

How to build a foundation of AI-based healthcare systems through language models?

Mehdi Bahrami, Ph.D., IEEE Senior Member
Artificial Intelligence Laboratory
Fujitsu Research of America,
Sunnyvale, CA



[Code & Slide](#)

Overview

- Why Language Model?
- What is a language model?
 - Word2Vec
 - GloVe
 - Transformer
 - BERT

Why language models?

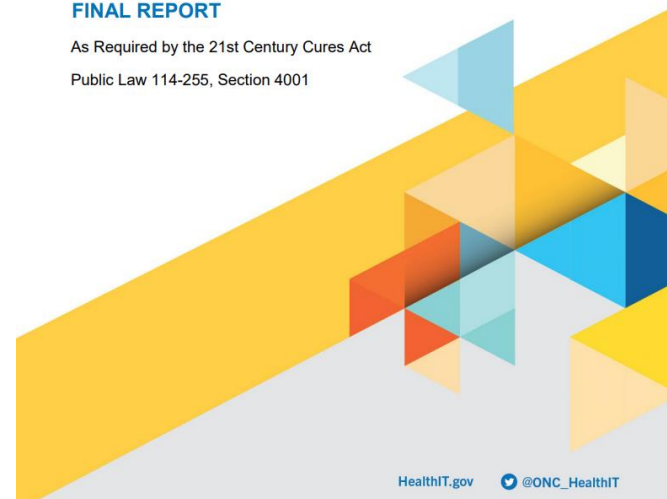
- Health Record System (EHR)
 - The Office of the National Coordinator for Health Information Technology
- Electronic Medical Record (EMR)
 - Electronic record of health-related information on an individual within one health care organization
- HHS Strategy on Reducing Regulatory and Administrative Burden Relating to the Use of Health IT and EHRs was released on February 21, 2020.



Strategy on Reducing Regulatory and Administrative Burden Relating to the Use of Health IT and EHRs

FINAL REPORT

As Required by the 21st Century Cures Act
Public Law 114-255, Section 4001

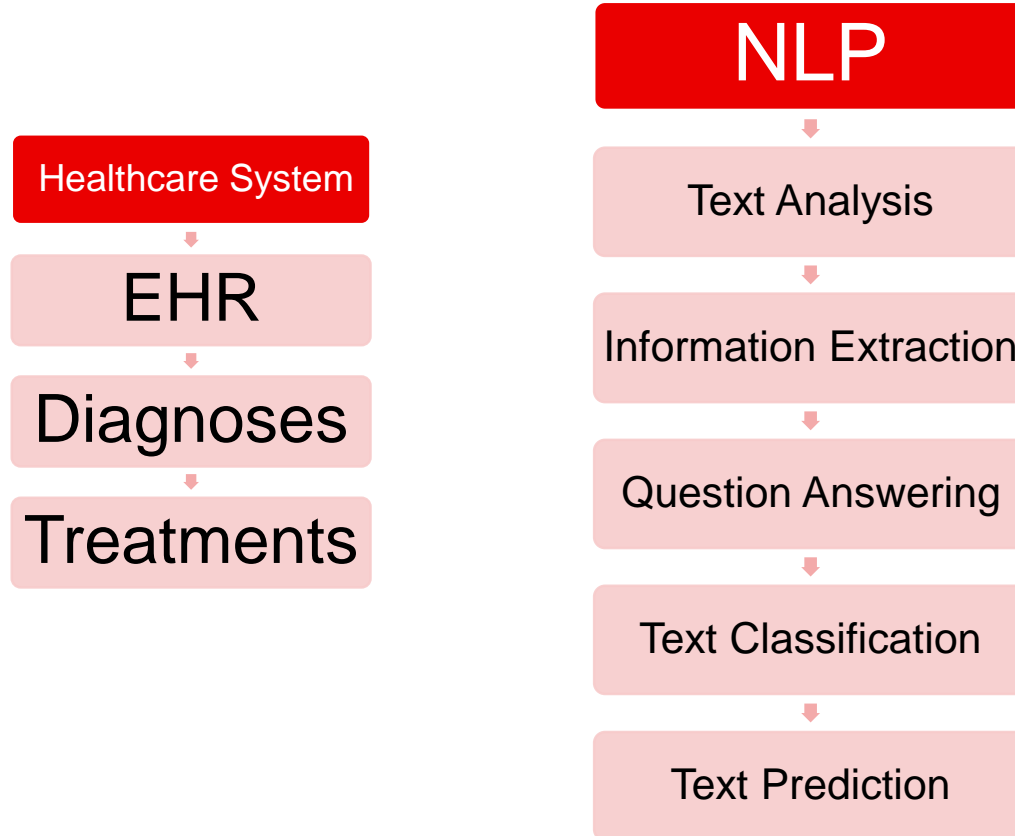


HealthIT.gov

@ONC_HealthIT

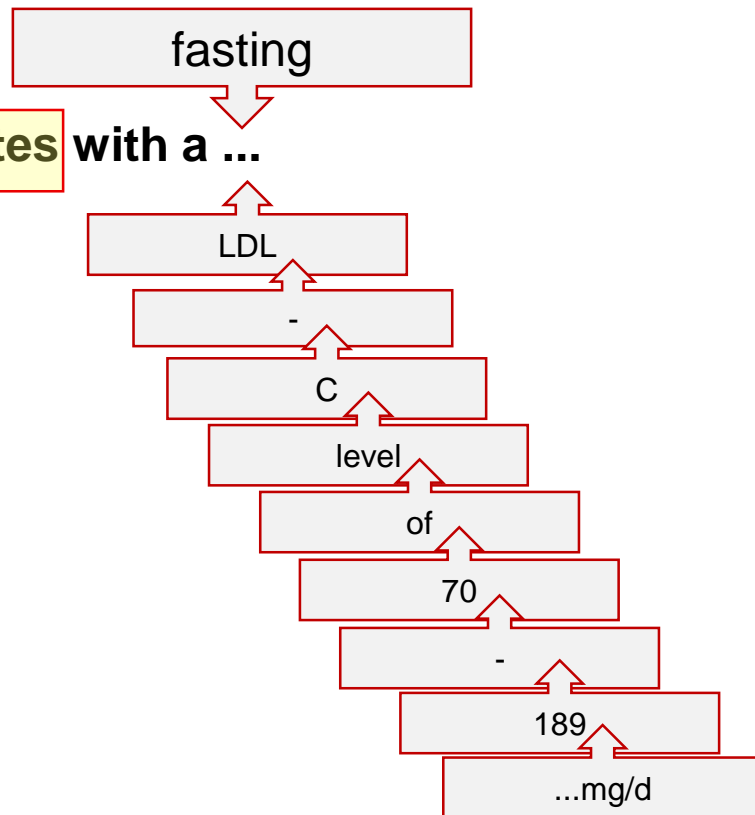
<https://www.healthit.gov/playbook/>

Utilizing NLP for Healthcare System



Language Modeling

Adults aged 40-75 years with a diagnosis of diabetes with a ...



Language Modeling

$k = 4$

Type 1 diabetes can develop at any age

However, some people with type 1 diabetes can develop insulin resistance.

type 1 diabetes can develop in people who have a particular HLA complex.

Detail about “diabetes” can be learned from the context (e.g., semantic similarity)

$$P(w_t | context) \forall t \in V$$

$$P(w_t | w_{t-k}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+k})$$

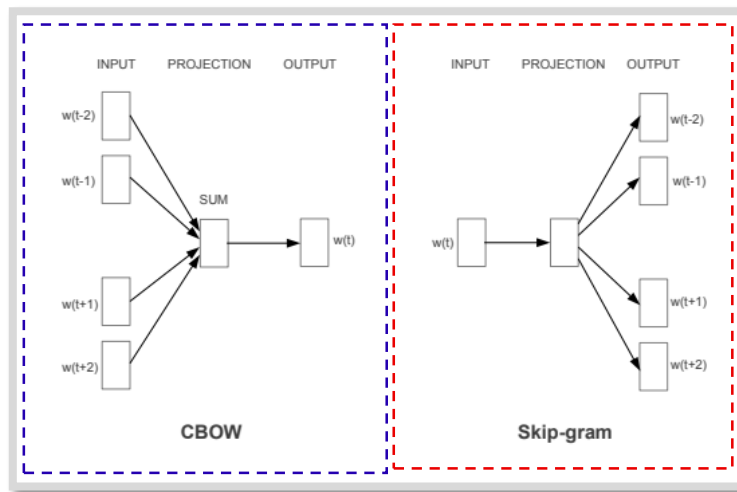
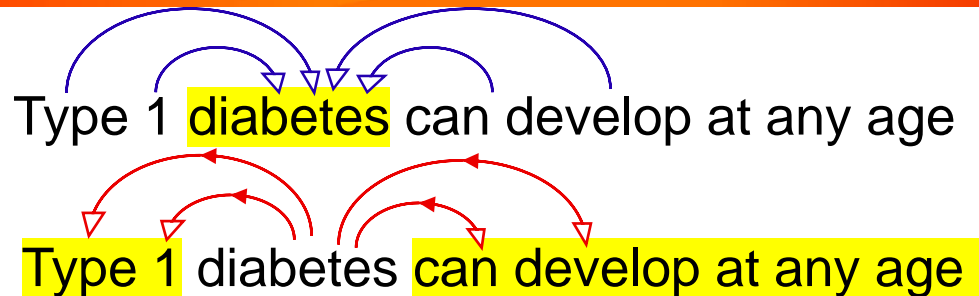
Language Modeling

○ Word2Vec

○ GloVe

○ Transformer

○ BERT



Mikolov, Tomas, et al. "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781* (2013).

Language Modeling

○ Word2Vec

○ GloVe

○ Transformer

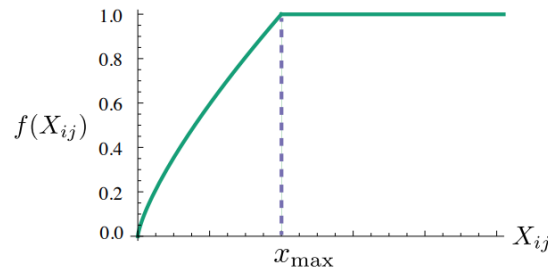
○ BERT

Frequent appearances
of **ice** and **solid**

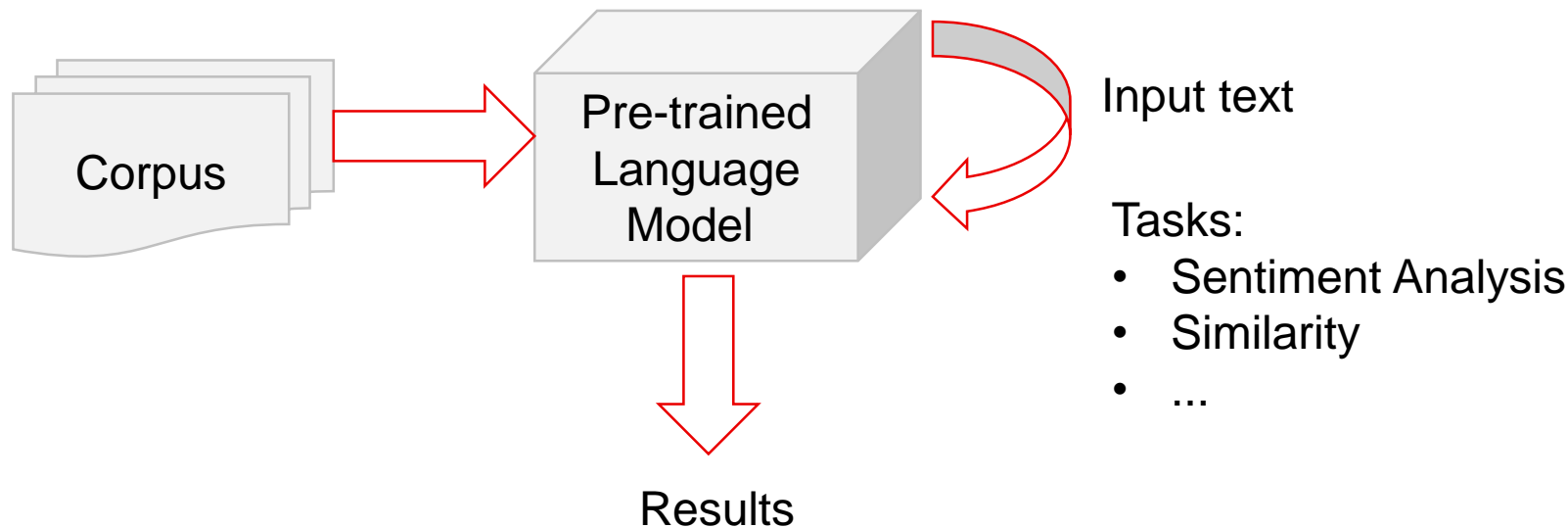
Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$	$k = \text{fashion}$
$P(k \text{ice})$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k \text{steam})$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k \text{ice})/P(k \text{steam})$	8.9	8.5×10^{-2}	1.36	0.96

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

Pennington, Jeffrey, Richard Socher, and Christopher D. Manning.
"Glove: Global vectors for word representation." *Proceedings of the
2014 conference on empirical methods in natural language
processing (EMNLP)*. 2014.

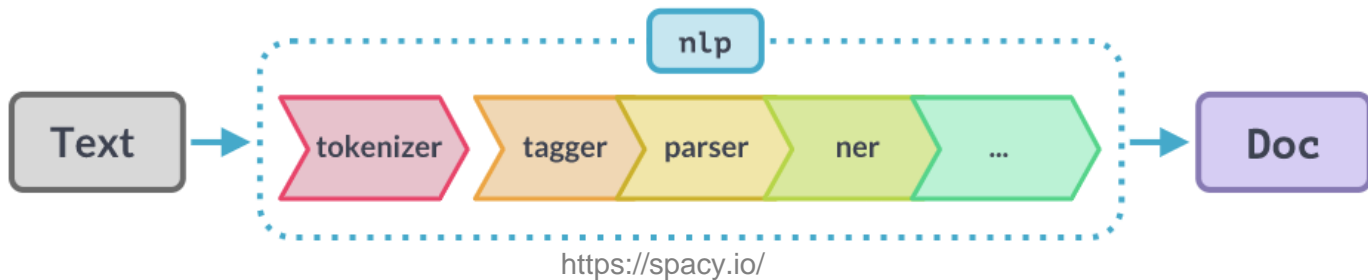


Text Analysis through Pre-trained Language Models



Text Analysis

- Text Cleaning
- Lemmatization
- Tokenizer
- Semantic



Code #1

jupyter Tutorial_NLP_Summit_2022_c1 Last Checkpoint: 8 hours ago (unsaved changes)



Logout

File Edit View Insert Cell Kernel Widgets Help

Trusted

Python 3

Run Stop Restart Clear Cell Toolbar

```
In [ ]: 1 import transformers
        2 import tensorflow as tf
        3 import keras as k
        4 import sys
        5 import seaborn as sns
        6 import numpy as np
        7 from spacy import displacy
        8 import matplotlib.pyplot as plt
        9 import logging
       10 logging.basicConfig(level=logging.ERROR)
       11 logger = logging.getLogger()
       12 logger.setLevel(logging.ERROR)
```

```
In [ ]: 1 print(f"transformers:{transformers.__version__}")
        2 print(f"tensorflow:{tf.__version__}")
        3 print(f"keras:{k.__version__}")
```

Sample EHR

```
In [ ]: 1 #Ref: https://medicalcodify.com/eh/webchart.cgi?f=layoutnouser&func=&module=&tabmodule=&name=RXDBmain&searchterm=
        2 corpus= ["Percentage of the following patients - all considered at high risk of cardiovascular events - who were
```

```
In [ ]: 1 corpus
```

Text Analysis

```
In [ ]: 1 import spacy
        2 try:
        3     nlp = spacy.load("en_core_web_lg", )
        4 except:
        5     !sys.executable -m spacy download en_core_web_lg
```

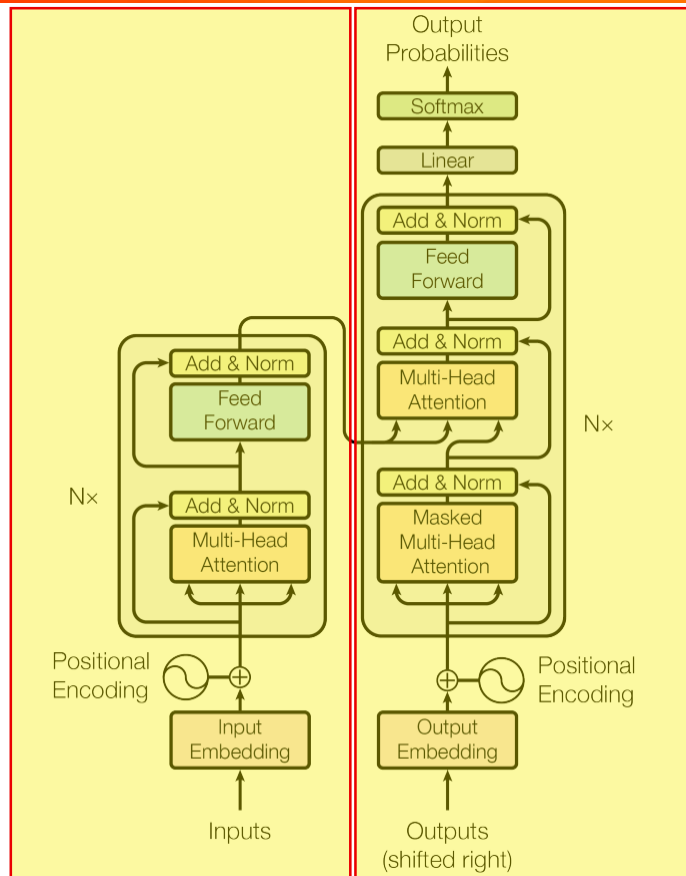
Language Modeling

○ Word2Vec

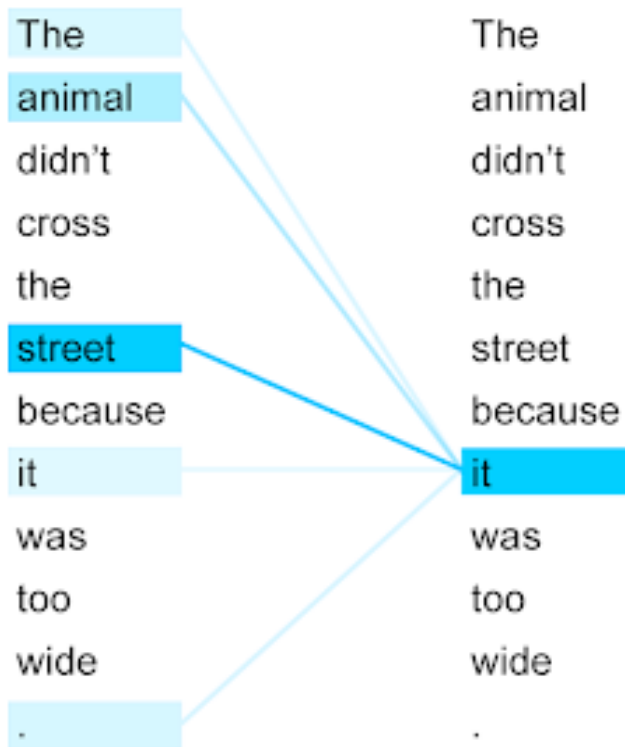
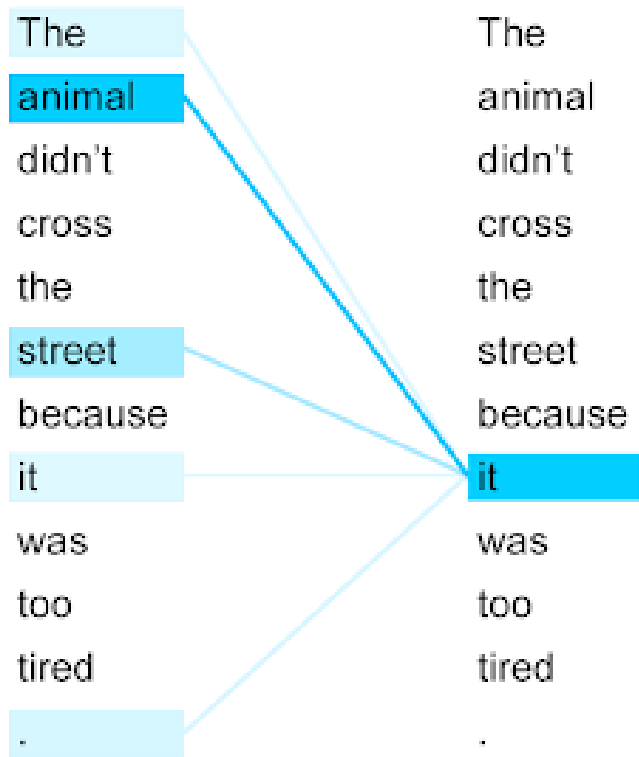
○ GloVe

○ **Transformer**

○ BERT



Mikolov, Tomas, et al. "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781* (2013).



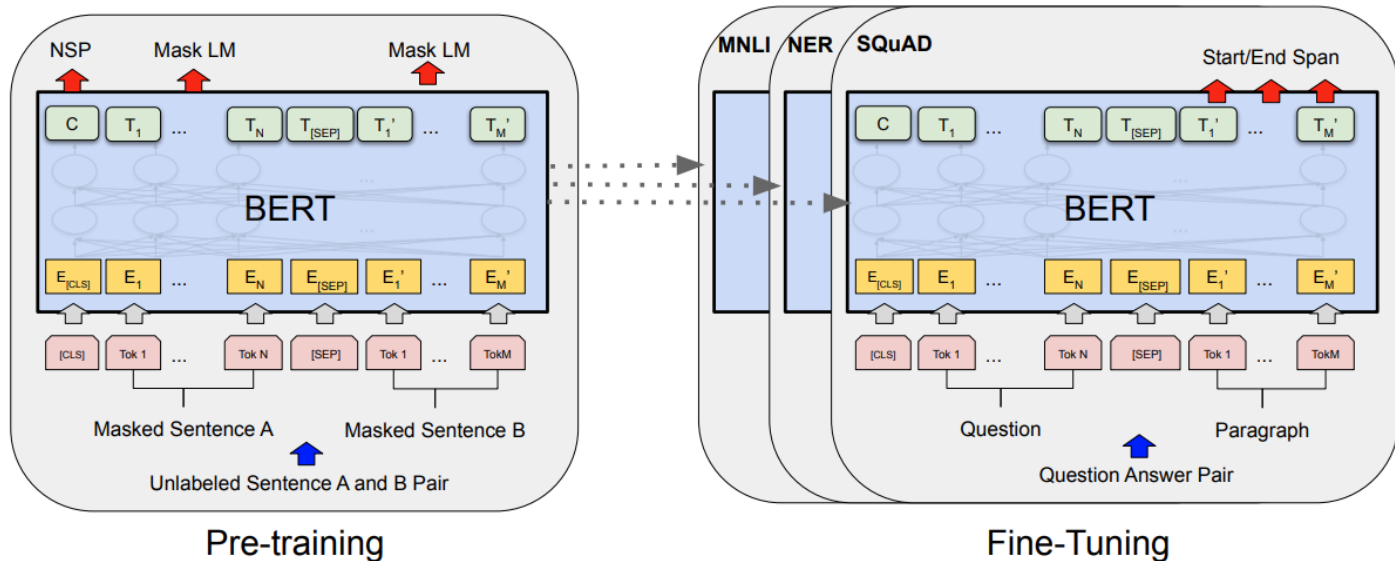
Language Modeling

○ Word2Vec

○ GloVe

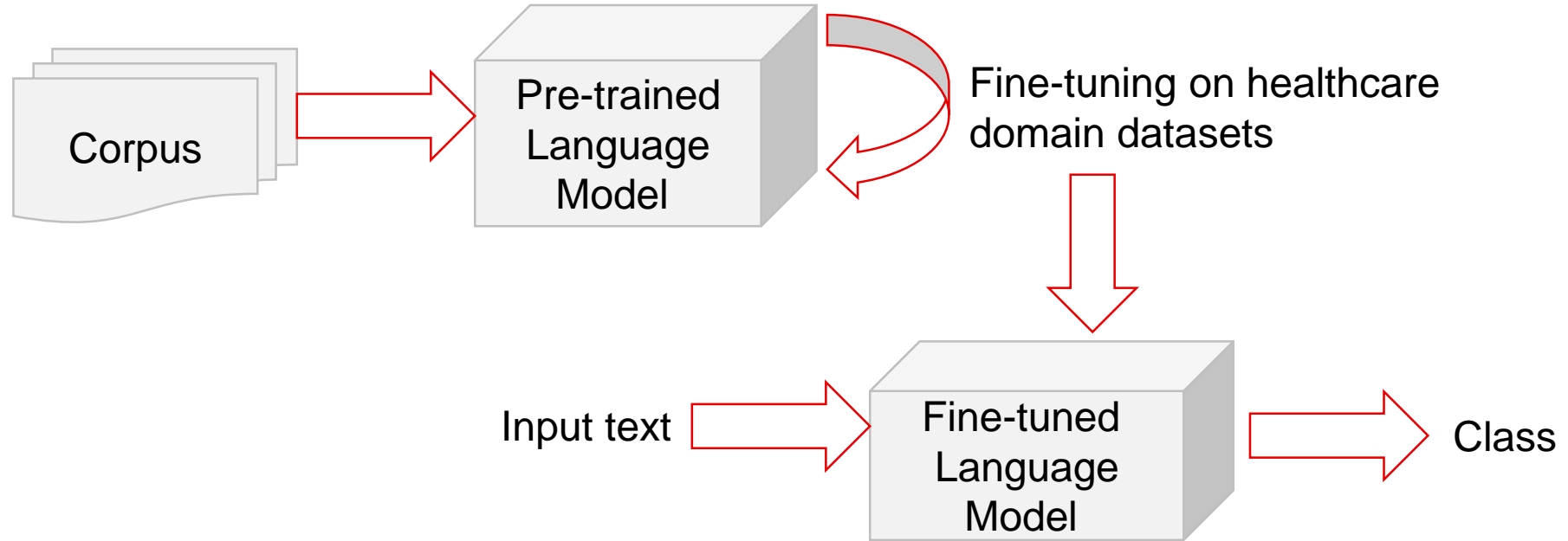
○ Transformer

○ BERT



Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).

Text Classification



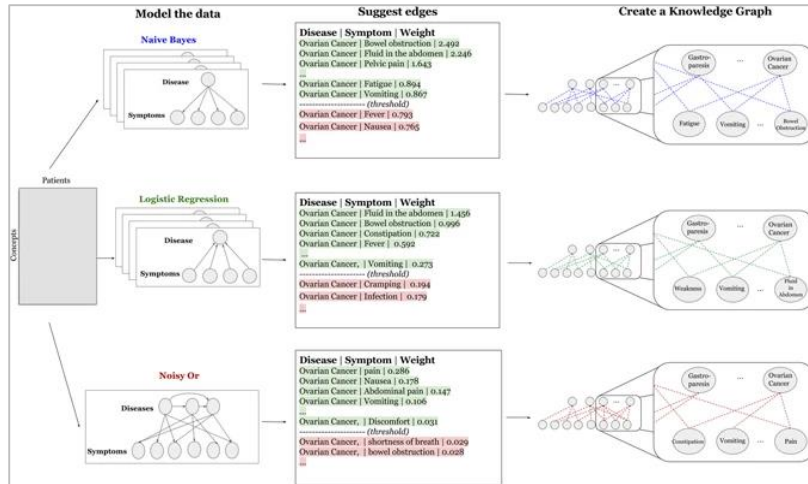
Question Answering

○ Structured Data

Y	filename	sex	age	age_years	corpus	group	child_TNW	child_TNS
1	922.cha	male	113	9.416666667	ENNI	SLI	1404	185
1	924.cha	male	112	9.333333333	ENNI	SLI	1162	125
1	926.cha	female	114	9.5	ENNI	SLI	729	80
1	928.cha	male	116	9.666666667	ENNI	SLI	592	63

Diagnose Specific Language Impairment in Children / Kaggle

Knowledge graph



Rotmensch, Maya, et al. "Learning a health knowledge graph from electronic medical records." *Scientific reports* 7.1 (2017): 1-11.

○ Unstructured Data

Description

Percentage of the following patients - all considered at high risk of cardiovascular events - who were prescribed or were on statin therapy during the measurement period: *Adults aged ≥ 21 years who were previously diagnosed with or currently have an active diagnosis of clinical atherosclerotic cardiovascular disease (ASCVD); OR *Adults aged ≥ 21 years who have ever had a fasting or direct low-density lipoprotein cholesterol (LDL-C) level ≥ 190 mg/dL or were previously diagnosed with or currently have an active diagnosis of familial or pure hypercholesterolemia; OR *Adults aged 40-75 years with a diagnosis of diabetes with a fasting or direct LDL-C level of 70-189 mg/dL

<https://medicalcodify.com/eh/webchart.cgi>

Question Answering

○ Open Generative Question Answering

Description

Percentage of the following patients - all considered at high risk of cardiovascular events - who were prescribed or were on statin therapy during the measurement period: *Adults aged ≥ 21 years who were previously diagnosed with or currently have an active diagnosis of clinical atherosclerotic cardiovascular disease (ASCVD); OR *Adults aged ≥ 21 years who have ever had a fasting or direct low-density lipoprotein cholesterol (LDL-C) level ≥ 190 mg/dL or were previously diagnosed with or currently have an active diagnosis of familial or pure hypercholesterolemia; OR *Adults aged 40-75 years with a diagnosis of diabetes with a fasting or direct LDL-C level of 70-189 mg/dL

"What are previous diagnoses?"

Source: <https://medicalcodify.com/eh/webchart.cgi>

Question Answering

○ Extractive Question Answering

Description

Percentage of the following patients - all considered at high risk of cardiovascular events - who were prescribed or were on statin therapy during the measurement period: *Adults aged ≥ 21 years who were previously diagnosed with or currently have an active diagnosis of clinical **atherosclerotic cardiovascular disease** (ASCVD); OR *Adults aged ≥ 21 years who have ever had a fasting or direct low-density lipoprotein cholesterol (LDL-C) level ≥ 190 mg/dL or were previously diagnosed with or currently have an active diagnosis of familial or pure hypercholesterolemia; OR *Adults aged 40-75 years with a diagnosis of diabetes with a fasting or direct LDL-C level of 70-189 mg/dL

Source: <https://medicalcodify.com/eh/webchart.cgi>

"What is the disease of the patients?"

Question Answering

○ Closed Generative Question Answering

Description

Percentage of the following patients - all considered at high risk of cardiovascular events - who were prescribed or were on statin therapy during the measurement period: *Adults aged ≥ 21 years who were previously diagnosed with or currently have an active diagnosis of clinical atherosclerotic cardiovascular disease (ASCVD); OR *Adults aged ≥ 21 years who have ever had a fasting or direct low-density lipoprotein cholesterol (LDL-C) level ≥ 190 mg/dL or were previously diagnosed with or currently have an active diagnosis of familial or pure hypercholesterolemia; OR *Adults aged 40-75 years with a diagnosis of diabetes with a fasting or direct LDL-C level of 70-189 mg/dL

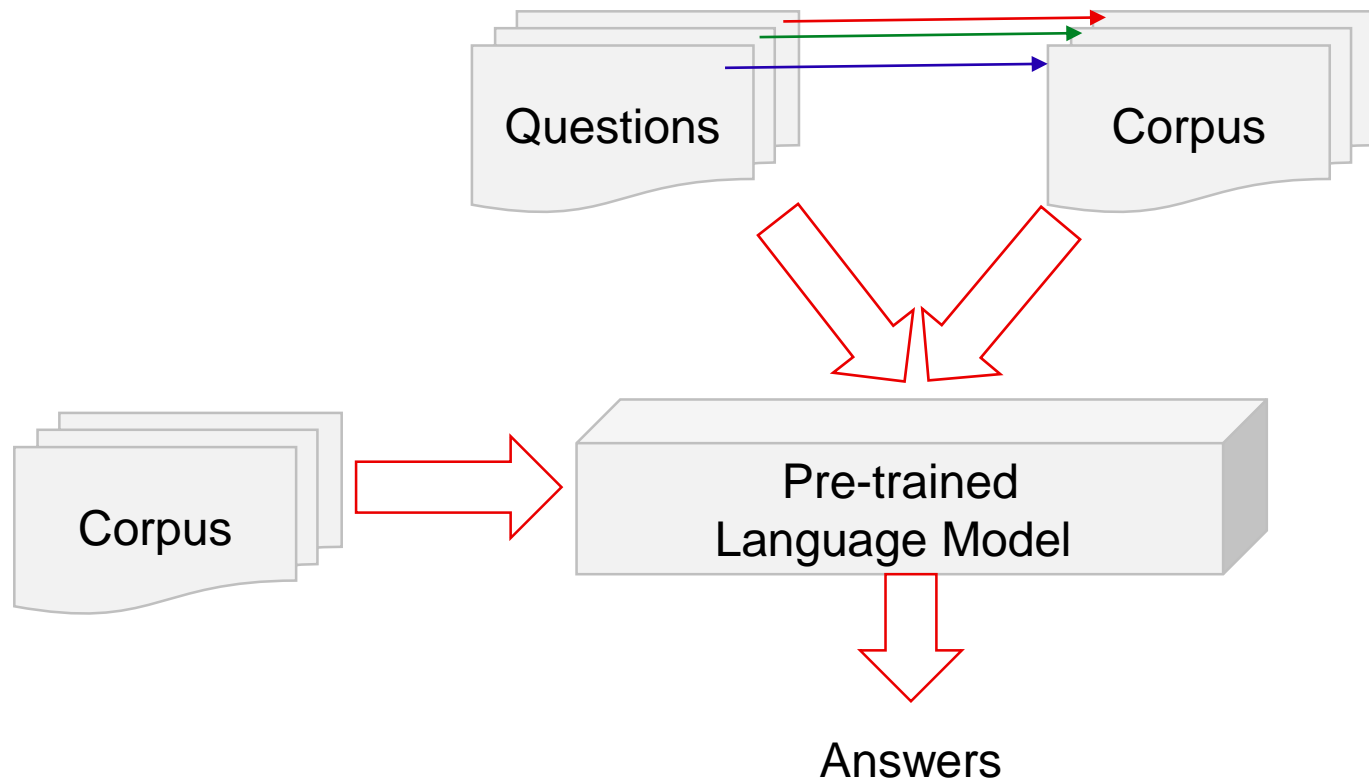
Source: <https://medicalcodify.com/eh/webchart.cgi>

"What are prevention methods for ASCVD?"

The screenshot displays two web pages side-by-side. The left page is the Wikipedia article for 'Cardiovascular disease', which includes a sidebar with navigation links and a main content area with introductory text. The right page is the Mayo Clinic website, specifically the 'Arteriosclerosis / atherosclerosis' section. This page features a search bar, navigation tabs for 'Symptoms & causes', 'Diagnosis & treatment', 'Doctors & departments', and 'Care at Mayo Clinic'. It also includes an 'Overview' section with text about the condition, a 'Request an Appointment' button, and a 'Products & Services' section at the bottom.

<https://www.mayoclinic.org/diseases-conditions/arteriosclerosis-atherosclerosis/symptoms-causes/syc-20350569>
https://en.wikipedia.org/wiki/Cardiovascular_disease

Question Answering



Code #2

Jupyter Tutorial_NLP_Summit_2022_c1 Last Checkpoint: 10 hours ago (unsaved changes)



Logout

File Edit View Insert Cell Kernel Widgets Help

Trusted

Python 3

Run Stop Restart Clear All Run and Restart Run and Clear

Markdown

0.40659, -0.73905, 0.44078, -0.012808], dtype=float32)

Question Answering

```
In [ ]: 1 from transformers import pipeline
```

```
In [ ]: 1 qa_model = pipeline("question-answering",  
2                               model="distilbert-base-cased-distilled-squad")  
3 #for more detail refer to: https://github.com/huggingface/notebooks/blob/master/examples/question\_answering.ipynb
```

```
In [ ]: 1 corpus[0]
```

```
In [ ]: 1 question_1="What are previous diagnosis?"  
2 answer_1 = qa_model(question = question_1, context = corpus[0])  
3 answer_1
```

```
In [ ]: 1 question_2="What is the disease of the patients?"  
2 answer_2 = qa_model(question = question_2, context = corpus[0])  
3 answer_2
```

```
In [ ]: 1 from sentence_transformers import SentenceTransformer, util  
2 import numpy as np
```

```
In [ ]: 1 #ref: https://huggingface.co/allenai/biomed\_roberta\_base  
2 model = SentenceTransformer("allenai/biomed_roberta_base")
```

```
In [ ]: 1 classifier = pipeline("text-classification", model = "roberta-large-mnli")
```

```
In [ ]: 1 mnli_context="Patients were previously diagnosed with atherosclerotic cardiovascular disease or were previously
```

```
In [ ]: 1 classifier(f"{mnli_context} Does patient previously diagnosed with diabetes?")
```

Text Classification

○ Entailment

Description

Percentage of the following patients - all considered at high risk of cardiovascular events - who were prescribed or were on statin therapy during the measurement period: *Adults aged ≥ 21 years who were previously diagnosed with or currently have an active diagnosis of clinical **atherosclerotic cardiovascular disease** (ASCVD); OR *Adults aged ≥ 21 years who have ever had a fasting or direct low-density lipoprotein cholesterol (LDL-C) level ≥ 190 mg/dL or were previously diagnosed with or currently have an active diagnosis of familial or pure hypercholesterolemia; OR *Adults aged 40-75 years with a diagnosis of diabetes with a fasting or direct LDL-C level of 70-189 mg/dL

Diagnoses Dataset

Heart Disease

Cancer

Diabetes

...

Code #3

Jupyter Tutorial_NLP_Summit_2022_c1 Last Checkpoint: 11 hours ago (autosaved)



Logout

File Edit View Insert Cell Kernel Widgets Help

Notebook saved

Trusted



Python 3

Run Code

QA & Text Classification

```
In [ ]: 1 question_2, answer_2["answer"]
```

```
In [ ]: 1 corpus[0]
```

```
In [ ]: 1 diseases = ["Heart Disease", "Cancer", "Diabetes"]
```

```
In [ ]: 1 embeddings=[]  
2 cosine_scores=[]  
3 answer_embedding = model.encode(answer_2["answer"], convert_to_tensor=True)  
4 for disease in diseases:  
5     target=model.encode(disease, convert_to_tensor=True)  
6     embeddings.append(target)  
7     cosine_scores.append(util.pytorch_cos_sim(target, answer_embedding))
```

```
In [ ]: 1 len(embeddings), embeddings[0].shape
```

```
In [ ]: 1 embeddings[0]
```

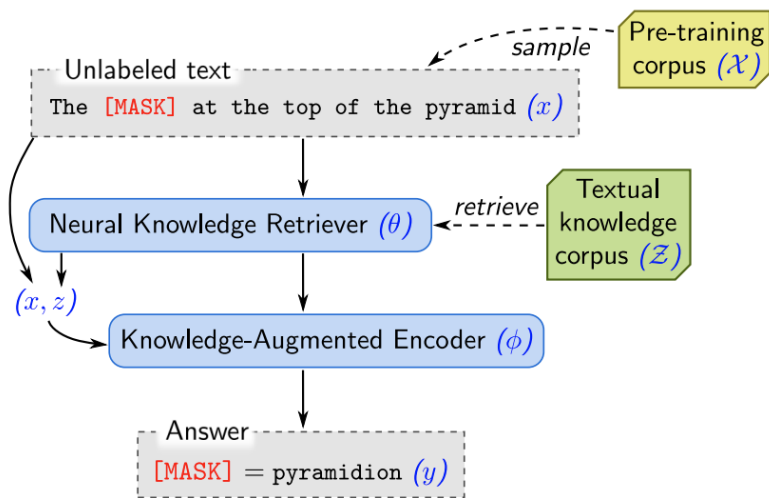
```
In [ ]: 1 cosine_scores
```

```
In [ ]: 1 for index in range(len(diseases)):  
2     print(f"{diseases[index]}:{cosine_scores[index]}")
```

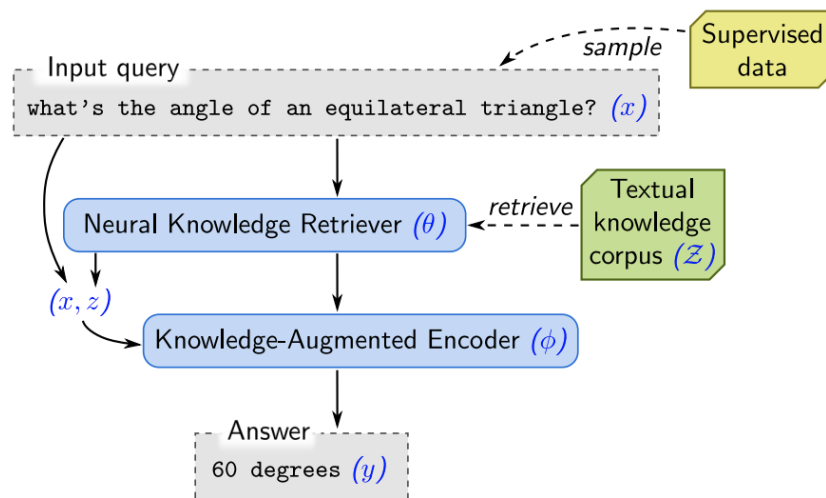
```
In [ ]: 1 top_related_answer= np.argmax(cosine_scores)  
2 diseases[top_related_answer]
```

```
In [ ]: 1
```

QA & Text Classification / Training



Unsupervised pre-training



Supervised fine-tuning

Code #4

jupyter Tutorial_NLP_Summit_2022_c1 Last Checkpoint: 12 hours ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

Out[53]: `tensor([[0.9217]]), tensor([[0.9137]]), tensor([[0.9208]])`

```
In [54]: 1 for index in range(len(diseases)):
          2     print(f"{diseases[index]}:{cosine_scores[index]}")

Heart Disease:tensor([[0.9217]])
Cancer:tensor([[0.9137]])
Diabetes:tensor([[0.9208]])
```

```
In [55]: 1 top_related_answer = np.argmax(cosine_scores)
          2 diseases[top_related_answer]
```

Out[55]: 'Heart Disease'

Masked Language Model

```
In [ ]: 1 from transformers import pipeline
          2 mlm = pipeline("fill-mask",
          3                 model="microsoft/BiomedNLP-PubMedBERT-base-uncased-abstract-fulltext")

In [ ]: 1 #original: all considered at high risk of cardiovascular
          2 predicted_tokens = mlm("all considered at high [MASK]")
          3 predicted_tokens = [_ for _ in predicted_tokens if _['token_str'] not in ['.', ',', ';']]
          4 predicted_tokens

In [90]: 1

In [ ]: 1
```

Conclusion

- Based on use case, different language models can be used such as Word2Vec, GloVe, BERT, GPT
- Pre-trained language models can be used easily and works well across different domains
- The latest Pre-trained language models are expensive
- The pre-trained models can be fine-tune on specific down-stream tasks

Thank you

?

?Question?

?

