

## State-of-the-art Natural Language Processing with **Spark NLP 4**



David Talby CTO, John Snow Labs

ORGANIZED BY 😂 databricks

### Agenda

- Why Spark NLP?
- What's New in 4.0?
- 10 Hidden Gems





### John Snow Labs is the team behind Spark NLP





### Spark NLP

State-of-the-art models Enterprise grade Natively Scalable Python, Java, and Scala Trainable & Tunable Releases every 2 weeks

Entity Recognition	Information Extraction They met Last week DATE -> 29-04-2020	Spellin abc Shet ✓ Shet	g & Grammar become the first> became the first	Text Classification
<b>Translation</b> (je t'aime -> i love you)	Summarization () → ()	Questi E	on Answering	Emotion Detection
<ul> <li>Split Text</li> <li>Sentence Detector</li> <li>Tokenizer</li> <li>Normalizer</li> <li>nGram Generator</li> <li>Word Segmentation</li> </ul>	<b>Clean Text</b> <ul> <li>Spell Checker</li> <li>Grammer Checker</li> <li>Writing Style Checker</li> <li>Stopword Cleaner</li> <li>Summarization</li> </ul>	6, Pre-tra Models BERT ALBERT	<b>DOOO+</b> ined Pipelines, & Transformers ELMO GloVe DeBERTa USE	250+ Languages
<ul> <li>Understand Grammar</li> <li>Stemmer</li> <li>Lemmatizer</li> <li>Part of Speech Tagger</li> <li>Dependency Parser</li> <li>Translation</li> </ul>	Find in Text <ul> <li>Text Matcher</li> <li>Regex Matcher</li> <li>Date Matcher</li> <li>Chunker</li> <li>Question Answering</li> </ul>	Longforr T5 N DistilBE XLM S-BER	ner ELECTRA IMT LaBSE IRT RoBERTa 1-RoBERTa T XLNet	
Trainable & Tunable	alable to a Cluster Fast I	nference ghtPipeline	Hardware Optim	ized Community

### Spark NLP

- 99 total releases
- Release every two weeks for the past 4 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
Stemming	Yes	Yes	Yes	Yes	No
Lemmatization	Yes	Yes	Yes	Yes	No
POS tagging	Yes	Yes	Yes	Yes	Yes
Entity recognition	Yes	Yes	Yes	Yes	Yes
Dependency parser	Yes	Yes	Yes	Yes	No
Text matcher	Yes	Yes	No	No	Νο
Date matcher	Yes	No	No	Νο	No
Sentiment detector	Yes	No	Yes	Yes	Yes
Text classification	Yes	Yes	Yes	Νο	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
Question Answering	Yes	Yes	No	No	Yes
Text Style Transfer	Yes	No	No	No	Yes
Translation	Yes	No	No	No	Yes
Summarization	Yes	No	Νο	No	Yes
250+ Languages supported	Yes	No	No	No	Yes
Text Generation (GPT, T5)	Yes	No	No	No	Yes

#### Text Classification with Word & Sentence Embeddings & Transformers in Spark NLP

Glove		ELMO	BE	RT	Electra	Long	ormer		Universal S		Sentence	e Encoder					
Albert	XLN	et Rol	Berta		DistilBert	XImRo	Berta		Bert Se		Bert Sent		Bert Sentence		beddings		
Word2Vec		Bert Toke Classifie	en er	Distil C	Bert Token lassifier	Albert T Classi	oken fier	F	RoBerta Clas	a Tok sifier	ken	RoBerta Clas	a Token sifier	DeBl	ERTa		
Token Classifier		XLNet Token Classifier		Lor C	ngFormer Token lassifier	Bert Seq Classi	Sequence assifier		DistilBert Sequence Classifier		DistilBert Sequence Classifier		DistilBer Sequenc Classifie				
Albert Sequence Classifier		RoBerta Sequence Classifie	e r	Xlm Se Cl	nRoBerta equence assifier	XLNe Sequer Classif	it ice ier	Γ	Longformer Sequence Classifier		Longformer Sequence Classifier		Nor Trai	n Transforme nsformer Em	r Embeddi oedding	ng	
RoBerta Sentence Embeddings			Х	lmRoBerta Se	entence Emb	eddin	gs				Trai	nsformer Sec nsformer Tok	uence Cla en Classifi	issifier er			
a stirring transportin beauty and apparently cutting roc daytime soa  1,999 the movie i filmmaking just experi	Senter , funny ng re ima the bea: reassem methoor ap  is undon methodo imental o	and finally agining of st and 1930s bled from the of any given 0  e by a logy that 's 1 enough		Alre	DistilBERT Pady (pre-)trained	0 1 	0 -0.215 -0.172  0.124	entence E 1 -0.1402 -0.144  0.014	anbeddings           and           and	<b>767</b> 0.201 0.371  0.274	label 1. 0.  1.		Sp C Su Mult Sequ Tok	oark NLF lassifierDL entimentDL tiClassifierD ence Classifi en Classifie	L ier r		

### **The Production-Grade NLP Models Hub**



### **Optimized, Tested, Supported Integrations**



### **NLU: The Simplicity of Python, the Power of Spark NLP**

#### **Powerful One-Liners**

Hundreds of NLP models in tens of languages are at your fingertips with just one line of code

#### **Elegant Python**

Directly read and write pandas dataframes for frictionless integration with other libraries and existing ML pipelines

#### 100% Open Source

Including pre-trained models & pipelines

lu.load("ner"	).predict(	"Angela	Market	is	from	Germany")	
---------------	------------	---------	--------	----	------	-----------	--

	++
entity	tag
	++
Angela Merkel	

Germany LOC

nlu.load('en.classify.fakenews').predict('Unicorns sighted on Mars!')

	++
fake	confidence
	++
FAKE	1.000000

XPATTERNS | BY ATIGEC

### **The Demos Hubs: 300+ Notebooks**

Spark NLP: English FREE

Detect Sentiment & Emotion

Analyze Spelling & Grammar

FREE

Infer Meaning & Intent

**Classify Documents** 

Spark NLP: World

Identify & Translate

European Languages

East Asian Languages

Middle Eastern Languages

Languages of India

African Languages

Languages

Languages

**Recognize Entities** 



**Run popular demos** 

https://nlp.johnsnowlabs.com/demos

### What's new in Spark NLP 4.0?



### **Extractive transformer-based Question answering**

#### What is Spark NLP?

Compute

Context

Spark NLP is a state-of-the-art Natural Language Processing library built on top of Apache Spark. It provides simple, performant & accurate NLP annotations for machine learning pipelines that scale easily in a distributed environment. Spark NLP comes with 5000+ pretrained pipelines and models in more than 200+ languages. It also offers tasks such as Tokenization, Word Segmentation, Part-of-Speech Tagging, Word and Sentence Embeddings, Named Entity Recognition, Dependency Parsing, Spell Checking, Text Classification, Sentiment Analysis, Token Classification, Machine Translation (+180 languages), Summarization & Question Answering, Text Generation, and many more NLP tasks.

Computation time on cpu: cached

a state-of-the-art Natural Language Processing library

### **Extractive transformer-based Question answering**

How many languages Spark NLP supports?

Compute

Context

Spark NLP is a state-of-the-art Natural Language Processing library built on top of Apache Spark. It provides simple, performant & accurate NLP annotations for machine learning pipelines that scale easily in a distributed environment. Spark NLP comes with 5000+ pretrained pipelines and models in more than 200+ languages. It also offers tasks such as Tokenization, Word Segmentation, Part-of-Speech Tagging, Word and Sentence Embeddings, Named Entity Recognition, Dependency Parsing, Spell Checking, Text Classification, Sentiment Analysis, Token Classification, Machine Translation (+180 languages), Summarization & Question Answering, Text Generation, and many more NLP tasks.

Computation time on cpu: 0.093 s

200+

### **Extractive transformer-based Question answering**

Who's your daddy?	Compute

Context

Spark NLP is a state-of-the-art Natural Language Processing library built on top of Apache Spark. It provides simple, performant & accurate NLP annotations for machine learning pipelines that scale easily in a distributed environment. Spark NLP comes with 5000+ pretrained pipelines and models in more than 200+ languages. It also offers tasks such as Tokenization, Word Segmentation, Part-of-Speech Tagging, Word and Sentence Embeddings, Named Entity Recognition, Dependency Parsing, Spell Checking, Text Classification, Sentiment Analysis, Token Classification, Machine Translation (+180 languages), Summarization & Question Answering, Text Generation, and many more NLP tasks.

Computation time on cpu: 0.111 s

Apache Spark

### 1,000+ Models for Modern extractive transformer-based Question answering (QA)

- BERT
- ELECTRA
- DistilBERT
- RoBERTa
- DeBERTa
- XLM-RoBERTa
- Longformer



### **More Accurate: NER**

### Highest F1 for an English Named Entity Recognition Model based on CoNLL03

- Spark NLP v3 achieved F1 95% on CoNLL03 Dev dataset by using BERT Large
- Spark NLP v4 achieves F1 96% on CoNLL03 Dev dataset by using DeBERTa Large

SYSTEM	YEAR	LANGUAGE	CONLL '03
Spark NLP v4	2022	Python/Scala/Java/R	93 (test F1) 96 (dev F1)
Spark NLP v3	2021	Python/Scala/Java/R	93 (test F1) 95 (dev F1)
spaCy v3	2021	Python	91.6
Stanza (StanfordNLP)	2020	Python	92.1
Flair	2018	Python	93.1
CoreNLP	2015	Java	89.6
SYSTEM	YEAR	LANGUAGE	ONTONOTES
Spark NLP v3	2021	Python/Scala/Java/R	90.0 (test F1) 92.5 (dev F1)
spaCy RoBERTa	2020	Python	89.8 (dev F1)
Stanza (StanfordNLP)	2020	Python	88.8 (dev F1)
Flair	2018	Python	89.7

### **More Accurate: Coreference Resolution**

John told Sally that she should come watch him play the violin.

John told Sally that she should come watch him play the violin.

#### **BERT for Coreference Resolution: Baselines and Analysis**

#### Mandar Joshi, Omer Levy, Daniel S. Weld, Luke Zettlemoyer

We apply BERT to coreference resolution, achieving strong improvements on the OntoNotes (+3.9 F1) and GAP (+11.5 F1) benchmarks. A qualitative analysis of model predictions indicates that, compared to ELMo and BERT-base, BERT-large is particularly better at distinguishing between related but distinct entities (e.g., President and CEO). However, there is still room for improvement in modeling document-level context, conversations, and mention paraphrasing. Our code and models are publicly available.

### **Faster! Optimized for the Latest Hardware**













### **Performance Improvements on GPU Devices**

Up to **8X performance** improvements by optimizing batch processing of rows on **GPU** devices Spark NLP 3.4 vs. Spark NLP 4.0 on GPU



### **Performance Improvements on CPU Devices**

Spark NLP 3.4 on CPU vs. Spark NLP 4.0 on CPU with oneDNN

Up to **1.97X performance** improvements on **Intel CPU** by using **oneAPI** Deep Neural Network Library (**oneDNN**)





Model

#### Spark NLP 3.x CPU/GPU vs. Spark NLP 4.0 CPU+oneDNN/GPU



### **Newly Supported Platforms**



### databricks

10.4 +LTS +ML +GPU 10.5 +LTS +ML +GPU



Spark 3.2.x PySpark 3.2.x



TensorFlow 2.7.1



Clusters without HDFS/DBFS/S3 Storage





EMR 6.6.0

Scala 2.12.15

### 10 Hidden Gems



### **1. Edit Sentences for Style**

nlu.load("en.t5.informal\_to\_formal\_styletransfer")
 .predict("Who gives a crap?")
>> Who cares?

#### Spark NLP also comes with a Slang normalizer:



### **2. Build Knowledge Graphs**

Automatically find relationships in free text, using dependency parsing or semantic relation extraction:



Then, use the GraphExtraction annotator to transform them into triplets ready to be loaded into a graph database:



### **3. Train in one language, predict in 100+ languages**

### 100+ Languages supported by Language-agnostic BERT Sentence Embedding (LABSE) and XLM-RoBERTa

#### •••

# Binary Class Classifier, 2 classes nlu.load('xx.embed\_sentence.labse train.sentiment').fit(train\_df).predict(test\_df)

# Multi Class Classifier, N classes nlu.load('xx.embed\_sentence.labse train.classifier').fit(train\_df).predict(test\_df)

# Multi Class Classifier with multiple labels example (i.e. Hashtags)
# N classes, where one row can be assigned up to N labels
nlu.load('xx.embed\_sentence.labse train.multi\_classifier').fit(train\_df).predict(test\_df)

ISO	NAME	ISO	NAME	ISO	NAME
af	AFRIKAANS	ht	HAITIAN_CREOLE	pt	PORTUGUESE
am	AMHARIC	hu	HUNGARIAN	го	ROMANIAN
ar	ARABIC	hy	ARMENIAN	ru	RUSSIAN
as	ASSAMESE	id	INDONESIAN	rw	KINYARWANDA
az	AZERBAIJANI	ig	IGBO	si	SINHALESE
be	BELARUSIAN	is	ICELANDIC	sk	SLOVAK
bg	BULGARIAN	it	ITALIAN	sl	SLOVENIAN
bn	BENGALI	ja	Japanese	sm	SAMOAN
bo	TIBETAN	jv	JAVANESE	sn	SHONA
bs	BOSNIAN	ka	GEORGIAN	so	SOMALI
ca	CATALAN	kk	KAZAKH	sq	ALBANIAN
ceb	CEBUANO	km	KHMER	SF	SERBIAN
co	CORSICAN	kn	KANNADA	st	SESOTHO
CS	CZECH	ko	KOREAN	su	SUNDANESE
cy	WELSH	ku	KURDISH	sv	SWEDISH
da	DANISH	ky	KYRGYZ	SW	SWAHILI
de	GERMAN	la	LATIN	ta	TAMIL
el	GREEK	Ib	LUXEMBOURGISH	te	TELUGU
en	ENGLISH	lo	LAOTHIAN	tg	TAJIK
co	ESPERANTO	lt	LITHUANIAN	th	THAI
cs	SPANISH	lv	LATVIAN	tk	TURKMEN
et	ESTONIAN	mg	MALAGASY	tl	TAGALOG
eu	BASQUE	mi	MAORI	tr	TURKISH
fa	PERSIAN	mk	MACEDONIAN	tt	TATAR
fi	FINNISH	ml	MALAYALAM	ug	UIGHUR
fr	FRENCH	mn	MONGOLIAN	uk	UKRAINIAN
fy	FRISIAN	mr	MARATHI	ur	URDU
ga	IRISH	ms	MALAY	uz	UZBEK
gd	SCOTS_GAELIC	mt	MALTESE	vi	VIETNAMESE
gl	GALICIAN	my	BURMESE	wo	WOLOF
gu	GUJARATI	ne	NEPALI	xh	XHOSA
ha	HAUSA	nl	DUTCH	yi	YIDDISH
haw	HAWAIIAN	no	NORWEGIAN	yo	YORUBA
he	HEBREW	ny	NYANJA	zh	Chinese
hi	HINDI	or	ORIYA	zu	ZULU
hmn	HMONG	pa	PUNJABI		
hr	CROATIAN	pl	POLISH		

### 4. Translate between 200+ languages

>											
Afrikaans	Arabic	Azeri	Bulgarian	Bislama	Bengali	Breton	Catalan	Czech	Welsh	Danish	German
af	ar	az	bg	bi	bn	br	ca	cs	су	da	de
			*	·6)							
Ewe	Greek	English	Esperanto	Spanish	Estonian	Basque	Farsi	Finnish	Fiji	French	Irish
e	el	en	eo	es	et	eu	fa	fi	fj	fr	ga
	*	•	x¢x	۲	24					-	
alician	Manx	Hausa	Hebrew	Hindi	Hiri	Haitian	Hungarian	Armenian	Indonesia	nIgbo	Icelandio
1	gv	ha	he	hi	Motu	ht	hu	hy	id	ig	is
		+ +			ho				*		
		+ +					380 J	<u> </u>			
tallan +	Japanese	Georgian	Kongo	Kuanyama	Greenland	isorean	Latin	Ganda	Lingala	Luba-	Latvian
	Ja	Ka	Ng		KI		Ia	Ig		lu	10
	•	$\geq \leq$		*	*		影			*	•
alagasy	Marshalle	SEYRO	Malayalam	Marathi	Maltese	Ndonga	Dutch	Norwegian	Chichewa	Oromoor	Punjabi
g	mh	Macedonia	nml	mr	mt	ng	nl	Bokmal	ny	om	ра
		mk					$\square$	no			
	۲			-	-						*
olish	Portugues	eKirundi	Romanian	Russian	Kinyarwan	dSangro	Slovak	Slovenian	Samoan	Shona	Somali
1	pt	IN	IO	IU	rw	sg	sk	sl	sm	sn	so
1998		*			<b>&gt;</b> *			+	C*		
lbanian	Siswati	Sesotho	Swedish	Thai	Tigrinya	Tagalog	Tswana	Tonga	Turkish	Tsonga	Twi
p	SS	st	sv	th	ti	tl	tn	to	tr	ts	tw
		C		*	45		88		ė	04	
ahitian	Ukrainian	Urdu	Venda	Vietnames	eWalloon	Xhosa	Yoruba	Chinese	Zulu	94 m	iore!
v	uk	UT	Ve	vi	wa	xh	VO	zh	711		

#### •••

# Use ISO standards for the languages nlu.load('<start\_language>.translate\_to.<target\_language>')

#Translate Turkish to English: nlu.load('tr.translate\_to.en')

#Translate English to French: nlu.load('en.translate\_to.fr')

#Translate French to Hebrew
nlu.load('fr.translate\_to.he')`

#Translate English to German
nlu.load('en.translate\_to.de')`

# **5. Import and Scale models from HF and TFHub**



#### **TF Hub to Spark NLP**

Spark NLP	TF Hub Notebooks	Colab
BertEmbeddings	TF Hub in Spark NLP - BERT	CO Open in Colab
BertSentenceEmbeddings	TF Hub in Spark NLP - BERT Sentence	CO Open in Colab
AlbertEmbeddings	TF Hub in Spark NLP - ALBERT	CO Open in Colab

https://nlp.johnsnowlabs.com/docs/en/transformers#example-notebooks

#### **HuggingFace to Spark NLP**

Spark NLP	HuggingFace Notebooks	Colab
BertEmbeddings	HuggingFace in Spark NLP - BERT	Open In Colab
BertSentenceEmbeddings	HuggingFace in Spark NLP - BERT Sentence	CO Open in Colab
DistilBertEmbeddings	HuggingFace in Spark NLP - DistilBERT	Open in Colab
RoBertaEmbeddings	HuggingFace in Spark NLP - RoBERTa	Open in Colab
XImRoBertaEmbeddings	HuggingFace in Spark NLP - XLM-RoBERTa	Open in Colab
AlbertEmbeddings	HuggingFace in Spark NLP - ALBERT	CO Open in Colab
XInetEmbeddings	HuggingFace in Spark NLP - XLNet	Open in Colab
LongformerEmbeddings	HuggingFace in Spark NLP - Longformer	CO Open in Colab
BertForTokenClassification	HuggingFace in Spark NLP - BertForTokenClassification	Open in Colab
DistilBertForTokenClassification	HuggingFace in Spark NLP - DistilBertForTokenClassification	CO Open in Colab
AlbertForTokenClassification	HuggingFace in Spark NLP - AlbertForTokenClassification	CO Open In Colab
RoBertaForTokenClassification	HuggingFace in Spark NLP - RoBertaForTokenClassification	CO Open in Colab
XlmRoBertaForTokenClassification	HuggingFace in Spark NLP - XImRoBertaForTokenClassification	CO Open in Colab
BertForSequenceClassification	HuggingFace in Spark NLP - BertForSequenceClassification	CC Open in Colab
DistilBertForSequenceClassification	HuggingFace in Spark NLP - DistilBertForSequenceClassification	CO Open n Colab
AlbertForSequenceClassification	HuggingFace in Spark NLP - AlbertForSequenceClassification	Open in Colab
RoBertaForSequenceClassification	HuggingFace in Spark NLP - RoBertaForSequenceClassification	Open in Colab
XImRoBertaForSequenceClassification	HuggingFace in Spark NLP - XImRoBertaForSequenceClassification	Open in Colab
XInetForSequenceClassification	HuggingFace in Spark NLP - XInetForSequenceClassification	00 Open in Colab

### 6. Unsupervised Keyword Extraction

<u>YAKE</u> Is Yet Another Keyword Extraction Algorithm that can extract keywords without any by leveraging statistical properties of ngrams

nlu.load('yake').predict("NLU is a Python Library for beginners and experts in NLP")

keywords_score_confidence	keywords	sentence
0.454232	[nlu, nlp, python library]	NLU is a Python Library for beginners and expe

### 7. Check Spelling Based on Context

```
# 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']

Word error rate of 8.09% versus 20.24% with the JamSpell library

### 8. Deep-Learning-Based Sentence Splitter

How do you correctly split text to sentences when there is no punctuation, bad punctuation, or unexpected line breaks? (common in OCR and ASR outputs)

```
text = """John loves Mary.Mary loves Peter
Peter loves Helen .Helen loves John;
Total: four people involved."""
for anno in sd model.fullAnnotate(text)[0]["sentences"]:
    print("{}\t{}\t{}\.format(
        anno.metadata["sentence"], anno.begin, anno.end, anno.result))
                       John loves Mary.
               15
0
       0
1
               32
                       Mary loves Peter
       16
                       Peter loves Helen .
2
       33
               51
3
       52
                       Helen loves John;
               68
4
       71
               98
                       Total: four people involved.
```

### 9. Identify 10+ Kinds of Toxic Content



nlu.load('en.classify.toxic').predict('You are to stupid')

### **10. Measure Semantic Similarity**

#### First sentence to compare

Sign up for our mailing list to get free offers and updates about our products!

#### Second sentence to compare

Subscribe to notifications to receive information about discounts and new offerings.

Detected similarity: 66%

These sentences are similar.

Useful for searching in knowledge bases, filtering content, or for unsupervised topic modelling.

### DATA+AI SUMMIT 2022

# Thank you!

![](_page_33_Picture_2.jpeg)

in In/davidtalby

@davidtalby

nlp.johnsnowlabs.com

github.com/johnsnowlabs

![](_page_33_Picture_7.jpeg)

SUMMIT 2022