



---

# Spark NLP for Data Scientists

Jan 26-27, 2021

Veysel Kocaman  
Lead Data Scientist  
[veysel@johnsnowlabs.com](mailto:veysel@johnsnowlabs.com)





# Spark NLP

## FOR DATA SCIENTISTS

Day 1: January 26, 2020, 12:00 PM – 4:00 PM EST

Day 2: January 27, 2020, 12:00 PM – 4:00 PM EST

Training Location: Online

Certification Exam: Online, January 28th to February 3rd



# Spark NLP

## for Healthcare Data Scientists

Day 1: January 28, 2020, 12:00 PM – 4:00 PM EST

Day 2: January 29, 2020, 12:00 PM – 4:00 PM EST

Training Location: Online

Certification Exam: Online, January 30th to February 5th

**4 days, 16 hrs live training**

**Certification exams in 2 weeks**

# Welcome

Day-1	60 min	- Intro to John Snow Labs and Spark NLP - Intro to NLP and Spark NLP Basics (colab)
	10 min	Break
	50 min	- Text preprocessing with Spark NLP Notebook
	10 min	Break
	50 min	- Pretrained Models in Spark NLP
Day-2	50 min	- Named Entity Recognition (NER) (part-1)
	10 min	Break
	60 min	- Text Classification with Spark NLP
	10 min	Break
	50 min	- T5 (Text-to-Text Transfer Transformer)

# Setup

 Open in Colab

## RUNNING CODE:

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

[How to set up Google Colab]

## BOOKMARK:

<https://nlp.johnsnowlabs.com/>  
<https://nlp.johnsnowlabs.com/docs/en/quickstart>  
[spark-nlp.slack.com](https://spark-nlp.slack.com)

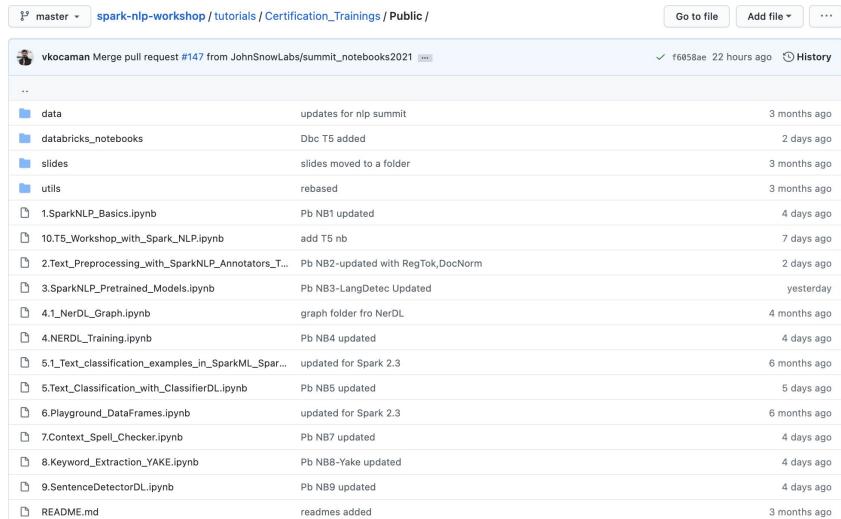
**spark-nlp==2.7.1**

```
[ ] import os

# Install java
! apt-get update -qq
! apt-get install -y openjdk-8-jdk-headless -qq > /dev/null

os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
os.environ["PATH"] = os.environ["JAVA_HOME"] + "/bin:" + os.environ["PATH"]
! java -version

# Install pyspark
! pip install --ignore-installed -q pyspark==2.4.4
! pip install --ignore-installed -q spark-nlp==2.7.1
```



The screenshot shows a GitHub pull request history for a repository named 'spark-nlp-workshop'. The pull request is titled 'Merge pull request #147 from JohnSnowLabs/summit\_notebooks2021' and was merged 22 hours ago by user 'vkocaman'. The history lists several commits and file changes:

- data: updates for nlp summit (3 months ago)
- databricks\_notebooks: Dbc T5 added (2 days ago)
- slides: slides moved to a folder (3 months ago)
- utils: rebased (3 months ago)
- 1.SparkNLP\_Basics.ipynb: Pb NB1 updated (4 days ago)
- 10.T5\_Workshop\_with\_Spark\_NLP.ipynb: add T5 nb (7 days ago)
- 2.Text\_Preprocessing\_with\_SparkNLP\_Annotators\_T...: Pb NB2-updated with RegTok,DocNorm (2 days ago)
- 3.SparkNLP\_Pretrained\_Models.ipynb: Pb NB3-LangDetec Updated (yesterday)
- 4.1\_NerDL\_Graph.ipynb: graph folder fro NerDL (4 months ago)
- 4.NERDL\_Training.ipynb: Pb NB4 updated (4 days ago)
- 5.1\_Text\_classification\_examples\_in\_SparkML\_Spar...: updated for Spark 2.3 (6 months ago)
- 5.Text\_Classification\_with\_ClassifierDL.ipynb: Pb NB5 updated (5 days ago)
- 6.Playground\_DataFrames.ipynb: updated for Spark 2.3 (6 months ago)
- 7.Context\_Spell\_Checker.ipynb: Pb NB7 updated (4 days ago)
- 8.Keyword\_Extraction\_YAKE.ipynb: Pb NB8-Yake updated (4 days ago)
- 9.SentenceDetectorDL.ipynb: Pb NB9 updated (4 days ago)
- README.md: readme added (3 months ago)

# Part - I

- ❖ Introducing JSL and Spark NLP
- ❖ Natural Language Processing (NLP) Basics
- ❖ Spark NLP Pretrained Pipelines
- ❖ Text Preprocessing with Spark NLP
- ❖ Spark NLP Pretrained Models



"John Snow Labs enables healthcare organizations to deploy state-of-the-art artificial intelligence (AI) platforms, models and data in production today."

# JOHN SNOW LABS



"John Snow Labs wows in both proven customer success and verifiable state-of-the-art technology – making it a natural winner of the highly competitive 2019

AI Platform of the Year Award."



"Keep an eye on this company – as it represents where the industry and data science are headed."

# Introducing Spark NLP

Daily ~ 12K  
Monthly ~ 360K

PyPI link

<https://pypi.org/project/spark-nlp>

Total downloads

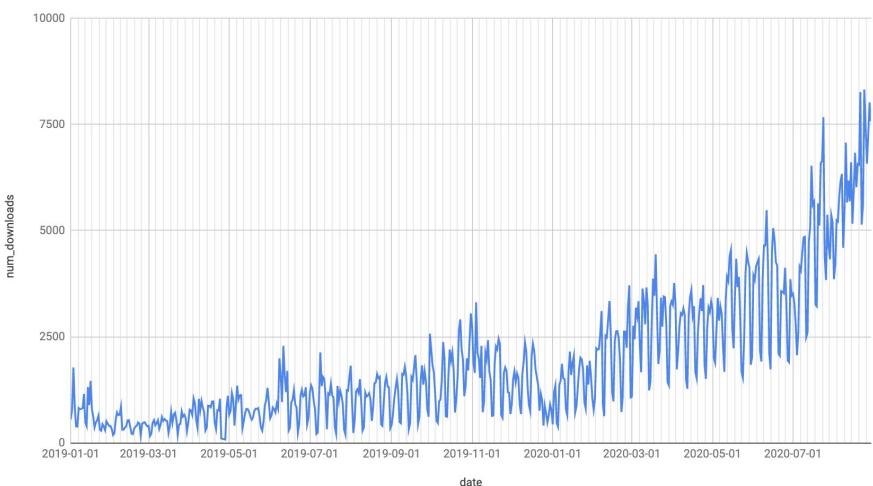
2,823,976

Total downloads - 30 days

360,631

Total downloads - 7 days

79,475



- 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 (26 new releases in 2020, 30 new releases in 2019.)

# Spark NLP

- 73 total releases
- Release every two weeks for the past 3 years
- A single unified library for all your NLP/NLU need
- Active community on Slack and GitHub

NLP Feature	Spark NLP	spaCy	NLTK	CoreNLP	Hugging Face
Tokenization	Yes	Yes	Yes	Yes	Yes
Sentence segmentation	Yes	Yes	Yes	Yes	No
Steeming	Yes	Yes	Yes	Yes	No
Lemmatization	Yes	Yes	Yes	Yes	No
POS tagging	Yes	Yes	Yes	Yes	No
Entity recognition	Yes	Yes	Yes	Yes	Yes
Dep parser	Yes	Yes	Yes	Yes	No
Text matcher	Yes	Yes	No	No	No
Date matcher	Yes	No	No	No	No
Sentiment detector	Yes	No	Yes	Yes	Yes
Text classification	Yes	Yes	Yes	No	Yes
Spell checker	Yes	No	No	No	No
Language detector	Yes	No	No	No	No
Keyword extraction	Yes	No	No	No	No
Pretrained models	Yes	Yes	Yes	Yes	Yes
Trainable models	Yes	Yes	Yes	Yes	Yes

# TRUSTED BY



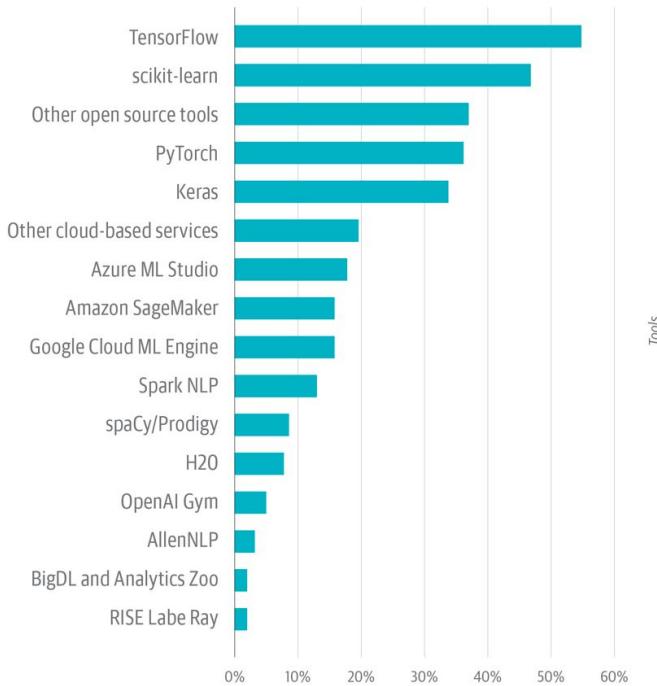
Imperial College  
London



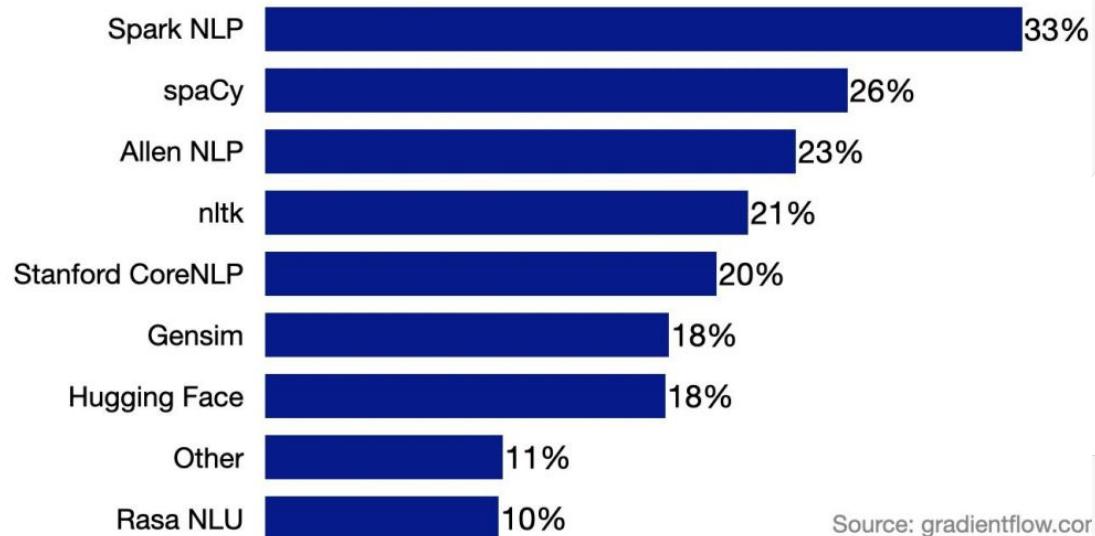
STANFORD  
UNIVERSITY

# Spark NLP in Industry

Which of the following AI tools do you use?



Which NLP libraries does your organization use?



Source: gradientflow.co

NLP Industry Survey by Gradient Flow,  
an independent data science research & insights company, September 2020

## OFFICIALLY SUPPORTED RUNTIMES



databricks

CLOUDERA



Azure



# Biomedical Named Entity Recognition at Scale

Veysel Kocaman  
John Snow Labs Inc.  
16192 Coastal Highway  
Lewes, DE , USA 19958  
veysel@johnsnowlabs.com

**Abstract**—Named entity recognition (NER) is a widely applicable natural language processing task and building block of question answering, topic modeling, information retrieval, etc. In the medical domain, NER plays a crucial role by extracting meaningful chunks from clinical notes and reports, which are then fed to downstream tasks like assertion status detection, entity resolution, relation extraction, and de-identification. Reimplementing a Bi-LSTM-CNN-Char deep learning architecture on top of Apache Spark, we present a single trainable NER model that obtains new state-of-the-art results on seven public biomedical benchmarks without using heavy contextual embeddings like BERT. This includes improving BC4CHEMD to 93.72% (4.1% gain), Species800 to 80.91% (4.6% gain), and JNLPBA to 81.29% (5.2% gain). In addition, this model is freely available within a production-grade code base as part of the open-source Spark NLP library; can scale up for training and inference in any Spark cluster; has GPU support and libraries for popular programming languages such as Python, R, Scala and Java; and can be extended to support other human languages with no code changes.

## I. INTRODUCTION

Electronic health records (EHRs) are the primary source of information for clinicians tracking the care of their patients. Information fed into these systems may be found in structured fields for which values are inputted electronically (e.g. laboratory test orders or results) [1] but most of the time information in these records is unstructured making it largely inaccessible

## Abstract

Named entity recognition (NER) is one of the most important building blocks of NLP tasks in the medical domain by extracting meaningful chunks from clinical notes and reports, which are then fed to downstream tasks like assertion status detection, entity resolution, relation extraction, and de-identification. Due to the growing volume of healthcare data in unstructured format, an increasingly important challenge is providing high accuracy implementations of state-of-the-art deep learning (DL) algorithms at scale. In this study, we introduce a production-grade clinical and biomedical NER algorithm based on a modified BiLSTM-CNN-Char DL architecture built on top of Apache Spark. This algorithm establishes new state-of-the-art accuracy on 7 of 8 well-known biomedical NER benchmarks and 3 clinical concept extraction challenges: 2010 i2b2/VA clinical concept extraction, 2014 n2c2 de-identification, and 2018 n2c2 medication extraction. Moreover, clinical NER models trained using this implemen-

Anonymous NAACL-HLT 2021 submission

# Spark NLP: Natural Language Understanding at Scale

Veysel Kocaman, David Talby

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

## Accurate Clinical and Biomedical Named Entity Recognition at Scale

## Improving Clinical Document Understanding on COVID-19 Research with Spark NLP

Veysel Kocaman, David Talby

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

## Abstract

Following the global COVID-19 pandemic, the number of scientific papers studying the virus has grown massively, leading to increased interest in automated literature review. We present a clinical text mining system that improves on previous efforts in three ways. First, it can recognize over 100 different entity types including social determinants of health, anatomy, risk factors, and adverse events in addition to other commonly used clinical and biomedical entities. Second, the text processing pipeline includes assertion status detection, to distinguish between clinical facts that are present, absent, conditional, or about someone other than the patient. Third, the deep learning models used are more accurate than previously available, leveraging an integrated pipeline of state-of-the-art pre-trained named entity recognition models, and improving on the previous best performing benchmarks for assertion status detection. We illustrate extracting trends and insights - e.g. most frequent disorders and symptoms, and most common vital signs and EKG findings – from the COVID-19 Open Research Dataset (CORD-19). The system is built using the Spark NLP library which natively supports scaling to use distributed clusters, leveraging GPU's, configurable and reusable NLP pipelines, healthcare-specific embeddings, and the ability to train models to support new entity types or human languages with no code changes.

be found in structured fields for which values are inputted electronically (e.g. laboratory test orders or results) (Liede et al. 2015) but most of the time information in these records is unstructured making it largely inaccessible for statistical analysis (Murdoch and Detsky 2013). These records include information such as the reason for administering drugs, previous disorders of the patient or the outcome of past treatments, and they are the largest source of empirical data in biomedical research, allowing for major scientific findings in highly relevant disorders such as cancer and Alzheimer's disease (Perera et al. 2014).

A primary building block in such text mining systems is named entity recognition (NER) - which is regarded as a critical precursor for question answering, topic modelling, information retrieval, etc (Yadav and Bethard 2019). In the medical domain, NER recognizes the first meaningful chunks out of a clinical note, which are then fed down the processing pipeline as an input to subsequent downstream tasks such as clinical assertion status detection (Uzuner et al. 2011), clinical entity resolution (Tzitzivacos 2007) and de-identification of sensitive data (Uzuner, Luo, and Szolovits 2007) (see Figure 1). However, segmentation of clinical and drug entities is considered to be a difficult task in biomedical NER systems because of complex orthographic structures of named entities

# Spark NLP Modules

Clinical Entity Recognition	Clinical Entity Linking	Assertion Status	Relation Extraction	Entity Recognition	Information Extraction	Sentiment Analysis	Document Classification						
<p>60 units <b>DOSAGE</b> of insulin glargine <b>DRUG</b> at night. <b>FREQUENCY</b></p>	<p>Suspect diabetes SNOMED-CT: 473122005 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>		<p>How Lucy <small>version</small></p>	<p>They met <small>date</small> → 29-04-2020</p>								
Algorithms				Content									
<b>Extract Knowledge</b> <ul style="list-style-type: none"> <li>Entity Linker</li> <li>Entity Disambiguator</li> <li>Document Classifier</li> <li>Contextual Parser</li> </ul>	<b>De-identify text</b> <ul style="list-style-type: none"> <li>Structured Data</li> <li>Unstructured Text</li> <li>Obfuscator</li> <li>Generalizer</li> </ul>	<b>Medical Transformers</b> <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>	<b>Linked Medical Terminologies</b> <ul style="list-style-type: none"> <li>SNOMED-CT</li> <li>CPT</li> <li>ICD-10-CM</li> <li>RxNorm</li> <li>ICD-10-PCS</li> <li>ICD-O</li> <li>LOINC</li> </ul>	<b>Split Text</b> <ul style="list-style-type: none"> <li>Sentence Detector</li> <li>Deep Sentence Detector</li> <li>Tokenizer</li> <li>nGram Generator</li> <li>Word Segmentation</li> </ul>	<b>Clean Text</b> <ul style="list-style-type: none"> <li>Spell Checking</li> <li>Spell Correction</li> <li>Normalizer</li> <li>Stopword Cleaner</li> <li>Summarization</li> </ul>	<b>Transformers</b> <ul style="list-style-type: none"> <li>BERT</li> <li>ELMO</li> <li>GloVe</li> <li>ALBERT</li> <li>XLNet</li> <li>USE</li> <li>Small BERT</li> <li>ELECTRA</li> <li>T5</li> <li>NMT</li> <li>LaBSE</li> </ul>	<b>200+ Languages</b> 						
<b>Split Text</b> <ul style="list-style-type: none"> <li>Sentence Detector</li> <li>Deep Sentence Detector</li> <li>Tokenizer</li> <li>nGram Generator</li> </ul>	<b>Clean Medical Text</b> <ul style="list-style-type: none"> <li>Spell Checking</li> <li>Spell Correction</li> <li>Normalizer</li> <li>Stopword Cleaner</li> </ul>	<b>75+ Pretrained Models</b> <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>Understand Grammar</b> <ul style="list-style-type: none"> <li>Stemmer</li> <li>Lemmatizer</li> <li>Part of Speech Tagger</li> <li>Dependency Parser</li> <li>Translation</li> </ul>	<b>Find in Text</b> <ul style="list-style-type: none"> <li>Text Matcher</li> <li>Regex Matcher</li> <li>Date Matcher</li> <li>Chunker</li> <li>Question Answering</li> </ul>	<b>Pre-trained Models</b> <p><b>700+</b></p>	<b>Pre-trained Pipelines</b> <p><b>400+</b></p>
<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 Grammar</b> <ul style="list-style-type: none"> <li>Stemmer</li> <li>Lemmatizer</li> <li>Part of Speech Tagger</li> <li>Dependency Parser</li> </ul>	<b>Find in Text</b> <ul style="list-style-type: none"> <li>Text Matcher</li> <li>Regex Matcher</li> <li>Date Matcher</li> <li>Chunker</li> </ul>		<b>Trainable &amp; Tunable</b>	<b>Scalable to a Cluster</b>	<b>Fast Inference</b>	<b>Hardware Optimized</b>	<b>Community</b>						
													

# Spark NLP Modules (Enterprise and Public)

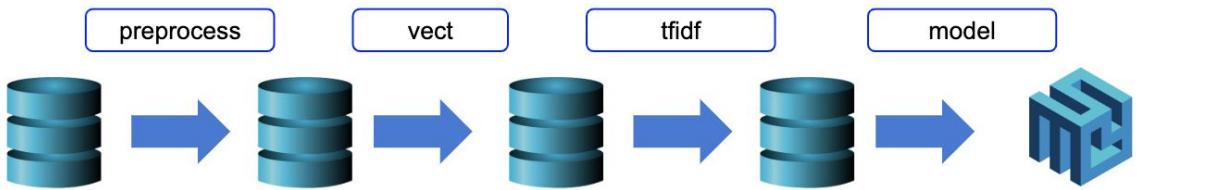
# Spark NLP: Apache License 2.0

## NLP Features

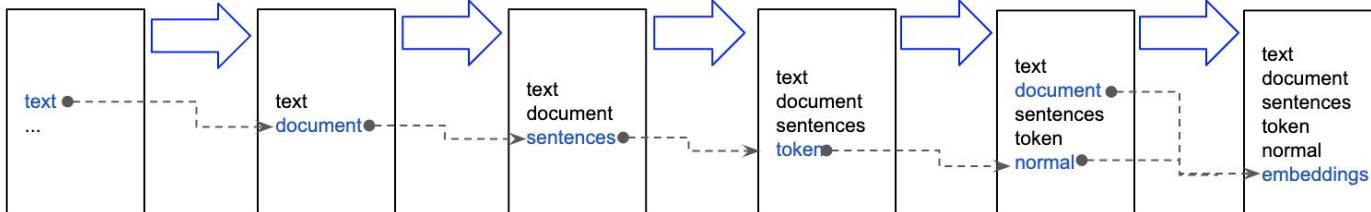
- Tokenization
- Word Segmentation
- Stop Words Removal
- Normalizer
- Stemmer
- Lemmatizer
- NGrams
- Regex Matching
- Text Matching
- Chunking
- Date Matcher
- **Part-of-speech** tagging
- Sentence Detector (DL models)
- **Dependency** parsing
- **Sentiment** Detection (ML models)
- **Spell** Checker (ML and DL models)
- Word Embeddings (**GloVe** and **Word2Vec**)
- **BERT** Embeddings
- **ELMO** Embeddings
- **Universal Sentence** Encoder
- **BERT Sentence** Embeddings
- **Sentence** Embeddings
- **Chunk** Embeddings
- Neural **Machine Translation**
- **Text-To-Text** Transfer Transformer (**Google T5**)
- Unsupervised **keywords extraction**
- Language **Detection & Identification** (up to 375 languages)
- Multi-class Text **Classification** (DL model)
- Multi-label Text **Classification** (DL model)
- Multi-class **Sentiment Analysis** (DL model)
- **Named entity** recognition (DL model)
- Easy **TensorFlow** integration
- **GPU** Support
- Full integration with **Spark ML** functions
- **710+** pre-trained **models** in **+192 languages!**
- **450+** pre-trained **pipelines** in **+192 languages!**

# 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).

# Natural Language Processing

Information Retrieval

Doc A



Doc 1

Doc 2

Doc 3

Sentiment Analysis



Information Extraction



Machine Translation



Question Answering



Human: When was Apollo sent to space?

Machine: First flight -  
AS-201,  
February 26,  
1966

# NLP Basics

## LEMMATIZATION

Find the **lemma** of each word:

- How does it show in the dictionary?

Uses a lookup from a full dictionary.

am, are, is → be

liver → liver

lives → live

## STEMMING

Find the **stem** of each word.

Uses rules: e.g, remove common suffixes.

Form	Suffix	Stem
studies	-es	studi
study <b>ing</b>	- <b>ing</b>	study
niñ <b>as</b>	- <b>as</b>	niñ
niñ <b>ez</b>	- <b>ez</b>	niñ

- The goal of both **stemming** and **lemmatization** is to reduce inflectional forms and sometimes derivationally related forms of a word to a common base form for normalization purposes.
- Lemmatization always returns real words, **stemming** doesn't.

# NLP Basics

## Remove stop words and apply stemming

it was a bright cold day in april  
and the clocks **were** striking  
thirteen winston smith **his** chin  
nuzzled **into his** breast in an  
effort **to** escape **the** vile wind  
slipped quickly **through the** glass  
doors **of** victory mansions though  
**not** quickly enough **to** prevent a  
swirl **of** gritty dust **from** entering  
along **with him**



bright cold day april clocks  
striking thirteen winston smith  
chin nuzzled breast effort  
escape vile wind slipped quickly  
glass doors victory mansions  
though quickly enough prevent  
swirl gritty dust entering along

- For tasks like text classification, where the text is to be classified into different categories, **stopwords** are **removed** or excluded from the given text so that more focus can be given to those words which define the meaning of the text.

### Stopwords

a  
able  
about  
above  
according  
accordingly  
across  
actually  
after  
afterwards  
again  
against  
ain  
all  
allow  
allows  
almost  
alone  
along  
already  
also

(520 stopwords)

# Spell Checking & Correction



```
val pipeline = PretrainedPipeline("spell_check_ml", "en")
val result = pipeline.annotate("Harry Potter is a graet muvie")

println(result("spell"))
/* will print Seq[String](..., "is", "a", "great", "movie") */
```

- 3 trainable approaches
- **Norvig Approach:**
  - Retrieves tokens and auto-corrects based on a given dictionary
- **Symmetric Delete:**
  - Uses distance metrics to find possible words
- **Context Aware:**
  - Most accurate: Judges words in context
  - Deep learning based

# Context Spell Checker

The Spell Checker can leverage the context of words for ranking different correction sequences. Let's take a look at some examples,

```
# check for the different occurrences of the word "siter"
example1 = ["I will call my siter.", \
            "Due to bad weather, we had to move to a different siter.", \
            "We travelled to three siter in the summer."]
beautify(lp.annotate(example1))
```

```
['I will call my sister .\n',
 'Due to bad weather , we had to move to a different site .\n',
 'We travelled to three sites in the summer .\n']
```

```
# check for the different occurrences of the word "ueather"
example2 = ["During the summer we have the best ueather.", \
            "I have a black ueather jacket, so nice.", \
            "I introduce you to my sister, she is called ueather."]
beautify(lp.annotate(example2))
```

```
['During the summer we have the best weather .\n',
 'I have a black leather jacket , so nice .\n',
 'I introduce you to my sister , she is called Heather .\n']
```

Notice that in the first example, 'siter' is indeed a valid English word,

<https://www.merriam-webster.com/dictionary/siter>

# NORMALIZATION

Remove or replace undesirable characters or regular expressions:

from: @Have a\$ #2great birth) day>!  
to: Have a great birth day!

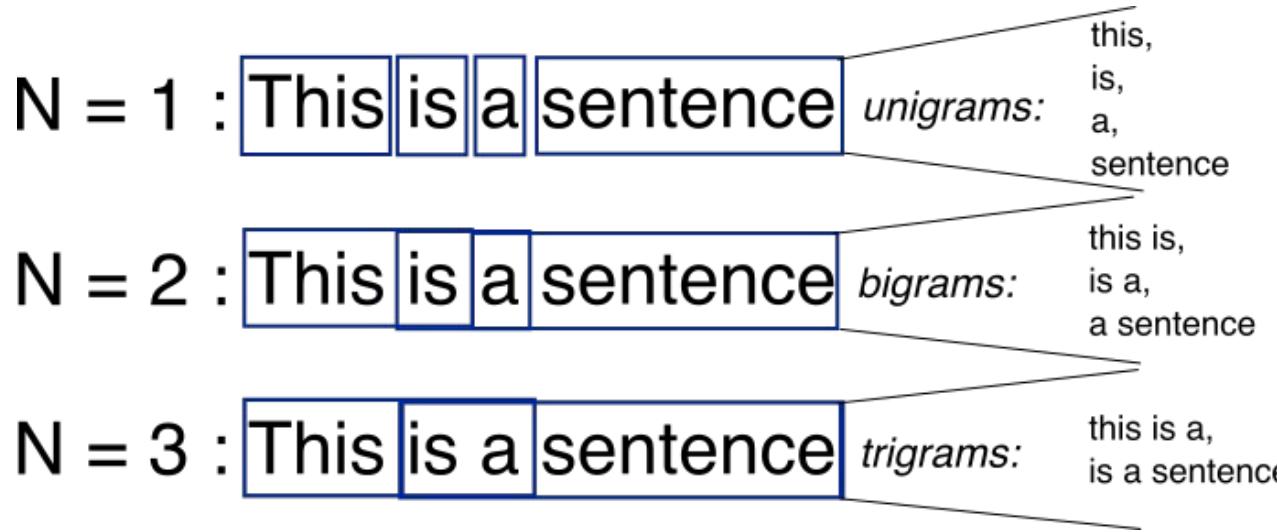
Spark NLP also comes with a Slang normalizer:

**Original tweet**  
@USER, r u cuming 2 MidCorner dis Sunday?

**Normalized tweet**

@USER, are you coming to MidCorner this Sunday?

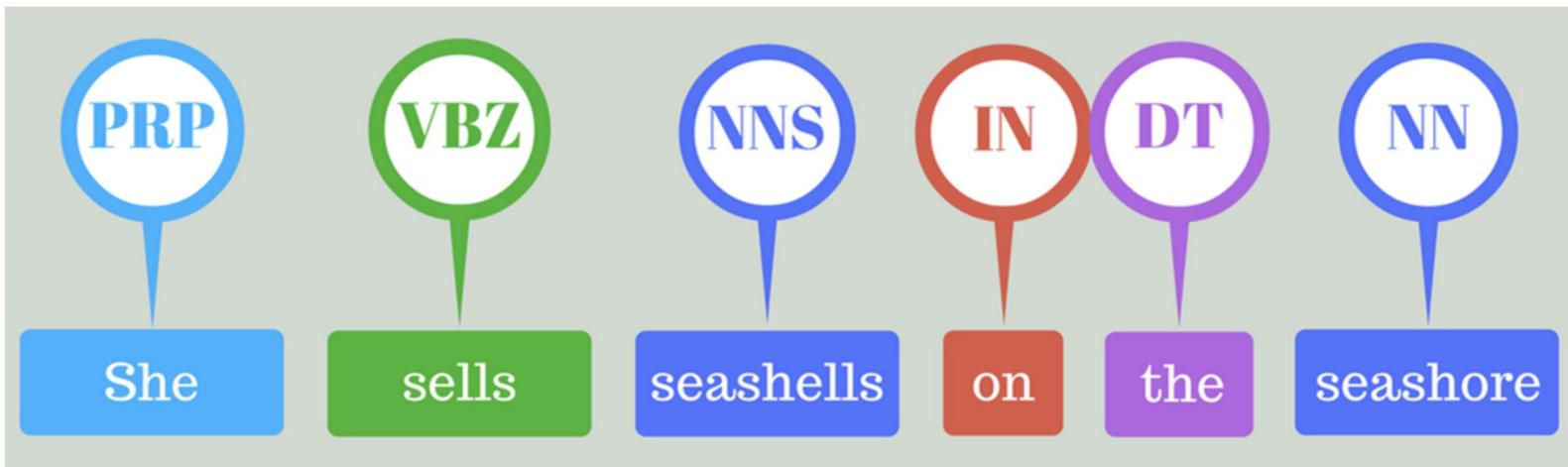
# N-gram Tokenization



- Kind of tokenizers which split words or sentences into several tokens
- Each token has certain number of characters
- Number of character depends on the type of ngram tokenizer
- Unigram, bigram, trigram, etc.

# PART OF SPEECH TAGGING

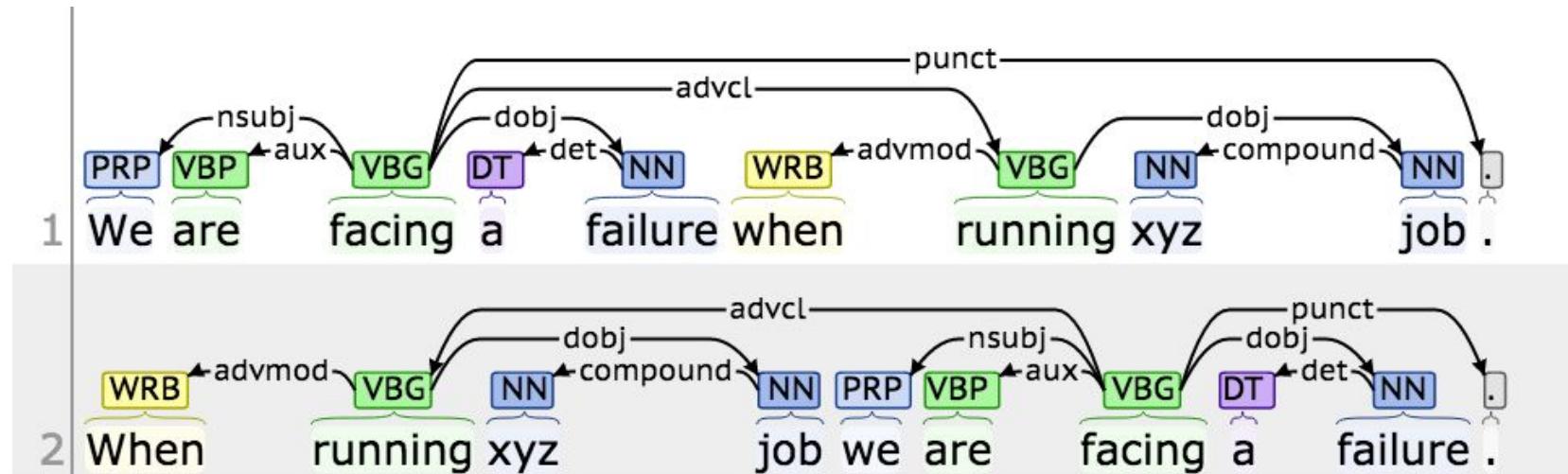
Often useful for recognizing named entities or word relationships.



A **POS tag** (or **part-of-speech tag**) is a special label assigned to each token (word) in a text corpus to indicate the **part of speech** and often also other grammatical categories such as tense, number (plural/singular), case etc.

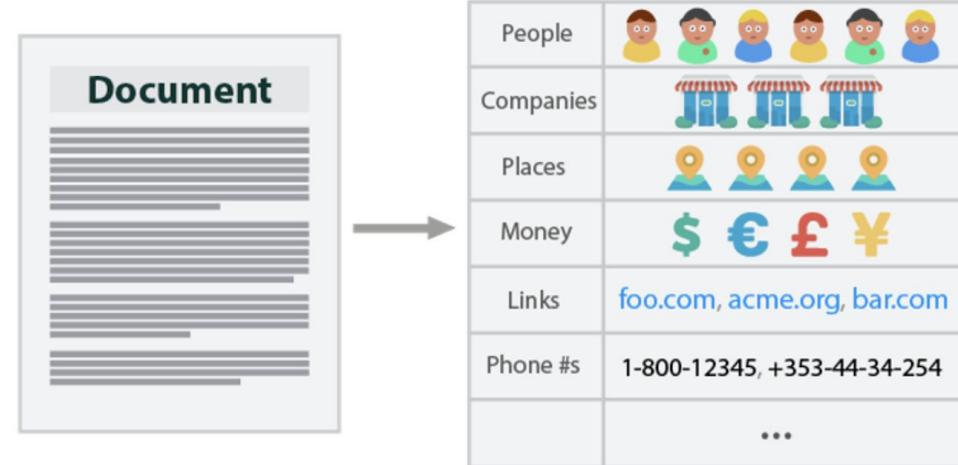
# DEPENDENCY PARSING

Useful for extracting relationships (i.e. building knowledge graphs):



# Named Entity Recognition (NER)

NER is a subtask of information extraction that seeks to **locate and classify named entity** mentioned in unstructured text into pre-defined categories such as **person names, organizations, locations, medical codes, time expressions, quantities, monetary values, percentages, etc.**



But Google **ORG** is starting from behind. The company made a late push into hardware, and Apple **ORG**'s Siri **PRODUCT**, available on iPhones **PRODUCT**, and Amazon **ORG**'s Alexa **PRODUCT** software, which runs on its Echo **PRODUCT** and Dot **PRODUCT** devices, have clear leads in consumer adoption.

# Word & Sentence Embeddings

## Vocabulary

index: Word:

0 aardvark  
1 able  
...

2409 black  
2410 bling  
...

3202 candid

3203 cast

3204 cat

...

5281 is

5282 island

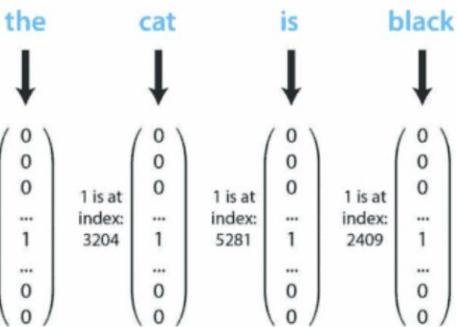
...

8676 the

8677 thing

...

9999 zombie



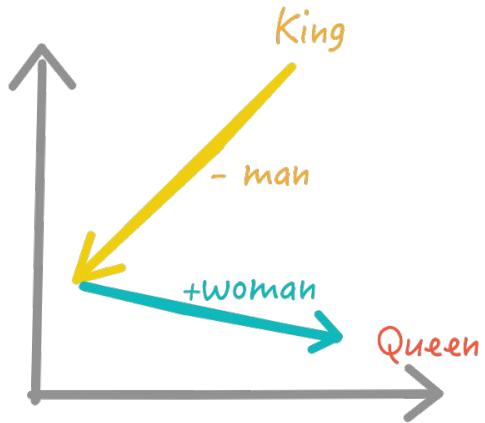
One-hot vector encoding for words in input sentence complete.

In [9]: doc[3].vector

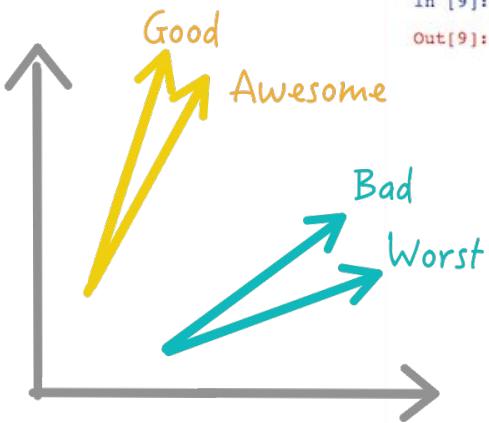
```
Out[9]: array([ 0.037103 , -0.31259 , -0.17857 ,  0.30001 ,  0.078154 ,
 0.17958 ,  0.12048 , -0.11879 , -0.20601 ,  1.2849 ,
-0.20409 ,  0.80613 ,  0.34344 , -0.19191 , -0.084511 ,
 0.17339 ,  0.042483 ,  2.0282 , -0.16278 , -0.60306 ,
-0.53766 ,  0.35711 ,  0.22882 ,  0.1171 ,  0.42983 ,
 0.16165 ,  0.407 ,  0.036476 ,  0.52636 , -0.13524 ,
-0.016897 ,  0.029259 , -0.079115 , -0.32305 ,  0.052255 ,
-0.3617 , -0.18355 , -0.34717 , -0.3691 ,  0.16881 ,
 0.21018 , -0.38376 , -0.096909 , -0.36296 , -0.37319 ,
 0.0021152,  0.32512 ,  0.063977 ,  0.36249 , -0.26935 ,
-0.59341 , -0.13625 ,  0.016425 , -0.2474 , -0.07498 ,
 0.034708 , -0.01476 , -0.11648 ,  0.25559 , -0.35002 ,
-0.52707 ,  0.21221 ,  0.062456 ,  0.26184 ,  0.53149 ,
 0.34957 , -0.22692 ,  0.44076 ,  0.4438 ,  0.6335 ,
-0.049757 , -0.08134 ,  0.65618 , -0.4716 ,  0.090675 ,
-0.084873 ,  0.31455 , -0.38495 , -0.19247 ,  0.48064 ,
 0.26688 ,  0.095743 ,  0.13024 ,  0.37023 ,  0.46269 ,
-0.32844 ,  0.17375 , -0.36325 ,  0.30672 , -0.075042 ,
-0.64684 , -0.49822 ,  0.12372 , -0.28547 ,  0.61811 ,
-0.19228 ,  0.0040473 ,  0.1774 ,  0.033154 , -0.54862 ,
 0.34695 , -0.53506 , -0.013381 ,  0.085712 , -0.054447 ,
-0.64673 ,  0.016749 ,  0.47676 ,  0.037803 , -0.10066 ,
-0.4165 , -0.20252 ,  0.2794 ,  0.10852 , -0.40154 ])
```

- Words that are used in similar contexts will be given similar representations. That is, words that are used in similar ways will be placed close together within the high-dimensional semantic space—these points will cluster together, and their distance to each other will be low.

# Word & Sentence Embeddings



a) Learns Analogy



b) Similar Words have same angles

```
In [9]: doc[3].vector
```

```
Out[9]: array([ 0.037103 , -0.31259 , -0.17857 ,  0.30001 ,  0.078154 ,  
  0.17958 ,  0.12048 , -0.11879 , -0.20601 ,  1.2849 ,  
 -0.20409 ,  0.80613 ,  0.34344 , -0.19191 , -0.084511 ,  
  0.17339 ,  0.042483 ,  2.0282 , -0.16278 , -0.60306 ,  
 -0.53766 ,  0.35711 ,  0.22882 ,  0.1171 ,  0.42983 ,  
  0.16165 ,  0.407 ,  0.036476 ,  0.52636 , -0.13524 ,  
 -0.016897 ,  0.029259 , -0.079115 , -0.32305 ,  0.052255 ,  
 -0.3617 , -0.18355 , -0.34717 , -0.3691 ,  0.16881 ,  
  0.21018 , -0.38376 , -0.096909 , -0.36296 , -0.37319 ,  
  0.0021152,  0.32512 ,  0.063977 ,  0.36249 , -0.26935 ,  
 -0.59341 , -0.13625 ,  0.016425 , -0.2474 , -0.07498 ,  
  0.034708 , -0.01476 , -0.11648 ,  0.25559 , -0.35002 ,  
 -0.52707 ,  0.21221 ,  0.062456 ,  0.26184 ,  0.53149 ,  
  0.34957 , -0.22692 ,  0.44076 ,  0.4438 ,  0.6335 ,  
 -0.049757 , -0.08134 ,  0.65618 , -0.4716 ,  0.090675 ,  
 -0.084873 ,  0.31455 , -0.38495 , -0.19247 ,  0.48064 ,  
  0.26688 ,  0.095743 ,  0.13024 ,  0.37023 ,  0.46269 ,  
 -0.32844 ,  0.17375 , -0.36325 ,  0.30672 , -0.075042 ,  
 -0.64684 , -0.49822 ,  0.12372 , -0.28547 ,  0.61811 ,  
 -0.19228 ,  0.0040473 ,  0.1774 ,  0.033154 , -0.54862 ,  
  0.34695 , -0.53506 , -0.013381 ,  0.085712 , -0.054447 ,  
 -0.64673 ,  0.016749 ,  0.47676 ,  0.037803 , -0.10066 ,  
 -0.4165 , -0.20252 ,  0.2794 ,  0.10852 , -0.40154 ])
```

- Deep-Learning-based natural language processing systems.
- They encode **words** and **sentences** in fixed-length dense vectors to drastically improve the processing of textual data.
- Based on **The Distributional Hypothesis**: Words that occur in the same contexts tend to have similar meanings.

# Word & Sentence Embeddings

Glove  
(100, 200, 300)

ELMO  
(512, 1024)

BERT  
(768d)

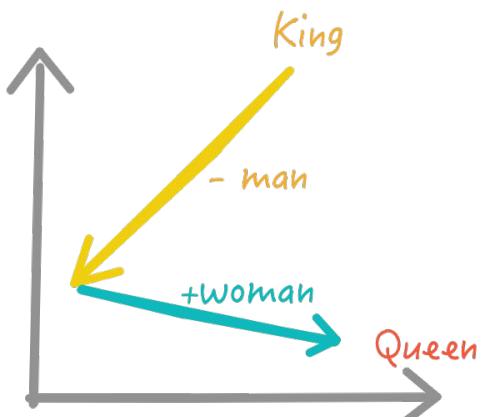
Universal Sentence Encoders  
(512)

Albert  
(768, 1024, 2048, 4096)

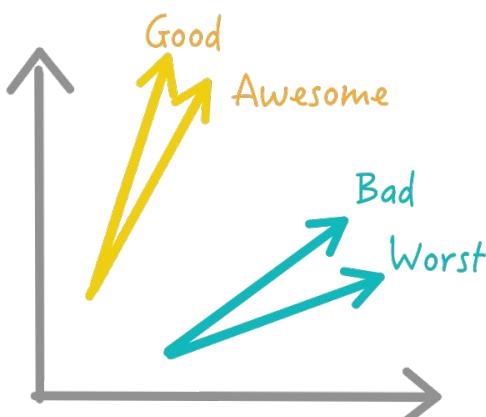
XLNet  
(768, 1024)

Electra  
(768)

Bert Sentence Embeddings  
(768)



a) Learns Analogy



b) Similar Words have same angles

- Deep-Learning-based natural language processing systems.
- They encode **words** and **sentences** in fixed-length dense vectors to drastically improve the processing of textual data.
- Based on **The Distributional Hypothesis**: Words that occur in the same contexts tend to have similar meanings.
- Elmo and Bert-family embeddings are context-aware.

# Text Classification with Word & Sentence Embeddings

Glove  
(100, 200, 300)

ELMO  
(512, 1024)

BERT  
(768d)

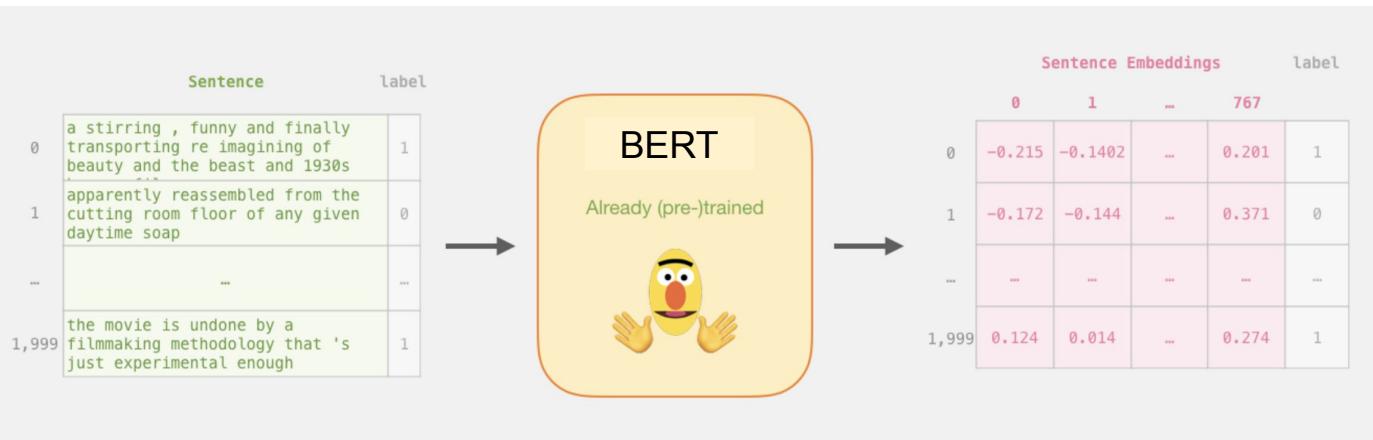
Universal Sentence Encoders  
(512)

Albert  
(768, 1024, 2048, 4096)

XLNet  
(768, 1024)

Electra  
(768)

Bert Sentence Embeddings  
(768)



Model	Name		BertSentenceEmbeddings	
WordEmbeddingsModel (GloVe)	glove_100d	small_bert_L2_512	BertSentenceEmbeddings	sent_small_bert_L6_128
BertEmbeddings	electra_small_uncased	small_bert_L4_512	BertSentenceEmbeddings	sent_small_bert_L8_128
BertEmbeddings	electra_base_uncased	small_bert_L6_512	BertSentenceEmbeddings	sent_small_bert_L10_128
BertEmbeddings	electra_large_uncased	small_bert_L8_512	BertSentenceEmbeddings	sent_small_bert_L12_128
BertEmbeddings	bert_base_uncased	small_bert_L10_512	BertSentenceEmbeddings	sent_small_bert_L2_256
BertEmbeddings	bert_base_cased	small_bert_L12_512	BertSentenceEmbeddings	sent_small_bert_L4_256
BertEmbeddings	bert_large_uncased	small_bert_L2_768	BertSentenceEmbeddings	sent_small_bert_L6_256
BertEmbeddings	bert_large_cased	small_bert_L4_768	BertSentenceEmbeddings	sent_small_bert_L8_256
BertEmbeddings	biobert_pubmed_base_cased	small_bert_L6_768	BertSentenceEmbeddings	sent_small_bert_L10_256
BertEmbeddings	biobert_pubmed_large_cased	small_bert_L8_768	BertSentenceEmbeddings	sent_small_bert_L12_256
BertEmbeddings	biobert_pmc_base_cased	small_bert_L10_768	BertSentenceEmbeddings	sent_small_bert_L2_512
BertEmbeddings	biobert_pubmed_pmc_base_cased	elmo	BertSentenceEmbeddings	sent_small_bert_L4_512
BertEmbeddings	biobert_clinical_base_cased	albert_base_uncased	BertSentenceEmbeddings	sent_small_bert_L6_512
BertEmbeddings	biobert_discharge_base_cased	albert_large_uncased	BertSentenceEmbeddings	sent_small_bert_L8_512
BertEmbeddings	covidbert_large_uncased	albert_xlarge_uncased	BertSentenceEmbeddings	sent_small_bert_L10_512
BertEmbeddings	small_bert_L2_128	albert_xxlarge_uncased	BertSentenceEmbeddings	sent_small_bert_L12_512
BertEmbeddings	small_bert_L4_128	xlnet_base_cased	BertSentenceEmbeddings	sent_small_bert_L2_768
BertEmbeddings	small_bert_L6_128	xlnet_large_cased	BertSentenceEmbeddings	sent_small_bert_L4_768
BertEmbeddings	small_bert_L8_128	tfhub_use	BertSentenceEmbeddings	sent_small_bert_L6_768
		tfhub_use_lg	BertSentenceEmbeddings	sent_small_bert_L8_768
			BertSentenceEmbeddings	sent_small_bert_L10_768
			BertSentenceEmbeddings	sent_small_bert_L12_768

# Word & Sentence Embeddings

albert\_base = [https://tfhub.dev/google/albert\\_base/3](https://tfhub.dev/google/albert_base/3) |  
768-embed-dim, 12-layer, 12-heads, 12M parameters

albert\_large = [https://tfhub.dev/google/albert\\_large/3](https://tfhub.dev/google/albert_large/3) |  
1024-embed-dim, 24-layer, 16-heads, 18M parameters

albert\_xlarge = [https://tfhub.dev/google/albert\\_xlarge/3](https://tfhub.dev/google/albert_xlarge/3) |  
2048-embed-dim, 24-layer, 32-heads, 60M parameters

albert\_xxlarge = [https://tfhub.dev/google/albert\\_xxlarge/3](https://tfhub.dev/google/albert_xxlarge/3) |  
4096-embed-dim, 12-layer, 64-heads, 235M parameters

XLNet-Large =  
[https://storage.googleapis.com/xlnet/released\\_models/cased\\_L-24\\_H-1024\\_A-16.zip](https://storage.googleapis.com/xlnet/released_models/cased_L-24_H-1024_A-16.zip) | 24-layer, 1024-hidden, 16-heads

XLNet-Base =  
[https://storage.googleapis.com/xlnet/released\\_models/cased\\_L-12\\_H-768\\_A-12.zip](https://storage.googleapis.com/xlnet/released_models/cased_L-12_H-768_A-12.zip) | 12-layer, 768-hidden, 12-heads.

<https://nlp.johnsnowlabs.com/api/>

**BERT** is a bi-directional transformer for pre-training over a lot of unlabeled textual data to learn a language representation that can be used to fine-tune for specific machine learning tasks. While BERT outperformed the NLP state-of-the-art on several challenging tasks, its performance improvement could be attributed to the bidirectional transformer, novel pre-training tasks of Masked Language Model and Next Structure Prediction along with a lot of data and Google's compute power.

**XLNet** is a large bidirectional transformer that uses improved training methodology, larger data and more computational power to achieve better than BERT prediction metrics on 20 language tasks.

To improve the training, XLNet introduces permutation language modeling, where all tokens are predicted but in random order. This is in contrast to BERT's masked language model where only the masked (15%) tokens are predicted.

**Albert** is a Google's new "ALBERT" language model and achieved state-of-the-art results on three popular benchmark tests for natural language understanding (NLU): GLUE, RACE, and SQuAD 2.0. ALBERT is a "lite" version of Google's 2018 NLU pretraining method BERT. Researchers introduced two parameter-reduction techniques in ALBERT to lower memory consumption and increase training speed.

# Coding ...

## 1. Spark NLP Basics

## 2. Text Preprocessing with Spark NLP

## 3. Spark NLP Pretrained Models

(click on Colab icon or open in a new tab)

[https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification\\_Trainings/Public/1.SparkNLP\\_Basics.ipynb](https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification_Trainings/Public/1.SparkNLP_Basics.ipynb)

[https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification\\_Trainings/Public/2.Text\\_Preprocessing\\_with\\_SparkNLP\(Annotators\\_Transformers\).ipynb](https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification_Trainings/Public/2.Text_Preprocessing_with_SparkNLP(Annotators_Transformers).ipynb)

[https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification\\_Trainings/Public/3.SparkNLP\\_Pretrained\\_Models.ipynb](https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification_Trainings/Public/3.SparkNLP_Pretrained_Models.ipynb)

### Spark NLP Basics and Pretrained Pipelines

- Colab Setup
- How to prevent Google Colab from disconnecting?
- Start Spark Session
- Using Pretrained Pipelines
  - Explain Document ML
  - Explain Document DL
  - Recognize Entities DL
  - Clean Stop Words
  - Clean Slang
  - Spell Checker
  - Spell Checker DL
  - Parsing a list of texts
  - Using fullAnnotate to get more details
  - Use pretrained match\_chunk Pipeline for Individual Noun Phrase
  - Extract exact dates from referential date phrases
  - Sentiment Analysis
    - Vivek algo
    - DL version (trained on imdb)
    - DL version (trained on twitter dataset)

### Text Preprocessing with Spark NLP

- Colab Setup
- Annotators and Transformer Concepts
- Create Spark Dataframe
- Transformers
- Document Assembler
- Sentence Detector
- Tokenizer
- Stacking Spark NLP Annotators in Spark ML Pipeline
- Normalizer
- Stopwords Cleaner
- Token Assembler
- Stemmer
- Lemmatizer
- NGram Generator
- TextMatcher
- RegexMatcher
- Text Cleaning with UDF
- Finisher
- LightPipeline

### Spark NLP Pretrained Models

- Colab Setup
- LemmatizerModel
- PerceptronModel (POS - Part of speech tags)
- Chunker
- Dependency Parser
- SpellChecker
  - Norvig Spell Checker
  - Context SpellChecker
- Language Detector
- Embeddings
  - Word Embeddings (Glove)
  - Using your own Word embeddings in Spark NLP
- Elmo Embeddings
- Bert Embeddings
- Albert Embeddings
- XlnetEmbeddings
- Chunk Embeddings
- UniversalSentenceEncoder
- Loading Models from local
- Getting Sentence Embeddings from word embeddings
  - Cosine similarity between two embeddings (sentence similarity)
- NERDL Model
  - Public NER (CoNLL 2003)
  - NerDL OntoNotes 100D
  - NER with Bert (Onto)
  - Getting the NER chunks with NER Converter
- Highlight the entities



# Spark NLP for Data Scientists

Session start at 13:30

Veysel Kocaman

Sr. Data Scientist

[veysel@johnsnowlabs.com](mailto:veysel@johnsnowlabs.com)

# Part - II

- ❖ Named Entity Recognition (NER) in Spark NLP

## CoNLL 2003 (English)

The CoNLL 2003 NER task consists of newswire text from the Reuters RCV1 corpus tagged with four different entity types (PER, LOC, ORG, MISC). Models are evaluated based on span-based F1 on the test set. \* used both the train and development splits for training.

Model	F1	Paper / Source	Code
CNN Large + fine-tune (Baevski et al., 2019)	93.5	Cloze-driven Pretraining of Self-attention Networks	
RNN-CRF+Flair	93.47	Improved Differentiable Architecture Search for Language Modeling and Named Entity Recognition	
LSTM-CRF+ELMo+BERT+Flair	93.38	Neural Architectures for Nested NER through Linearization	Official
Flair embeddings (Akbik et al., 2018)*	93.09	Contextual String Embeddings for Sequence Labeling	Flair framework
BERT Large (Devlin et al., 2018)	92.8	BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding	
CVT + Multi-Task (Clark et al., 2018)	92.61	Semi-Supervised Sequence Modeling with Cross-View Training	Official
BERT Base (Devlin et al., 2018)	92.4	BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding	
BILSTM-CRF+ELMo (Peters et al., 2018)	92.22	Deep contextualized word representations	AllenNLP Project AllenNLP GitHub
Peters et al. (2017) *	91.93	Semi-supervised sequence tagging with bidirectional language models	
CRF + AutoEncoder (Wu et al., 2018)	91.87	Evaluating the Utility of Hand-crafted Features in Sequence Labelling	Official
Bi-LSTM-CRF + Lexical Features (Ghadhar and Langlais 2018)	91.73	Robust Lexical Features for Improved Neural Network Named-Entity Recognition	Official
BILSTM-CRF + IntNet (Xin et al., 2018)	91.64	Learning Better Internal Structure of Words for Sequence Labeling	
Chiu and Nichols (2016) *	91.62	Named entity recognition with bidirectional LSTM-CNNs	

# NER-DL in Spark NLP

SYSTEM	YEAR	LANGUAGE	ACCURACY
Spark NLP v2.4	2020	Python/Scala/Java/R	93.3 (test F1) - 95.9 (dev F1)
Spark NLP v2.x	2019	Python/Scala/Java/R	93
Spark NLP v1.x	2018	Python/Scala/Java/R	92
spaCy v2.x	2017	Python/Cython	92.6
spaCy v1.x	2015	Python/Cython	91.8
ClearNLP	2015	Java	91.7
CoreNLP	2015	Java	89.6
MATE	2015	Java	92.5
Turbo	2015	C++	92.4

The best NER score in production

93.3 %  
Test Set

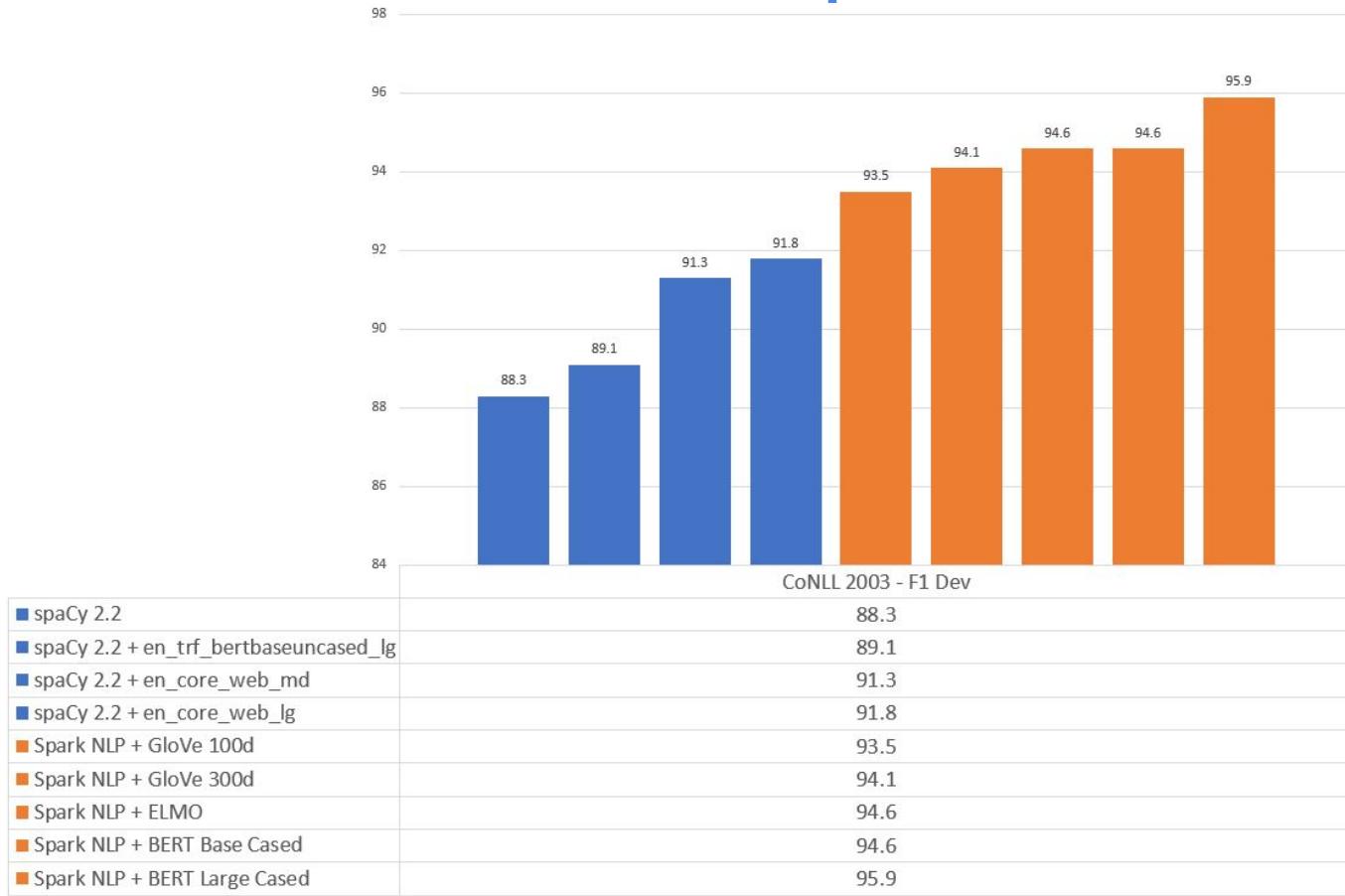


Bert



NerDLApproach

# NER-DL in Spark NLP



# NER Systems

Feature-engineered machine learning systems	Dict	SP	DU	EN	GE
Carreras et al. (2002) binary AdaBoost classifiers	Yes	81.39	77.05	-	-
Malouf (2002) - Maximum Entropy (ME) + features	Yes	73.66	68.08	-	-
Li et al. (2005) SVM with class weights	Yes	-	-	88.3	-
Passos et al. (2014) CRF	Yes	-	-	90.90	-
Ando and Zhang (2005a) Semi-supervised state of the art	No	-	-	89.31	75.27
Agerri and Rigau (2016)	Yes	<b>84.16</b>	<b>85.04</b>	<b>91.36</b>	<b>76.42</b>
Feature-inferring neural network word models					
Collobert et al. (2011) Vanilla NN +SLL / Conv-CRF	No	-	-	81.47	-
Huang et al. (2015) Bi-LSTM+CRF	No	-	-	84.26	-
Yan et al. (2016) Win-BiLSTM (English), FF (German) (Many fets)	Yes	-	-	88.91	<b>76.12</b>
Collobert et al. (2011) Conv-CRF (SENNNA+Gazetteer)	Yes	-	-	89.59	-
Huang et al. (2015) Bi-LSTM+CRF+ (SENNNA+Gazetteer)	Yes	-	-	<b>90.10</b>	-
Feature-inferring neural network character models					
Gillick et al. (2015) – BTS	No	<b>82.95</b>	<b>82.84</b>	<b>86.50</b>	<b>76.22</b>
Kuru et al. (2016) CharNER	No	82.18	79.36	84.52	70.12
Feature-inferring neural network word + character models					
Yang et al. (2017)	Yes	85.77	<b>85.19</b>	91.26	-
Luo (2015)	Yes	-	-	91.20	-
Chiu and Nichols (2015)	Yes	-	-	<b>91.62</b>	-
Ma and Hovy (2016)	No	-	-	91.21	-
Santos and Guimaraes (2015)	No	82.21	-	-	-
Lample et al. (2016)	No	85.75	81.74	90.94	<b>78.76</b>
Bharadwaj et al. (2016)	Yes	<b>85.81</b>	-	-	-
Dernoncourt et al. (2017)	No	-	-	90.5	-
Feature-inferring neural network word + character + affix models					
Re-implementation of Lample et al. (2016) (100 Epochs)	No	85.34	85.27	90.24	78.44
Yadav et al. (2018)(100 Epochs)	No	86.92	87.50	90.69	78.56
Yadav et al. (2018) (150 Epochs)	No	<b>87.26</b>	<b>87.54</b>	90.86	<b>79.01</b>

## 1. Classical Approaches (rule based)

## 2. ML Approaches

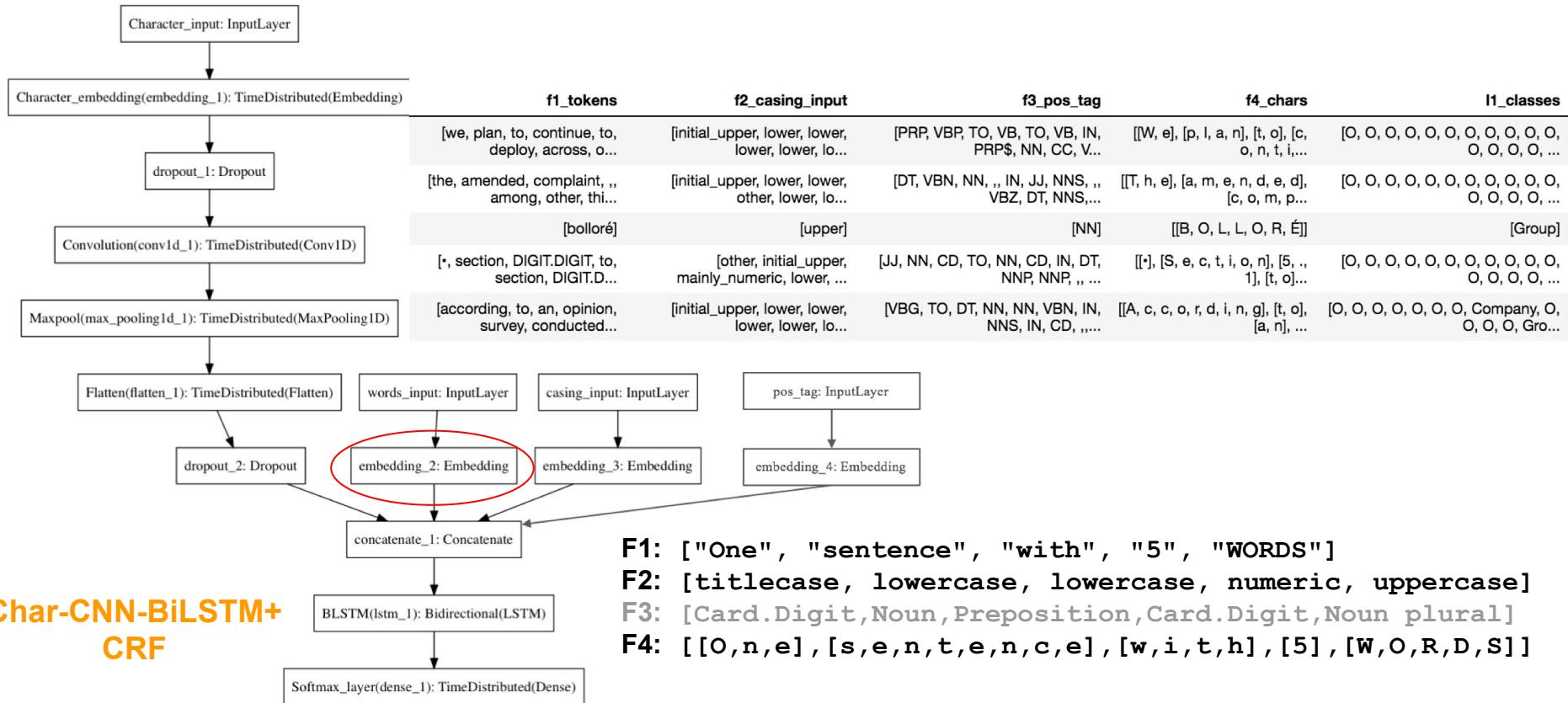
- Multi-class classification
- Conditional Random Field (CRF)

## 3. DL Approaches

- Bidirectional LSTM-CRF
- Bidirectional LSTM-CNNs
- Bidirectional LSTM-CNNS-CRF
- Pre-trained language models  
(Bert, Elmo)

## 4. Hybrid Approaches (DL + ML)

# NER-DL in Spark NLP



# NER-DL in Spark NLP

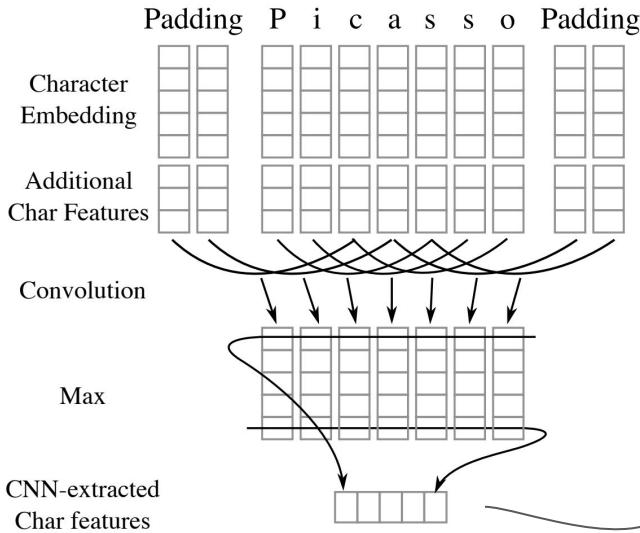


Figure 2: The convolutional neural network extracts character features from each word. The character embedding and (optionally) the character type feature vector are computed through lookup tables. Then, they are concatenated and passed into the CNN.

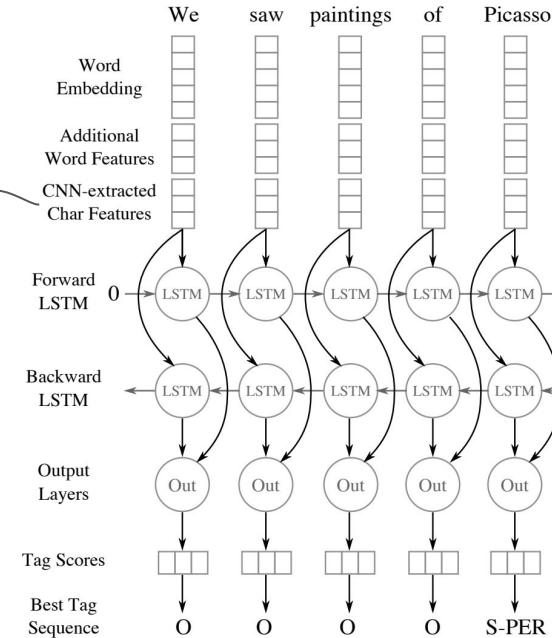


Figure 1: The (unrolled) BLSTM for tagging named entities. Multiple tables look up word-level feature vectors. The CNN (Figure 2) extracts a fixed length feature vector from character-level features. For each word, these vectors are concatenated and fed to the BLSTM network and then to the output layers (Figure 3).

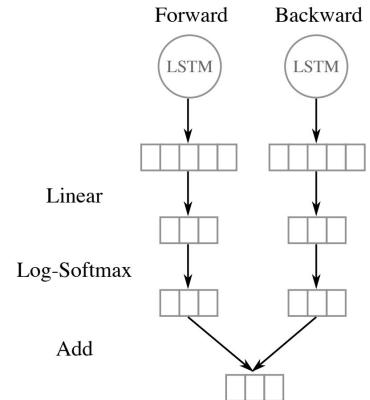


Figure 3: The output layers (“Out” in Figure 1) decode output into a score for each tag category.

## Char-CNN-BiLSTM

# Clinical Named Entity Recognition

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, associated with an acute hepatitis, and obesity with a body mass index (BMI) of 33.5 kg/m<sup>2</sup>, presented with a one-week history of polyuria, polydipsia, poor appetite, and vomiting. Two weeks prior to presentation, she was treated with a five-day course of amoxicillin for a respiratory tract infection. 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, bicarbonate 18 mmol/L, anion gap 20, creatinine 0.4 mg/dL, triglycerides 508 mg/dL, total cholesterol 122 mg/dL, glycated hemoglobin (HbA1c) 10%, and venous pH 7.27. Serum lipase was normal at 43 U/L. Serum acetone levels could not be assessed as blood samples kept hemolyzing due to significant lipemia. The patient was initially admitted for starvation ketosis, as she reported poor oral intake for three days prior to admission. However, serum chemistry obtained six hours after presentation revealed her glucose was 186 mg/dL, the anion gap was still elevated at 21, serum bicarbonate was 16 mmol/L, triglyceride level peaked at 2050 mg/dL, and lipase was 52 U/L. The β-hydroxybutyrate level was obtained and found to be elevated at 5.29 mmol/L - the original sample was centrifuged and the chylomicron layer removed prior to analysis due to interference from turbidity caused by lipemia again.

## Clinical NER

Color codes: PROBLEM, TREATMENT, TEST,

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 from 3 months.

## Posology NER

Color codes: FREQUENCY, DOSAGE, DURATION, DRUG, FORM, STRENGTH,

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.

## Anatomy NER

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

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

# NER-DL in Spark NLP

## CoNLL2003 format

All data files contain one word per line with empty lines representing sentence boundaries. At the end of each line there is a tag which states whether the current word is inside a named entity or not. The tag also encodes the type of named entity. Here is an example sentence:

U.N.	NNP	I-NP	I-ORG
official	NN	I-NP	O
Ekeus	NNP	I-NP	I-PER
heads	VBZ	I-VP	O
for	IN	I-PP	O
Baghdad	NNP	I-NP	I-LOC
.	.	O	O

\* Each line contains four fields: the word, its part-of-speech tag, its chunk tag and its named entity tag.

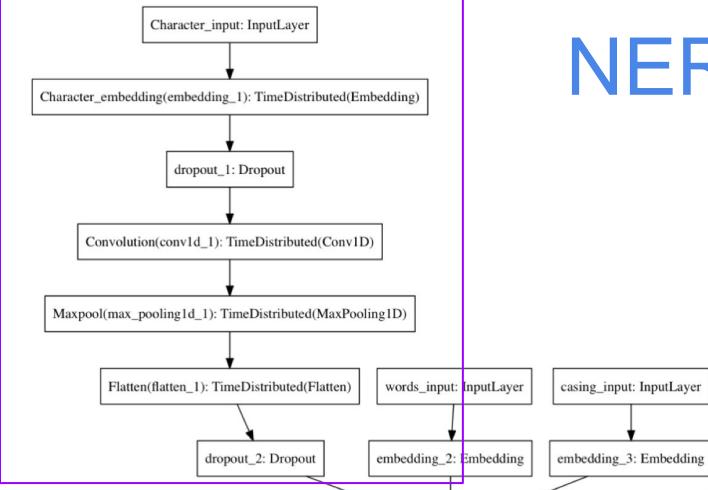
\* CoNLL: Conference on Computational Natural Language Learning

## BIO schema

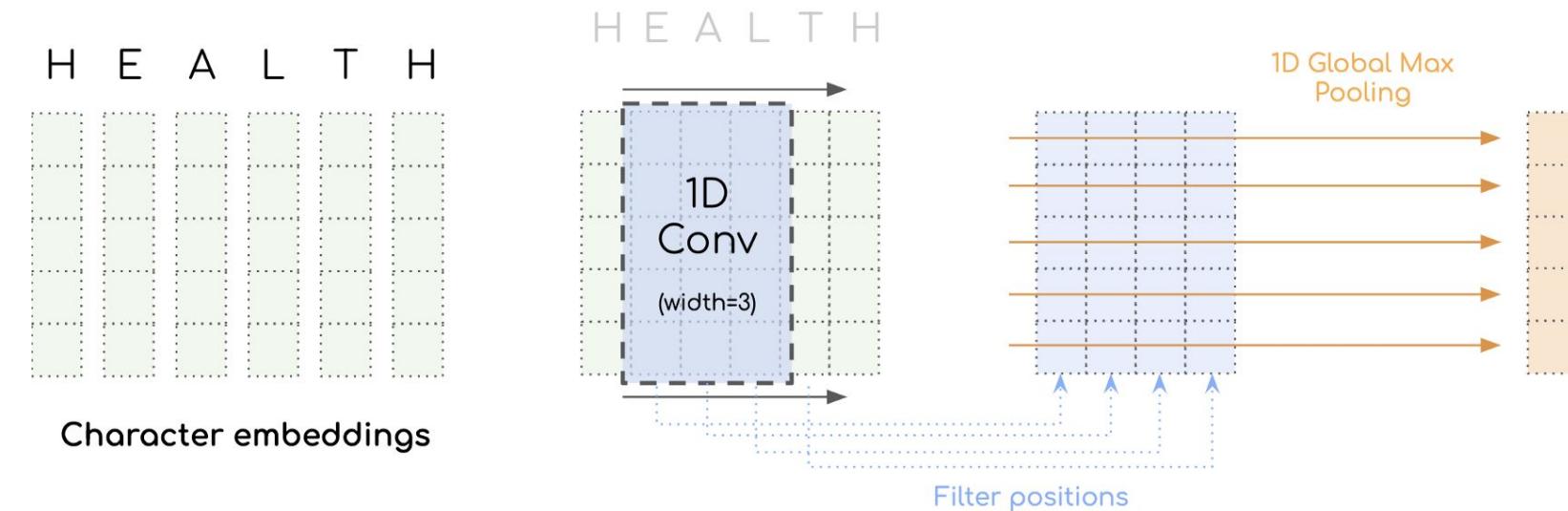
John	B-PER
Smith	I-PER
lives	O
in	O
New	B-LOC
York	I-LOC

John Smith ⇒ PERSON  
New York ⇒ LOCATION

# NER-DL in Spark NLP



Char-CNN process, e.g. on the world “HEALTH”



# 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

# NER-DL in Spark NLP

## Classification

	The	company	XYZ	Private	Limited	works	in	the	health	sector	in	Europe
GROUND TRUTH	O	O	Company	Company	Company	O	O	O	Activity	Activity	O	Location
PREDICTION	O	O	O	Company	Company	O	Group	O	O	Activity	O	Location

Class	O	Company	Location	Activity	Group
O	91969	546	295	1069	251
Company	569	3084	69	43	129
Location	137	48	1677	1	4
Activity	735	28	0	1329	0
Group	98	95	8	0	1185

Recall	Precision	F1
97,7 %	98,4 %	98 %
79,2 %	81,1 %	80,2 %
89,8 %	81,8 %	85,6 %
63,5 %	54,4 %	58,6 %
85,5 %	75,5 %	80,2 %

# Coding ...

Open 4. NERDL Training notebook in Colab

(click on Colab icon or open in a new tab)

[https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification\\_Trainings/Public/4.NERDL\\_Training.ipynb](https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification_Trainings/Public/4.NERDL_Training.ipynb)

## Test set evaluation

```
In [ ]: import pyspark.sql.functions as F  
  
predictions.select(F.explode(F.arrays_zip('token.result','label.result','ner.result')).alias("cols")) \  
.select(F.expr("cols['0']").alias("token"),  
        F.expr("cols['1']").alias("ground_truth"),  
        F.expr("cols['2']").alias("prediction")).show(truncate=False)
```

token	ground_truth	prediction
CRICKET	o	o
-	o	o
LEICESTERSHIRE	B-ORG	B-ORG
TAKE	o	o
OVER	o	o
AT	o	o
TOP	o	o
AFTER	o	o
INNINGS	o	o
VICTORY	o	o
.	o	o
LONDON	B-LOC	B-LOC
1996-08-30	o	o
West	B-MISC	B-MISC
Indian	I-MISC	I-MISC
all-rounder	o	o
Phil	B-PER	B-PER
Simmons	I-PER	I-PER
took	o	o
four	o	o

only showing top 20 rows

```
In [ ]: from sklearn.metrics import classification_report  
  
preds_df = predictions.select(F.explode(F.arrays_zip('token.result','label.result','ner.result')).alias("cols")) \  
.select(F.expr("cols['0']").alias("token"),  
        F.expr("cols['1']").alias("ground_truth"),  
        F.expr("cols['2']").alias("prediction")).toPandas()  
  
print (classification_report(preds_df['ground_truth'], preds_df['prediction']))
```

	precision	recall	f1-score	support
B-LOC	0.88	0.93	0.90	1837
B-MISC	0.80	0.82	0.81	922
B-ORG	0.92	0.73	0.81	1341
B-PER	0.94	0.95	0.95	1842



# Spark NLP for Data Scientists

Session start at 13:30

Veysel Kocaman

Sr. Data Scientist

[veysel@johnsnowlabs.com](mailto:veysel@johnsnowlabs.com)

# Part - III

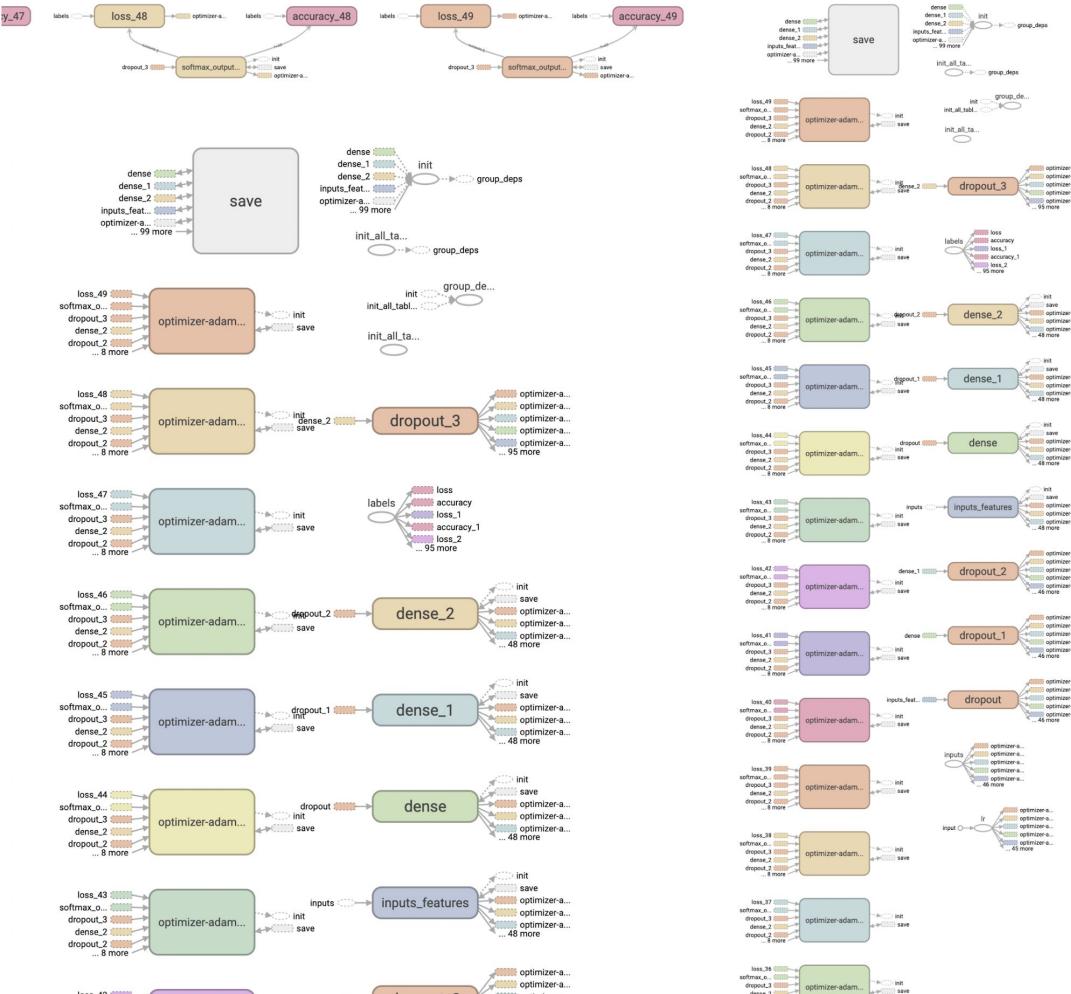
- ❖ Text Classification with Classifier DL in Spark NLP

# SentimentDL, ClassifierDL, and MultiClassifierDL

- BERT
  - Small BERT
  - BioBERT
  - CovidBERT
  - LaBSE
  - ALBERT
  - ELECTRA
  - XLNet
  - ELMO
  - Universal Sentence Encoder
  - GloVe
- 2 classes (positive/negative)
  - 3 classes (0, 1, 2)
  - 4 classes (Sports, Business, etc.)
  - 5 classes (1.0, 2.0, 3.0, 4.0, 5.0)
  - ... 100 classes!

- 100 dimensions
  - 200 dimensions
  - 128 dimensions
  - 256 dimensions
  - 300 dimensions
  - 512 dimensions
  - 768 dimensions
  - 1024 dimensions
- tfhub\_ues
  - tfhub\_use\_lg
  - glove\_6B\_100
  - glove\_6B\_300
  - glove\_840B\_300
  - bert\_base\_cased
  - bert\_base\_uncased
  - bert\_large\_cased
  - bert\_large\_uncased
  - bert\_multi\_uncased
  - electra\_small\_uncased
  - elmo
  - ... 90+ Word & Sentence models

# Classifier DL Tensorflow Architecture



# Coding ...

## Open 5. Text Classification with ClassifierDL notebook in Colab

(click on Colab icon or open in a new tab)

[https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification\\_Trainings/Public/4.NERDL\\_Training.ipynb](https://github.com/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/Certification_Trainings/Public/4.NERDL_Training.ipynb)

### ClassifierDL with Universal Sentence Embeddings

```
In [ ]: # actual content is inside description column
document = DocumentAssembler() \
    .setInputCol("description") \
    .setOutputCol("document")

# we can also use sentence detector here if we want to train on and get predictions for each sentence

use = UniversalSentenceEncoder.pretrained() \
    .setInputCols(["document"]) \
    .setOutputCol("sentence_embeddings")

# the classes/labels/categories are in category column
classifierdl = ClassifierDLApproach() \
    .setInputCols(["sentence_embeddings"]) \
    .setOutputCol("class") \
    .setLabelColumn("category") \
    .setMaxEpochs(5) \
    .setEnableOutputLogs(True)

use_clf_pipeline = Pipeline(
    stages = [
        document,
        use,
        classifierdl
    ]
)

tfhub_use.download started this may take some time.
Approximate size to download 923.7 MB
[OK!]

In [ ]: use_pipelineModel = use_clf_pipeline.fit(trainDataset)
# 5 epochs takes around 10 min

In [ ]: !cd ~/annotator_logs && ls -l
total 8
-rw-r--r-- 1 root root 533 Apr  7 17:14 ClassifierDLApproach_a0fdc9e970b.log
-rw-r--r-- 1 root root 976 Apr  7 16:59 ClassifierDLApproach_f222663dfb2c.log

In [ ]: !cat ~/annotator_logs/ClassifierDLApproach_a0fdc9e970b.log
Training started - total epochs: 5 - learning rate: 0.005 - batch size: 64 - training examples: 120000
Epoch 0/5 - 34.868680102%.2fs - loss: 1620.7466 - accuracy: 0.8803833 - batches: 1875
Epoch 1/5 - 35.627811455%.2fs - loss: 1604.4518 - accuracy: 0.8915333 - batches: 1875
Epoch 2/5 - 34.687788982%.2fs - loss: 1597.8773 - accuracy: 0.8966333 - batches: 1875
Epoch 3/5 - 34.629944711%.2fs - loss: 1593.4987 - accuracy: 0.900275 - batches: 1875
Epoch 4/5 - 34.714090256%.2fs - loss: 1590.3165 - accuracy: 0.90335834 - batches: 1875
```

# T5: Text-To-Text Transfer Transformer

President Franklin <M> born <M> January 1882.

Lily couldn't <M>. The waitress had brought the largest <M> of chocolate cake <M> seen.

Our <M> hand-picked and sun-dried <M> orchard in Georgia.

**T5**

D. Roosevelt was <M> in

believe her eyes <M> piece <M> she had ever

peaches are <M> at our

*Pre-training*

*Fine-tuning*

President Franklin D.  
Roosevelt was born  
in January 1882.

When was Franklin D.  
Roosevelt born?

**T5**

1882

1. Text Summarization
2. Question Answering
3. Translation
4. Sentiment analysis
5. Natural Language Inference
6. Coreference Resolution
7. Sentence Completion
8. Word Sense Disambiguation

# T5: Text-To-Text Transfer Transformer

1. Text Summarization
2. Question Answering
3. Translation
4. Sentiment analysis
5. Natural Language Inference
6. Coreference Resolution
7. Sentence Completion
8. Word Sense Disambiguation

Task Name	Explanation
<a href="#">1.CoLA</a>	Classify if a sentence is grammatical correct
<a href="#">2.RTE</a>	Classify whether if a statement can be deduced from a sentence
<a href="#">3.MNLI</a>	Classify for a hypothesis and premise whether they contradict or contradict each other or neither of both (3 class).
<a href="#">4.MRPC</a>	Classify whether a pair of sentences is a re-phrasing of each other (semantically equivalent)
<a href="#">5.QNLI</a>	Classify whether the answer to a question can be deduced from an answer candidate.
<a href="#">6.QQP</a>	Classify whether a pair of questions is a re-phrasing of each other (semantically equivalent)
<a href="#">7.SST2</a>	Classify the sentiment of a sentence as positive or negative
<a href="#">8.STSB</a>	Classify the sentiment of a sentence on a scale from 1 to 5 (21 Sentiment classes)
<a href="#">9.CB</a>	Classify for a premise and a hypothesis whether they contradict each other or not (binary).
<a href="#">10.COPA</a>	Classify for a question, premise, and 2 choices which choice the correct choice is (binary).
<a href="#">11.MultiRc</a>	Classify for a question, a paragraph of text, and an answer candidate, if the answer is correct (binary),
<a href="#">12.WiC</a>	Classify for a pair of sentences and a disambiguous word if the word has the same meaning in both sentences.
<a href="#">13.WSC/DPR</a>	Predict for an ambiguous pronoun in a sentence what it is referring to.
<a href="#">14.Summarization</a>	Summarize text into a shorter representation.
<a href="#">15.SQuAD</a>	Answer a question for a given context.
<a href="#">16.WMT1.</a>	Translate English to German
<a href="#">17.WMT2.</a>	Translate English to French
<a href="#">18.WMT3.</a>	Translate English to Romanian

# 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-enterprise/>

<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>



---

# How to find your way in Spark NLP Universe?

A guidebook that contains all the helpful blogs, GitHub repos, documents, and demos of Spark NLP.

# The Table of Contents

---

- Models Hub
- Streamlit Demos
- GitHub Repositories
- Content of the Certification Training
- Slack Channels
- Docs
- Blog Posts
- Youtube Channel

# Models Hub



→ Main Page

→ Models Hub

→ Description Pages

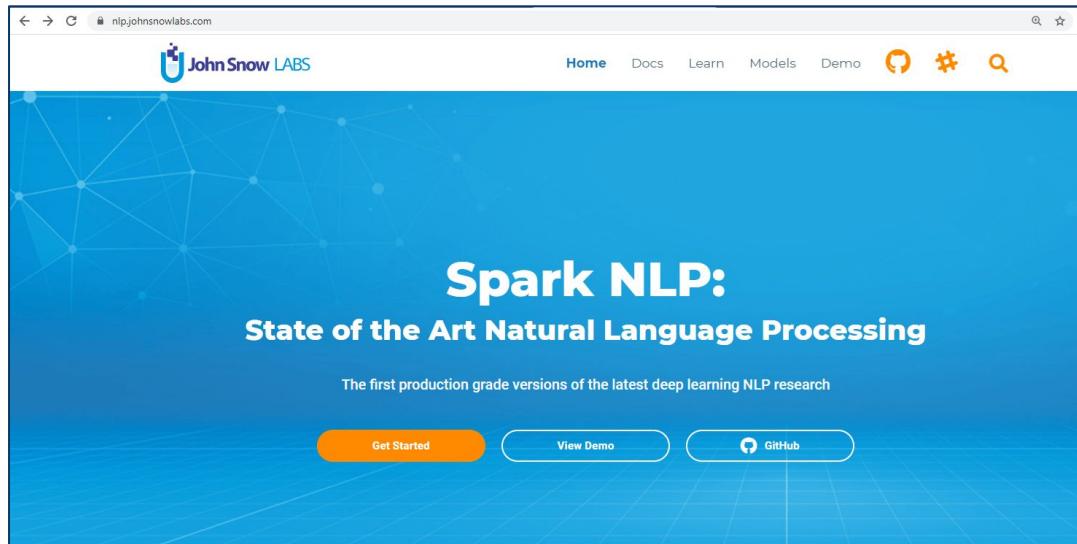
→ How to Use

→ Demo Pages

→ Colab Notebooks

→ Sample Code Snippets

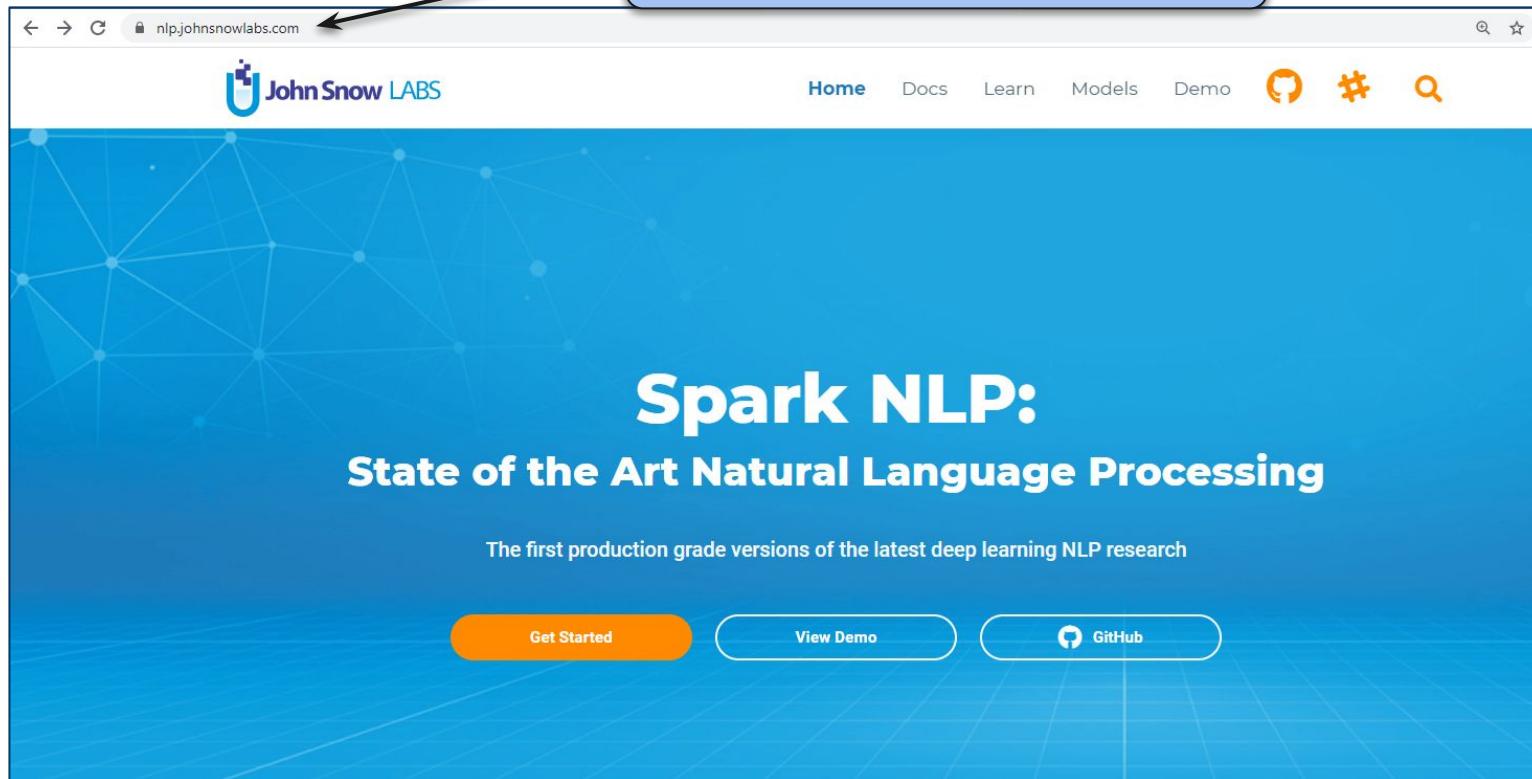
→ Benchmarking (Metrics)





# Main Page

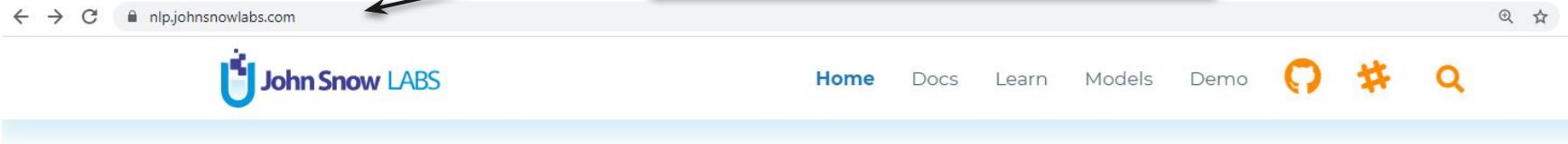
http://nlp.johnsnowlabs.com/



A screenshot of a web browser displaying the John Snow LABS main page. The URL "nlp.johnsnowlabs.com" is visible in the address bar. The page features a blue header with the "John Snow LABS" logo and navigation links for Home, Docs, Learn, Models, and Demo. Below the header is a large blue background image with a network graph. The main content area has a white background and displays the text "Spark NLP: State of the Art Natural Language Processing". Below this, a subtext reads "The first production grade versions of the latest deep learning NLP research". At the bottom are three orange call-to-action buttons: "Get Started", "View Demo", and "GitHub".

# Models Hub

http://nlp.johnsnowlabs.com/



## Right Out of The Box

Spark NLP ships with many **NLP features**, pre-trained **models** and **pipelines**

Models Hub

Models Hub

# Models Hub

Available Models and Pipelines

Show All 1036

aav<sup>4</sup> af<sup>5</sup> afa<sup>4</sup> alv<sup>4</sup> am<sup>2</sup> ar<sup>9</sup> art<sup>2</sup> ase<sup>2</sup> assertion<sup>2</sup> az<sup>4</sup> bat<sup>4</sup> bcl<sup>4</sup> bem<sup>4</sup> ber<sup>4</sup>  
bg<sup>7</sup> bh<sup>2</sup> bho<sup>2</sup> bi<sup>4</sup> bn<sup>5</sup> bnt<sup>4</sup> br<sup>3</sup> bzs<sup>4</sup> ca<sup>7</sup> cau<sup>2</sup> ccs<sup>2</sup> ceb<sup>4</sup> cel<sup>4</sup> chk<sup>4</sup> classifier<sup>19</sup>  
clinical<sup>101</sup> cn<sup>7</sup> cpf<sup>4</sup> cpp<sup>4</sup> crs<sup>4</sup> cs<sup>7</sup> cus<sup>4</sup> cy<sup>4</sup> da<sup>6</sup> de<sup>12</sup> deidentify<sup>4</sup> dra<sup>4</sup> ee<sup>4</sup> efi<sup>4</sup> el<sup>5</sup>  
embeddings<sup>55</sup> en<sup>809</sup> entity\_resolution<sup>29</sup> eo<sup>5</sup> es<sup>22</sup> et<sup>6</sup> eu<sup>7</sup> euq<sup>4</sup> fa<sup>5</sup> fi<sup>7</sup> fiu<sup>4</sup> fij<sup>4</sup> fr<sup>7</sup> ga<sup>7</sup>  
gaa<sup>4</sup> gem<sup>4</sup> gil<sup>4</sup> gl<sup>7</sup> gmq<sup>4</sup> gmw<sup>4</sup> grk<sup>4</sup> guw<sup>4</sup> gv<sup>4</sup> ha<sup>5</sup> he<sup>7</sup> hi<sup>7</sup> hil<sup>4</sup> ho<sup>4</sup> ht<sup>4</sup> hu<sup>7</sup>  
hy<sup>7</sup> id<sup>7</sup> ig<sup>4</sup> iir<sup>4</sup> ilo<sup>4</sup> inc<sup>4</sup> ine<sup>4</sup> is<sup>4</sup> iso<sup>4</sup> it<sup>7</sup> itc<sup>4</sup> ja<sup>7</sup> jap<sup>4</sup> ka<sup>2</sup> kab<sup>2</sup> kg<sup>4</sup> kj<sup>4</sup> kl<sup>2</sup>  
ko<sup>6</sup> kqn<sup>4</sup> kwn<sup>4</sup> kwy<sup>4</sup> la<sup>3</sup> language\_detection<sup>16</sup> lemmatizer<sup>41</sup> lg<sup>4</sup> licenced<sup>1</sup> licensed<sup>82</sup> ln<sup>4</sup> loz<sup>4</sup>  
lu<sup>4</sup> lua<sup>4</sup> lue<sup>4</sup> lun<sup>4</sup> luo<sup>4</sup> lus<sup>4</sup> lv<sup>5</sup> map<sup>2</sup> mfe<sup>4</sup> mg<sup>4</sup> mh<sup>4</sup> mk<sup>4</sup> mkh<sup>4</sup> ml<sup>4</sup> mos<sup>4</sup> mr<sup>7</sup>  
mt<sup>4</sup> mul<sup>4</sup> nb<sup>2</sup> ner<sup>79</sup> ng<sup>4</sup> nic<sup>4</sup> niu<sup>4</sup> nl<sup>9</sup> nn<sup>1</sup> nso<sup>4</sup> ny<sup>4</sup> nyk<sup>4</sup> om<sup>4</sup> open\_source<sup>810</sup> pa<sup>2</sup>  
pag<sup>4</sup> pap<sup>4</sup> phi<sup>4</sup> pipeline<sup>351</sup> pis<sup>4</sup> pl<sup>8</sup> pon<sup>4</sup> pos<sup>46</sup> poz<sup>2</sup> pqe<sup>4</sup> pqw<sup>2</sup> pt<sup>6</sup> question\_answering<sup>1</sup>  
re<sup>3</sup> relation\_extraction<sup>2</sup> relation\_extraction<sup>4</sup> rn<sup>4</sup> rnd<sup>4</sup> ro<sup>5</sup> roa<sup>4</sup> ru<sup>10</sup> run<sup>4</sup> rw<sup>4</sup> sal<sup>4</sup> sem<sup>4</sup>  
sentence\_detection<sup>5</sup> sentiment<sup>10</sup> seq2seq<sup>649</sup> sg<sup>4</sup> sit<sup>2</sup> sk<sup>7</sup> sl<sup>3</sup> sla<sup>4</sup> sm<sup>4</sup> sn<sup>4</sup> so<sup>1</sup> sq<sup>4</sup> srm<sup>2</sup>

# Description Page



2020/05/10/wikiner\_6B\_300\_pt.html

[Home](#) [Docs](#) [Learn](#) [Models](#) [Demo](#)

## Description

WikiNER is a Named Entity Recognition (or NER) model, meaning it annotates text to find features like the names of people, places, and organizations. This NER model does not read words directly but instead reads word embeddings, which represent words as points such that more semantically similar words are closer together. WikiNER 6B 300 is trained with GloVe 6B 300 word embeddings, so be sure to use the same embeddings in the pipeline.

## Predicted Entities

Persons- PER , Locations- LOC , Organizations- ORG , Miscellaneous- MISC .

[Live Demo](#)

[Open in Colab](#)

[Download](#)

May 18 Detect Persons, Locations, Organizations and Misc Entities in Polish (WikiNER 6B 300)

May 18 Detect Persons, Locations, Organizations and Misc Entities in Portuguese (WikiNER 6B 300)

May 18 Detect Persons, Locations, Organizations and Misc Entities in Dutch (WikiNER 840B 300)

May 18 Detect Persons, Locations, Organizations and Misc Entities in Polish (WikiNER 840B 300)

May 18 Detect Persons, Locations, Organizations and Misc Entities in Portuguese (WikiNER 840B 300)

... [View All Recent Posts](#)



## How to Use

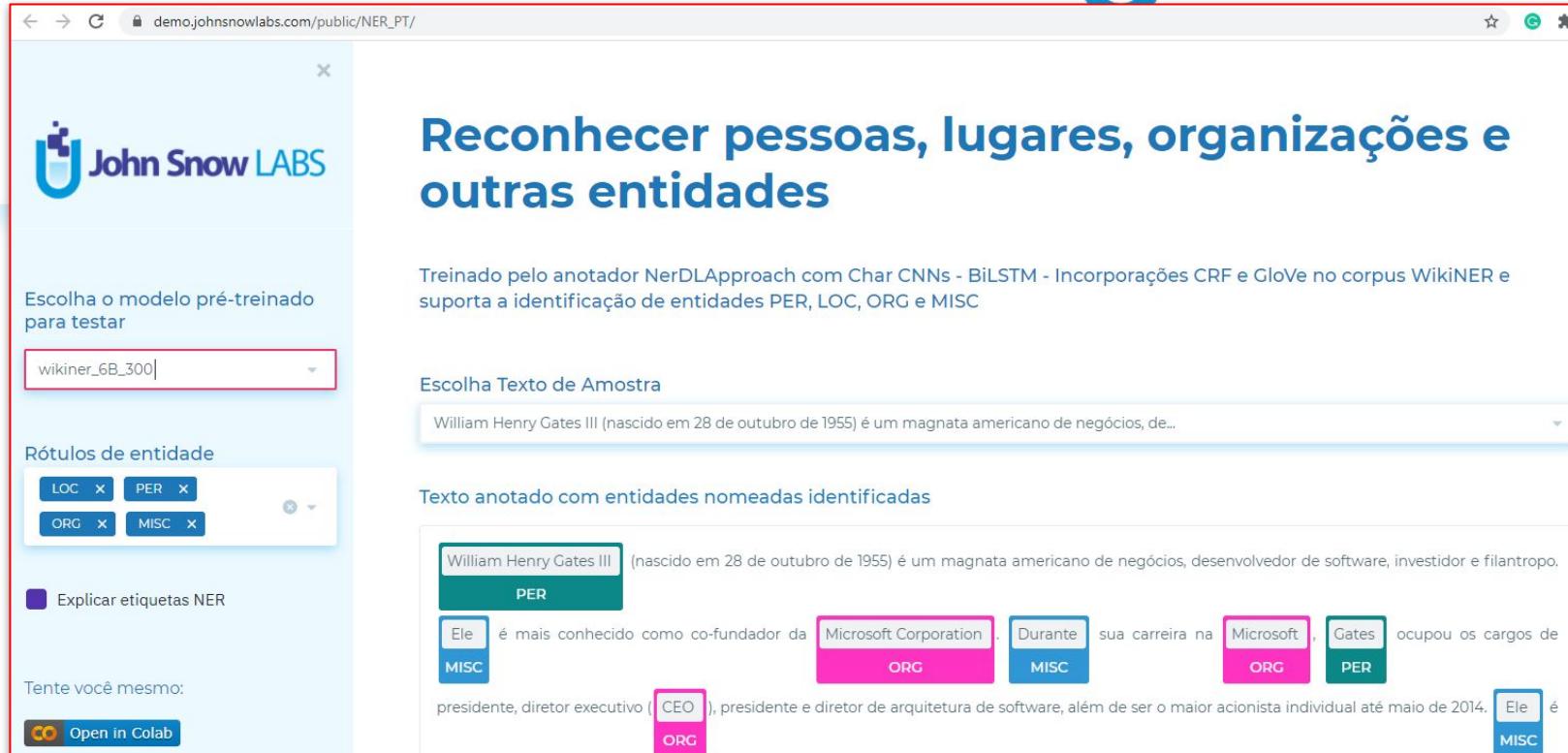
## How to

# Python

Scala

## How to use

demo.johnsnowlabs.com/public/NER\_PT/



Escolha o modelo pré-treinado para testar  
wikiner\_6B\_300

Rótulos de entidade  
LOC X PER X  
ORG X MISC X

Explicar etiquetas NER

Tente você mesmo:  
[Open in Colab](#)

Reconhecer pessoas, lugares, organizações e outras entidades

Treinado pelo anotador NerDLApproach com Char CNNs - BiLSTM - Incorporações CRF e GloVe no corpus WikiNER e suporta a identificação de entidades PER, LOC, ORG e MISC

Escolha Texto de Amostra

William Henry Gates III (nascido em 28 de outubro de 1955) é um magnata americano de negócios, de...

Texto anotado com entidades nomeadas identificadas

William Henry Gates III (nascido em 28 de outubro de 1955) é um magnata americano de negócios, desenvolvedor de software, investidor e filantropo.

Ele é mais conhecido como co-fundador da Microsoft Corporation. Durante sua carreira na Microsoft, Gates ocupou os cargos de presidente, diretor executivo (CEO), presidente e diretor de arquitetura de software, além de ser o maior acionista individual até maio de 2014. Ele é

Live Demo  Open in Colab  Download

Live Demo

 Open in Colab

 Download

# Colab Notebooks



colab.research.google.com/github/JohnSnowLabs/spark-nlp-workshop/blob/master/tutorials/streamlit\_notebooks/NER\_PT.ipynb

File Edit View Insert Runtime Tools Help

Table of contents

Detect entities in Portuguese text

Colab Setup

Start the Spark session

Select the DL model

Some sample examples

Define Spark NLP pipeline

Run the pipeline

Visualize results

+ Section

John Snow LABS

Open in Colab

## ▼ Detect entities in Portuguese text

### ▼ 1. Colab Setup

```
[ ] # Install Java
! apt-get update -qq
! apt-get install -y openjdk-8-jdk-headless -qq > /dev/null
! java -version

# Install pyspark
! pip install --ignore-installed -q pyspark==2.4.4
```

Live Demo

Open in Colab

Download

# Benchmarking (Metrics)

## Benchmarking

label	tp	fp	fn	prec	rec	f1
B-LOC	2081	157	142	0.9298481	0.93612236	0.9329747
I-ORG	1292	220	152	0.8544974	0.8947368	0.87415427
I-LOC	293	81	66	0.78342247	0.81615597	0.79945433
I-PER	1578	127	99	0.9255132	0.940966	0.9331757
B-ORG	1846	185	145	0.9089119	0.9271723	0.91795135
B-PER	3043	186	206	0.942397	0.93659586	0.9394875

tp: 10133 fp: 956 fn: 810 labels: 6

Macro-average prec: 0.890765, rec: 0.9086249, f1: 0.8996063

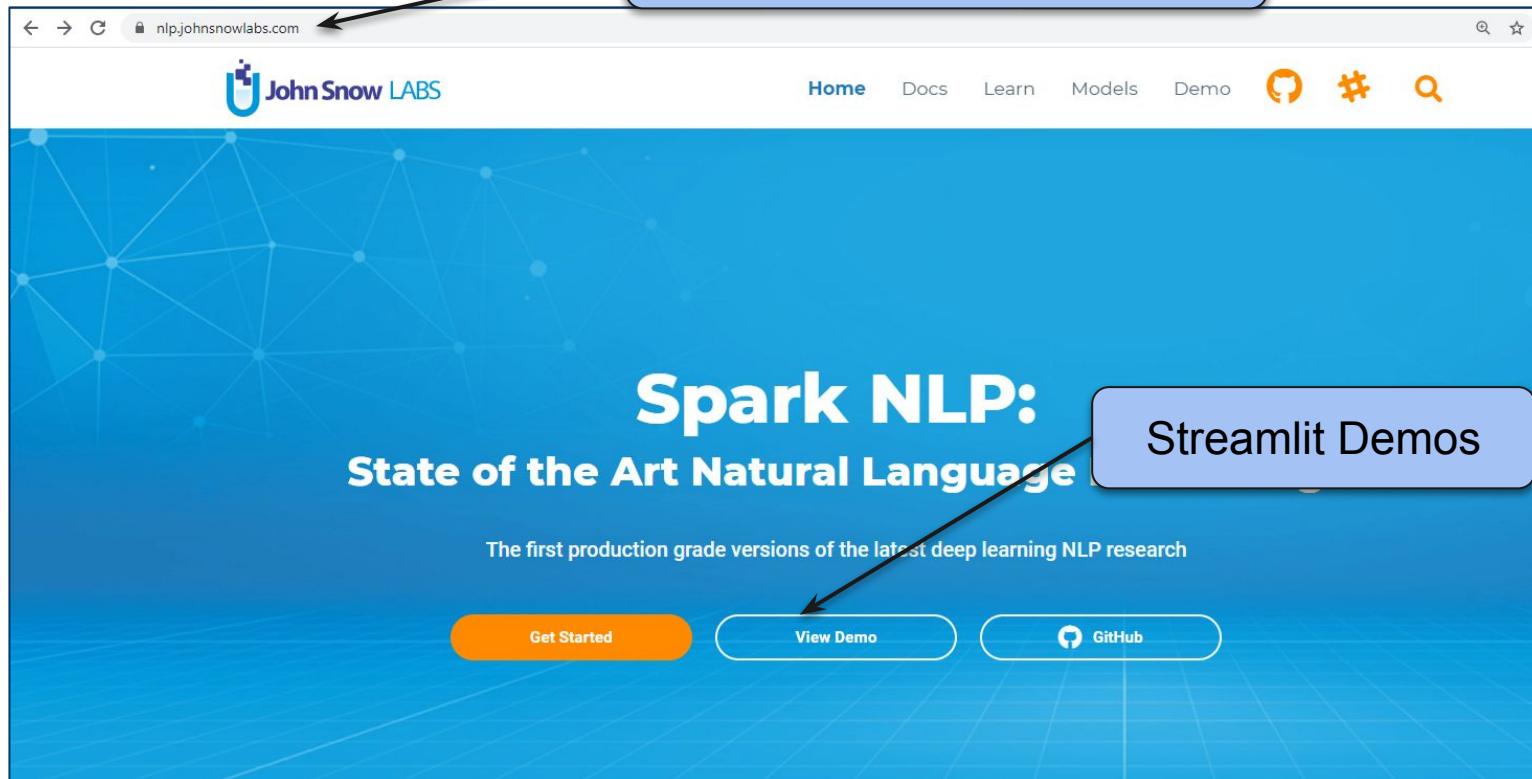
Micro-average prec: 0.91378844, rec: 0.9259801, f1: 0.9198439

# Streamlit Demos

- 
- Open Source
  - Languages
  - Healthcare
  - Spark OCR
  - De-Identification

# Main Page

http://nlp.johnsnowlabs.com/



The screenshot shows the main page of the John Snow LABS website at <http://nlp.johnsnowlabs.com/>. The page features a blue header with the logo and navigation links: Home, Docs, Learn, Models, Demo, GitHub, Hash, and Search. The main content area has a blue background with a network graph pattern. It prominently displays the text "Spark NLP: State of the Art Natural Language" and describes it as "The first production grade versions of the latest deep learning NLP research". Three calls-to-action are visible: an orange "Get Started" button, a white "View Demo" button, and a white "GitHub" button. A callout box highlights the "Streamlit Demos" link next to the "View Demo" button.

# Streamlit Demos

---

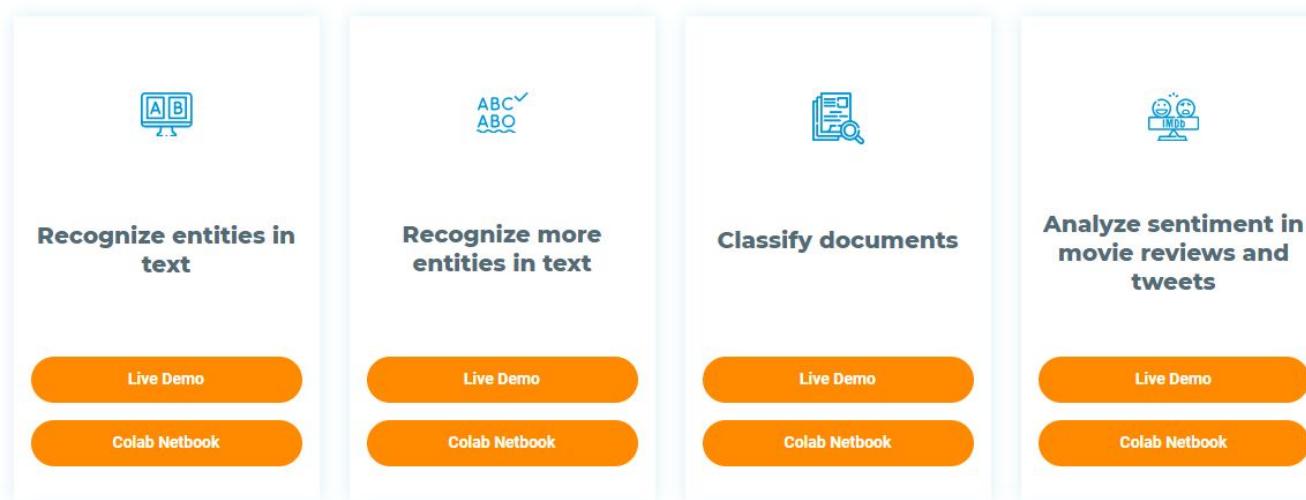
## Spark NLP in Action



# Open Source



## Open Source



**Recognize entities in text**

**Recognize more entities in text**

**Classify documents**

**Analyze sentiment in movie reviews and tweets**

[Live Demo](#)

[Colab Netbook](#)

[Live Demo](#)

[Colab Netbook](#)

[Live Demo](#)

[Colab Netbook](#)

[Live Demo](#)

[Colab Netbook](#)

# Languages

Open Source FREE

Languages FREE

Healthcare

Spark OCR

De-identification

## Languages



Detect language



Recognize entities in  
English text



Recognize entities in  
French text



Recognize entities in  
German text

[Live Demo](#)

[Live Demo](#)

[Live Demo](#)

[Live Demo](#)

[Colab Netbook](#)

[Colab Netbook](#)

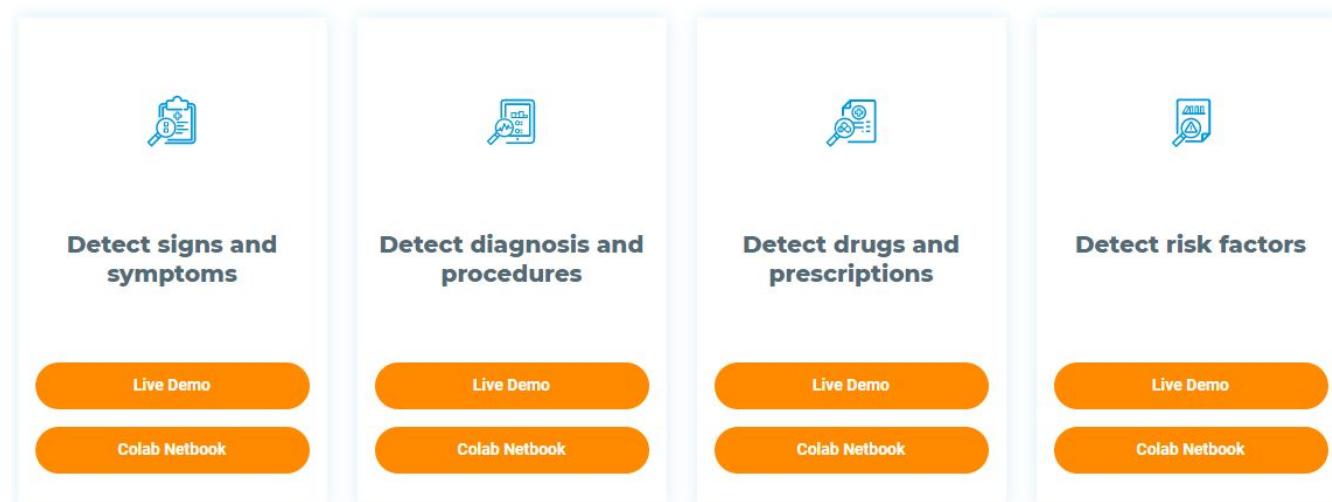
[Colab Netbook](#)

[Colab Netbook](#)

# Healthcare



## Healthcare



The page is divided into four main sections, each represented by a white card with rounded corners. Each card contains an icon of a magnifying glass over a clipboard or document, followed by a title and two orange buttons at the bottom.

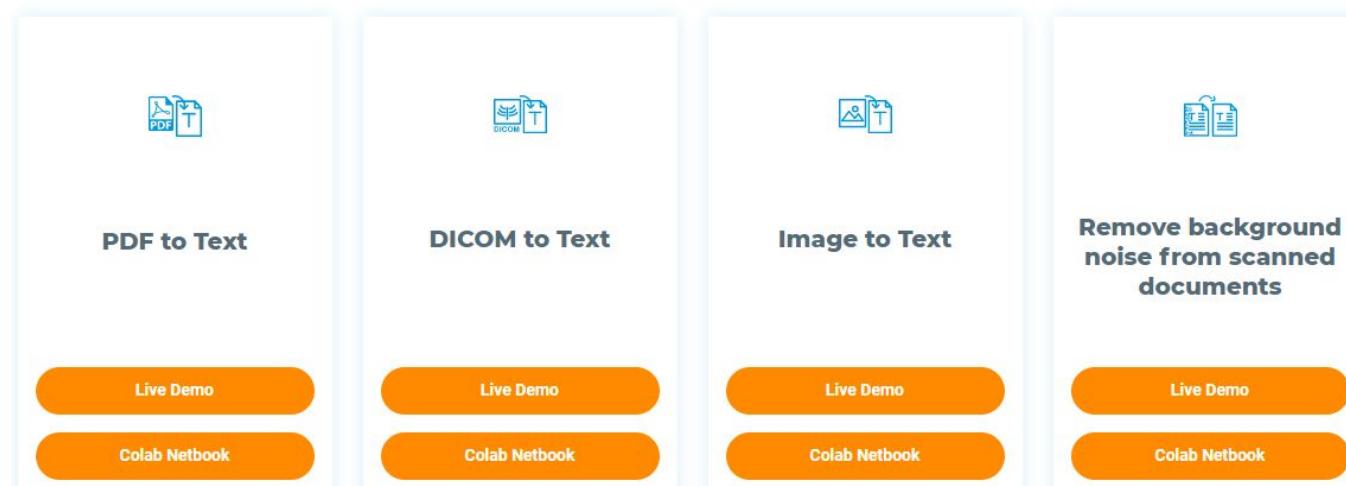
- Detect signs and symptoms**
- Detect diagnosis and procedures**
- Detect drugs and prescriptions**
- Detect risk factors**

**Live Demo** and **Colab Netbook** buttons are present in each section.

# Spark OCR



## Spark OCR



The page displays four main features of the Spark OCR service:

- PDF to Text**: Converts PDF files into text. Includes "Live Demo" and "Colab Netbook" buttons.
- DICOM to Text**: Converts DICOM files into text. Includes "Live Demo" and "Colab Netbook" buttons.
- Image to Text**: Converts images into text. Includes "Live Demo" and "Colab Netbook" buttons.
- Remove background noise from scanned documents**: A feature for cleaning scanned documents. Includes "Live Demo" and "Colab Netbook" buttons.

# De-Identification



## De-identification



**Deidentify structured data**



**Deidentify free text documents**



**Deidentify DICOM documents**



**De-identify PDF documents - HIPAA Compliance**

[Live Demo](#)

[Colab Netbook](#)

[Live Demo](#)

[Colab Netbook](#)

[Live Demo](#)

[Colab Netbook](#)

[Live Demo](#)

[Colab Netbook](#)

# GitHub Repositories



→ spark-nlp

→ spark-nlp-workshop

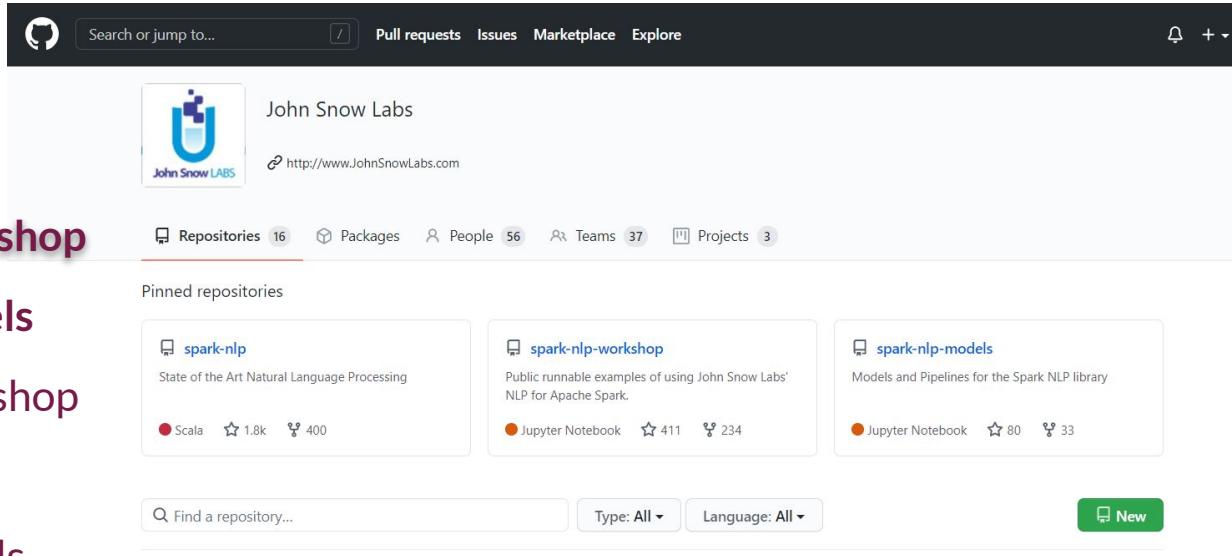
→ spark-nlp-models

→ spark-ocr-workshop

→ nlu

→ spark-nlp-models

→ spark-nlp-display

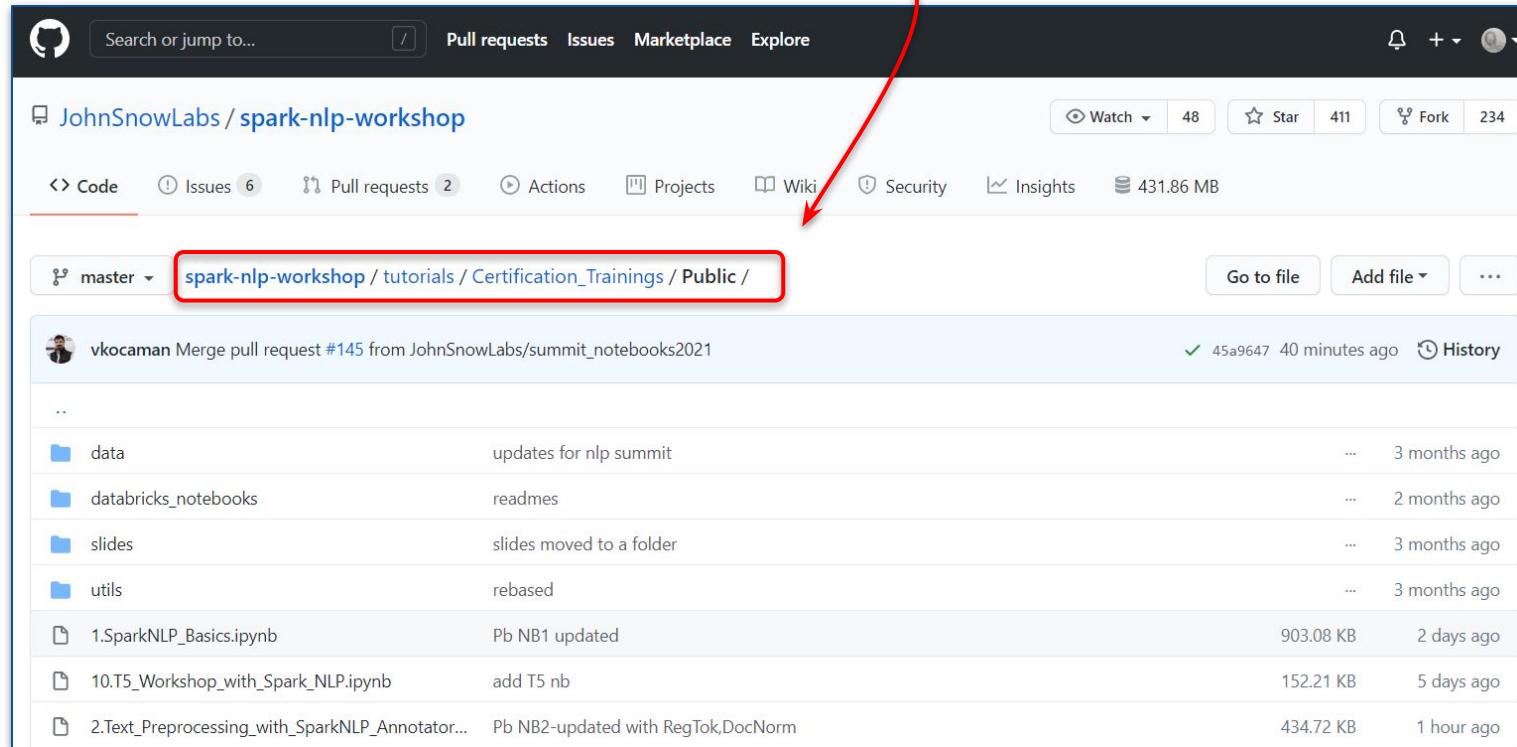


The screenshot shows the GitHub profile of 'John Snow Labs'. At the top, there's a search bar, a pull requests button, an issues button, a marketplace button, and an explore button. Below the header, the 'John Snow Labs' organization profile is displayed with its logo, name, and website link (<http://www.JohnSnowLabs.com>). The main content area shows three pinned repositories:

- spark-nlp**: State of the Art Natural Language Processing. Scala, 1.8k stars, 400 forks.
- spark-nlp-workshop**: Public runnable examples of using John Snow Labs' NLP for Apache Spark. Jupyter Notebook, 411 stars, 234 forks.
- spark-nlp-models**: Models and Pipelines for the Spark NLP library. Jupyter Notebook, 80 stars, 33 forks.

At the bottom of the page, there's a search bar with placeholder text 'Find a repository...', dropdown menus for 'Type: All' and 'Language: All', and a green 'New' button.

# Certification Training



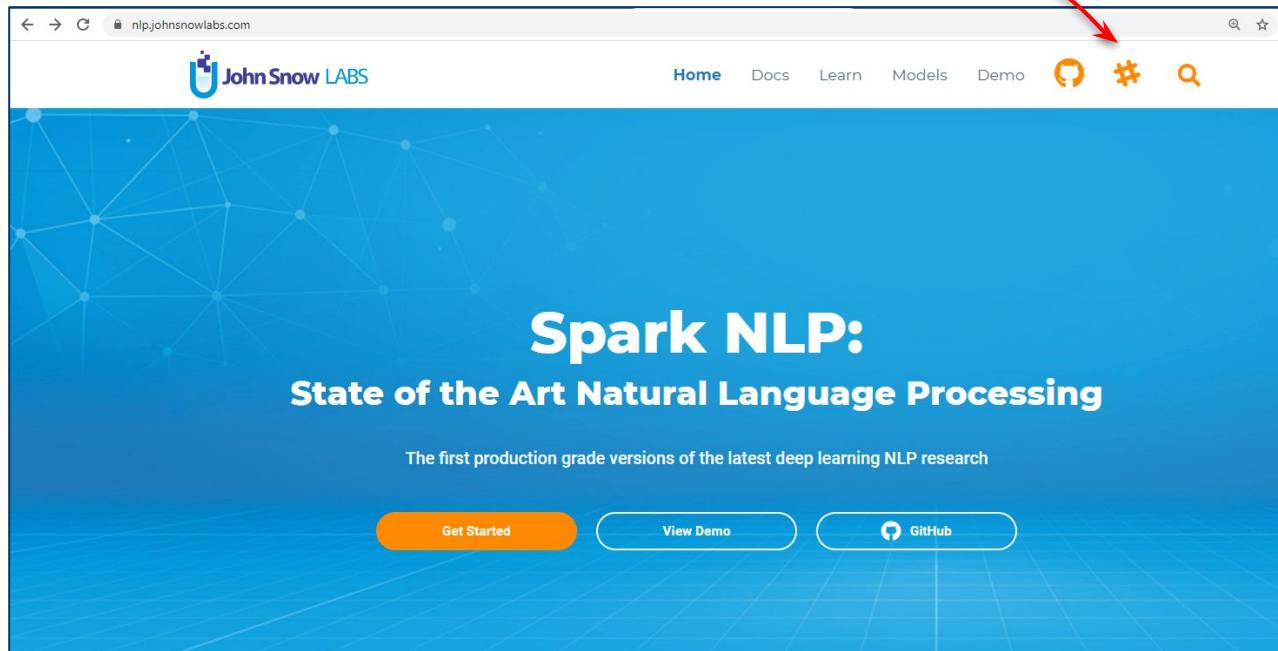
A screenshot of a GitHub repository page for 'JohnSnowLabs / spark-nlp-workshop'. The repository has 48 pull requests, 411 stars, and 234 forks. A red arrow points from the top right towards the breadcrumb navigation bar, which shows the path: 'spark-nlp-workshop / tutorials / Certification\_Trainings / Public /'. Below the path, a list of files and commits is displayed.

File/Commit	Description	Last Modified
..		40 minutes ago
data	updates for nlp summit	3 months ago
databricks_notebooks	readmes	2 months ago
slides	slides moved to a folder	3 months ago
utils	rebased	3 months ago
1.SparkNLP_Basics.ipynb	Pb NB1 updated	2 days ago
10.T5_Workshop_with_Spark_NLP.ipynb	add T5 nb	5 days ago
2.Text_Preprocessing_with_SparkNLP_Annotator...	Pb NB2-updated with RegTok,DocNorm	1 hour ago

# Spark-NLP Slack Channels



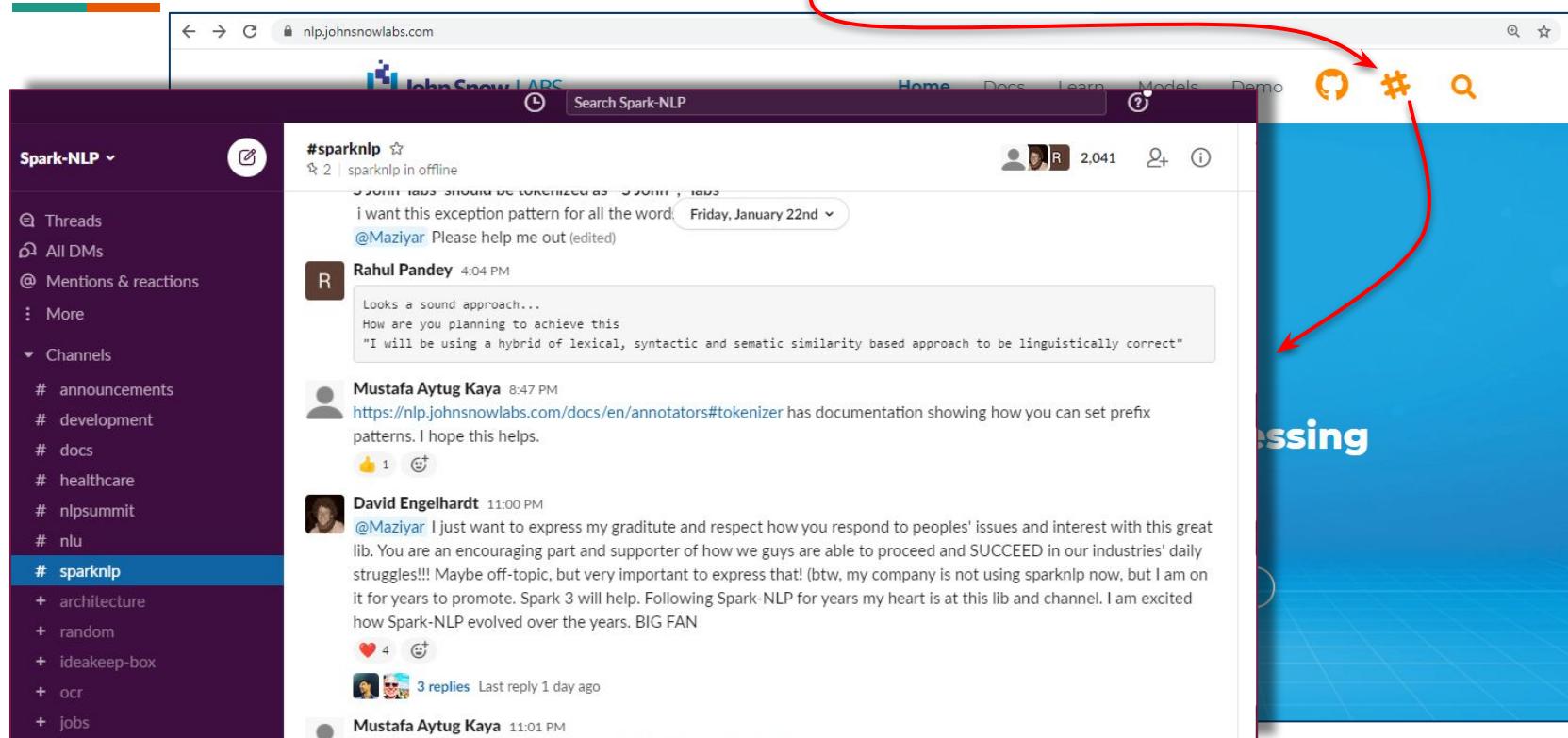
- #sparknlp
- #healthcare
- #nlu
- #ocr
- #jobs



The screenshot shows the homepage of the John Snow LABS NLP platform. The header features the 'John Snow LABS' logo and navigation links for Home, Docs, Learn, Models, and Demo. Below the header is a large blue banner with the text 'Spark NLP: State of the Art Natural Language Processing'. At the bottom of the page are three call-to-action buttons: 'Get Started', 'View Demo', and 'GitHub'.

# Spark-NLP Slack Channels

spark-nlp.slack.com



The screenshot shows a web browser displaying the URL [nlp.johnsnowlabs.com](https://nlp.johnsnowlabs.com). The page has a purple header with the John Snow LABS logo and navigation links for Home, Docs, Learn, Models, and Demo. A red arrow points from the text "spark-nlp.slack.com" at the top right to the Slack interface. The main content area is a Slack channel for "#sparknlp". The sidebar on the left lists various channels, with "# sparknlp" being the active one, indicated by a blue bar at the bottom. The channel feed shows messages from users like Maziyar, Rahul Pandey, Mustafa Aytug Kaya, and David Engelhardt. The interface includes standard Slack features like a search bar, a hashtag icon, and a magnifying glass icon.

#sparknlp ☆  
2 | sparknlp is offline

John Snow LABS should be tokenized as "John", "Snow", "LABS"  
I want this exception pattern for all the words. Friday, January 22nd  
@Maziyar Please help me out (edited)

Rahul Pandey 4:04 PM  
Looks a sound approach...  
How are you planning to achieve this  
"I will be using a hybrid of lexical, syntactic and semantic similarity based approach to be linguistically correct"

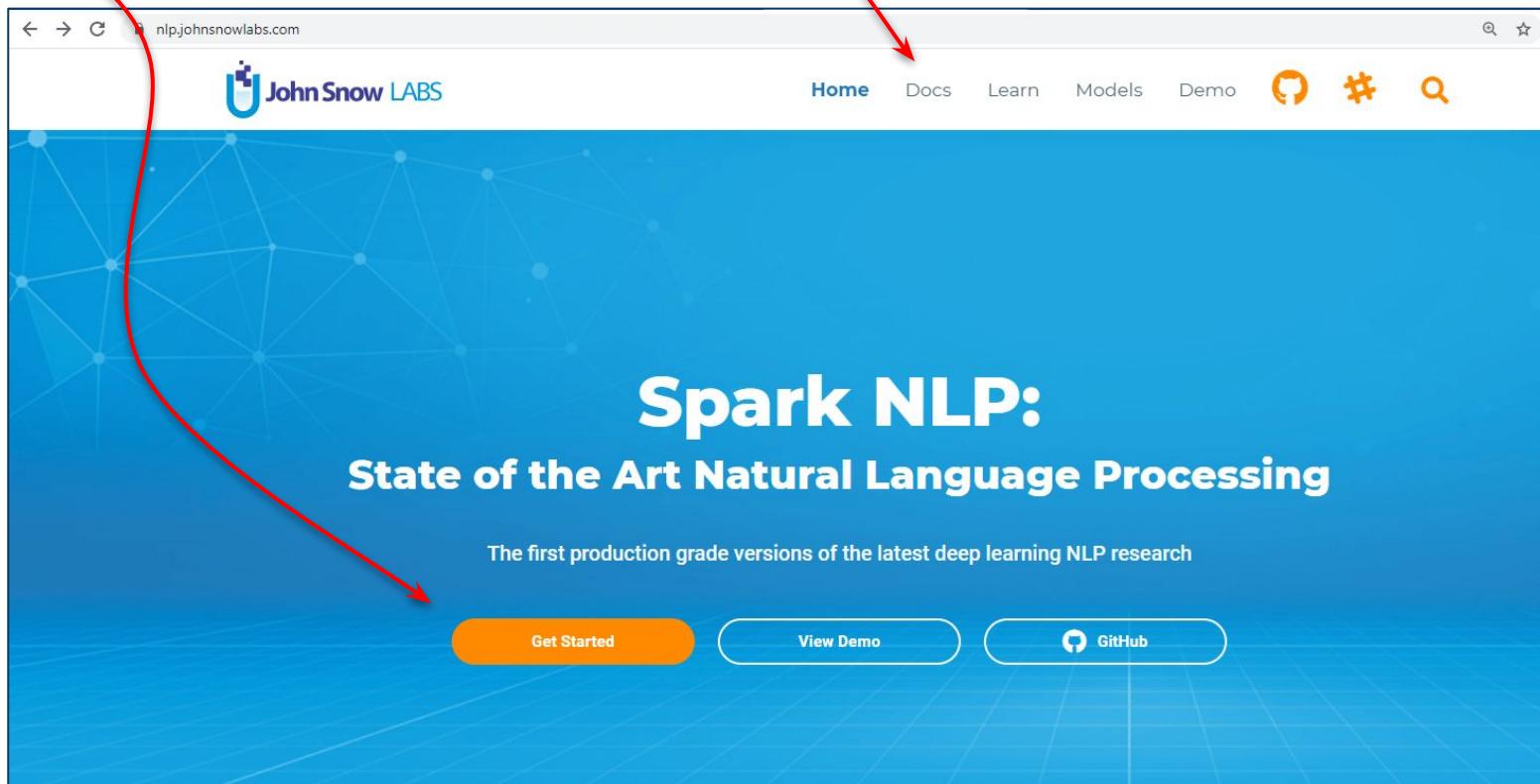
Mustafa Aytug Kaya 8:47 PM  
<https://nlp.johnsnowlabs.com/docs/en/annotators#tokenizer> has documentation showing how you can set prefix patterns. I hope this helps.  
1 thumbs up

David Engelhardt 11:00 PM  
@Maziyar I just want to express my gratitude and respect how you respond to people's issues and interest with this great lib. You are an encouraging part and supporter of how we guys are able to proceed and SUCCEED in our industries' daily struggles!!! Maybe off-topic, but very important to express that! (btw, my company is not using sparknlp now, but I am on it for years to promote. Spark 3 will help. Following Spark-NLP for years my heart is at this lib and channel. I am excited how Spark-NLP evolved over the years. BIG FAN)  
4 hearts, 1 thumbs up

3 replies Last reply 1 day ago

Mustafa Aytug Kaya 11:01 PM

## Docs



The screenshot shows the homepage of the [nlp.johnsnowlabs.com](https://nlp.johnsnowlabs.com) website. The page has a blue header with the John Snow LABS logo and navigation links for Home, Docs, Learn, Models, Demo, and a search bar. Below the header is a large blue banner featuring a network graph background and the text "Spark NLP: State of the Art Natural Language Processing". A red curved arrow points from the "Docs" heading in the top left towards the "Docs" link in the header. Another red arrow points from the "Docs" link in the header down to the "Docs" link in the banner.

nlp.johnsnowlabs.com

John Snow LABS

Home Docs Learn Models Demo

Spark NLP:  
State of the Art Natural Language Processing

The first production grade versions of the latest deep learning NLP research

Get Started View Demo GitHub



# Docs



← → C nlp.johnsnowlabs.com/docs/en/quickstart

John Snow LABS Home Docs Learn Models Demo

Spark NLP

- Getting Started
- Install Spark NLP
- General Concepts
- Transformers
- Annotators
- Helpers
- Pipelines
- Models
- Training
- Scaladoc
- Spark NLP Display
- Developers
- Release Notes

Annotation Lab

- Getting Started
- Start Page
- Project Setup
- Import Documents
- Tasks
- Annotate

## Quick Start

### Requirements & Setup

Spark NLP is built on top of **Apache Spark**

**2.4.4.** For using Spark NLP you need:

- Java 8
- Apache Spark 2.4.x (or Apache Spark 2.3.x)

It is recommended to have basic knowledge of the framework and a working environment before using Spark NLP. Please refer to Spark [documentation](#) to get started with Spark.

Install Spark NLP in

- Python
- Scala and Java
- Databricks

John Snow LABS Home Docs Learn Models Demo

Join our Slack channel

Join our channel, to ask for help and share your feedback. Developers and users can help each other getting started here.

Spark NLP Slack

Spark NLP in Action

Make sure to check out our demos built by Streamlit to showcase Spark NLP in action:

Spark NLP Demo

Spark NLP Workshop

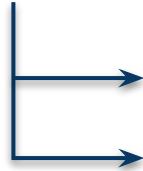
If you prefer learning by example, check this repository:

Spark NLP Workshop

It is full of fresh examples and even a docker container if you want to skip installation.

Below, you can follow into a more theoretical and thorough quick start guide.

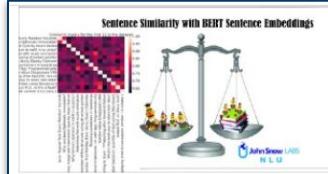
# Articles-Blog Posts



Medium-Spark NLP (<https://medium.com/spark-nlp>)

Blog Posts

# Medium-Spark NLP



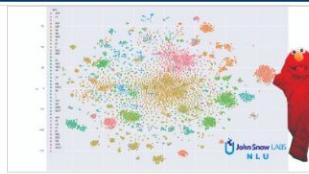
**Easy sentence similarity with BERT Sentence Embeddings using John Snow Labs NLU**

1 Python line to Bert Sentence Embeddings and 5 more for Sentence similarity. Leverage your data to answer questions!



Christian Kasim Loan

Nov 20, 2020 · 8 min read



**1 Python Line for ELMo Word Embeddings with John Snow Labs' NLU**

Including Part of Speech, Named Entity Recognition, Emotion, and Sentiment Classification in the same line! With Bonus t-SNE plots and...



Christian Kasim Loan

Oct 24, 2020 · 7 min read



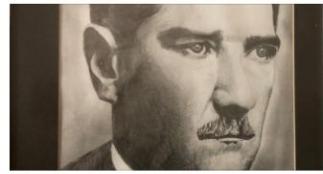
**Turkish NER Model Training using Spark NLP, with the help of NLU.**

The NLU miracle allows us to produce a perfect CoNLL file and a perfect CoNLL file makes the Turkish NER model perfect. This is the first...

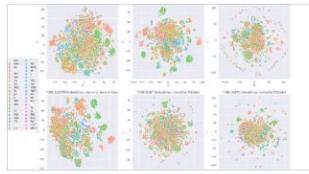


Murat Gunay

Oct 24, 2020 · 10 min read



**Classification of Texts Written in Turkish Language Using Spark NLP**



**1 line of code for BERT, ALBERT, ELMO, ELECTRA, XLNET, GLOVE, Part of Speech with...**



**1 line to BERT Word Embeddings with NLU**

Including Part of Speech, Named Entity

Medium

# Blog Posts

---

an intro article for Spark NLP:

<https://towardsdatascience.com/introduction-to-spark-nlp-foundations-and-basic-components-part-i-c83b7629ed59>

How to start Spark NLP in 2 weeks:

<https://towardsdatascience.com/how-to-get-started-with-sparknlp-in-2-weeks-cb47b2ba994d>

<https://towardsdatascience.com/how-to-wrap-your-head-around-spark-nlp-a6f6a968b7e8>

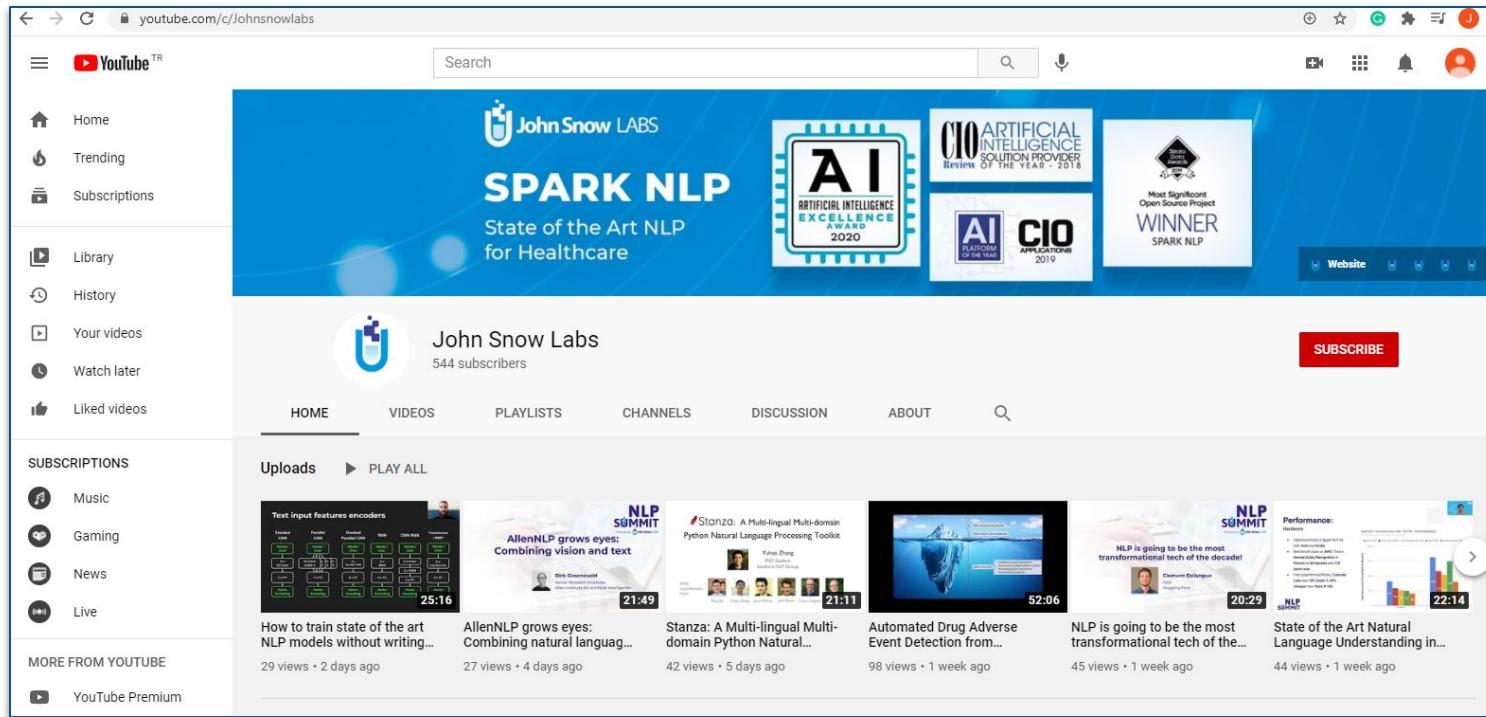
Article for NER and text classification in Spark NLP

<https://towardsdatascience.com/named-entity-recognition-ner-with-bert-in-spark-nlp-874df20d1d77>

<https://medium.com/spark-nlp/named-entity-recognition-for-healthcare-with-sparknlp-nerdl-and-nercrf-a7751b6ad571>

<https://towardsdatascience.com/text-classification-in-spark-nlp-with-bert-and-universal-sentence-encoders-e644d618ca32>

# Youtube Channel



The screenshot shows the YouTube channel page for "John Snow Labs". The channel banner features the text "SPARK NLP" and "State of the Art NLP for Healthcare". It also displays several awards: "CIO ARTIFICIAL INTELLIGENCE SOLUTION PROVIDER OF THE YEAR - 2018", "AI EXCELLENCE AWARD 2020", "CIO PLATFORM OF THE YEAR 2019", and "CIO APPLICATIONS 2019". The channel has 544 subscribers and a "SUBSCRIBE" button.

**HOME**

**VIDEOS**

**PLAYLISTS**

**CHANNELS**

**DISCUSSION**

**ABOUT**

**SUBSCRIPTIONS**

**Uploads** ► **PLAY ALL**

Video Title	Length	Views	Last Updated
Text input features encoders	25:16	29 views	2 days ago
AllenNLP grows eyes: Combining vision and text	21:49	27 views	4 days ago
Stanza: A Multi-lingual Multi-domain Python Natural Language Processing Toolkit	21:11	42 views	5 days ago
Automated Drug Adverse Event Detection from...	52:06	98 views	1 week ago
NLP is going to be the most transformational tech of the...	20:29	45 views	1 week ago
Performance: State of the Art Natural Language Understanding in...	22:14	44 views	1 week ago

**MORE FROM YOUTUBE**

**YouTube Premium**