

# USE CASE Healthcare-Specific Large Language Models in Action



**LUCA MARTIAL**  
Senior Data Scientist  
**John Snow Labs**



# *Healthcare-Specific Large Language Models in Action*



**LUCA MARTIAL**  
Senior Data Scientist  
**John Snow Labs**



JohnSnowLabs / **spark-nlp** Public

State of the Art Natural Language Processing

[nlp.johnsnowlabs.com](https://nlp.johnsnowlabs.com)

Apache-2.0 license

3k stars 608 forks

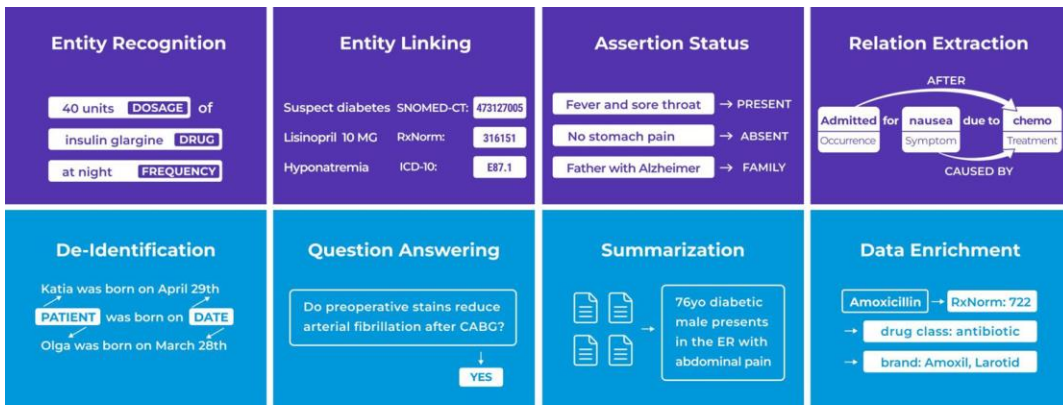
<b>Entity Recognition</b> I love <b>Lucy</b> <b>PERSON</b>	<b>Text Classification</b> 	<b>Spelling &amp; Grammar</b> abc She become the first... ✓ She became the first	<b>Information Extraction</b> They met <b>Last week</b> <b>DATE</b> -> 29-04-2020
<b>Question Answering</b> 	<b>Speech to Text</b> 	<b>Image Classification</b> 	<b>Reading Comprehension</b> 
<b>Translation</b> [je t'aime -> i love you]	<b>Summarization</b> 	<b>Paraphrasing</b> You bet! > For sure.	<b>Emotion Detection</b> 

<b>Split Text</b> <ul style="list-style-type: none"> <li>Sentence Detector</li> <li>Tokenizer</li> <li>Normalizer</li> <li>nGram Generator</li> <li>Word Segmentation</li> </ul>	<b>Clean Text</b> <ul style="list-style-type: none"> <li>Spell Checker</li> <li>Grammar Checker</li> <li>Writing Style Checker</li> <li>Stopword Cleaner</li> <li>Summarization</li> </ul>	<b>20,000+</b> <b>Pre-trained Pipelines, Models &amp; Transformers</b> <table border="1"> <tr> <td>BERT</td> <td>ELMO</td> <td>TAPAS</td> </tr> <tr> <td>ALBERT</td> <td>DeBERTa</td> <td>USE</td> </tr> <tr> <td>Longformer</td> <td>ELECTRA</td> <td></td> </tr> <tr> <td>T5</td> <td>NMT</td> <td>ViT</td> </tr> <tr> <td>DistilBERT</td> <td>RoBERTa</td> <td></td> </tr> <tr> <td colspan="3">XLM-RoBERTa</td> </tr> <tr> <td>Wav2Vec2</td> <td>XLNet</td> <td></td> </tr> </table>	BERT	ELMO	TAPAS	ALBERT	DeBERTa	USE	Longformer	ELECTRA		T5	NMT	ViT	DistilBERT	RoBERTa		XLM-RoBERTa			Wav2Vec2	XLNet		<b>250+</b> <b>Languages</b>
BERT	ELMO	TAPAS																						
ALBERT	DeBERTa	USE																						
Longformer	ELECTRA																							
T5	NMT	ViT																						
DistilBERT	RoBERTa																							
XLM-RoBERTa																								
Wav2Vec2	XLNet																							
<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>Trainable &amp; Tunable</b> 	<b>Scalable</b> 	<b>Fast Inference</b> 	<b>Hardware Optimized</b> 	<b>Community</b> 
------------------------------------	---------------------	---------------------------	-------------------------------	----------------------



### Algorithms

### Content

#### Information Extraction

- Document Classification
- Entity Disambiguation
- Contextual Parsing
- Patient Risk Scoring

#### Data Obfuscation

- Name Consistency
- Gender Consistency
- Age Group Consistency
- Format Consistency

#### Medical Language Models



#### Medical Terminologies



#### Clinical Grammar

- Deep Sentence Detector
- Medical Spell Checking
- Medical Part of Speech
- Terminology Mapping

#### Zero-Shot Learning

- Entities by Prompt
- Relations by Prompt
- Classification by Prompt
- Relative Data Extraction

### 1,000+ Pretrained Models

#### Clinical Text

Signs, Symptoms, Treatments, Findings, Procedures, Drugs, Tests, Labs, Vitals, Sections, Adverse Effects, Risk Factors, Anatomy, Social Determinants, Vaccines, Demographics, Sensitive Data

#### Biomedical Text

Clinical Trial Design, Protocols, Objectives, Results; Research Summary & Outcomes; Organs, Cell Lines, Organisms, Tissues, Genes, Variants, Expressions, Chemicals, Phenotypes, Proteins, Pathogens

### Trainable & Tunable

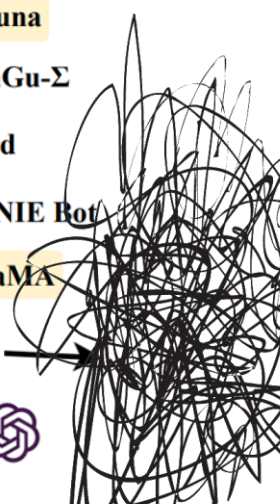
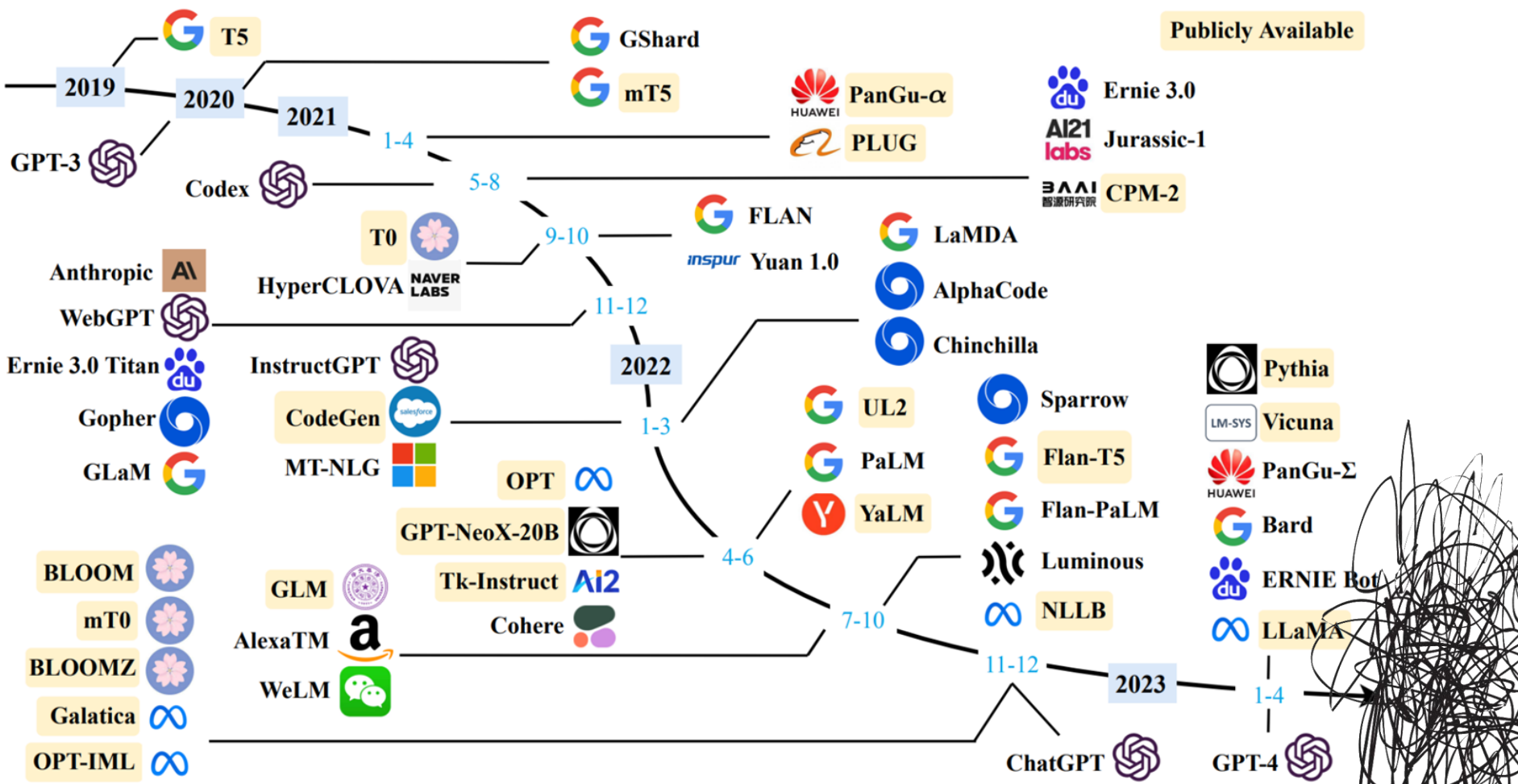
### Scalable

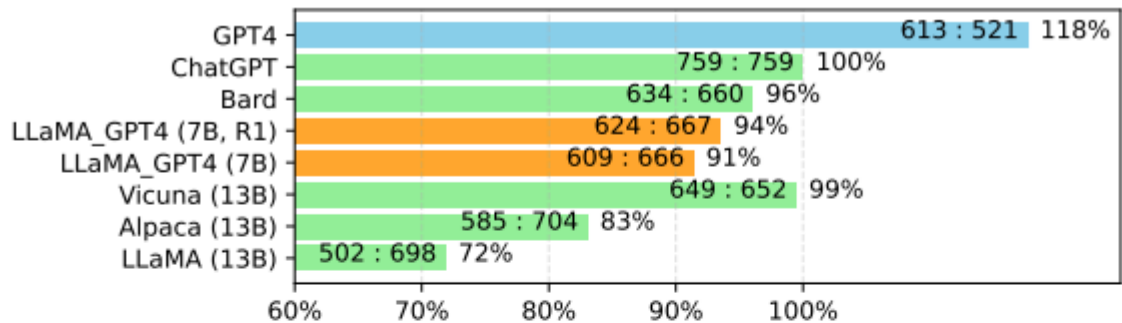
### Fast Inference

### Hardware Optimized

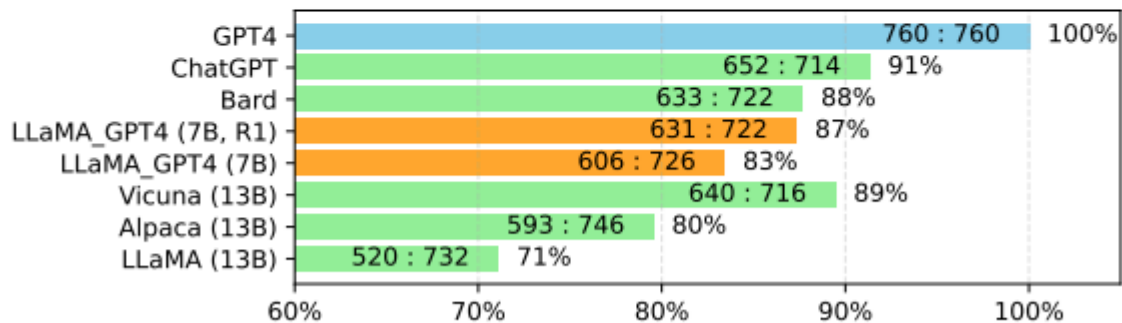
### Community







(c) All chatbots against ChatGPT



(d) All chatbots against GPT-4

- GPT4all
- OpenChatKit
- Alpaca
- Lit-LLaMA
- Dolly
- **Vicuna**
- medAlpaca
- ColossalChat
- Cerebras-GPT
- Koala
- Baize





What about the  
medical domain?





# Capabilities of GPT-4 on Medical Challenge Problems

Dataset	Component	<b>GPT-4 (5 shot)</b>	GPT-4 (zero shot)	GPT-3.5 (5 shot)	GPT-3.5 (zero shot)	Flan-PaLM 540B* (few shot)
MedQA	Mainland China	<b>75.31</b>	71.07	44.89	40.31	–
	Taiwan	<b>84.57</b>	82.17	53.72	50.60	–
	United States (5-option)	<b>78.63</b>	74.71	47.05	44.62	–
	United States (4-option)	<b>81.38</b>	78.87	53.57	50.82	60.3**
PubMedQA	Reasoning Required	74.40	75.20	60.20	71.60	<b>79.0</b>
MedMCQA	Dev	<b>72.36</b>	69.52	51.02	50.08	56.5
MMLU	Clinical Knowledge	<b>86.42</b>	86.04	68.68	69.81	77.00
	Medical Genetics	<b>92.00</b>	91.00	68.00	70.00	70.00
	Anatomy	<b>80.00</b>	<b>80.00</b>	60.74	56.30	65.20
	Professional Medicine	<b>93.75</b>	93.01	69.85	70.22	83.80
	College Biology	93.75	<b>95.14</b>	72.92	72.22	87.50
	College Medicine	76.30	<b>76.88</b>	63.58	61.27	69.90

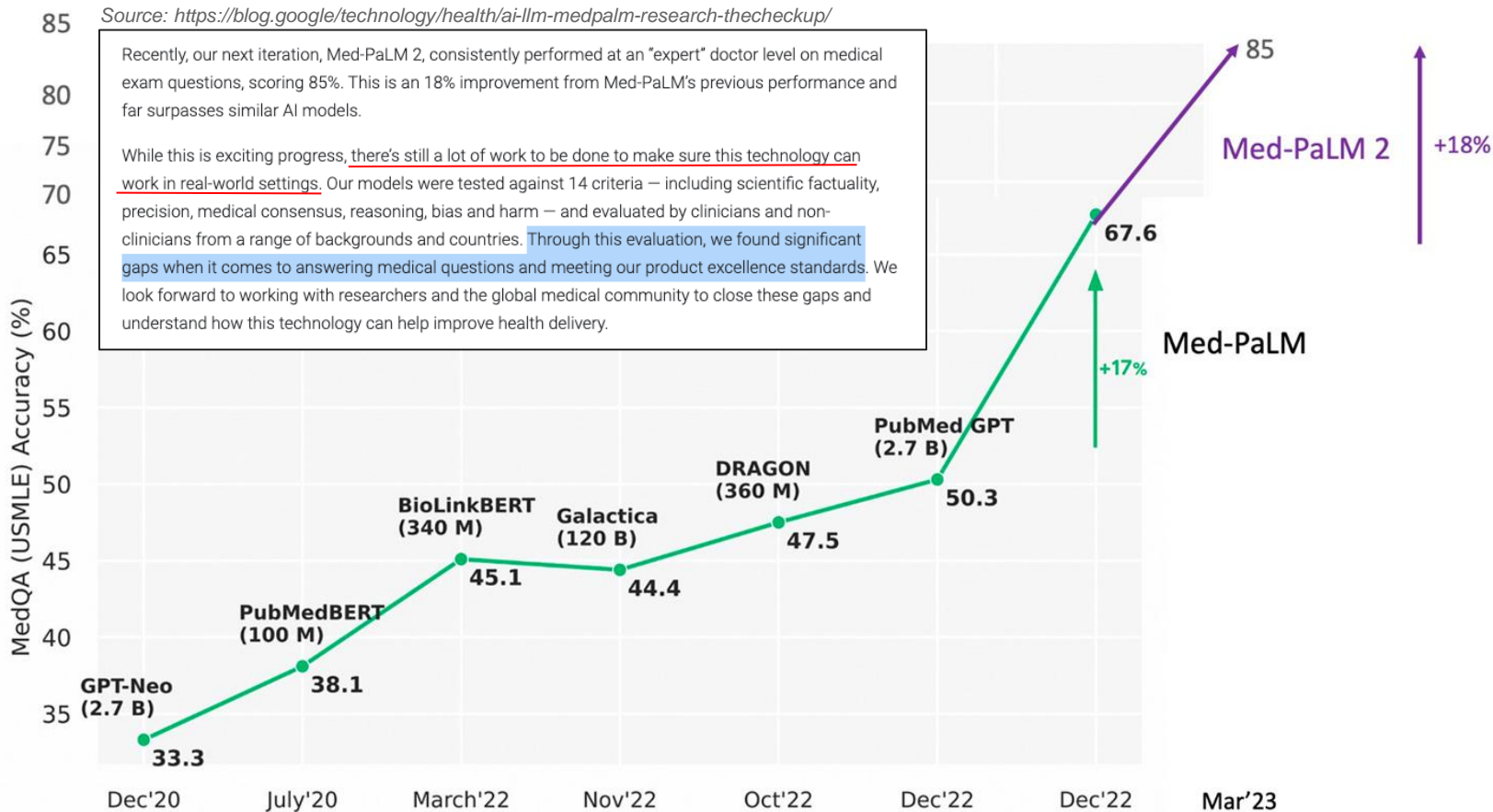
# Google's MedPaLM-2 on USMLE (Medical License Exam)

Source: <https://blog.google/technology/health/ai-llm-medpalm-research-thecheckup/>

Recently, our next iteration, Med-PaLM 2, consistently performed at an "expert" doctor level on medical exam questions, scoring 85%. This is an 18% improvement from Med-PaLM's previous performance and far surpasses similar AI models.

While this is exciting progress, there's still a lot of work to be done to make sure this technology can work in real-world settings. Our models were tested against 14 criteria — including scientific factuality, precision, medical consensus, reasoning, bias and harm — and evaluated by clinicians and non-clinicians from a range of backgrounds and countries. Through this evaluation, we found significant gaps when it comes to answering medical questions and meeting our product excellence standards. We look forward to working with researchers and the global medical community to close these gaps and understand how this technology can help improve health delivery.

Medical Question Answering





Choose the Task :

Summarizer ▾

- Summarizer
- Question Answering
- Text Generation

Try it yourself:

 [Open in Colab](#)

# Explore Medical Large Language Models

## Clinical Text Summarization

This model is specifically trained on clinical data for text summarization.

 [Select an example](#)

Medical Specialty: Allergy / Immunology, Sample Name: Allerg... ▾

Text

Medical Specialty: Allergy / Immunology, Sample Name: Allergic Rhinitis  
Description: A 23-year-old white female presents with complaint of allergies. (Medical Transcription Sample Report)

## Clinical Entity Recognition

40 units **DOSAGE** of

insulin glargine **DRUG**

at night **FREQUENCY**

## Clinical Entity Linking

Suspect diabetes SNOMED-CT: **473127005**

Lisinopril 10 MG RxNorm: **316151**

Pyonatremia ICD-10: **E87.1**

## Assertion Status

Fever and sore throat → PRESENT

No stomach pain → ABSENT

Father with Alzheimer → FAMILY

## De-Identification

Ora **NAME**, a 25 **AGE** yo

cashier **PROFESSION** from

Morocco **LOCATION**

## Relation Extraction



## *Name Entity Recognition*

*Please identify Person, Organization, Location and Miscellaneous Entity from the given text.*

**Text:** *All four teams are level with one point each from one game.*

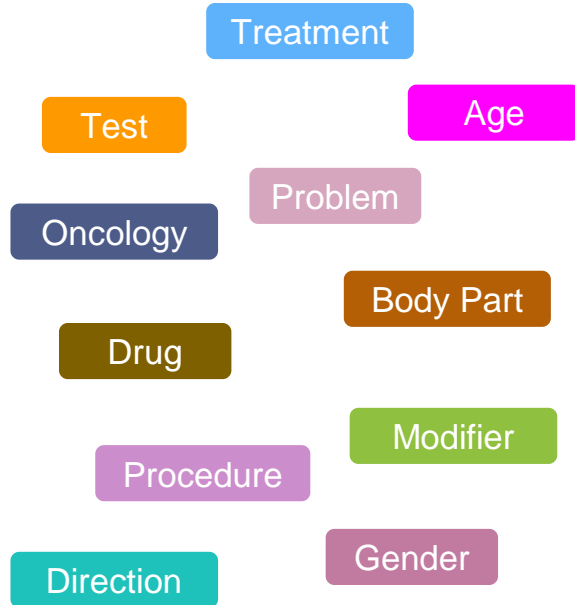
**Entity:**

## NER

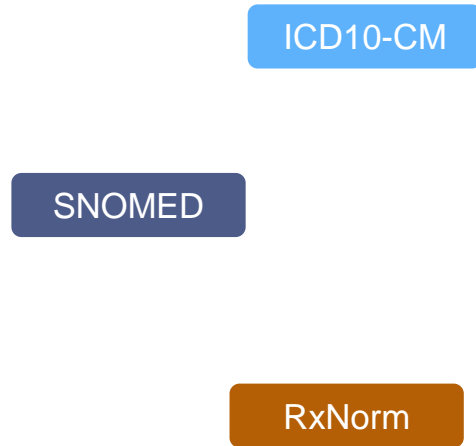
<b>Model</b>	<i>Zero-Shot</i>		<i>Fine-Tuned</i>		
	ChatGPT	GPT-3.5	Flair	LUKE	ACE
<b>All</b>	<b>53.7</b>	53.5	93.0	93.9	<b>94.6</b>
<b>Loc</b>	<b>72.2</b>	67.1	94.0	-	-
<b>Per</b>	<b>81.4</b>	78.0	97.4	-	-
<b>Org</b>	45.1	<b>50.0</b>	91.9	-	-
<b>Misc</b>	4.5	<b>4.8</b>	83.0	-	-

# Scope of Experiments

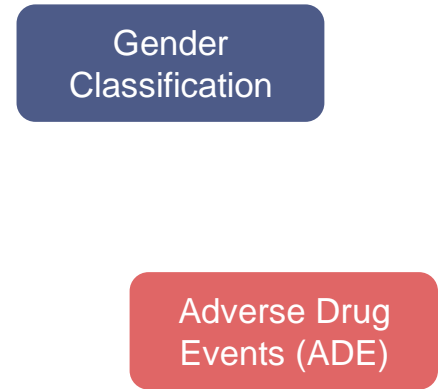
## Named Entity Recognition




## Entity Resolution (Medical Terminologies)



## Text Classification



# In-Depth Comparison of Spark NLP for Healthcare and ChatGPT on Clinical Named Entity Recognition


 **Veysel Kocaman** · Follow  
Published in John Snow Labs · 5 min read · Apr 18

55

🔖 🔄 📄 ⋮

Spark NLP for Healthcare NER models outperform ChatGPT by 10–45% on key medical concepts, resulting in half the errors compared to ChatGPT.

# Comparing Spark NLP for Healthcare and ChatGPT in Extracting ICD10-CM Codes from Clinical Notes


 **Veysel Kocaman** · Follow  
Published in John Snow Labs · 6 min read · Apr 16

69

🔖 🔄 📄 ⋮

In assigning ICD10-CM codes, Spark NLP for Healthcare achieved a 76% success rate, while GPT-3.5 and GPT-4 had overall accuracies of 26% and 36% respectively.

# A Comprehensive Comparison of ChatGPT and Spark NLP for Healthcare in De-Identification of Sensitive Data (PHI)

 **Veysel Kocaman** · Follow  
Published in John Snow Labs · 5 min read · Apr 18

109

🔖 🔄 📄 ⋮

Spark NLP for Healthcare De-Identification module demonstrates superior performance with a 93% accuracy rate compared to ChatGPT's 60% accuracy on detecting PHI entities in clinical notes.

[spark-nlp-workshop](#) / [tutorials](#) / [academic](#) / [LLMs\\_in\\_Healthcare](#) / [benchmarks](#) / 

Add file ▾ ⋮

 **aydinmyilmaz** add missing deid prompt b38305b · last month 🕒 History

Name	Last commit message	Last commit date
..		
config	add benchmarks	last month
data	add deid sentences	last month
workbench	add missing deid prompt	last month
README.md	add benchmarks	last month
requirements.txt	add benchmarks	last month



# Extracting Medical Problems

100 sentence, ~800 entities

## Prompt

You are a highly experienced, skilled and helpful medical annotator who have been working on medical texts to label medical entities.

I will provide you some entity types with sample chunks and I want you to find similar entities from given texts.

- Entity Type: Problem
- 1. Example chunks for Problem Type: feels weak, shortness of breath, backache
- 2. Example chunks for Problem Type: gastroparesis, gastritis, allergies, pneumonitis
- 3. Example chunks for Problem Type: spine fractures, ligature strangulation, abrasions
- 4. Example chunks for Problem Type: depression, bipolar disorder, psychosis
- 5. Example chunks for Problem Type: colon cancer, mesothelioma , brachial plexus tumor
- 6. Example chunks for Problem Type: depression, anxiety, bipolar disorder, psychosis
- 7. Example chunks for Problem Type: coronary artery disease, CAD, cardiomyopathy
- 8. Example chunks for Problem Type: renal disease, nephrolithiasis, hydronephrosis
- 9. Example chunks for Problem Type: overweight
- 10. Example chunks for Problem Type: DM Type II, diabetic
- 11. Example chunks for Problem Type: obese
- 12. Example chunks for Problem Type: wandering atrial pacemaker, multifocal atrial tachycardia, frequent APCs, bradycardia
- 13. Example chunks for Problem Type: tuberculosis, sexually transmitted diseases, HIV
- 14. Example chunks for Problem Type: increased attenuation, T1 hypointensity, opacity in apex right lung
- 15. Example chunks for Problem Type: stroke, TIA
- 16. Example chunks for Problem Type: increased cholesterol, hypercholesterolemia
- 17. Example chunks for Problem Type: tachycardic, afebrile
- 18. Example chunks for Problem Type: high blood pressure, HTN

I want you to extract Problem type of entities from the given text and label them as Problem

Task :

Find entities in the given sentence.

Answer value must be as given (valid JSON) for the given sentence as example:

```
{{"given_sentence": "Patient feels weak.", "list_of_entities": [{"entity_type": "Problem", "chunk": "feels weak"}]}}
```

Now I want you to find the Problem entities in the given sentence:

76%

GPT 3.5

The patient denies chest pain , irregular heartbeats , sudden changes in heartbeat or palpitation , shortness of breath , difficulty breathing at night , swollen legs or feet , heart murmurs , high blood pressure , cramps in his legs with walking , pain in his feet or toes at night or varicose veins .

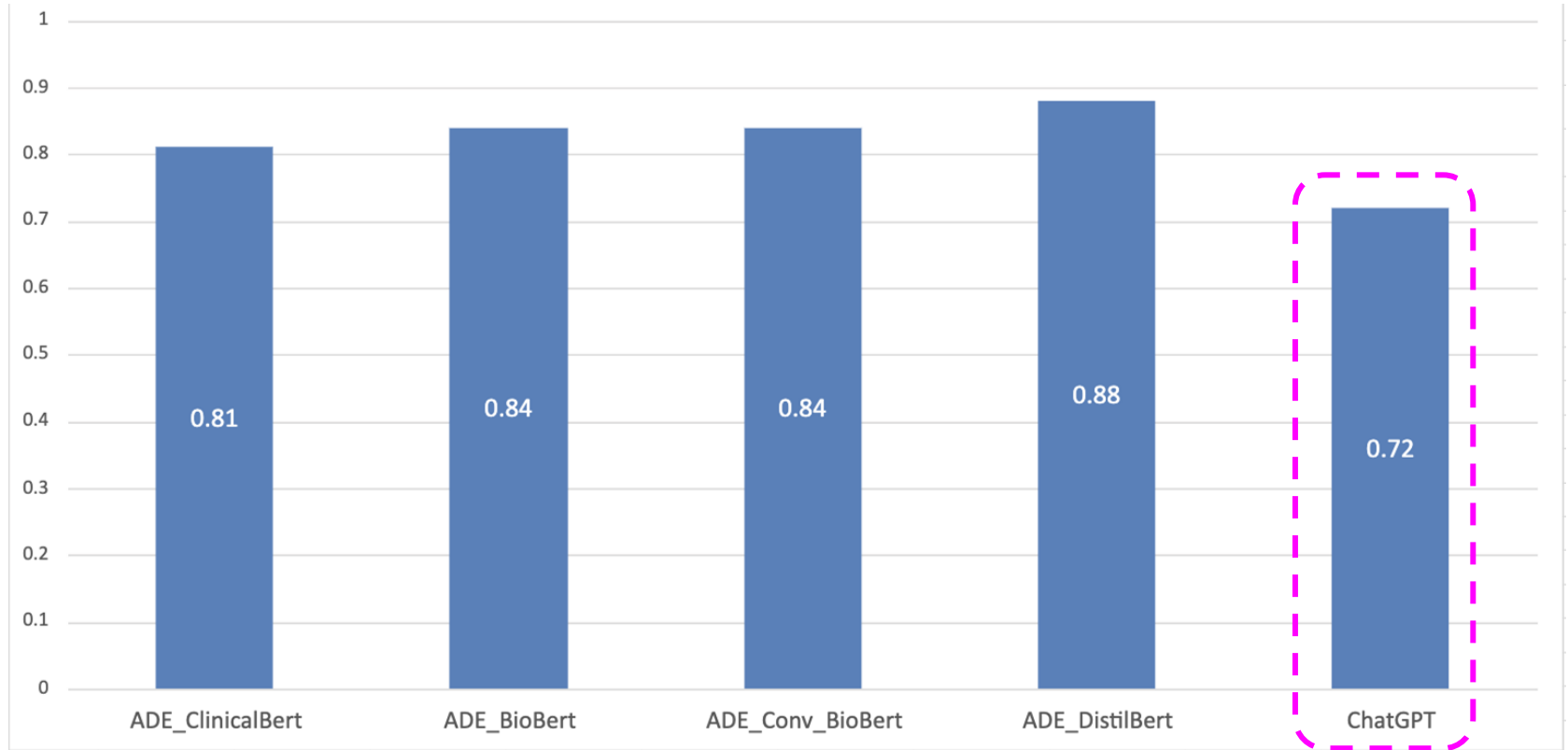
89%

Spark NLP  
(ner\_jsl\_reduced)

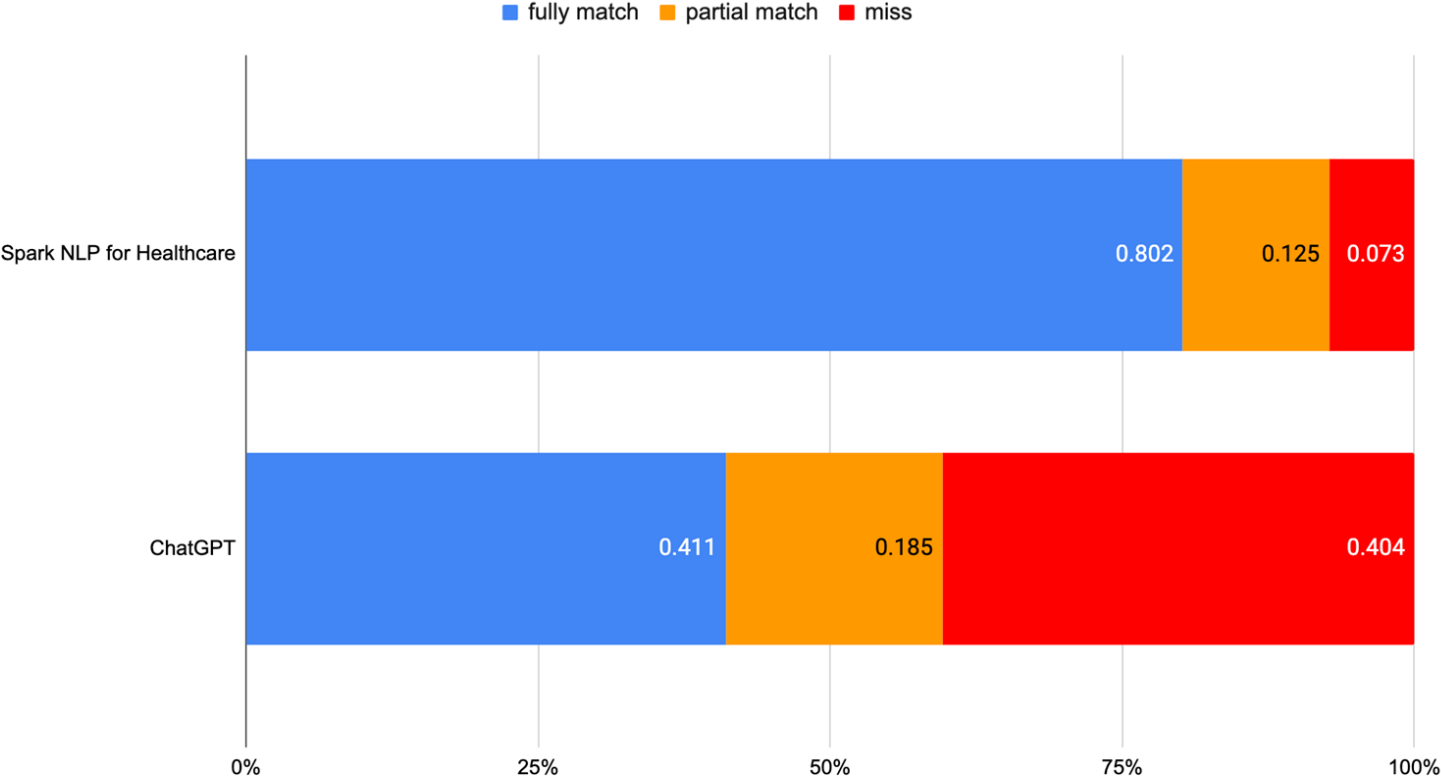
The patient denies chest pain , irregular heartbeats , sudden changes in heartbeat or palpitation , shortness of breath , difficulty breathing at night , swollen legs or feet , heart murmurs , high blood pressure , cramps in his legs with walking , pain in his feet or toes at night or varicose veins .

\* lenient metrics (partially overlapping chunks counted as hit)

# Detecting Adverse Drug Events

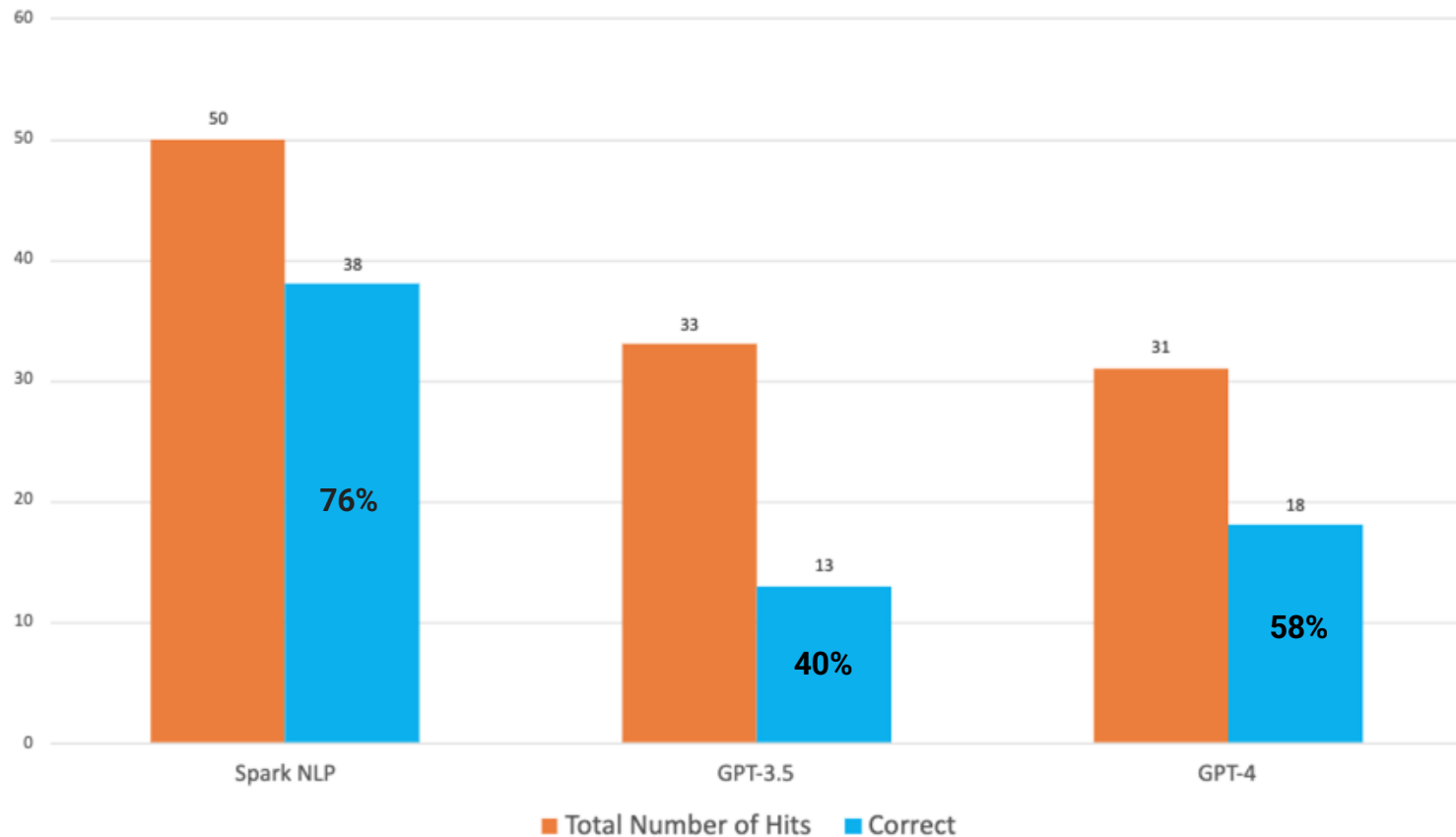


# De-Identifying PHI Data



A Comprehensive Comparison of ChatGPT and Spark NLP for Healthcare in De-Identification of Sensitive Data (Medium)

# Extracting ICD10-CM Codes



# LLMs: The Good, The Bad




- ✓ Tasks making HCP lives easier (note summarizing, patient profiling, etc)
- ✓ Indexing, querying databases
- ✓ Agent-based integrations




- ✗ Unreliable data abstractors
- ✗ On-prem deployment is not possible
- ✗ Hallucinating and fabricating incorrect results w/ high confidence
- ✗ Domain & task-specific fine tuning will be required (expertise, time, money)

# Solutions Moving Forward

 Granular data abstraction by high-precision models; repetitive & laborious tasks by LLMs

 LLMs guard-railed by explainable DL/ML models, knowledge graphs, rule-based systems

 LLMs as smart assistants (convert natural language to structured queries (SQL, Cypher))





# INTELLIGENT HEALTH UK 2023

Breaking down the barriers  
between tech and healthcare