

# Spark NLP for Healthcare Data Scientists

July 20-21, 2022

**Veysel Kocaman**  
**Head of Data Science**  
[veysel@johnsnowlabs.com](mailto:veysel@johnsnowlabs.com)



# Agenda

Day	Dur.	Topic	Host
Day-1	50 min	<ul style="list-style-type: none"><li>- Intro to John Snow Labs and Spark NLP</li><li>- Healthcare NLP in Spark NLP</li></ul>	Veysel
	50 min	<ul style="list-style-type: none"><li>- Clinical Named Entity Recognition (<i>nb 1, 1.5 and 7</i>)</li></ul>	Veysel
	50 min	<ul style="list-style-type: none"><li>- Clinical Pretrained Pipelines (<i>nb 1, 1.5 and 7</i>)</li></ul>	Veysel
	50 min	<ul style="list-style-type: none"><li>- Clinical Assertion Status Model (<i>nb 2</i>)</li></ul>	Hasham
Day-2	50 min	<ul style="list-style-type: none"><li>- Clinical Entity Resolution (<i>nb 3, 3.1, 13, 13.1</i>)</li></ul>	Veysel
	50 min	<ul style="list-style-type: none"><li>- Clinical Entity Resolution (<i>nb 3, 3.1, 13, 13.1</i>)</li><li>- Misc (<i>ADE, ChunkMerger, ChunkMapper, ChunkKeyPhraseExtractor, ChunkSentenceSplitter, Generic Clf, Gender Clf</i>) (<i>nb 8, 9, 18, 16, 21</i>)</li></ul>	Veysel
	50 min	<ul style="list-style-type: none"><li>- Clinical Relation Extraction Model (<i>nb 10, 10.1, 10.2, 10.3</i>)</li></ul>	Hasham
	50 min	<ul style="list-style-type: none"><li>- De-Identification and Obfuscation of PHI (<i>nb 4</i>)</li></ul>	Luca

# Colab Notebooks

 Open in Colab

## RUNNING CODE:

[https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification\\_Trainings/Healthcare](https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification_Trainings/Healthcare)

[How to set up Google Colab]

## BOOKMARK:

<https://nlp.johnsnowlabs.com/models>  
[https://nlp.johnsnowlabs.com/docs/en/quickstart\\_spark-nlp.slack.com](https://nlp.johnsnowlabs.com/docs/en/quickstart_spark-nlp.slack.com)

### Clinical Named Entity Recognition

- [1.Clinical\\_Named\\_Entity\\_Recognition\\_Model.ipynb](#)
- [1.2.Contextual\\_Parser\\_Rule\\_Based\\_NER.ipynb](#)
- [1.3.prepare\\_CoNLL\\_from\\_annotations\\_for\\_NER.ipynb](#)
- [1.4.Resume\\_MedicalNer\\_Model\\_Training.ipynb](#)
- [1.5.BertForTokenClassification\\_NER\\_SparkNLP\\_with\\_Transformers.ipynb](#)
- [7.Clinical\\_NER\\_Chunk\\_Merger.ipynb](#)
- [16.Adverse\\_Drug\\_Event\\_ADE\\_NER\\_and\\_Classifier.ipynb](#)

### Clinical Assertion

- [2.Clinical\\_Assertion\\_Model.ipynb](#)
- [2.1.Scope\\_window\\_tuning\\_assertion\\_status\\_detection.ipynb](#)

### Clinical Entity Resolution

- [3.Clinical\\_Entity\\_Resolvers.ipynb](#)
- [3.1.Calculate\\_Medicare\\_Risk\\_Adjustment\\_Score.ipynb](#)
- [3.2.Sentence\\_Entity\\_Resolvers\\_with\\_EntityChunkEmbeddings.ipynb](#)
- [13.Snomed\\_Entity\\_Resolver\\_Model\\_Training.ipynb](#)
- [13.1.Finetuning\\_Sentence\\_Entity\\_Resolver\\_Model.ipynb](#)
- [22.CPT\\_Entity\\_Resolver.ipynb](#)
- [24.Improved\\_Entity\\_Resolvers\\_in\\_SparkNLP\\_with\\_sBert.ipynb](#)
- [24.1.Improved\\_Entity\\_Resolution\\_with\\_SentenceChunkEmbeddings.ipynb](#)

**spark-nlp==4.0.0**  
**spark-nlp-jsl==4.0.0**

 master spark-nlp-workshop / tutorials / Certification\_Trainings / Healthcare /

[Go to file](#) [Add file](#) [...](#)

### Clinical Relation Extraction

- [10.Clinical\\_Relation\\_Extraction.ipynb](#)
- [10.1.Clinical\\_Relation\\_Extraction\\_BodyParts\\_Models.ipynb](#)
- [10.2.Clinical\\_REL\\_Knowledge\\_Graph\\_with\\_Neo4j.ipynb](#)
- [10.3.ZeroShot\\_Clinical\\_Relation\\_Extraction.ipynb](#)

### Clinical De-identification

- [4.Clinical\\_Deidentification.ipynb](#)
- [4.1.Clinical\\_Deidentification\\_in\\_German.ipynb](#)
- [4.2.Clinical\\_Deidentification\\_in\\_Spanish.ipynb](#)
- [4.3.Clinical\\_Deidentification\\_SparkNLP\\_vs\\_Cloud\\_Providers\\_Comparison.ipynb](#)
- [4.4.Clinical\\_Deidentification\\_SparkNLP\\_vs\\_SpaCy\\_Comparison.ipynb](#)
- [4.5.Clinical\\_Deidentification\\_in\\_French.ipynb](#)
- [4.6.Clinical\\_Deidentification\\_in\\_Italian.ipynb](#)

### Optical Character Recognition with Spark OCR

- [5.Spark\\_OCR.ipynb](#)
- [5.1.Spark\\_OCR\\_Multi\\_Modals.ipynb](#)
- [5.2.Spark\\_OCR\\_Deidentification.ipynb](#)

### Clinical Pipelines

- [11.Pretrained\\_Clinical\\_Pipelines.ipynb](#)
- [11.1.Healthcare\\_Code\\_Mapping.ipynb](#)
- [11.2.Pretrained\\_NER\\_Profiling\\_Pipelines.ipynb](#)

### Classifiers

- [8.Generic\\_Classifier.ipynb](#)
- [21.Gender\\_Classifier.ipynb](#)

### Normalizers

- [23.Drug\\_Normalizer.ipynb](#)
- [25.Date\\_Normalizer.ipynb](#)

### Auxillary Notebooks

- [6.Clinical\\_Context\\_Spell\\_Checker.ipynb](#)
- [9.Chunk\\_Key\\_Phrase\\_Extraction.ipynb](#)
- [12.Named\\_Entity\\_Disambiguation.ipynb](#)
- [14.German\\_Healthcare\\_Models.ipynb](#)
- [15.German\\_Licensed\\_Models.ipynb](#)
- [17.Graph\\_builder\\_for\\_DL\\_models.ipynb](#)
- [18.Chunk\\_Sentence\\_Splitter.ipynb](#)
- [19.Financial\\_Contract\\_NER.ipynb](#)
- [20.SentenceDetectorDL\\_Healthcare.ipynb](#)
- [26.Chunk\\_Mapping.ipynb](#)

[johnsnowlabs.com/spark-nlp-in-action/](https://johnsnowlabs.com/spark-nlp-in-action/)

Home Products Solutions Customers Company Learn Install Software Schedule a Call

**Spark NLP: English**

- Entity Mention & Intent
- Classify Documents
- Recognize Entities
- Detect Sentiment & Emotion
- Analyze Spelling & Grammar

**Spark NLP: World Languages**

- Identify & Translate Languages
- European Languages
- East Asian Languages
- Languages of India
- Middle Eastern Languages
- Languages of Africa

**Spark NLP for Healthcare**

- Recognize Clinical Entities
- Recognize Biomedical Entities
- De-Identification
- Resolve Entities to Codes
- Recognize Disease Symptoms
- Extract Relationships
- Split & Clean Medical Text
- Analyze non-English Medical Text

**Spark OCR**

- Extract Text from Documents
- Enhance Low-Quality Images
- Convert Tables & Structured Data
- Analyze non-English Medical Text

# Spark NLP in Action

Spark NLP – English → Recognize Entities



Recognize entities in text

[Live Demo](#) [Create Notebook](#)



Recognize more entities in text

[Live Demo](#) [Create Notebook](#)



Detect Key Phrases

[Live Demo](#) [Create Notebook](#)



Find a text in document

[Live Demo](#) [Create Notebook](#)



Detect and normalize dates

[Live Demo](#) [Create Notebook](#)

http://nlp.johnsnowlabs.com/docs/en/licensed\_install



Home Docs Learn Models Demo [Get Started](#) [Documentation](#)

**Spark NLP for Healthcare**

Getting Started  
Installation  
Annotators  
Training  
Models  
Evaluation  
Scalability API (Scalable)  
Python API (Spacy)  
Version Compatibility  
Release Notes  
Cluster Speed Benchmarks

Documentation > Spark NLP for Healthcare

# Spark NLP for Healthcare

## Getting started

Spark NLP for Healthcare is a commercial extension of Spark NLP for clinical and biomedical text mining. If you don't have a Spark NLP for Healthcare subscription yet, you can ask for a free trial by clicking on the button below.

[Try Now](#)

Spark NLP for Healthcare provides healthcare-specific annotators, pipelines, models, and embeddings for:

- Clinical entity recognition
- Clinical Entity Linking
- Entity normalization
- Attention Status Detection
- De-identification
- Relative Entropy
- Named Clinical Entity Extraction

After you are going to use any premained NER model, you don't need to install licensed library. As long as you have the AWS keys and license keys in your environment, you will be able to use licensed NER models with Spark NLP public library. For the other licensed premained models like AssertionOrC, DeIdentification, EntityResolvers and Relation Extraction models, you will need to install Spark NLP Enterprise as well.

The library offers access to several clinical and biomedical transformers: JS\_BERT\_Clinical, BioBERT, CliniceBert, Glue-Med, Glue-ICD-10. It also includes over 50 pre-trained healthcare models, that can recognize the following entities (any many more):

- Clinical - support Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections
- Drugs - support Name, Dosage, Strength, Route, Duration, Frequency

Getting started  
Install Spark NLP for Healthcare  
Setup AWS CLI Credentials for Amazon SageMaker  
Start Spark NLP for Healthcare Sessions from Python  
Spark NLP for Healthcare Cheat Sheet  
Install Spark NLP for Healthcare Databricks  
Install Spark NLP for Healthcare GCP Notebook  
Google Colab Notebook

The screenshot shows the homepage of the John Snow Labs NLP Models Hub. The header features the John Snow LABS logo and navigation links for Home, Docs, Learn, Models, Demo, and Support. A prominent search bar at the top right allows users to "Upload Your Model!" or search for specific models and pipelines. Below the search bar, a large banner reads "NLP Models Hub" with the subtitle "A place for sharing and discovering Spark NLP models and pipelines". A secondary search bar below the main one is used to find "models & pipelines". The main content area displays a grid of supported NLP models, each with a "Supported" badge, a title, and a brief description. The first model listed is "Sentence Entity Resolver for UMLS CUI Codes (Clinical Drug)", followed by "Task: Entity Resolution", "Language: English", and "Edition: Spark NLP for Healthcare 3.2.3". Other models shown include "Sentence Entity Resolver for RxNorm (RxNorm\_Case\_Cased\_mllm\_embeddings)" and "Text: Entity Resolution" for "Task: Entity Resolution", "Language: English", and "Edition: Spark NLP for Healthcare 3.2.3". Each model entry includes a date (e.g., Date: 10.2021) and a small icon.

The screenshot shows a browser window displaying the `sparknlp_jsl.annimator` module from the `nlp.johnsnowlabs.com` website. The page includes a header with the John Snow Labs logo and navigation links for "Getting Started" and "API Reference". The main content area features a search bar and a table of contents titled "sparknlp\_jsl.annimator". Below the table of contents, several sub-sections are listed, each with a brief description and a link to the full documentation. The sections include:

- sparknlp\_jsl.annimator**
- sparknlp\_jsl.annimator.AssertionOrApproach**
- sparknlp\_jsl.annimator.AssertionOrModel**
- sparknlp\_jsl.annimator.AssertionOrFilter**
- sparknlp\_jsl.annimator.AssertionLogitApproach**
- sparknlp\_jsl.annimator.AssertionLogitModel**
- sparknlp\_jsl.annimator.AssertionEmbeddings**
- sparknlp\_jsl.annimator.BertSentenceChunkEmbed**
- sparknlp\_jsl.annimator.Chunk2Token**
- sparknlp\_jsl.annimator.ChunkConverter**
- sparknlp\_jsl.annimator.ChunkFilterer**
- sparknlp\_jsl.annimator.ChunkFiltererApproach**
- sparknlp\_jsl.annimator.ChunkMergeApproach**
- sparknlp\_jsl.annimator.ChunkMergeModel**
- sparknlp\_jsl.annimator.ContextualPaserParam**
- sparknlp\_jsl.annimator.ContextualPaserApproach**
- sparknlp\_jsl.annimator.DataNormalizer**
- sparknlp\_jsl.annimator.Delimiter**
- sparknlp\_jsl.annimator.DelimiterClassificationModel**
- sparknlp\_jsl.annimator.DocumentLogregClassif**
- sparknlp\_jsl.annimator.DocumentLogregClassifApproach**
- sparknlp\_jsl.annimator.DrugNormalizer**
- sparknlp\_jsl.annimator.GenericClassifierApproach**
- sparknlp\_jsl.annimator.GenericClassifierModel**
- sparknlp\_jsl.annimator.IOTagger**
- sparknlp\_jsl.annimator.MedicalApproach**
- sparknlp\_jsl.annimator.MedicalModel**

JohnSnowLabs / spark-nlp-workshop		Public	
		Code	Watch
<a href="#">Code</a>	<a href="#">Issues</a>	<a href="#">Pull requests</a>	<a href="#">Discussions</a>
<a href="#">Actions</a>	<a href="#">Projects</a>	<a href="#">Wiki</a>	<a href="#">Security</a>
<a href="#">master</a> · <a href="#">75 branches</a> · <a href="#">7 tags</a>		<a href="#">Go to file</a>	<a href="#">Add file</a> · <a href="#">Code</a>
 galiph Merge pull request #400 from JohnSnowLabs/galiph · 2 hours ago		✓ 1873387 · 2 hours ago	1501 commits
<a href="#">data</a>	Add files via upload	2 months ago	
<a href="#">databricks</a>	Update benchmark.md	11 days ago	
<a href="#">java</a>	add java examples	3 months ago	
<a href="#">jupyter</a>	Update docker-compose.yaml	3 days ago	
<a href="#">nlu</a>	nlu notebooks updated	4 months ago	
<a href="#">platforms</a>	added more info for troubleshooting	2 months ago	
<a href="#">scala</a>	update scala codes	5 months ago	
<a href="#">tutorials</a>	Notebooks updated with v3.0	2 hours ago	
<a href="#">zeppelin</a>	clean notebook	3 years ago	
<a href="#">.gitattributes</a>	Ignore html from linguist-vendored	3 years ago	
<a href="#">.gitignore</a>	removed outdated notebooks	13 months ago	
<a href="#">ISSUE_TEMPLATE.md</a>	Create ISSUE_TEMPLATE.md	3 years ago	
<a href="#">LICENSE</a>	Initial commit	3 years ago	
<a href="#">README.md</a>	Update README.md	26 days ago	
<a href="#">colab_setup.sh</a>	Update colab_setup.sh	4 months ago	
<a href="#">jd_colab_setup.sh</a>	Update jd_colab_setup.sh	6 months ago	
<a href="#">jd_colab_setup_with_OCR.sh</a>	Update jd_colab_setup_with_OCR.sh	2 months ago	
<a href="#">jd_sagemaker_setup.sh</a>	corrected paths for binarys	4 months ago	
<a href="#">jd_sagemaker_setup_30.1.sh</a>	added script	4 months ago	
<a href="#">jd_sagemaker_setup_with_OCR.sh</a>	Update jd_sagemaker_setup_with_OCR.sh	5 months ago	

The screenshot shows the ScalaDoc interface for the `MedicalNerApproach` class. The top navigation bar includes links for Home, Back, Forward, Help, and Update. The URL is `nlp.johnsnowlabs.co...`. The search bar contains the query "Search". The title is "Spark NLP 3.0.3 ScalaDoc". On the right side, there's a sidebar with the package structure:

```
com.johnsnowlabs.nlp.annotators  
  MedicalNerApproach  
    Companion object  
    MedicalNerApproach  
      class MedicalNerApproach extends  
        AnnotatorApproach[Hed[calNodeModel]] with  
        NerApproach[Hed[calNodeApproach]] with Logging with  
        ParamsAndFeaturesWritable with CheckLicense
```

The main content area displays the `Linear Supertypes` for `MedicalNerApproach`, which includes `IOBTagger`, `MedicalNerApproach`, `MedicalNerModel`, `NamedEntityConfidence`, `NaiveChunker`, `NerConverterInternal`, `NerTaggerInternal`, and `TagsEncodingInternal`. Below this, the `anno` section lists two annotations:

- `val InputAnnotatorTypes: Array[String]`  
Input annotator types: DOCUMENT, TOKEN, WORD\_EMBEDDINGS
- `val outputAnnotatorType: String`  
Input annotator types: NAMED\_ENTITY

The bottom section, `getParam`, lists several parameters:

- `def getBatchSize: Int`  
Batch size
- `def getConfigProtoBytes: Option[Array[Byte]]`  
ConfigProto from tensorflow, serialized into byte array.
- `def getDropout: Float`  
Dropout coefficient
- `def getEnableMemoryOptimizer: Boolean`  
Memory Optimizer
- `def getEnableOutputLogs: Boolean`  
Whether to output to annotators log folder

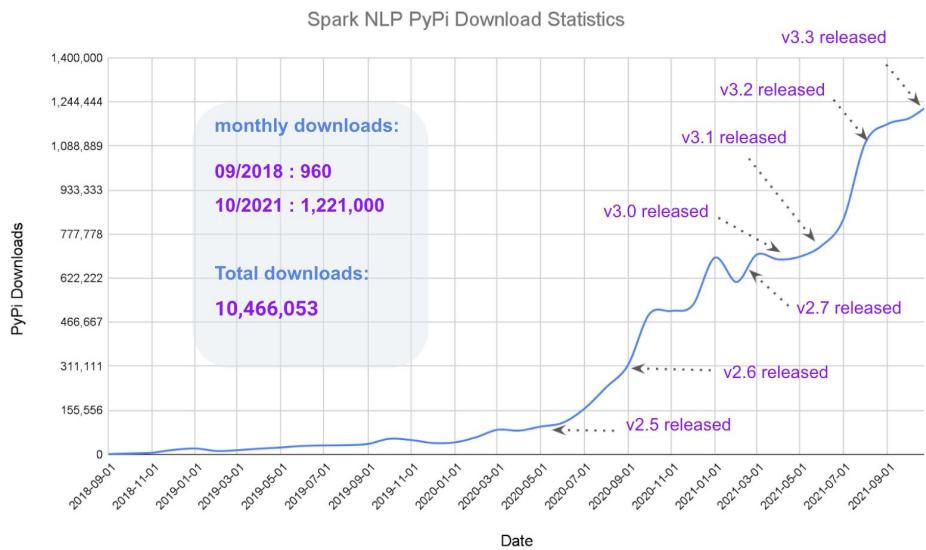
# Part - I

- ❖ Overview and key concepts in Spark NLP
- ❖ NLP basics & review
- ❖ Common medical NLP use cases
- ❖ Clinical named entity recognition

# Introducing Spark NLP

Daily ~ 80K  
Monthly ~ 2.3M

PyPI link	<a href="https://pypi.org/project/spark-nlp">https://pypi.org/project/spark-nlp</a>
Total downloads	30,067,117
Total downloads - 30 days	2,331,332
Total downloads - 7 days	543,740



- 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 <b>DOSAGE</b> of insulin glargine <b>DRUG</b> at night <b>FREQUENCY</b></p>	<p>Suspect diabetes SNOMED-CT: <b>473127005</b> Lisinopril 10 MG RxNorm: <b>316151</b> Hyponatremia ICD-10: <b>E87.1</b></p>	<p>Fever and sore throat → PRESENT No stomach pain → ABSENT Father with Alzheimer → FAMILY</p>							
Algorithms		Content							
<p><b>Extract Knowledge</b></p> <ul style="list-style-type: none"> <li>• Entity Linker</li> <li>• Entity Disambiguator</li> <li>• Document Classifier</li> <li>• Contextual Parser</li> </ul>		<p><b>De-identify text</b></p> <ul style="list-style-type: none"> <li>• Structured Data</li> <li>• Unstructured Text</li> <li>• Obfuscator</li> <li>• Generalizer</li> </ul>							
<p><b>Split Text</b></p> <ul style="list-style-type: none"> <li>• Sentence Detector</li> <li>• Deep Sentence Detector</li> <li>• Tokenizer</li> <li>• nGram Generator</li> </ul>		<p><b>Clean Medical Text</b></p> <ul style="list-style-type: none"> <li>• Spell Checking</li> <li>• Spell Correction</li> <li>• Normalizer</li> <li>• Stopword Cleaner</li> </ul>							
<p><b>Clinical Grammar</b></p> <ul style="list-style-type: none"> <li>• Stemmer</li> <li>• Lemmatizer</li> <li>• Part of Speech Tagger</li> <li>• Dependency Parser</li> </ul>		<p><b>Find in Text</b></p> <ul style="list-style-type: none"> <li>• Text Matcher</li> <li>• Regex Matcher</li> <li>• Date Matcher</li> <li>• Chunker</li> </ul>							
Trainable & Tunable	Scalable to a Cluster	Fast Inference	Hardware Optimized						
			 						
Community									
<a href="#">Get Started</a>		<a href="#">View Documentation</a>							
<p><b>600+ Pretrained Models</b></p> <table border="1"> <tr> <td> <b>Clinical:</b>            Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections         </td><td> <b>Anatomy:</b>            Organ, Subdivision, Cell, Structure, Organism, Tissue, Gene, Chemical         </td></tr> <tr> <td> <b>Drugs:</b>            Name, Dosage, Strength, Route, Duration, Frequency, Poisons, Adverse Effects         </td><td> <b>Demographics:</b>            Age, Gender, Height, Weight, Race, Ethnicity, Marital Status, Vital Signs         </td></tr> <tr> <td> <b>Risk Factors:</b>            Smoking, Obesity, Diabetes, Hypertension, Substance Abuse         </td><td> <b>Sensitive Data:</b>            Patient Name, Address, Phone, Email, Dates, Providers, Identifiers         </td></tr> </table>				<b>Clinical:</b> Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections	<b>Anatomy:</b> Organ, Subdivision, Cell, Structure, Organism, Tissue, Gene, Chemical	<b>Drugs:</b> Name, Dosage, Strength, Route, Duration, Frequency, Poisons, Adverse Effects	<b>Demographics:</b> Age, Gender, Height, Weight, Race, Ethnicity, Marital Status, Vital Signs	<b>Risk Factors:</b> Smoking, Obesity, Diabetes, Hypertension, Substance Abuse	<b>Sensitive Data:</b> Patient Name, Address, Phone, Email, Dates, Providers, Identifiers
<b>Clinical:</b> Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections	<b>Anatomy:</b> Organ, Subdivision, Cell, Structure, Organism, Tissue, Gene, Chemical								
<b>Drugs:</b> Name, Dosage, Strength, Route, Duration, Frequency, Poisons, Adverse Effects	<b>Demographics:</b> Age, Gender, Height, Weight, Race, Ethnicity, Marital Status, Vital Signs								
<b>Risk Factors:</b> Smoking, Obesity, Diabetes, Hypertension, Substance Abuse	<b>Sensitive Data:</b> 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 Agarwal<sup>1</sup>, Sutanay Choudhury<sup>1\*</sup>, Sindhu Tipirneni<sup>2</sup>, Pritam Mukherjee<sup>3</sup>, Colby Ham<sup>1</sup>, Suzanne Tamang<sup>1</sup>, Matthew Baker<sup>4</sup>, Siyi Tang<sup>5</sup>, Veysel Kocaman<sup>7</sup>, Olivier Gevaert<sup>1,5</sup>, Robert Rallo<sup>1</sup>, and Chandan K Reddy<sup>1</sup>

<sup>1</sup>Pacific Northwest National Laboratory, Richland, 99354, USA

<sup>2</sup>Department of Computer Science, Virginia Tech, Arlington, 22203, USA

<sup>3</sup>Stanford Center for Biomedical Informatics Research, Department of Medicine, School of Medicine, Stanford University, Stanford, 94305, USA

<sup>4</sup>Department of Biomedical Data Science, Stanford University, Stanford, 94305, USA

<sup>5</sup>Department of Electrical Engineering, Stanford University, Stanford, 94305, USA

<sup>6</sup>Division of Immunology and Rheumatology, Department of Medicine, Stanford University, Stanford, 94305, USA

<sup>7</sup>John Snow Labs, Delaware City, 19968, USA

stanford.cs@stanford.edu



American Medical  
Informatics  
Association

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

Sutanay Choudhury<sup>1</sup>, Khushbu Agarwal<sup>1</sup>, Colby Ham<sup>1</sup>, Pritam Mukherjee<sup>2</sup>, Siyi Tang<sup>3</sup>, Sindhu Tipirneni<sup>3</sup>, Chandan Reddy<sup>4</sup>, Suzanne Tamang<sup>2</sup>, Robert Rallo<sup>1</sup>, Veysel Kocaman<sup>7</sup>,  
<sup>1</sup>Pacific Northwest National Laboratory; <sup>2</sup>Stanford University; <sup>3</sup>Virginia 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 studies<sup>1</sup>, 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-



## 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)



Health  
Intelligence  
(W3PHIAI-22)  
at AAAI

### 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,

# TRUSTED BY



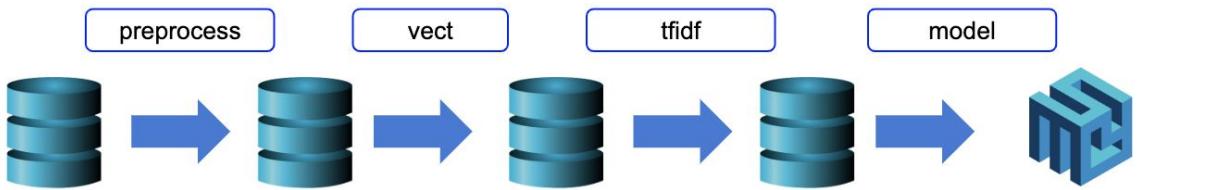
Imperial College  
London



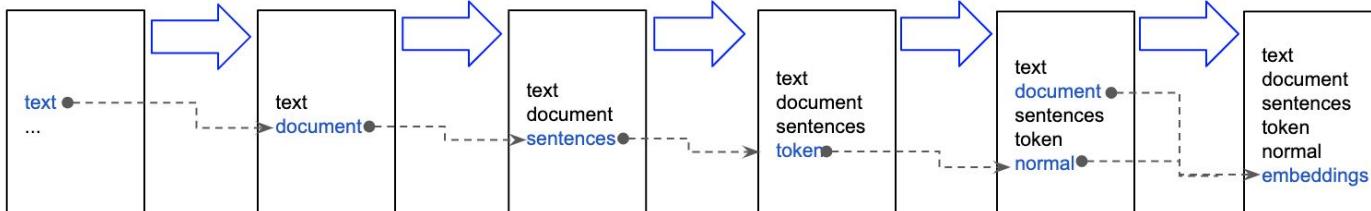
STANFORD  
UNIVERSITY

# Introducing Spark NLP

## 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)
```

# Introducing Spark NLP



Faster inference

```
from sparknlp.base import LightPipeline  
LightPipeline(someTrainedPipeline).annotate(someStringOrArray)
```

Spark is like a [locomotive](#) racing a [bicycle](#). The [bike](#) will win if the load is light, it is quicker to accelerate and more agile, but with a heavy load the [locomotive](#) might take a while to get up to speed, but [it's](#) going to be faster in the end.

**LightPipelines** are Spark ML pipelines converted into a single machine but multithreaded task, becoming more than 10x times faster for smaller amounts of data (small is relative, but 50k sentences is roughly a good maximum).

# Spark 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, ...

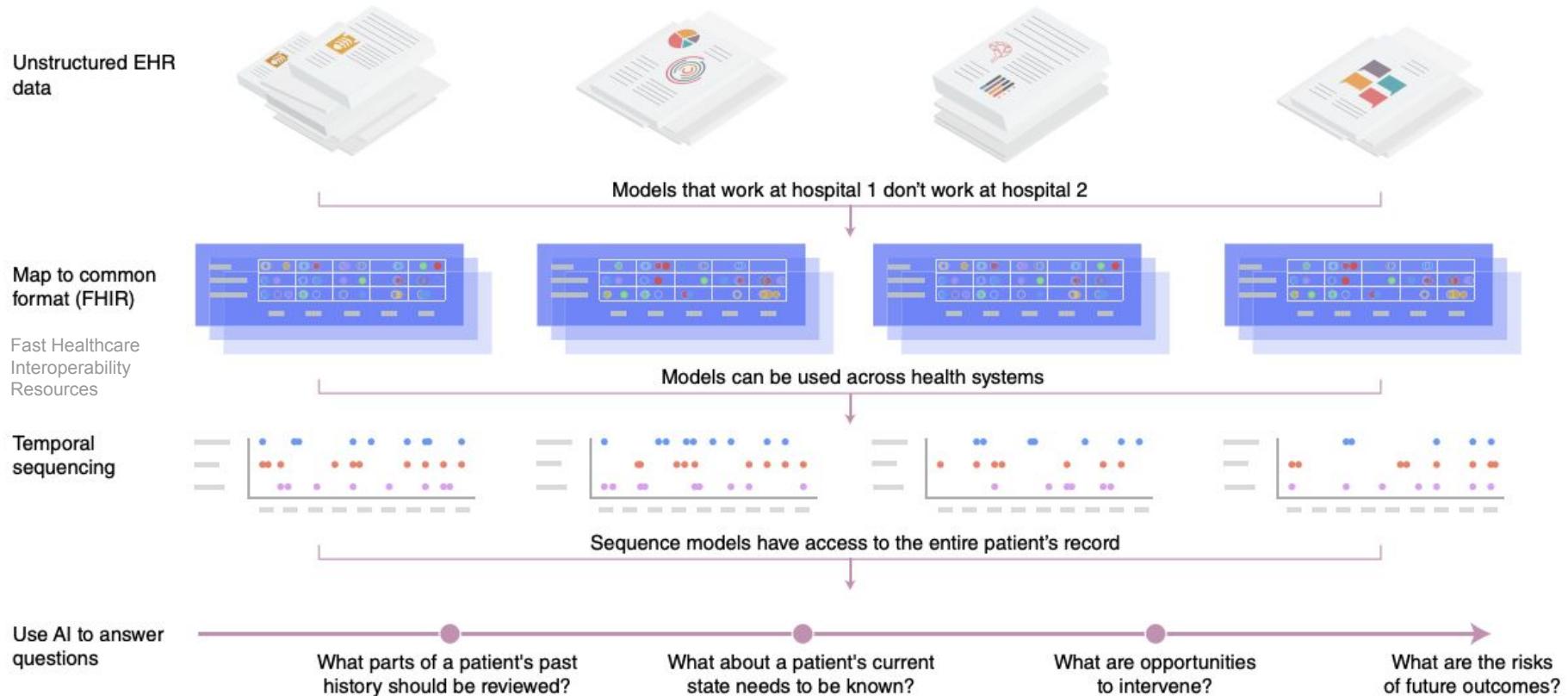
### 3. ML & DL Models

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."



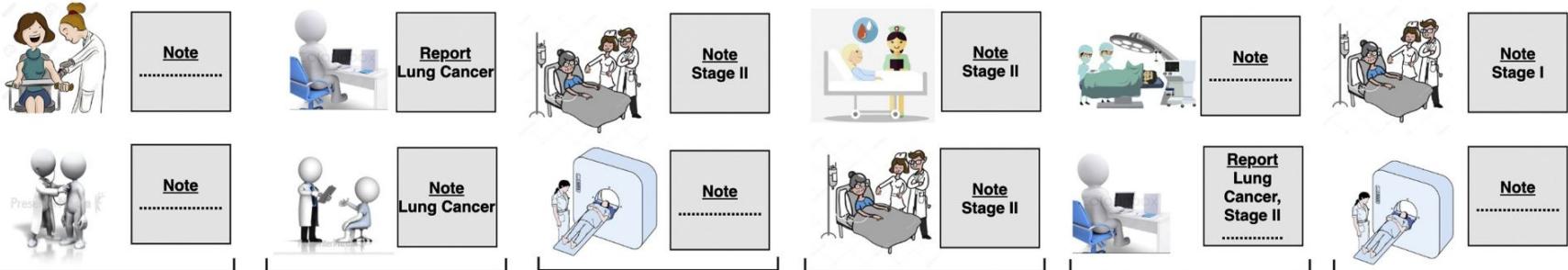
Features	
Type of Pain	Symptoms
Intensity of Pain	Onset of symptoms
Body part of region	Attempted home remedy



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

# Putting the clinical facts on a timeline

Natural History



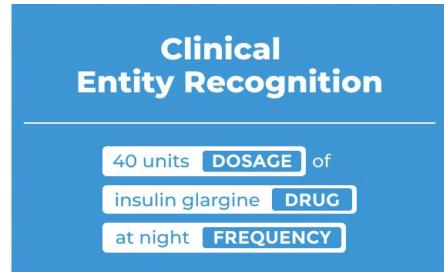
Medical Timeline

Lung Cancer  
Diagnosis

Tumor Stage II

Tumor Stage I

# 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

# 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

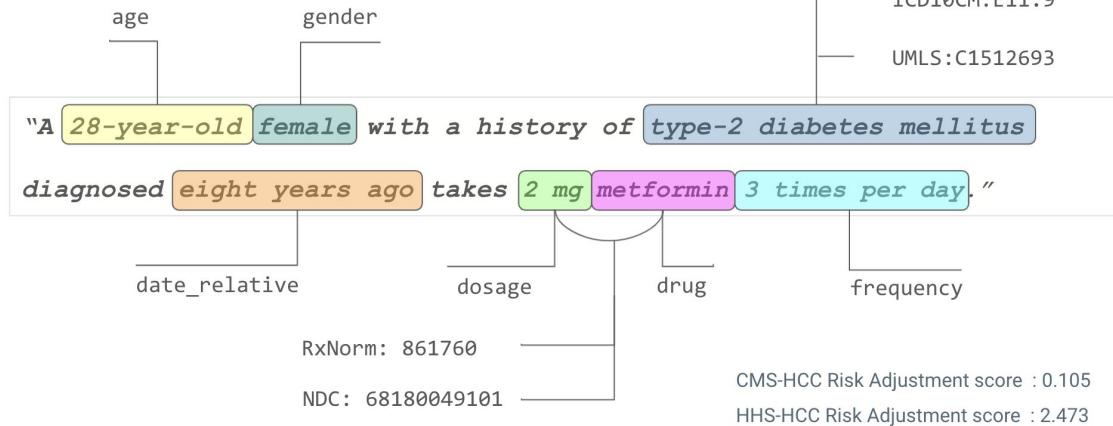
UMLS:C1512693

Sentence Splitter

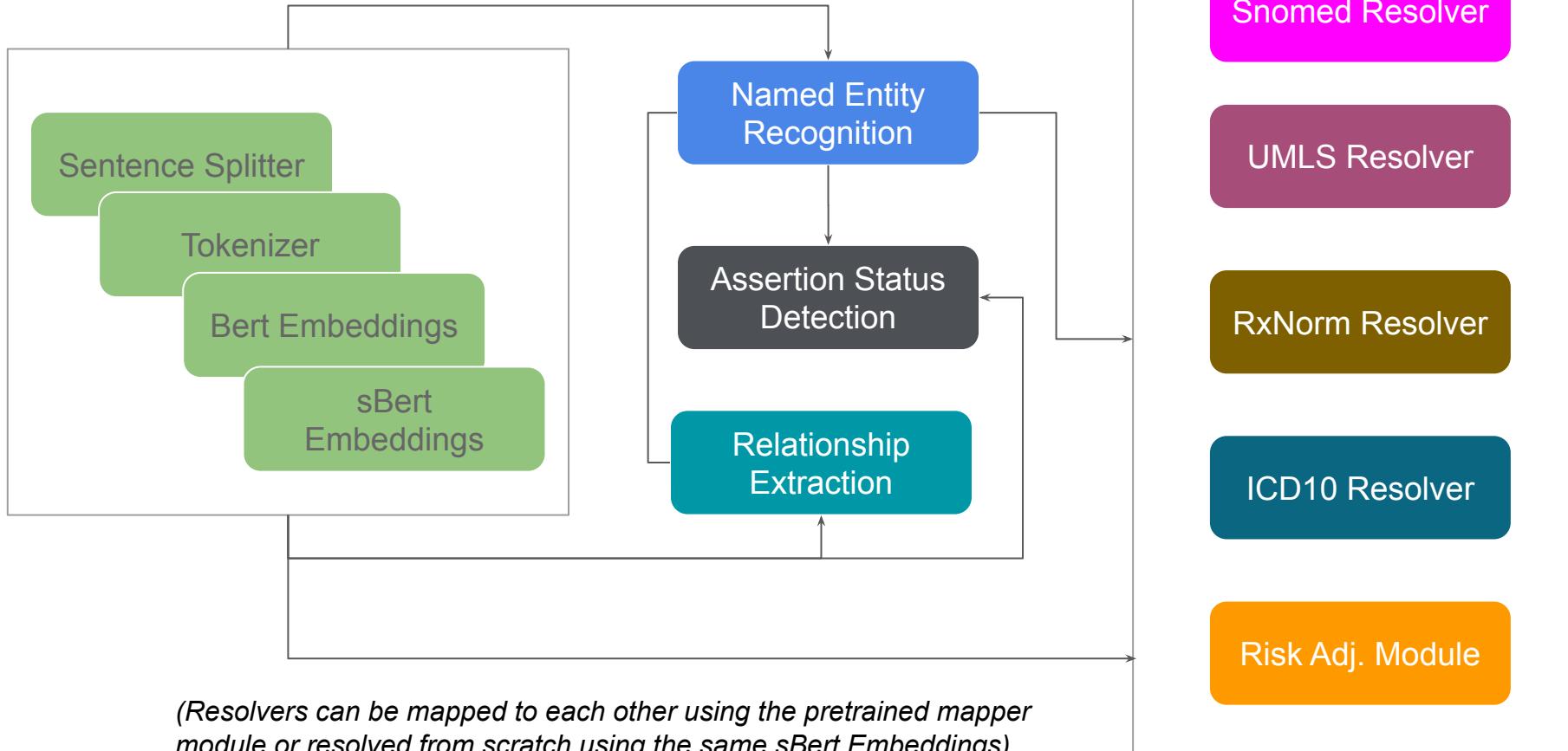
Tokenizer

Bert Embeddings

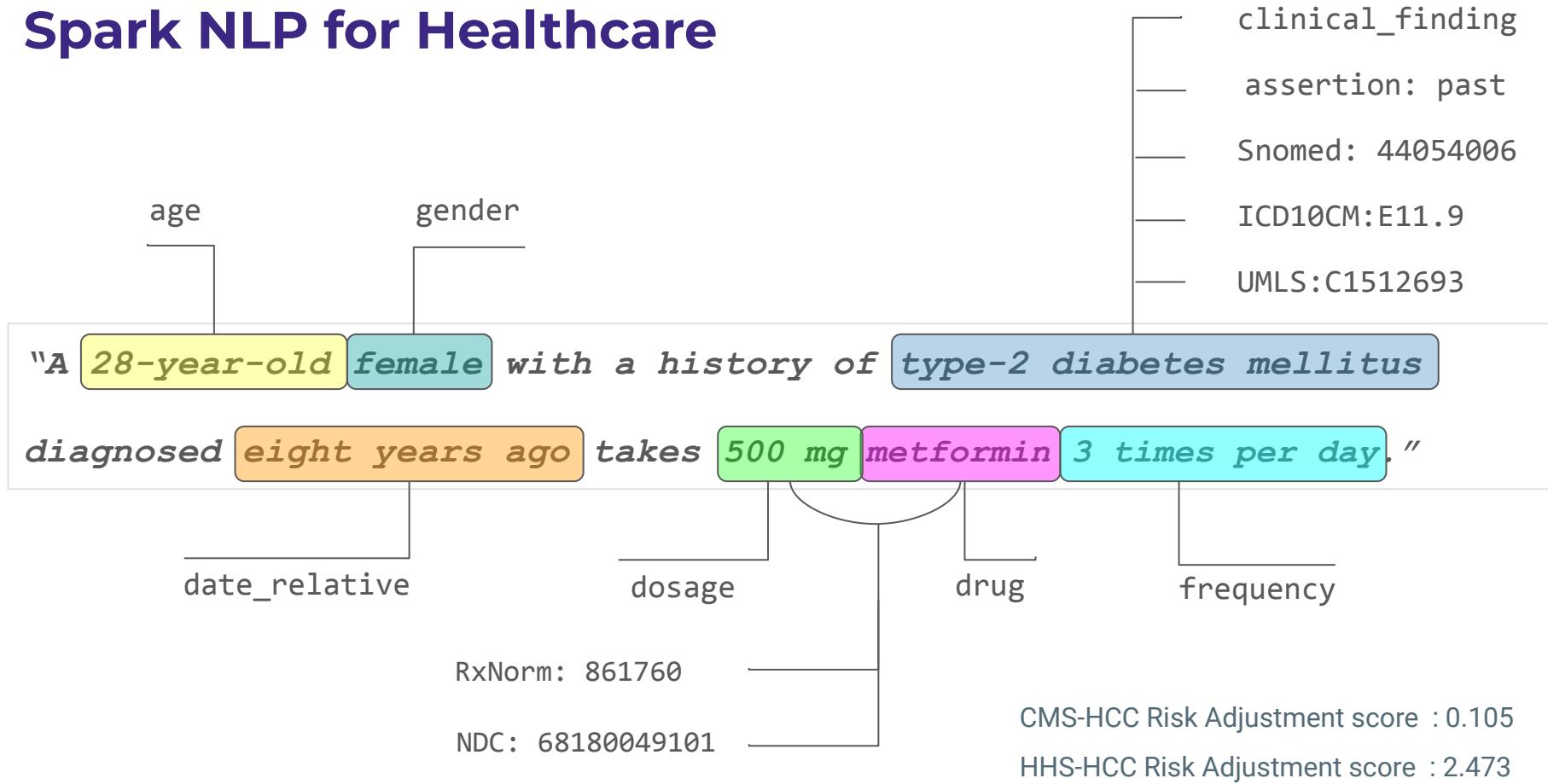
sBert Embeddings



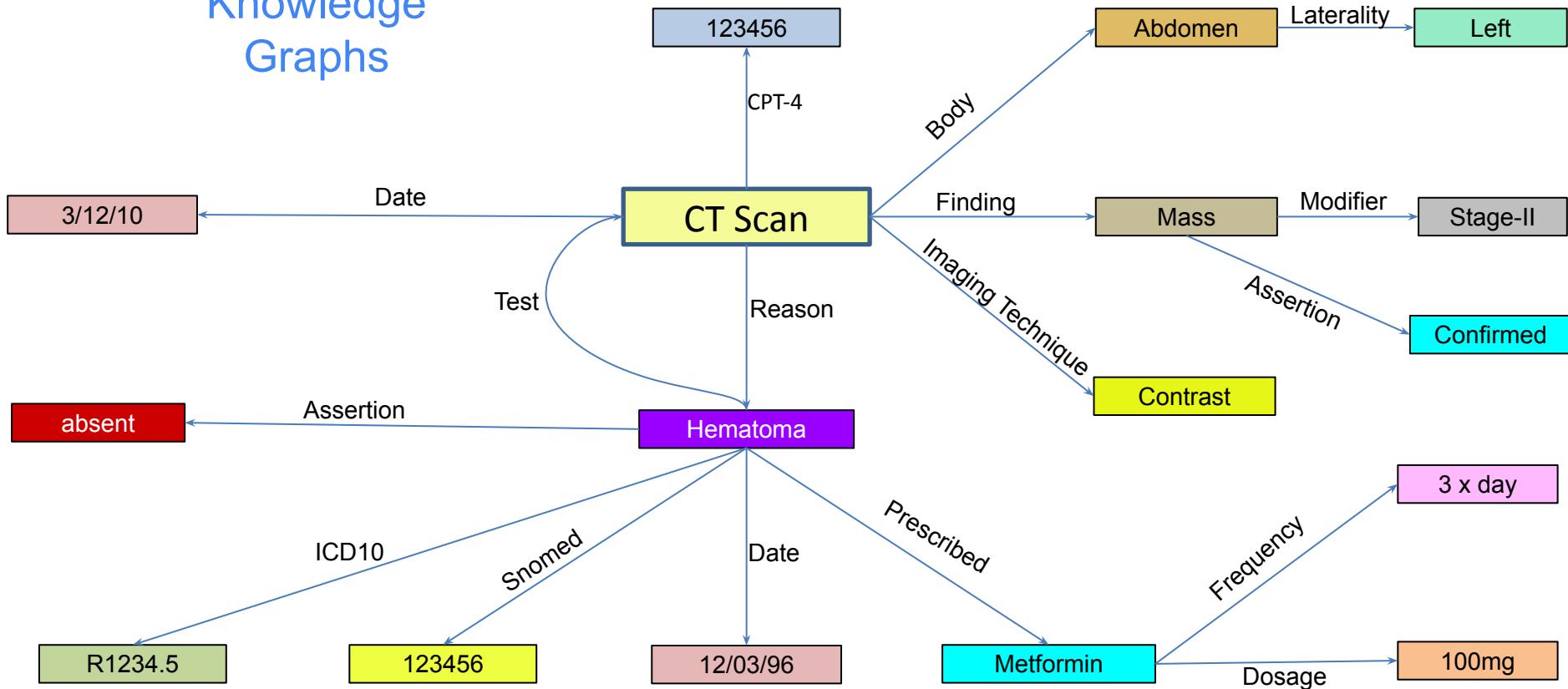
# Spark NLP for Healthcare

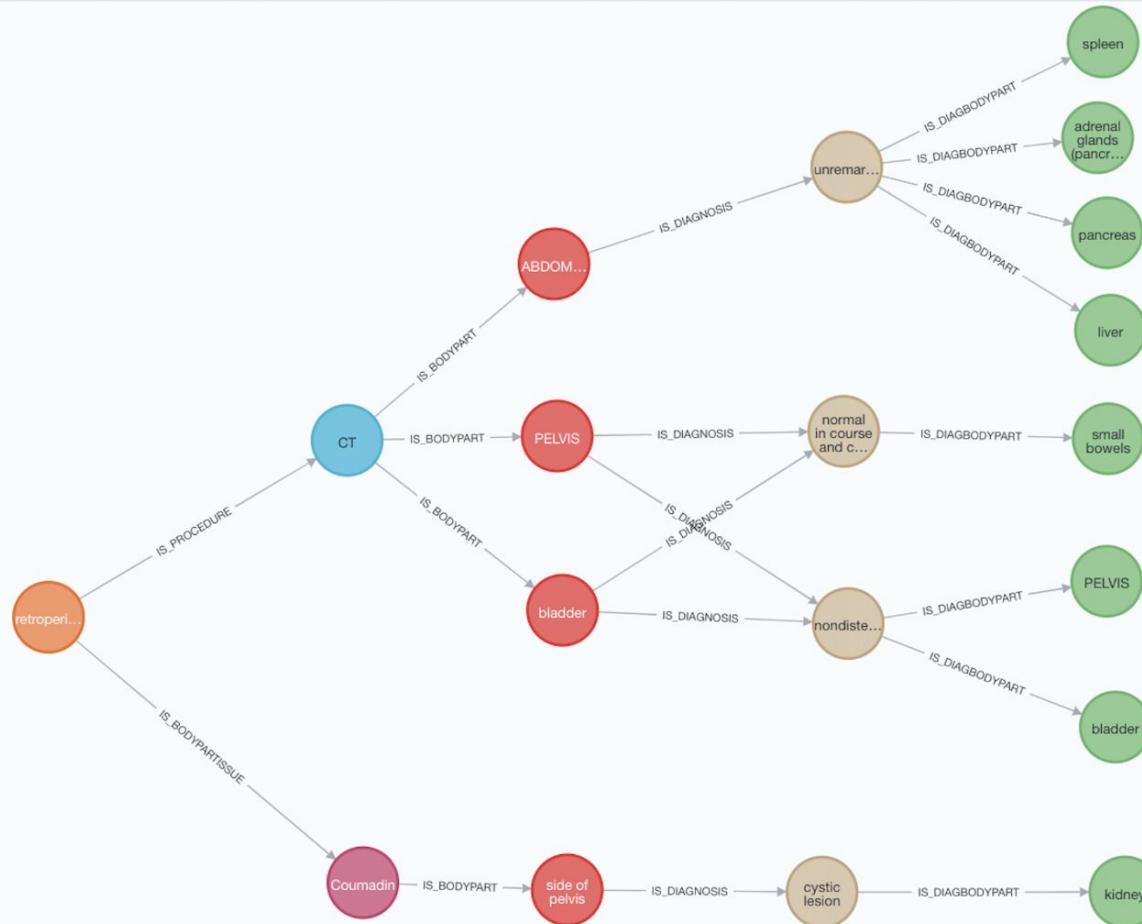


# Spark NLP for Healthcare



# Knowledge Graphs





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.

# NLP in Healthcare

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



Annotate your own data and train a custom NER model

Recently diagnosed, stage 4 adenocarcinoma of both lungs with metastasis to bone.  
CT scan shows no indication of mets on brain.

Extract named entities with Spark NLP *NERDL* and assign assertion statuses with *AssertionDL* model

Feature extraction & engineering

Prediction: Bone metastasis on June 2018



Text embeddings thru clinical word embeddings

- The frequency of clinical visits using document dates
- The number of positive site (organ) entities (hits by window)
- The number of Radiology/Oncology/Pathology reports in the last x days
- The number of tests applied in the last x days
- The number of diseases detected in the last x days
- Family history and social health determinants, etc.
- Date extraction and normalization

# Clinical Named Entity Recognition (NER)

- Extract structured data from free text
- Automate record keeping & abstraction process
- Feeding downstream tasks
- Features for ML models

Clinical Entity Recognition	Clinical Entity Linking	Assertion Status	Relation Extraction
40 units <b>DOSAGE</b> of insulin glargine <b>DRUG</b> at night <b>FREQUENCY</b>	Suspect diabetes SNOMED-CT: A73122005 Lisinopril 10 MG RxNorm: 316151 Hyponatremia ICD-10: E87.1	Fever and sore throat → PRESENT No stomach pain → ABSENT Father with Alzheimer → FAMILY	AFTER Admitted Occurrence for nausea Symptom due to chemo Treatment CAUSED BY
Algorithms		Content	
Extract Knowledge		Medical Transformers	Linked Medical Terminologies
<ul style="list-style-type: none"><li>• Entity Linker</li><li>• Entity Disambiguator</li><li>• Document Classifier</li><li>• Contextual Parser</li></ul>		<ul style="list-style-type: none"><li>• Structured Data</li><li>• Unstructured Text</li><li>• Obfuscator</li><li>• Generalizer</li></ul>	<ul style="list-style-type: none"><li>JSL-BERT-Clinical</li><li>BioBERT</li><li>ClinicalBERT</li><li>GloVe-Med</li><li>GloVe-ICD-O</li><li>BlueBERT</li></ul>
Split Text		Clean Medical Text	75+ Pretrained Models
<ul style="list-style-type: none"><li>• Sentence Detector</li><li>• Deep Sentence Detector</li><li>• Tokenizer</li><li>• nGram Generator</li></ul>		<ul style="list-style-type: none"><li>• Spell Checking</li><li>• Spell Correction</li><li>• Normalizer</li><li>• Stopword Cleaner</li></ul>	<p><b>Clinical:</b> Signs, Symptoms, Treatments, Procedures, Tests, Labs, Sections</p> <p><b>Anatomy:</b> Organ, Subdivision, Cell, Structure, Organism, Tissue, Gene, Chemical</p> <p><b>Drugs:</b> Name, Dosage, Strength, Route, Duration, Frequency, Poisons, Adverse Effects</p> <p><b>Risk Factors:</b> Smoking, Obesity, Diabetes, Hypertension, Substance Abuse</p> <p><b>Demographics:</b> Age, Gender, Height, Weight, Race, Ethnicity, Marital Status, Vital Signs</p> <p><b>Sensitive Data:</b> Patient Name, Address, Phone, Email, Dates, Providers, Identifiers</p>
Clinical Grammar		Find in Text	Trainable & Tunable Scalable to a Cluster Fast Inference Hardware Optimized Community
<ul style="list-style-type: none"><li>• Stemmer</li><li>• Lemmatizer</li><li>• Part of Speech Tagger</li><li>• Dependency Parser</li></ul>		<ul style="list-style-type: none"><li>• Text Matcher</li><li>• Regex Matcher</li><li>• Date Matcher</li><li>• Chunker</li></ul>	    

# Pretrained NER Models

- Clinical NER Models

index	model_name	index	model_name	index	model_name	index	model_name
1	jsl_ner_wip_clinical	17	ner_chexpert	33	ner_deid_subentity (German)	49	ner_jsl_greedy
2	jsl_ner_wip_greedy_clinical	18	ner_clinical	34	ner_diseases_large	50	ner_jsl_slim
3	jsl_ner_wip_modifier_clinical	19	ner_clinical_icdem	35	ner_drugs	51	ner_measurements_clinical
4	jsl_rd_ner_wip_greedy_clinical	20	ner_clinical_large	36	ner_drugs_greedy	52	ner_medmentions_coarse
5	ner_ade_clinical	21	ner_clinical_large_en	37	ner_drugs_large	53	ner_posology
6	ner_ade_clinicalalbert	22	ner_deid_augmented	38	ner_events_admission_clinical	54	ner_posology_experimental
7	ner_ade_healthcare	23	ner_deid_enriched	39	ner_events_clinical	55	ner_posology_greedy
8	ner_anatomy	24	ner_deid_generic_augmented	40	ner_events_healthcare	56	ner_posology_healthcare
9	ner_anatomy_coarse	25	ner_deid_generic (German)	41	ner_genetic_variants	57	ner_posology_large
10	ner_bacterial_species	26	ner_deid_large	42	ner_healthcare	58	ner_posology_small
11	ner_bionlp	27	ner_deid_sd	43	ner_human_phenotype_gene_clinical	59	ner_profiling_clinical
12	ner_cancer_genetics	28	ner_deid_sd_large	44	ner_human_phenotype_go_clinical	60	ner_radiology
13	ner_cellular	29	ner_deid_subentity_augmented	45	ner_jsl	61	ner_radiology_wip_clinical
14	ner_chemicals	30	ner_deid_synthetic	46	ner_jsl_enriched	62	ner_risk_factors
15	ner_chemprot_clinical	31	ner_deidentify_dl	47	ner_nihss	63	ner biomarker
16	ner_abbreviation_clinical	32	ner_deid_subentity_augmented_i2b2	48	ner_diseases	64	ner_drugprot_clinical

- BioBert NER Models

index	model_name	index	model_name	index	model_name	index	model_name
1	jsl_ner_wip_greedy_biobert	7	ner_cellular_biobert	13	ner_events_biobert	18	ner_jsl_greedy_biobert
2	jsl_rd_ner_wip_greedy_biobert	8	ner_chemprot_biobert	14	ner_human_phenotype_gene_biobert	19	ner_posology_biobert
3	ner_ade_biobert	9	ner_clinical_biobert	15	ner_human_phenotype_go_biobert	20	ner_posology_large_biobert
4	ner_anatomy_biobert	10	ner_deid_biobert	16	ner_jsl_biobert	21	ner_profiling_biobert
5	ner_anatomy_coarse_biobert	11	ner_deid_enriched_biobert	17	ner_jsl_enriched_biobert	22	ner_risk_factors_biobert
6	ner_bionlp_biobert	12	ner_diseases_biobert				

- BertForTokenClassification Clinical NER models

model_name
1 bert_token_classifier_ner_ade
2 bert_token_classifier_ner_clinical
3 bert_token_classifier_ner_deid
4 bert_token_classifier_ner_drugs
5 bert_token_classifier_ner_jsl
6 bert_token_classifier_ner_jsl_slim
7 bert_token_classifier_ner_bionlp
8 bert_token_classifier_ner_bacteria
9 bert_token_classifier_ner_anatomy
10 bert_token_classifier_ner_cellular
11 bert_token_classifier_ner_chemprot
12 bert_token_classifier_ner_chemicals
13 bert_token_classifier_drug_development_trials

Approach	embeddings	# of models
BiLSTM-CNN-Char	Clinical (glove)	70+
BiLSTM-CNN-Char	Biobert	20+
Bert for Token Cls.	Biobert	10+
Total		100+

# NER JSL

Let's show an example of `ner_js1` model that has about 80 clinical entity labels by changing just only the model name.

## Entities

Injury_or_Poisoning	Direction	Test	Admission_Discharge	Death_Entity
Relationship_Status	Duration	Respiration	Hyperlipidemia	Birth_Entity
Age	Labour_Delivery	Family_History_Header	BMI	Temperature
Alcohol	Kidney_Disease	Oncological	Medical_History_Header	Cerebrovascular_Disease
Oxygen_Therapy	O2_Saturation	Psychological_Condition	Heart_Disease	Employment
Obesity	Disease_Syndrome_Disorder	Pregnancy	ImagingFindings	Procedure
Medical_Device	Race_Ethnicity	Section_Header	Symptom	Treatment
Substance	Route	Drug_Ingredient	Blood_Pressure	Diet
External_body_part_or_region	LDL	VS_Finding	Allergen	EKG_Findings
Imaging_Technique	Triglycerides	RelativeTime	Gender	Pulse
Social_History_Header	Substance_Quantity	Diabetes	Modifier	Internal_organ_or_component
Clinical_Dept	Form	Drug_BrandName	Strength	Fetus_NewBorn
RelativeDate	Height	Test_Result	Sexually_Active_or_Sexual_Orientation	Frequency
Time	Weight	Vaccine	Vital_Signs_Header	Communicable_Disease
Dosage	Overweight	Hypertension	HDL	Total_Cholesterol
Smoking	Date			

# 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	<b>0.814</b>	0.862	0.711	<b>0.88</b>	0.71
ner_jsl_slim	0.87	0.766	0.86	<b>0.778</b>	<b>0.89</b>	0.75
ner_deid	<b>0.94</b>	<b>0.77</b>	0.93	0.77	0.75	0.63
ner_drug	0.964	0.964	0.912	0.911	<b>1</b>	<b>0.98</b>
ner_ade	0.84	0.807	0.839	0.819	<b>0.89</b>	<b>0.84</b>

\* 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.

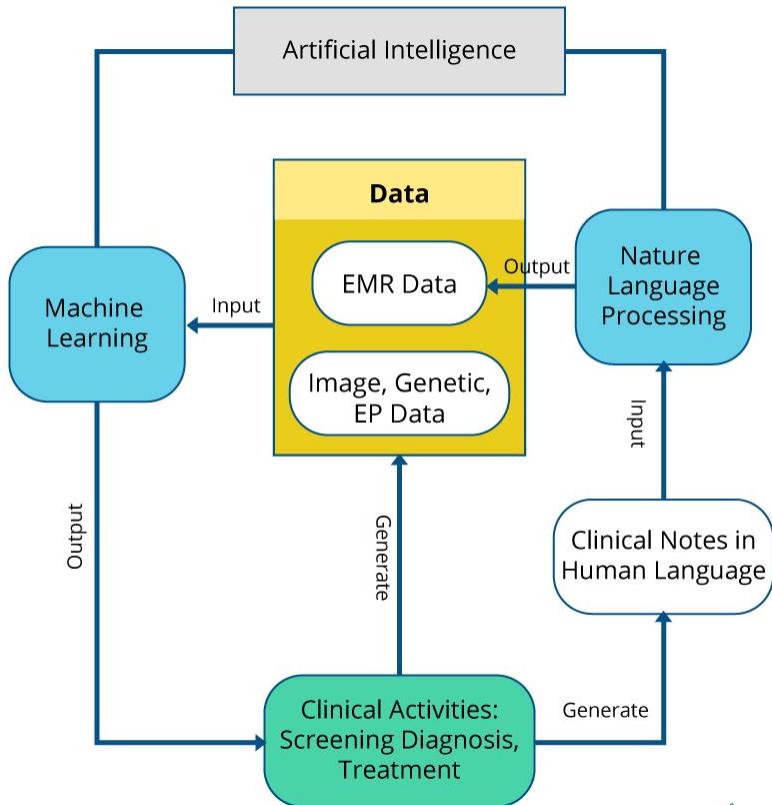
# NLP in Healthcare

"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** .

# Clinical Named Entity Recognition (NER)



The patient was prescribed 1 capsule of Advil for 5 days . He was seen by the endocrinology service and she was discharged on 40 units of insulin glargine at night , 12 units of insulin lispro with meals , and metformin 1000 mg two times a day . It was determined that all SGLT2 inhibitors should be discontinued indefinitely fro 3 months .

Color codes:FREQUENCY, DOSAGE, DURATION, DRUG, FORM, STRENGTH, **Posology NER**

No findings in urinary system , skin color is normal , brain CT and cranial checks are clear . Swollen fingers and eyes . Extensive stage small cell lung cancer . Chemotherapy with carboplatin and etoposide . Left scapular pain status post CT scan of the thorax .

Color codes:Organ, Organism\_subdivision, Organism\_substance, PathologicalFormation, Anatomical\_system,

**Anatomy NER**

A . Record date : 2093-01-13 , David Hale , M.D . , Name : Hendrickson , Ora MR . # 7194334 Date : 01/13/93 PCP : Oliveira , 25 years-old , Record date : 2079-11-09 . Cocke County Baptist Hospital . 0295 Keats Street

Color codes:STREET, DOCTOR, AGE, HOSPITAL, PATIENT, DATE, MEDICALRECORD,

**PHI NER**

# Clinical Named Entity Recognition (NER)

A 28-year-old female with a history of gestational diabetes mellitus diagnosed eight years prior to presentation and subsequent type two diabetes mellitus (T2DM), one prior episode of HTG-induced pancreatitis three years prior to presentation, and associated with an acute hepatitis, presented with a one-week history of polyuria, poor appetite, and vomiting. She was on metformin, glipizide, and dapagliflozin for T2DM and atorvastatin and gemfibrozil for HTG. She had been on dapagliflozin for six months at the time of presentation. Physical examination on presentation was significant for dry oral mucosa ; significantly , her abdominal examination was benign with no tenderness, guarding, or rigidity. Pertinent laboratory findings on admission were: serum glucose 111 mg/dL, creatinine 0.4 mg/dL, triglycerides 508 mg/dL, total cholesterol 122 mg/dL, and venous pH 7.27.

D

A 28-year-old female with a history of gestational diabetes mellitus diagnosed eight years prior to presentation and subsequent type two diabetes mellitus (T2DM), one prior episode of HTG-induced pancreatitis three years prior to presentation, and associated with an acute hepatitis, presented with a one-week history of polyuria, poor appetite, and vomiting. She was on metformin, glipizide, and dapagliflozin for T2DM and atorvastatin and gemfibrozil for HTG. She had been on dapagliflozin for six months at the time of presentation. Physical examination on presentation was significant for dry oral mucosa ; significantly , her abdominal examination was benign with no tenderness, guarding, or rigidity. Pertinent laboratory findings on admission were: serum glucose 111 mg/dL, creatinine 0.4 mg/dL, triglycerides 508 mg/dL, total cholesterol 122 mg/dL, and venous pH 7.27.

ner\_clinical

A 28-year-old female with a history of gestational diabetes mellitus diagnosed eight years prior to presentation and subsequent type two diabetes mellitus (T2DM), one prior episode of HTG-induced pancreatitis three years prior to presentation, and associated with an acute hepatitis, presented with a one-week history of polyuria, poor appetite, and vomiting. She was on metformin, glipizide, and dapagliflozin for T2DM and atorvastatin and gemfibrozil for HTG. She had been on dapagliflozin for six months at the time of presentation. Physical examination on presentation was significant for dry oral mucosa ; significantly , her abdominal examination was benign with no tenderness, guarding, or rigidity. Pertinent laboratory findings on admission were: serum glucose 111 mg/dL, creatinine 0.4 mg/dL, triglycerides 508 mg/dL, total cholesterol 122 mg/dL, and venous pH 7.27.

ner\_jsl

# Clinical Named Entities – Spark NLP vs Others

## Spark NLP

Google	Azure	AWS
PROBLEM	DIAGNOSIS SYMPTOM_OR_SIGN ALLERGEN	MEDICAL_CONDITION_DX_NA MEDICAL_CONDITION_SIGN MEDICAL_CONDITION_SYMPTOM
PROCEDURE	TREATMENT_NAME	PROCEDURE_NAME TREATMENT_NAME
MEDICINE	MEDICATION_CLASS MEDICATION_NAME	MEDICATION_BRAND_NAME MEDICATION_GENERIC_NAME
ANATOMICAL_STRUCTURE	BODY_STRUCTURE	SYSTEM_ORGAN_SITE
LABORATORY_DATA BODY_MEASUREMENT	EXAMINATION_NAME	TEST_NAME
SEVERITY	CONDITION_QUALIFIER CONDITION_SCALE	MEDICAL_CONDITION_ACUITY
MED_DOSE MED_TOTALDOSE MED_STRENGTH MED_UNIT	DOSAGE	MEDICATION_DOSAGE MEDICATION_STRENGTH MEDICATION_RATE
MED_FREQUENCY	FREQUENCY	MEDICATION_FREQUENCY
MED_FORM	MEDICATION_FORM	MEDICATION_FORM
MED_ROUTE	MEDICATION_ROUTE	MEDICATION_ROUTE_OR_MODE
MED_DURATION	TIME	MEDICATION_DURATION
LAB_VALUE MED_VALUE	MEASUREMENT_VALUE	TEST_VALUE
LAB_UNIT BM_UNIT	MEASUREMENT_UNIT	TEST_UNIT



Symptom, Disease\_Syndrome\_Disorder, Heart\_Disease, VS\_Finding, Communicable\_Disease, Hypertension, Diabetes, Kidney\_Disease, Cerebrovascular\_Disease, Injury\_or\_Poisoning, Psychological\_Condition, Total\_Cholesterol, Hyperlipidemia, Obesity, Oncological, Pregnancy, EKG\_Findings, Death\_Entity, ImagingFindings, Female\_Reproductive\_Status, Fetus\_NewBorn, Pregnancy\_Delivery\_Puerperium, Overweight, Puerperium



Test, Test\_Result, Treatment, Pulse, Imaging\_Technique, Labour\_Delivery, Temperature, Blood\_Pressure, Oxygen\_Therapy, Weight, LDL, O2\_Saturation, BMI, Vaccine, Respiration, Triglycerides

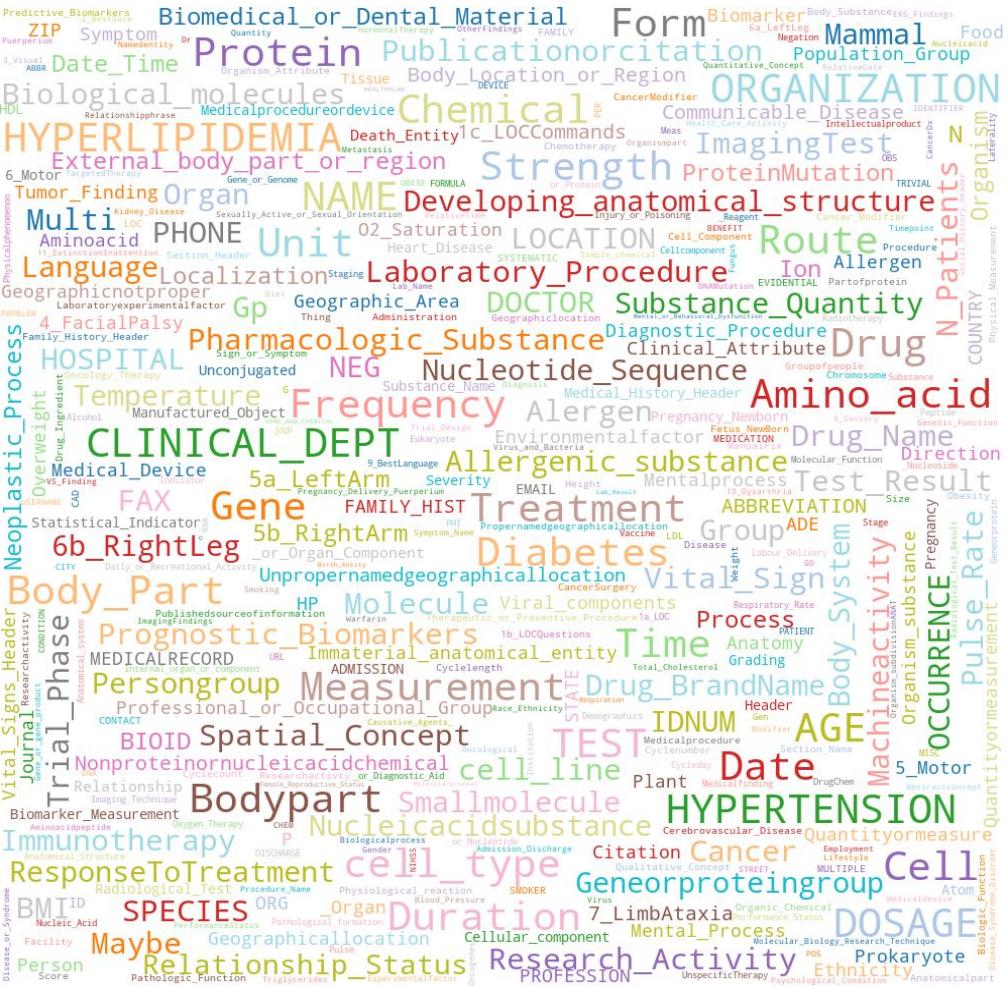
Mother with a lung cancer ONCOLOGICAL , a pregnant PREGNANCY patient is diagnosed as breast cancer ONCOLOGICAL in 1991 DATE and then admitted ADMISSION\_DISCHARGE to Mayo Clinic CLINICAL\_DEPT in Oct 2000 DATE , went under chemo TREATMENT for 6 months DURATION , discharged ADMISSION\_DISCHARGE in April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

# Clinical Named Entities

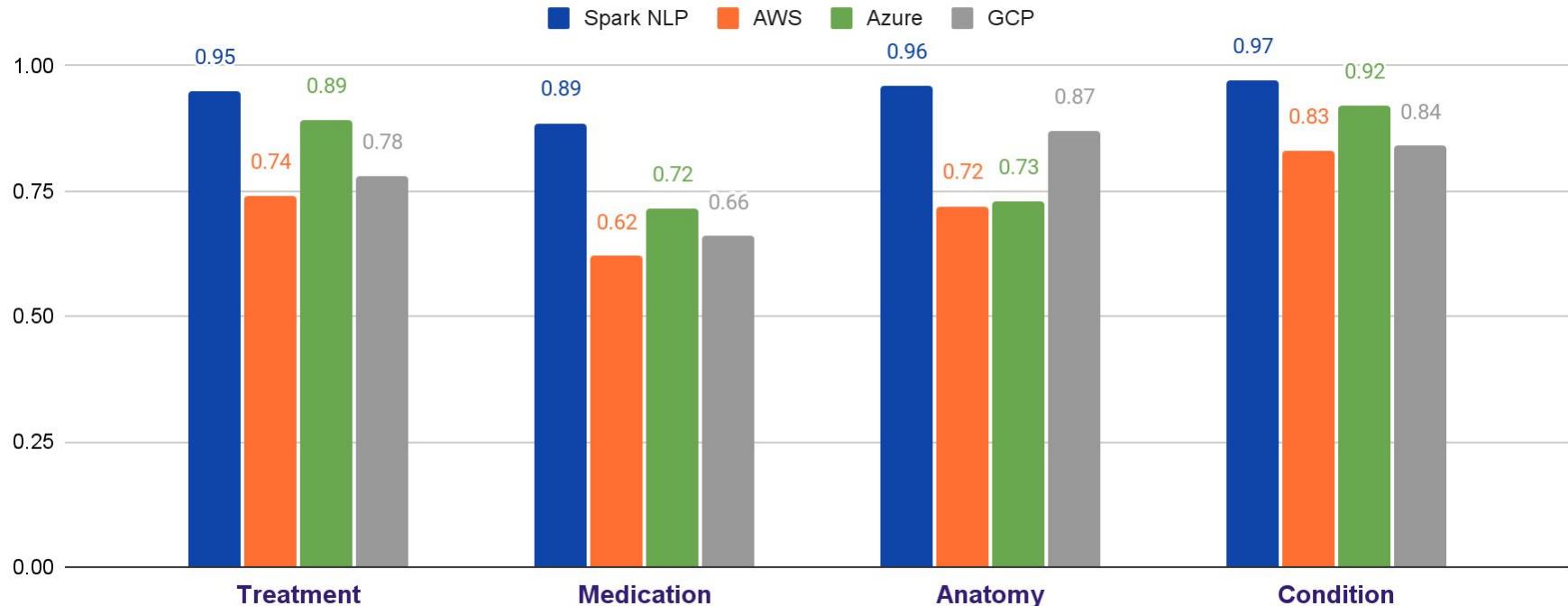
## Spark NLP vs Others

Google	Azure	AWS
PROBLEM	DIAGNOSIS SYMPTOM_OR_SIGN ALLERGEN	MEDICAL_CONDITION_DX_NA MEDICAL_CONDITION_SIGN MEDICAL_CONDITION_SYMPTOM
PROCEDURE	TREATMENT_NAME	PROCEDURE_NAME TREATMENT_NAME
MEDICINE	MEDICATION_CLASS MEDICATION_NAME	MEDICATION_BRAND_NAME MEDICATION_GENERIC_NAME
ANATOMICAL_STRUCTURE	BODY_STRUCTURE	SYSTEM_ORGAN_SITE
LABORATORY_DATA BODY_MEASUREMENT	EXAMINATION_NAME	TEST_NAME
SEVERITY	CONDITION_QUALIFIER CONDITION_SCALE	MEDICAL_CONDITION_ACUITY
MED_DOSE MED_TOTALDOSE MED_STRENGTH MED_UNIT	DOSAGE	MEDICATION_DOSAGE MEDICATION_STRENGTH MEDICATION_RATE
MED_FREQUENCY	FREQUENCY	MEDICATION_FREQUENCY
MED_FORM	MEDICATION_FORM	MEDICATION_FORM
MED_ROUTE	MEDICATION_ROUTE	MEDICATION_ROUTE_OR_MODE
MED_DURATION	TIME	MEDICATION_DURATION
LAB_VALUE MED_VALUE	MEASUREMENT_VALUE	TEST_VALUE
LAB_UNIT BM_UNIT	MEASUREMENT_UNIT	TEST_UNIT

# 400+ entities from 100+ models



# NER Benchmarks



# Spark NLP vs AWS vs GCP vs Academic

		Spark NLP	Competition Best	Last Best
Clinical Concept Extraction	2010 i2b2/VA	<b>0.876</b>	0.852	0.862
De-Identification	2014 n2c2	<b>0.961</b>	0.936	0.955
Medication Extraction	2018 n2c2	<b>0.899</b>	0.896	0.896

Entity	Sample	Spark NLP Clinical Models			AWS Medical Comprehend			GCP Healthcare API		
		Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Problem	4891	0.726	0.585	<b>0.648</b>	0.539	0.478	<b>0.507</b>	0.850	0.516	0.642
Test	5903	0.782	0.662	<b>0.717</b>	0.594	0.703	<b>0.644</b>	0.576	0.461	0.512
Drug	10284	0.946	0.882	0.913	0.815	0.910	<b>0.860</b>	0.962	0.885	<b>0.922</b>
Avg. F1				<b>0.759</b>			<b>0.670</b>			0.692

# Biomedical Named Entity Recognition

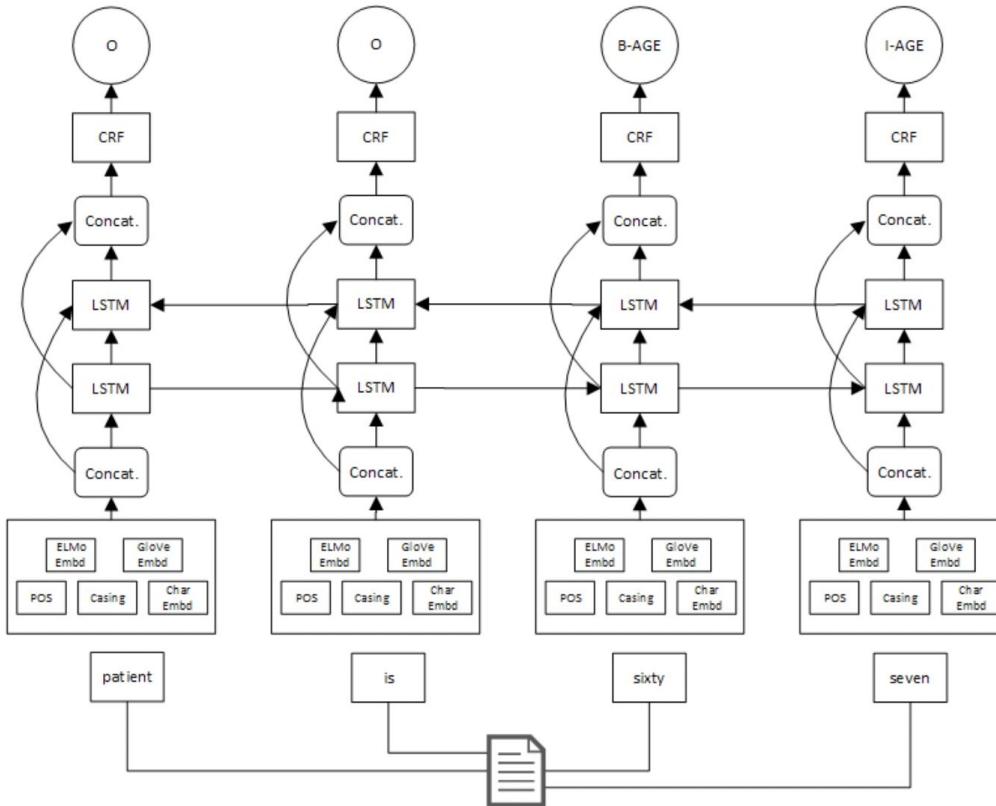
## Spark NLP vs Spacy vs Stanza

Dataset	Entities	Spark - Biomedical	Spark - GloVe 6B	Stanza	SciSpacy
NBCI-Disease	Disease	<b>89.13</b>	87.19	87.49	81.65
BC5CDR	Chemical, Disease	<b>89.73</b>	88.32	88.08	83.92
BC4CHEMD	Chemical	<b>93.72</b>	92.32	89.65	84.55
Linnaeus	Species	86.26	85.51	<b>88.27</b>	81.74
Species800	Species	<b>80.91</b>	79.22	76.35	74.06
JNLPBA	5 types in cellular	<b>81.29</b>	79.78	76.09	73.21
AnatEM	Anatomy	<b>89.13</b>	87.74	88.18	84.14
BioNLP13-CG	16 types in Cancer Genetics	<b>85.58</b>	84.3	84.34	77.6

Benchmarks on BioMedical NER Datasets

# NER Architecture

Label Prediction  
CRF Layer  
Bi-LSTM Layer  
Embedding Layer  
Sentence Input  
Document Input



Char-CNN-BiLSTM

John	B-PER
Smith	I-PER
lives	0
in	0
New	B-LOC
York	I-LOC

John Smith  $\Rightarrow$  PERSON  
New York  $\Rightarrow$  LOCATION

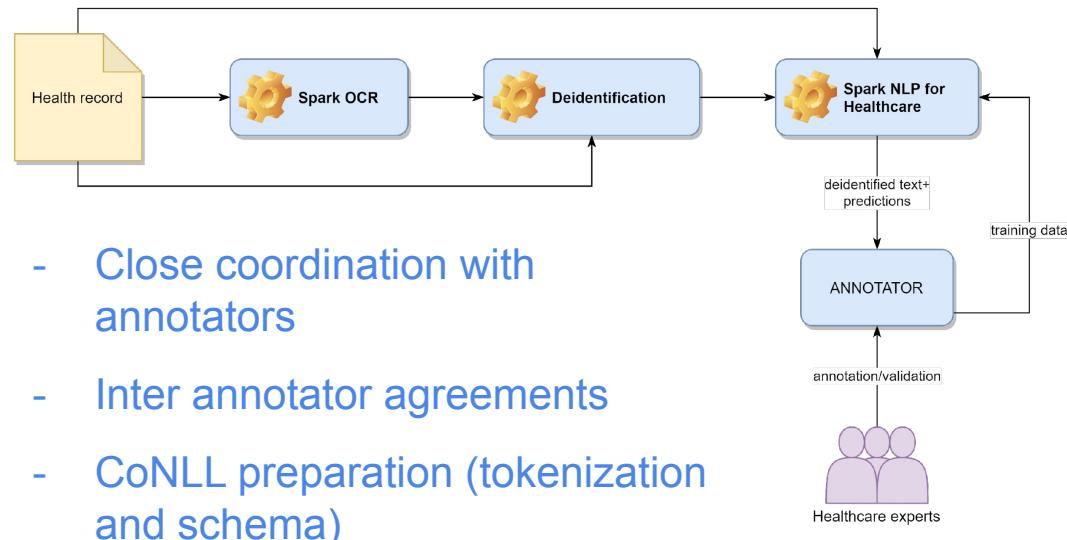
word	POS_tag	chunk_tag	NER_tag
She	PRP	O	B-person
presented	VBD	B-VP	O
with	IN	B-VP	O
left	JJ	B-NP	B-problem
upper	JJ	I-NP	I-problem
quadrant	NN	I-NP	I-problem
pain	NN	I-NP	I-problem
as	RB	O	O
well	RB	O	O
as	IN	B-VP	O
nausea	NN	B-NP	B-problem

John	B-PER
Smith	I-PER
lives	O
in	O
New	B-LOC
York	I-LOC

John Smith  $\Rightarrow$  PERSON  
 New York  $\Rightarrow$  LOCATION

word	POS_tag	chunk_tag	NER_tag
She	PRP	O	B-person
presented	VBD	B-VP	O
with	IN	B-VP	O
left	JJ	B-NP	B-problem
upper	JJ	I-NP	I-problem
quadrant	NN	I-NP	I-problem
pain	NN	I-NP	I-problem
as	RB	O	O
well	RB	O	O
as	IN	B-VP	O
nausea	NN	B-NP	B-problem

# NER in Healthcare



- Close coordination with annotators
- Inter annotator agreements
- CoNLL preparation (tokenization and schema)

She returns today for ongoing evaluation of her EGFR mutated, stage 4 lung cancer with metastasis to her L2 vertebrae and her lungs bilaterally.

Bone negative for metastatic disease.

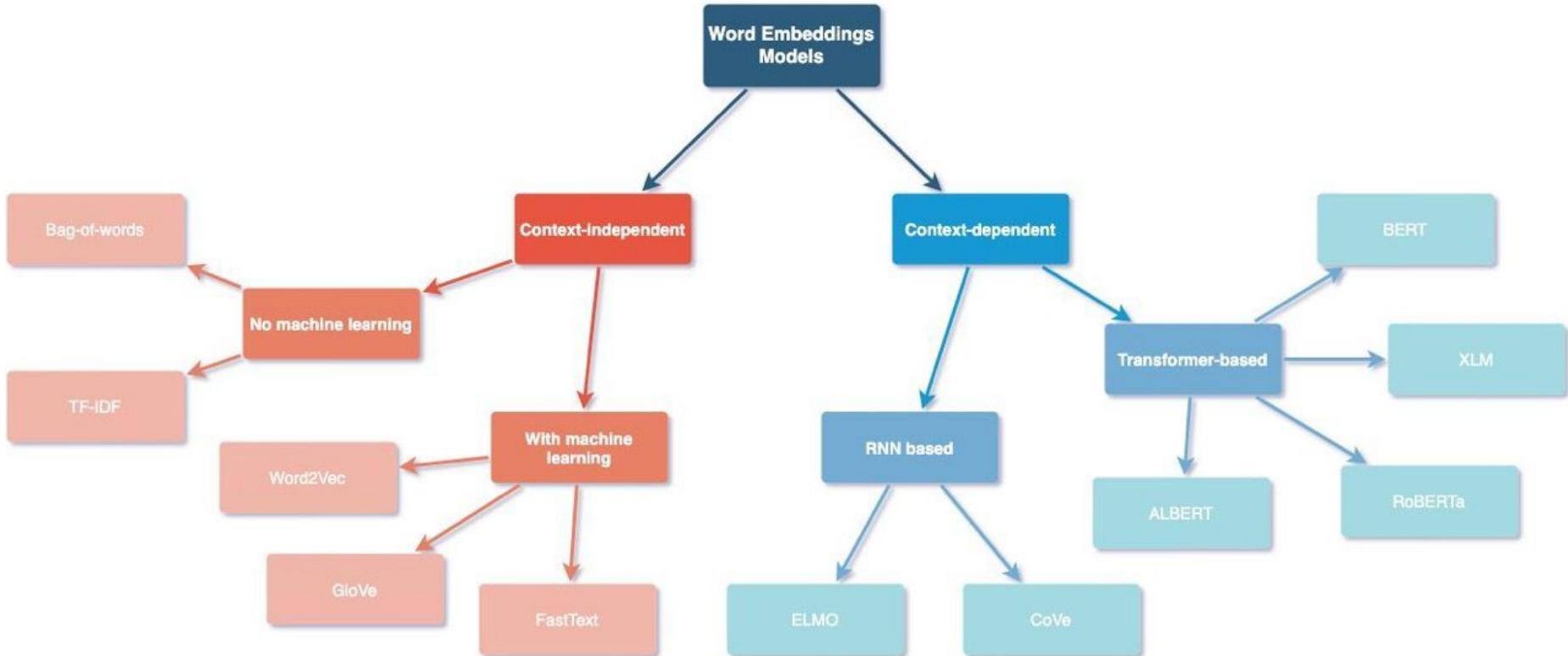
Patient denies any family history of cancer.

# NER-DL in Spark NLP

## Char-CNN-BiLSTM

	F1 : Tokens	F2 : Casing	F3 : POS	F4 : Char CNN	Labels
The					O
company					O
XYZ					Company
Private					Company
Limited					Company
works					O
in					O
the					O
health					Activity
sector					Activity
in					O
Europe					Location

# Clinical Word/Sentence Embeddings



# AN NLP TIMELINE AND THE TRANSFORMER FAMILY

## BAG OF WORDS (BOW)

Count the occurrences of each word in the documents and use them as features.

1954

## TF-IDF

The BOW scores are modified so that rare words have high scores and common words have low scores.

1972

## WORD2VEC

Each word is mapped to a high-dimensional vector called word embedding, which captures its semantic. Word embeddings are learned by a neural network looking for word correlations on a large corpus.

2013

## RNN

RNNs compute document embeddings leveraging word context in sentences, which was not possible with word embeddings alone.

### LSTM

Capture long term dependencies.

1997

### Bidirectional RNN

Capture left-to-right and right-to-left dependencies.

1997

### Encoder-decoder RNN

An RNN creates a document embedding (i.e. the encoder) and another RNN decodes it into text (i.e. the decoder).

2014

1986

## TRANSFORMER

An encoder-decoder model that leverages attention mechanisms to compute better embeddings and to better align output to input.

2017

## BERT

Bidirectional Transformer pretrained using a combination of Masked Language Modeling and Next Sentence Prediction objectives. It uses global attention.

2018

## GPT

The first autoregressive model based on the Transformer architecture.

### GPT-2

A bigger and optimized version of GPT, pre-trained on WebText.

2019

### GPT-3

A bigger and optimized version of GPT-2, pre-trained on Common Crawl.

2020

2018

## CTRL

Similar to GPT but with control codes for conditional text generation.

2019

## TRANSFORMER-XL

It's an autoregressive Transformer which can reuse previously computed hidden-states to attend to longer context.

2019

## ALBERT

A lighter version of BERT, where (1) Next Sentence Prediction is replaced by Sentence Order Prediction, and (2) parameter-reduction techniques are used for lower memory consumption and faster training.

2019

## ROBERTA

Better version of BERT, where (1) the Masked Language Modeling objective is dynamic, (2) the Next Sentence Prediction objective is dropped, (3) the BPE tokenizer is employed, and (4) better hyperparameters are used.

2019

## XLM

Transformer pre-trained on a corpus of several languages using objectives like Causal Language Modeling, Masked Language Modeling, and Translation Language Modeling.

2019

## XLNET

Transformer-XL with a generalized autoregressive pre-training method that enables learning bidirectional dependences.

2019

## PEGASUS

A bidirectional encoder and a left-to-right decoder pre-trained with Masked Language Modeling and Gap Sentence Generation objectives.

2019

## DISTILBERT

Same as BERT but smaller and faster, while preserving over 95% of BERT's performances. Trained by distillation of the pre-trained BERT model.

2019

## XLM-ROBERTA

RoBERTa trained on a multilingual corpus with the Masked Language Modeling objective.

2019

## BART

A bidirectional encoder and a left-to-right decoder trained by corrupting text with an arbitrary noising function and learning a model to reconstruct the original text.

2019

## CONVBERT

Better version of BERT, where self-attention blocks are replaced with new ones that leverage convolutions to better model global and local context.

2019

## FUNNEL TRANSFORMER

A type of Transformer that gradually compresses the sequence of hidden states to a shorter one and hence reduces the computation cost.

2020

## REFORMER

A more efficient Transformer thanks to local-sensitive hashing attention, axial position encoding and other optimizations.

2020

## T5

A bidirectional encoder and a left-to-right decoder pre-trained on a mix of unsupervised and supervised tasks.

2020

## LONGFORMER

A Transformer model replacing the attention matrices with sparse matrices for higher training efficiency.

2020

## PROPHETNET

A Transformer model trained with the Future N-gram Prediction objective and with a novel self-attention mechanism.

2020

## ELECTRA

Same as BERT but lighter and better. The model is trained with the Replaced Token Detection objective.

2020

## SWITCH TRANSFORMER

A sparsely-activated expert Transformer model that aims to simplify and improve over Mixtute of Experts.

2021



NLPLANET

The community of  
NLP enthusiasts!



<https://www.linkedin.com/company/nlplanet>



<https://medium.com/nlplanet>



[@nlplanet\\_](https://twitter.com/nlplanet_)

By Fabio Chiusano

# Clinical Word/Sentence Embeddings

Clinical Glove  
(200d)

PubMed + PMC

ICDO Glove  
(200d)

PubMed + ICD10  
UMLS + MIMIC III

Sent BERT

BioBert finetuned on  
NLI and MedNLI

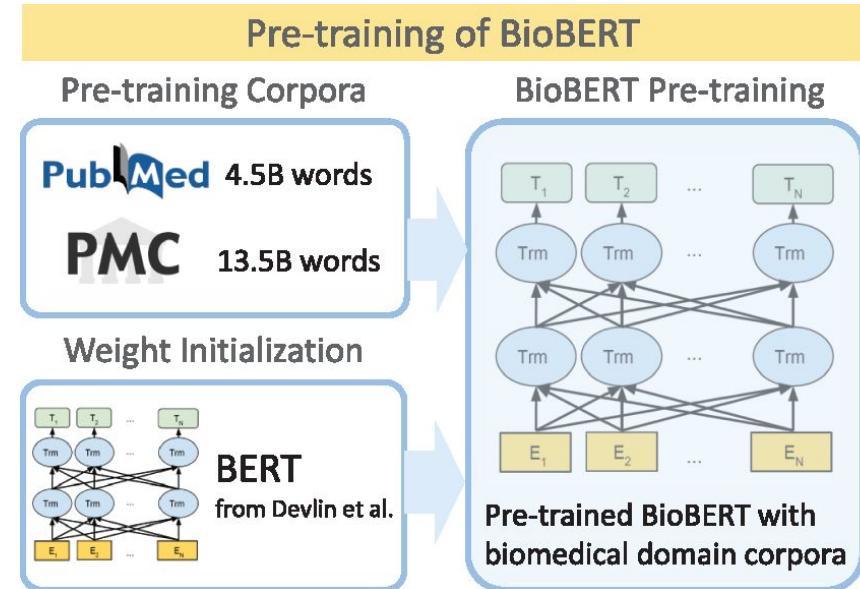
Bio/Clinical BERT

Fine tuned Pubmed + PMC + Discharge summaries



PubMed abstracts and PMC full-text articles

<https://www.nlm.nih.gov/bsd/difference.html>



# Part - II

- ❖ Assertion Status detection

# Assertion Status Detection

"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

```
clinical_assertion = AssertionDLModel\  
    .pretrained("assertion_dl", "en", "clinical/models")\  
    .setInputCols(["sentence", "ner_chunk", "embeddings"]) \  
    .setOutputCol("assertion")
```

Classify the assertions made on given medical concepts as being

- present,
- absent,
- possible,
- conditionally present under certain circumstances,
- hypothetically present at some future point, mentioned in the patient report but associated with someone else.

# Assertion Status Detection

- The deep neural network architecture for assertion status detection in Spark NLP is based on a Bi-LSTM framework, and is a modified version of the architecture proposed by Federico Fancellu, Adam Lopez and Bonnie Webber ([Neural Networks For Negation Scope Detection](#)).
- In the proposed implementation, input units depend on the target tokens (a named entity) and the neighboring words that are explicitly encoded as a sequence using word embeddings.
- Similar to paper mentioned above, it is observed that that 95% of the scope tokens (neighboring words) fall in a window of 9 tokens to the left and 15 to the right of the target tokens in the same dataset. Therefore, the same window size was implemented,
- following parameters were used: learning rate 0.0012, dropout 0.05, batch size 64 and a maximum sentence length 250.
- The model has been implemented within Spark NLP as an annotator called AssertionDLModel. After training 20 epoch and measuring accuracy on the official test set, this implementation exceeds the latest state-of-the-art accuracy benchmarks

Assertion Label	Spark NLP	Latest Best
Absent	0.944	0.937
Someone-else	0.904	0.869
Conditional	0.441	0.422
Hypothetical	0.862	0.890
Possible	0.680	0.630
Present	0.953	0.957
micro F1	0.939	0.934

Mother with a lung cancer,

# AssertionDLModel

	<b>model_name</b>	<b>Predicted Entities</b>
1	assertion_dl	Present, Absent, Possible, Planned, Someoneelse, Past, Family, None, Hypothetical
2	assertion_dl_biobert	absent, present, conditional, associated_with_someone_else, hypothetical, possible
3	assertion_dl_healthcare	absent, present, conditional, associated_with_someone_else, hypothetical, possible
4	assertion_dl_large	hypothetical, present, absent, possible, conditional, associated_with_someone_else
5	assertion_dl_radiology	Confirmed, Suspected, Negative
6	assertion_jsl	Present, Absent, Possible, Planned, Someoneelse, Past, Family, None, Hypothetical
7	assertion_jsl_large	present, absent, possible, planned, someoneelse, past
8	assertion_ml	Hypothetical, Present, Absent, Possible, Conditional, Associated_with_someone_else

# Spark NLP for Healthcare Data Scientists

July 20-21, 2022

**Veysel Kocaman**  
**Head of Data Science**  
[veysel@johnsnowlabs.com](mailto:veysel@johnsnowlabs.com)



# Part - III

- ❖ Entity Resolution (ICD1-, RxNorm, Snomed, etc.)

# Entity Resolution in Spark NLP for Healthcare

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

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

# Clinical Entity Resolution

## ICD10CM

- sbiobertresolve\_icd10cm\_augmented
- sbiobertresolve\_icd10pcs
- sbiobertresolve\_icd10cm\_augmented\_billable\_hcc
- sbiobertresolve\_icd10cm
- sbiobertresolve\_icd10cm\_slim\_normalized
- sbiobertresolve\_icd10cm\_slim\_billable\_hcc
- sbertrresolve\_icd10cm\_slim\_billable\_hcc\_med
- sbiobertresolve\_icd10cm\_generalised

## CPT

- sbiobertresolve\_cpt
- sbiobertresolve\_cpt\_procedures\_augmented
- sbiobertresolve\_cpt\_augmented
- sbiobertresolve\_cpt\_procedures\_measurements\_augmented

## Snomed

- sbiobertresolve\_snomed\_auxConcepts\_int
- sbiobertresolve\_snomed\_findings
- sbiobertresolve\_snomed\_findings\_int
- sbiobertresolve\_snomed\_auxConcepts
- sbertrsolve\_snomed\_bodyStructure\_med
- sbiobertresolve\_snomed\_bodyStructure
- sbiobertresolve\_snomed\_findings\_aux\_concepts
- sbertrsolve\_snomed\_conditions

## RxNorm

- sbiobertresolve\_rxnorm
- demo\_sbiobertresolve\_rxnorm
- sbiobertresolve\_rxnorm\_dispo
- sbiobertresolve\_rxnorm\_disposition
- sbertrsolve\_rxnorm\_disposition
- sbiobertresolve\_rxnorm\_ndc

## LOINC

- sbluebertresolve\_loinc
- sbiobertresolve\_loinc

and more ...

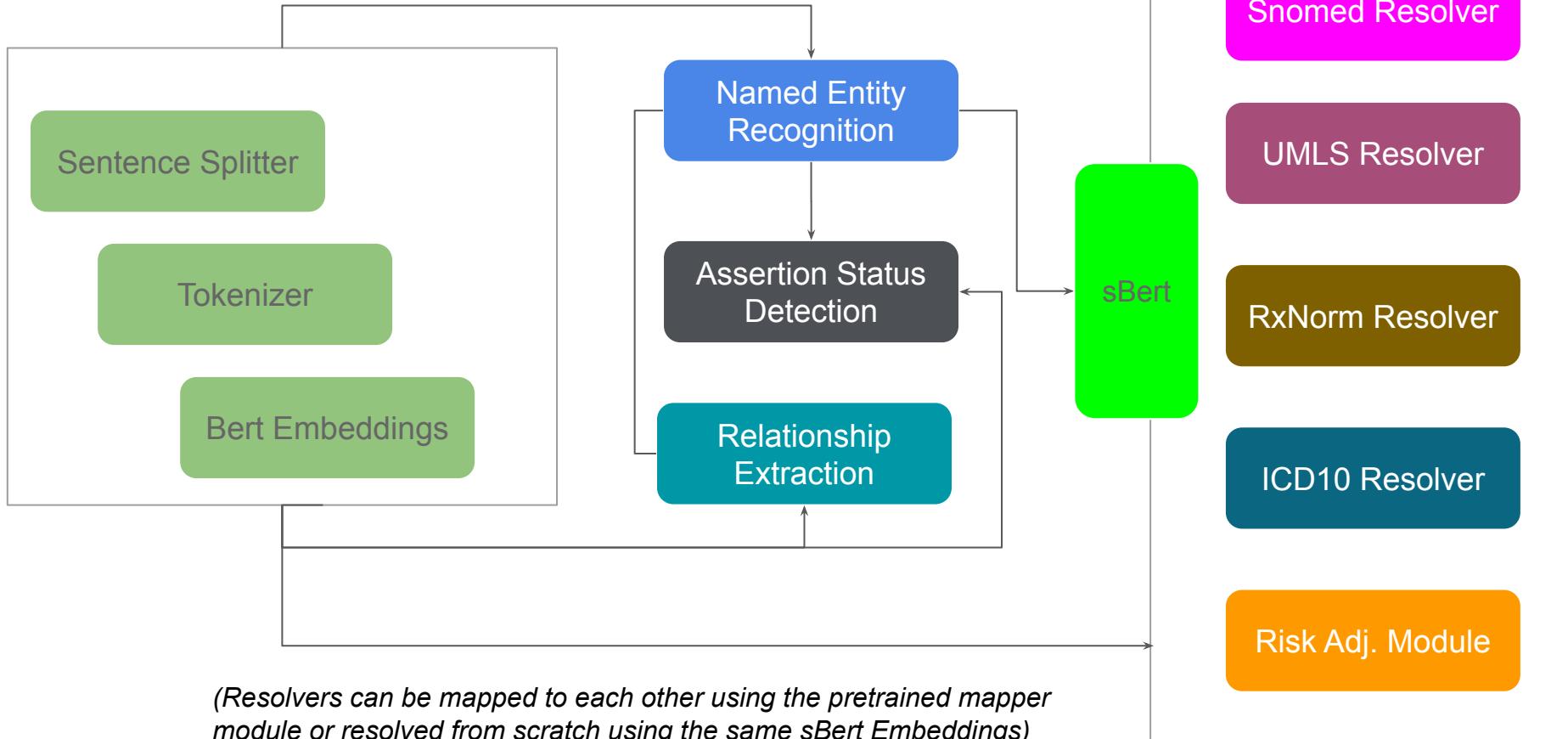
## UMLS

- sbiobertresolve\_umls\_findings
- sbiobertresolve\_umls\_major\_concepts
- sbiobertresolve\_umls\_disease\_syndrome
- sbiobertresolve\_umls\_clinical\_drugs

## mapping

- icd10cm\_snomed\_mapping : ICD10 Codes to Snomed Codes
- snomed\_icd10cm\_mapping : Snomed Codes to ICD Codes
- icd10cm\_umls\_mapping : ICD Codes to UMLS Codes
- snomed\_umls\_mapping : Snomed Codes to UMLS Codes
- rxnorm\_umls\_mapping : RxNorm Codes to UMLS Codes
- mesh\_umls\_mapping : MeSH Codes to UMLS Codes
- rxnorm\_mesh\_mapping : RxNorm Codes to MeSH Codes

# Clinical Entity Resolution



# Entity Resolution with SentenceEntityResolver

```
[ ] documentAssembler = DocumentAssembler()\n    .setInputCol("text")\\n    .setOutputCol("ner_chunk")\n\nsbert_embedder = BertSentenceEmbeddings.pretrained('sbiobert_base_cased_mli', 'en','clinical/models')\\n    .setInputCols(["ner_chunk"])\n    .setOutputCol("sentence_embeddings")\\n    .setCaseSensitive(False)\n\nrxnorm_resolver = SentenceEntityResolverModel.pretrained("sbiobertresolve_rxnorm_augmented","en", "clinical/models") \\n    .setInputCols(["ner_chunk", "sentence_embeddings"]) \\n    .setOutputCol("rxnorm_code")\\n    .setDistanceFunction("EUCLIDEAN")\n\nrxnorm_pipelineModel = PipelineModel(\n    stages = [\n        documentAssembler,\n        sbert_embedder,\n        rxnorm_resolver])\n\nrxnorm_lp = LightPipeline(rxnorm_pipelineModel)
```

```
sbiobert_base_cased_mli download started this may take some time.\nApproximate size to download 384.3 MB\n[OK!]\nsbiobertresolve_rxnorm_augmented download started this may take some time.\n[OK!]
```

```
▶ text = 'aspirin 10 meq/ 5 ml oral sol'
```

```
*time p.predict(text)[cols]
```

```
CPU times: user 92.7 ms, sys: 14.4 ms, total: 107 ms\nWall time: 1.53 s
```

```
resolution_rxnorm_code_target_text resolution_rxnorm_code resolution_rxnorm_code_k_codes resolution_rxnorm_code_k_resolution resolution_rxnorm_code_k_cos_dis
```

```
[memantine hydrochloride 2 MG/ML Oral\nSolution, dicyclomine hydrochloride 2 MG/ML\nOral Solution, codeine phosphate 2 MG/ML\nOral Solution, guaifenesin 10 MG/ML Oral\nSolution, prednisolone 2 MG/ML Oral Solution,\nhydroxyzine hydrochloride 2 MG/ML Oral
```

```
[ ] # Annotator that transforms a text column from dataframe into an Annotation ready for NLP
documentAssembler = DocumentAssembler()\n    .setInputCol("text")\n    .setOutputCol("document")\n\n# Sentence Detector DL annotator, processes various sentences per line\nsentenceDetectorDL = SentenceDetectorDLModel.pretrained("sentence_detector_dl_healthcare", "en", 'clinical/models') \\n    .setInputCols(["document"])\n    .setOutputCol("sentence")\n\n# Tokenizer splits words in a relevant format for NLP\ntokenizer = Tokenizer()\n    .setInputCols(["sentence"])\n    .setOutputCol("token")\n\n# WordEmbeddingsModel pretrained "embeddings_clinical" includes a model of 1.7Gb that needs to be downloaded\nword_embeddings = WordEmbeddingsModel.pretrained("embeddings_clinical", "en", "clinical/models")\\n    .setInputCols(["sentence", "token"])\n    .setOutputCol("word_embeddings")\n\n# Named Entity Recognition for clinical concepts.\nclinical_ner = MedicalNerModel.pretrained("ner_clinical", "en", "clinical/models") \\n    .setInputCols(["sentence", "token", "word_embeddings"])\n    .setOutputCol("ner")\n\nner_converter_icd = NerConverterInternal() \\n    .setInputCols(["sentence", "token", "ner"])\n    .setOutputCol("ner_chunk")\n    .setWhiteList(['PROBLEM'])\n    .setPreservePosition(False)\n\nc2doc = Chunk2Doc() \\n    .setInputCols("ner_chunk")\n    .setOutputCol("doc_ner_chunk")\n\nsbert_embedder = BertSentenceEmbeddings.pretrained('sbiobert_base_cased_mli', 'en','clinical/models')\\n    .setInputCols("doc_ner_chunk")\n    .setOutputCol("sentence_embeddings")\n    .setCaseSensitive(False)\n\nicd_resolver = SentenceEntityResolverModel.pretrained("sbiobertresolve_icd10cm_augmented_billable_hcc","en", "clinical/models") \\n    .setInputCols(["ner_chunk", "sentence_embeddings"])\n    .setOutputCol("icd10cm_code")\n    .setDistanceFunction("EUCLIDEAN")\n\n# Build up the pipeline\nresolver_pipeline = Pipeline(\n    stages = [\n        documentAssembler,\n        sentenceDetectorDL,\n        tokenizer,
```

# Entity Resolution with ChunkMapper

```
[ ] documentAssembler = DocumentAssembler()\
.setInputCol("text")\
.setOutputCol("chunk")

chunkerMapper = ChunkMapperModel.pretrained("rxnorm_mapper", "en", "clinical/models")\
.setInputCols(["chunk"])\
.setOutputCol("rxnorm")\
.setRels(["rxnorm_code"])

pipeline = Pipeline().setStages([documentAssembler,
chunkerMapper])

model = pipeline.fit(spark.createDataFrame([[ '' ]]).toDF('text'))

lp = LightPipeline(model)

rxnorm_mapper download started this may take some time.
[OK!]
CPU times: user 4 µs, sys: 1 µs, total: 5 µs
Wall time: 18.1 µs
[{'chunk': [Annotation(document, 0, 8, metformin, {})],
 'rxnorm': [Annotation(labeled_dependency, 0, 8, 6809, {'entity': 'metformin', 'relation': 'rxnorm_code', 'all_relations': ''})]}]

[ ] %%time

lp.fullAnnotate("metformin")

CPU times: user 6.9 ms, sys: 803 µs, total: 7.71 ms
Wall time: 19.3 ms
[{'chunk': [Annotation(document, 0, 8, metformin, {})],
 'rxnorm': [Annotation(labeled_dependency, 0, 8, 6809, {'entity': 'metformin', 'relation': 'rxnorm_code', 'all_relations': ''})]}]
```

```

documentAssembler = DocumentAssembler()\
    .setInputCol('text')\
    .setOutputCol('document')

sentence_detector = SentenceDetector()\
    .setInputCols(["document"])\
    .setOutputCol("sentence")

tokenizer = Tokenizer()\
    .setInputCols("sentence")\
    .setOutputCol("token")

ner = MedicalBertForTokenClassifier.pretrained("bert_token_classifier_drug_development_trials", "en", "clinical/models")\
    .setInputCols("token", "sentence")\
    .setOutputCol("ner")

nerconverter = NerConverter()\
    .setInputCols("sentence", "token", "ner")\
    .setOutputCol("ner_chunk")

#drug_action_treatment_mapper with "action" mappings
chunkerMapper = ChunkMapperModel().pretrained("drug_action_treatment_mapper", "en", "clinical/models")\
    .setInputCols(["ner_chunk"])\
    .setOutputCol("action_mappings")\
    .setRels(["action"])

pipeline = Pipeline().setStages([documentAssembler,
                                sentence_detector,
                                tokenizer,
                                ner,
                                nerconverter,
                                chunkerMapper])

text = [
    """The patient was female and patient of Dr. X. and she was given Dermovate, Aspagan"""
]

test_data = spark.createDataFrame(text).toDF("text")
res = pipeline.fit(test_data).transform(test_data)

```

bert\_token\_classifier\_drug\_development\_trials download started this may take some time.

[OK!]

drug\_action\_treatment\_mapper download started this may take some time.

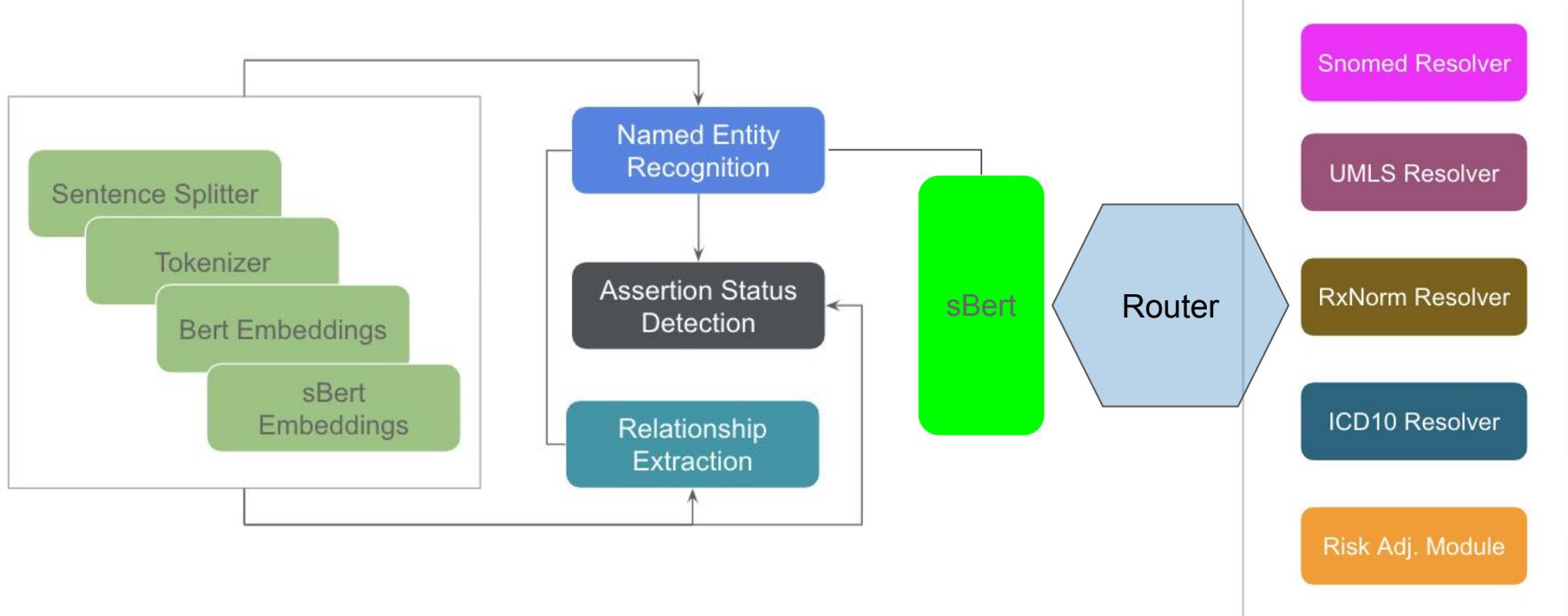
[OK!]

Chunks detected by ner model

```
] : res.select(F.explode('ner_chunk.result').alias("chunks")).show(truncate=False)
```

chunks
Dermovate
Aspagan

# Clinical Entity Resolution





```
ner = MedicalNerModel().pretrained("ner_clinical", "en", "clinical/models") \
    .setInputCols(["sentence", "token", "word_embeddings"]) \
    .setOutputCol("ner")

ner_chunk = NerConverter() \
    .setInputCols("sentence", "token", "ner") \
    .setOutputCol("ner_chunk") \
    .setWhiteList(["PROBLEM", "TREATMENT"])

chunk2doc = Chunk2Doc().setInputCols("ner_chunk").setOutputCol("doc_jsl_chunk")

sbiobert_embeddings = BertSentenceEmbeddings.pretrained("sbiobert_base_cased_mli", "en", "clinical/models") \
    .setInputCols(["doc_jsl_chunk"]) \
    .setOutputCol("sbert_embeddings") \
    .setCaseSensitive(False)

router_sentence_icd10 = Router() \
    .setInputCols("sbert_embeddings") \
    .setFilterFieldsElements(["PROBLEM"]) \
    .setOutputCol("problem_embeddings")

router_sentence_rxnorm = Router() \
    .setInputCols("sbert_embeddings") \
    .setFilterFieldsElements(["TREATMENT"]) \
    .setOutputCol("drug_embeddings")

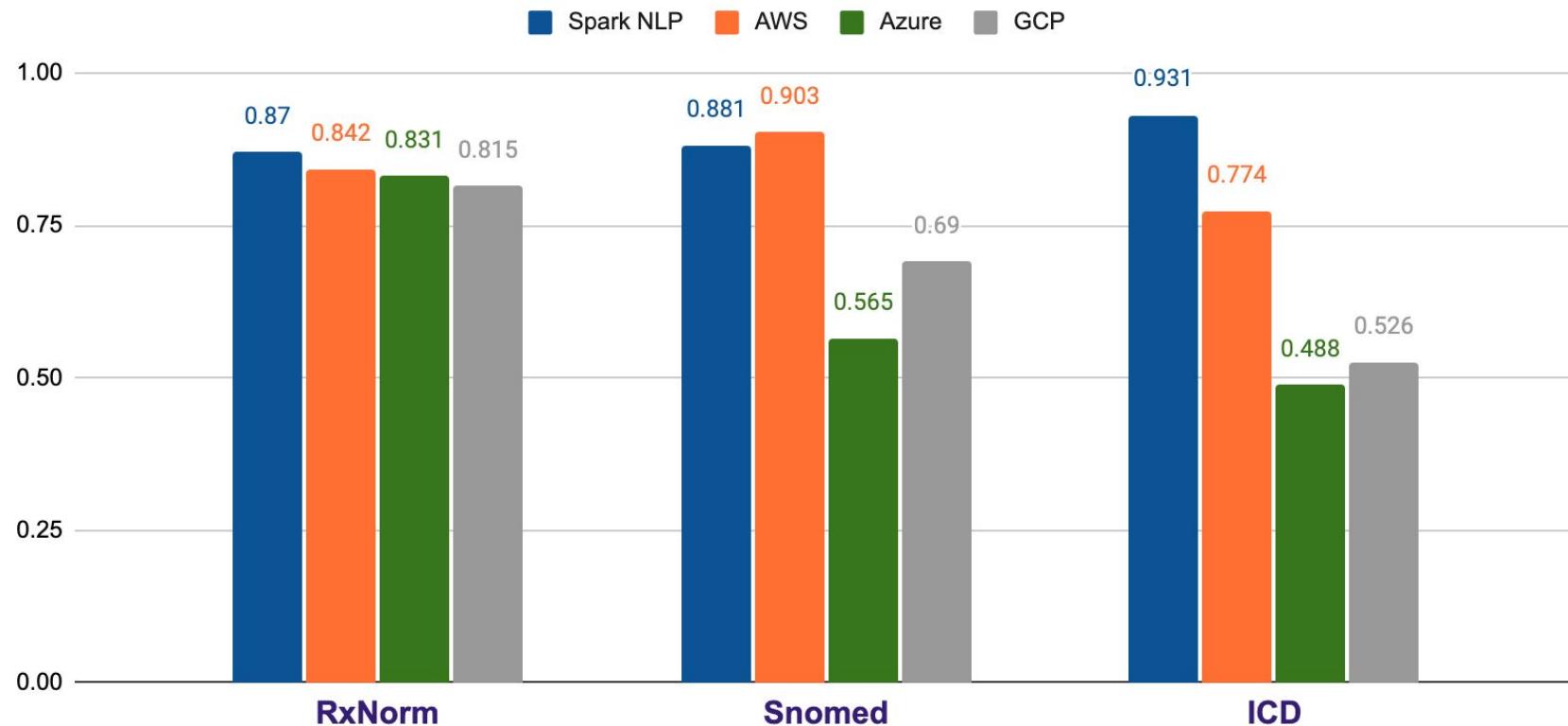
icd_resolver = SentenceEntityResolverModel.pretrained("sbiobertresolve_icd10cm_slim_billable_hcc", "en", "clinical/models") \
    .setInputCols(["problem_embeddings"]) \
    .setOutputCol("icd10cm_code") \
    .setDistanceFunction("EUCLIDEAN")

rxnorm_resolver = SentenceEntityResolverModel.pretrained("sbiobertresolve_rxnorm_augmented", "en", "clinical/models") \
    .setInputCols(["drug_embeddings"]) \
    .setOutputCol("rxnorm_code") \
    .setDistanceFunction("EUCLIDEAN")

resolver_pipeline = Pipeline(stages=[

    documentAssembler,
    sentenceDetector,
    tokenizer,
    word_embeddings,
    ner,
    ner_chunk,
    chunk2doc,
    sbiobert_embeddings,
    router_sentence_icd10,
    router_sentence_rxnorm,
    icd_resolver,
    rxnorm_resolver
])
```

# Top - 5 Results

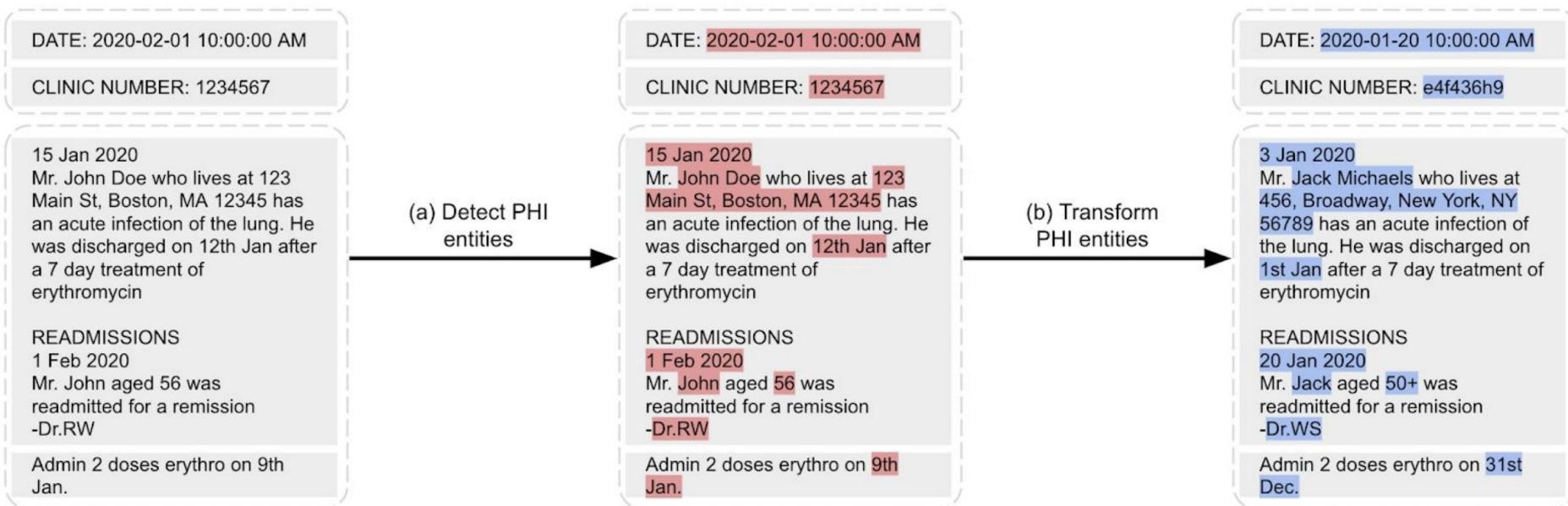


# Part - IV

- ❖ De-Identification and Obfuscation of PHI data

# De-Identification

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



# De-Identification

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

Group Name	Included Entities
A (defined by the HIPAA Safe Harbor Implementation)	Age over 89, Phone/Fax numbers, Email addresses, Websites and URLs, IP Addresses, Dates, Social security numbers, Medical record numbers, Vehicle/Device numbers, Account/Certificate/License numbers, Health plan numbers, Biometric identifiers, Street addresses, City, Zip code, Employer names, Personal names of patients and family members
B	Group A, Doctor names, User IDs (of care providers), State
C	Group B, Hospital names, Country
D	Group C, Holidays, Days of the week, Occupations

# De-Identification

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

```
A . Record date : 2093-01-13 , David Hale , M.D . , Name : Hendrickson , Ora MR . # 7194334  
Date : 01/13/93 PCP : Oliveira , 25 month years-old , Record date : 2079-11-09 . Cocke  
County Baptist Hospital . 0295 Keats Street
```

Color codes: DOCTOR, HOSPITAL, DATE, STREET, MEDICALRECORD, PATIENT,

## Deidentified Text

```
['A .',  
 'Record date : <DATE> , <DOCTOR> , M.D .',  
 ', Name : <PATIENT> , <PATIENT> MR .',  
 '# <MEDICALRECORD> Date : <DATE> PCP : <DOCTOR> , 25  
month years-old , Record date : <DATE> .',  
 '<HOSPITAL> .',  
'<STREET>']
```

```
def get_deidentify_model():  
  
    custom_ner_converter = NerConverter()\  
        .setInputCols(["sentence", "token", "ner"])\\  
        .setOutputCol("ner_chunk")  
        #.setWhiteList(entity_types)  
  
    deidentify_pipeline = Pipeline(  
        stages = [  
            documentAssembler,  
            sentenceDetector,  
            tokenizer,  
            word_embeddings,  
            clinical_ner,  
            custom_ner_converter,  
            deidentification_rules  
        ])  
  
    empty_data = spark.createDataFrame([[""]]).toDF("text")  
  
    model_deidentify = deidentify_pipeline.fit(empty_data)  
  
    return model_deidentify
```

# 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

DATE: 2020-01-20 10:00:00 AM

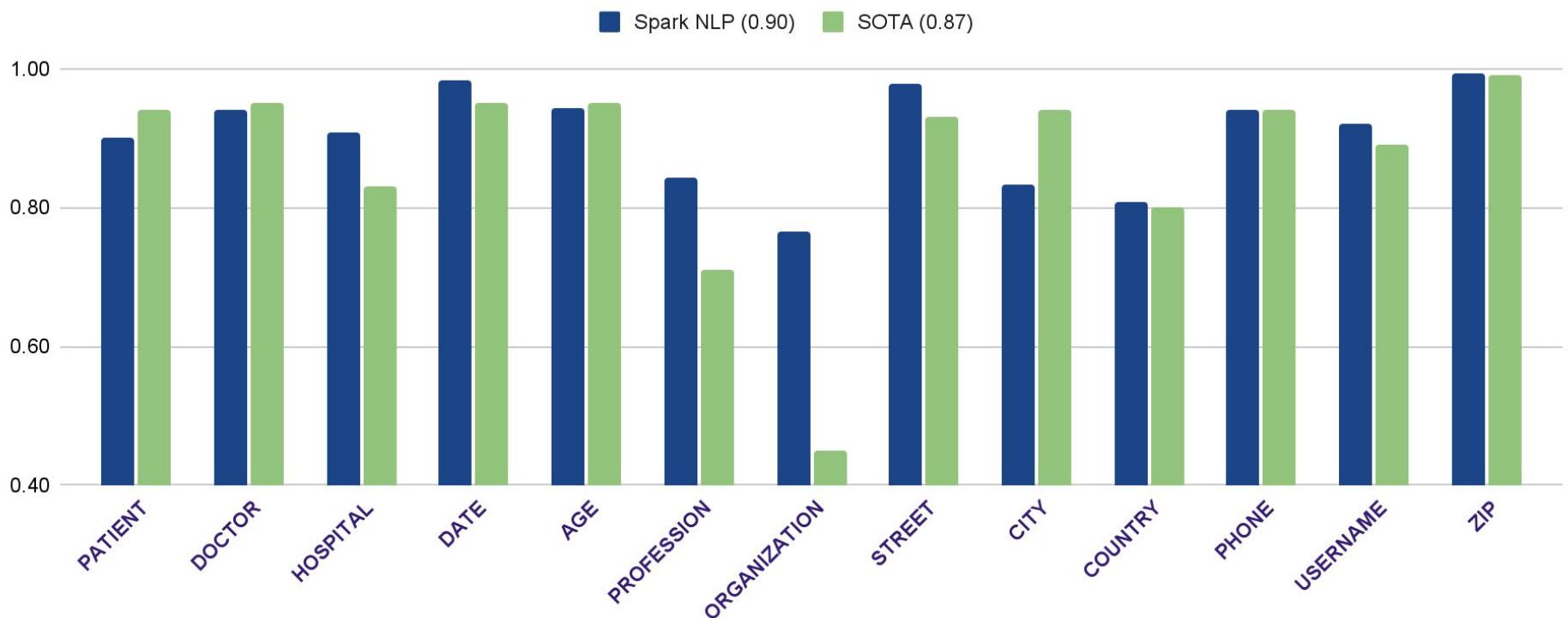
CLINIC NUMBER: e4f436h9

3 Jan 2020  
Mr. Jack Michaels who lives at  
456, Broadway, New York, NY  
56789 has an acute infection of  
the lung. He was discharged on  
1st Jan after a 7 day treatment of  
erythromycin

READMISSIONS  
20 Jan 2020  
Mr. Jack aged 50+ was  
readmitted for a remission  
-Dr.WS

Admin 2 doses erythro on 31st  
Dec.

# Deidentification Benchmarks



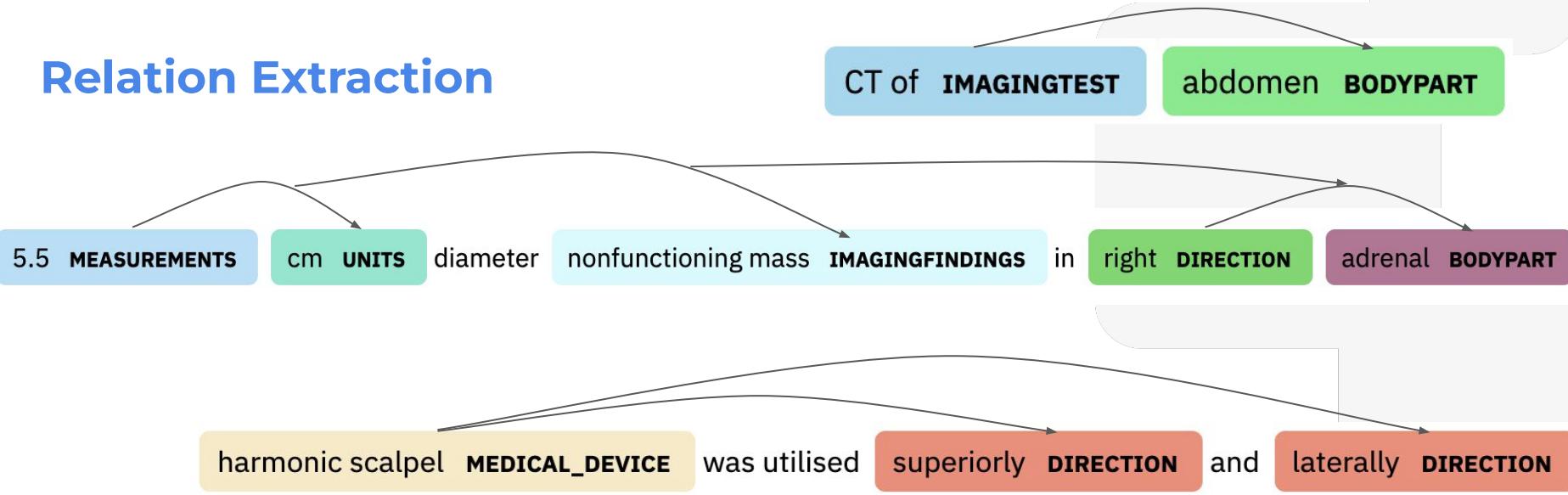
# Part - V

- ❖ Relation Extraction

# Clinical 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."

## Relation Extraction



# Clinical Relation Extraction

## model\_name

0	re_ade_biobert
1	re_ade_clinical
2	re_bodypart_directions
3	re_bodypart_problem
4	re_bodypart_proceduretest
5	re_chemprot_clinical
6	re_clinical
7	re_date_clinical
8	re_drug_drug_interaction_clinical
9	re_human_phenotype_gene_clinical
10	re_temporal_events_clinical
11	re_temporal_events_enriched_clinical
12	re_test_problem_finding
13	re_test_result_date

14	redl_ade_biobert
15	redl_bodypart_direction_biobert
16	redl_bodypart_problem_biobert
17	redl_bodypart_procedure_test_biobert
18	redl_chemprot_biobert
19	redl_clinical_biobert
20	redl_date_clinical_biobert
21	redl_drug_drug_interaction_biobert
22	redl_human_phenotype_gene_biobert
23	redl_temporal_events_biobert

Relation	Recall	Precision	F1	SOTA
DRUG-ADE	0.66	1.00	<b>0.80</b>	0.76
DRUG-DOSAGE	0.89	1.00	<b>0.94</b>	0.91
DRUG-DURATION	0.75	1.00	<b>0.85</b>	0.92
DRUG-FORM	0.88	1.00	<b>0.94</b>	0.95*
DRUG-FREQUENCY	0.79	1.00	<b>0.88</b>	0.90
DRUG-REASON	0.60	1.00	<b>0.75</b>	0.70
DRUG-ROUTE	0.79	1.00	<b>0.88</b>	0.95*
DRUG-STRENGTH	0.95	1.00	<b>0.98</b>	0.97

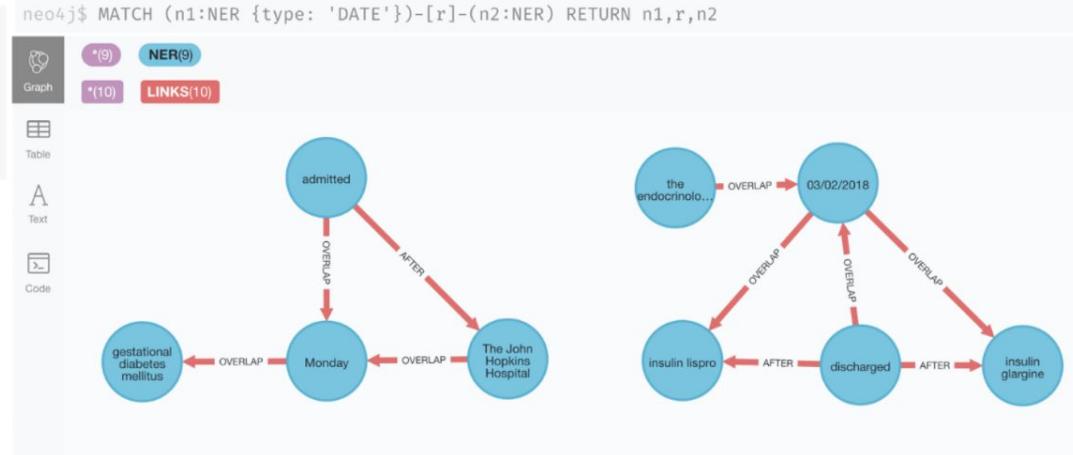
Relation	Recall	Precision	F1
OVERLAP	0.81	0.73	<b>0.77</b>
BEFORE	0.85	0.88	<b>0.86</b>
AFTER	0.38	0.46	<b>0.43</b>

# Clinical Relation Extraction

*She is admitted to The John Hopkins Hospital on Monday with a history of gestational diabetes mellitus diagnosed. She was seen by the endocrinology service and she was discharged on 03/02/2018 on 40 units of insulin glargin and 12 units of insulin lispro.*

```
1 query = """
2 | MATCH (n1:NER {type: 'DATE'})-[r]-(n2:NER)
3 | RETURN n1.name AS date, r.relation AS relation, n2.name AS event
4 """
5
6 df = pd.DataFrame([dict(_) for _ in conn.query(query)])
7 df
```

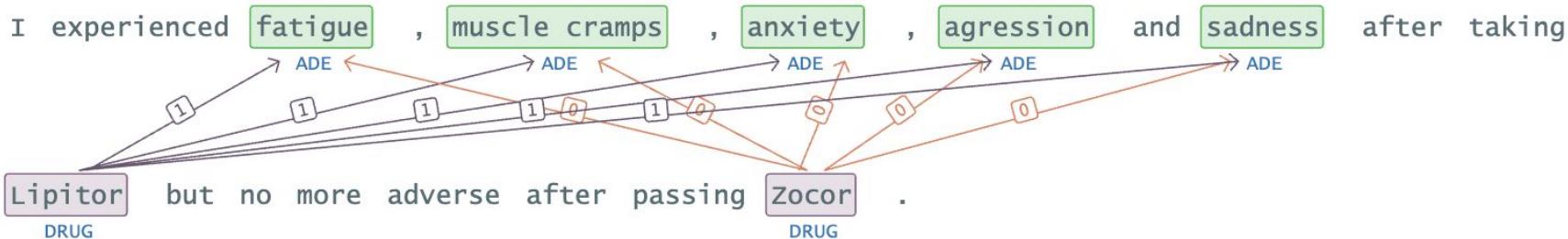
	date	relation	event
0	Monday	OVERLAP	gestational diabetes mellitus
1	Monday	OVERLAP	The John Hopkins Hospital
2	Monday	OVERLAP	admitted
3	03/02/2018	OVERLAP	insulin lispro
4	03/02/2018	OVERLAP	insulin glargin
5	03/02/2018	OVERLAP	discharged
6	03/02/2018	OVERLAP	the endocrinology service



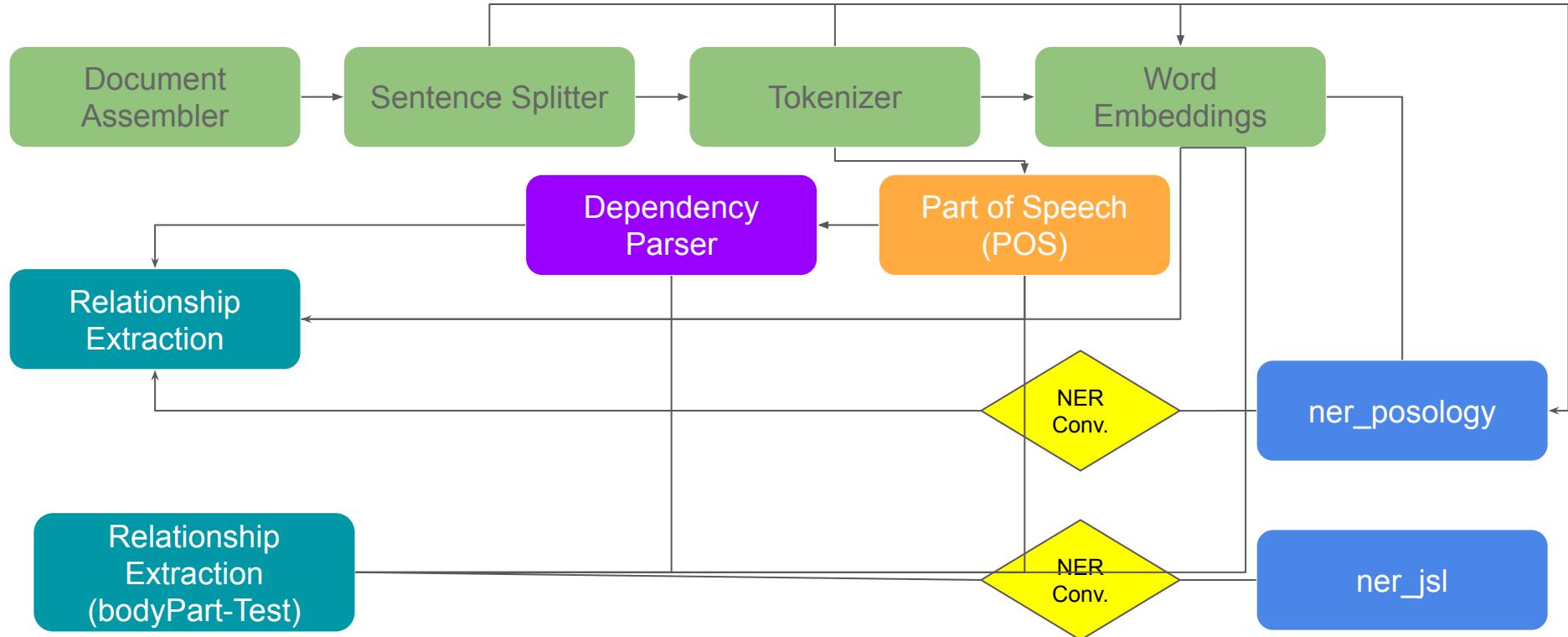
# Clinical Relation Extraction

	relation	entity1	entity1_begin	entity1_end	chunk1	entity2	entity2_begin	entity2_end	chunk2	confidence
0	1	ADE	14	20	fatigue	DRUG	82	88	Lipitor	0.9996617
1	0	ADE	14	20	fatigue	DRUG	124	128	Zocor	0.9952187
2	1	ADE	23	35	muscle cramps	DRUG	82	88	Lipitor	0.9999827
3	0	ADE	23	35	muscle cramps	DRUG	124	128	Zocor	0.91462934
4	1	ADE	38	44	anxiety	DRUG	82	88	Lipitor	0.7636133
5	0	ADE	38	44	anxiety	DRUG	124	128	Zocor	0.9999691
6	1	ADE	47	55	agression	DRUG	82	88	Lipitor	0.99999833
7	0	ADE	47	55	agression	DRUG	124	128	Zocor	0.99781835
8	1	ADE	61	67	sadness	DRUG	82	88	Lipitor	1.0
9	0	ADE	61	67	sadness	DRUG	124	128	Zocor	0.9999572

I experienced fatigue, muscle cramps, anxiety, agression and sadness after taking Lipitor but no more adverse after passing Zocor.



# Clinical Relation Extraction

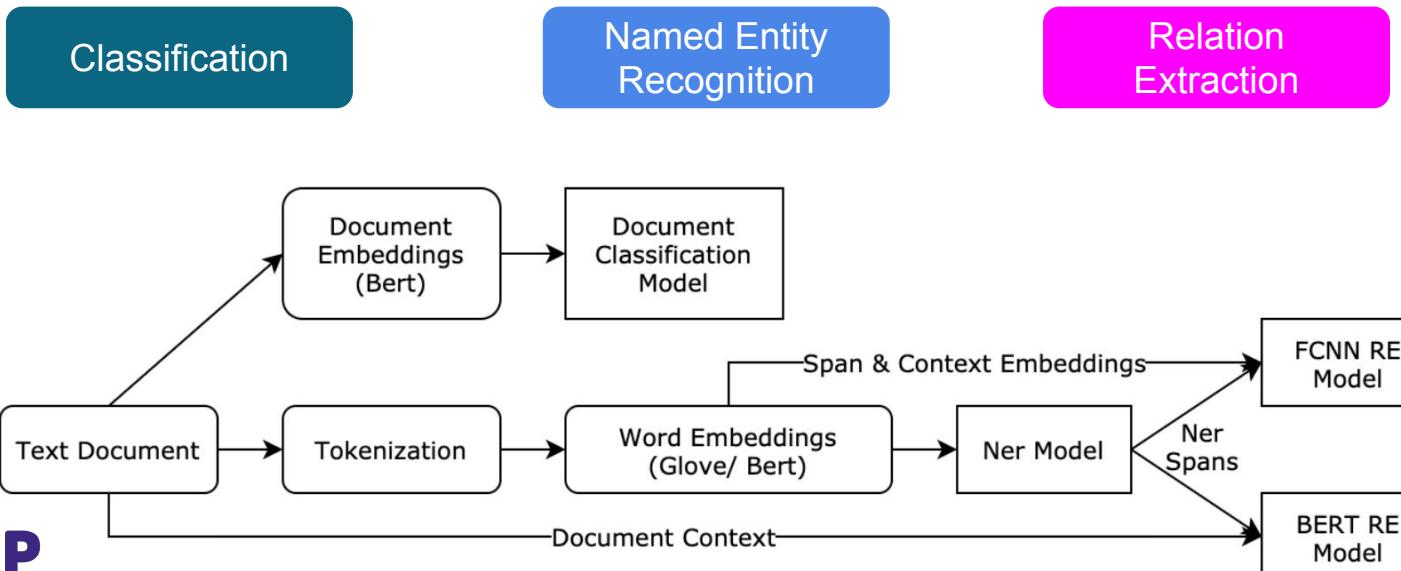


# Part - VI

## ❖ Adverse Drug Reactions (ADR)

# Adverse Drug Reactions (ADR)

Document	Class	ADE Entity	Drug Entity	relation
I feel a bit drowsy & have a little blurred vision after taking insulin.	ADE	drowsy blurred vision	insulin insulin	Positive Positive
@yho fluvastatin gave me cramps, but lipitor suits me!	ADE	cramps cramps	fluvastatin lipitor	Positive Negative
I just took advil and haven't had any gastric problems so far.	NEG	-	-	-



# Adverse Drug Reactions (ADR) Benchmark



Dataset	GLoVe Embeddings						BERT Embeddings						SOTA	
	Precision		Recall		F1		Precision		Recall		F1		F1	
	strict	relax	strict	relax	strict	relax	strict	relax	strict	relax	strict	relax		
ADE	88.32	93.77	89.26	94.80	88.78	94.27	90.0	94.47	93.56	98.22	<b>91.75</b>	96.31	<b>91.3</b>	
	87.81	93.59	88.81	94.66	88.30	94.12	89.6	94.37	93.18	98.13	91.36	96.21		
CADEC	78.14	89.04	77.14	88.01	77.62	88.50	78.53	88.63	79.03	89.32	<b>78.76</b>	88.95	<b>71.9</b>	
	71.87	86.36	71.67	86.13	71.75	86.23	72.38	86.14	73.64	87.66	72.99	86.88		
SMM4H	81.43	90.33	72.17	78.51	76.01	83.41	78.5	86.76	75.23	82.42	<b>76.73</b>	84.41	<b>67.81</b>	
	83.66	91.34	71.31	77.86	76.99	84.06	79.13	87.09	74.33	81.81	76.65	84.36		

Table 2: NER metrics on benchmark datasets. For each dataset, macro and micro averaged scores are displayed on first and second row respectively. SOTA metrics for ADE, CADEC, and SMM4H are obtained from (Yan et al. 2021), (Stanovsky, Gruhl, and Mendes 2017), and (Ge et al. 2020) respectively, and are macro-averaged.

## Named Entity Recognition (NER)

# Adverse Drug Reactions (ADR) Benchmark



Dataset	GLoVe (Avg.) Embeddings			BERT (Avg.) Embeddings			BERT Sentence Embeddings			SOTA
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	F1
ADE	75.96	79.53	76.86	76.91	84.96	79.37	87.41	84.72	<b>85.96</b>	<b>87.0</b>
	86.84	81.22	83.43	88.13	84.38	85.38	90.97	91.20	91.03	
CADEC	85.29	84.24	84.71	86.50	86.11	86.30	87.13	86.32	<b>86.69</b>	<b>81.5</b>
	85.99	86.10	86.0	87.38	87.43	87.40	87.78	87.86	87.79	

Table 1: Classification Metrics on benchmark datasets. For each dataset, Macro and Micro averaged scores are displayed on first and second row respectively. SOTA metrics for ADE and CADEC datasets are obtained from (Huynh et al. 2016) and (Alimova and Tutubalina 2019) respectively.

Dataset	Base (FCNN) RE			BERT RE			SOTA
	Precision	Recall	F1	Precision	Recall	F1	F1
ADE Corpus	69.11	86.22	<b>74.70</b>	81.31	79.03	<b>80.10</b>	<b>83.74</b>
ADE Enriched with n2c2	89.01	89.44	<b>89.22</b>	89.19	90.93	<b>90.02</b>	

Table 3: Relation Extraction performance on the ADE benchmark dataset. The test set was kept standard for a fair comparison, and all scores are macro-averaged due to high class imbalance. SOTA metrics for RE on ADE corpus as reported by (Crone 2020)

# Part - VII

## ❖ Key Chunk Phrase Extractor

# Chunk Key Phrase Extractor (KCPE)

key_phrase	source	DocumentSimilarity	MMRScore	sentence
type two diabetes mellitus	NER	0.7639750686118073	0.4583850593816694	0
subsequent type diabetes	ngrams	0.7503709443591438	0.08298243928224425	0
HTG-induced pancreatitis years	ngrams	0.6817062970203589	0.11246275270031031	0
hepatitis obesity	ngrams	0.6666053470245074	0.1177052008980295	0
mellitus diagnosed years	ngrams	0.6389213391545323	0.08129479185432026	0
history gestational diabetes	ngrams	0.6219876368539883	0.0950104202982544	0
vomiting	ngrams	0.5824238088130589	0.14864183399720493	0
admitted starvation ketosis	ngrams	0.5789875069392564	0.12008073486190007	0
five-day amoxicillin respiratory	ngrams	0.5330653868257814	0.09428153526023508	0
28-year-old female history	ngrams	0.38613601247069695	0.12987678861407687	0

## YAKE

keyword	score
years prior presentation	0.006335399690627251
prior presentation	0.011644010991495998
prior presentation subsequent	0.020272229518351368
weeks prior presentation	0.020272229518351368
respiratory tract infection	0.02568455658449274
anion gap	0.025965846371439553
physical examination presentation	0.02840600503736659
obtained hours presentation	0.028532992974589392
examination presentation significant	0.028532992974589392
prior	0.029673513395379065
years prior	0.03008818777992058
anion gap elevated	0.031568192739369824

## CKPE

key_phrase_candidate	DocumentSimilarity
pancreatitis years prior	0.6491587146812722
diagnosed years prior	0.38594469396979897
respiratory tract infection	0.34452766290310755
patient treated insulin	0.3413457416284759
serum	0.3371024001999838
presentation revealed glucose	0.31458360368143906
examination presentation significant	0.29099950377907047
prior analysis due	0.22501711661945623
prior	0.21634008371261446
physical examination presentation	0.19165189487112474

key_phrase_candidate	source
28-year-old female history	ngram
28-year-old	NER
female history gestational	ngram
female	NER
history gestational diabetes	ngram
gestational diabetes mellitus	NER
gestational diabetes mellitus	ngram
diabetes mellitus diagnosed	ngram
mellitus diagnosed years	ngram
diagnosed years prior	ngram
eight years prior	NER
years prior presentation	ngram
prior presentation subsequent	ngram
presentation subsequent type	ngram
subsequent type diabetes	ngram
type diabetes mellitus	ngram
type two diabetes mellitus	NER
diabetes mellitus (	ngram
mellitus ( T2DM	ngram
( T2DM ),	ngram
T2DM ), prior	ngram
T2DM	NER
), prior episode	ngram
prior episode HTG-induced	ngram
episode HTG-induced pancreatitis	ngram
HTG-induced pancreatitis years	ngram
HTG-induced pancreatitis	NER
pancreatitis years prior	ngram
three years prior	NER
years prior presentation	ngram
prior presentation ,	ngram
presentation , acute	ngram
, acute hepatitis	ngram



# Thank you !

**Veysel Kocaman**  
Principal Data Scientist  
John Snow Labs



# Spark NLP Resources

Spark NLP Official page

Spark NLP Workshop Repo

JSL Youtube channel

JSL Blogs

Introduction to Spark NLP: Foundations and Basic Components (Part-I)

Introduction to: Spark NLP: Installation and Getting Started (Part-II)

Named Entity Recognition with Bert in Spark NLP

Text Classification in Spark NLP with Bert and Universal Sentence Encoders

Spark NLP 101 : Document Assembler

Spark NLP 101: LightPipeline

<https://www.oreilly.com/radar/one-simple-chart-who-is-interested-in-spark-nlp/>

<https://blog.dominodatalab.com/comparing-the-functionality-of-open-source-natural-language-processing-libraries/>

<https://databricks.com/blog/2017/10/19/introducing-natural-language-processing-library-apache-spark.html>

<https://databricks.com/fr/session/apache-spark-nlp-extending-spark-ml-to-deliver-fast-scalable-unified-natural-language-processing>

<https://medium.com/@saif1988/spark-nlp-walkthrough-powered-by-tensorflow-9965538663fd>

<https://www.kdnuggets.com/2019/06/spark-nlp-getting-started-with-worlds-most-widely-used-nlp-library-enterprise.html>

<https://www.forbes.com/sites/forbestechcouncil/2019/09/17/why-spark-nlp-is-the-most-widely-used-nlp-library-enterprise/>

<https://medium.com/hackernoon/mueller-report-for-nerds-spark-meets-nlp-with-tensorflow-and-bert-part-1-32490a8f8f12>

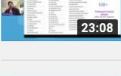
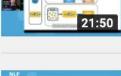
<https://www.analyticsindiamag.com/5-reasons-why-spark-nlp-is-the-most-widely-used-library-in-enterprises/>

<https://www.oreilly.com/ideas/comparing-production-grade-nlp-libraries-training-spark-nlp-and-spacy-pipelines>

<https://www.oreilly.com/ideas/comparing-production-grade-nlp-libraries-accuracy-performance-and-scalability>

<https://www.infoworld.com/article/3031690/analytics/why-you-should-use-spark-for-machine-learning.html>

# Healthcare NLP Summit '22

1	 End-to-End No-Code Development of NER model for Text with Annotation Lab John Snow Labs  4:42	Using Spark NLP in R: A Drug Standardization Case Study John Snow Labs  15:00	11	Natural Language Technologies: Current Status and Future Evolution John Snow Labs  19:54
2	 How to Build a Foundation of AI-based Healthcare Systems Through Language Models? John Snow Labs  34:12	How to sign up for a free trial and how to use your license on COLAB John Snow Labs  1:09	12	The Quest for Proactive and Reactive Healthcare John Snow Labs  15:31
3	 Automated Patient Risk Adjustment and Medicare HCC Coding from Clinical Notes John Snow Labs  28:19	Few-Shot Text Classification in the Real-World John Snow Labs  25:16	13	The Unified NLP Platform John Snow Labs  26:44
4	 End-to-End No-Code Development of Visual NER Models for PDFs and Images John Snow Labs  6:31	Medical NLP: Domain Expertise and Data Quality are Vital For Success John Snow Labs  13:50	14	Industry Survey Analysis: AI in Healthcare 2022 John Snow Labs  19:59
5	 A Hierarchical Approach for Automated ICD-10 Coding Using Phrase-level Attention John Snow Labs  29:04	Automatic mining of adverse drug reactions from social media posts and unstructured chats John Snow Labs  19:24	15	Using NLP at Scale to Process Patient Charts for Identifying Patient Encounters John Snow Labs  23:54
6	 End-to-End No-Code Development of AI Models for Text and Images John Snow Labs  14:41	Data Centric AI for Healthcare John Snow Labs  26:23	16	Transfer Learning From Existing Diseases Via Hierarchical Multi-Modal BERT Models to Predict COVID19 John Snow Labs  24:19
7	 Opportunities and Challenges of Applying Advances in NLP to Healthcare John Snow Labs  20:38	Radiology Report Summarization John Snow Labs  36:15	17	Supporting Mental Healthcare Delivery Using NLP John Snow Labs  25:33
8	 Current State-of-the-Art Accuracy for Key Medical Natural Language Processing Benchmarks John Snow Labs  23:08	Entropy and Sentiment: the Anna Karenina Principle in Patient Experience Data John Snow Labs  31:06	18	Using Spark NLP to De-Identify Doctor Notes in the German Language John Snow Labs  28:50
9	 Machine Reading for Precision Medicine John Snow Labs  21:50	1 Line of Code to Use 600+ State-of-the-Art Clinical & Biomedical NLP Models John Snow Labs  31:31	19	Accelerating Use of the Digital Medical Record with Optimized Structured & Unstructured Data John Snow Labs  30:20
10	 Collaborative Healthcare NLP: Customisable NLP platforms for health and related research John Snow Labs  18:10	Cleanlab: Making AI Work with Messy, Real-World Healthcare and NLP Data John Snow Labs  23:40	20	Is it Enough to Simply Apply Language Model for Optimal Text Classification? John Snow Labs  25:57