



# The Attack of the Muppets!

Data Management in the Era of  
Pretrained Transformers

**Jimmy Lin**

David R. Cheriton School of Computer Science  
University of Waterloo

ICDE 2021 Keynote  
Thursday, April 22, 2021



UNIVERSITY OF  
**WATERLOO**

# Message



(Historically)  
NLP and database  
researchers don't have much  
to say to each other



**Changes are coming!**



One word: **BERT!**



SEARCH

# Understanding searches better than ever before

Pandu Nayak  
Google Fellow and Vice President, Search

Published Oct 25, 2019

If there's one thing I've learned over the 15 years working on Google Search, it's that people's curiosity is endless. We see billions of searches every day, and 15 percent of those queries are ones we haven't seen before--so we've built ways to return results for queries we can't anticipate.

When people like you or I come to Search, we aren't always quite sure about the best way to formulate a query. We might not know the right words to use, or how to spell something, because often times, we come to Search looking to learn--we don't necessarily have the knowledge to begin with.

At its core, Search is about understanding language. It's our job to figure out what you're searching for and surface helpful information from the web, no matter how you spell or combine the words in your query. While we've continued to improve our language understanding capabilities over the years, we sometimes still don't quite get it right, particularly with complex or conversational queries. In fact, that's one of the reasons why people often use "keyword-ese," typing strings of words that they think we'll understand, but aren't actually how they'd naturally ask a question.

With the latest advancements from our research team in the science of language understanding--made possible by machine learning--we're making a significant improvement to how we understand queries, representing the biggest leap forward in the past five years, and one of the biggest leaps forward in the history of Search.

## Applying BERT models to Search

Last year, we introduced and open-sourced a neural network-based technique for natural language processing (NLP) pre-training called Bidirectional Encoder Representations from Transformers, or as we call it--BERT, for short. This technology enables anyone to train their own state-of-the-art question answering system.

This breakthrough was the result of Google research on [transformer models](#) that process words in relation to all the other words in a sentence, rather than in a one-by-one order. BERT models can therefore consider the full context of a word by looking at the words that come before and after it--particularly useful for understanding the intent behind search queries.

# BERT!

<https://blog.google/products/search/search-language-understanding-bert/>



# Bing delivers its largest improvement in search experience using Azure GPUs

Posted on November 18, 2019



[Jeffrey Zhu](#), Program Manager, Bing Platform

Over the last couple of years, deep learning has become widely adopted across the Bing search stack and powers a vast number of our intelligent features. We use natural language models to improve our core search algorithm's understanding of a user's search intent and the related webpages so that Bing can deliver the most relevant search results to our users. We rely on deep learning computer vision techniques to enhance the discoverability of billions of images even if they don't have accompanying text descriptions or summary metadata. We leverage machine-based reading comprehension models to retrieve captions within larger text bodies that directly answer the specific questions users have. All these enhancements lead toward more relevant, contextual results for web search queries.

Recently, there was a breakthrough in natural language understanding with a type of model called transformers (as popularized by Bidirectional Encoder Representations from Transformers, [BERT](#)). Unlike previous deep neural network (DNN) architectures that processed words individually in order, transformers understand the context and relationship between each word and all the words around it in a sentence. Starting from April of this year, we used large transformer models to deliver the largest quality improvements to our Bing customers in the past year. For example, in the query "what can aggravate a concussion", the word "aggravate" indicates the user wants to learn about actions to be taken after a concussion and not about causes or symptoms. Our search powered by these models can now understand the user intent and deliver a more useful result. More importantly, these models are now applied to every Bing search query globally making Bing results more relevant and intelligent.

## BEFORE

The screenshot shows the Bing search interface with the query "what can aggravate a concussion". The results page displays a list of links, with the top result being a PDF titled "Facts About Concussion and Brain Injury". The description of this result includes text about concussions occurring from falls or blows to the head, and mentions that they are usually not life-threatening.

## AFTER

The screenshot shows the same search query "what can aggravate a concussion" on the Bing search results page after the BERT update. The results are more refined, showing links that specifically address what can aggravate a concussion, such as a link to a Cleveland Clinic article titled "Suspect a Concussion? How to Protect Your Recovery ...". The description of this result emphasizes avoiding activities like exercise, which can worsen symptoms.

**BERT!**

NLP

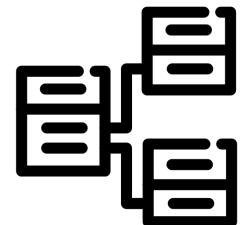
DB

(Historically)

NLP and database  
researchers don't have much  
to say to each other



Unstructured  
Text



Structured  
Relations



Source: <https://www.flickr.com/photos/mathiasappel/27065384106/>



Source: <https://www.flickr.com/photos/leszekleszczynski/15447205451/>

NLP

DB

Gartner®

Blog home > Blog post



# Big Content: The Unstructured Side of Big Data

By Darin Stewart | May 01, 2013 |

1 Comment

Search all blog posts



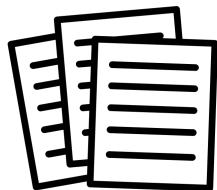
Search

Open Data

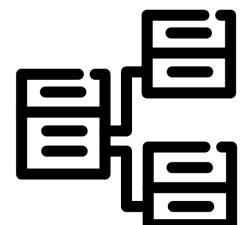
Enterprise Content Management

Big Content

The age of information overload is slowly drawing to a close. The enterprise is finally getting comfortable with managing massive amounts of data, content and information. The pace of information creation continues to accelerate, but the ability of infrastructure and information management to keep pace is coming within sight. Big data is now considered a blessing rather than a curse. Even so, managing information is not the same as fully exploiting information. While 'Big Data' technologies and techniques are unlocking secrets previously hidden in enterprise data, the largest source of potential insight remains largely untapped. Unstructured content represents as much as eighty percent of an organizations total information assets. While Big Data technologies and techniques are well suited to exploring unstructured information, this 'Big Content' remains grossly underutilized and its potential largely unexplored.



Unstructured  
Text



Structured  
Relations



# NLP

# DB

## A productive dialogue...

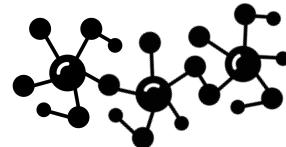
Hey look, we can extract entities, relations, and facts from free text!

A knowledge graph!

And this is something we know how to build!



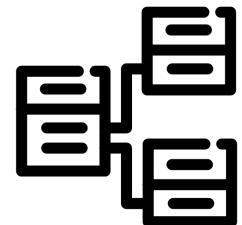
Unstructured  
Text



Knowledge  
Graphs

A knowledge graph!

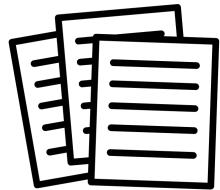
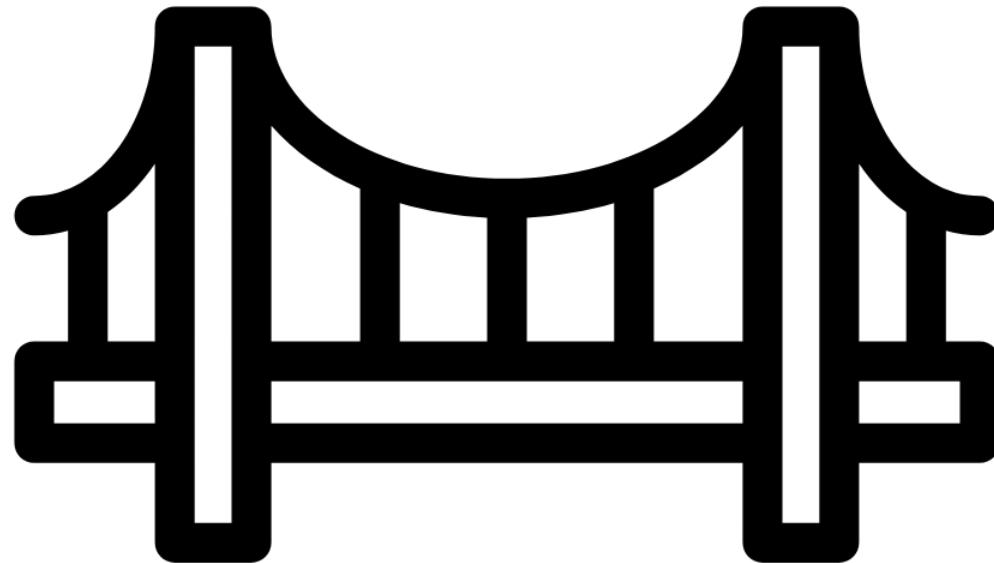
Cool, this is something we can work with!



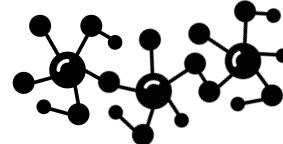
Structured  
Relations

# NLP

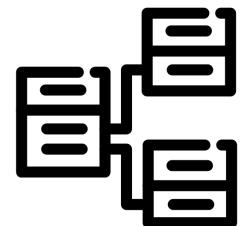
# DB



Unstructured  
Text



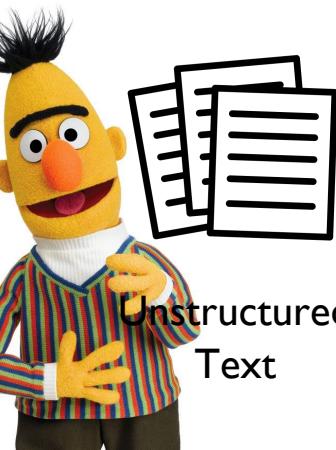
Knowledge  
Graphs



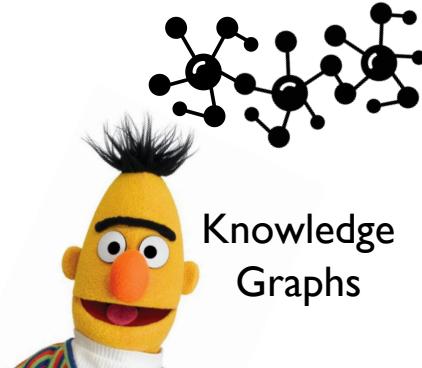
Structured  
Relations

# NLP

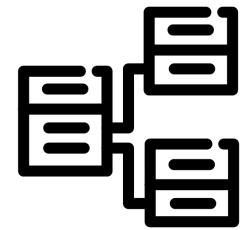
# DB



Unstructured  
Text



Knowledge  
Graphs



Structured  
Relations

NLP

DB

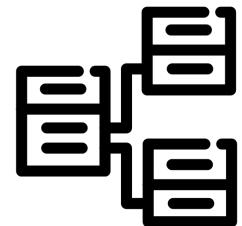
What if BERT took over  
everything?



Unstructured  
Text



Knowledge  
Graphs



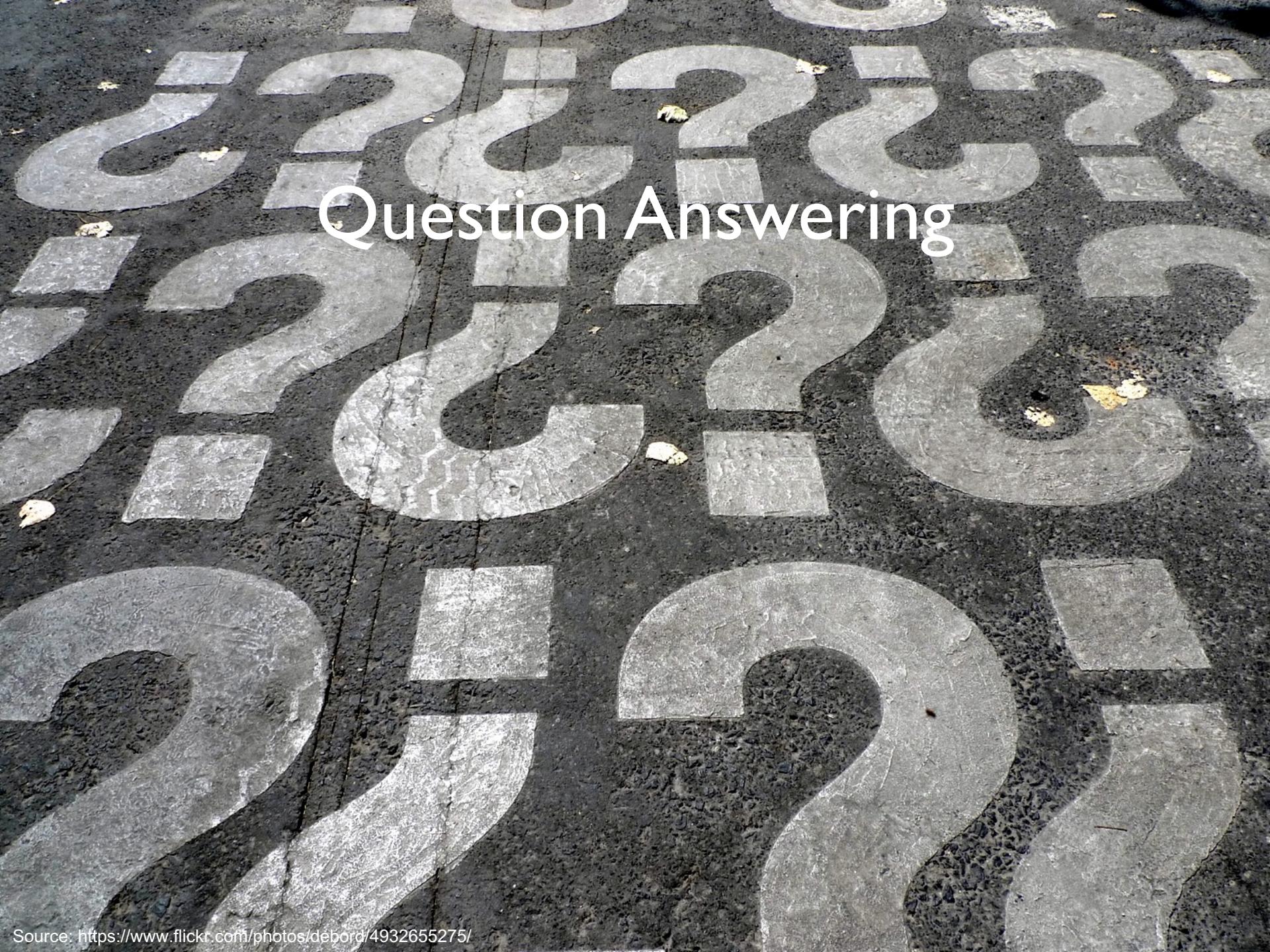
Structured  
Relations

# Message



## Messenger





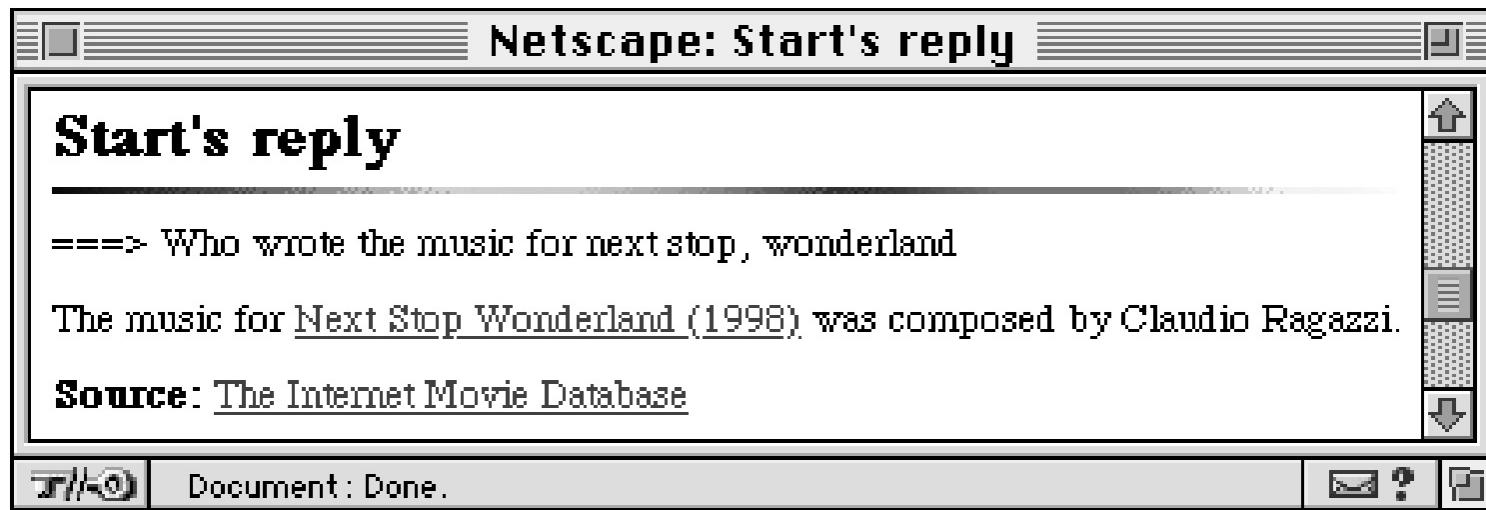
# Question Answering

# 1997: My journey begins

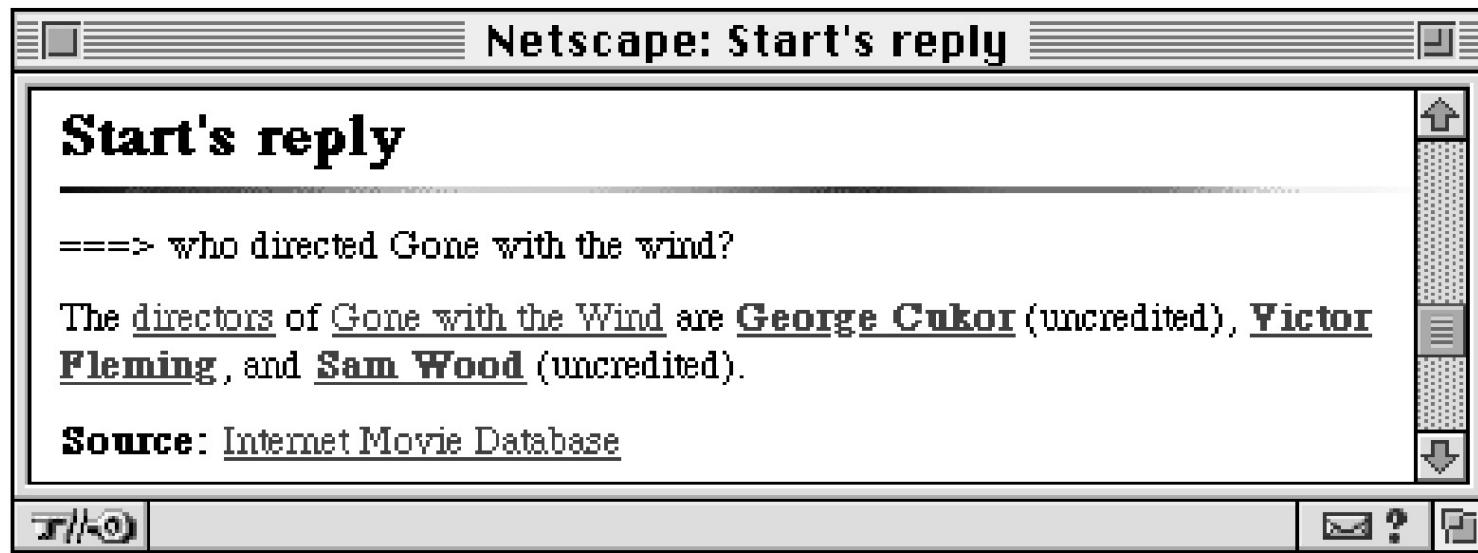


# 1993: The START System

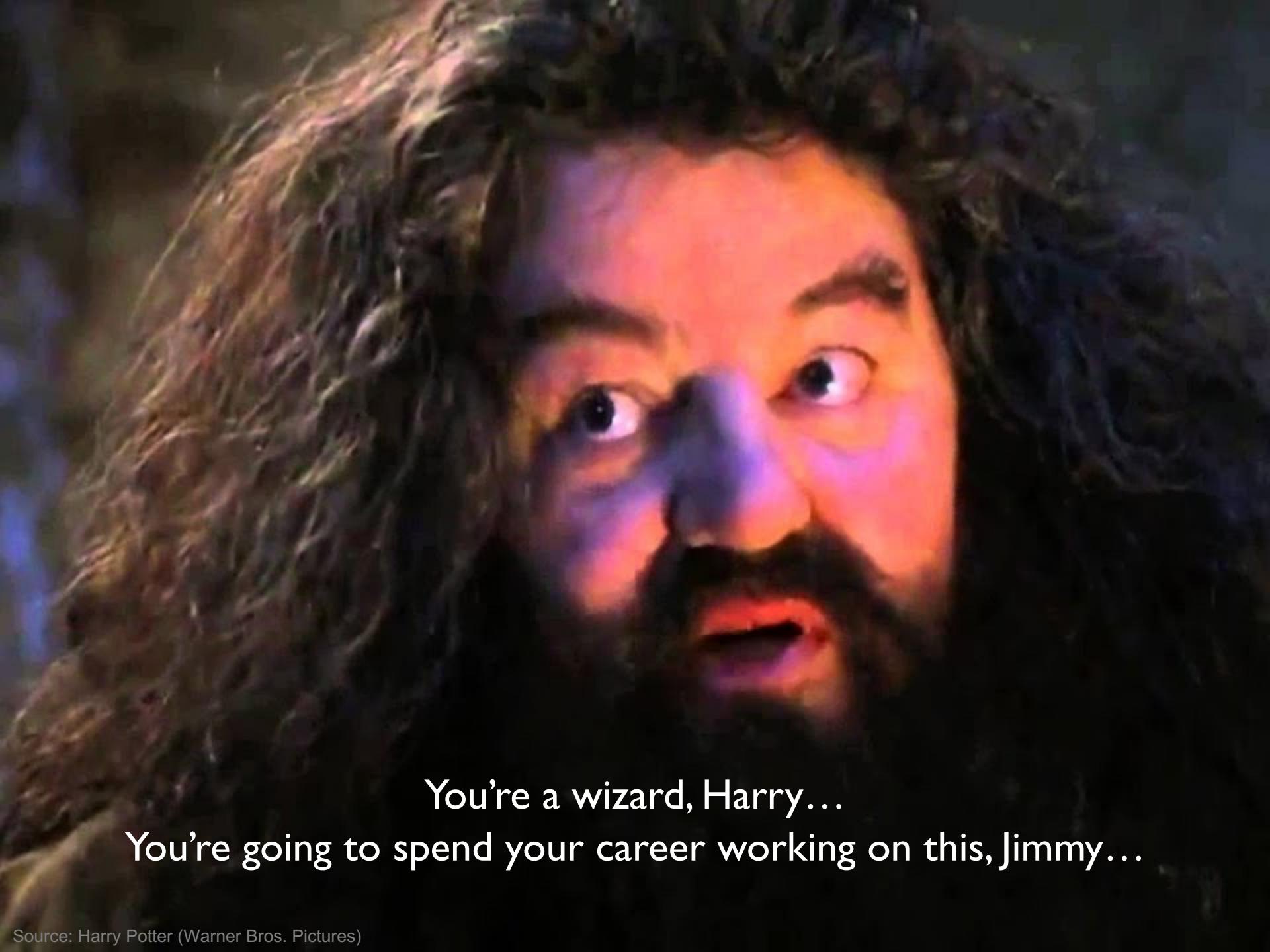
## First QA system on the web!



<http://start.csail.mit.edu/>



<http://start.csail.mit.edu/>



You're a wizard, Harry...  
You're going to spend your career working on this, Jimmy...

# My career-long quest...

Connecting users with relevant information



# My career-long quest...

Connecting users with relevant information

What? **text**, speech, images, graphs, semi-structured data, relational data, log data...

Who? **general information seekers**, domain experts, legal scholars, historians, data scientists, etc.

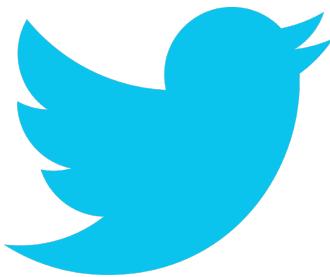
## Information Access

(question answering, document retrieval, summarization, ...)

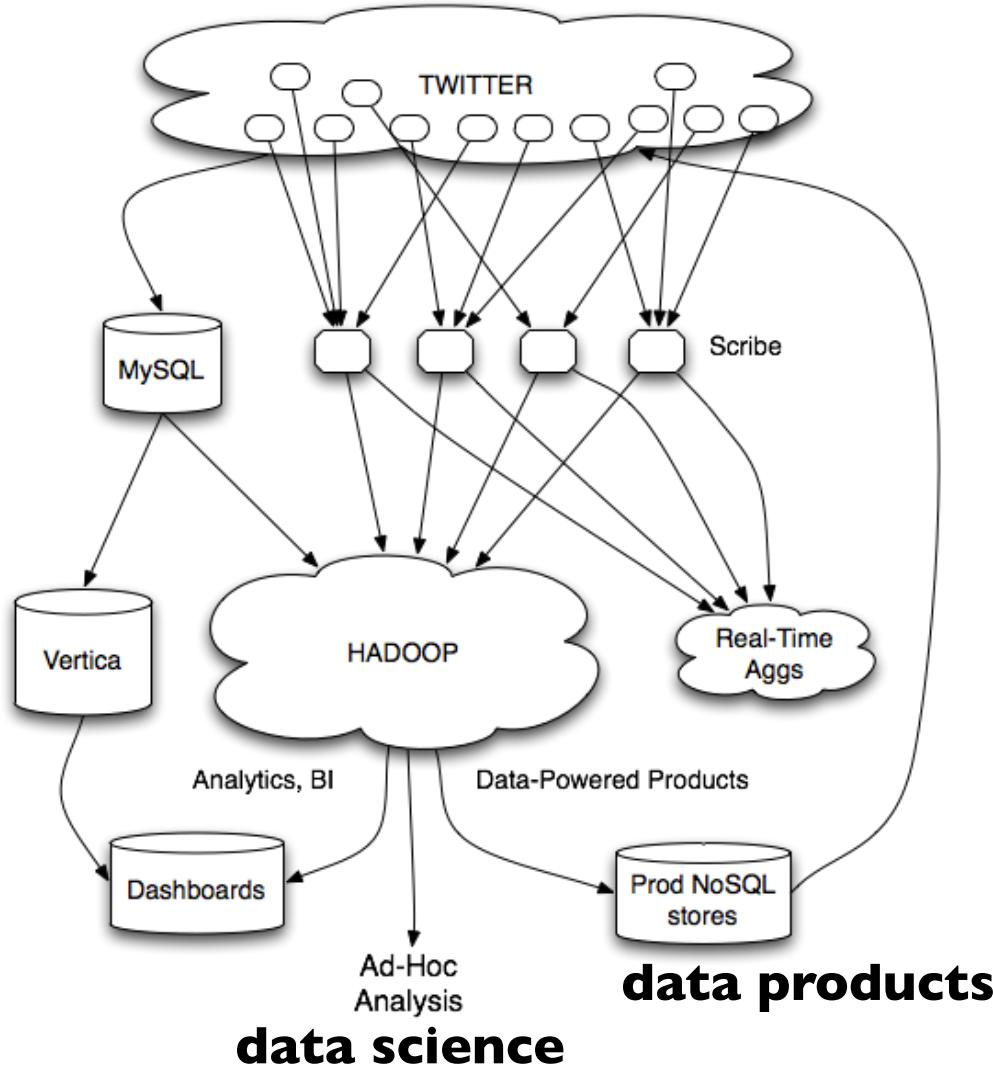


# Question Answering

I also dabble in data management



Circa 2012



I worked on...

- **analytics infrastructure to support data science**
- **data products to connect users to relevant content**



## Tweets

Mishne et al. Fast Data in the Era of Big Data: Twitter's Real-Time Related Query Suggestion Architecture. SIGMOD 2013.



Struggling with complex data  
of Data Science 2/20 to rehi

► Promoted by Cloudera

TWEETS FOLLOWING FOLLOWERS  
1,64

Leibert et al. Automatic Management of Partitioned,  
Replicated Search Services. SoCC 2011

Compose new Tweet...

Who to follow · Refresh · View all

plotly @plotlygraphs

+ Follow

► Promoted

Brad Anderson @boorad

Followed by Florian Leibert ...

+ Follow

Sheila Morrissey @sheilaMorr

+ Follow

Popular accounts · Find friends

## Trends

+ Change

#Olymp

I worked on...

Ukraine

#ConfessYourUnpopulOpinion

Venny

#PremioLoNuestro



Clinton Paquin @clintonpaquin  
Simply stated, "The only problem is muscle memory" @TheChange

► View conversation



The Hill @thehill · 1h

Republicans take debt ceiling fight to Senate

► View summary



Retweeted by Alex Feinberg

Popehat @Popehat · 10h

In a world in which few things are certain, this feed does.

Expand



#Sochi2014

#SochiProblems

Sochi

#SochiFail

Sochi 2014 ✅ @Sochi2014

Sochi Olympics 2014 @2014Sochi

Игры Сочи 2014 ✅ @sochi2014\_ru

Sochi Problems @SochiProblem

NYT Olympics @SochiNYT

Sochi Problems @SochiProblems

Search all people for sochi

Reply Retweet Favorite More

Gupta et al. WTF: The Who to Follow Service at Twitter. WWW 2013

Lin and Kolcz. Large-Scale Machine Learning at Twitter. SIGMOD 2012

@mollycooper and @BBCsasak

► View summary

... More

Expand

... More

# What is BERT?

# Why should you care?

# What is BERT?



# What's BERT?

a transformer encoder pretrained with a mask language model (MLM) objective that generates contextualized representations and can be fine-tuned for a wide range of NLP tasks.

## What's a monad?

a monad is a monoid in the category of endofunctors

What's the problem?





# What's BERT?

❶ a transformer encoder pretrained  
with a mask language model  
(MLM) objective that generates  
❷ contextualized representations  
and can be fine-tuned for a wide  
range of NLP tasks.❸

- ❶ What's a transformer?
- ❷ What are contextualized representations?
- ❸ What's pretraining?
- ❹ What's fine-tuning?

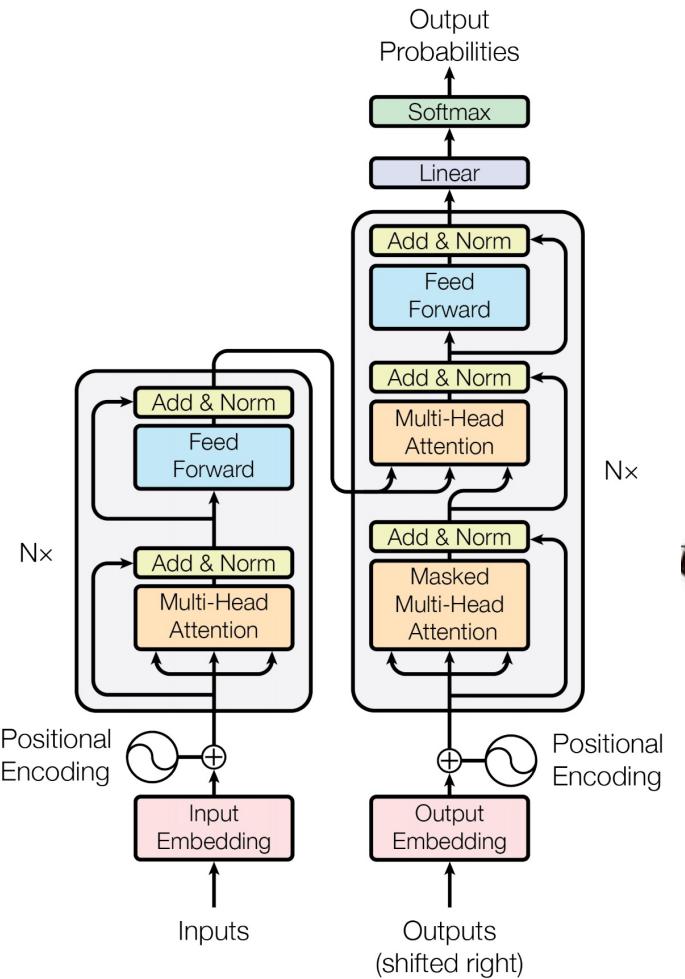


# What's BERT?

1 a **transformer encoder** pretrained  
with a mask language model  
(MLM) objective that generates  
contextualized representations  
and can be fine-tuned for a wide  
range of NLP tasks. 4

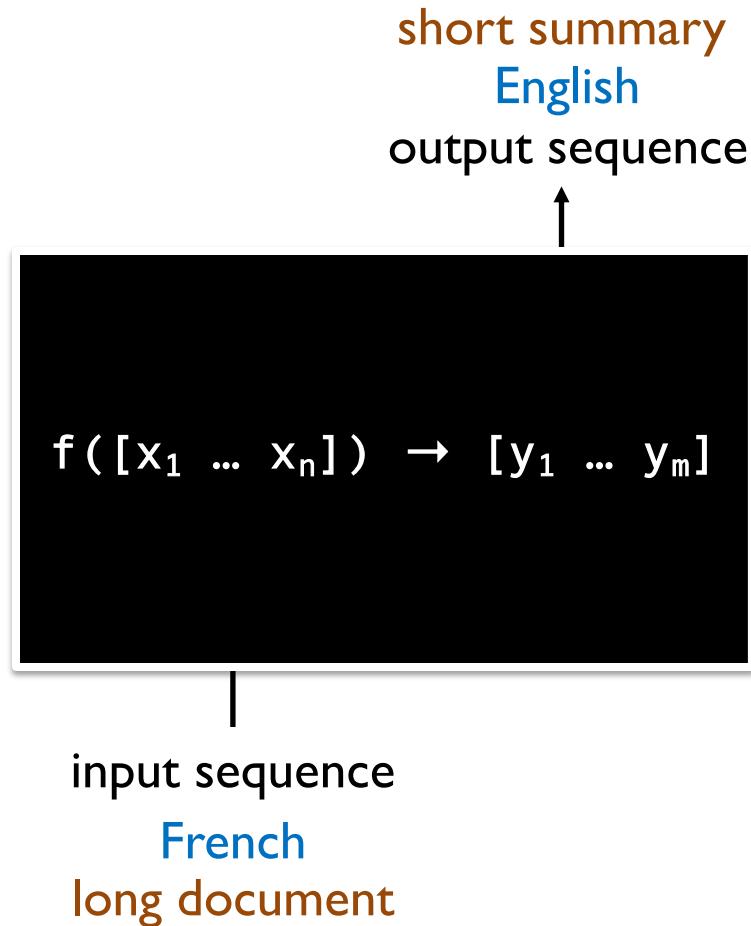
- 1 What's a transformer?
- 2 What are contextualized representations?
- 3 What's pretraining?
- 4 What's fine-tuning?

# What's a transformer?



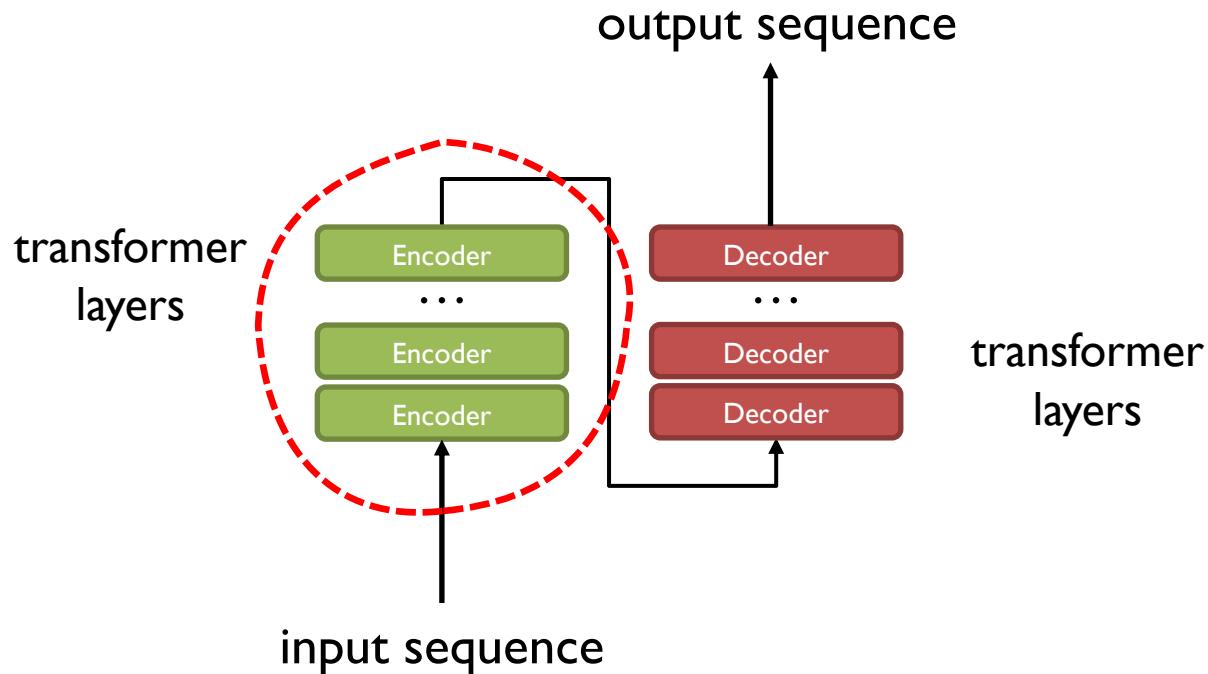
# What's a transformer?

(sequence-to-sequence model)



# What's a transformer?

(sequence-to-sequence model)



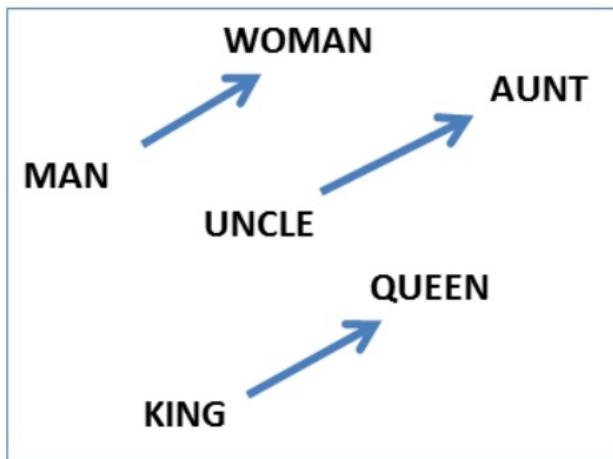


# What's BERT?

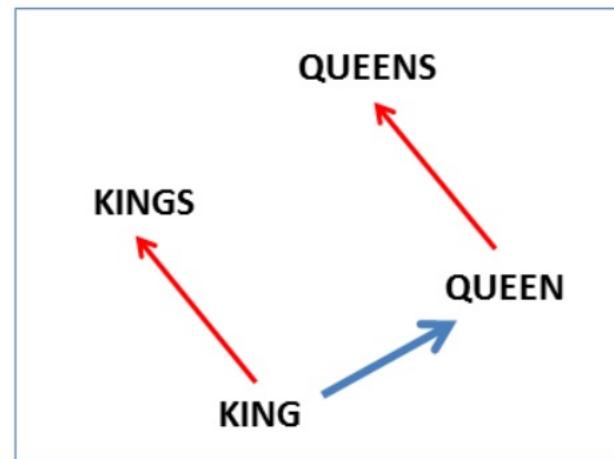
❶ a transformer encoder pretrained  
with a mask language model  
(MLM) objective that **generates**  
**❷ contextualized representations**  
and can be fine-tuned for a wide  
range of NLP tasks.❸

- ❶ What's a transformer?
- ❷ **What are contextualized representations?**
- ❸ What's pretraining?
- ❹ What's fine-tuning?

# Static Representations



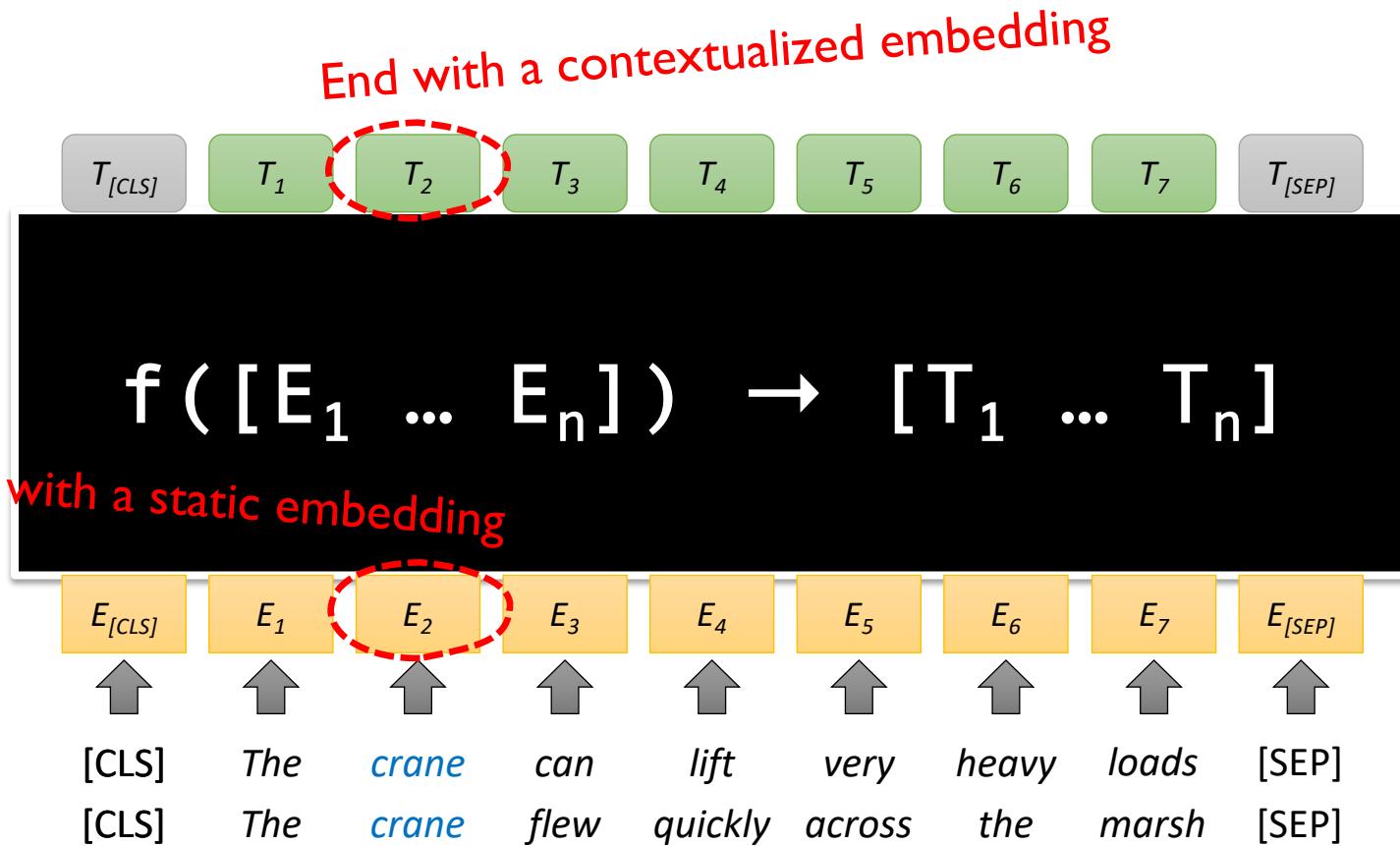
$$\text{king} - \text{man} + \text{woman} = \text{queen}$$



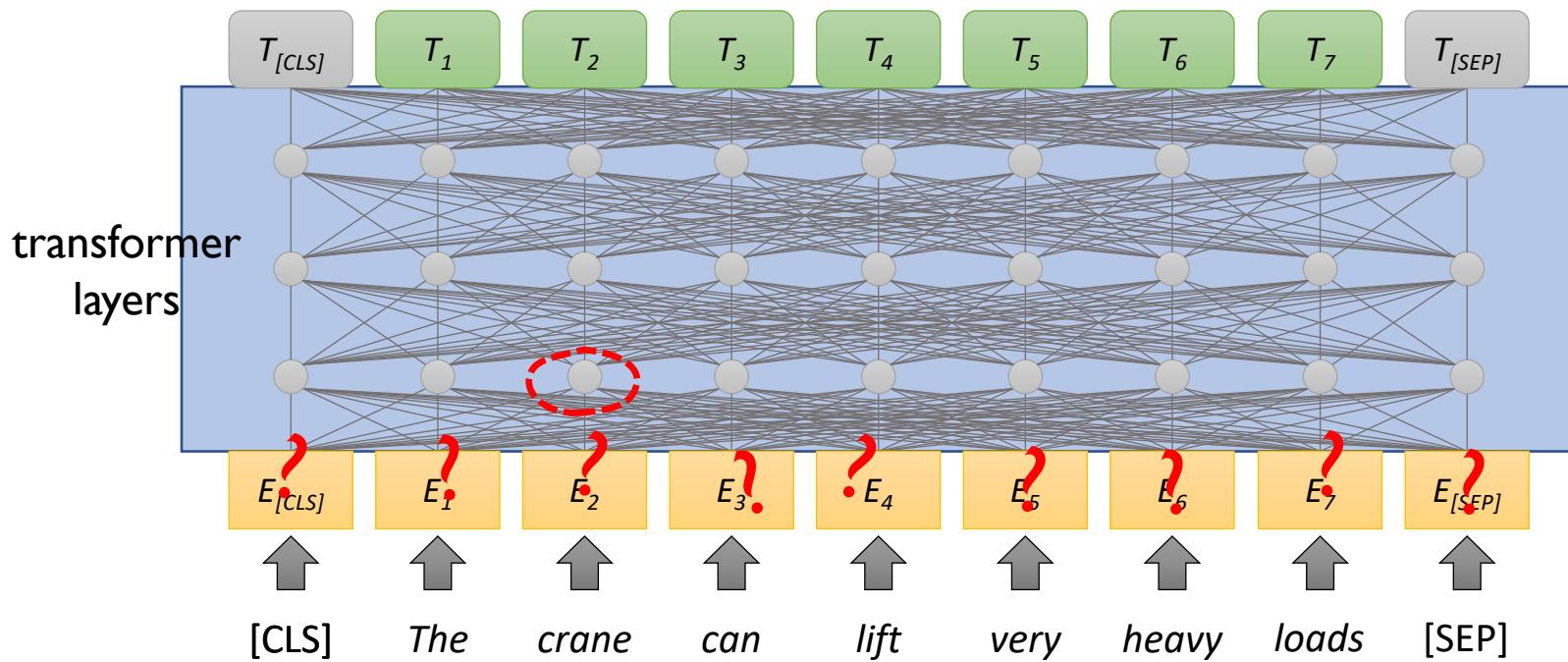
$$\text{kings} - \text{king} = \text{queens} - \text{queen}$$

**tokens → dense vectors  
(embeddings)**

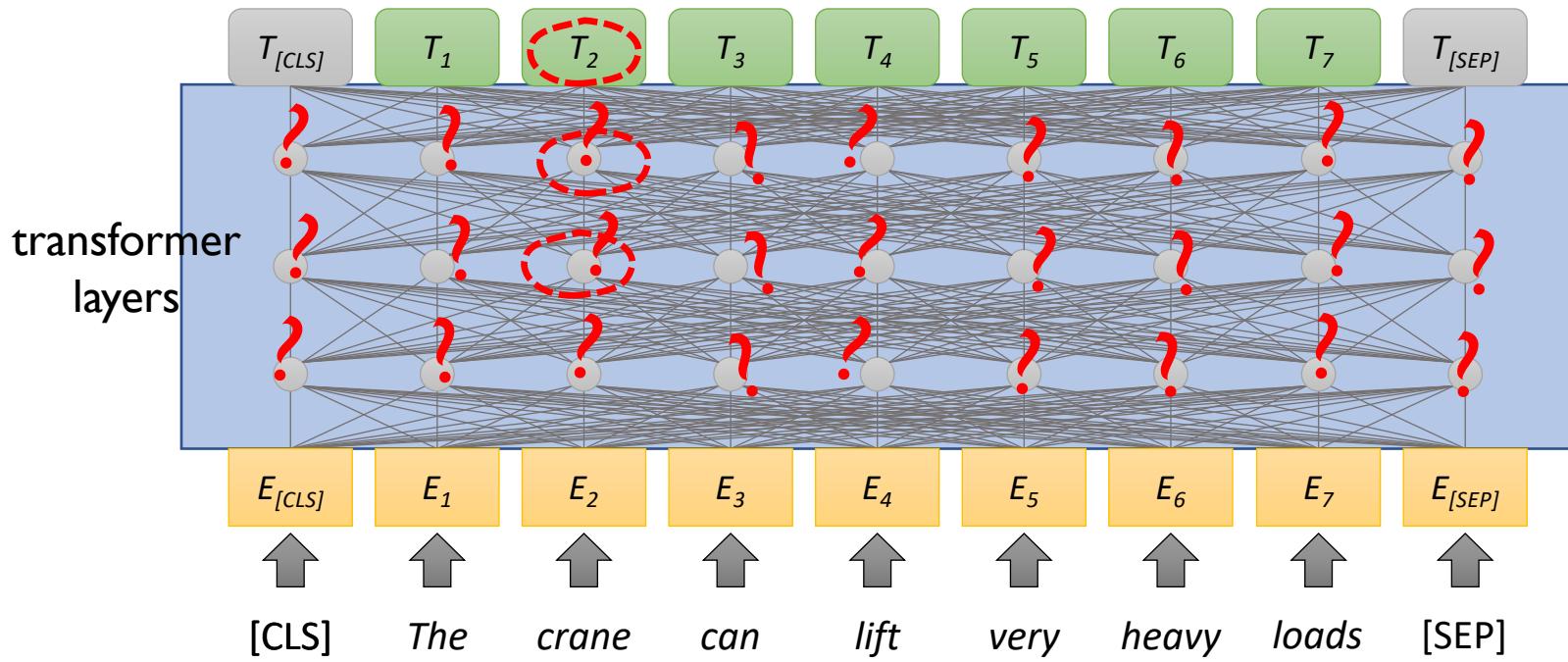
# Contextualized Representations?



# Attention!



# Attention!





# What's BERT?

❶ a transformer encoder ❸ pretrained  
with a mask language model  
**(MLM) objective** that generates ❷  
contextualized representations ❸  
and can be fine-tuned for a wide  
range of NLP tasks. ❹

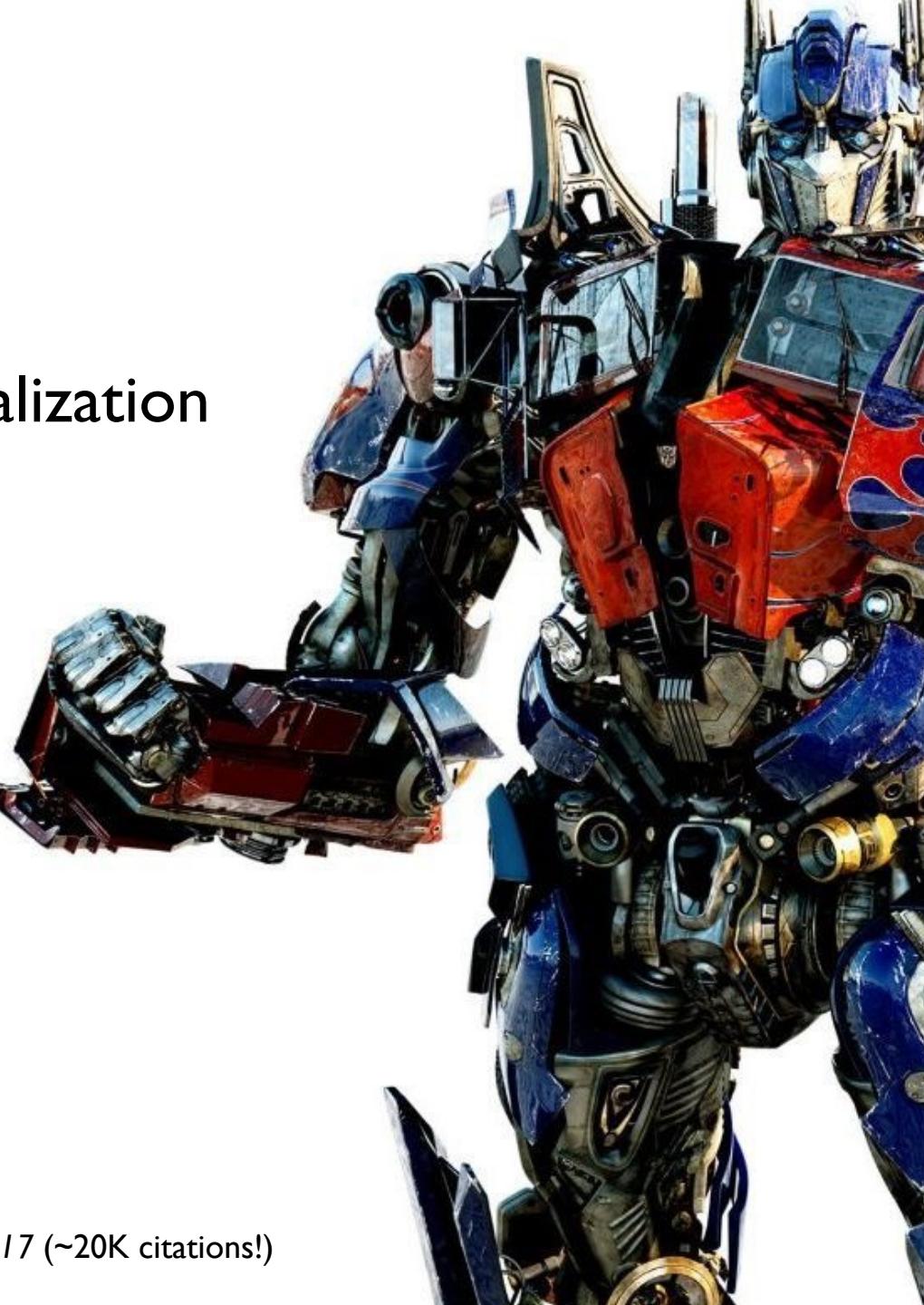
- ❶ What's a transformer?
- ❷ What are contextualized representations?
- ❸ **What's pretraining?**
- ❹ What's fine-tuning?

# *Just do it!*

Start with random initialization

Train with labeled data

# Supervised Training!





Unsupervised  
Training?

You can never have enough labeled data!

# Self-Supervised Training



(language modeling objective)

You don't need labels... just text!

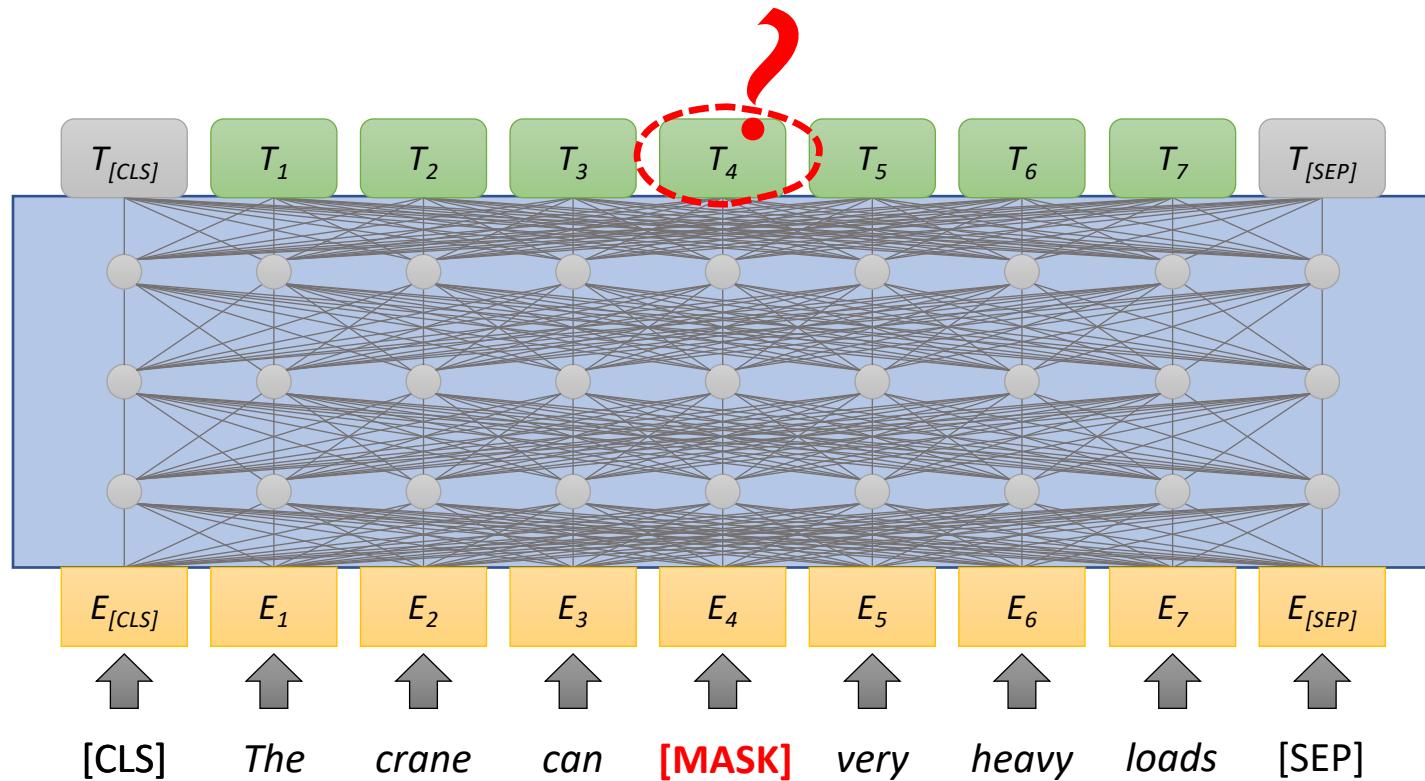
# But wait! I don't actually care!



5+4=9  
Doesn't  
Matter!



# Masked Language Model = MLM





Start with this...



Do this...

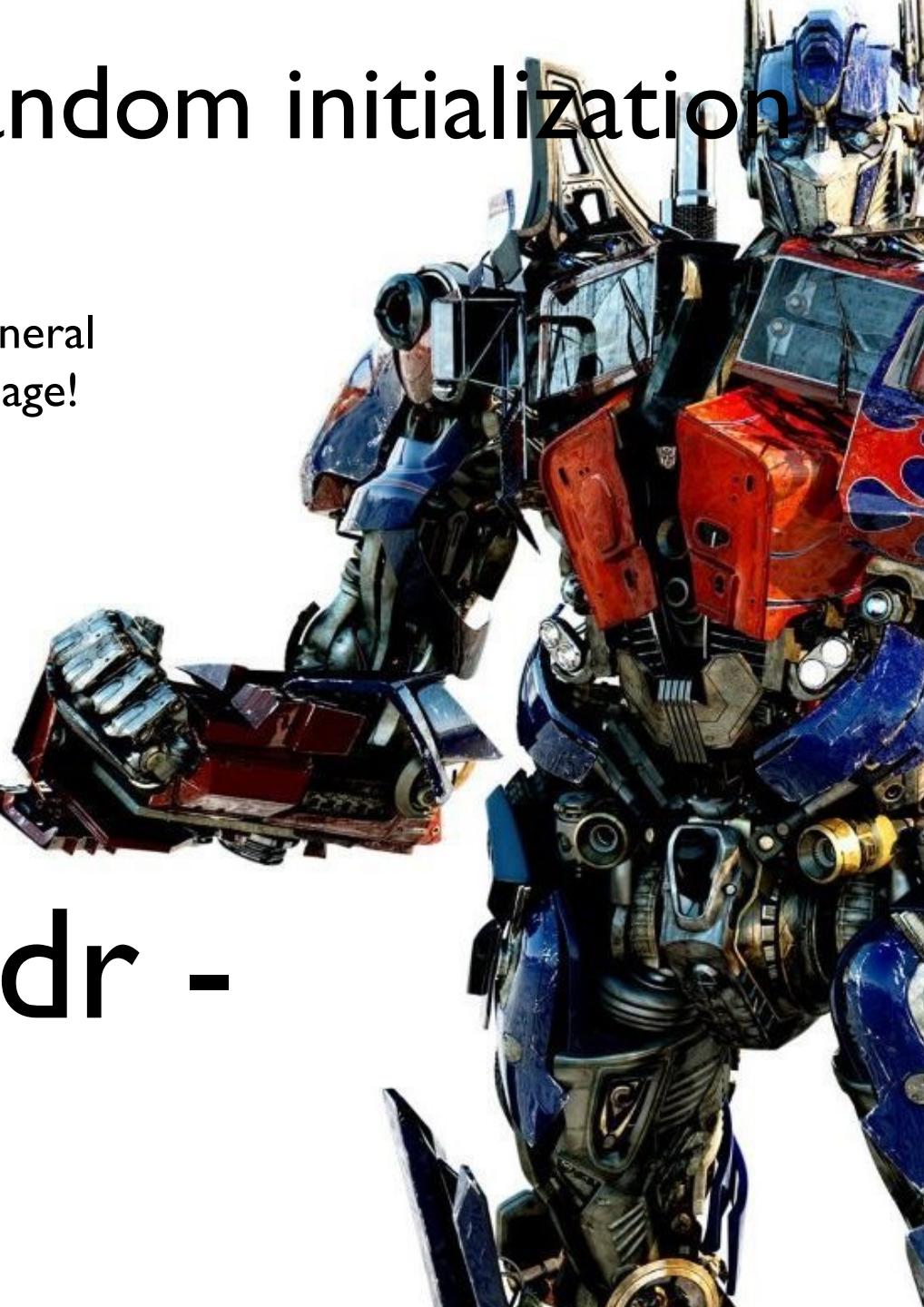
Many many times...  
With lots of GPUs!



# Random initialization

## MLM pretraining

Already knows general things about language!



tl;dr -



# What's BERT?

❶ a transformer encoder pretrained  
with a mask language model  
(MLM) objective that generates  
❷ contextualized representations  
and can be fine-tuned for a wide  
range of NLP tasks.❸

- ❶ What's a transformer?
- ❷ What are contextualized representations?
- ❸ What's pretraining?
- ❹ What's fine-tuning?



# Pretraining

Play peek-a-boo!

Self-Supervised

How do we get BERT to  
something more interesting?

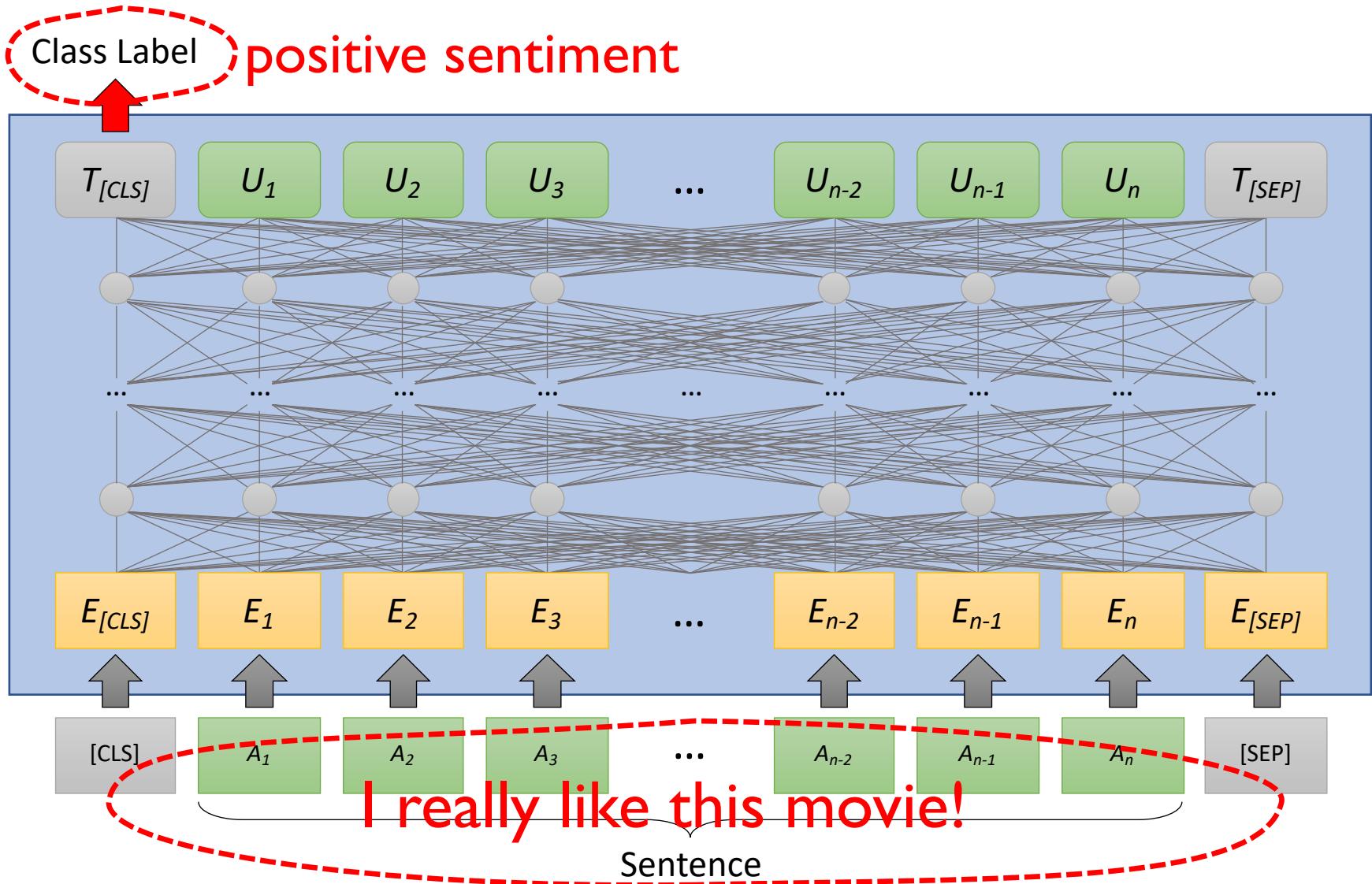
# Fine-Tuning

Adjust model parameters using  
labeled data for the target task

Supervised

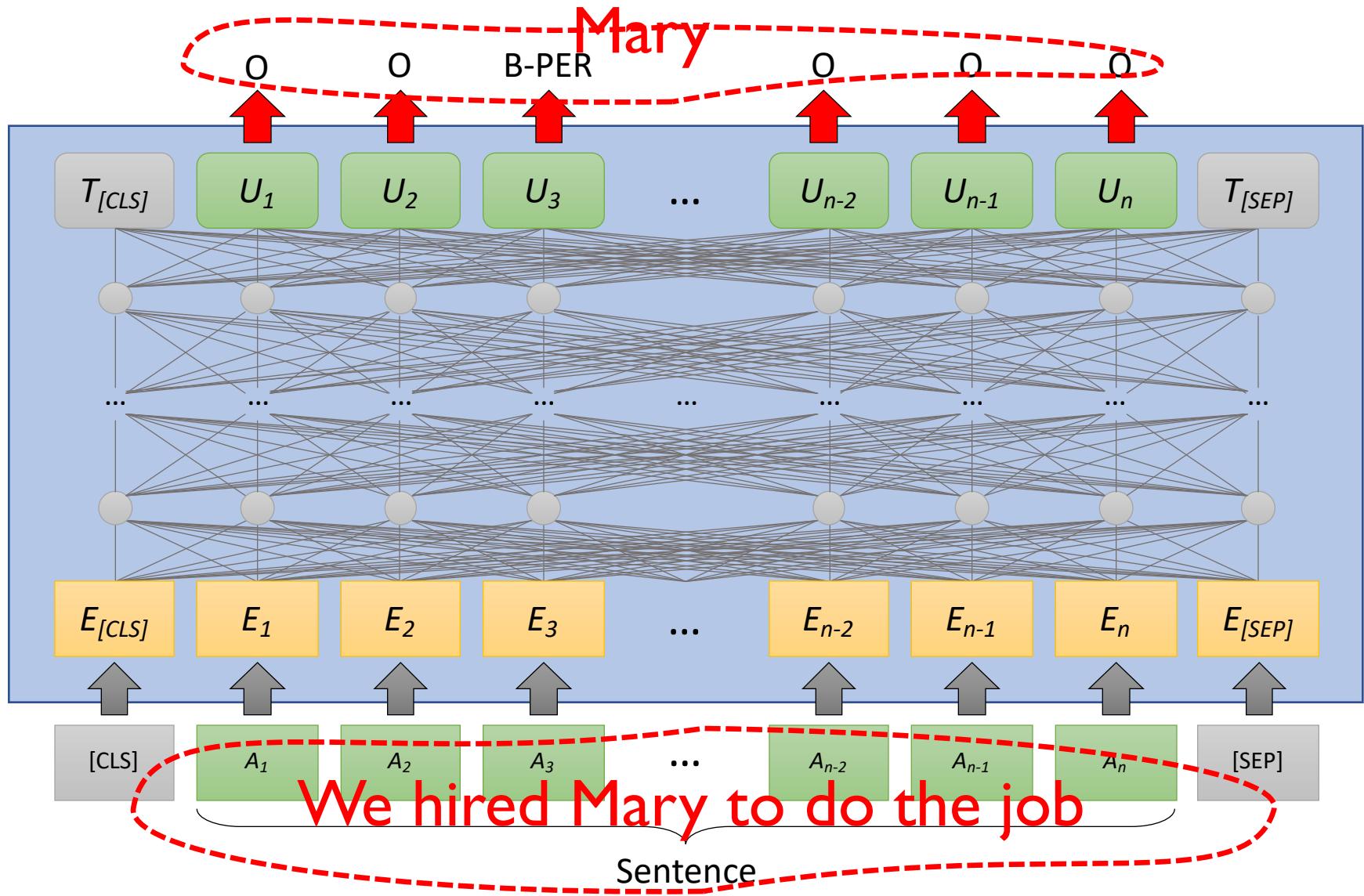
$$f([A_1, A_2, \dots, A_n]) \rightarrow y$$

Single-Input Classification (e.g., sentiment analysis)



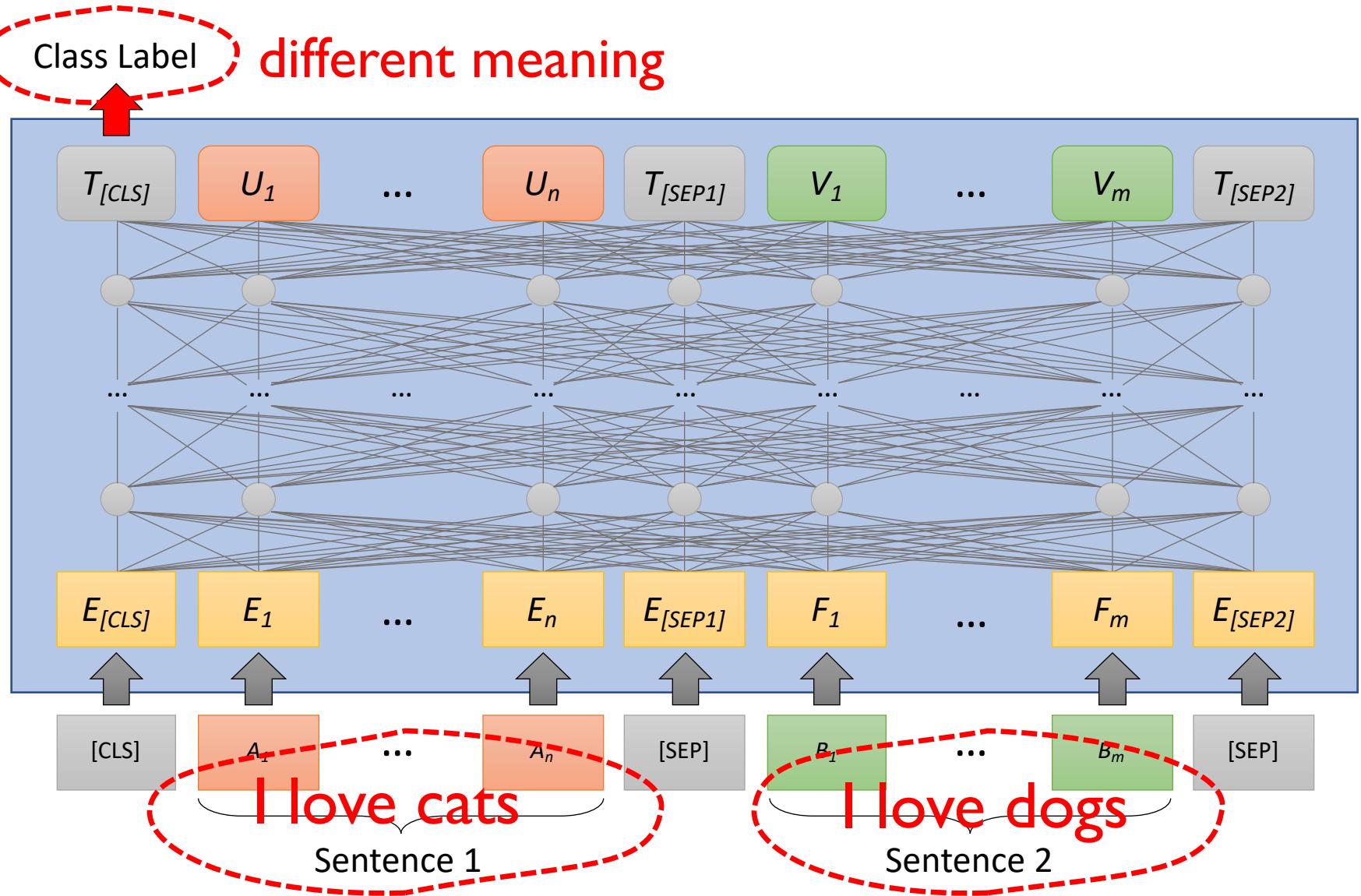
$$f([A_1, A_2, \dots A_n]) \rightarrow [y_1 \dots y_n]$$

Single-Input Sequence Labeling (e.g., named-entity recognition)



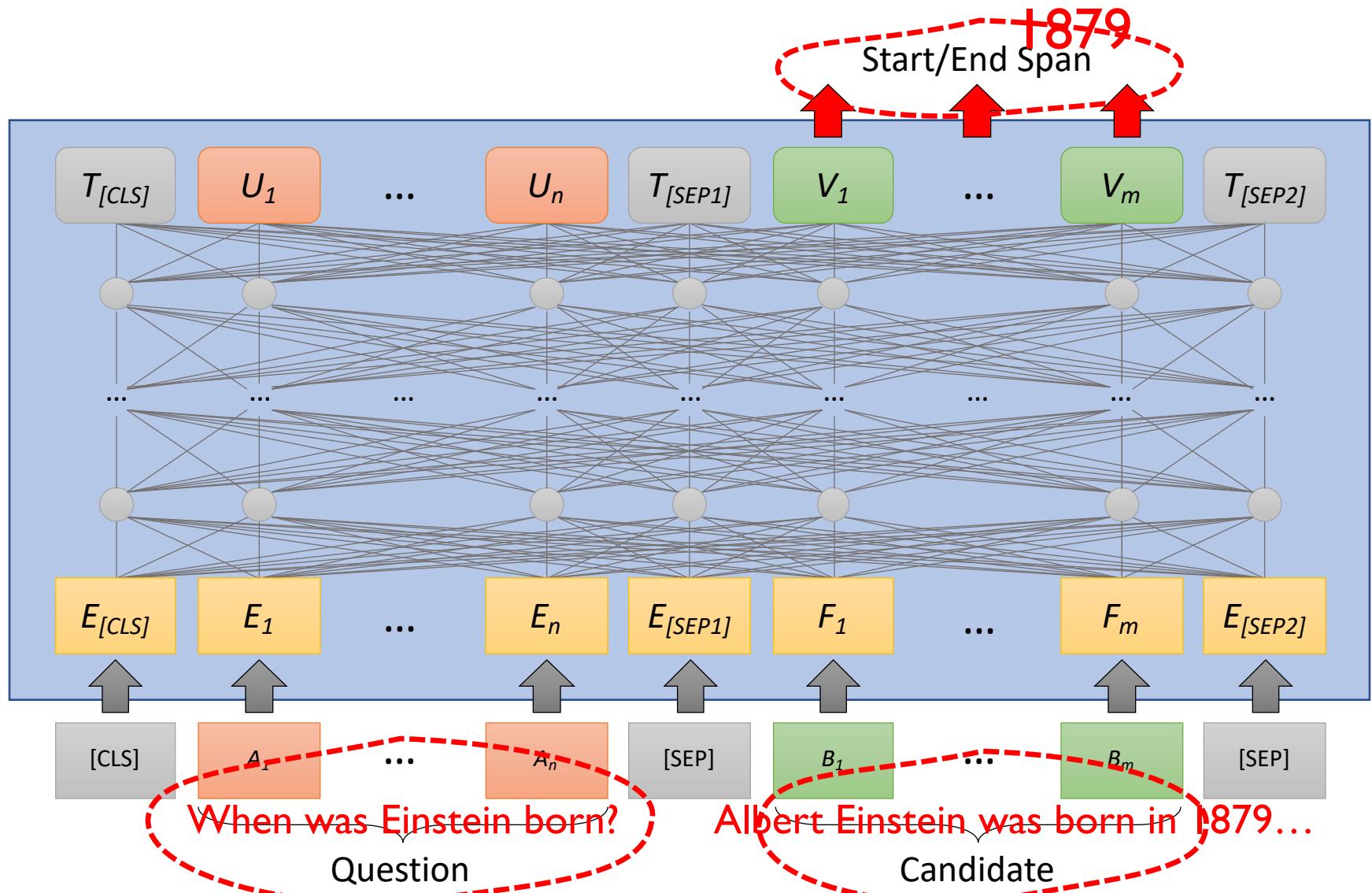
$$f([A_1 \dots A_n], [B_1 \dots B_m]) \rightarrow y$$

Two-Input Classification (e.g., paraphrase detection)



$$f([A_1 \dots A_n], [B_1 \dots B_m]) \rightarrow [y_1 \dots y_m]$$

Two-Input Sequence Labeling (e.g., question answering)





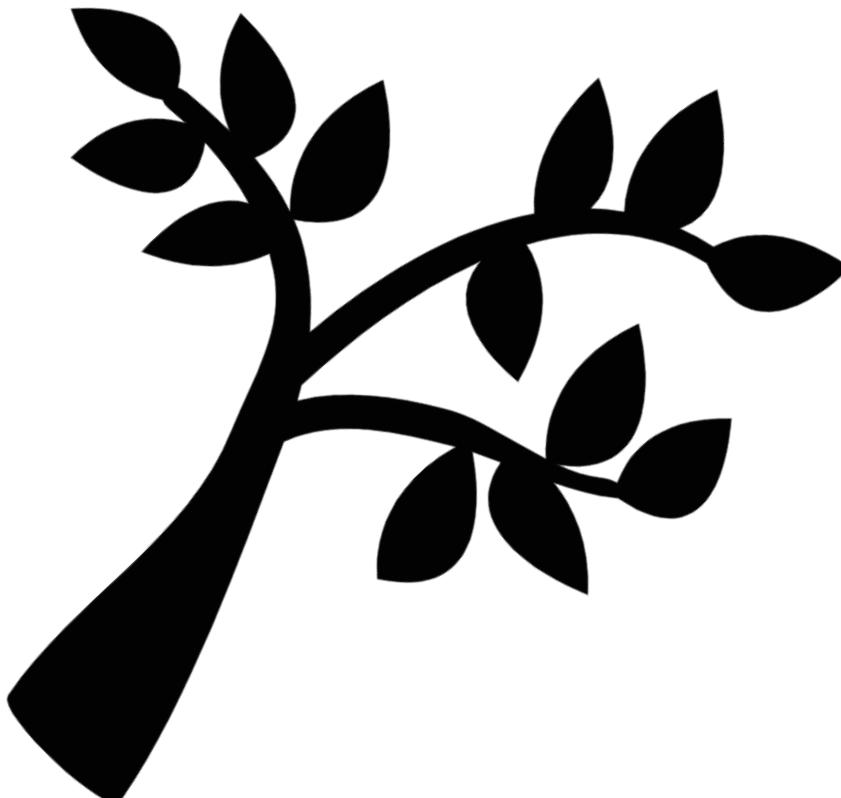
# What's BERT?

a transformer encoder pretrained with a mask language model (MLM) objective that generates contextualized representations and can be fine-tuned for a wide range of NLP tasks.

Now you know what  
BERT is!



# BERT's family tree



Devlin et al. BERT. NAACL 2019 (~18K citations!)

# What's a muppet?

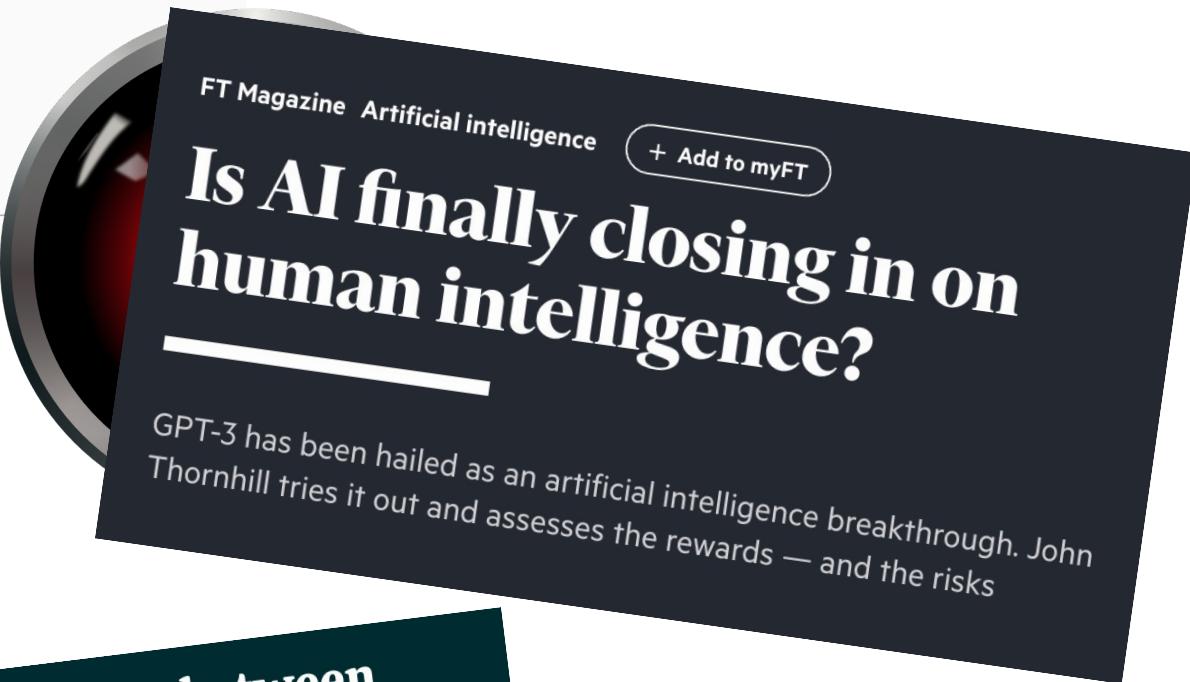


# What's GPT-3?

What Is GPT-3 And Why Is It  
Revolutionizing Artificial  
Intelligence?



Bernard Marr Contributor   
Enterprise Tech

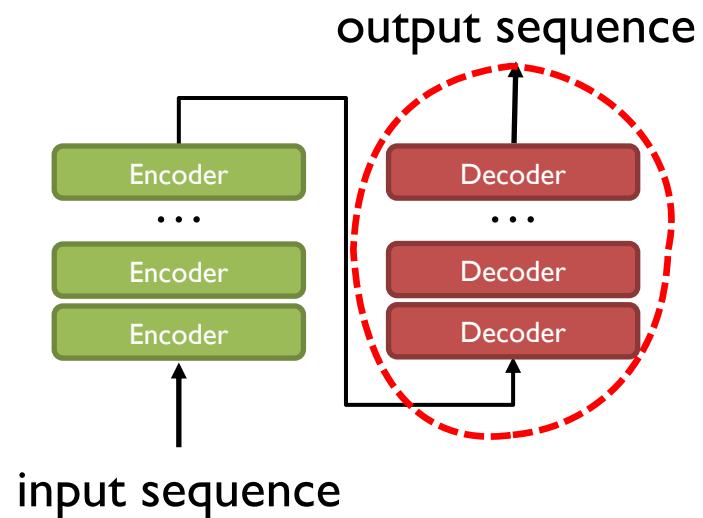


The image shows a dark-colored digital magazine cover from FT Magazine. At the top left, it says "FT Magazine Artificial intelligence". To the right is a button with "+ Add to myFT". The main title, "Is AI finally closing in on human intelligence?", is displayed prominently in large, bold, white letters. Below the title is a subtitle: "GPT-3 has been hailed as an artificial intelligence breakthrough. John Thornhill tries it out and assesses the rewards — and the risks". The overall design is modern and professional.

We're entering the AI twilight zone between  
narrow and general AI

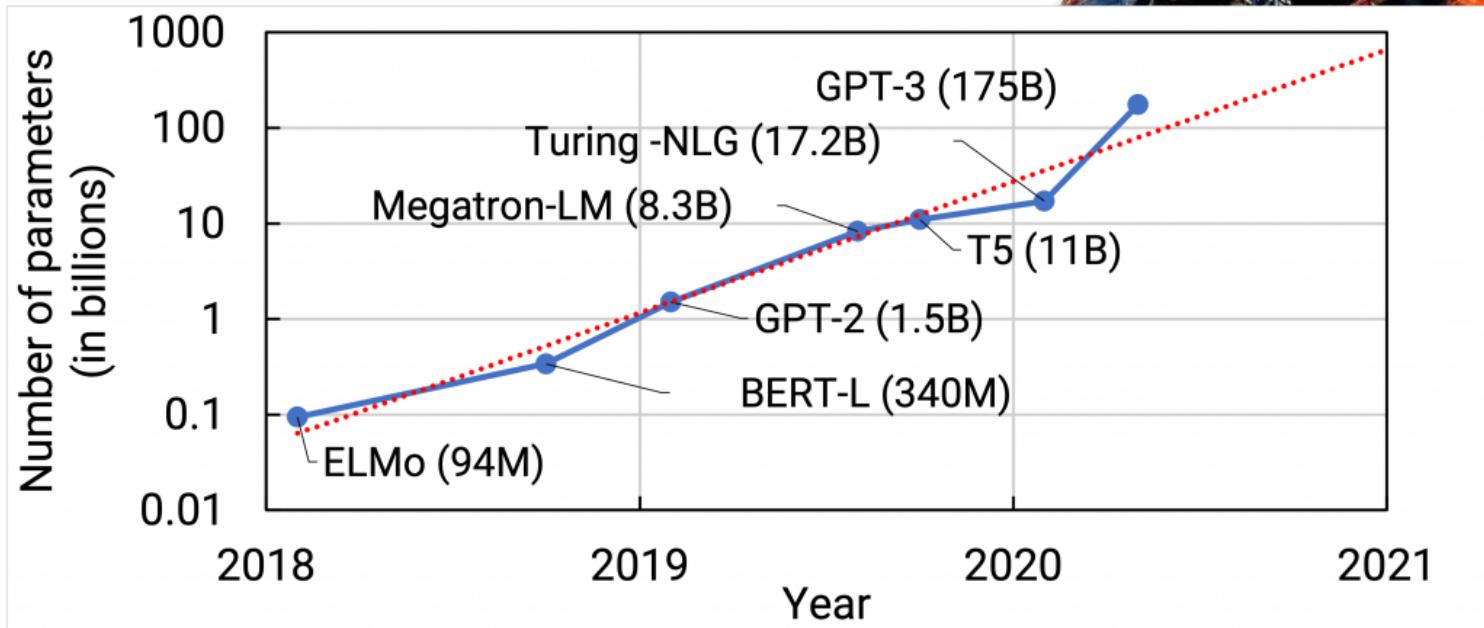
f

# What's GPT-3?



# What's GPT-3?

## pretrained decoder

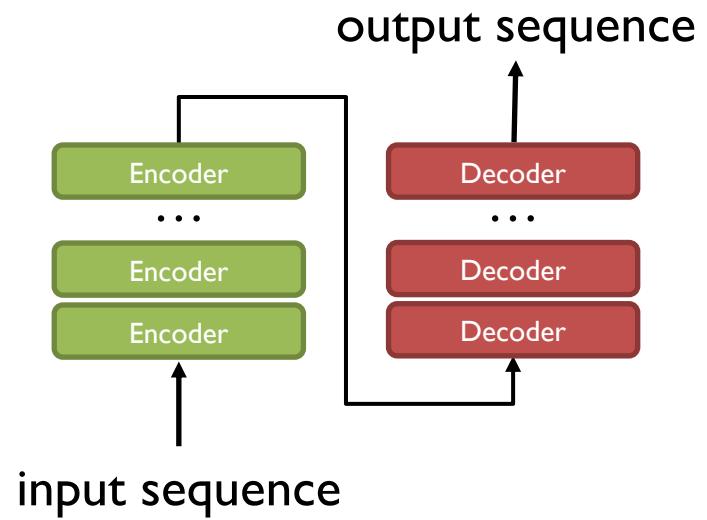


# 175 BILLION parameters!

*Before:* training a GPT-3 model ... would take 36 years on eight V100 GPUs, or seven months with 512 V100 GPUs.

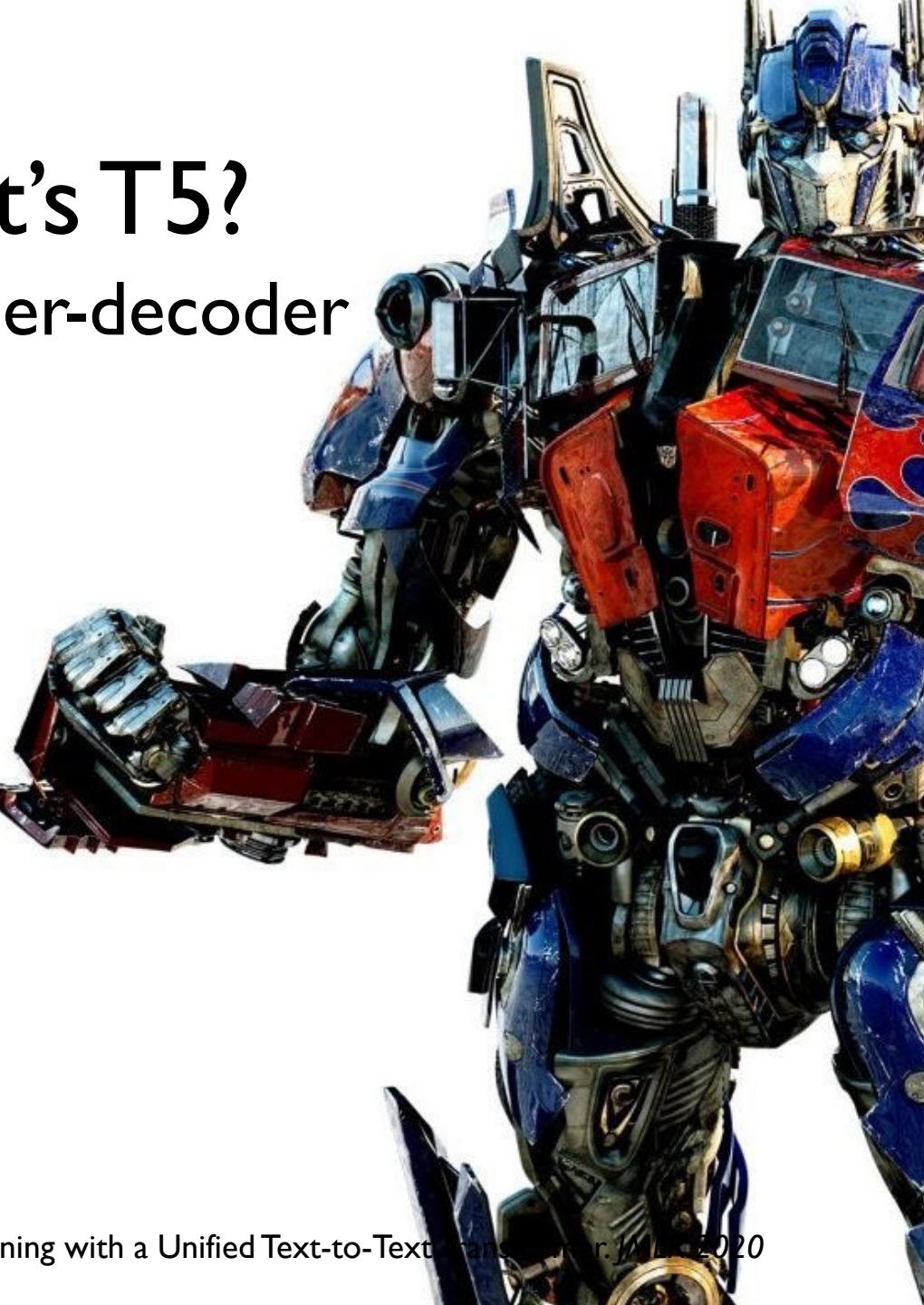
*Now:* The GPT-3 model ... requires just over a month to train using 1024 A100 GPUs!

# What's T5?



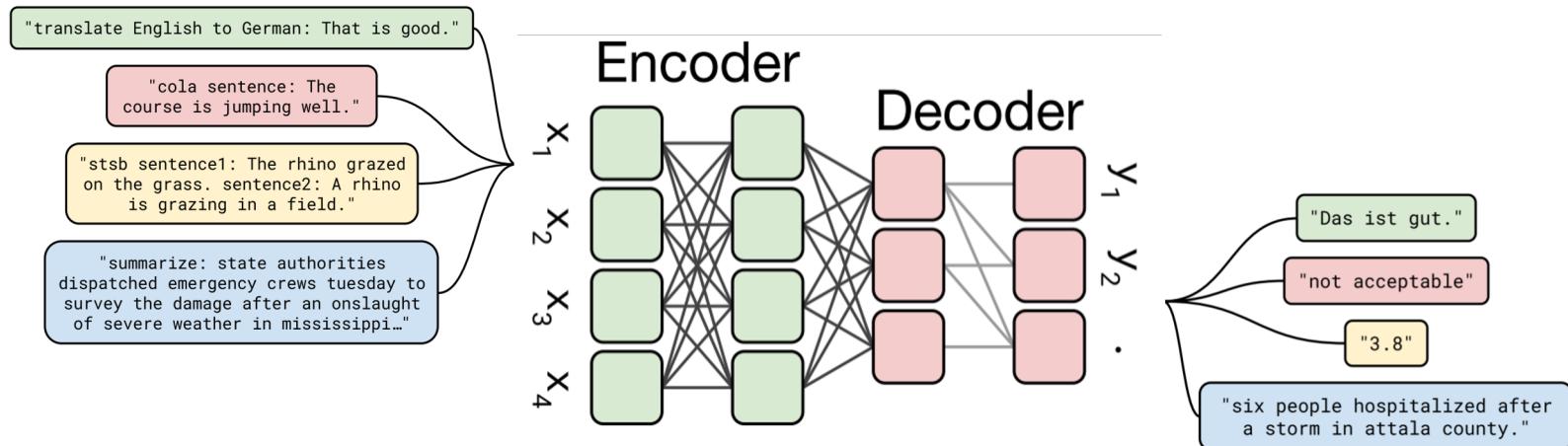
# What's T5?

## pretrained encoder-decoder

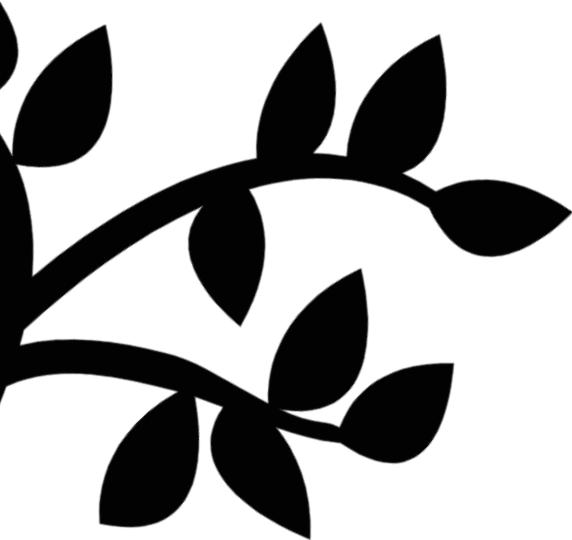


# What's T5?

## pretrained encoder-decoder



Everything's a sequence-to-sequence task!



**2017 – “vanilla” transformers**  
(no pretraining)

**Mid-2018 – GPT-2**  
(decoder, pretrained w/ “predict next word”)

**Late-2018 – BERT**  
(encoder, pretrained w/ “peek-a-boo”)

**Late-2019 – T5**  
(encoder–decoder, pretrained)

**Mid-2020 – GPT-3**  
(like GPT-2, but larger!)

# What is BERT? (and his muppet friends)

Why should you care?

# BERT is good at text ranking

*what is a lobster roll?*

A Lobster Roll is a bread roll filled with bite-sized chunks of lobster meat. Lobster Rolls are made on the Atlantic coast of North America, from the New England area of the United States on up into the Maritimes areas of Canada.



# BERT is good at text ranking



Select some  
promising texts

Rerank  
selected texts



 *lucene*

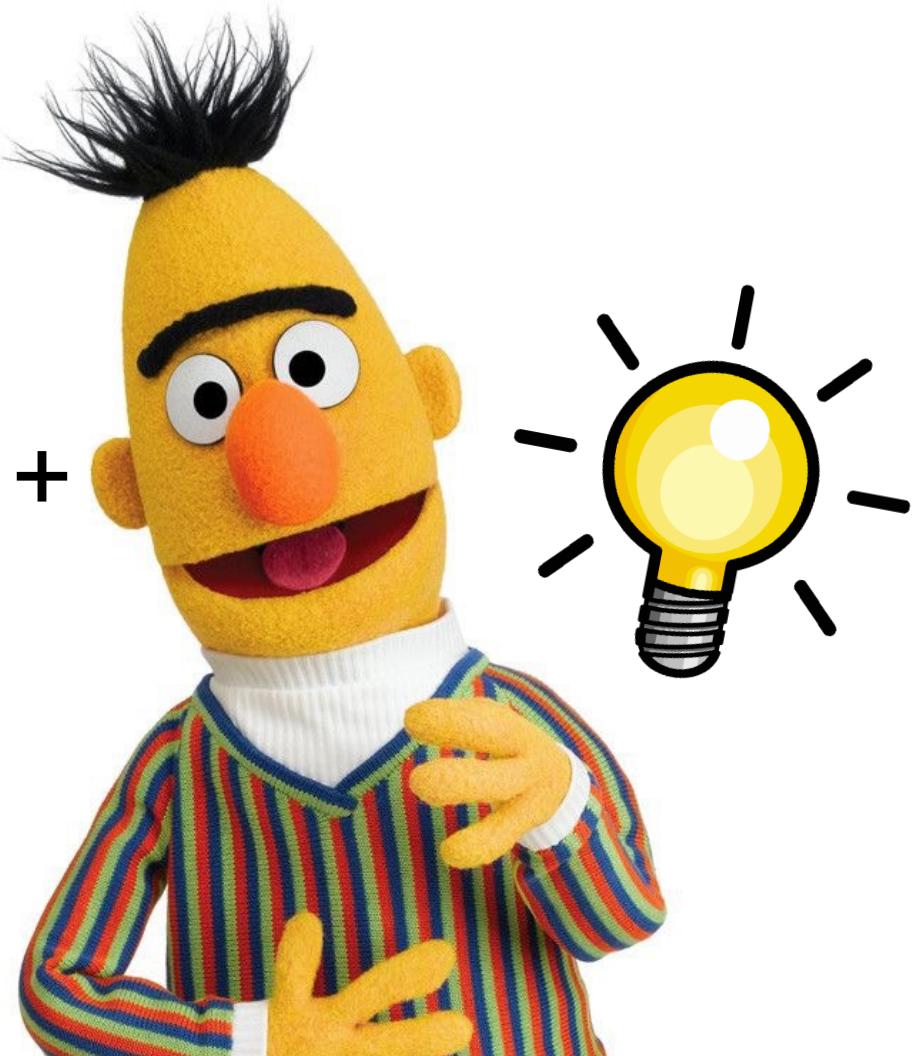
The Lucene logo, featuring the word "lucene" in a green, italicized, lowercase font. A stylized green "L" with three curved lines extending from its top right corner is positioned to the left of the text.

# BERT is good at text ranking



Select some  
promising texts

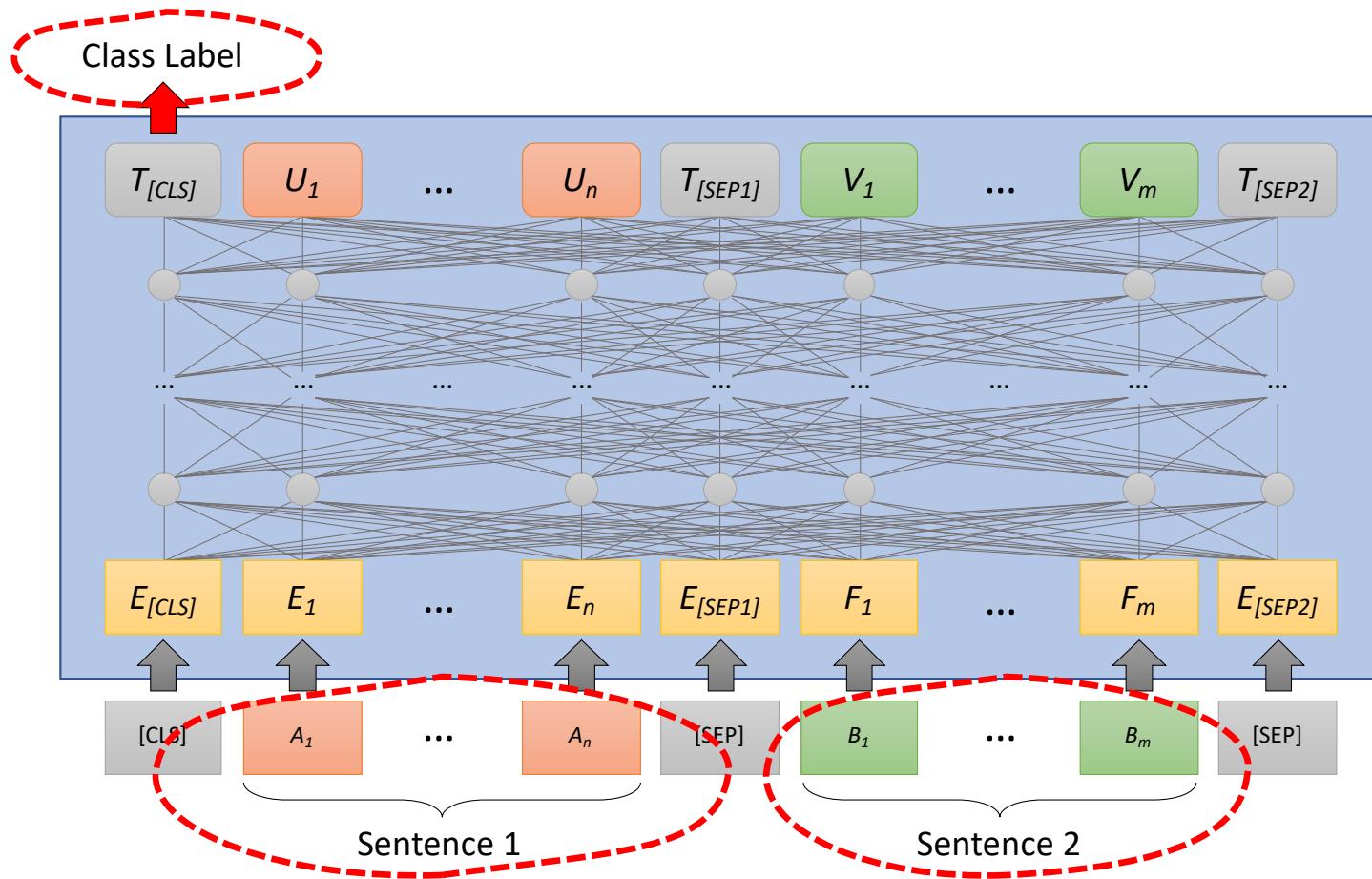
The Lucene logo, featuring the word "Lucene" in a stylized green font where the letters are interconnected by horizontal lines, with a small green swoosh icon preceding the "L".



# BERT is good at text ranking

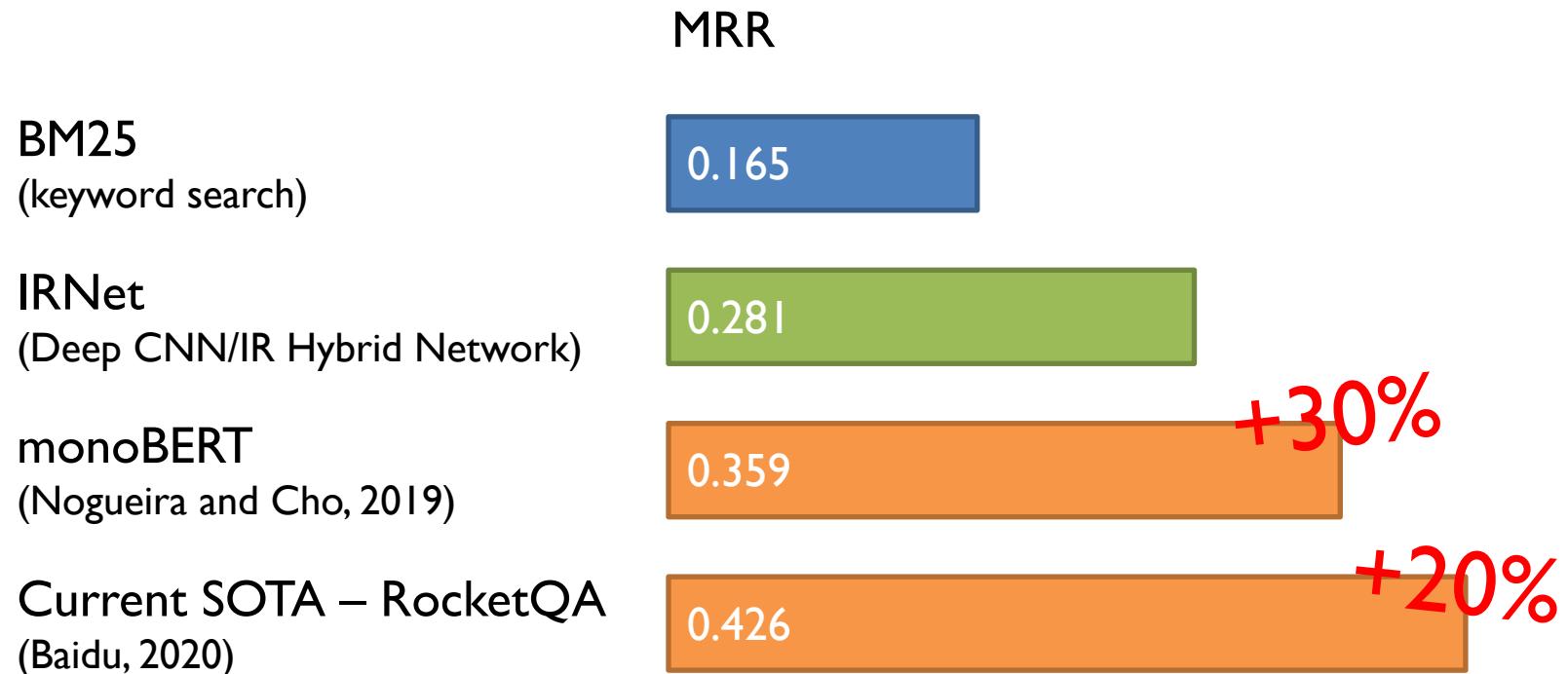
$$f([A_1 \dots A_n], [B_1 \dots B_m]) \rightarrow y$$

query                    candidate text    prob. relevance



*Like how good?*

# BERT is good at text ranking



MS MARCO passage ranking  
(8.8M passages, queries from Bing logs)

# BERT is good at QA

*who lives in the imperial palace in tokyo?*

The Tokyo Imperial Palace ( 皇居 , Kōkyō , literally " Imperial Residence " ) is the primary residence of the Emperor of Japan . It is a large park - like area located in the Chiyoda ward of Tokyo and contains buildings including the main palace ( 宮殿 , Kyūden ) , the private residences of the Imperial Family , an archive , museums and administrative offices.

# BERT is good at QA



Select some  
promising texts

Extract correct  
answer span



# BERT is good at QA



Select some  
promising texts

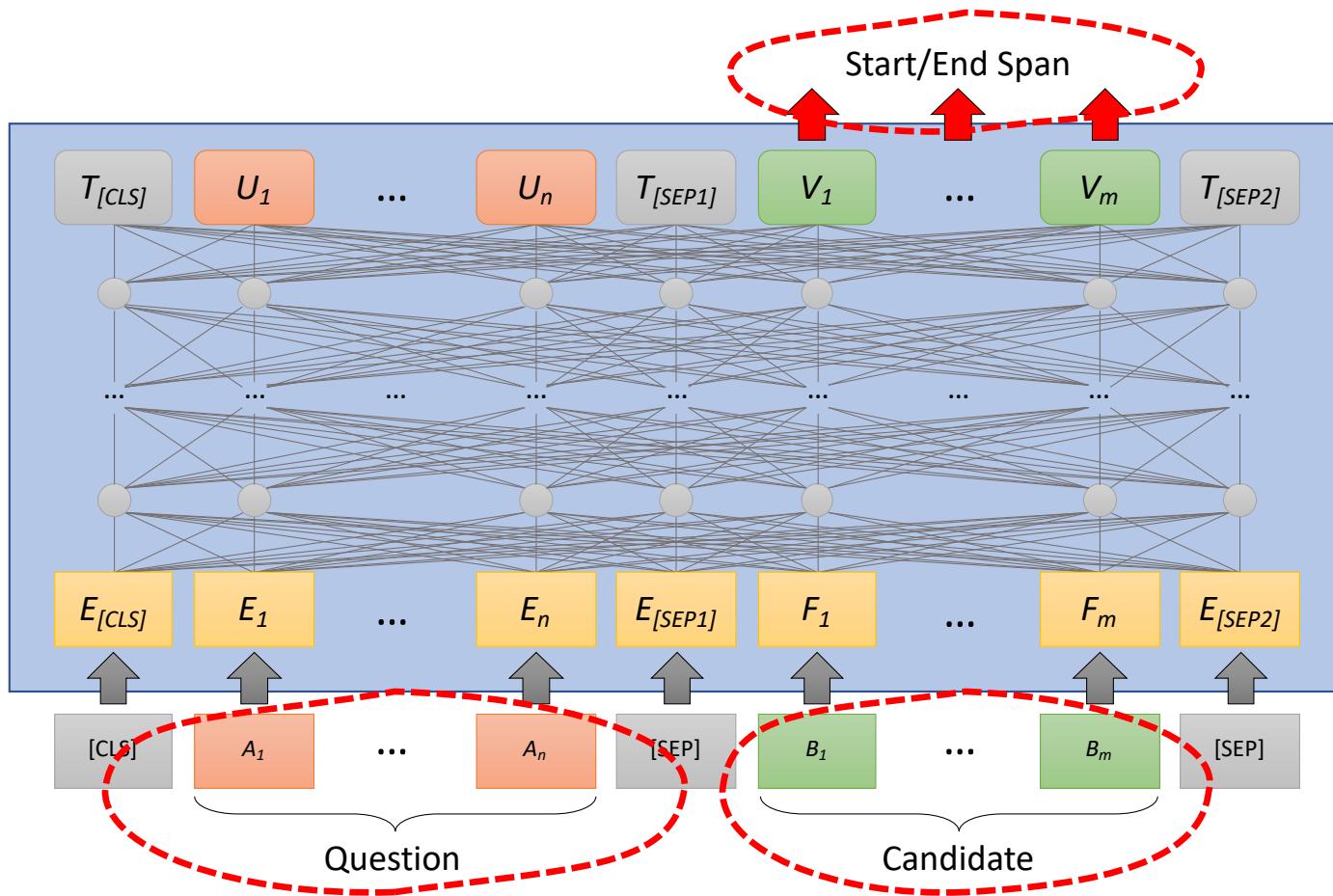
The Lucene logo, featuring the word "Lucene" in a green, stylized font where the letters are interconnected by horizontal lines, with a small green swoosh icon preceding the "L".



# BERT is good at QA

$$f([A_1 \dots A_n], [B_1 \dots B_m]) \rightarrow [y_1 \dots y_m]$$

question                    candidate answer                    prob. start/end pos



# BERT is good at QA

*Like how good?*

Exact Match score

Dr.QA  
(Chen et al. 2017)

0.271

BERTserini  
(Yang et al. 2019)

0.386

Current SOTA  
(circa 2021)

~0.57

SQuAD questions against ~5M Wikipedia articles



# What's QA? select-project queries!

*who lives in the imperial palace in tokyo?*

```
select "residents" from location  
      where city = "tokyo" and  
            name = "imperial palace"
```

# BERT can do multi-hop QA

*Which former member of the Pittsburgh Pirates was nicknamed “The Cobra”?*

Paragraph 1: Several current and former members of the Pittsburgh Pirates – ... John Milner, [Dave Parker](#), and Rod Scurry...

Paragraph 2: [David Gene Parker](#), nicknamed “The Cobra”, is an American former player in Major League Baseball...



# What's multi-hop QA? select-project-join queries!

*Which former member of the Pittsburgh Pirates was  
nicknamed “The Cobra”?*

```
select roster.name from roster, player  
  where roster.name = player.name and  
roster.team = "Pittsburgh Pirates" and  
player.nickname = "The Cobra"
```

# BERT is good at QA

“Open-Book” Question Answering



Select some  
promising texts

 *lucene*

The Lucene logo, which consists of a stylized green 'L' shape followed by the word 'lucene' in a lowercase, italicized, green sans-serif font.

# BERT is good at QA

“Closed-Book” Question Answering



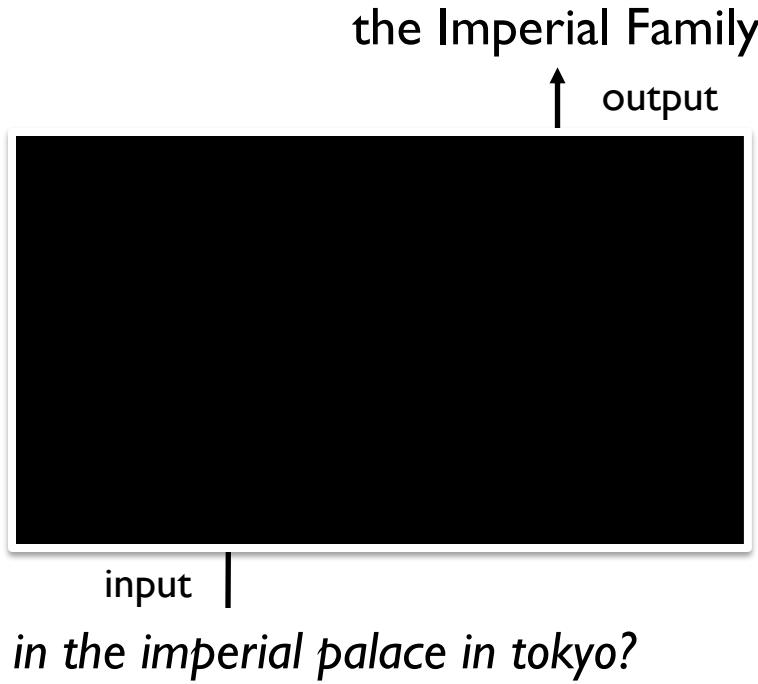
# T5 is good at QA

“Closed-Book” Question Answering

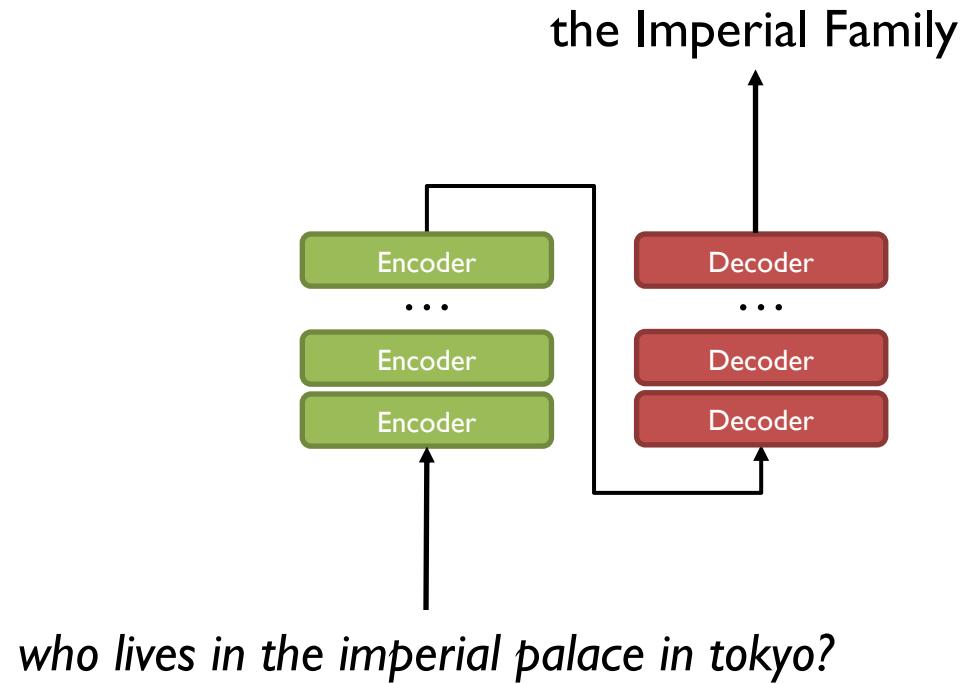


No more corpus!  
(Just the model)

# Closed-Book QA with T5



# Closed-Book QA with T5



No more corpus!  
(Just a transformer)

A close-up photograph of a young boy with a shaved head, looking directly at the camera with a neutral expression. A large, dark, semi-transparent shadow of a hand is cast over the left side of his face, obscuring his eye and nose. The background is dark and out of focus.

There is no spoon.

(Just a transformer)

# There is no corpus

Pretraining  
Fine-Tuning



A close-up photograph of a young boy with a shaved head, looking directly at the camera with a neutral expression. He is wearing a light-colored t-shirt and holding a dark book or magazine open with both hands, showing the pages. The background is slightly blurred, showing what appears to be a shelf with more books.

There is no fine-tuning.  
(Just a pretrained transformer)

# What does BERT know?







# What does BERT know?

Wikidata P19: place of birth

Francesco Bartolomeo  
Conti was born in \_\_\_\_\_

Accuracy: 16%

select birthplace from person  
where name = "Francesco Bartolomeo"

Accuracy: 100%

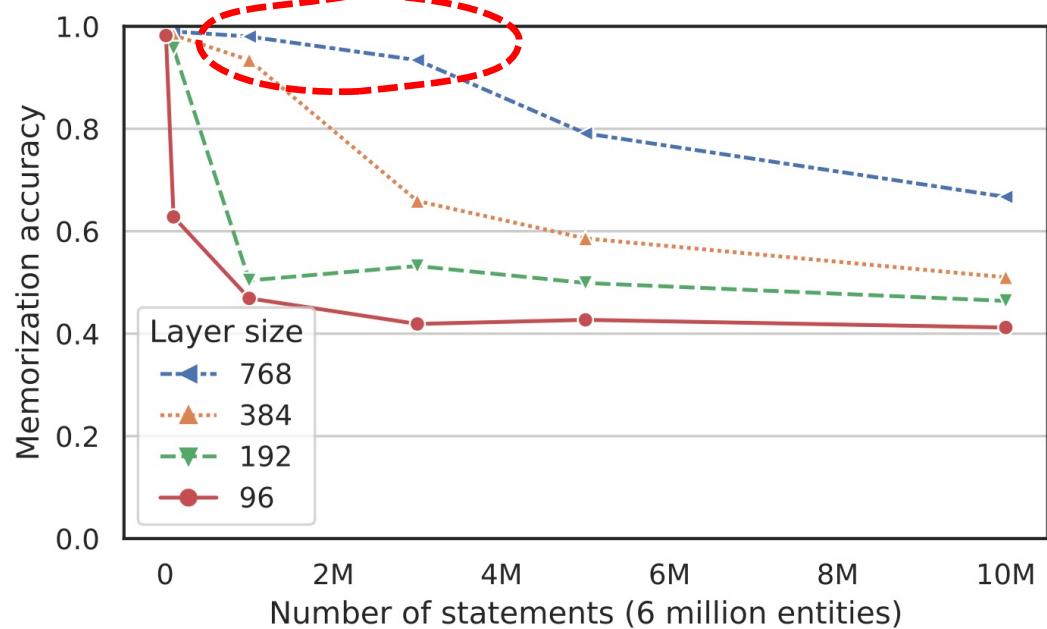


# Let's teach BERT!





# Can BERT memorize?





What's fact memorization?  
Inserts!

What's QA?  
select-project queries!

What's multi-hop QA?  
select-project-join queries!

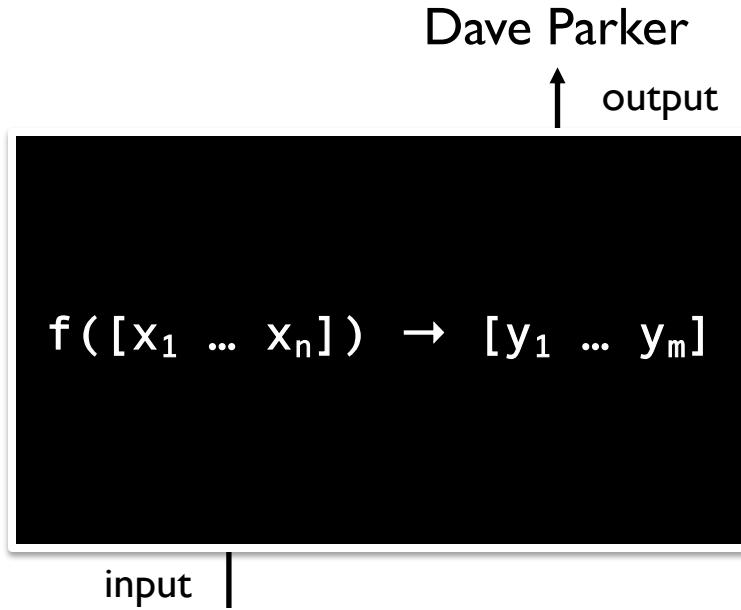
What's fact memorization?  
Inserts!

*and we've gotten rid of the data!*

# Implications?

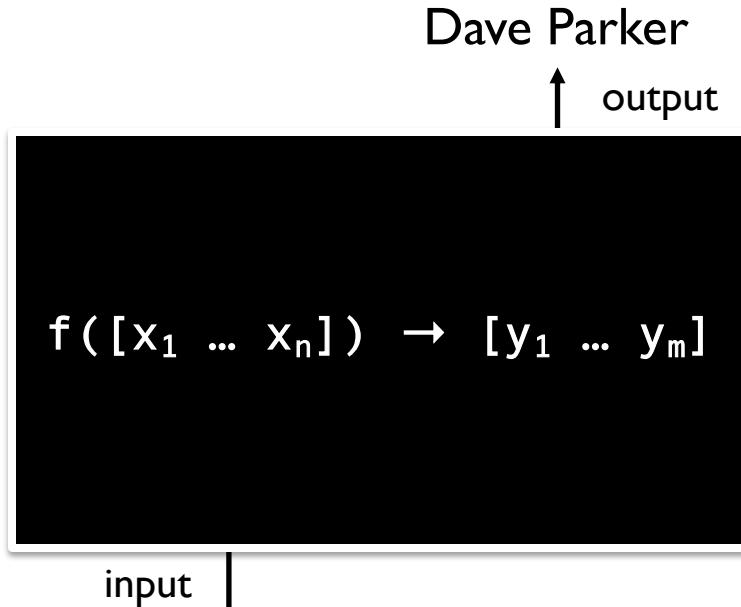
# What's an RDBMS?

(sequence-to-sequence model)



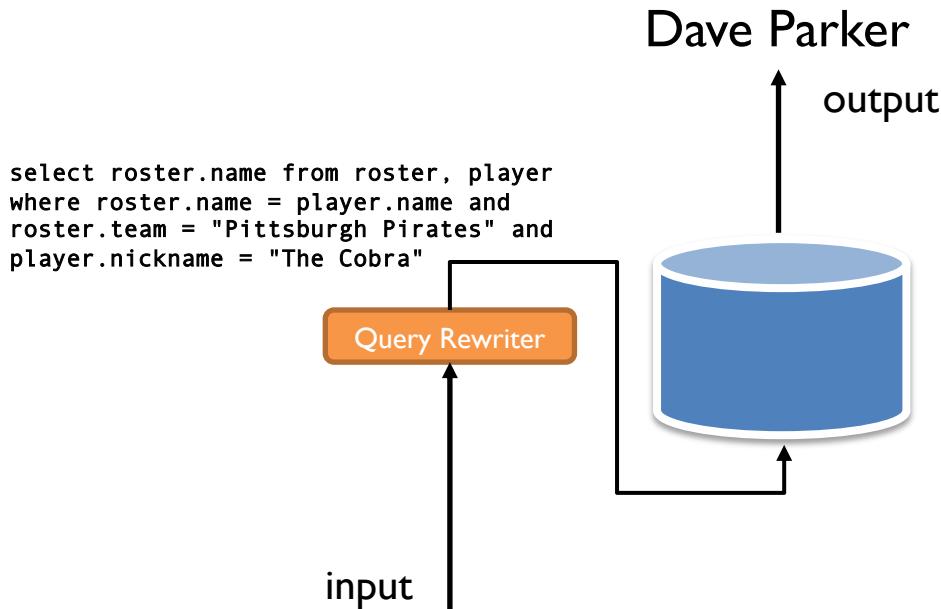
```
select roster.name from roster, player  
where roster.name = player.name and  
roster.team = "Pittsburgh Pirates" and  
player.nickname = "The Cobra"
```

# What's a QA system? (sequence-to-sequence model)



Which former member of the Pittsburgh  
Pirates was nicknamed “The Cobra”?

# What's a QA system? (sequence-to-sequence model)



Which former member of the Pittsburgh  
Pirates was nicknamed “The Cobra”?

# Already being done!

Which countries in Europe have at least 3 car manufacturers?

```
SELECT T1.country_name
FROM countries AS T1 JOIN continents
AS T2 ON T1.continent = T2.cont_id
JOIN car_makers AS T3 ON
T1.country_id = T3.country
WHERE T2.continent = 'Europe'
GROUP BY T1.country_name
HAVING COUNT(*) >= 3
```

# Spider 1.0



Yale Semantic Parsing and Text-to-SQL Challenge

# Already being done!

What does 175 billion parameters get you?  
GPT-3 can write code based on natural language descriptions!

Sharif Shameem @sharifshameem · Jul 13, 2020

This is mind blowing.

With GPT-3, I built a layout generator where you just describe any layout you want, and it generates the JSX code for you.

W H A T

Describe a layout.

Just describe any layout you want, and it'll try to render below!

a button for every color of the rainbow

Generate

<div style={{backgroundColor: 'red', padding: 20}>Red</div><div style={{backgroundColor: 'orange', padding: 20}>Orange</div><div style={{backgroundColor: 'yellow', padding: 20}>Yellow</div><div style={{backgroundColor: 'green', padding: 20}>Green</div><div style={{backgroundColor: 'blue', padding: 20}>Blue</div><div style={{backgroundColor: 'purple', padding: 20}>Purple</div><div style={{backgroundColor: 'violet', padding: 20}>Violet</div>

▶ 1.8M views 0:36 / 2:00

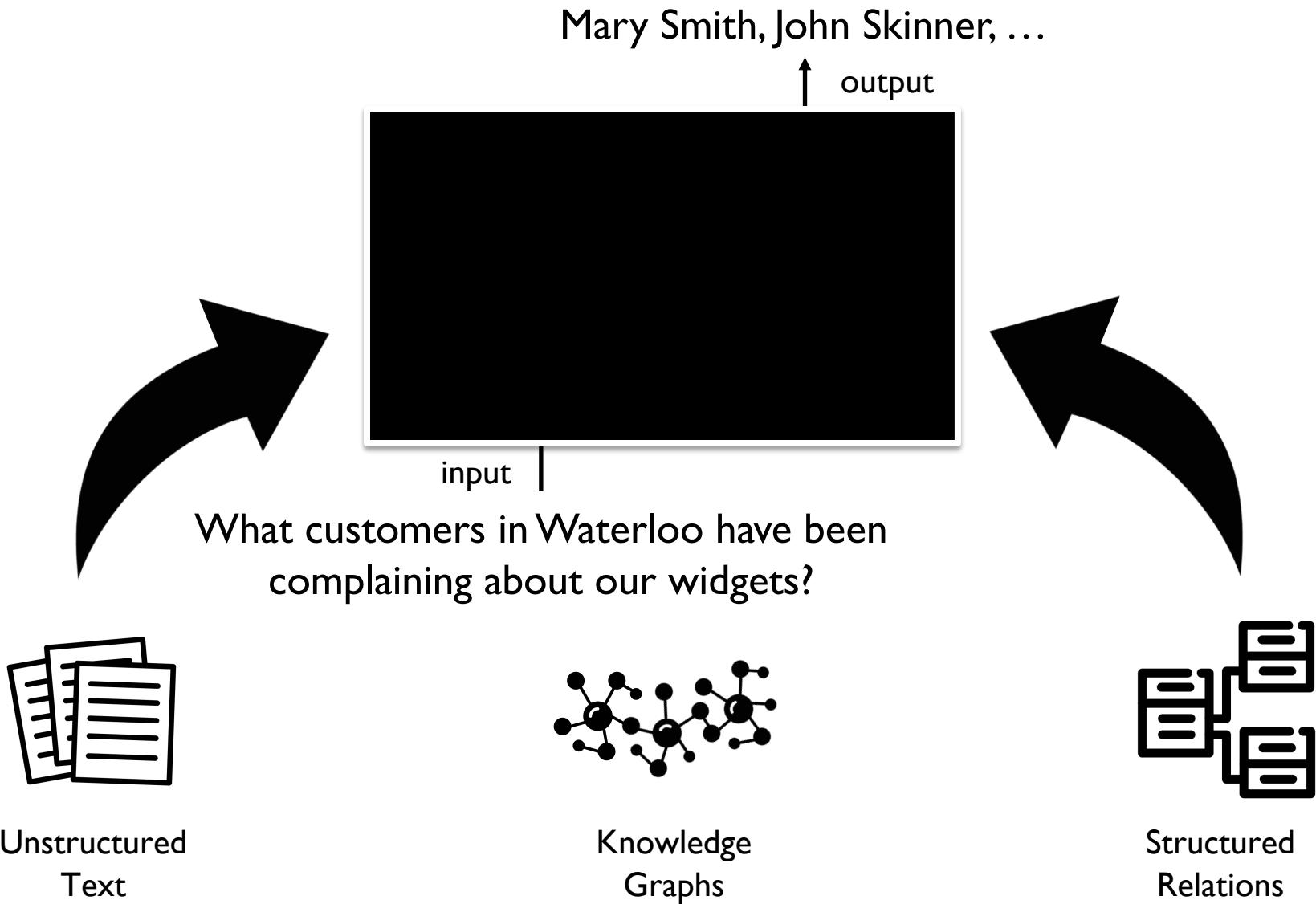
693 14.3K 42.2K

<https://twitter.com/sharifshameem/status/1282676454690451457>

NLP

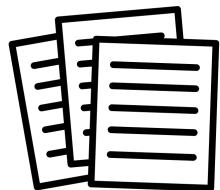
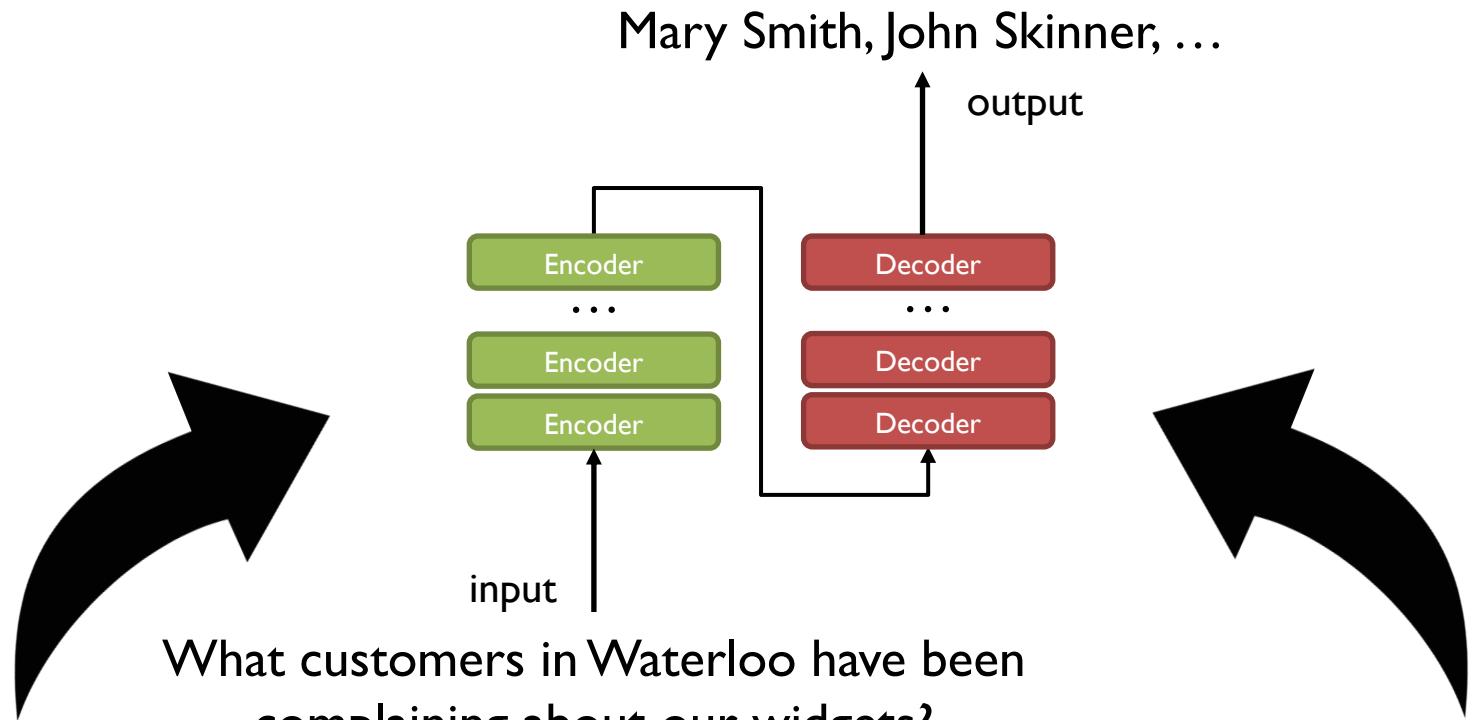
# Natural language access to all your organization's data!

DB

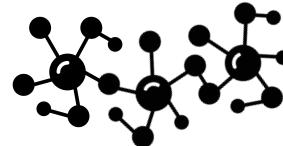


# NLP

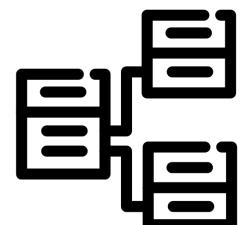
# DB



Unstructured  
Text



Knowledge  
Graphs



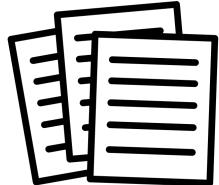
Structured  
Relations

NLP

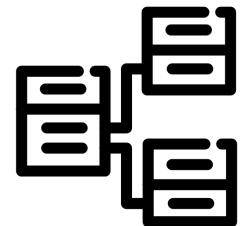
DB



What if BERT took over  
everything?



Unstructured  
Text



Structured  
Relations

NLP

DB



What if transformers took  
over everything?



Unstructured  
Text



Structured  
Relations

A close-up, shallow depth-of-field photograph of pistachio nuts. The background is filled with out-of-focus pistachios, creating a warm, golden-yellow hue. In the foreground, several pistachios are visible, their light brown shells glistening. One pistachio is split open, revealing its vibrant green, fleshy pit. The lighting is soft, highlighting the texture of the nut shells.

# Nuts?

# From Natural Language Processing to Neural Databases

James Thorne  
University of Cambridge  
Facebook AI  
jt719@cam.ac.uk

Fabrizio Silvestri  
Facebook AI  
fsilvestri@fb.com

Majid Yazdani  
Facebook AI  
myazdani@fb.com

Sebastian Riedel  
Facebook AI  
University College London  
riedel@fb.com

Marzieh Saeidi  
Facebook AI  
marzieh@fb.com

Alon Halevy  
Facebook AI  
ayh@fb.com

## ABSTRACT

In recent years, neural networks have shown impressive performance gains on long-standing AI problems, such as answering queries from text and machine translation. These advances raise the question of whether neural nets can be used at the core of query processing to derive answers from facts, even when the facts are expressed in natural language. If so, it is conceivable that we could relax the fundamental assumption of database management, namely, that our data is represented as fields of a pre-defined schema. Furthermore, such technology would enable combining information from text, images, and structured data seamlessly.

This paper introduces *neural databases*, a class of systems that use NLP transformers as localized answer derivation engines. We ground the vision in NEURALDB, a system for querying facts represented as short natural language sentences. We demonstrate that recent natural language processing models, specifically transformers, can answer select-project-join queries if they are given a set of relevant facts. However, they cannot scale to non-trivial databases nor answer set-based and aggregation queries. Based on these insights, we identify specific research challenges that are needed to build neural databases. Some of the challenges require drawing upon literature in data management, and others pose new research opportunities to the NLP community. Finally, we show that with preliminary solutions, NEURALDB can already answer queries over thousands of sentences with very high accuracy.

### PVLDB Reference Format:

James Thorne, Majid Yazdani, Marzieh Saeidi, Fabrizio Silvestri, Sebastian Riedel, and Alon Halevy. From Natural Language Processing to Neural Databases. PVLDB, 14(6): 1033-1039, 2021.  
doi:10.14778/3447689.3447706

## 1 INTRODUCTION

Researchers have long considered the application of neural nets to data management problems, including learning indices [16], query optimization, data cleaning and entity matching [20, 23, 32]. In applying neural networks to data management, research has so far assumed that the data was modeled by a database schema.

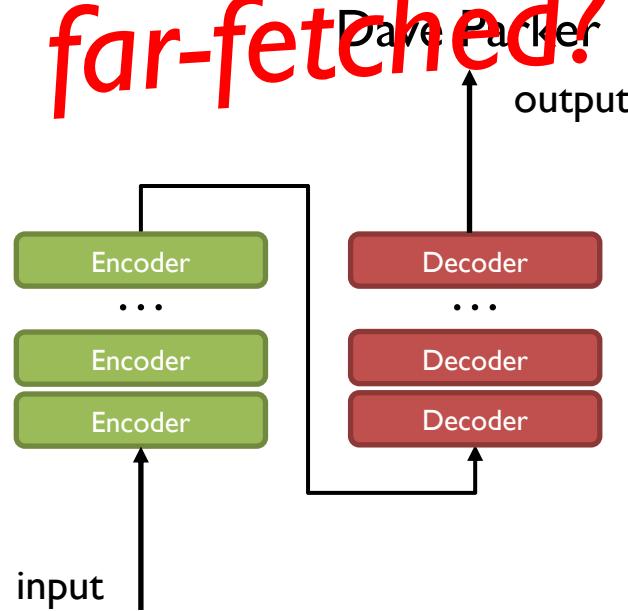
The success of neural networks in processing unstructured data such as natural language and images raises the question of whether their use can be extended to a point where we can relax the fundamental assumption of database management, which is that the data we process is represented as fields of a pre-defined schema. What if, instead, data and queries can be represented as short natural language sentences, and queries can be answered from these sentences? Furthermore, what if relevant data from images can be seamlessly combined with text and structured data?

This paper describes a vision for neural databases and preliminary empirical evidence of its potential. Neural databases offer several benefits that database systems have struggled to support for decades. The first, and most important benefit, is that a neural database has no pre-defined schema. Therefore, the scope of the database does not need to be defined in advance, and any data that becomes relevant as the application is used can be stored and queried. The second benefit is that updates and queries can be posed in a variety of natural language forms, as is convenient to any user. In contrast, a traditional database query needs to be based on the database schema. Even when the data is modeled with a more flexible formalism such as RDF, there is still a single name for any given relation, and that name needs to be used in updates and queries. Third, with recent advances in machine translation, the language of queries and answers can be different from the language of the data in the neural database. A final benefit comes from the fact

NEURALDB

# RDBMS based on Neural Networks

*Is this really that  
far-fetched?*



```
select roster.name from roster, player
where roster.name = player.name and
roster.team = "Pittsburgh Pirates" and
player.nickname = "The Cobra"
```

# RDBMS based on Neural Networks

*Is this really that  
far-fetched?*

*This isn't fundamentally different from  
model-based (approximate) query processing!*

Deshpande and Madden. MauveDB: Supporting Model-Based User Views  
in Database Systems. *SIGMOD 2006*

Mühleisen et al. Capturing the Laws of (Data) Nature. *CIDR 2015*

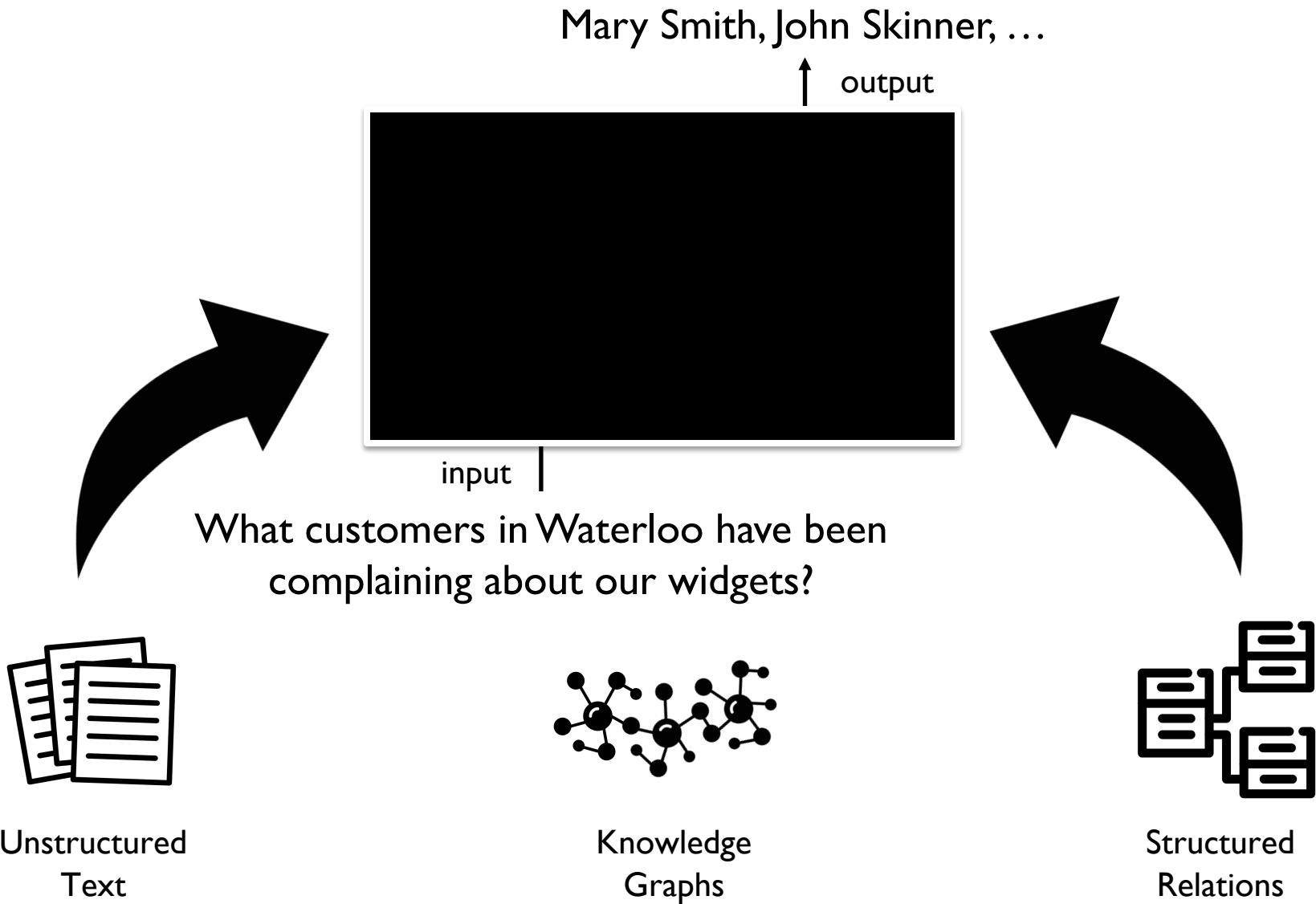
Kulessa et al. Model-based Approximate Query Processing. *arXiv:1811.06224*

Vorona et al. DeepSPACE: Approximate Geospatial Query Processing  
with Deep Learning. *SIGSPATIAL 2019*

NLP

# Natural language access to all your organization's data!

DB



**TRANS  
FORMERS**

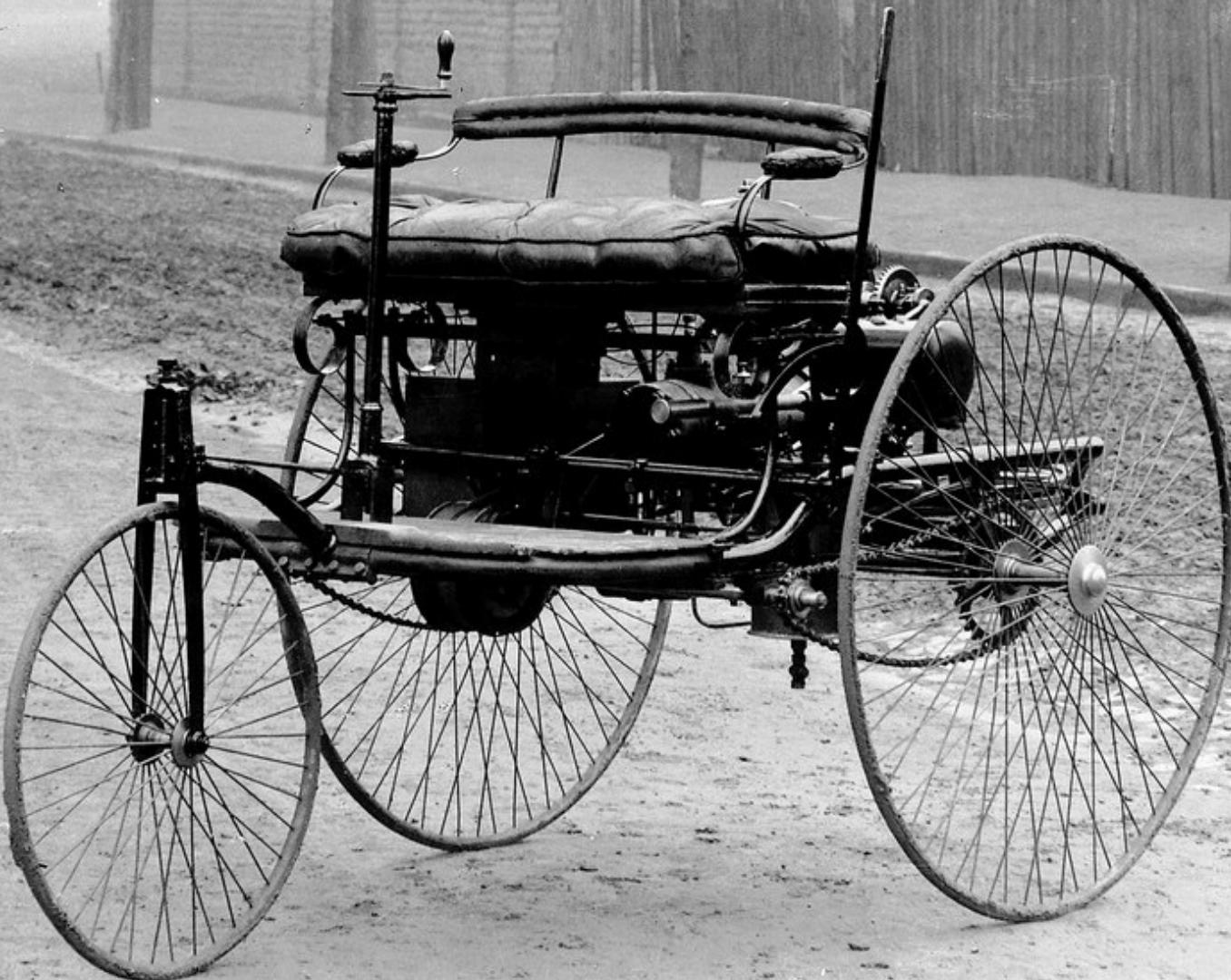
What's QA?  
select-project queries!

What's multi-hop QA?  
select-project-join queries!

Some of the pieces  
What's fact memorization?  
already exist!  
Inserts!

and we've gotten rid of the data!

# We've begun!







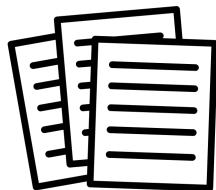
NLP

DB

# TRANSFORMERDB

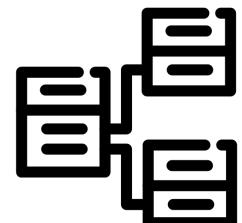
A neural network that ingests all your organization's data and allows you to query it using natural language.

[transformerdb.ai](https://transformerdb.ai)



Unstructured  
Text

(looking for a seed round)



Structured  
Relations