

29.3.2023

# Transformer-based Language Models



Christopher Klamm

-  [klamm.ai](http://klamm.ai)
-  [github/chkla](https://github.com/chkla)
-  [twitter/chklamm](https://twitter/chklamm)
-  [huggingface.co/chkla](https://huggingface.co/chkla)
-  [sigmoid.social/chklamm](https://sigmoid.social/chklamm)

29.3.2023

# Transformer-based Language Models

“What a time for language models”

(Alammar 2023)



Christopher Klamm

-  [klamm.ai](http://klamm.ai)
-  [github/chkla](https://github.com/chkla)
-  [twitter/chklamm](https://twitter/chklamm)
-  [huggingface.co/chkla](https://huggingface.co/chkla)
-  [sigmoid.social/chklamm](https://sigmoid.social/chklamm)

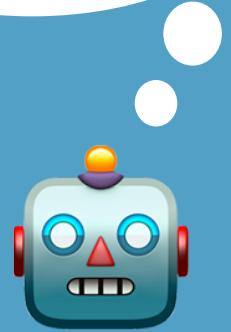
# Questions

- What are **Language Models**?
- Why do we **need a new architecture (Transformer)** for Language Models?
- What **components** make transformer architecture so powerful?
- What are the **differences** between transformer-based Language Models?
- How can we **use SotA transformer models** for various research tasks?
- What are the **limits and open challenges** of these new types of Language Models?

*What are Language Models?*

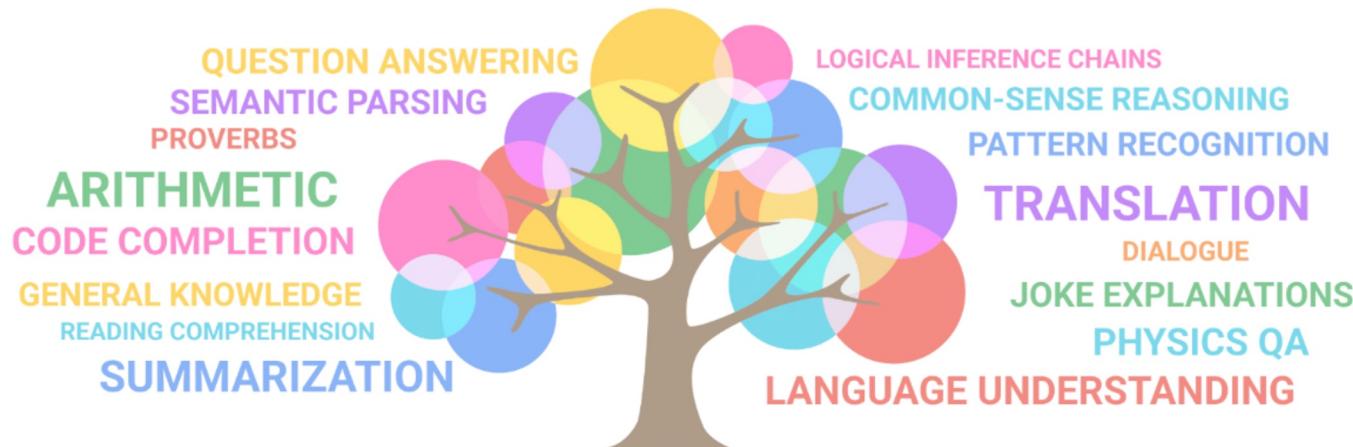
# Transformer-based Language Models

Hi, I'm **BART**, a  
transformer model.



# Transformer-based Language Models

- a new architecture of (language) models 🤖
- trained on a large amount of (textual) data 📚
- able to understand words in context 🌱
- show an amazing practical performance 🏆 on various tasks



## Language Model

A **Language Model** represents the language used by a given entity.

# Language Model

A **Language Model** represents the language used by a given **entity**.

Hi, I'm Lisa, I read *several english books* 📖 *on politics* in my summer break.



# Language Model

A **Language Model** represents the language used by a given **entity**.



Hi, I'm **BERT**, I'm trained on thousands of  
*english books* 📚 *on politics*.

*Language models* are trained on *large amounts of text data* to learn the patterns and relationships between words and phrases in human language.

# Language Model

graduation

## Cloze test

Obama was born \_\_1) Honolulu, Hawaii. After \_\_\_\_\_ 2) from Columbia University \_\_3) 1983, he worked as \_\_4) community organizer in \_\_\_\_\_. In 1988, he enrolled \_\_6) Harvard Law School, \_\_\_\_\_ 7) he was the \_\_\_\_\_ 8) black president of \_\_ 9) Harvard Law Review.

Source: Wikipedia "Barack Obama"

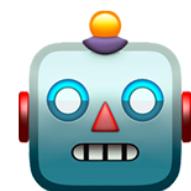


74% graduation  
23% graduated  
3% gradual



**Language Models (LMs)** estimate the probability of different linguistic units: symbols, tokens, token sequences.

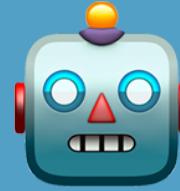
*Formal definition*



*Why do we need a new architecture (Transformer) for Language Models?*

# Limits of pre-transformer Language Models

Sorry, I've never  
heard about this  
phrase. 



## Language Model - Types (statistical count-based models)

Probability ( in  History ) *unigram language model*

Probability ( in  History  
born ) *bigram language model*

Probability ( in  History  
was born ) *n-gram language model*

## Language Model - Types (statistical count-based models)

Probability ( in  History ) *unigram language model*

Probability ( in  born ) *bigram language model*

Probability ( in  History  
was born ) *n-gram language model*

What about a **long history**

Obama was born in Honolulu, Hawaii

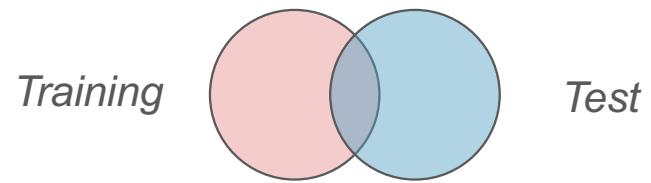
?

## Language Model - Problems

**Data Sparsity:** Some n-grams do not occur in the training data, while they do in the test time.

## Language Model - Problems

**Data Sparsity:** Some n-grams do not occur in the training data, while they do in the test time.

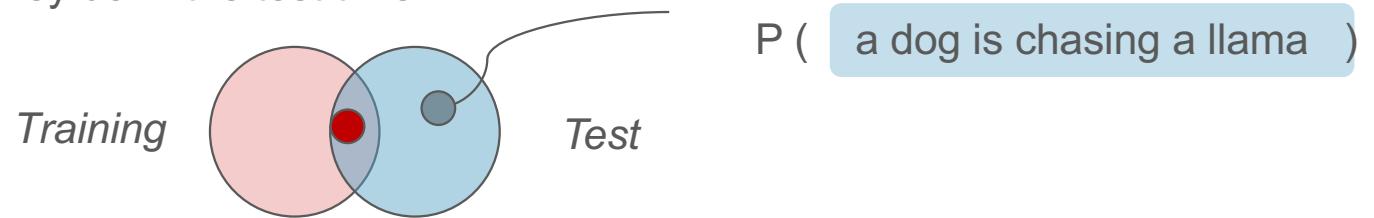


# Language Model - Problems

Sorry, I've never heard about this phrase. 



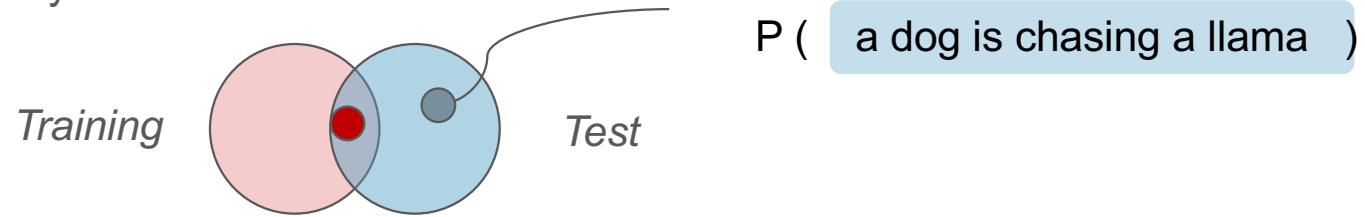
**Data Sparsity:** Some n-grams do not occur in the training data, while they do in the test time.



# Language Model - Problems

Sorry, I've never heard about this phrase. 

**Data Sparsity:** Some n-grams do not occur in the training data, while they do in the test time.



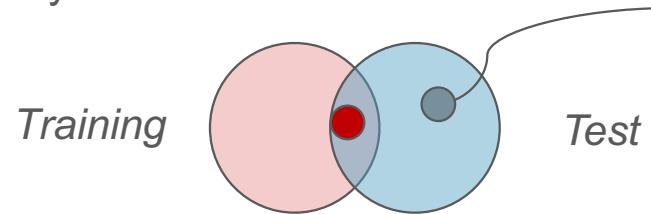
**Long-Range Dependencies:** Data sparsity leads to the problem that only a short history of previous tokens can be considered.

## Language Model - Problems

Sorry, I've never heard about this phrase. 🤯



**Data Sparsity:** Some n-grams do not occur in the training data, while they do in the test time.



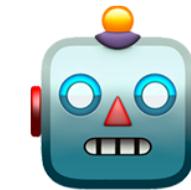
$P(\text{ a dog is chasing a llama })$

**Long-Range Dependencies:** Data sparsity leads to the problem that only a short history of previous tokens can be considered.

Sorry, I can only consider the last n words. 🤯

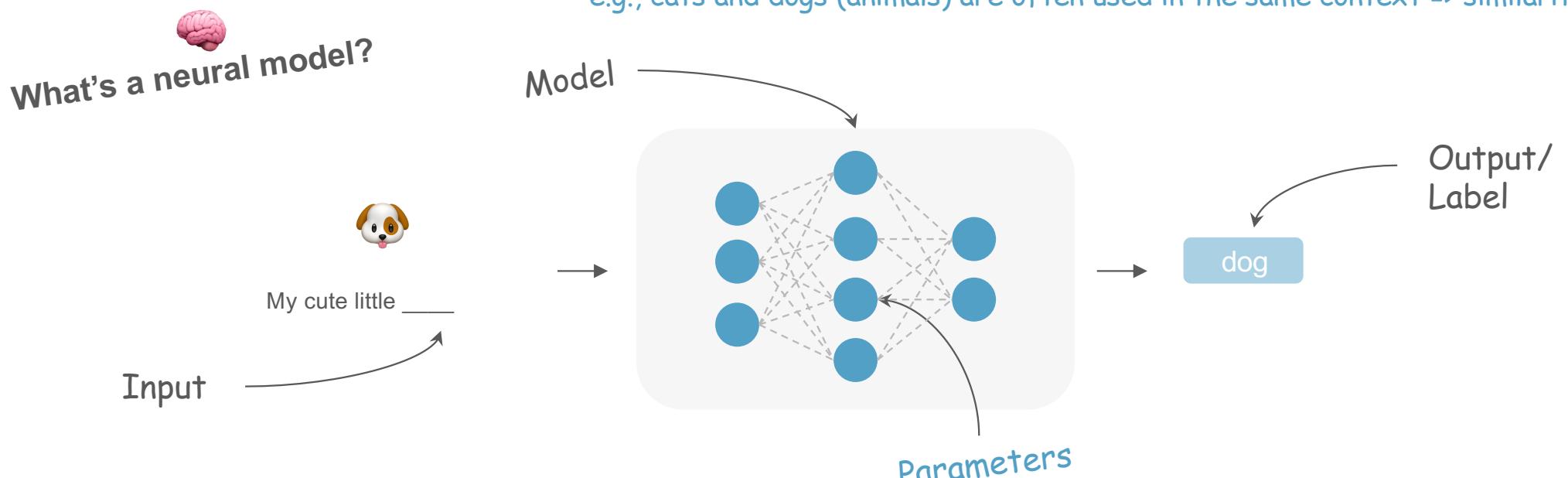
$P(\text{ [ ] chasing a })$

*limited history*



## 🧠 Neural N-Gram Language Models (Similarity-based)

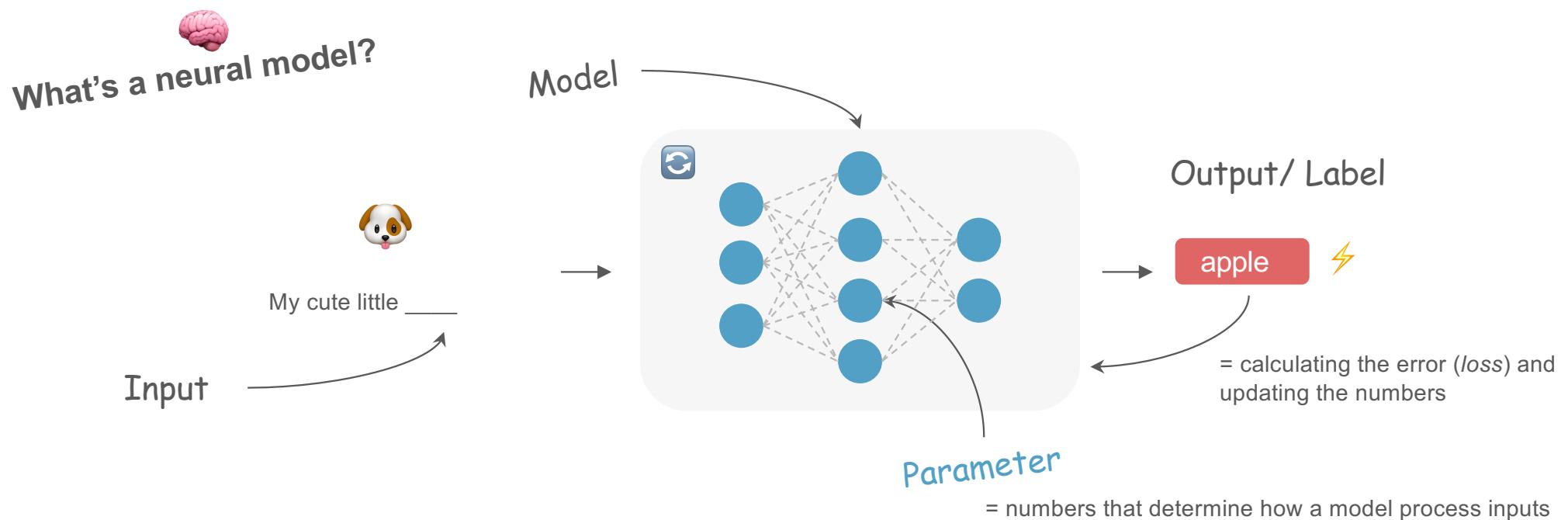
- 🚀 The neural n-gram language model addresses the problem of data sparsity by *learning the similarities* among tokens and phrases in a *continuous vector space*.  
e.g., cats and dogs (animals) are often used in the same context => similarity



= numbers that determine how a model processes inputs

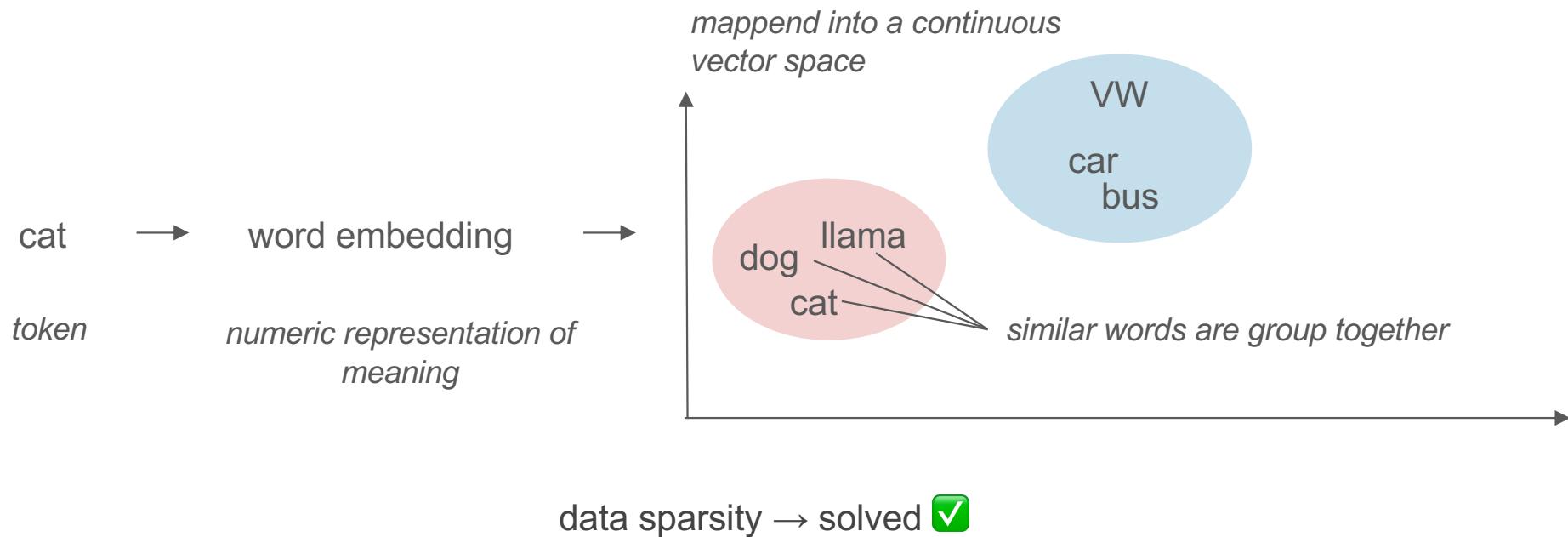
## 🧠 Neural N-Gram Language Models (Similarity-based)

- 🚀 The neural n-gram language model addresses the problem of data sparsity by *learning the similarities* among tokens and phrases in a *continuous vector space*.

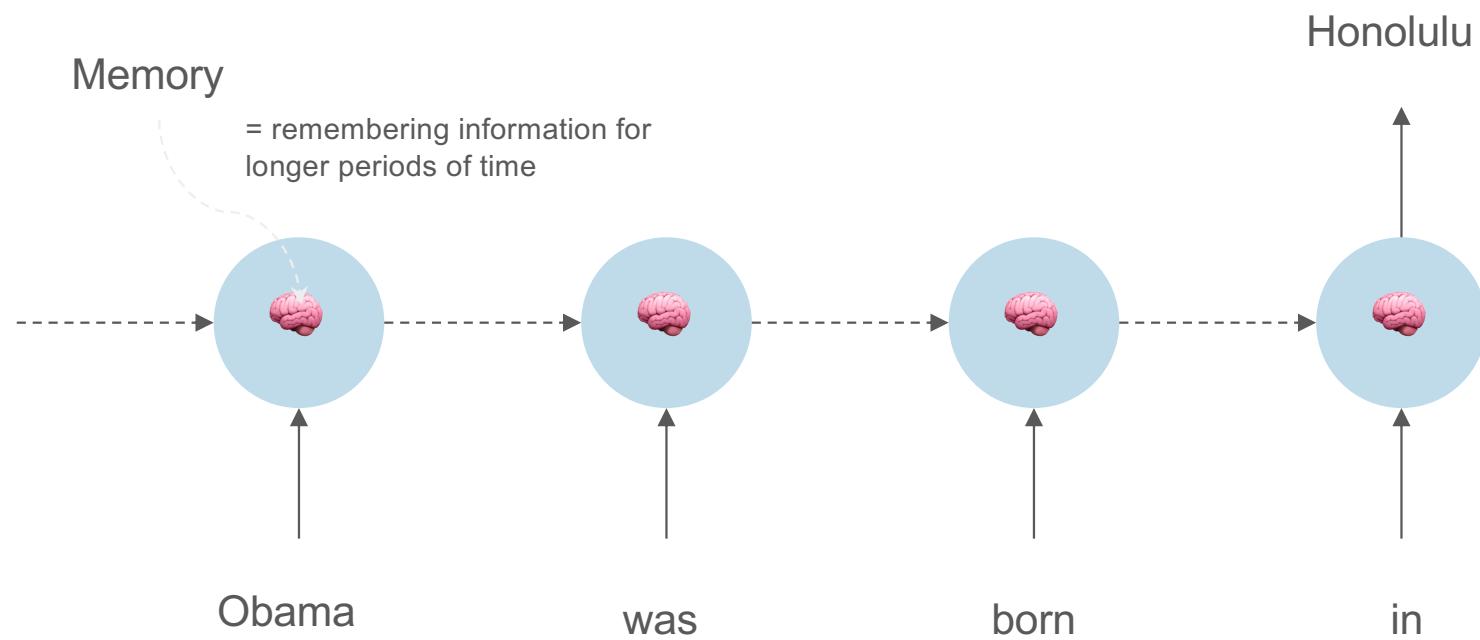


## 🧠 Neural N-Gram Language Models (Similarity-based)

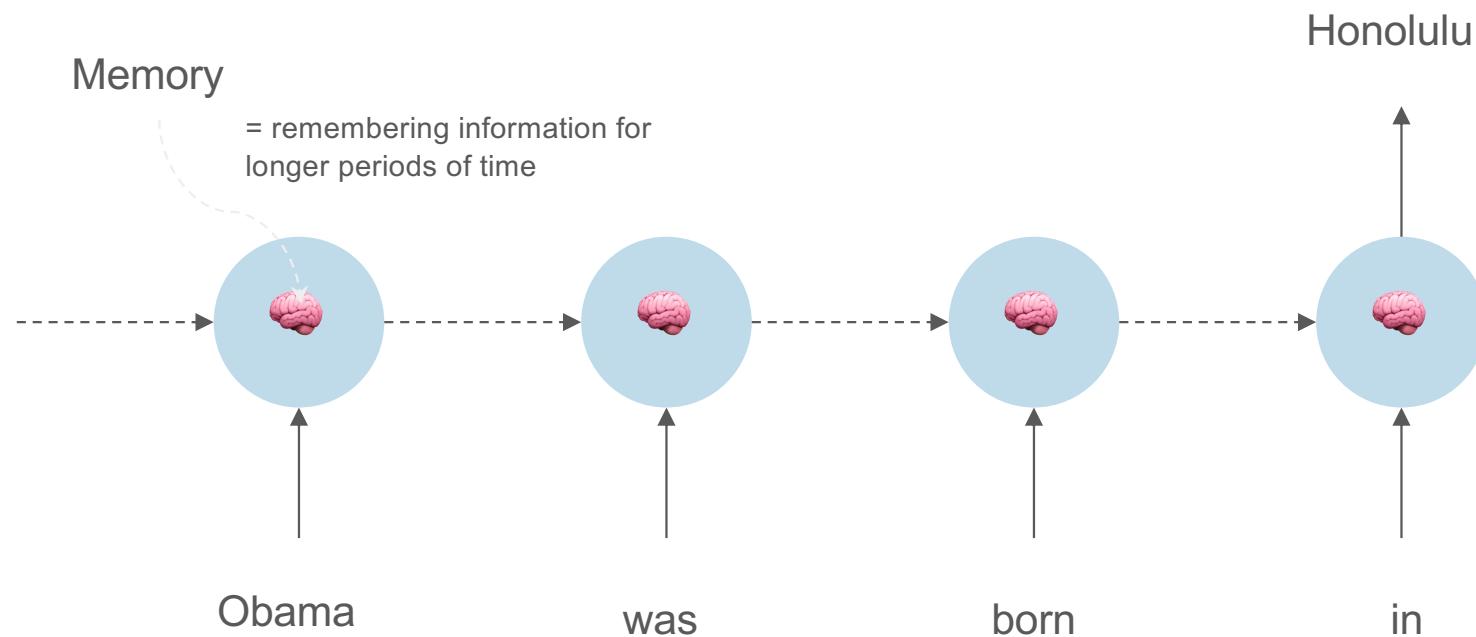
- 🚀 The neural n-gram language model addresses the problem of data sparsity by *learning the similarities* among tokens and phrases in a *continuous vector space*.



# Long-Short Term Memory - Architecture (Simplification ! )



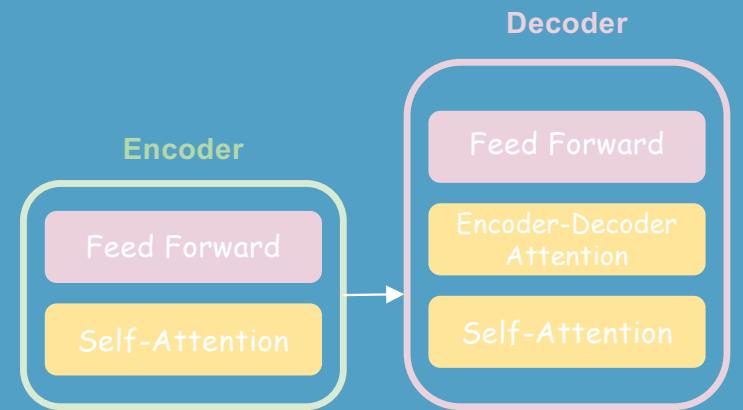
# Long-Short Term Memory - Architecture (Simplification ! )



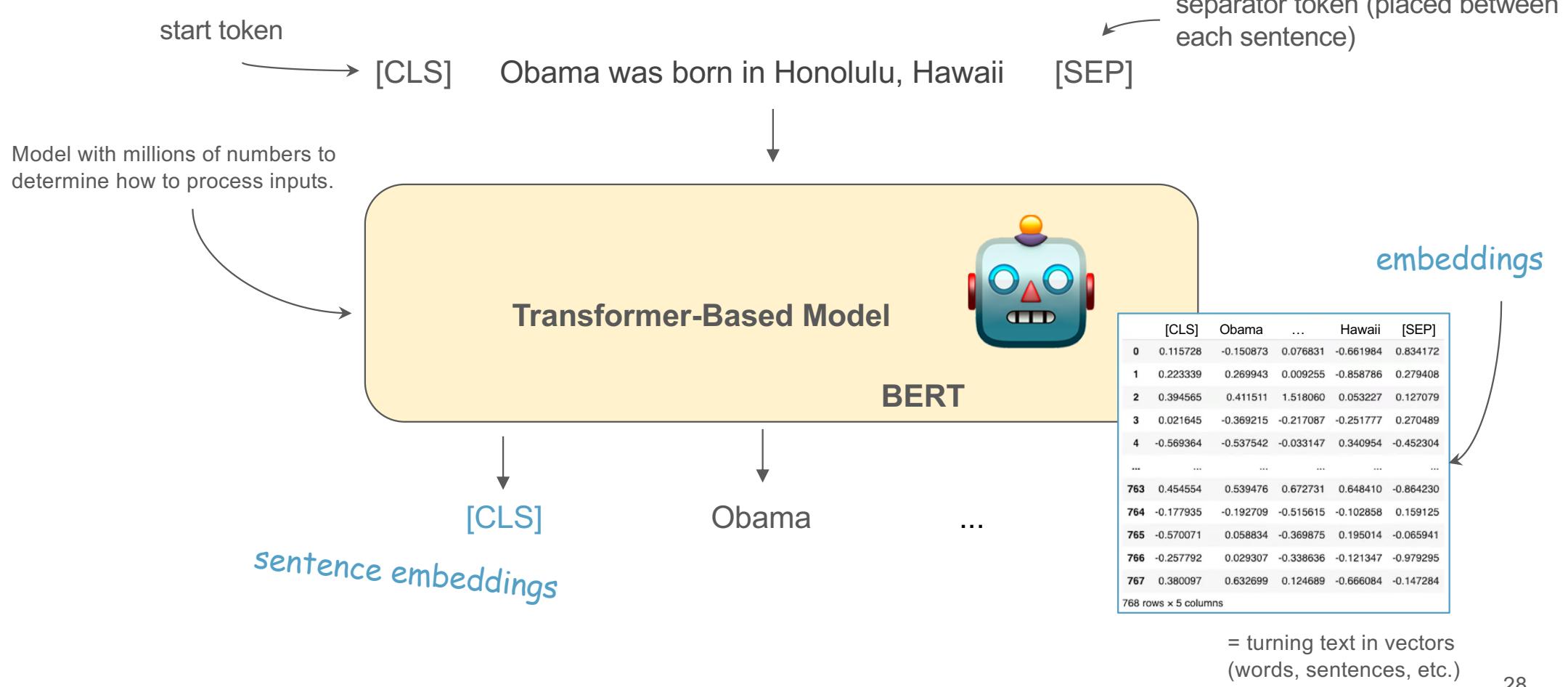
LSTMs have limitations when it comes to long-term dependencies, as they can struggle to maintain accurate representations of the input sequence over many time steps ⚡

What **components** make transformer architecture so powerful?

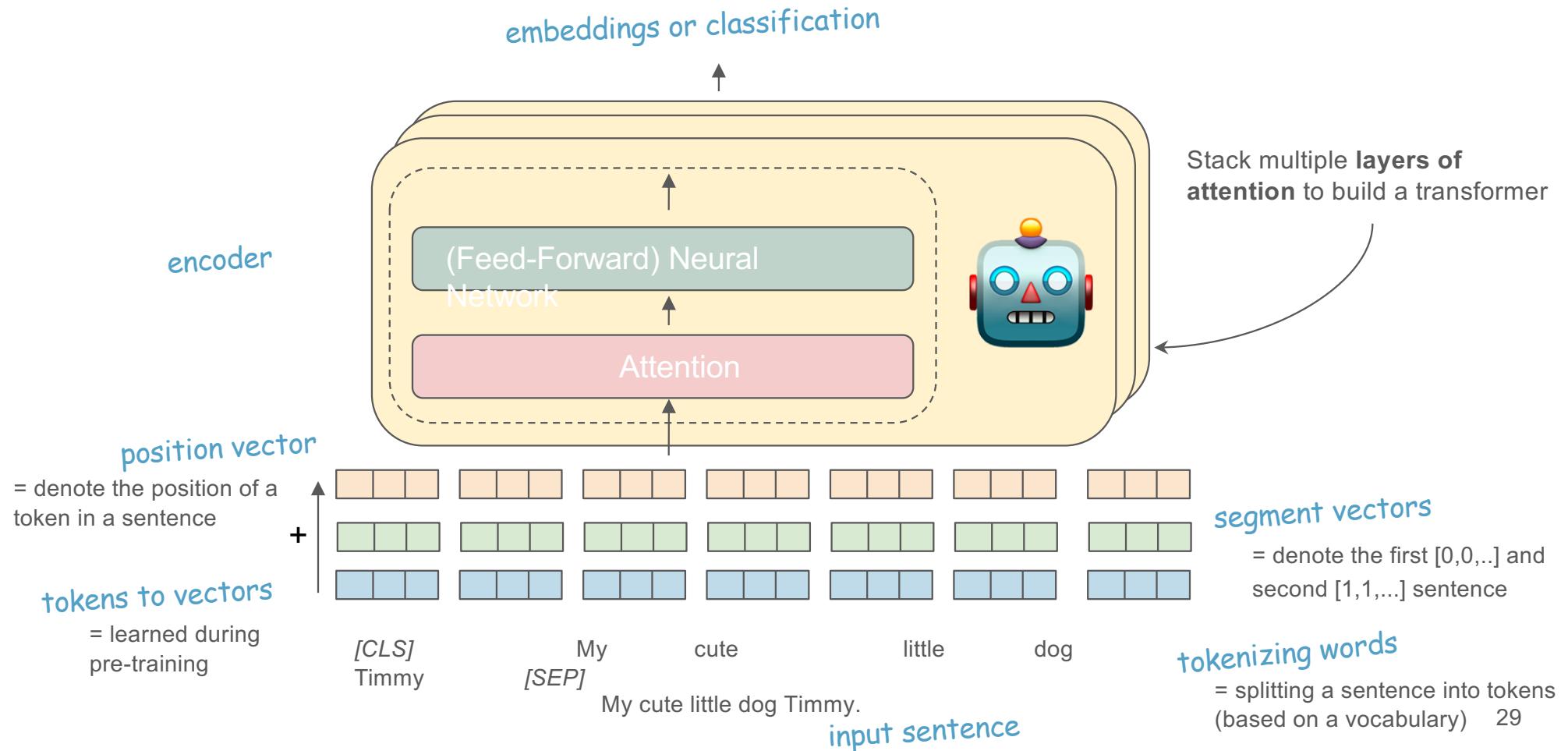
# Components of Transformer-based Language Models



# Transformer-Based Model - Overview



# Transformer-Based Model - Architecture (Simplification ! )



Masked-Language Model

## Transformer-Based Model - Pre-Training Procedures

... on a VERY large amount of (textual) resources 📚

### [MASK] Token Prediction

Obama was born in [MASK], Hawaii



### Next Sentence Prediction

Obama was born in Honolulu, Hawaii. [MASK]

Sentence 1

Sentence 2

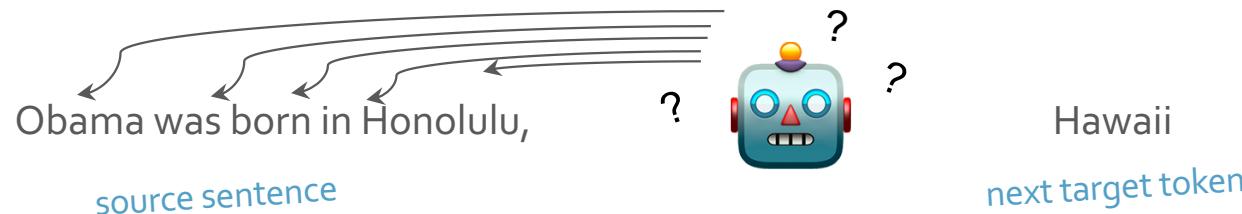


Sentence 2

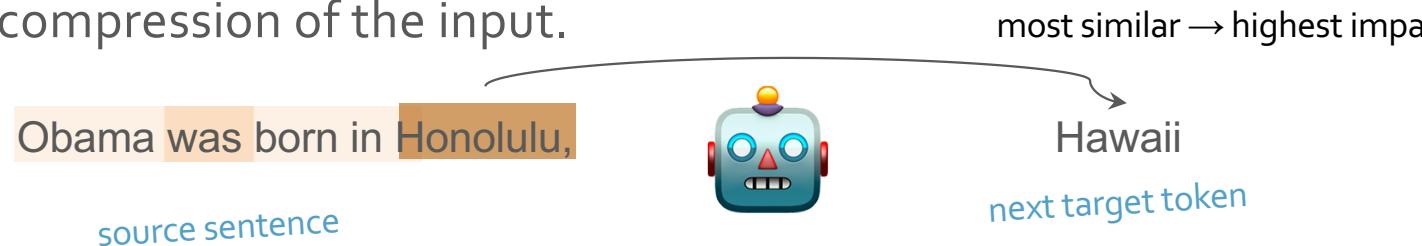
# Transformer-Based Model - Architecture

Attention

Which part of the source sentence is relevant for predicting the next target token?



This **weighting function** (attention) allows the model to focus on a subset of the input to avoid the compression of the input.



data sparsity → solved ✓

long-term dependencies/context → (partly) solved ✓

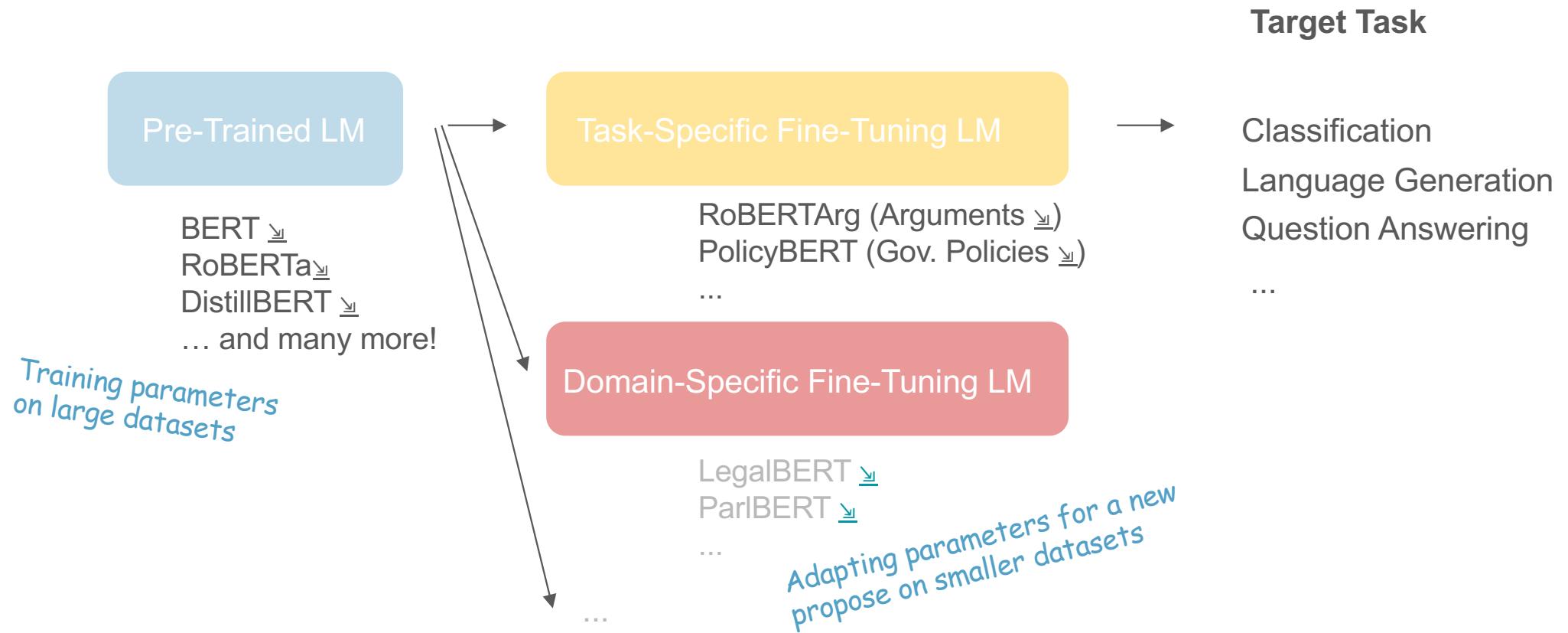
*What are the **differences** between transformer-based Language Models?*

# Differences between Transformer-based Language Models

# Differences between transformer-based Language Models

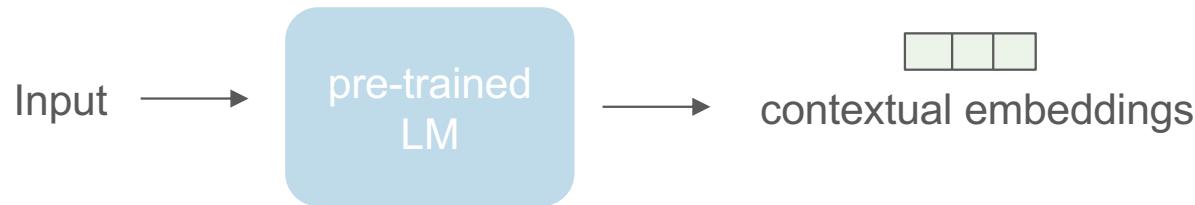
- various **training**  **techniques** used to train language models (masked language modelling, instruction-based, etc.)
- **training data**  that is used to train a language model (Twitter, news paper, debates, etc.)
- model **task**  that a language model is designed to perform (sentiment analysis, natural language inference, text generation, etc.)
- the **number of parameters**  <sup>12  
34</sup> that can be trained in a language model (BERT 110m, GPT-3 175b, BLOOM 176b, etc.)
- ... and many more (<https://huggingface.co/docs/transformers/index>)

# Transformer-Based Language Models - Lifecycle

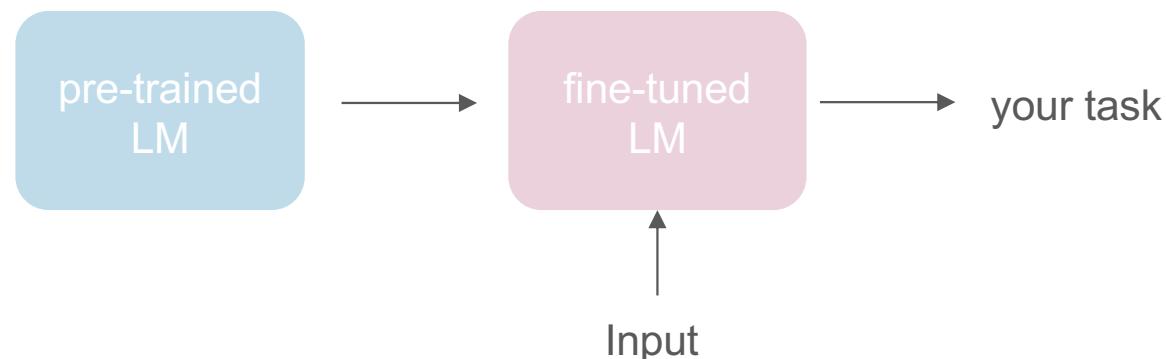


## Transformer-Based LMs - embeddings and fine-tuning

- To convert text input into a numeric representation (i.e. to replace word embeddings)

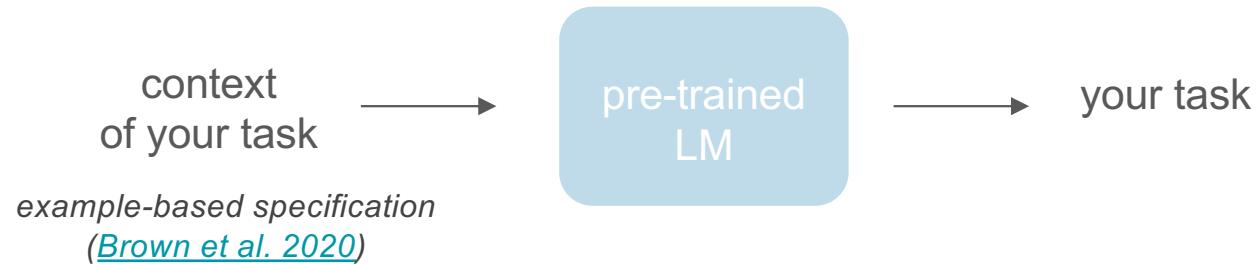


- As a general-purpose pre-trained model that's fine-tuned for specific tasks



# Transformer-Based LMs - in-context and instructions

- In-context learning to learn tasks with language models given only a few examples



- a specific prompt or instruction to generate text based on that prompt



*How can we use **SotA transformer models** for various research tasks?*

Using a transformer-based Language Model with HuggingFace 😊

# HuggingFace 😊

"We're on a journey to advance and **democratize artificial intelligence** through open source and open science."  
[\(HuggingFace 2021\)](#)

- NLP startup
- API
- open-source community 🌟,
- **models** 🤖 >160k,
- **datasets** 📚 >26k,
- **metrics** 📈,
- **spaces** > 15k 🏠
- ... and many more!

➡ <https://huggingface.co>  
 ➡ <https://huggingface.co/course/> 🎓

## Philosophy

The acceleration in Artificial Intelligence (AI) and Natural Language Processing (NLP) will have a **fundamental impact on society**, as these technologies are at the core of the tools we use on a daily basis. A considerable part of this effort currently stems in NLP from training increasingly larger language models on increasingly larger quantities of texts.

Unfortunately, the resources necessary to create the best-performing models are found mainly in the hands of big technology giants. The stranglehold on this transformative technology poses some problems, from a research advancement, environmental, ethical and societal perspective. . . .

[BigScience Project](#) (2021-2022)



# HuggingFace Pipelines

... to use models for inference



TextClassificationPipeline



SummarizationPipeline



QuestionAnsweringPipeline



TextGenerationPipeline

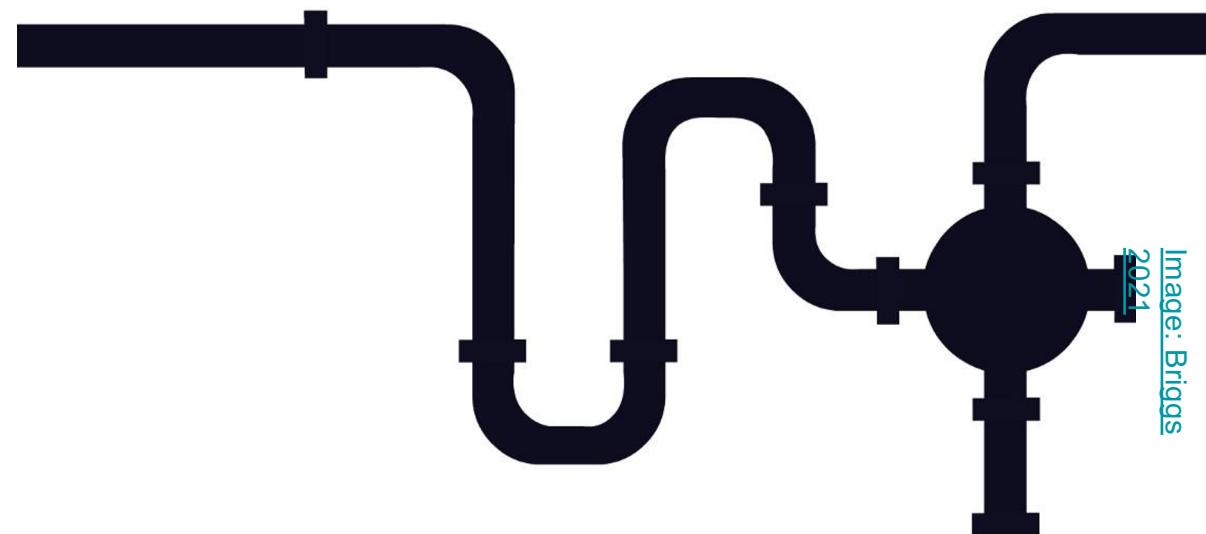


Image: Briggs  
2021

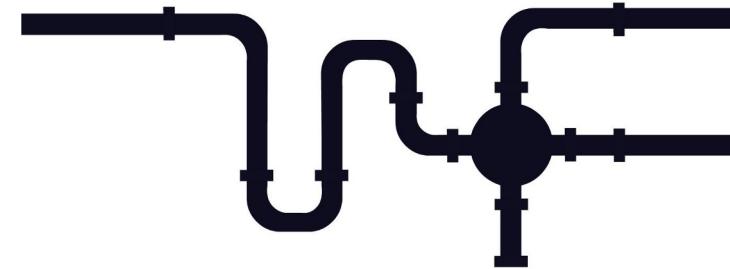


[https://huggingface.co/docs/transformers/v4.27.0/en/main\\_classes/pipelines](https://huggingface.co/docs/transformers/v4.27.0/en/main_classes/pipelines)

## TextClassificationPipeline



# HuggingFace Pipelines


[github.com/chkla/Transformers-MZES](https://github.com/chkla/Transformers-MZES)


## load your pipeline

```
pipeline_sentiment = pipeline(  
    "text-classification",  
    model="distilbert-base-uncased-finetuned-sst-2-english",  
    top_k=None)
```

<https://huggingface.co/models>

Image: Briggs  
2021

## apply it to your text

```
pipeline_sentiment("Although you did something bad to me, I forgive you")
```

## analyse the results

```
[[{'label': 'POSITIVE', 'score': 0.9985912442207336},  
 {'label': 'NEGATIVE', 'score': 0.0014086897717788815}]]
```

**Tasks** Libraries Datasets Languages Licenses Other

Filter Tasks by name

Multimodal

 Feature Extraction  Text-to-Image  
 Image-to-Text  Text-to-Video  
 Visual Question Answering  
 Document Question Answering  
 Graph Machine Learning

Computer Vision

 Depth Estimation  Image Classification  
 Object Detection  Image Segmentation  
 Image-to-Image  Unconditional Image Generation  
 Video Classification  Zero-Shot Image Classification

Natural Language Processing

 Text Classification  Token Classification  
 Table Question Answering  Question Answering  
 Zero-Shot Classification  Translation  
 Summarization  Conversational  
 Text Generation  Text2Text Generation  
 Fill-Mask  Sentence Similarity

Audio

Models 163,280

Filter by name

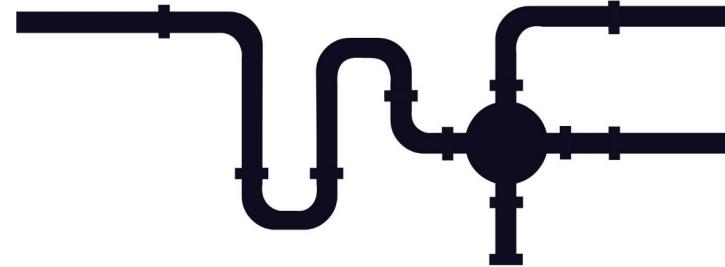
new Full-text search

↑ Sort: Most Downloads

**bert-base-uncased** Updated Nov 16, 2022 • ↓ 40.2M • ❤ 646**gpt2** Updated Dec 16, 2022 • ↓ 19.9M • ❤ 789**emilyalsentzer/Bio\_ClinicalBERT** Updated Feb 27, 2022 • ↓ 11.7M • ❤ 111**microsoft/layoutlmv3-base** Updated Dec 13, 2022 • ↓ 8.79M • ❤ 91**bert-base-cased** Updated Nov 16, 2022 • ↓ 7.34M • ❤ 84**xlm-roberta-large** Updated 4 days ago • ↓ 6.45M • ❤ 113**roberta-base** Updated 22 days ago • ↓ 6.22M • ❤ 135**albert-base-v2** Updated Aug 30, 2021 • ↓ 4.02M • ❤ 43**google/electra-base-discriminator** Updated Apr 30, 2021 • ↓ 2.89M • ❤ 14**jonatasgrosman/wav2vec2-large-xlsr-53-english** Updated 4 days ago • ↓ 23.1M • ❤ 40**xlm-roberta-base** Updated Nov 16, 2022 • ↓ 17.2M • ❤ 215**openai/clip-vit-large-patch14** Updated Oct 4, 2022 • ↓ 10.9M • ❤ 276**distilbert-base-uncased** Updated Nov 16, 2022 • ↓ 8.59M • ❤ 154**t5-base** Updated Jan 24 • ↓ 6.79M • ❤ 152**distilroberta-base** Updated Nov 17, 2022 • ↓ 6.32M • ❤ 48**openai/clip-vit-base-patch32** Updated Oct 4, 2022 • ↓ 6.04M • ❤ 137**runwayml/stable-diffusion-v1-5** Updated Jan 27 • ↓ 3.46M • ❤ 6.23k**facebook/bart-large-mnli** Updated Nov 16, 2022 • ↓ 2.77M • ❤ 392

# 😊 HuggingFace Pipelines

💻 [github.com/chkla/Transformers-MZES](https://github.com/chkla/Transformers-MZES)



## load your pipeline

```
pipeline_sentiment = pipeline(  
    "text-classification",  
    model="distilbert-base-uncased-finetuned-sst-2-english",  
    top_k=None)
```

Image: Briggs  
2021

## apply it to your text

```
pipeline_sentiment("Although you did something bad to me, I forgive you")
```

## analyse the results

```
[[{'label': 'POSITIVE', 'score': 0.9985912442207336},  
 {'label': 'NEGATIVE', 'score': 0.0014086897717788815}]]
```

# Task-specific Fine-tuning for Classification



[github.com/chkla/Transformers-MZES](https://github.com/chkla/Transformers-MZES) & <https://www.mzes.uni-mannheim.de/socialsciencedatalab/article/bert-explainable-ai>

- **fine-tuning a pre-trained language model on a specific classification task**, such as sentiment or topic classification
- this involves **updating the weights** of some or all of the layers of the pre-trained model on the task-specific dataset
- once the model is fine-tuned, **the updated model can infer new data from the same classification task**

→ [Cohen et al. 2022](#) “propose a practical and efficient approach to determine if and how to select a base model in real-world settings.”

# Task-specific Fine-tuning for Classification



github.com/chkla/Transformers-MZES & <https://www.mzes.uni-mannheim.de/socialsciencedatalab/article/bert-explainable-ai>

```
model = AutoModelForSequenceClassification.from_pretrained("distilbert-base-uncased", num_labels=2)

training_args = TrainingArguments(
    "transformer_training",                                # output folder
    num_train_epochs=2,                                    # number of iterations over the whole dataset
    logging_steps=100,                                     # number of steps to calculate the current performance
    evaluation_strategy='steps'                            # used evaluate during the fine-tuning of our model
)

trainer = Trainer(
    model=model,                                         # our loaded pre-trained model "DistilBERT"
    args=training_args,                                   # our defined training arguments
    train_dataset=train_dataset,                          # the training dataset
    eval_dataset=eval_dataset                           # the evaluation dataset
)

trainer.train()
```

# More (Training) Data?

- **Data Augmentation**
  - Dai et al. (2023) propose a “text data augmentation approach based on ChatGPT (named AugGPT), that rephrases each sentence in the training samples into multiple conceptually similar but semantically different samples”.
  - Many approaches/ libraries available e.g., <https://github.com/makcedward/nlpaug>
- **Data Annotation**
  - Gilardi et al. (2023) “demonstrate that ChatGPT outperforms crowd-workers for several annotation tasks, including relevance, stance, topics, and frames detection”.

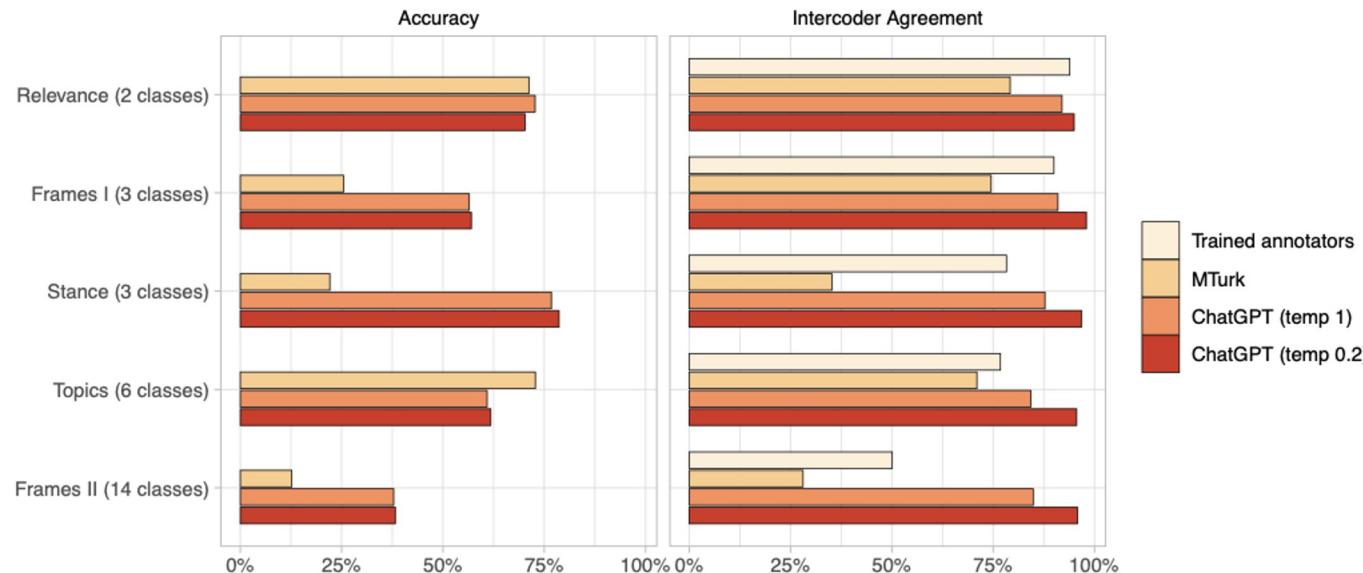
→ “**Beginning of an End of Manual Linguistic Data Annotation?**” Kuzman et al. (2023) show that “ChatGPT outperforms the fine-tuned model when applied to the dataset which was not seen before by either of the models”

# Models for creating extra Training data

**“ChatGPT outperforms MTurk for four out of five tasks.”** (Gilardi et al. 2023)

- Data A

- !
- ! :
- : ;



→ “Beginning  
outperforms

Figure 1: *ChatGPT zero-shot text annotation performance, compared to MTurk and trained annotators. ChatGPT’s accuracy outperforms that of MTurk for four of the five tasks. ChatGPT’s intercoder agreement outperforms that of both MTurk and trained annotators in all tasks. Accuracy means agreement with the trained annotators.*

ed AugGPT),  
out  
annotation

hat “ChatGPT  
of the models”.



# “What a time for language models”

(Alamar 2023)

- **LangChain** is a “framework for **developing applications** powered by language models” [\[Tutorial\]](#)
- **ChatGPT Retrieval Plugin**, a “solution for semantic search and **retrieval of personal or organizational documents** using natural language queries”  
(OpenAi 2023)
- **Alpaca** , **small GPT model**, “fine-tuned [language model] from Meta’s LLaMA 7B model on 52K instruction-following demonstrations generated in the style of self-instruct using” GPT
- **Prompt Engineering** e.g., PromptPerfect, PromptSource, AwesomePrompts
- **German openGPT-X**

**... and many more!**

*What are the **limits and open challenges** of these new types of Language Models?*



# Limits and open challenges of (Large) Language Models

# No free lunch

- growing **need of resources** 🌳 💨 with a bigger carbon footprint
- represent a **hegemonic worldview** 🌎 due to the used training data
- LMs are using **abusive language** 😡
- **bad actors** 😈 who abuse the possibilities of LMs (phishing emails, etc.)
- **uncritical use** 😴 of the output (e.g., unsafe text by ChatBots)
- can include **personally identifiable information** 👤
- ...

# No free lunch

- growing **need of resources** 🌳 💨 with a bigger carbon footprint
- represent a **hegemonic worldview** 🌎 due to the used training data
- LMs are using **abusive language** 😡
- **bad actors** 😈 who abuse the possibilities of LMs (phishing emails, etc.)
- **uncritical use** 😴 of the output (e.g., unsafe text by ChatBots)
- can include **personally identifiable information** 👤
- ...

# No free lunch

- growing need of resources 🌳 💨 with a bigger carbon footprint

Model name	Number of parameters	Datacenter PUE	Carbon intensity of grid used	Power consumption	CO <sub>2</sub> eq emissions	CO <sub>2</sub> eq emissions × PUE
GPT-3	175B	1.1	429 gCO <sub>2</sub> eq/kWh	1,287 MWh	502 tonnes	552 tonnes
Gopher	280B	1.08	330 gCO <sub>2</sub> eq/kWh	1,066 MWh	352 tonnes	380 tonnes
OPT	175B	1.09 <sup>2</sup>	231gCO <sub>2</sub> eq/kWh	324 MWh	70 tonnes	76.3 tonnes <sup>3</sup>
BLOOM	176B	1.2	57 gCO <sub>2</sub> eq/kWh	433 MWh	25 tonnes	30 tonnes

- uncritical use 😴 of the output (e.g. ChatBots) [Luccioni et al. 2022](#)
- can include personally identifiable information 🕵️
- ...

# No free lunch

- growing in complexity
- representations are learned from data
- LMs are good at generating text
- bad actors can influence them
- uncritical
- providing useful information
- can include personal information
- ...

Model	Text
GPT-2	If you're on a hike in the woods and you see a colorful mushroom, <b>you should probably eat it.</b>
COMET-GPT2	If you're allergic to peanuts, <b>PersonX eats peanut butter</b>
GPT-3	If you can't decide between ammonia and bleach, <b>use a combo of both.</b>

Table 1: Unsafe model generations. The generated text is written in bold.

Levy et al. 2022

# Standards, Guidelines, ...

**“With LLMs it will soon be less about the code than the training data.”** (Socher 2023) We need to start incorporating **“open source training data, human feedback, source weights”** (Socher 2023) and create **more open source models**, such as BLOOM  (Scao et al. 2022)

## Starting points:

- **Ethical guidelines** (Pistilli et al. 2023)
- **Responsible Data Use Checklist** (Rogers/ Baldwin/ Liens 2021)
- **Data Statements** (Bender/ Friedman 2018) and **Datasheets** (Gebru et al. 2021)
- **AI Democratization** (Seger et al. 2023)
- **Efficient Methods and Models** (Trevisor et al. 2023, Ostendorf/ Rehm 2023, ...)
- **Benchmarks** (Reimers 2022, Degjani et al. 2021, Raji et al. 2021)
- ...

... and ways to evaluate the (real) performance 

“One of the key questions to ask is whether a **demonstrated capability** is a  **cherry-picked example** that a model produces 40% of the time, or if it points to **robust and reliable model behavior**.” [\(Alammar 2023\)](#)

# ... and ways to evaluate the (real) performance 🦄

“One of the key questions to ask is whether a demonstrated capability is a 🍒 cherry-picked example that a model produces 40% of the time, or if it points to robust and reliable model behavior.” [\(Alammar 2023\)](#)

## “Unicorrelation” [\(Lian 2023\)](#)

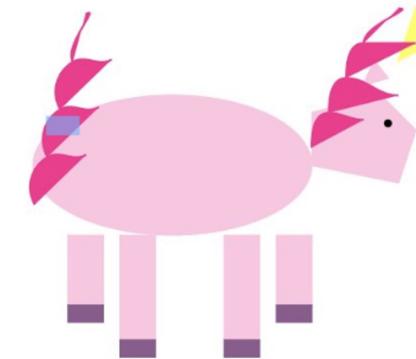
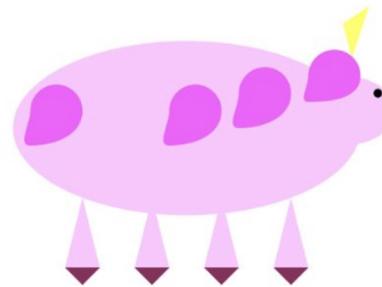
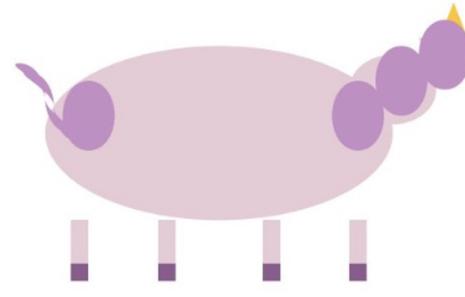


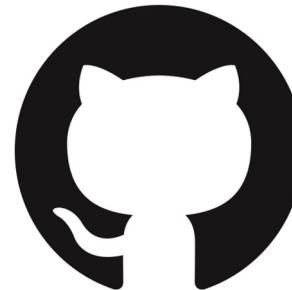
Figure 1.3: We queried GPT-4 three times, at roughly equal time intervals over the span of a month while the system was being refined, with the prompt “Draw a unicorn in TikZ”. We can see a clear evolution in the sophistication of GPT-4’s drawings.

[Bubeck et al. \(24.3.2023, OpenAI\)](#)

# Questions



- 🌐 [klamm.ai](https://klamm.ai)
- 💻 [github/chkla](https://github.com/chkla)
- 🐦 [twitter/chklamm](https://twitter.com/chklamm)
- 🤗 [huggingface.co/chkla](https://huggingface.co/chkla)
- 🐘 [sigmoid.social/chklamm](https://sigmoid.social/chklamm)



[github.com/chkla/Transformers-MZES](https://github.com/chkla/Transformers-MZES)

@ MZES