

Spark NLP for Healthcare: Modular Approach to Solve Problems at Scale in Healthcare NLP

April 20, 2021

Colab notebook : https://bit.ly/spark_nlp_colab

License keys : https://bit.ly/spark_nlp_license

Slides : https://bit.ly/spark_nlp_slides

ODSC EAST

BOSTON | APRIL 19-21

Veysel Kocaman

Lead Data Scientist
John Snow Labs



The Challenges we face today;



Healthcare AI practitioners usually need to make unacceptable trade-offs between

- delivering state-of-the-art accuracy,
- generalizing over unseen data points,
- preventing the sharing of personal data or intellectual property.

NLP in Healthcare

Clean & structured data



Raw & unstructured data



Healthcare data



- Less than **50% of the structured data** and less than **1% of the unstructured data** is being leveraged for decision making in companies (HBR). This is even worse in healthcare.
- NLP is ultra domain specific, so train your own models.

Why is language understanding hard?

Human Language is:

- Nuanced
- Fuzzy
- Contextual
- Medium specific
- Domain specific

Healthcare specific needs:

1. Core Annotators

Part of speech, spell checking, ...

2. Vocabulary

Ontologies, relationships, word embeddings, ...

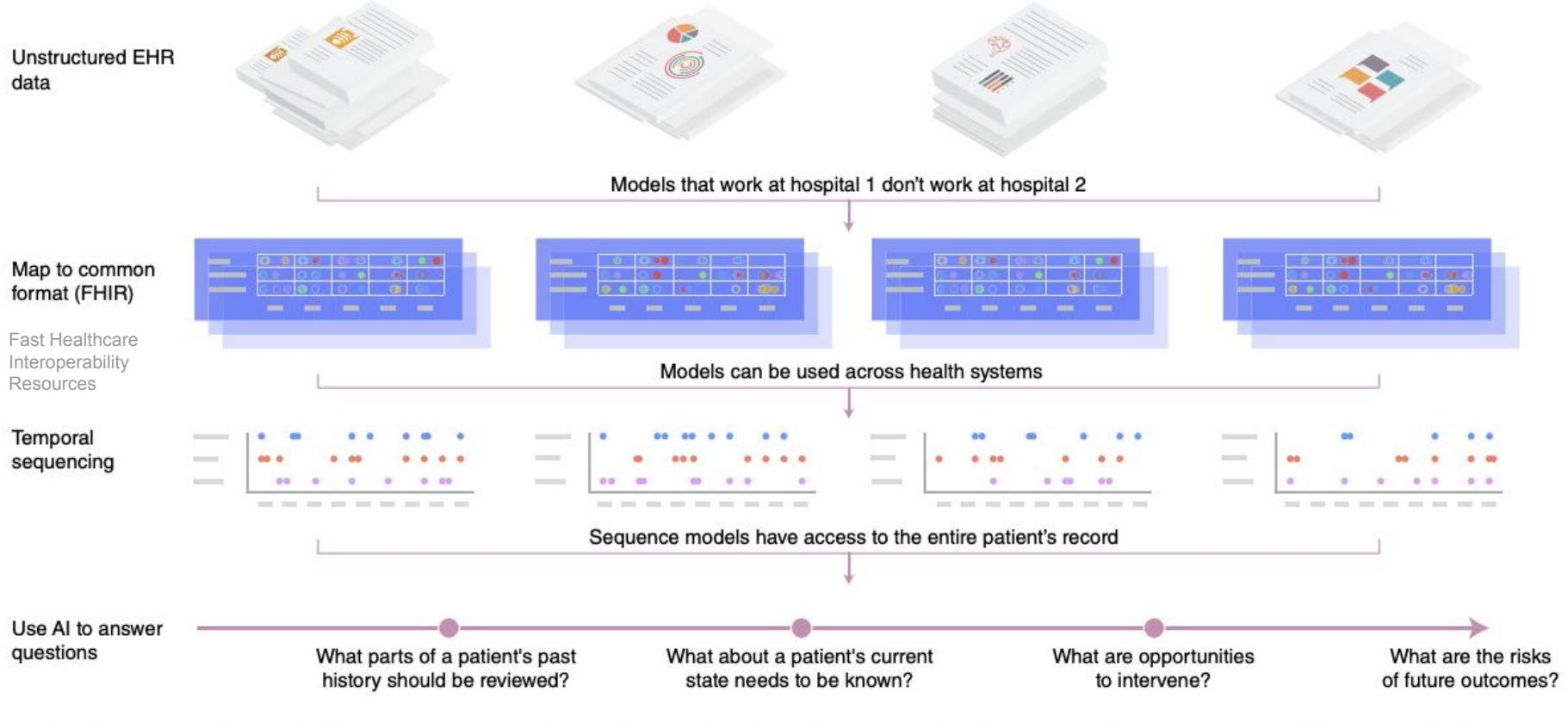
3. ML & DL Models

Named entity recognition, entity resolution, ...

| ED Triage Notes |
|--|
| states started last night, upper abd, took alka seltzer approx 0500, no relief. nausea no vomiting |
| Since yesterday 10/10 "constant Tylenol 1 hr ago. +nausea. diaphoretic. Mid abd radiates to back |
| Generalized abd radiating to lower x 3 days accompanied by dark stools. Now with bloody stool this am. Denies dizzy, sob, fatigue. Visiting from Japan on business." |



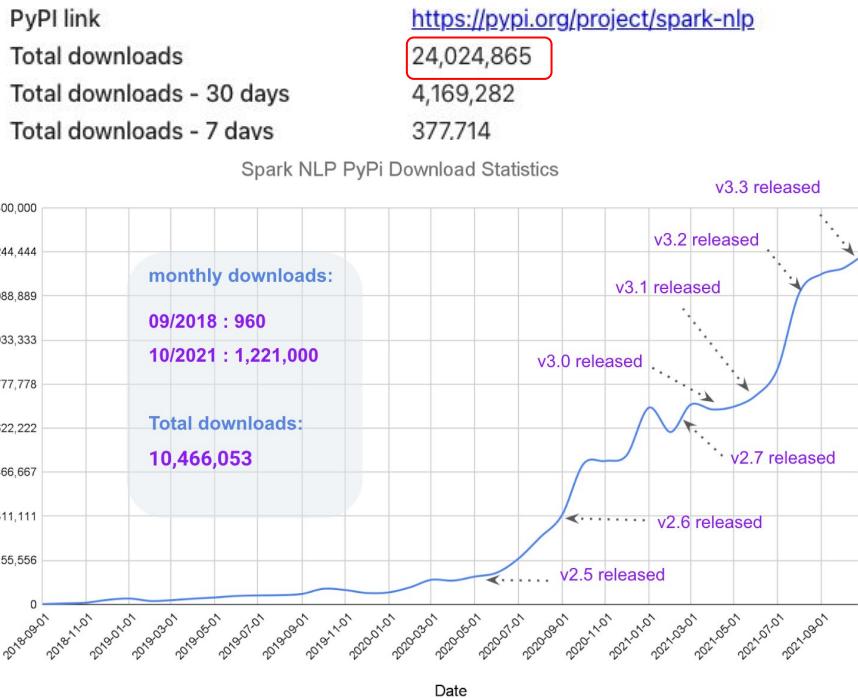
| Features | |
|---------------------|-----------------------|
| Type of Pain | Symptoms |
| Intensity of Pain | Onset of symptoms |
| Body part of region | Attempted home remedy |



“Systems used to generate health data are designed for operations, not to organize data effectively for research or analytics.”

Introducing Spark NLP

Daily ~ 50K
Monthly ~ 4.2M



- Spark NLP is an open-source natural language processing library, built on top of Apache Spark and Spark ML. (initial release: Oct 2017)
 - A single unified solution for all your NLP needs
 - Take advantage of transfer learning and implementing the latest and greatest SOTA algorithms and models in NLP research
 - The most widely used NLP library in industry (3 yrs in a row)
 - Delivering a mission-critical, enterprise grade NLP library (used by multiple Fortune 500)
 - Full-time development team (a new release every other week)

Spark NLP for Healthcare

Spark NLP for Healthcare provides

- accurate,
- scalable,
- private,
- tunable,
- modular

software library that helps healthcare & pharma organizations build longitudinal patient records and knowledge graphs on real-world EHR data.

| Clinical Entity Recognition | Clinical Entity Linking | Assertion Status | Relation Extraction | | | | | | |
|--|--|--|---|--|---|--|---|---|--|
| <p>40 units DOSAGE of insulin glargine DRUG at night FREQUENCY</p> | <p>Suspect diabetes SNOMED-CT: 473127005 Lisinopril 10 MG RxNorm: 316151 Hyponatremia ICD-10: E87.1</p> | <p>Fever and sore throat → PRESENT No stomach pain → ABSENT Father with Alzheimer → FAMILY</p> |  | | | | | | |
| Algorithms | | Content | | | | | | | |
| <p>Extract Knowledge</p> <ul style="list-style-type: none"> • Entity Linker • Entity Disambiguator • Document Classifier • Contextual Parser | | <p>De-identify text</p> <ul style="list-style-type: none"> • Structured Data • Unstructured Text • Obfuscator • Generalizer <p>Medical Transformers</p> <table border="1"> <tr> <td>JSL-BERT-Clinical</td> <td>BioBERT</td> </tr> <tr> <td>ClinicalBERT</td> <td>GloVe-Med</td> </tr> <tr> <td>GloVe-ICD-O</td> <td>BlueBERT</td> </tr> </table> | | JSL-BERT-Clinical | BioBERT | ClinicalBERT | GloVe-Med | GloVe-ICD-O | BlueBERT |
| JSL-BERT-Clinical | BioBERT | | | | | | | | |
| ClinicalBERT | GloVe-Med | | | | | | | | |
| GloVe-ICD-O | BlueBERT | | | | | | | | |
| <p>Split Text</p> <ul style="list-style-type: none"> • Sentence Detector • Deep Sentence Detector • Tokenizer • nGram Generator | | <p>Clean Medical Text</p> <ul style="list-style-type: none"> • Spell Checking • Spell Correction • Normalizer • Stopword Cleaner | | | | | | | |
| <p>Clinical Grammar</p> <ul style="list-style-type: none"> • Stemmer • Lemmatizer • Part of Speech Tagger • Dependency Parser | | <p>Find in Text</p> <ul style="list-style-type: none"> • Text Matcher • Regex Matcher • Date Matcher • Chunker | | | | | | | |
| Trainable & Tunable | Scalable to a Cluster | Fast Inference | Hardware Optimized | | | | | | |
|  |  |  |   | | | | | | |
| Community | |  | | | | | | | |
| Get Started | | View Documentation | | | | | | | |
| <p>600+ Pretrained Models</p> <table border="1"> <tr> <td>Clinical: Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections</td> <td>Anatomy: Organ, Subdivision, Cell, Structure Organism, Tissue, Gene, Chemical</td> </tr> <tr> <td>Drugs: Name, Dosage, Strength, Route, Duration, Frequency Poisons, Adverse Effects</td> <td>Demographics: Age, Gender, Height, Weight, Race, Ethnicity, Marital Status, Vital Signs</td> </tr> <tr> <td>Risk Factors: Smoking, Obesity, Diabetes, Hypertension, Substance Abuse</td> <td>Sensitive Data: Patient Name, Address, Phone, Email, Dates, Providers, Identifiers</td> </tr> </table> | | | | Clinical: Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections | Anatomy: Organ, Subdivision, Cell, Structure Organism, Tissue, Gene, Chemical | Drugs: Name, Dosage, Strength, Route, Duration, Frequency Poisons, Adverse Effects | Demographics: Age, Gender, Height, Weight, Race, Ethnicity, Marital Status, Vital Signs | Risk Factors: Smoking, Obesity, Diabetes, Hypertension, Substance Abuse | Sensitive Data: Patient Name, Address, Phone, Email, Dates, Providers, Identifiers |
| Clinical: Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections | Anatomy: Organ, Subdivision, Cell, Structure Organism, Tissue, Gene, Chemical | | | | | | | | |
| Drugs: Name, Dosage, Strength, Route, Duration, Frequency Poisons, Adverse Effects | Demographics: Age, Gender, Height, Weight, Race, Ethnicity, Marital Status, Vital Signs | | | | | | | | |
| Risk Factors: Smoking, Obesity, Diabetes, Hypertension, Substance Abuse | Sensitive Data: Patient Name, Address, Phone, Email, Dates, Providers, Identifiers | | | | | | | | |

Academic Activities & Benchmarks



Preparing for the Next Pandemic: Transfer Learning from Existing Diseases via Hierarchical Multi-Modal BERT Models to Predict COVID-19 Outcomes

Khushbu Agarwal¹, Sutanay Choudhury^{1*}, Sindhu Tipirneni², Pritam Mukherjee³, Colby Ham¹, Suzanne Tamang¹, Matthew Baker⁴, Siyi Tang⁵, Veysel Kocaman⁷, Olivier Gervais^{3,4}, Robert Rallo¹, and Chandan K Reddy¹

¹Pacific Northwest National Laboratory, Richland, 99354, USA

²Department of Computer Science, Virginia Tech, Arlington, 22203, USA

³Stanford Center for Biomedical Informatics Research, Department of Medicine, School of Medicine, Stanford University, Stanford, 94305, USA

⁴Department of Biomedical Data Science, Stanford University, Stanford, 94305, USA

⁵Department of Electrical Engineering, Stanford University, Stanford, 94305, USA

⁶Division of Immunology and Rheumatology, Department of Medicine, Stanford University, Stanford, 94305, USA

⁷John Snow Labs, Delaware City, 19968, USA

stanford.cs@stanford.edu; gervais@pnnl.gov



INFORMATICS PROFESSIONALS. LEADING THE WAY.

American Medical
Informatics
Association

Tracking the Evolution of COVID-19 via Temporal Comorbidity Analysis from Multi-Modal Data

Sutanay Choudhury¹, Khushbu Agarwal¹, Colby Ham¹, Pritam Mukherjee², Siyi Tang³, Sindhu Tipirneni³, Chandan Reddy⁴, Suzanne Tamang², Robert Rallo¹, Veysel Kocaman⁷,
¹Pacific Northwest National Laboratory; ²Stanford University; ³Virginia Tech;

John Snow Labs

Introduction

We aim to characterize the evolution in the effectiveness of treatment for different patient groups over the course of the COVID-19 pandemic. In contrast to most existing studies¹, we study the evolution of patient trajectories based on unique sets of frequent comorbid conditions discovered from the data. Further, we study the association between frequent co-morbid conditions to the length of stay (LOS) as a measure of treatment efficacy, for poor COVID-19 related outcomes.

Journal of Biomedical Semantics

SOFTWARE

Accurate Clinical and Biomedical Named Entity Recognition at Scale

Veysel Kocaman* and David Talby

*Correspondence:
veysel@johnsnowlabs.com
John Snow Labs, Lewes, DE, USA
Full list of author information is available at the end of the article

Scientific Document Understanding (SDU) at AAAI

Deeper Clinical Document Understanding Using Relation Extraction

Hasham Ul Haq, Veysel Kocaman, David Talby

John Snow Labs Inc.
16192 Coastal Highway
Lewes, DE, USA 19958
{hasham, veysel, david}@johnsnowlabs.com

Abstract

The surging amount of biomedical literature & digital clinical records presents a growing need for text mining techniques that can not only identify but also semantically enrich

publications and literature are growing rapidly, there still lacks structured knowledge that can be easily processed by computer programs. Relation Extraction becomes even more pertinent in biomedical research as it can provide the criti-



2021

New State-of-the-art (SOTA) Benchmarks



- ✓ 6 academic publications & events and 1 patent application, 20+ medium blogposts
- ✓ new SOTA benchmarks on Clinical NER challenges (i2b2 2010 Clinical, i2b2 2014 Deid, n2c2 2018 Medication)
- ✓ new SOTA benchmarks on Adverse Drug Reaction NER datasets (ADE, CADEC, SMM4H)
- ✓ new SOTA benchmarks on Adverse Drug Reaction classification datasets (ADR, CADEC)
- ✓ new SOTA benchmarks on Clinical Relation Extraction datasets (i2b2, temporal, ADE, Posology, PGR – 5 out of 7)



Mining Adverse Drug Reactions from Unstructured Mediums at Scale

Hasham Ul Haq, Veysel Kocaman, David Talby

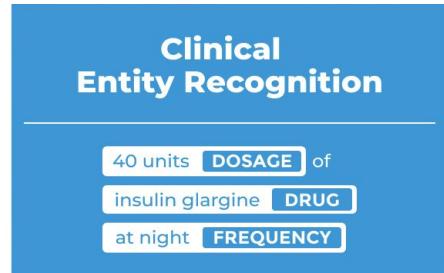
John Snow Labs Inc.
16192 Coastal Highway
Lewes, DE, USA 19958
{hasham, veysel, david}@johnsnowlabs.com

ADR's has been estimated to cost \$156 billion each year in the United States alone (van Der Hoof et al. 2006).

Finding all ADR's of a drug before it is marketed is not practical for several reasons. First, the number of human subjects going through clinical trials is often too small to detect rare ADR's. Second, many clinical trials are short-lasting while some ADR's take time to manifest. Third,

Health Intelligence (W3PHIAI-22) at AAAI

NLP in Healthcare



Clinical Entity Linking

Suspect diabetes SNOMED-CT: 473127005

Lisinopril 10 MG RxNorm: 316151

Pyponatremia ICD-10: E87.1

Assertion Status

Fever and sore throat → PRESENT

No stomach pain → ABSENT

Father with Alzheimer → FAMILY

De-Identification

Ora **NAME**, a **25 AGE** yo
cashier **PROFESSION** from
Morocco **LOCATION**

Relation Extraction

AFTER

Admitted for **nausea** due to **chemo**

Occurrence Symptom Treatment

CAUSED BY

NLP in Healthcare

"Mother with a lung cancer, a patient is diagnosed as breast cancer in 1991 and then admitted to Mayo Clinic in Oct 2000, went under chemo for 6 months, discharged in April 2001 with a prescription of 2 mg metformin 3 times per day."

Named Entities

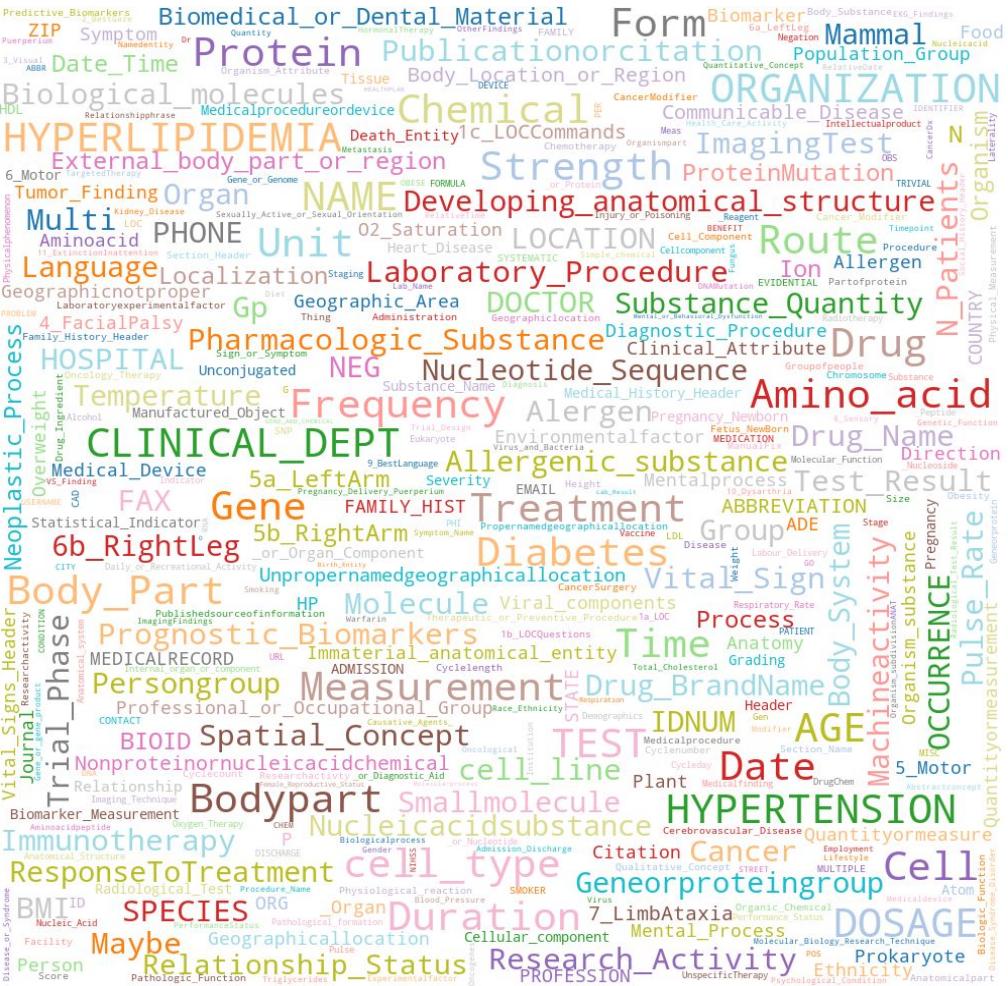
Mother with a lung cancer **ONCOLOGICAL** , a pregnant **PREGNANCY** patient is diagnosed as breast cancer **ONCOLOGICAL** in **1991 DATE** and then admitted **ADMISSION_DISCHARGE** to Mayo Clinic **CLINICAL_DEPT** in Oct **2000 DATE** , went under chemo **TREATMENT** for 6 months **DURATION** , discharged **ADMISSION_DISCHARGE** in **April 2001 DATE** with a prescription of **2 mg STRENGTH** metformin **DRUG_INGREDIENT** **3 times per day FREQUENCY** .

Pretrained NER Models

ner_ade_clinical
ner_posology_greedy
ner_risk_factors
jsl_ner_wip_clinical
ner_human_phenotype_gene_clinical
jsl_ner_wip_greedy_clinical
ner_cellular
ner_cancer_genetics
jsl_ner_wip_modifier_clinical
ner_drugs_greedy
ner_deid_sd_large
ner_diseases
nerdl_tumour_demo
ner_deid_subentity_augmented
ner_jsl_enriched
ner_genetic_variants
ner_bionlp
ner_measurements_clinical
ner_diseases_large
ner_radiology
ner_deid_augmented
ner_anatomy
ner_chemprot_clinical

| | |
|-----------------------------------|---------------------------------|
| ner_posology_greedy | ner_drugs |
| ner_risk_factors | ner_deid_sd |
| jsl_ner_wip_clinical | ner_posology_large |
| ner_human_phenotype_gene_clinical | ner_deid_large |
| jsl_ner_wip_greedy_clinical | ner_posology |
| ner_cellular | ner_deidentify_dl |
| ner_cancer_genetics | ner_deid_enriched |
| jsl_ner_wip_modifier_clinical | ner_bacterial_species |
| ner_drugs_greedy | ner_drugs_large |
| ner_deid_sd_large | ner_clinical_large |
| ner_diseases | jsl_rd_ner_wip_greedy_clinical |
| nerdl_tumour_demo | ner_medmentions_coarse |
| ner_deid_subentity_augmented | ner_radiology_wip_clinical |
| ner_jsl_enriched | ner_clinical |
| ner_genetic_variants | ner_chemicals |
| ner_bionlp | ner_deid_synthetic |
| ner_measurements_clinical | ner_events_clinical |
| ner_diseases_large | ner_posology_small |
| ner_radiology | ner_anatomy_coarse |
| ner_deid_augmented | ner_human_phenotype_go_clinical |
| ner_anatomy | ner_jsl_slim |
| ner_chemprot_clinical | ner_jsl |
| | ner_jsl_greedy |
| | ner_events_admission_clinical |

400+ entities from 100+ models



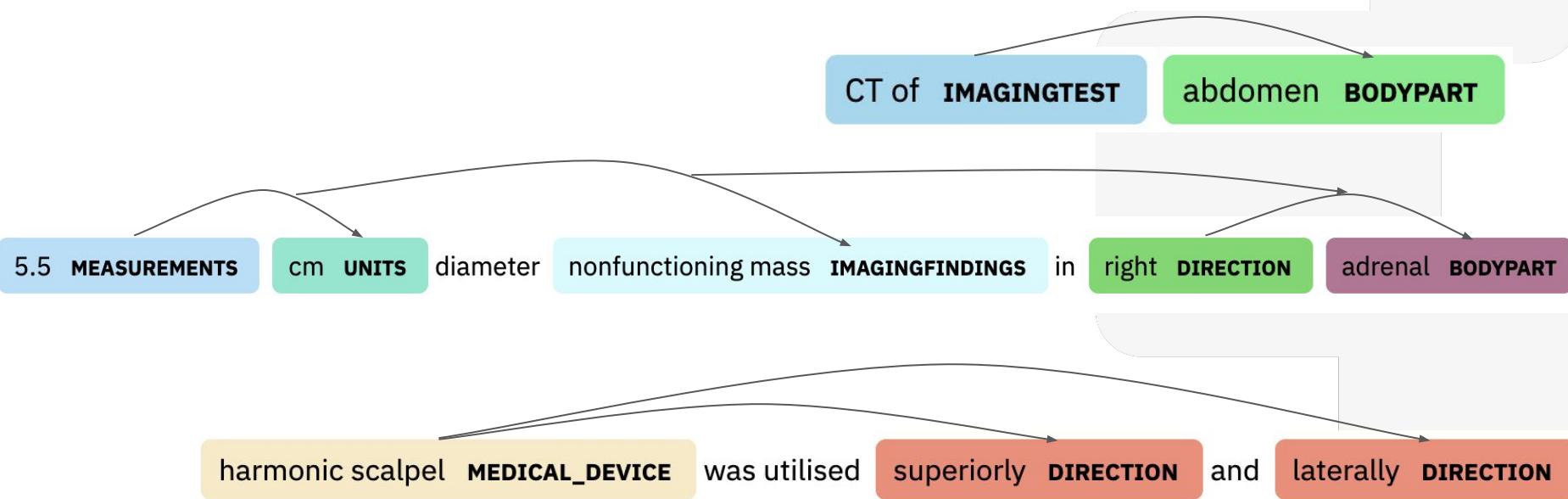
Assertion Status

"Mother with a lung cancer, a patient is diagnosed as breast cancer in 1991 and then admitted to Mayo Clinic in Oct 2000, went under chemo for 6 months, discharged in April 2001 with a prescription of 2 mg metformin 3x per day. No sign of gynecological disorder but she suffers from acute cramps if she doesn't take her drug."

| Chunk | Entity | Assertion |
|------------------------|-------------|-------------|
| lung cancer | Oncological | Family |
| breast cancer | Oncological | Past |
| chemo | Treatment | Past |
| gynecological disorder | Disorder | Absent |
| acute cramps | Disorder | Conditional |

Relation Extraction

"This is a 52-year-old inmate with a 5.5 cm diameter nonfunctioning mass in his right adrenal shown by CT of abdomen. During the umbilical hernia repair, the harmonic scalpel was utilised superiorly and laterally."



Entity Resolution

This is a 52-year-old AGE inmate with a 5.5 MEASUREMENTS cm UNITS diameter nonfunctioning mass SYMPTOM in his GENDER right DIRECTION adrenal BODYPART shown by CT of IMAGINGTEST abdomen BODYPART . During the umbilical hernia repair PROCEDURE , the harmonic scalpel MEDICAL_DEVICE was utilised superiorly DIRECTION and laterally DIRECTION .

Entity Resolution

ICD10CM, Snomed,

RxNorm, CPT-4,

ICD10CPS, RxCUI, ICDO,

UMLS, ATC, HPO,

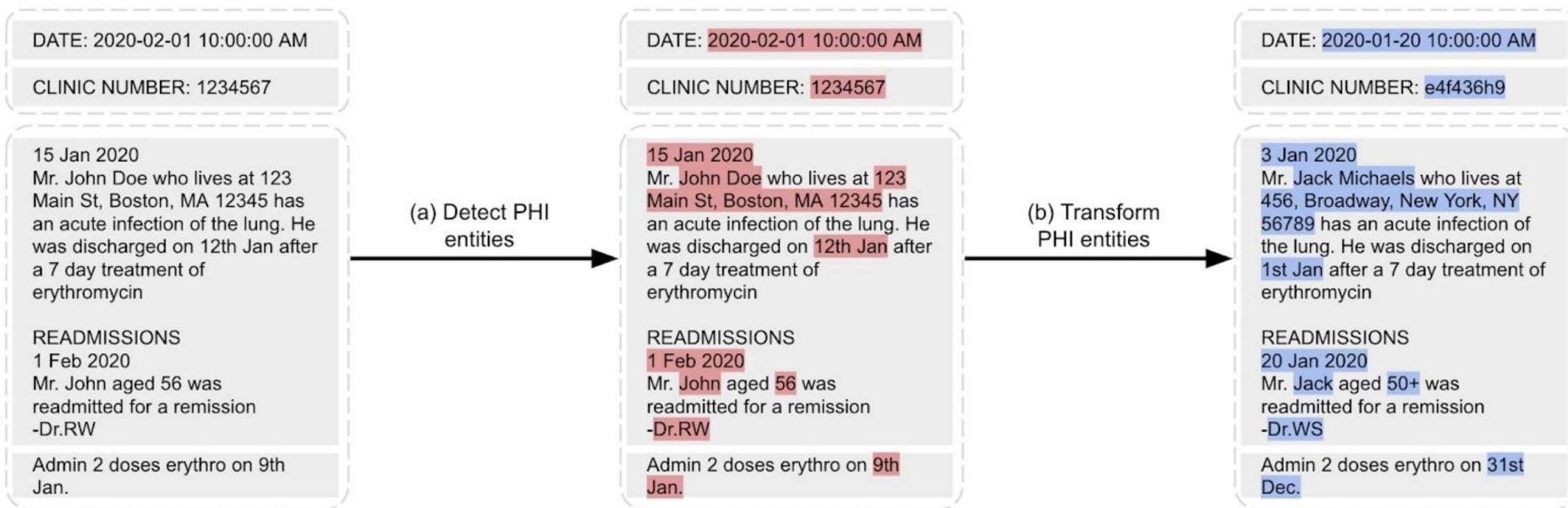
| Term | Vocab | Code | Explanation (ground truth) |
|---------------|-------|-------|---|
| CT | CPT-4 | 76497 | Unlisted computed tomography procedure |
| CT of abdomen | CPT-4 | 74150 | Computed tomography, abdomen; without contrast material |

weighted Sentence Chunk Embeddings (after 3.2.0)

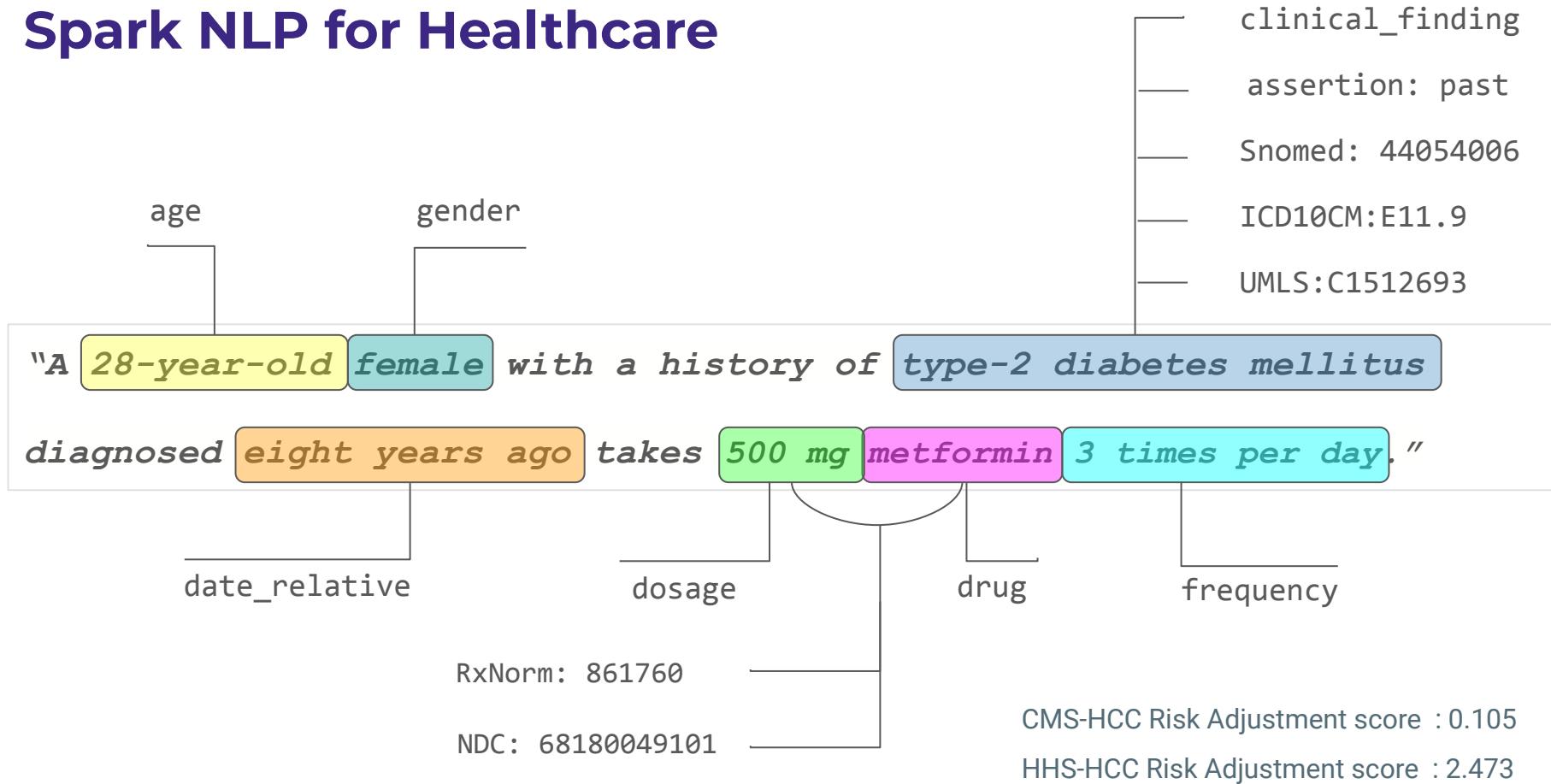
| Term | Vocab | Code | Explanation (ground truth) |
|------|-------|-------|---|
| CT | CPT-4 | 74150 | Computed tomography, abdomen; without contrast material |

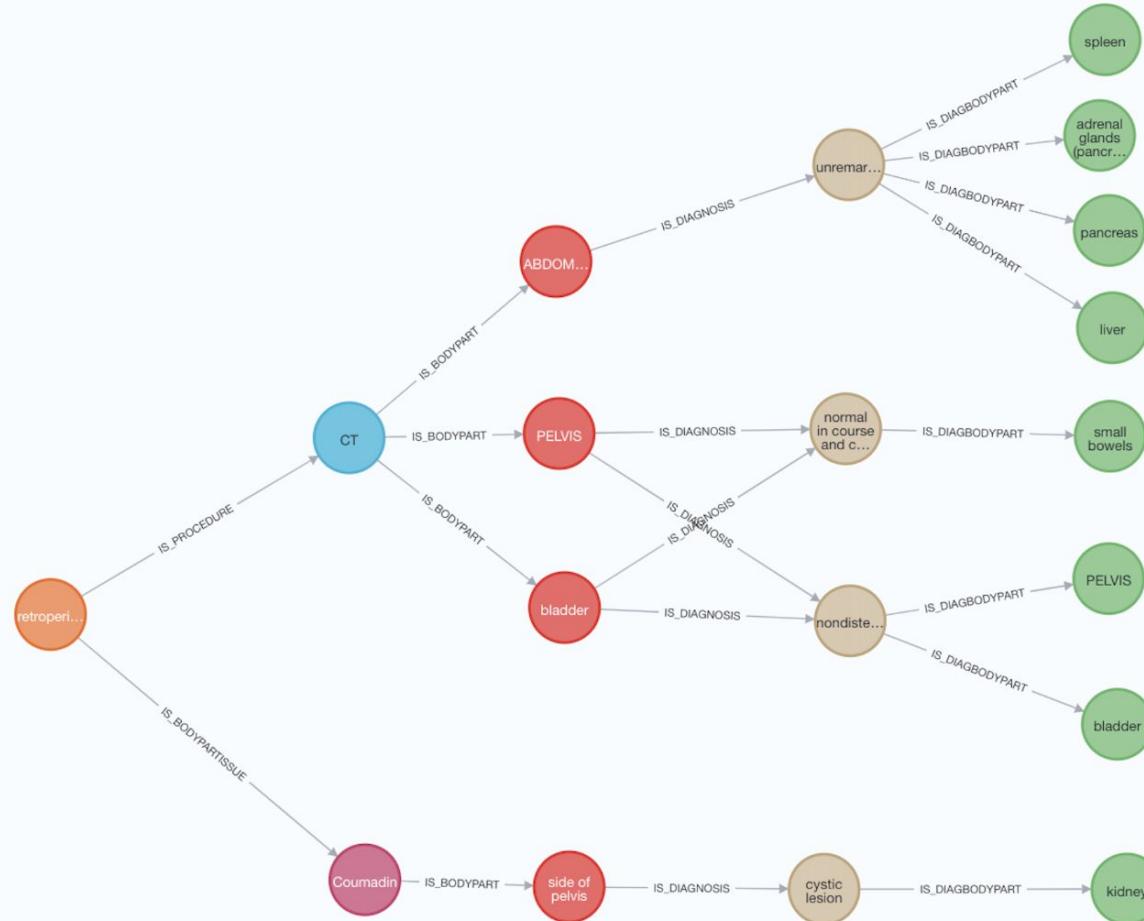
De-Identification

- * Identifies potential pieces of content with personal information about patients and remove them by replacing with semantic tags.



Spark NLP for Healthcare





REASON FOR EXAM: Evaluate for retroperitoneal hematoma on the right side of pelvis, the patient has been following, is currently on Coumadin.

CT ABDOMEN: There is no evidence for a retroperitoneal hematoma.

The liver, spleen, adrenal glands, and pancreas are unremarkable.

Within the superior pole of the left kidney, there is a 3.9 cm cystic lesion.

A 3.3 cm cystic lesion is also seen within the inferior pole of the left kidney.

No calcifications are noted. The kidneys are small bilaterally.

CT PELVIS: Evaluation of the bladder is limited due to the presence of a Foley catheter, the bladder is nondistended.

The large and small bowels are normal in course and caliber. There is no obstruction.

Spark NLP for Healthcare

Named Entity Recognition

ICD10 Resolver

Snomed Resolver

UMLS Resolver

Assertion Status Detection

Risk Adj. Module

RxNorm Resolver

Relationship Extraction

clinical_finding

Snomed: 44054006

ICD10CM:E11.9

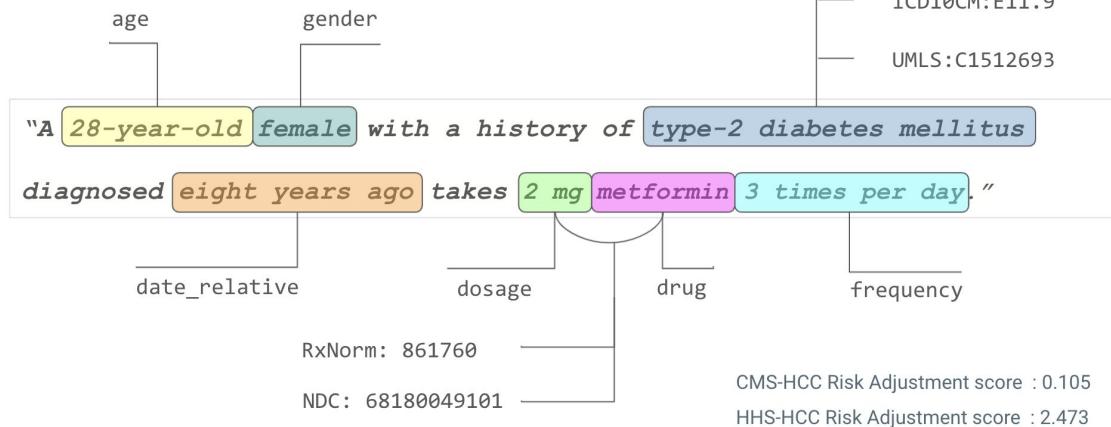
UMLS:C1512693

Sentence Splitter

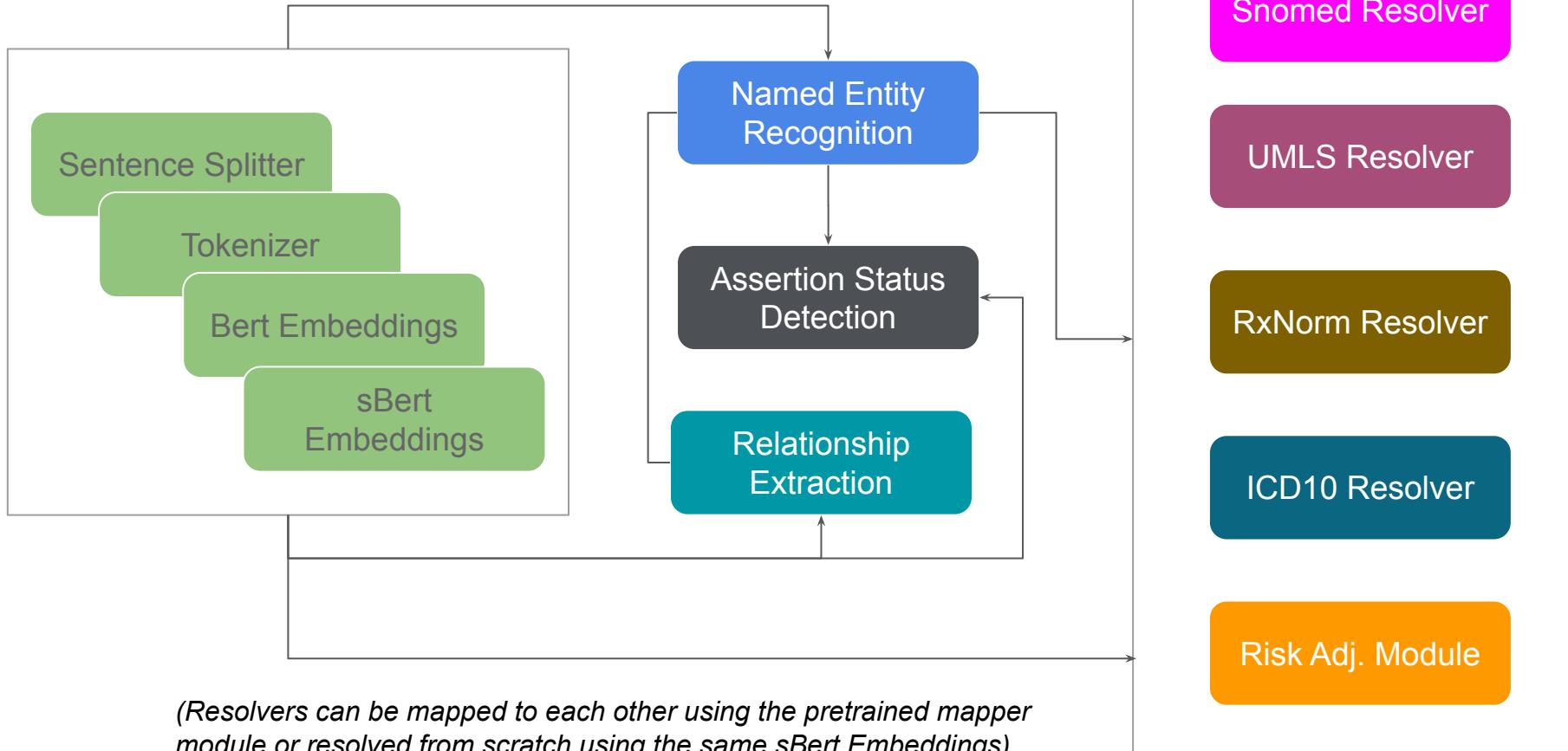
Tokenizer

Bert Embeddings

sBert Embeddings

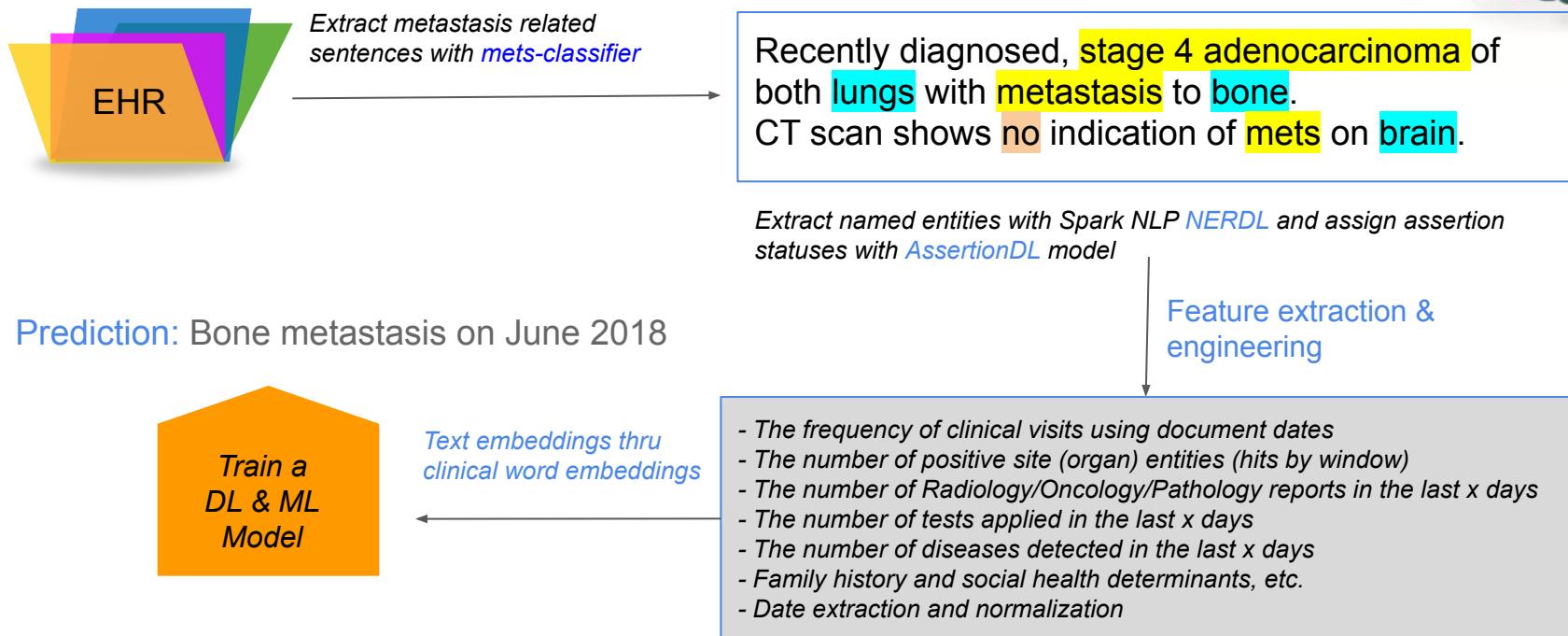


Spark NLP for Healthcare



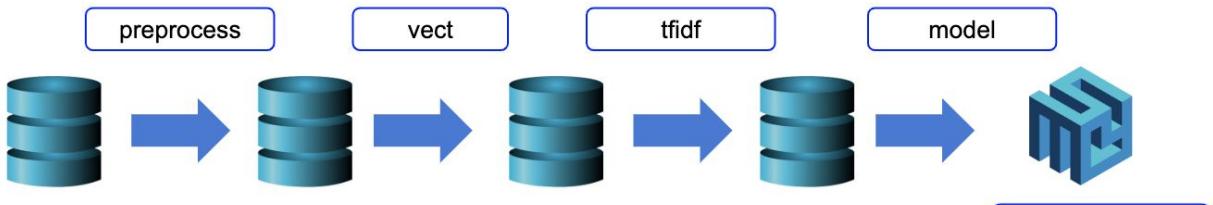
NLP in Healthcare

Case: Predicting if a patient would develop a metastasis on certain sites.

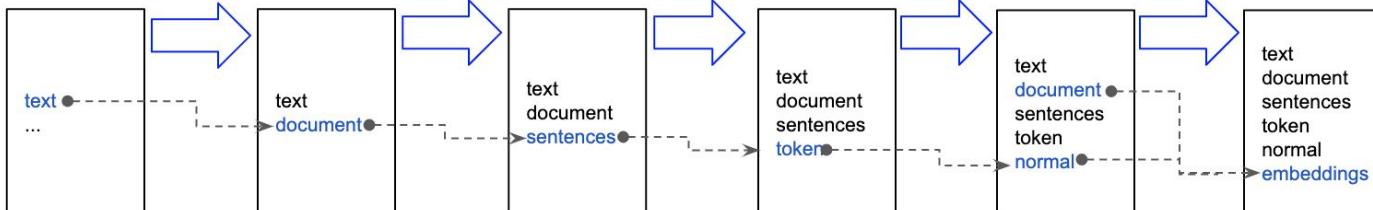


Spark NLP for Healthcare

Pipeline of annotators



DocumentAssembler() SentenceDetector() Tokenizer() Normalizer() WordEmbeddings()



DataFrame

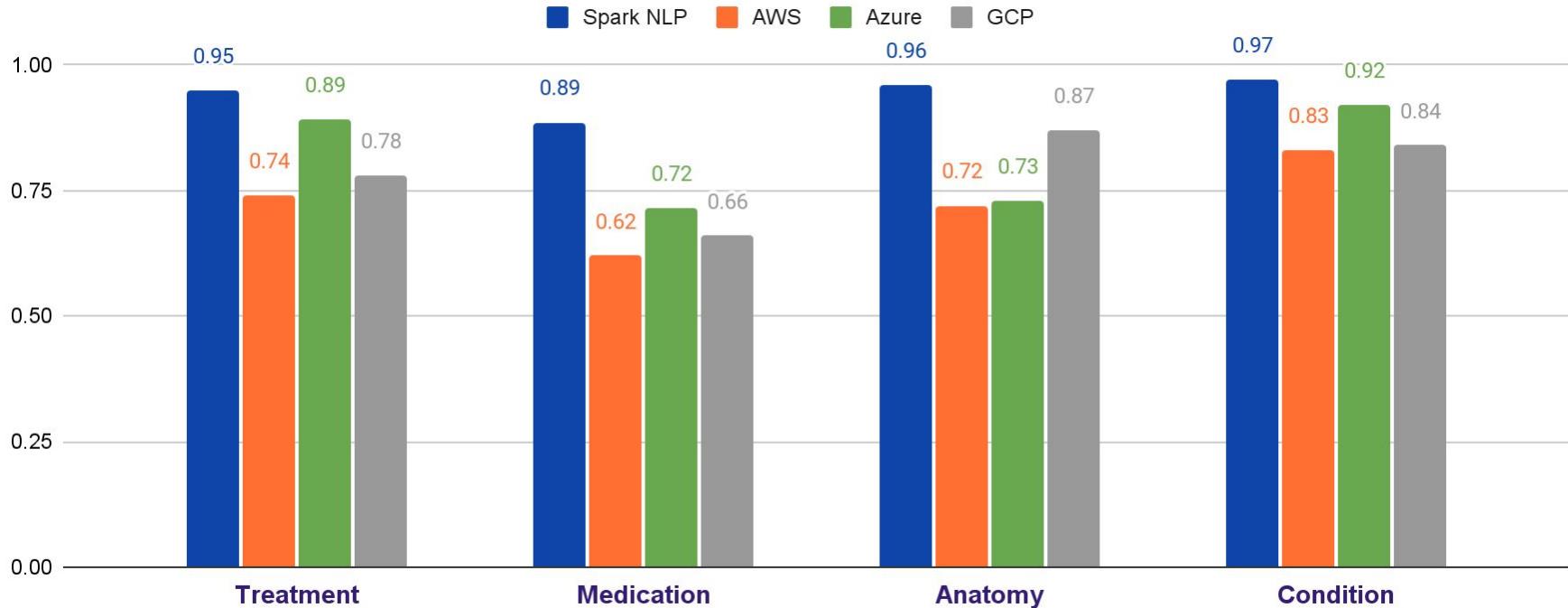
```
from pyspark.ml import Pipeline
document_assembler = DocumentAssembler()\
    .setInputCol("text")\
    .setOutputCol("document")
sentenceDetector = SentenceDetector()\
    .setInputCols(["document"])\
    .setOutputCol("sentences")
tokenizer = Tokenizer() \
    .setInputCols(["sentences"]) \
    .setOutputCol("token")
normalizer = Normalizer()\
    .setInputCols(["token"])\
    .setOutputCol("normal")
word_embeddings=WordEmbeddingsModel.pretrained()\
    .setInputCols(["document","normal"])\ 
    .setOutputCol("embeddings")
nlpPipeline = Pipeline(stages=[document_assembler,
    sentenceDetector,
    tokenizer,
    normalizer,
    word_embeddings,
])
nlpPipeline.fit(df).transform(df)
```

Attention (aka Bert) is all you need ?

| ner model | embeddings_clinical (BLSTM-CNN-Char) | | biobert (BLSTM-CNN-Char) | | BertForTokenClassification (SOTA) | |
|--------------|---|--------------|-----------------------------|--------------|--------------------------------------|-------------|
| | micro | macro | micro | macro | micro | macro |
| ner_jsl | 0.878 | 0.814 | 0.862 | 0.711 | 0.88 | 0.71 |
| ner_jsl_slim | 0.87 | 0.766 | 0.86 | 0.778 | 0.89 | 0.75 |
| ner_deid | 0.94 | 0.77 | 0.93 | 0.77 | 0.75 | 0.63 |
| ner_drug | 0.964 | 0.964 | 0.912 | 0.911 | 1 | 0.98 |
| ner_ade | 0.84 | 0.807 | 0.839 | 0.819 | 0.89 | 0.84 |

* On average, the GLoVe embeddings are 30% faster during training compared to BERT embeddings, and more than 5x faster during inference, while being on-par in terms of F1 score.

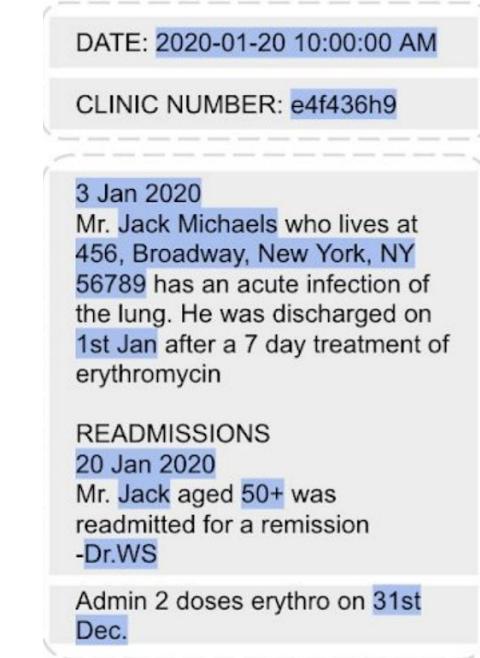
NER Benchmarks



~ 8K Sentences from MtSamples.com

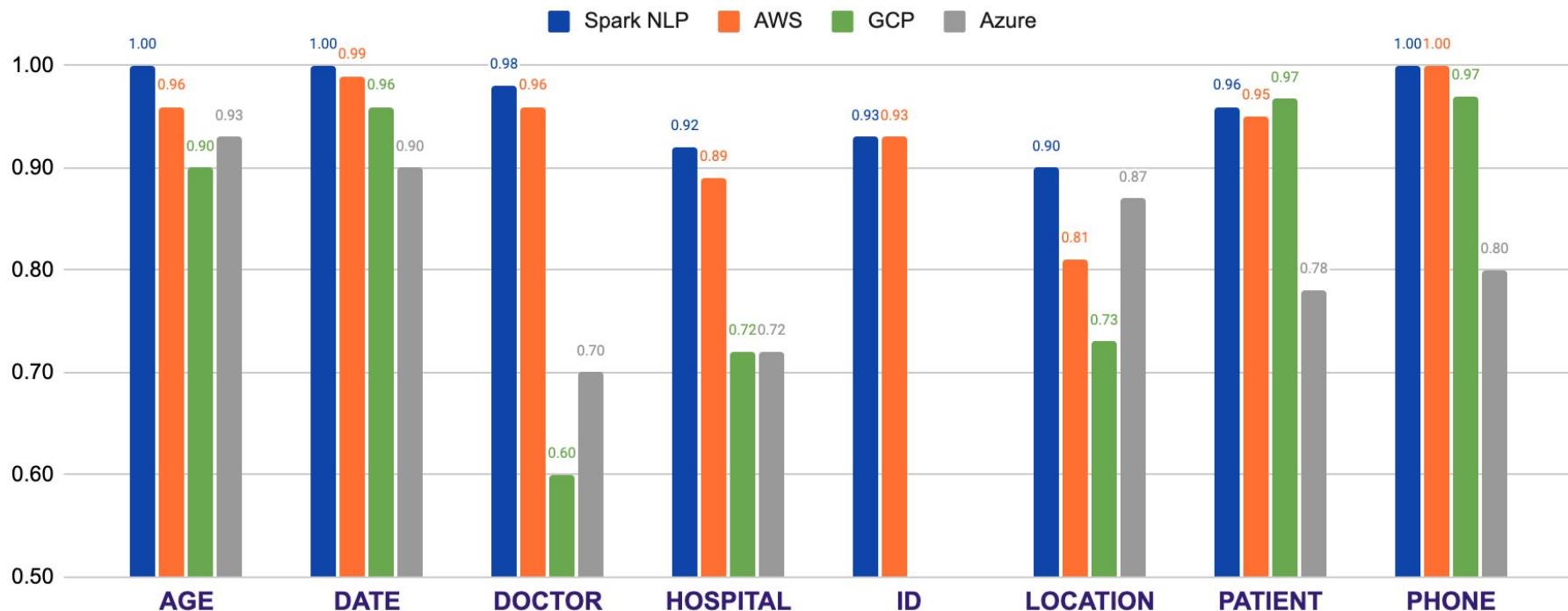
Spark NLP for Healthcare - Deidentification

| | English | German | French | Spanish | Italian |
|--------------|---------|--------|--------|---------|---------|
| PATIENT | 0.90 | 0.97 | 0.94 | 0.92 | 0.91 |
| DOCTOR | 0.94 | 0.98 | 0.99 | 0.92 | 0.92 |
| HOSPITAL | 0.91 | 1.00 | 0.94 | 0.86 | 0.90 |
| DATE | 0.98 | 1.00 | 0.98 | 0.99 | 0.98 |
| AGE | 0.94 | 0.99 | 0.86 | 0.98 | 0.98 |
| PROFESSION | 0.84 | 1.00 | 0.81 | 0.91 | 0.89 |
| ORGANIZATION | 0.77 | 0.94 | 0.77 | 0.83 | 0.74 |
| STREET | 0.9794 | 0.9802 | 0.8986 | 0.9448 | 0.9754 |
| CITY | 0.8331 | 0.9874 | 0.8643 | 0.8377 | 0.9678 |
| COUNTRY | 0.8083 | 0.9823 | 0.8983 | 0.8662 | 0.9262 |
| PHONE | 0.9412 | 0.882 | 0.9785 | 0.9027 | 0.9815 |
| USERNAME | 0.9215 | 1 | 0.9239 | 0.7407 | 0.9091 |
| ZIP | 0.9928 | - | 1 | 0.9895 | 0.9867 |

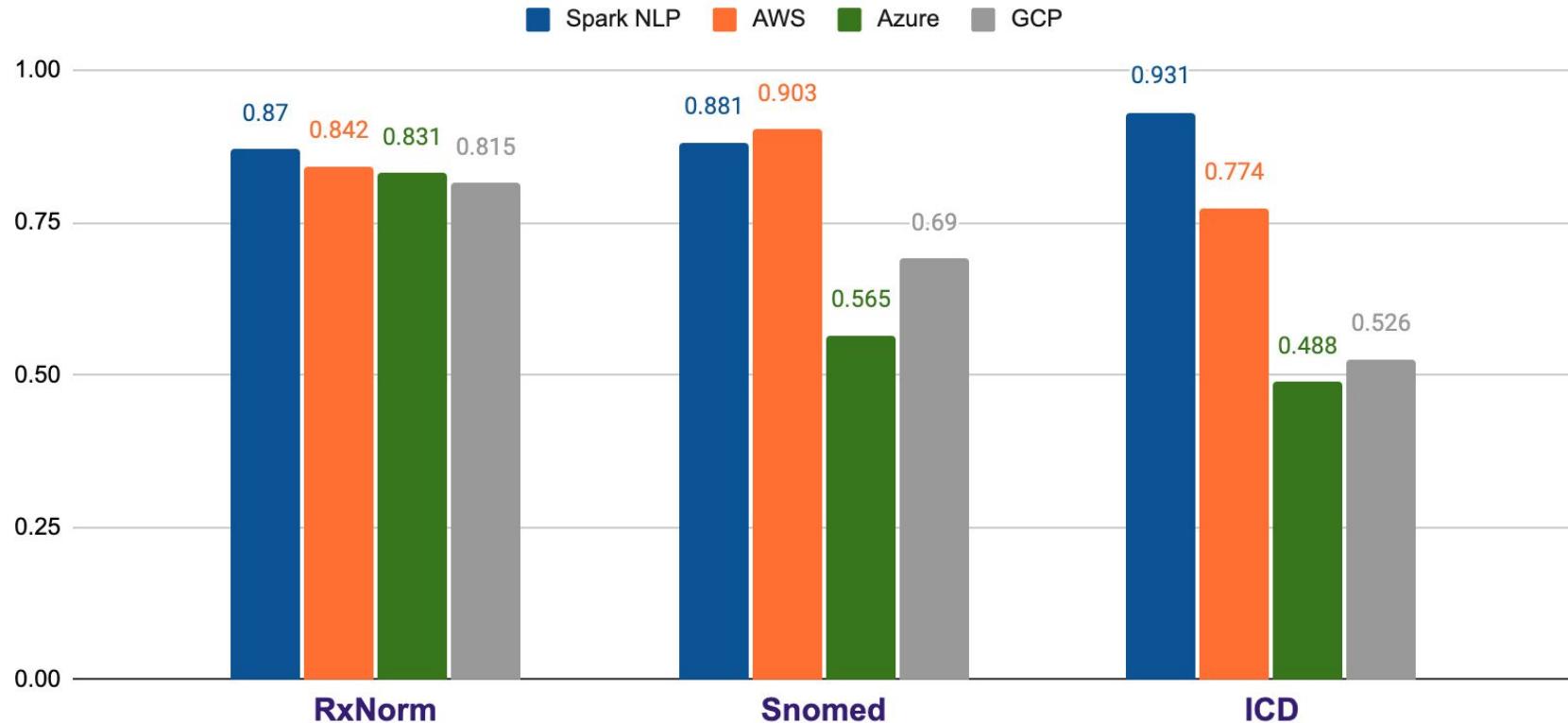


Sub-entity
(13-entity)

Deidentification Benchmarks



Top - 5 Results



~ 8K Sentences from MtSamples.com



Thank you !

Veysel Kocaman
Principal Data Scientist
John Snow Labs

