

# Why are (L)LMs powerful?



Size



Attention mechanism



Layers



Etc.



**Transfer Learning**

# Transfer Learning

## What:

Store 'language knowledge' and 'task knowledge' in the parameters of a model.

## Why:

Learn new tasks faster & better.

# Classical Algorithms

## (e.g. Regression, SVM)

No prior  
**‘language knowledge’**

No knowledge of semantic similarities between words like "attack", "war" and "tree".

No prior  
**‘task knowledge’**

No knowledge of tasks like "*Classify this text into ‘activist’ or ‘conservative’ rhetoric*".

# Word Embeddings (e.g. Word2Vec)

Provide 'language knowledge':

- Represents "*attack*", "*war*" in similar static vectors

Attack  $\approx$

0.1	0.8	...	-0.5	0.1
-----	-----	-----	------	-----

War  $\approx$

0.2	0.8	...	-0.4	0.3
-----	-----	-----	------	-----

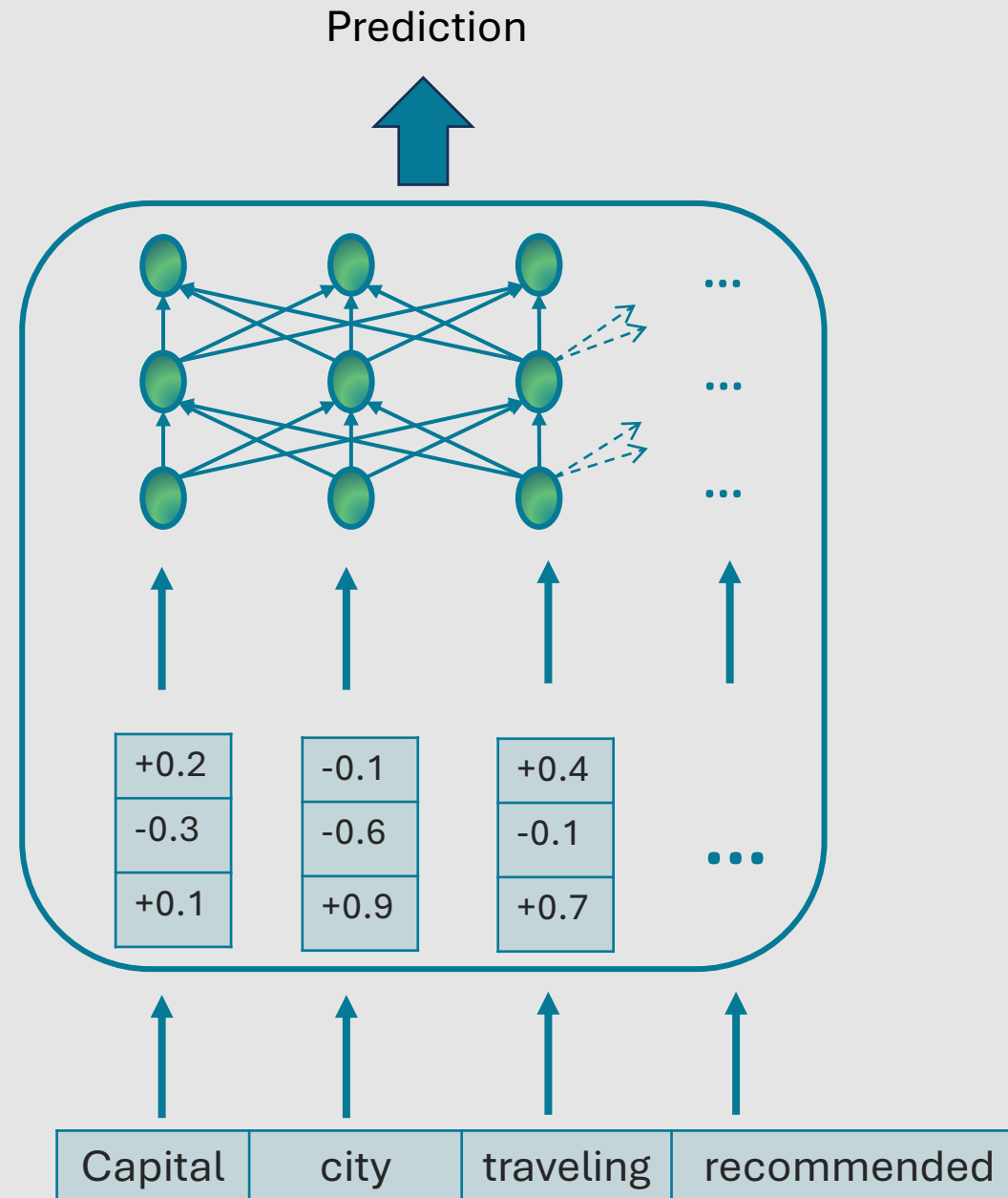
Tree  $\approx$

0.7	-0.4	...	0.1	-0.7
-----	------	-----	-----	------

# Transformers (e.g. BERT-base)

Prior 'language knowledge':

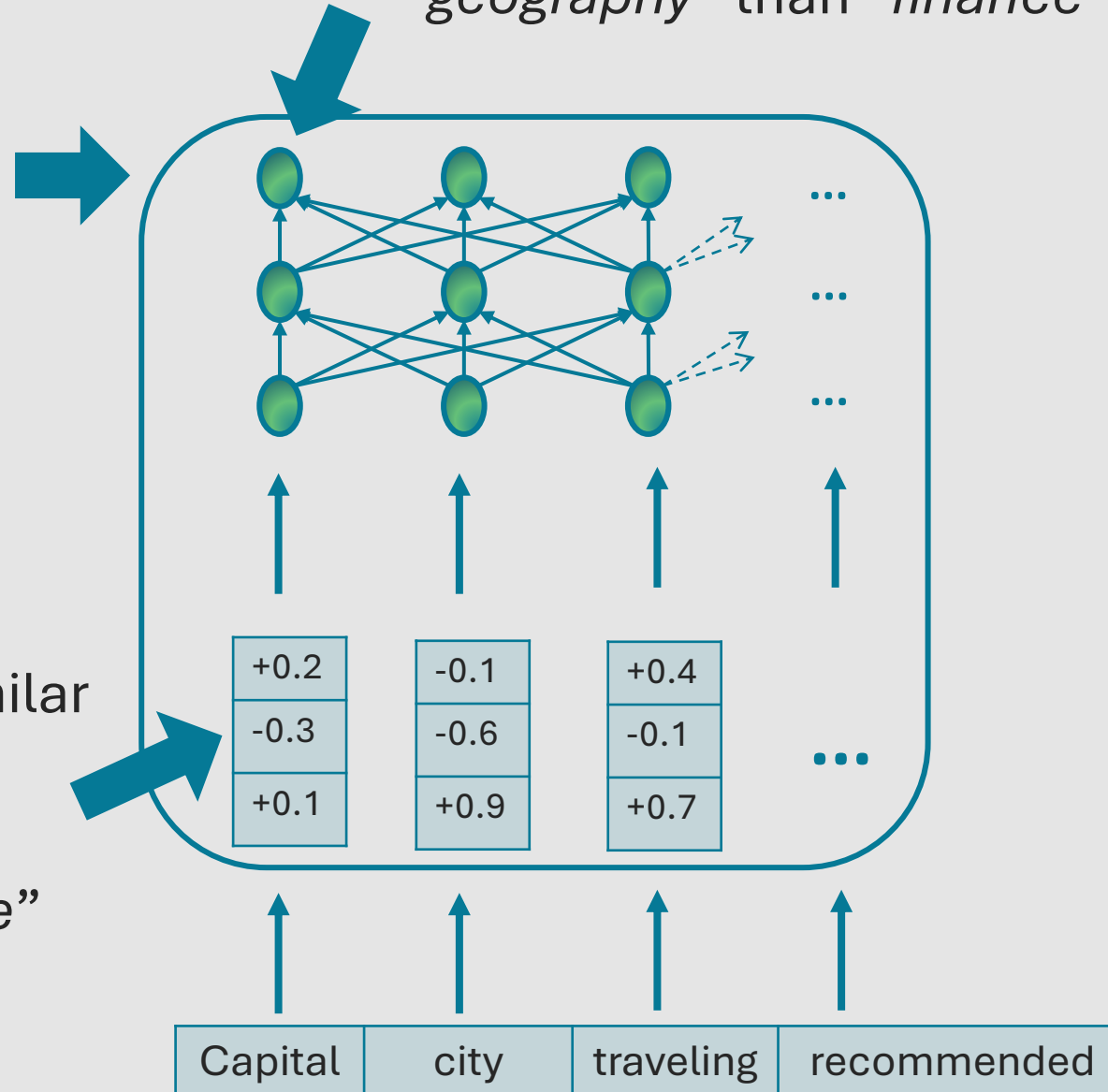
- Represents "*attack*", "*war*" in similar vectors
- Represents "*capital*" differently in context of "*city*" or "*investment*" or "*crime*"



Word vector of “*capital*” now closer to “*geography*” than “*finance*” or “*crime*”

Last layer: Contextualised representation of input

Word vector of “*capital*” with similar distance to “*geography*” / “*finance*” / “*crime*”



# Transformers (e.g. BERT-base)

Prior 'language knowledge':

- Represents "*attack*", "*war*" in similar vectors
- Represents "*capital*" differently in context of "*city*" or "*investment*" or "*crime*"

Learned through **simple, self-supervised task**:

- Masked Language Modelling



# How BERT acquires language knowledge

## ORIGINAL TEXT

“Capital punishment, also known as the death penalty and formerly called judicial homicide, is ...”



## MASKED TEXT

“Capital [MASK], also known as the death [MASK] and formerly called [MASK] homicide, is ...”

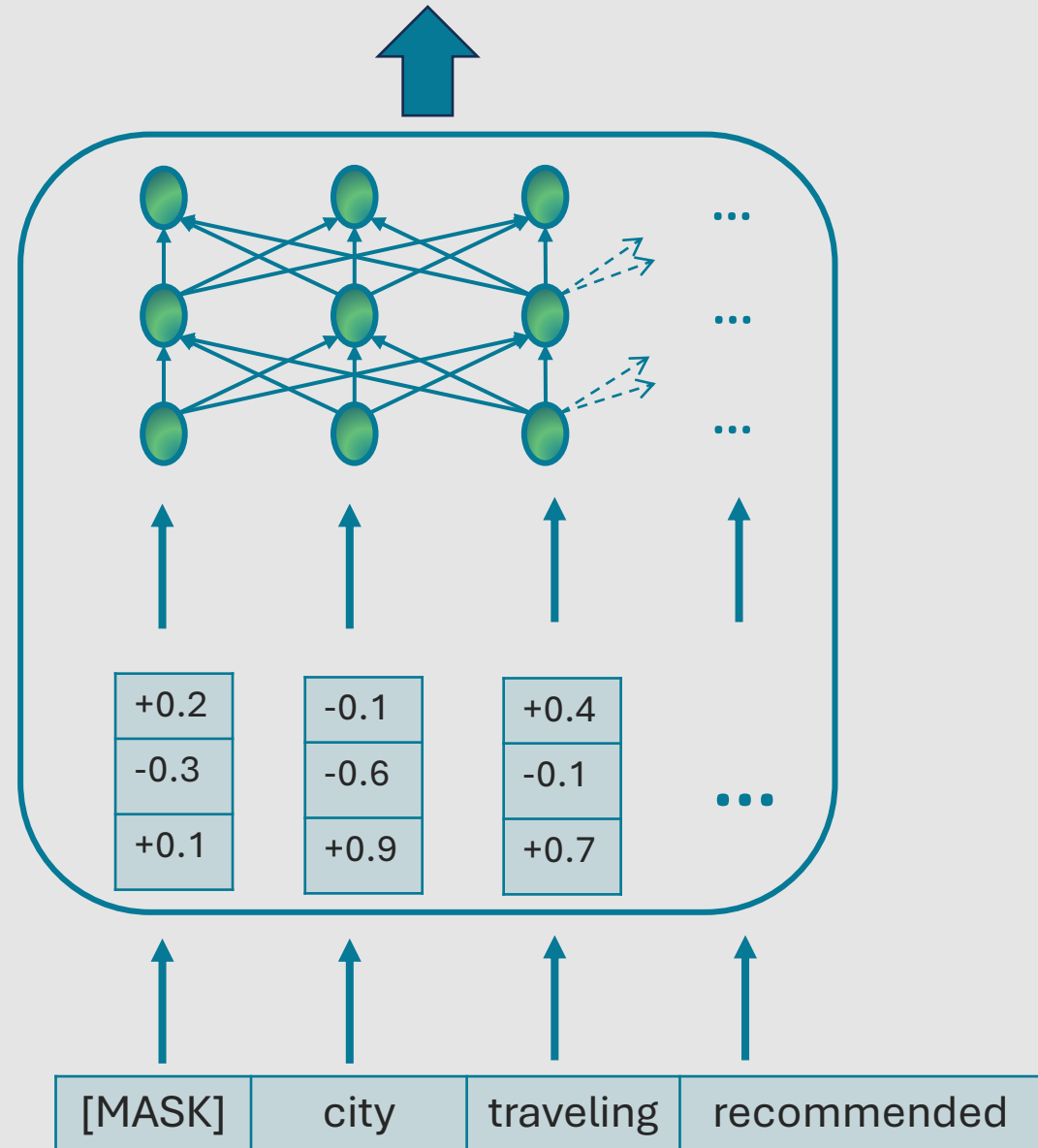
The algorithm learns to predict the correct word behind the [MASK] token.

This creates general **'language knowledge'**

# BERT-base

After millions of training iterations:

Many parameters (vectors),  
which are very good at  
predicting hidden words.



# Disadvantages of BERT-base

**‘Task knowledge’ from MLM is not useful:**

- BERT-base only knows how to predict hidden words (MLM task)
- We are actually interested in other tasks like classification, summarization etc.
- BERT-base needs to learn new, useful tasks from scratch

# Reusing more 'prior knowledge'

## Universal tasks

- Natural Language Inference (BERT-NLI)
- Next-word-prediction (GPT)

# 2

## Natural Language Inference (NLI)

The EU betrayed its partners during the negotiations on Sunday.

context-sentence

The EU is trustworthy.

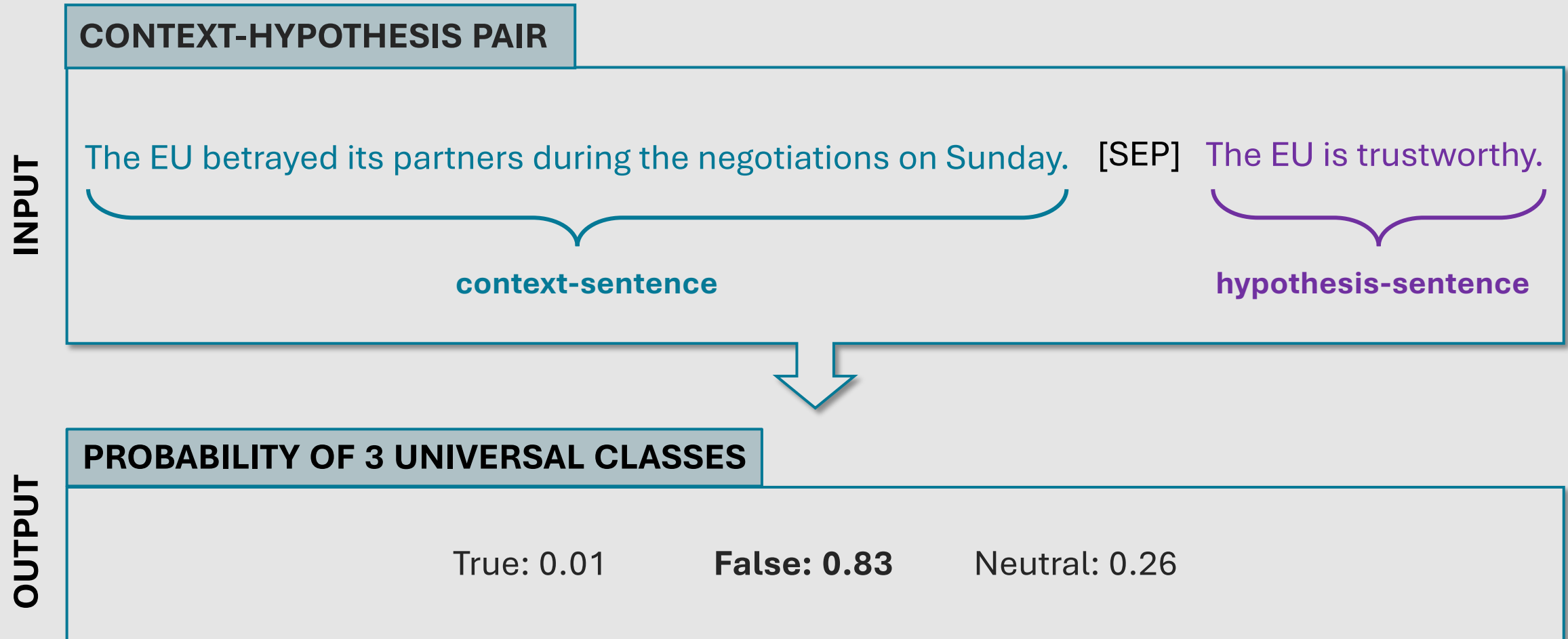
hypothesis-sentence

### NLI Task

Is the **hypothesis** true/false/neutral, given the **context-sentence**?

→ NLI is the task of determining whether a given statement (**hypothesis**) can logically be inferred from another statement (**context-sentence**).

# Natural Language Inference (NLI)



# Universal Task & Label Verbalisation

## Example task:

Identifying texts that indicate that the economy is performing well / badly

## Task reformulated for NLI:

### *NLI-input*

### *NLI-Output*

{context-sentence from news} [SEP] {hypothesis-sentence verbalising label}	Most "True" label
“The Rubel plummeted amidst a surge of investors withdrawing from Russia [SEP] The economy is performing badly”	<u>0,61 True</u> 0,13 False 0,26 Neutral
“The Rubel plummeted amidst a surge of investors withdrawing from Russia [SEP] The economy is performing well”	<u>0,19 True</u> 0,38 False 0,43 Neutral

# Universal Task & Label Verbalisation

*NLI-input*

*NLI-Output*

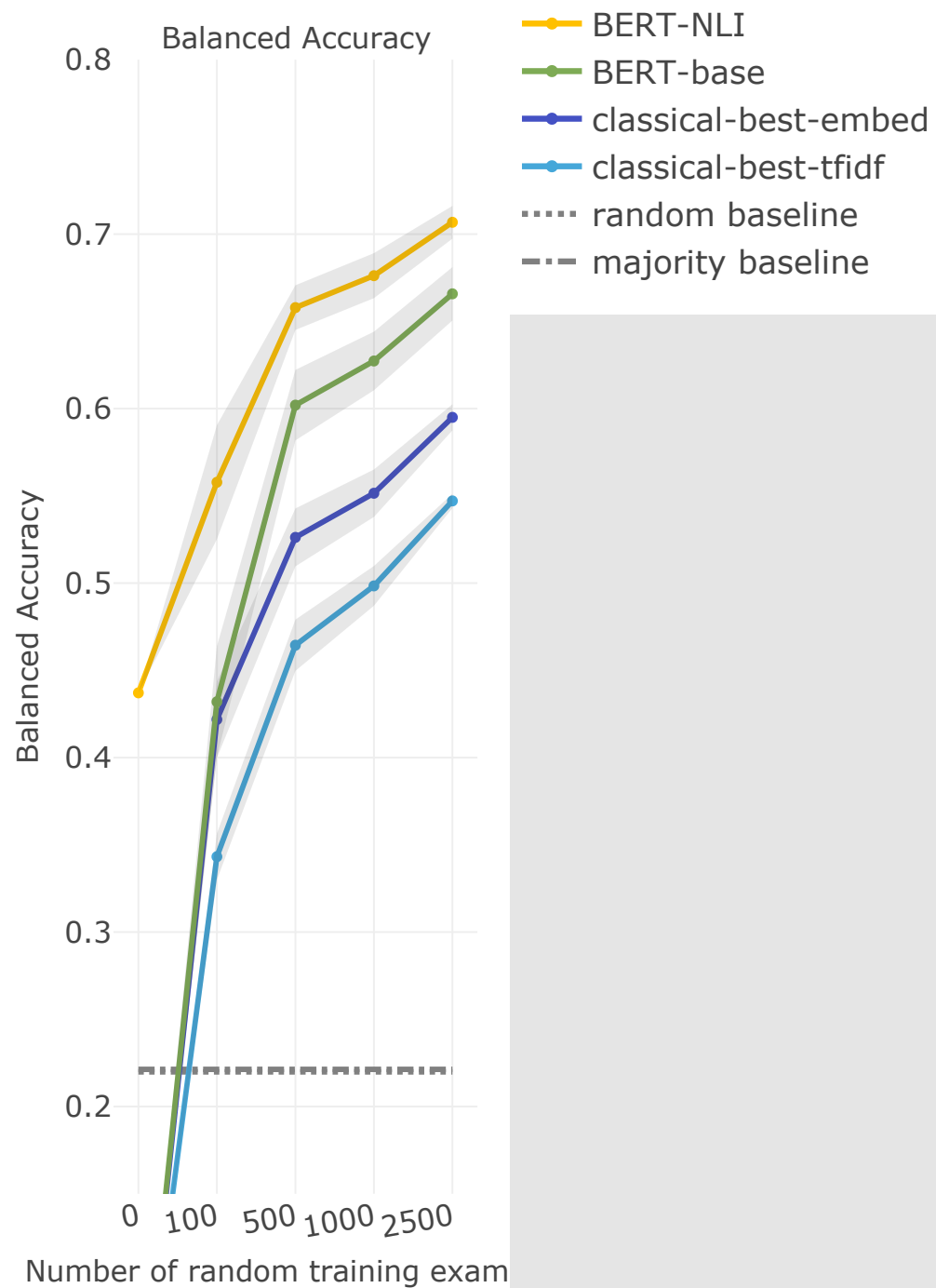
{context-sentence from news} [SEP] {hypothesis-sentence verbalising label}	Most "True" label
“The politicians were bribed by lobbyists. [SEP] It is about corruption.”	<u>0,61</u>
“The politicians were bribed by lobbyists. [SEP] It is about peace.”	0,01
“The politicians were bribed by lobbyists. [SEP] It is about free market.”	0,06
“The politicians were bribed by lobbyists. [SEP] It is about equality.”	0,04
...	...



# NLI: Data-Rich Task

## NLI is a data rich task

- Many NLI datasets with over 1 million annotated sentence pairs from different domains exist.
  - Examples: *SNLI* (570k examples, Bowman et al. 2015), *MultiNLI* (433k, Williams et al. 2018), *ANLI* (162k, Nie et al. 2020)
- **Helps address the issue of data scarcity**



**Average performance  
across eight tasks vs.  
training data size**

[Laurer et. al 2023](#)

# Limitations of NLI

- Usefulness decreases with training data size. If there is enough data to learn the new task ( $> 1000$  texts), BERT-base is better.
- BERT-NLI can only do classification tasks.
- No summarization, translation, information extraction ...

# A more universal task

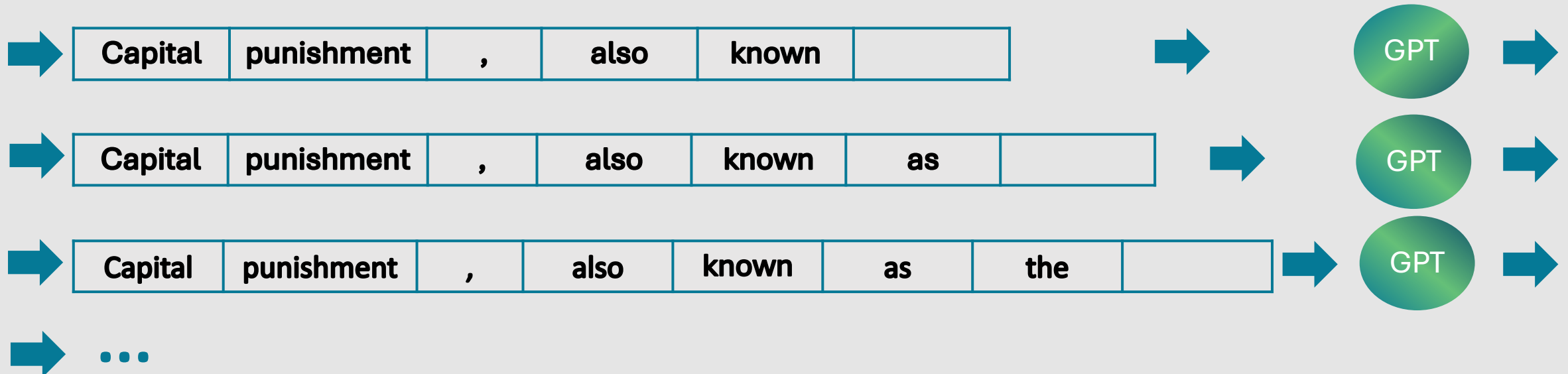
## Next-token-prediction

- Main pre-training task of GPT models
- Is a self-supervised task (no manual annotations)

# Next-token-prediction task

**Original text:**

*“Capital punishment, also known as the death penalty and formerly called judicial homicide, is ...”*



# Universal task

## Sentiment classification

[long text ...]	Is	this	text	positive	or
negative	?				



[long text ...]	Is	this	text	positive	or
negative	?	<b>positive</b>			

# Universal task

## Information extraction

[long text ...]	Extract	all	countries	from	the
text	:				



[long text ...]	Extract	all	countries	from	the
text	:	<b>Germany</b>			



[long text ...]	Extract	all	countries	from	the
text	:	<b>Germany</b>	,		




...


# Universal task

## Summarisation

[long text ...]	A	summary	of	the	preceding
text	:				



[long text ...]	A	summary	of	the	preceding
text	:	<b>The</b>			



[long text ...]	A	summary	of	the	preceding
text	:	<b>The</b>	<b>main</b>		



...





# Reflect and Q&A

- Q1: What is the difference between word embeddings and BERT?
- Q2: What are universal tasks and why are they useful?
- Q3: In your own words, try to define transfer learning.

**Write your responses on a piece of paper / notebook.  
Ask any questions about the slides in the chat.**