotators were first trained using Mayo notes to gain better understanding of the annotation guidelines. Inter-annotator agreement (IAA) was calculated after annotation and corresponding discrepancies were resolved by discussions between the two annotators to generate a final gold standard dataset.

NLP algorithm development and evaluation. Using the annotated corpus, we developed both a single-site and multi-site NLP algorithm using a regular expression-based matching method, which has been widely adopted for information extraction in clinical settings. Specifically, for the Mayo data, we randomly chose 101 notes out of the 313 annotated notes as development set, 105 notes as validation set, and the remaining 107 notes were used as the testing set. For the UKen data, 10 notes were used for training and 10 for testing. For the UMN data, 18 was used for training and 18 for testing. Single-site algorithm was developed using the development set and validation set from Mayo, tested on the Mayo testing set and all data from UKen and UMN. Multi-site algorithm was generated through further refinement of the single-site algorithm using training sets from UKen and UMN and then tested on testing sets from all sites.

We evaluated the performance of single-site and multi-site algorithms using precision, recall, and F1-score for the annotated concept mentions, without and with certainty. A span can be represented from the start position to the end position of the concept mention. Certainty is an attribute of the concept mention including positive, negated, hypothetical and possible. For the mention-level evaluation without certainty, when there are overlaps between the gold standard mention span and the NLP detected mention span while the concept type (i.e., the specific sign/symptom such as fever, cough) is the same, it is considered a true positive (TP). If a concept mention exists in the

gold standard annotation but not detected by the NLP algorithm, or spans overlap but the concept type is not matched, it is considered as a false negative (FN). If a concept mention is detected by the algorithm but does not exist in the gold standard annotation, the concept is considered as a false positive (FP). For the mention-level span and certainty evaluation, certainty match needs to be considered when calculating TP, FN and FP. The precision, recall and F1-score are then calculated as follows. We further manually analyzed errors from multi-site algorithm mention-level evaluation without certainty.

$$Precision = \frac{TP}{TP + FN}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$

Acknowledgment

This research was possible because of the patients whose information is included within the data and the organizations and scientists who have contributed to the on-going development of this community resource <https://doi.org/10.1093/jamia/ocaa196>. The analyses described in this publication were conducted with data or tools accessed through the NCATS N3C Data Enclave <https://covid.cd2h.org> and N3C Attribution & Publication Policy v1.2-2020-08-25b and supported by NCATS U24 TR002306. This task was made possible by the National Center for Advancing Translational Sciences of the National Institutes of Health under award number U01TR02062 and the Bill & Melinda Gates Foundation. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

Code Availability Statement

Framework components can be found on GitHub:

- ETL Pipeline: <https://github.com/OHNLP/Backbone>
- NLP Implementation: <https://github.com/OHNLP/MedTagger>
- Web Rule Editor Front-end: <https://github.com/OHNLP/OHNLPTK>
- MedTator annotation tool: <https://github.com/OHNLP/MedTator>

The developed NLP ruleset can be found at

<https://github.com/OHNLP/covid19ruleset/tree/main/covid19>

Data Availability Statement

- A detailed annotation guideline outlining the goals of the NLP task and how the corpora were annotated can be found at <https://github.com/OHNLP/N3C-NLP-Documentation/wiki/Annotation-guideline-for-COVID-19-concepts>
- The sample deidentified synthetic corpus used as part of this study can be found at https://github.com/OHNLP/N3C-NLP-Documentation/blob/master/n3c_omop_sample.csv

Competing Interests Declaration

MAH has a founding interest in Pryzm Health. HX and The University of Texas Health Science Center at Houston have financial related interests at Melax Technologies Inc.

Author contributions

Project Conceptualization: SL, AWen, LW, HH, SF, HL

Data curation: SL, AWen, LW, HH, RM, AWilliams, DH, RK, ML, NA, MR, RZ, JDO, JHS

Data integration: AWen, LW, HH, RM, DH, RK, ML, NA, RZ, TS, YH, EP, SSH, CGC, JHS

Data analysis: SL, AWen, LW, RM, RK, NA, RZ, MR, TS

Software development: SL, AWen, HH, MR, TS

Data quality assurance: SL, RM, DH, RK, LM, NA, RZ, TS, JDO, HY, EP, TD, PS, HX

Draft the manuscript: SL, AWen, LW, HH, RM, RZ, HL, SF

Critical revision of the manuscript for important intellectual content: AWilliams, RK, ML, NA, YH

Project evaluation: LW, SF, RM, AWilliams, DH, RK, ML, NA, RZ, MR, TS

Project management: SL, LW, RZ, TS, EP, JHS, RF, HL

Regulatory oversight / admin: EP, CGC, HL

Database / Information systems admin: RM, DH, NA, RZ, TD

Biological subject matter expertise: LW, YH, UT, MAH

Funding acquisition: MAH, CGC, PS, HX, HL

Reference

1. Health NIo. Clinical and Translational Science Awards (CTSA) Program. <https://ncats.nih.gov/ctsa>.
2. Welcome | Center for Data to Health. <https://cd2h.org/mission>.
3. Gottesman O, Kuivaniemi H, Tromp G, Faucett WA, Li R, Manolio TA, Sanderson SC, Kannry J, Zinberg R, Basford MA, Brilliant M, Carey DJ, Chisholm RL, Chute CG, Connolly JJ, Crosslin D, Denny JC, Gallego CJ, Haines JL, Hakonarson H, Harley J, Jarvik GP, Kohane I, Kullo IJ, Larson EB, McCarty C, Ritchie MD, Roden DM, Smith ME, Böttinger EP, Williams MS, and The e MN. The Electronic Medical Records and Genomics (eMERGE) Network: past, present, and future. *Genetics in Medicine* **15**, 761-771 (2013).
4. Selby JV, Beal AC, Frank L. The Patient-Centered Outcomes Research Institute (PCORI) national priorities for research and initial research agenda. *Jama* **307**, 1583-1584 (2012).
5. Investigators AoURP. The “All of Us” research program. *New England Journal of Medicine* **381**, 668-676 (2019).
6. Hripcsak G, Duke JD, Shah NH, Reich CG, Huser V, Schuemie MJ, Suchard MA, Park RW, Wong ICK, Rijnbeek PR. Observational Health Data Sciences and Informatics (OHDSI): opportunities for observational researchers. *Studies in health technology and informatics* **216**, 574 (2015).
7. Hripcsak G, Ryan PB, Duke JD, Shah NH, Park RW, Huser V, Suchard MA, Schuemie MJ, DeFalco FJ, Perotte A. Characterizing treatment pathways at scale using the OHDSI network. *Proceedings of the National Academy of Sciences* **113**, 7329-7336 (2016).

8. Hripcsak G, Schuemie MJ, Madigan D, Ryan PB, Suchard MA. Drawing reproducible conclusions from observational clinical data with OHDSI. *Yearbook of medical informatics* **30**, 283-289 (2021).
9. Hripcsak G, Wilcox A. Reference standards, judges, and comparison subjects: roles for experts in evaluating system performance. *J Am Med Inform Assoc* **9**, 1-15 (2002).
10. Rosenbloom ST, Denny JC, Xu H, Lorenzi N, Stead WW, Johnson KB. Data from clinical notes: a perspective on the tension between structure and flexible documentation. *Journal of the American Medical Informatics Association* **18**, 181-186 (2011).
11. Cannon J, Lucci S. Transcription and EHRs. Benefits of a blended approach. *Journal of AHIMA/American Health Information Management Association* **81**, 36 (2010).
12. Blease C, Kaptchuk TJ, Bernstein MH, Mandl KD, Halamka JD, DesRoches CM. Artificial Intelligence and the Future of Primary Care: Exploratory Qualitative Study of UK General Practitioners' Views. *J Med Internet Res* **21**, e12802 (2019).
13. Fu S, Chen D, He H, Liu S, Moon S, Peterson KJ, Shen F, Wang L, Wang Y, Wen A. Clinical concept extraction: a methodology review. *Journal of Biomedical Informatics*, 103526 (2020).
14. Rongali S, Soldaini L, Monti E, Hamza W. Don't parse, generate! a sequence to sequence architecture for task-oriented semantic parsing. In: *Proceedings of The Web Conference 2020* (2020).
15. Goleva SB, Lake AM, Torstenson ES, Haas KF, Davis LK. Epidemiology of functional seizures among adults treated at a university hospital. *JAMA network open* **3**, e2027920-e2027920 (2020).

16. Dowell A, Darlow B, Macrae J, Stubbe M, Turner N, McBain L. Childhood respiratory illness presentation and service utilisation in primary care: a six-year cohort study in Wellington, New Zealand, using natural language processing (NLP) software. *BMJ open* **7**, e017146 (2017).
17. Nunes AP, Seeger JD, Stewart A, Gupta A, McGraw T. Cardiovascular Outcome Risks in Patients With Erectile Dysfunction Co-Prescribed a Phosphodiesterase Type 5 Inhibitor (PDE5i) and a Nitrate: A Retrospective Observational Study Using Electronic Health Record Data in the United States. *The Journal of Sexual Medicine* **18**, 1511-1523 (2021).
18. Chan L, Beers K, Yau AA, Chauhan K, Duffy Á, Chaudhary K, Debnath N, Saha A, Pattharanitima P, Cho J, Kotanko P, Federman A, Coca SG, Van Vleck T, Nadkarni GN. Natural language processing of electronic health records is superior to billing codes to identify symptom burden in hemodialysis patients. *Kidney Int* **97**, 383-392 (2020).
19. Hylan TR, Von Korff M, Saunders K, Masters E, Palmer RE, Carrell D, Cronkite D, Mardekian J, Gross D. Automated prediction of risk for problem opioid use in a primary care setting. *The Journal of Pain* **16**, 380-387 (2015).
20. Blumenthal KG, Lai KH, Huang M, Wallace ZS, Wickner PG, Zhou L. Adverse and hypersensitivity reactions to prescription nonsteroidal anti-inflammatory agents in a large health care system. *The Journal of Allergy and Clinical Immunology: In Practice* **5**, 737-743. e733 (2017).
21. Haug CJ. From patient to patient—sharing the data from clinical trials. *New England Journal of Medicine* **374**, 2409-2411 (2016).
22. Kent DM, Leung LY, Zhou Y, Luetmer PH, Kallmes DF, Nelson J, Fu S, Zheng C, Liu H, Chen W. Association of silent cerebrovascular disease identified using natural language processing and future ischemic stroke. *Neurology* **97**, e1313-e1321 (2021).
23. Goodlett D, Hung A, Feriozzi A, Lu H, Bekelman JE, Mullins CD. Site engagement for multi-site clinical trials. *Contemporary Clinical Trials Communications* **19**, 100608 (2020).

24. McGraw D, Leiter AB. Pathways to success for multi-site clinical data research. *Pathways* **9**, 19-2013 (2013).
25. Glynn EF, Hoffman MA. Heterogeneity introduced by EHR system implementation in a de-identified data resource from 100 non-affiliated organizations. *JAMIA open* **2**, 554-561 (2019).
26. Fu S, Leung LY, Raulli A-O, Kallmes DF, Kinsman KA, Nelson KB, Clark MS, Luetmer PH, Kingsbury PR, Kent DM. Assessment of the impact of EHR heterogeneity for clinical research through a case study of silent brain infarction. *BMC medical informatics and decision making* **20**, 1-12 (2020).
27. Morris AH, Orme Jr J, Truwit JD, Steingrub J, Grissom C, Lee KH, Li GL, Thompson BT, Brower R, Tidswell M. A replicable method for blood glucose control in critically ill patients. *Critical care medicine* **36**, 1787-1795 (2008).
28. Chiang MF, Hwang JC, Alexander CY, Casper DS, Cimino JJ, Starren J. Reliability of SNOMED-CT coding by three physicians using two terminology browsers. In: *AMIA Annual Symposium Proceedings*. American Medical Informatics Association (2006).
29. Leung LY, Fu S, Luetmer PH, Kallmes DF, Madan N, Weinstein G, Lehman VT, Rydberg CH, Nelson J, Liu H. Agreement between neuroimages and reports for natural language processing-based detection of silent brain infarcts and white matter disease. *BMC neurology* **21**, 1-5 (2021).
30. Kennedy PJ, Leathley CM, Hughes CF. Clinical practice variation. *Medical Journal of Australia* **193**, S97-S99 (2010).
31. Sohn S, Wang Y, Wi C-I, Krusemark EA, Ryu E, Ali MH, Juhn YJ, Liu H. Clinical documentation variations and NLP system portability: a case study in asthma birth cohorts across institutions. *Journal of the American Medical Informatics Association* **25**, 353-359 (2018).

32. Liu H, Bielinski SJ, Sohn S, Murphy S, Waghlikar KB, Jonnalagadda SR, Ravikumar K, Wu ST, Kullo IJ, Chute CG. An information extraction framework for cohort identification using electronic health records. *AMIA Summits on Translational Science Proceedings* **2013**, 149 (2013).
33. Wen A, Fu S, Moon S, El Wazir M, Rosenbaum A, Kaggal VC, Liu S, Sohn S, Liu H, Fan J. Desiderata for delivering NLP to accelerate healthcare AI advancement and a Mayo Clinic NLP-as-a-service implementation. *NPJ digital medicine* **2**, 1-7 (2019).
34. Rando HM, Bennett TD, Byrd JB, Bramante C, Callahan TJ, Chute CG, Davis HE, Deer R, Gagnier J, Korashy FM, Liu F, McMurry JA, Moffitt RA, Pfaff ER, Reese JT, Relevo R, Robinson PN, Saltz JH, Solomonides A, Sule A, Topaloglu U, Haendel MA. Challenges in defining Long COVID: Striking differences across literature, Electronic Health Records, and patient-reported information. *medRxiv*, 2021.2003.2020.21253896 (2021).
35. Sharafeldin N, Bates B, Song Q, Madhira V, Yan Y, Dong S, Lee E, Kuhrt N, Shao YR, Liu F, Bergquist T, Guinney J, Su J, Topaloglu U. Outcomes of COVID-19 in Patients With Cancer: Report From the National COVID Cohort Collaborative (N3C). *Journal of Clinical Oncology*, JCO.21.01074 (2021).
36. Haendel MA, Chute CG, Bennett TD, Eichmann DA, Guinney J, Kibbe WA, Payne PRO, Pfaff ER, Robinson PN, Saltz JH, Spratt H, Suver C, Wilbanks J, Wilcox AB, Williams AE, Wu C, Blacketer C, Bradford RL, Cimino JJ, Clark M, Colmenares EW, Francis PA, Gabriel D, Graves A, Hemadri R, Hong SS, Hripscak G, Jiao D, Klann JG, Kostka K, Lee AM, Lehmann HP, Lingrey L, Miller RT, Morris M, Murphy SN, Natarajan K, Palchuk MB, Sheikh U, Solbrig H, Visweswaran S, Walden A, Walters KM, Weber GM, Zhang XT, Zhu RL, Amor B, Girvin AT, Manna A, Qureshi N, Kurilla MG, Michael SG, Portilla LM, Rutter JL, Austin CP, Gersing KR, the NCC. The National COVID Cohort Collaborative (N3C): Rationale, design, infrastructure, and deployment. *Journal of the American Medical Informatics Association* **28**, 427-443 (2021).
37. Github - OHNLP/MedTagger. <https://github.com/OHNLP/MedTagger>.

38. Harkema H, Dowling JN, Thornblade T, Chapman WW. ConText: an algorithm for determining negation, experiencer, and temporal status from clinical reports. *Journal of biomedical informatics* **42**, 839-851 (2009).
39. Merkel D. Docker: lightweight Linux containers for consistent development and deployment. *Linux J* **2014**, Article 2 (2014).
40. He H, Fu S, Wang L, Liu S, Wen A, Liu H. MedTator: a serverless annotation tool for corpus development. *Bioinformatics*, (2022).
41. Liu S, Fu S, Liu H. Best practices of annotating clinical texts for information extraction tasks. In: *Informatics Playbook* (ed Wu C) (2021).
42. He Y, Yu H, Ong E, Wang Y, Liu Y, Huffman A, Huang H-h, Beverley J, Hur J, Yang X, Chen L, Omenn GS, Athey B, Smith B. CIDO, a community-based ontology for coronavirus disease knowledge and data integration, sharing, and analysis. *Scientific Data* **7**, 181 (2020).
43. Mungall CJ, McMurry JA, Köhler S, Balhoff JP, Borromeo C, Brush M, Carbon S, Conlin T, Dunn N, Engelstad M, Foster E, Gourdine JP, Jacobsen JOB, Keith D, Laraway B, Lewis SE, NguyenXuan J, Shefchek K, Vasilevsky N, Yuan Z, Washington N, Hochheiser H, Groza T, Smedley D, Robinson PN, Haendel MA. The Monarch Initiative: an integrative data and analytic platform connecting phenotypes to genotypes across species. *Nucleic Acids Research* **45**, D712-D722 (2017).
44. Stenetorp P, Pyysalo S, Topić G, Ohta T, Ananiadou S, Tsujii Ji. brat: a Web-based Tool for NLP-Assisted Text Annotation.). Association for Computational Linguistics (2012).
45. Köhler S, Gargano M, Matentzoglou N, Carmody LC, Lewis-Smith D, Vasilevsky NA, Danis D, Balagura G, Baynam G, Brower AM, Callahan TJ, Chute CG, Est JL, Galer PD, Ganesan S, Griese M, Haimel M, Pazmandi J, Hanauer M, Harris NL, Hartnett MJ, Hastreiter M, Hauck F, He Y, Jeske T, Kearney H, Kindle G, Klein C, Knoflach K, Krause R, Lagorce D, McMurry JA, Miller

JA, Munoz-Torres MC, Peters RL, Rapp CK, Rath AM, Rind SA, Rosenberg AZ, Segal MM, Seidel MG, Smedley D, Talmy T, Thomas Y, Wiafe SA, Xian J, Yüksel Z, Helbig I, Mungall CJ, Haendel MA, Robinson PN. The Human Phenotype Ontology in 2021. *Nucleic Acids Res* **49**, D1207-d1217 (2021).

46. Bodenreider O. The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Research* **32**, 267D-270 (2004).
47. SNOMED-CT. <https://www.nlm.nih.gov/healthit/snomedct/index.html>.
48. Medicine NLo. Medical Subject Headings. <https://www.nlm.nih.gov/mesh/meshhome.html>.
49. Brown EG, Wood L, Wood S. The medical dictionary for regulatory activities (MedDRA). *Drug Saf* **20**, 109-117 (1999).

Figures

Figure 1. N3C deidentification and synthetic text generation workflow

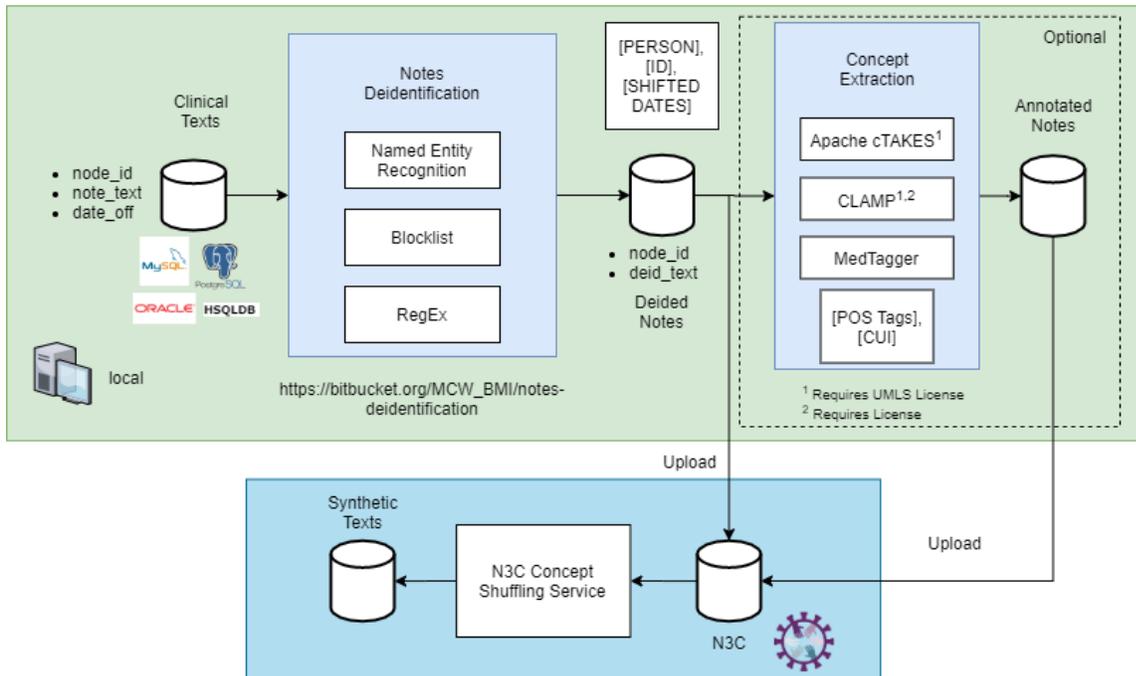


Figure 2. Screenshots of the Web GUI

(a)

The screenshot displays the N3C NLP Engine Demo interface. At the top, a dark navigation bar contains the site name and links for COVID-19 Demo, Dictionary Builder, Rule Editor, Wiki, user profile (guest), and GitHub. The main content area is divided into several sections:

- Input Text:** A text area with a 3,000-character limit containing the sentence: "The patient had a dry cough and fever or chills yesterday. He is also experiencing new loss of taste today and three days ago." Below this is a "Document Date" field set to "2021-02-18" and a "Run MedTagger" button.
- Visualization:** A tabbed interface with "Visualization" selected and "Raw Output" as an alternative. The "Brat visualization" shows the input text with green bars above and below it, indicating identified entities and their spans.
- Concept/Term List:** A list of medical terms categorized into two columns:
 - Dry cough
 - Fever
 - Lymphopenia
 - Sore Throat
 - Ground Glass Infiltrates
 - Elevated LDH
 - Diarrhea
 - Nasal Congestion
 - Loss of Appetite
 - Fatigue
 - Dyspnea
 - Headache
 - Myalgia
 - Abdominal Pain
 - Patchy Infiltrates
 - Elevated PT Time
 - InfluenzaA "Ruleset »" button is located at the bottom left of this section.
- COVID-19 Severe Case:** A section with the heading "COVID-19 Severe Case" and a sub-heading "To identify people at higher risk for severe illness using structured and unstructured data according to the CDC guideline." A "Wiki »" button is positioned below the text.

(c)

N3C NLP Engine Demo COVID-19 Demo Dictionary Builder Rule Editor Wiki guest GitHub

Dictionary Builder

Main

Ontology CIDO Load Data Extract Selection Upload to Server

Ontology Database Build Dictionary Deployment

Ontology Tree

The tree is built on 6818 classes

- All
- sequence_feature
- entity
 - continuant
 - occurrent
 - process
 - host response to coronavirus infection
 - disorder prevention
 - infectious disease epidemic
 - vaccine-induced host response
 - exposure event or process
 - Physiological Effects [PE]
 - life cycle
 - host exposure to infectious agent
 - coronaviral process to host
 - life-death temporal boundary
 - Cellular or Molecular Interactions [MoA]
 - immunization against infectious agent
 - immune response
 - host-coronavirus interaction
 - immunization
 - life cycle stage
 - Clinical Kinetics [PK]
 - disease course

Term List

continuant[An entity that exists in full at any time in which it exists at all, persists through time while maintaining its identity and has no temporal parts.]An entity that exists in full at any time in which it exists at all, persists through time while maintaining its identity and has no temporal parts.[continuant]OBJC

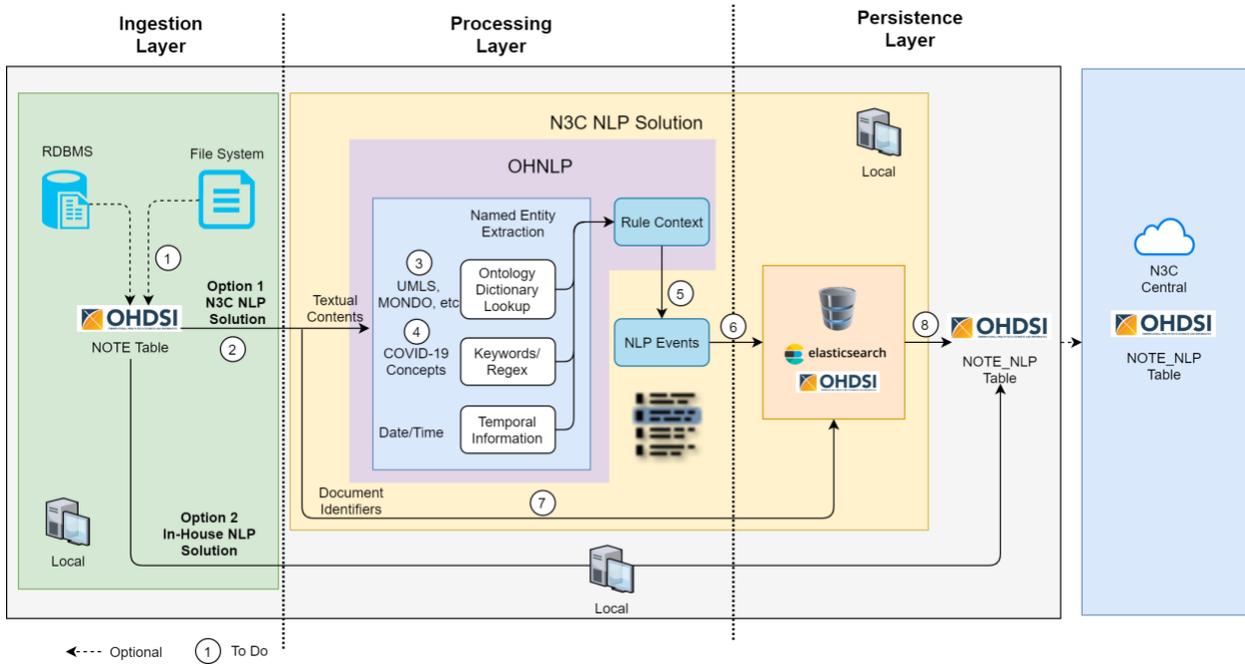
occurent[An entity that has temporal parts and that happens, unfolds or develops through time.]occurent[OBJC process]p is a process = Def. p is an occurrent that has temporal proper parts and for some time t, p s-depends_on some material entity at t. (axiom label in BFO2 Reference: [083-003])An occurrent that has temporal proper parts and for some time t, p s-depends_on some material entity at t.[process]OBJC

host response to coronavirus infection[A coronavirus-host interaction by which the host responds to viral infection.]host response to coronavirus infection[OBJC

disorder prevention[disorder prevention is a processual entity that prevents a disorder that is the physical basis of a disease.]disorder prevention[OBJC

infectious disease epidemic[A process of infectious disease realizations and for which there is a statistically significant increase in the infectious disease incidence of a population.]infectious disease epidemic[OBJC

Figure 3. Diagram of N3C NLP Solution



Tables

Table 1. Annotation corpora statistics. (Mayo: Mayo Clinic, UKen: University of Kentucky, UMN: University of Minnesota)

Concepts	Mayo (313 notes)	UKen (20 notes)	UMN (36 notes)
Abdominal_pain	59	2	3
Chest_pain	62	2	11
Chill	51	6	6
Cough	104	14	43
Cyanosis	9	4	17
Delirium	38	2	10
Diarrhea	92	5	11
Dyspnea	199	19	46
Fatigue	61	15	13
Fever	148	25	53
Headache	43	6	15
Hypersomnia	6	0	14
Loss_of_appetite	41	2	4
Loss_of_smell	23	4	6
Loss_of_taste	19	2	5
Myalgia	21	6	8
Nasal_obstruction	16	6	14
Nausea	87	7	10
Sore_throat	16	4	17
Vomiting	86	6	14

Table 2. Performance of the single-site NLP algorithm (Mayo: Mayo Clinic, UKen: University of Kentucky, UMN: University of Minnesota)

Dataset	Span			Span+Certainty		
	Precision	Recall	F1	Precision	Recall	F1
Mayo	0.882	0.869	0.876	0.789	0.639	0.706
UKen	0.698	0.714	0.706	0.664	0.643	0.653
UMN	0.658	0.735	0.694	0.534	0.438	0.481

Table 3. Performance of the multi-site NLP algorithm (Mayo: Mayo Clinic, UKentucky: University of Kentucky, UMN: University of Minnesota)

Dataset	Span			Span+Certainty		
	Precision	Recall	F1	Precision	Recall	F1
Mayo	0.863	0.908	0.884	0.824	0.681	0.746
UKen	0.696	0.859	0.769	0.662	0.734	0.696
UMN	0.718	0.918	0.806	0.562	0.456	0.504

Table 4. Error analysis of the multi-site algorithm mention-level evaluation without certainty for Mayo site

Error types of FP	No. FP (%)	Error types of FN	No. FN (%)
Annotation error: missing annotation	17 (26%)	NLP algorithm not complete	21 (66%)
Hard to judge if are COVID signs/symptoms	15 (23%)	Annotation error	8 (25%)
Not COVID signs/symptoms - instruction/patient education	10 (15%)	Tokenization error due to input format	2 (6%)
Not COVID signs/symptoms - adverse events/indication of treatment	7 (11%)	Template	1 (3%)
Not COVID signs/symptoms - others (anesthesia plan, symptoms of other disease, etc.)	5 (7%)		
Not COVID signs/symptoms - clinical goal/precaution	4 (6%)		
Not COVID signs/symptoms - template	5 (7%)		
NLP algorithm not precise	2 (3%)		

Table 5. Error analysis of the multi-site algorithm mention-level evaluation without certainty for UMN

Error types of FP	No. FP (%)	Error types of FN	No. FN (%)
Not COVID signs/symptoms - instruction/patient education	30 (61%)	NLP algorithm not complete	11 (85%)
Not COVID signs/symptoms - (substance use history; overdose; due to surgery, other comorbidity, etc.)	7 (14%)	Annotation error	2 (15%)
Annotation error: missing annotation	6 (12%)		
Not COVID signs/symptoms - template	4 (8%)		
Hard to judge if are COVID signs/symptoms	2 (3%)		

Table 6. Error analysis of the multi-site algorithm mention-level evaluation without certainty for UKen

Error types of PF	No. FP (%)	Error types of FN	No. FN (%)
Not COVID signs/symptoms - (substance use history; overdose; due to surgery, other comorbidity, etc.)	8 (33%)	NLP algorithm not complete	7 (78%)
Annotation error: missing annotation	5 (21%)	Annotation error	2 (22%)
Not COVID signs/symptoms - instruction/patient education	5 (21%)		
Hard to judge if are COVID signs/symptoms	3 (13%)		
Not COVID signs/symptoms - template	2 (8%)		
NLP algorithm not precise	1 (4%)		