# **Extracting OMOP Concepts from Clinical Narratives to Empower Clinical Research**

## January 25, 2022





## Outline

- N3C Overview
- Open Health Natural Language Processing (OHNLP)
  - Ethical AI framework
  - Scientific rigor
  - OHNLP Toolkits
  - PASC NLP algorithms
- Join the journey for NLP-empowered RECOVER
  - Contribute to the development effort
  - Implement the N3C NLP Infrastructure
  - Augment data with NLP for N3C
  - Collaboration and Partnership



## N3C Overview

#### **Chris Chute**

## N3C Enclave: The largest public HIPAA- limited longitudinal-EHR data set in US history



National COVID Cohort Collaborative

(1/20/22 release)





## N3C EHR Data Overview





#### Row level data:

- Person
- Drug
- Procedure
- Condition (diagnoses)
- Measurement (labs)

### Types of encounters:

- Inpatient, ED & Outpatient
- Longitudinal back to 2018



National COVID Cohort

## Each of the 69 sites has a pipeline with 2M+ transformations







The provenance between 5000 syntactic transformations across the 69 sites is automatically tracked.

#### This enables:

- pipeline developers to very quickly identify the root cause of data quality issues
- data pipelines can be refreshed in <20 minutes whenever the source data updates



#### National COVID Cohort Collaborative

## **Open Health Natural Language Processing**

## Hongfang Liu



## OHNLP - 2009 (Mayo and IBM)

## 03 Apr Mayo Clinic and IBM to move beyond EMR's to deliver knowledge at the point of care

Posted at 07:26h in Health Technology, IT by WTN News • 0 Comments

## 2009

https://wtnnews.com/2009/04/03/5840/ Madison – Biomedical informatics researchers at Mayo Clinic and IBM have launched a Web site for the newly founded Open Health Natural Language Processing (NLP) Consortium. The consortium is establishing the open-source space to promote past and current development efforts, including participation in information extraction from electronic medical records. Mayo Clinic and IBM Healthcare released clinical NLP technologies into the public domain. The site will allow the approximately 2,000 researchers and developers working on clinical language systems worldwide to contribute code and further develop the systems. Additionally, the VA Boston Healthcare System and Seattle Group Health have strongly indicated their support of he concept according to IBM.

"We are inviting our international colleagues to help continue development of these valuable tools," says Christopher Chute, M.D., Dr.P.H., Mayo Clinic bioinformatics expert and senior consultant on the project. "By making it an open-source initiative, we hope to enable wide use of these NLP tools so medical advancements can happen faster and more efficiently."

**cTAKES** functionality recognizes whether a clinical concept is negated, relevant to the patient or to the patient's family, which are attributes critical to understanding patient-centered medical **IBM's medKAT systems** (medical Knowledge Analysis Tool) is a system to extract structured information from unstructured data sources, such as pathology reports, clinical notes, discharge summaries and medical





## OHNLP - 2009 to 2012 (ONC SHARP)

Strategic Health IT Advanced Research Projects (SHARP) Research Focus Area 4 - Secondary Use of EHR Data Increasing efficiency of patient care through electronic healthcare records



Face to Face Group Photo June 11, 2012



## OHNLP - 2013 - 2016 (NIGMS)

- Research on implementability of open source NLP systems (<u>https://github.com/nlpie/nlp-adapt</u>) (led by Dr. Serguei Pakhomov)
- Research on usability to deliver production-ready NLP (led by Hua Xu)
- Research on NLP utility and generalizability for clinical research and practice (led by Hongfang Liu)

## OHNLP - 2017 - present







## Human-centric Ethical AI Framework



# The Use of AI for Health Guiding Principles

Human-centric Ethical AI Framework



- Protecting human autonomy
- Promoting human well-being and safety and the public interest
- Ensuring transparency, explainability and intelligibility
- Fostering responsibility and accountability
- Ensuring inclusiveness and equity
- Promoting AI that is responsive and sustainable

Opinion

## Too many AI researchers think real-world problems are not relevant

The community's hyperfocus on novel methods ignores what's really important.

by Hannah Kerner

August 18, 2020

#### IBM is selling off its Watson Health assets



By <u>Clare Duffy</u>, <u>CNN Business</u> Updated 2:29 PM ET, Fri January 21, 2022

- Language plays an unique role in humanity and society.
- Need to follow a social and ethical framework when working on health language.



## **Federated Development and Evaluation**



About 4,240,000,000 results (0.80 seconds)

The concept behind "data is the new oil" is that just like oil, raw data isn't valuable in and of itself, but, rather, the value is **created** when it is gathered completely and accurately, connected to other relevant data, and done so in a timely manner. ... COVID-19 related data is being generated quickly. Apr 27, 2020

https://www.kenwayconsulting.com > blog > data-is-the-n...

Is Data Really "The New Oil"? | Kenway Consulting



## Reproducibility and Methodology Rigor

#### **Scientific Rigor**

- How is NLP evaluation conducted in clinical research?
- What's the current reporting practice of the NLP component in the traditional clinical research?
- What is the degree of granularity when reporting NLP methodology and evaluation?
- What are the potential barriers to adopt clinical NLP to clinical research?
- What are the lessons learned from the existing evaluation and reporting practices





## Reproducibility and Methodology Rigor

**Scientific Rigor** 



We observed high heterogeneity in the reporting practice: 14% of studies did not report NLP methodology and evaluation, 22% of studies did not report evaluation design, and 10% of studies did not report NLP methodology. A few studies claimed that 'NLP has been evaluated prior to release'. However, no additional details can be found to justify the validity of the NLP results and clinical findings.



## Reproducibility and Methodology Rigor

**Scientific Rigor** 

#### **Evaluation Challenges**

- Only mention cohort duration but not NLP evaluation duration
- EHR system migration, change of definition, ETL process
- Study setting and EHR environment
- Definition variation between cohort and disease
- Commercial vs open source
- Evaluation performance/matrix

## **Implication to Clinical Research**

- Most data elements identified by NLP are either comprised of risk factors (48%) or outcome (42%), which have strong implication to the outcome of the study.
- The validity of the study is thus dependent on the rigor of NLP methodology and evaluation.



## A Collaborative Framework for NLP Development

**OHNLP** Toolkits





## A Federated Framework for NLP Evaluation

**OHNLP** Toolkits

13 Entity Tags 1 Link Tags	A_Chill_	_ <b>07.txt_hh.xml</b> 80 chars 14 tags	Save Save as	<ul> <li>Document</li> <li>I Sentences</li> </ul>	C 🖍 Rich Text	Smart O Off Accept All	Show Links	Sample About	a
Schema DTD File (.dtd)	Annota	ation File (.xml)	Save	Display Mode	Mark Mode	Hint Mode	Link Marks	Help	
20 files arge_01.bt_hh.xml a.C.chill_072.txt_hh.xml Dizziness_17.txt_hh.xml Dizziness_17.txt_hh.xml Dizziness_12.txt_hh.xml Myalgia_09.txt_hh.xml Myalgia_09.txt_hh.xml Myalgia_01.txt_hh.xml Alarge_01.txt_mm.xml Alarge_01.txt_mm.xml	Image: Constraint of the	B During th 11:00 PM the remai 2 Tried co 3 7:00AM - 4 Slept, fi 5 M Advil 6 2/21/21 -	e next 60 hour - 2 hours of f Inder of the n Advil , then headache, nally that nichelped the head Left with boo	rs I had the for incontrollable ight. In tylenol - o PAL Body Aches ght. ad and PA Body dy D Fatigue,	bllowing reacti shivering and didn't tough a , correction and , correction and , aches finally mild correction and a	ons, though not all headache 2/20/21 c chills, mild c	at the same ti 2:AM Extreme To diarrhea an	me: 02/19	<pre>//21: 10:00PM No Nausea and extreme headac (like Migraine), Nausea, PAC Body aching 1 VAX 2 PYREXIA 3 CHILL 4 COUGH 3 FATIGUE 4 PAIN 4 HEADACHE 5 SORENESS 5 SORENESS</pre>
Chill_47.txt_mm.xml Dizziness_17.txt_mm.xml Fatigue_01.txt_mm.xml Headache_02.txt_mm.xml Myalgia_09.txt_mm.xml	11 0 13 0 4 0 2 0 3 0	7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t	eats for most of energy back, good, except this whole orde	mild <b>DI</b> diarr left with <b>DI</b> eal, I <b>W2 const</b>	nild diarrhea	Leary of taking an ter , almost drinking	nything for the	<b>DI</b> diarrhe	ea; as I tend a MEDICATION
_Chill_47.txt_mm.xml Dizziness_17.txt_mm.xml Fatigue_01.txt_mm.xml Headache_02.txt_mm.xml Myalgia_09.txt_mm.xml Nausea_11.txt_mm.xml	11 0 13 0 4 0 2 0 3 0 2 0	7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t	eats for most of energy back, good, except this whole orde	of that night. mild <b>pl</b> diarr left with [pr] eal, I [c] consu	hea left. mild diarrhea umed tons of wa	Leary of taking an ter , almost drinki	nything for the ng nonstop due	to the thir	ea;, as I tend a MEDICATION
Chill_42.txt_mm.xml Dizziness_17.txt_mm.xml Fatigue_01.txt_mm.xml Myalgia_09.txt_mm.xml Nausea_11.txt_mm.xml Nausea_11.txt_mm.xml VAX	11 0 13 0 2 0 3 0 2 0 14	7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t MEDICATION	eats for most of energy back, good, except this whole order M2	of that night. mild D1 diarr left with 011 eal, I 00 const 814~836 co	mild diarrhea umed tons of water	Leary of taking an ter , almost drinking comment	nything for the	diarrhe	ea; as I tend a MEDICATION
Chill_42.txt_mm.xml Dizziness_17.txt_mm.xml Fatigue_01.txt_mm.xml Headache_02.txt_mm.xml Mayligia_03.txt_mm.xml Nausea_11.txt_mm.xml At! Tags VAX PYREXIA	11 0 13 0 4 0 2 0 3 0 2 0 14 0 1	7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t D D	ents for most of energy back, good, except this whole order M2	of that night. mild Di diarr left with Di eal, I C consu 814-836 cc 405-413 di	nild diarrhea umed tons of water	Leary of taking ar ter , almost drinkin comment	nything for the	DI diarrhe	ea, as I tend a MEDICATION
Chill_42.txt_mm.xml Dizziness_17.txt_mm.xml Fleadache_02.txt_mm.xml Headache_02.txt_mm.xml Nausea_11.txt_mm.xml Nausea_11.txt_mm.xml At! Tags VAX PYREXIA S CHILL	11 0 13 0 2 0 3 0 2 0 14 0 1	7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t DiarRHEA	entry back, good, except this whole order M2 DIO	of that night. mild diarr left with ori al, I co consi 814-836 co 405-413 di	thea left. mild diarrhea umed tons of wa onsumed tons of water	Leary of taking ar ter , almost drinkin comment comment	nything for the	<b>Di</b> diarrhe to the thir	ea, as I tend a MEDICATION
CHILL CLCHH		7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t DIARRHEA 9 DIARRHEA	eats for most of energy back, or good, except this whole order m M2 m Dio	of that night. mild 011 diarr left with 011 eal, I 92 const 814~836 cc 405-413 di 703-716 m	thea left. mild diarrhea umed tons of water arrhea wild diarrhea	Leary of taking ar	hything for the	<b>DI</b> diarrhe to the thir	ea, as I tend d MEDICATION
hill_42tst_mm.xml Dizziness_17.txt_mm.xml atigue_01.txt_mm.xml leadache_02.txt_mm.xml Ayaigia_09.txt_mm.xml Ayaigia_09.txt_mm.xml All Tags VAX PYREXIA VAX PYREXIA COUGH S FATIGUE		7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t MEDICATION DIARRHEA 9 DIARRHEA	eats for most of energy back, or good, except this whole order M2 DIO DIO DIO	of that night. mild indiarr left with [31] eal, I [32] const 814-836 cc 405-413 di 703-716 m	thea left. mild diarrhea umed tons of wa onsumed tons of water arrhea ald diarrhea	Leary of taking an ter , almost drinkin comment comment	hything for the	to the thir	as I tend a MEDICATION
Chill_42.txt_mm.xml Dizziness_17.txt_mm.xml Fatigue_01.txt_mm.xml Headache_02.txt_mm.xml Mayalgia_00.txt_mm.xml Nausea_11.txt_mm.xml Atl Tags VAX PYREXIA CHILL COUGH FATIGUE FATIGUE PATIGUE PATIGUE		7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t DIARRHEA DIARRHEA MEADACHE	eats for most of energy back, good, except this whole order to be an except of the second sec	of that night. mild of diarr left with original 814-836 cc 405-413 di 703-716 m 341-349 bb	thea left. mild diarrhea umed tons of wa onsumed tons of water iarrhea iild diarrhea eadache	Leary of taking ar	hything for the	to the thir	ea, as I tend a MEDICATION rst.
Chill_42txt_mm.xml Dizziness_17.txt_mm.xml Dizziness_17.txt_mm.xml Headache_02.txt_mm.xml Mausea_11.txt_mm.xml Atl Tags VAX  PYREXIA CUUGH CUUGH FATIGUE FATIGUE PAIN HEADACHE	111 0 13 0 2 0 3 0 2 0 14 0 1 1 0 2 2 2 2 2 2 2 2 2 2 2 2 2	7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t DIARRHEA © DIARRHEA © DIARRHEA	eats for most of energy back, good, except this whole order M2 DIO H1 H1	of that night. mild diarr left with ori al, I co consi 814-836 cc 405-413 di 703-716 m 341-349 ht	thea left. mild diarrhea umed tons of wa onsumed tons of water arrhea ild diarrhea eadache	Leary of taking ar ter , almost drinkin comment comment comment	hything for the ng nonstop due	to the thir	ea, as I tend a MEDICATION
Chill A2 txt_mm.xml Dizziness_17.txt_mm.xml Jiezidache_0.2.txt_mm.xml Myalgia_0.9.txt_mm.xml Musea_11.bt_mm.xml Mill Tags V vAx PYREXIA 3 CHILL 4 COUGH 5 FATIGUE 4 PAIN 9 HEADACHE 1 DIZTIMESS	111 0 13 0 2 0 2 0 14 0 1 1 0 0 2 2 2 2 2 2 2 2 2 2 2 2 2	7 Night swe 8 2/22/21 - 9 2/23/21 - 10 Through t DIARRHEA DIARRHEA DIARRHEA HEADACHE	eats for most of energy back, good, except this whole order M M2 M DIO M DIO M H1 M TO TO TO	of that night. mild 011 diarr left with 011 eal, I 02 const 814-836 co 405-413 di 703-716 m 341-349 ho	thea left. mild diarrhea umed tons of water arrhea sild diarrhea eadache	Leary of taking ar	hything for the ng nonstop due	DI diarrhe	ea, as I tend a MEDICATION



#### National COVID Cohort Collaborative A Minimal Viable Product (MVP) for NLP Deployment

**OHNLP** Toolkits







- Incorporate dictionary lookup for concept mentions
  - Default config includes COVID/ PASC relevant concepts
  - Dictionary resource for comprehensive NLP is available by request
    - A curated dictionary, MedLex, with OMOP Concept Identifiers mapped available
- Work on documentation for onboarding



> AMIA Annu Symp Proc. 2012;2012:568-76. Epub 2012 Nov 3.

# Towards a semantic lexicon for clinical natural language processing



Hongfang Liu<sup>1</sup>, Stephen T Wu, Dingcheng Li, Siddhartha Jonnalagadda, Sunghwan Sohn, Kavishwar Wagholikar, Peter J Haug, Stanley M Huff, Christopher G Chute

Affiliations + expand

PMID: 23304329 PMCID: PMC3540492 Paperpile
Access Options

Free PMC article

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3540492/



- 287 unique clinical findings mapped to Human Phenotype Ontology (HPO) terms with synonyms and Plain-language synonyms, (Deer, Rachel R., Madeline A. Rock, Nicole Vasilevsky, Leigh Carmody, Halie Rando, Alfred J. Anzalone, Marc D. Basson et al. "Characterizing long COVID: deep phenotype of a complex condition." *EBioMedicine* 74 (2021): 103722.)
- 355 high-level long COVID symptoms consolidated from 1520 UMLS concepts of 16,466 synonyms. (Wang, Liqin, Dinah Foer, Erin MacPhaul, Ying-Chih Lo, David W. Bates, and Li Zhou. "PASCLex: A comprehensive post-acute sequelae of COVID-19 (PASC) symptom lexicon derived from electronic health record clinical notes." *Journal of biomedical informatics* 125 (2022): 103951.)



- Map MedLex to the concept identifiers from the latest OMOP Vocabulary (through standard vocabulary)
- All terms from the two PASC articles were processed and mapped to the concept identifiers from the latest OMOP Vocabulary.
- For PASC terms could not be mapped, we manually reviewed and aligned with the closest OMOP concept identifiers if available or map to the HPO identifiers.



**Evaluation and Refinement** 

**OHNLP** Toolkits

## (a single site - preliminary result)

- Refinement with 97 long hauler notes
  - Add 23 concepts not included in the initial PASC dictionary but mappable to OMOP or HPO Vocabulary
  - Add 4 concepts not mappable to OMOP or HPO Vocabulary Ο
- The comparison with a human annotator on the PASC symptom section of long hauler sections. Among a total of 1500 annotations:
  - 1,022 (179 unique normalized text strings) were annotated by both.
  - 414 (251 unique normalized text strings) annotated by NLP but missed by Ο human. Some are human omissions. Some are contextual interpretation.
  - 62 (48 unique) annotated by human but missed by NLP. About 8 out of the Ο 48 unique annotations were not in either OMOP or the original PASC concepts. The remaining were caused by lexical variants or multi-span concept mentions (smell was complete lost).



# Join the journey for NLP-empowered RECOVER

## Hongfang Liu and Emily Plaff



## N3C Enclave NLP Data Contribution

#jointhejourney



https://github.com/National-COVID-Cohort-Collaborative/Phenotype\_Data\_Acquisition/wiki/NLP-Submission-Process



Sites who are interested in deploying the infrastructure, reach out to Katelyn Cordie

(cordie.katelyn@mayo.edu)

Sites who are interested in contributing to the NLP releases, reach out to Sijia Liu

(liu.sijia@mayo.edu)

N3C domain teams who have new NLP needs, reach out to Rafael Fuentes

(rafael.fuentes@nih.gov)

All other questions, reach out to Hongfang Liu

(liu.hongfang@mayo.edu)



## **Planned Development**

- Partnering with TriNetX
- NLP Knowledge Authoring and Concept Mapping Service empowered by Terminology Service and Large Language Models
- Incorporating Open Source NLP Algorithms available for:
  - Family History
  - Genetic Information
  - Context Classifier
  - Precision Oncology
  - SDoH Variables

#jointhejourney



## **Ongoing Community Collaboration**

#### **OHNLP** Collaboration

• Open source tools and resources

#### **Refinement and Evaluation**

• Continuous improvement

#### **Context Classification**

• Documentation heterogeneity

#### **Data Quality and Machine Learning**

• Evaluate the contribution of NLP to data quality and research



## Acknowledgement

- National Institute of Health
- The Mayo NLP Program
- The N3C NLP Community
- The OHNLP Community
- The CD2H Community
- The CTSA/iEC Community
- The OHDSI Community

Join the Journey in Translating Technology Innovation to Empower Clinical Research towards Better Health for Everyone.



## Join the journey



NATIONAL CENTER FOR DATA TO HEALTH



National COVID Cohort Collaborative

#### National Center for Advancing Translational Sciences



NATIONAL CENTER FOR DATA TO HEALTH Onboarding to N3C: <u>bit.ly/cd2h-onboarding-form</u>

**Joining Workstreams: N3C Data Ingestion & Harmonization Workstream Slack Channel Harmonization Google Group Harmonization** 

> N3C Phenotype & Data Acquisition Workstream Slack Channel Phenotype Google Group Phenotype

**N3C Collaborative Analytics Workstream Slack Channel Analytics Google Group Analytics** 

N3C Data Partnership & Governance Workstream Slack Channel Governance **Google Group Governance** 

N3C Synthetic Clinical Data Workstream Slack Channel Synthetic **Google Group Synthetic** 

Additional Information:

Onboarding N3C, Slack, Google | Finding and Joining a Google Group