## (5) Transformers vs LLMs

o Not all transformers are LLMs

o Transformers can also be used for computer vision

vision transformer → can be used for image classification, like CNN

o Not all LLMs are transformers
a can be LSTM or recurrent NN

---

### Closer look at GPT

① Zero-shot vs few shot learning

↳ Zero-shot = Generalize to completely unseen tasks w/ out prior specific examples

o Generative pre-training (GPT)
  o text used is not labelled
  o Dont need to provide labels
  o Generating next word
  o Transformers + unsupervised pretraining

o Few-shot → Learning from a minimum number of examples which user provides as input

(OpenAI)

o GPT 3  175B parameters

o

zero-shot: model predicts answer given only a natural
language description of task. No gradient updates are performed

Translate english to french (desc)
cheese ->                        GPT-4 can do
                    (prompt)     zero-shot too

__One-shot__: in addition to task description, model sees a
single example of the task. No gradient updates are performed

Translate English -> French: (task desc)
  Sea otter => Loutre de mer (example)
  cheese ->              (prompt)

Few-Shot: sees a few examples               GPT-4 is better w/
  sea otter -> loutre de mer                  __few shot__
  peppermint -> menthe poivrée  } examples
  plush giraffe -> girafe peluche   GPT-3 is few shot learner
  cheese ->        give few examples for translation

Token -> unit of text which model reads forming continue of token/word

model pre-training cost for GPT-3 is $4.6 million

* pretrained models are base/foundational models which can be
used for further finetuning
* many pretrained LLM's are available as open-source models ->
can be used as general purpose tools to write, extract and
edit text which was not part of training data

③ GPT Architecture (only decoder block, no encoder)

o simply trained on 'next word' prediction tasks
  ↳ the lion roams in the __jungle__
                            next word

* w/ this training, they can do a wide card of other tasks like
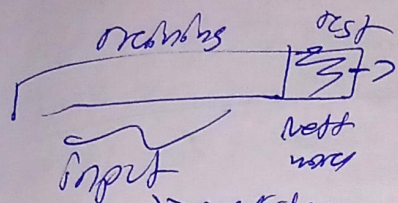translation, spelling correction etc.      (output)(labels)
ex: second law of robotics:                 a
  second law of robotics: a                 robot
  second law of robotics: a robot           must

* Next word prediction; Self-supervised learning
  ↳ self-labelling



* we don't collect labels for training data, but use structure of data itself
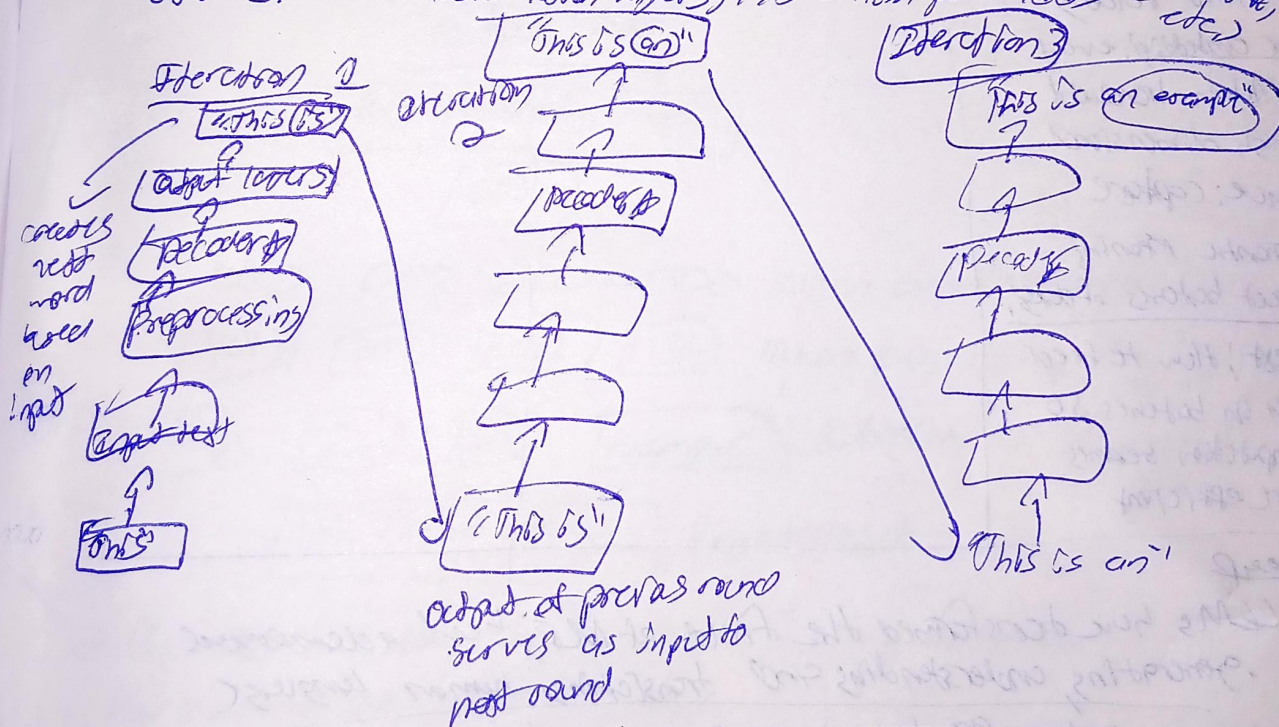  ↳ Next word / sentence is used as label ∴ Auto regressive model;
  use previous outputs as inputs for future predictions

① Pretraining
  ↓
unsupervised content
  observed)
Auto Regressive models

* compared to original transformer architecture, GPT arch is simpler

* There is no encoder, we just take the decoder

Original Transformer: 6 encoder-decoder blocks

GPT-3: 96 Transformation layers, 175 billion parameters (decoder, etc.)    96 transformer layers



Iteration 1
"This is"
[output layers]
[Decoder]
[Preprocessing]
[input text]
"This"
predicts next word based on input

attention
"This is an"
[Decoder]
"This is"
output of previous round serves as input to next round

Iteration 3
"This is an example"
[Decoder]
"This is an"

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

"Emergent Behavior" → ability of a model to perform tasks that
the model wasn't explicitly trained to perform
How to explain? Still isn't known
    ○ Good research topic

# Stages of Building LLM

## stage 1



| data prep and sampling | Attention mechanism | LLM arch |

↓ Building an LLM

Implement data sampling;
understand basic
mechanism

tokenization; break
text into tokens

vector embedded word

## stage 2

| Training Loop | model eval. | load pretrained weights |

↓ Foundation Model

pretrain on LLM on
unlabeled data

## stage 3

Finetuning → Classifier

Finetuning → personal assistant

Pretrained pretrained LLM