

Lecture-5: A closer look at GPT

1. Zero shot Vs Few shot learning

Zero-Shot: Ability to generalize to completely unseen tasks without any pair.

- The model predict the answer given only a natural language description of the task. No gradient updates are performed.

specific examples:

- Translate English to French:
 - Cheese \Rightarrow French
- task description.

One-Shot:

- In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

eg.

- Translate English to French:
 - Sea otter \Rightarrow loutre de mer
 - Cheese \leftarrow prompt
- task description
example

Few shot: Learning from a minimum number of examples which the user provides as input.

- In addition to the task description, the model see a few examples of the task. No gradient updates are performed.

eg.

- Translate English to French:
 - sea otter \leftarrow example
 - peppermint \leftarrow
 - plush giraffe
 - Cheese \leftarrow prompt
- task descriptor

Zero-shot Vs Few shot

	Input	Output
Complete task without example. ↳ Zero-shot	Translate English to French	petit - déjeuner
Complete task with a few examples. ↳ few-shot	goat → goat shoe → shoe phone →	phone

2. Utilizing Large Datasets

Let's us look at pretraining datasets of GPT-3 LLM.

Dataset name	Dataset description	Number of tokens	Proportion in training data
CommonCrawl (filtered)	Web crawl data	410 billion	60%
WebText2	Web crawl data	19 billion	22%
Books1	Internet-based book corpus	12 billion	8%
Books2	Internet-based book corpus	55 billion	8%
Wikipedia	High-quality text	3 billion	3%

Here take reference from GPT-3 main paper

A token is a unit of text which the model reads. For now, you can think of 1 token = 1 word.

* The total pre-training cost for GPT-3 \approx 4.6 Billions dollars.

* These pretrained models are base / foundational models which can be used for further finetuning.

* Many pretrained LLMs are available as open-source models \rightarrow can be used as general purpose tools to write, extract and edit text which was not part of the training data.

3. GPT Architecture

GPT \rightarrow Generative pretrained Transformer

Original
Paper
2018

\rightarrow GPT-3 is a scaled up version of this model, implemented on a larger dataset

* GPT models are simply trained on "next-word" prediction tasks.

The lion roams in the jungle.
next-word

* with this training, they can do a wide range of other tasks like translation, spelling correction etc.

* Next-word prediction: self-supervised learning
 \downarrow
Self-labelling

* we don't collect labels for training data, but use the structure of of the data itself.

next word in sentence is used as label

Auto regressive model.
use previous outputs as input for future prediction

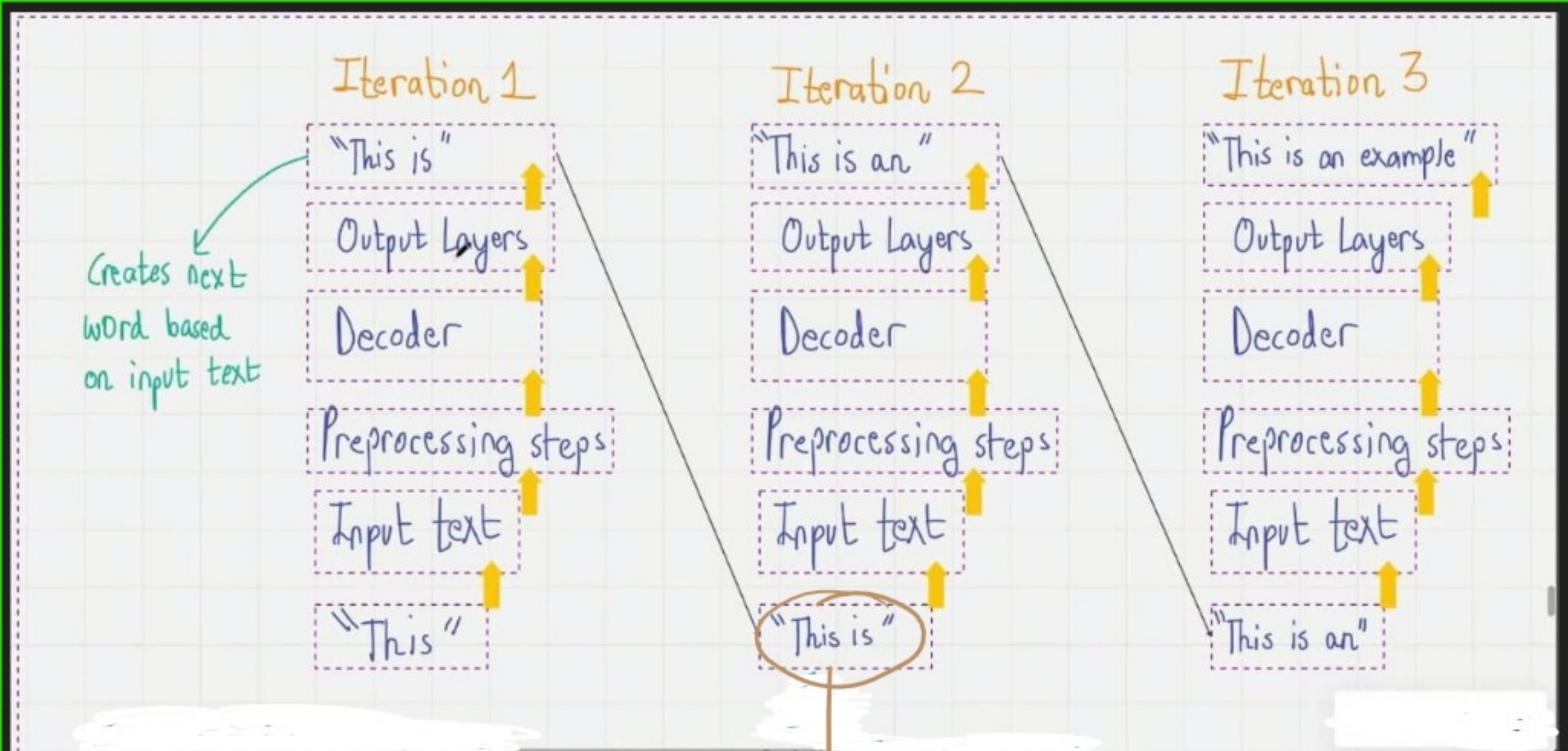
*** pretraining part of GPT model
 \downarrow
Supervised
 \downarrow
Auto-regressive model

* compared to original transformer architecture GPT architecture is simpler.

* In GPT architecture: there is no encoder we just have the decoder.

original transformer = 6 encoder-decoder block

GPT-3: 96 transformer layers,
175 billions parameters



GPT Architecture (decoder only)

Output of previous round serve as input to next round.

Although trained only for next word prediction GPT model can perform other tasks like language translation.

This is called "emergent
behaviour"

Ability of a model to perform tasks that the model wasn't explicitly trained to perform.