

Natural language processing & language modelling



Content



How do we model:

- Language
- Words
- Sequences
- Context

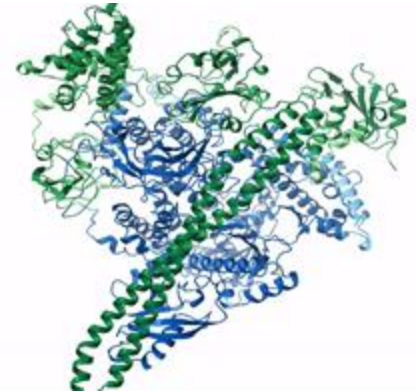
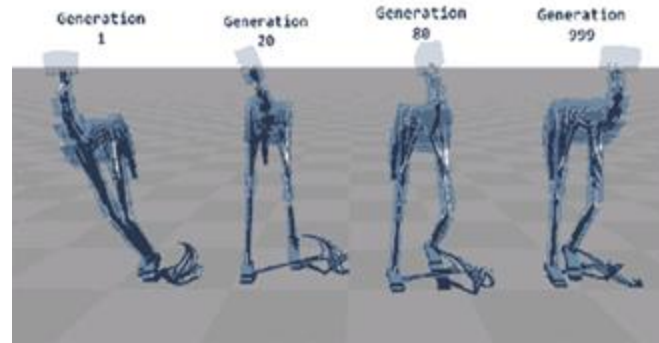
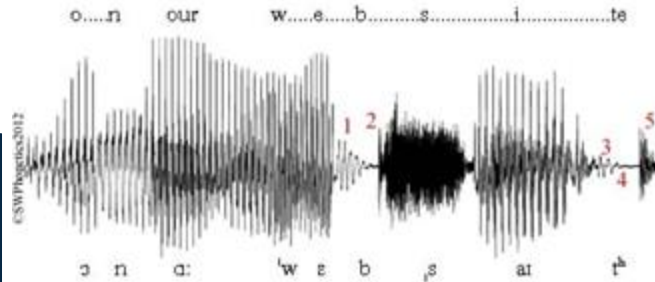
Transformers

LLMs



Sequential data

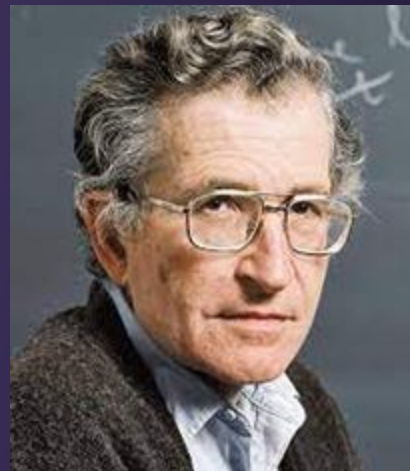
Universe seems to be inherently sequential



How can we model this?

Is there one (sequential) model to rule them all?

How does modelling language differ from other sequences?



“

*Language is a process of free creation; its laws and principles are **fixed**, but the manner in which the principles of generation are used is free and **infinitely varied**.*

~ Noam Chomsky

Sequential data

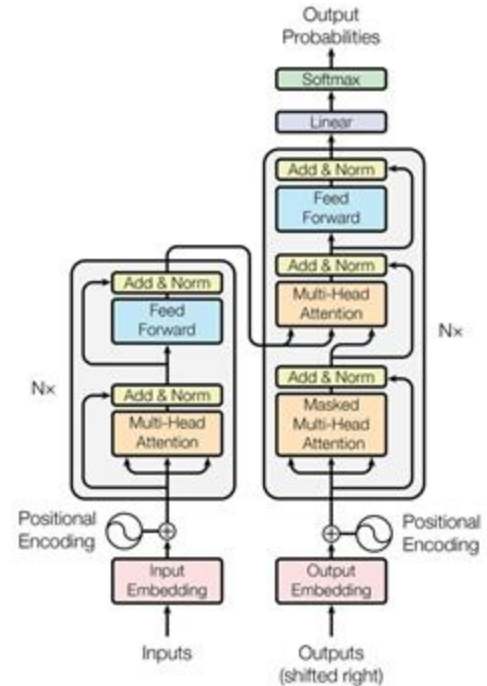
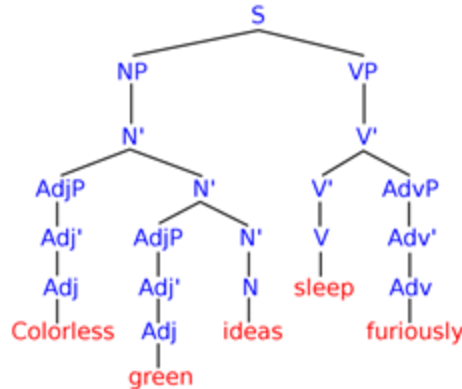
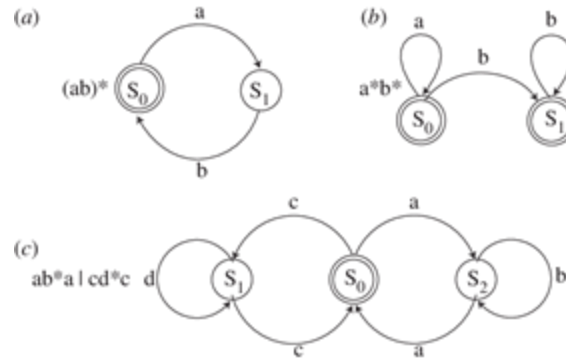
Modelling *natural language*

We have the innate ability for language (Universal Grammar)
Syntactic Structures

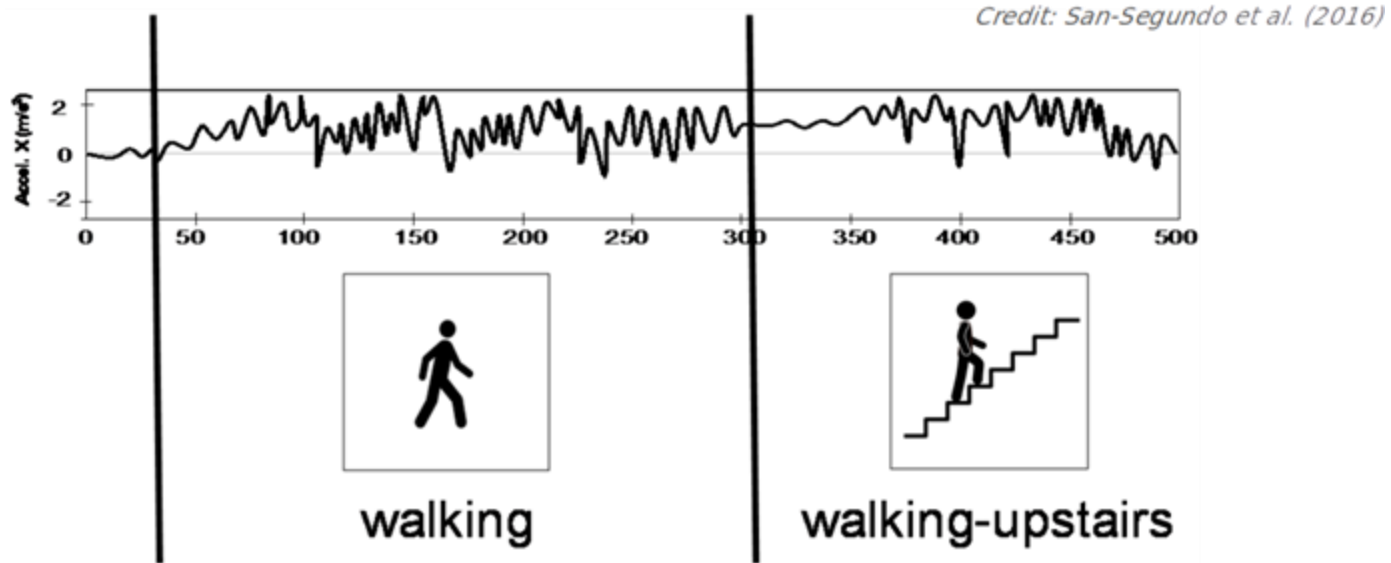
(Chomsky, N. 1957)

Different ways to model language:

- Finite State Automata
- Generative Grammars
- Deep Learning + Data
-



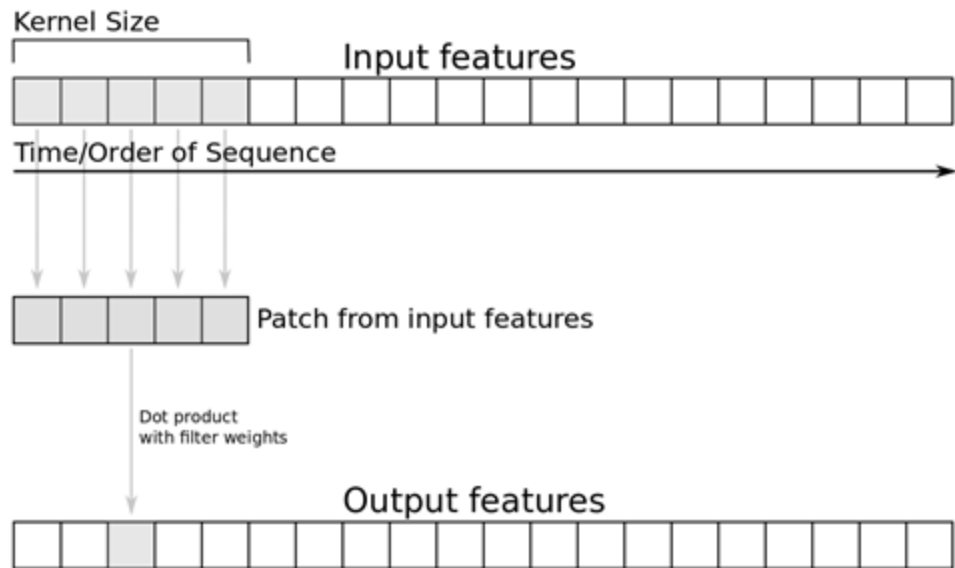
Human Activity Recognition



Human Activity Recognition

Can we use a CNN?

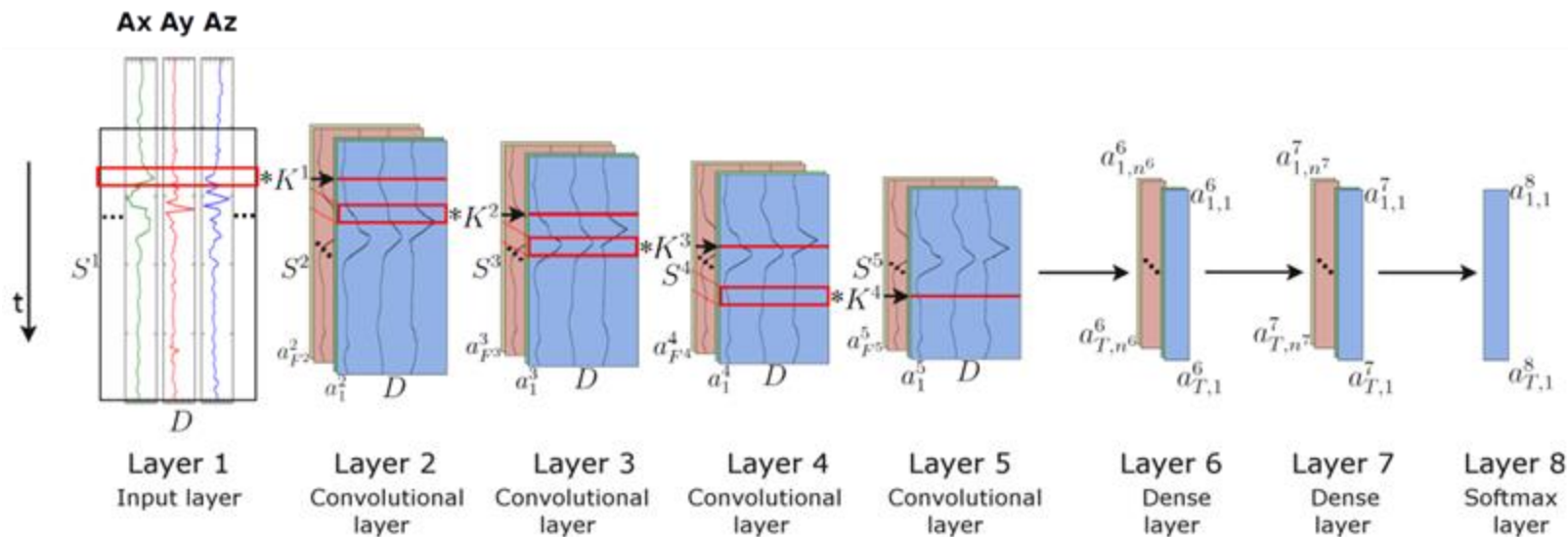
CNNs work here because they take **nearby** tokens into account



Human Activity Recognition



Can we use a CNN?

CNNs work here because they take **nearby** tokens into account



How to model language?



Language model objective:

-  Model the **probability** of the next symbol given the previous context by creating a probability distribution over language
-  Language models like ChatGPT model language *left to right*: one word at a time

$$P(w_1, w_2, \dots w_n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1, w_2) \dots P(w_n \mid w_1, \dots w_{n-1})$$

How to model language?

Language model objective:

-  Model the **probability** of the next symbol given the previous context by creating a probability distribution over language
-  Language models like ChatGPT model language *left to right*: one word at a time

$$P(w_1, w_2, \dots w_n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1, w_2) \dots P(w_n \mid w_1, \dots w_{n-1})$$

The cat

How to model language?

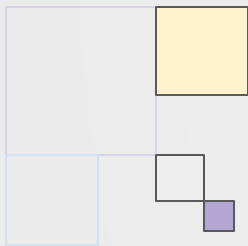
Use CNNs for modelling language?

CNN's context size is small and does **not remember** past context

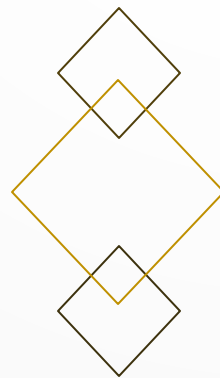
Language is very noisy and heavily (long-distance) context dependent

$$P(w_1, w_2, \dots w_n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1, w_2) \dots P(w_n \mid w_1, \dots w_{n-1})$$

The cat



Neural representations



How do we represent words?

Language is **discrete**!

“

One cat ate a tuna, one
cat ate a chicken, one
weird cat ate a human.

How do we represent words?

Language is **discrete**!

“

One cat ate a tuna, one
cat ate a chicken, one
weird cat ate a human.

Vocabulary

A vocabulary is the set of unique
tokens in your language (dataset)

How do we represent words?

Language is **discrete**!

“

One cat ate a tuna, one
cat ate a chicken, one
weird cat ate a human.

Vocabulary

A vocabulary is the set of unique
tokens in your language (dataset)

one
cat
ate
tuna
chicken
weird
human
a

How do we represent words?

Language is **discrete**!

“

One cat ate a tuna, one
cat ate a chicken, one
weird cat ate a human.

Vocabulary

A vocabulary is the set of unique
tokens in your language (dataset)

one
cat
ate
tuna
chicken
weird
human
a

$\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \end{bmatrix}$

How do we represent words?

Language is **discrete**!

“

One cat ate a tuna, one
cat ate a chicken, one
weird cat ate a human.

Vocabulary

A vocabulary is the set of unique
tokens in your language (dataset)

one
cat
ate
tuna
chicken
weird
human
a

$\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \end{bmatrix}$

● ————— ● $[0.122, 1.54, -0.193]$

How do we represent words?

Language is **discrete**!

“

One cat ate a tuna, one
cat ate a chicken, one
weird cat ate a human.

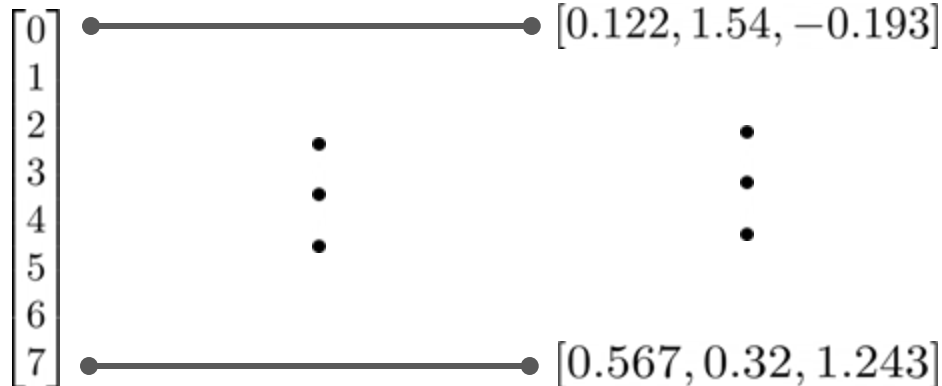
Easy to **compose** and easy to
vectorize

Vocabulary

A vocabulary is the set of unique
tokens in your language (dataset)

Embedding

one
cat
ate
tuna
chicken
weird
human
a



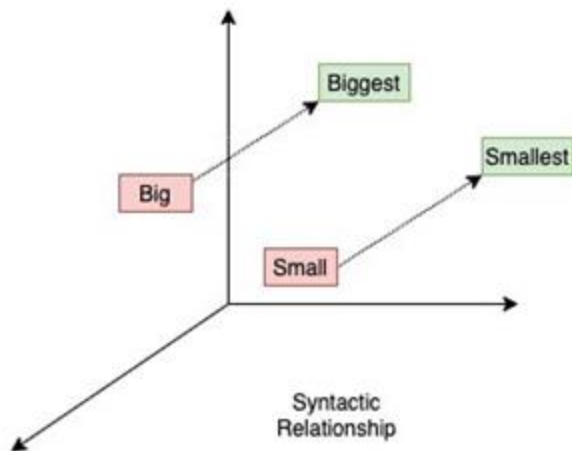
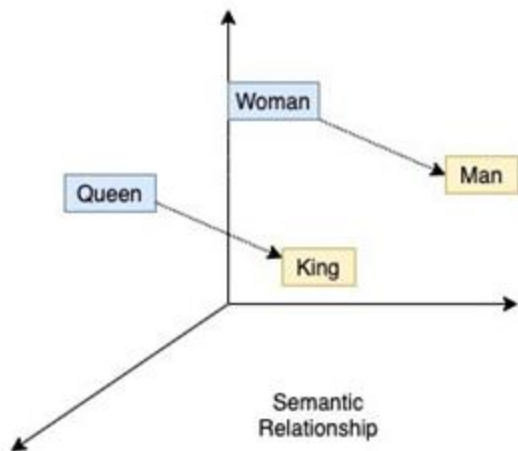
How do we represent words?

- Most basic: one hot encoding $[0, 0, 1, 0, \dots, 0, 0]$
✗ No meaningful distance between vectors
- Word2Vec

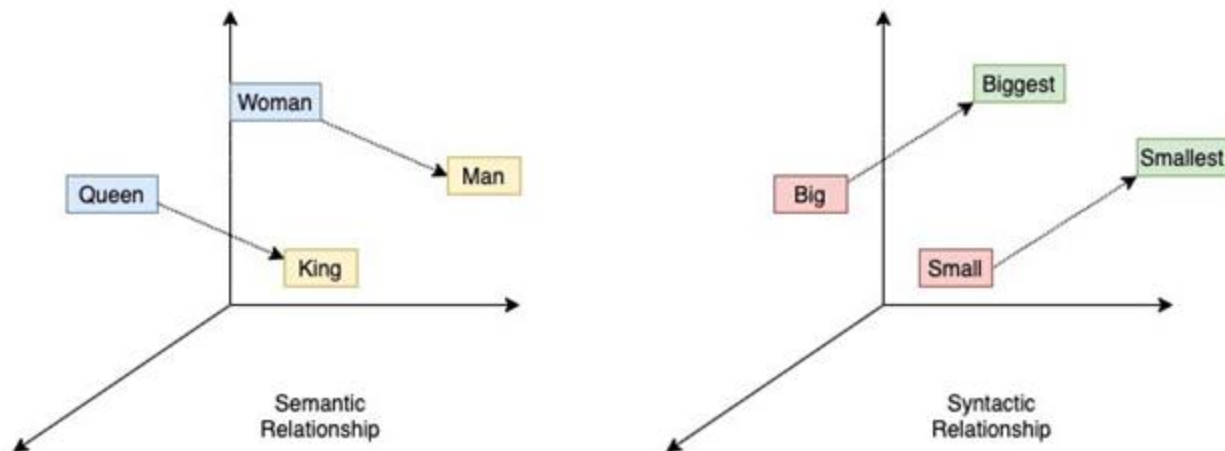
Distributional hypothesis

- The **context** surrounding a given word provides information about its meaning
- Words are similar if they share similar **contexts** (semantic + syntactic)
- Semantic similarity \approx distributional similarity

How do we represent words?



How do we represent words?



$$\text{King} - \text{Man} + \text{Woman} = \text{Queen}$$

Each word has a numerical vector representation learnt by a language model

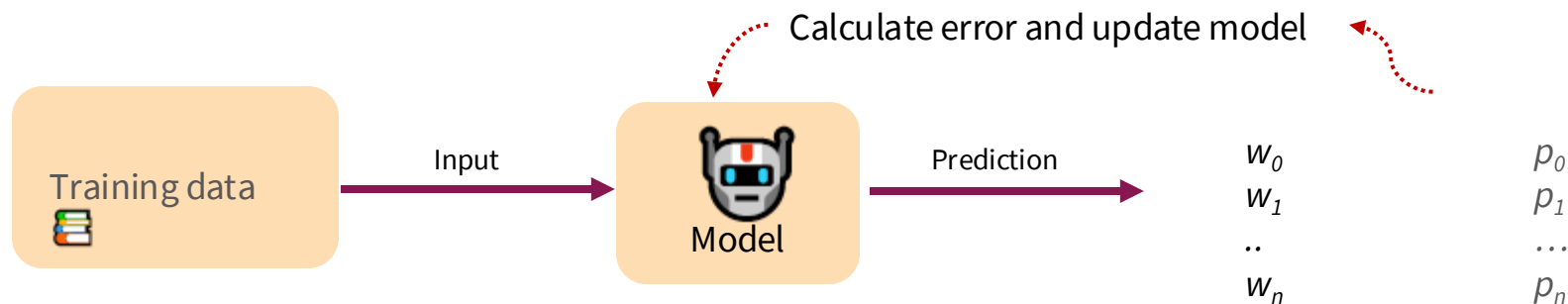
How do we represent words?

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Mikolov et al. 2013. *Efficient Estimation of Word Representations in Vector Space*

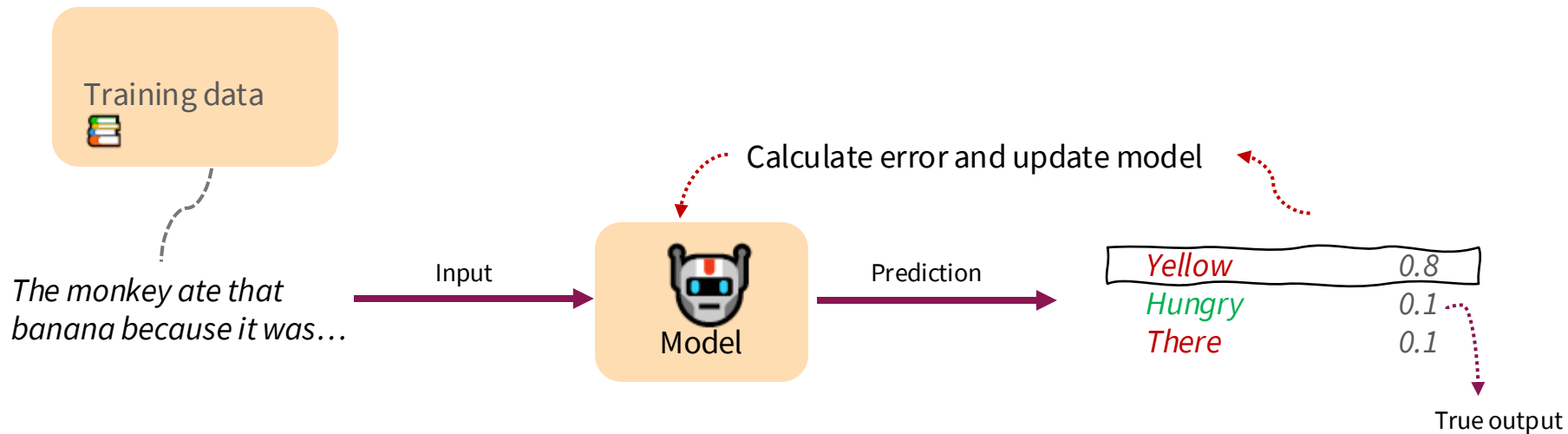
How do we model sequences?

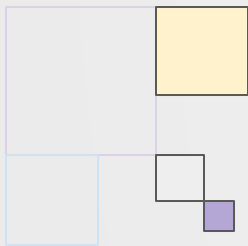
Unsupervised learning



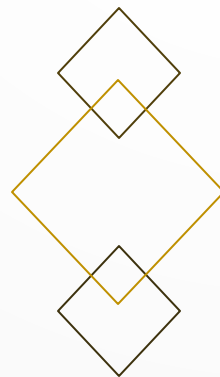
How do we model sequences?

Unsupervised learning





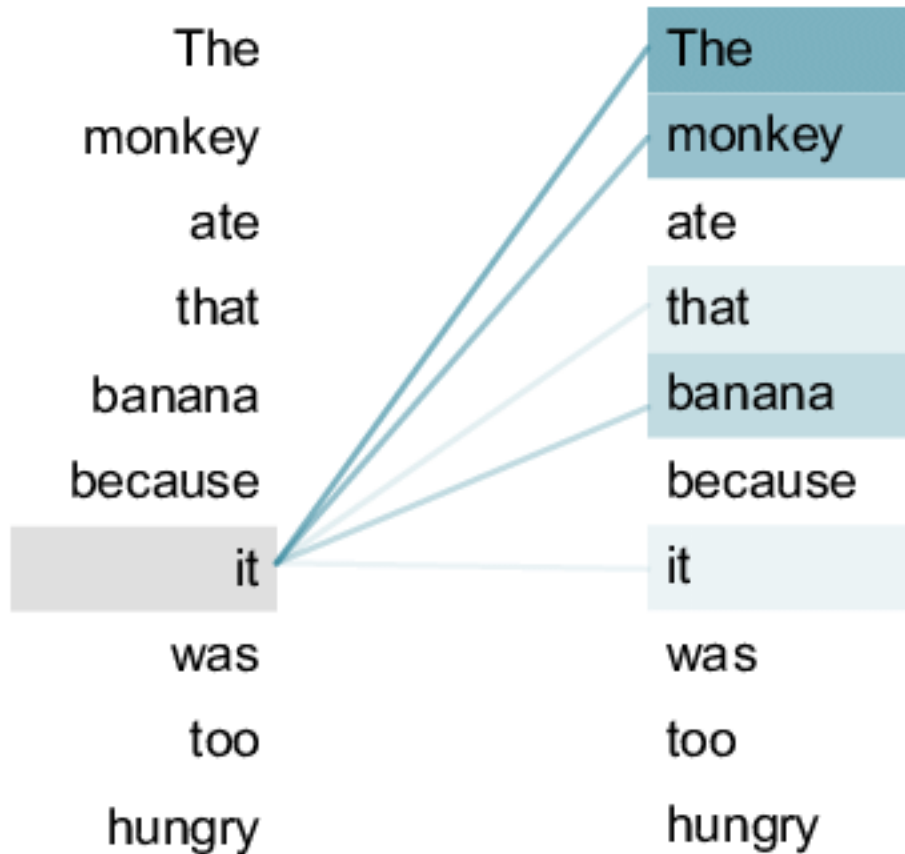
The Transformer



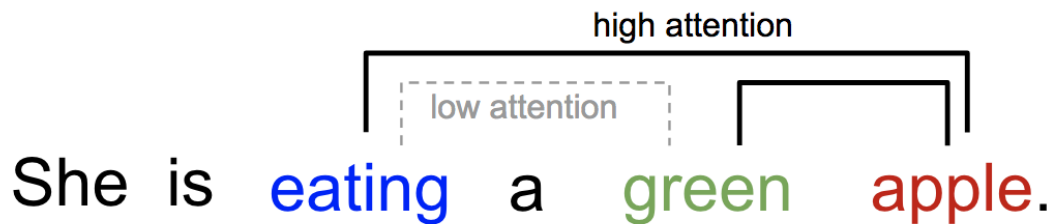
High Performance Machine Learning Group

Attention

Weighted average over
hidden states



Attention is all you need: self-attention

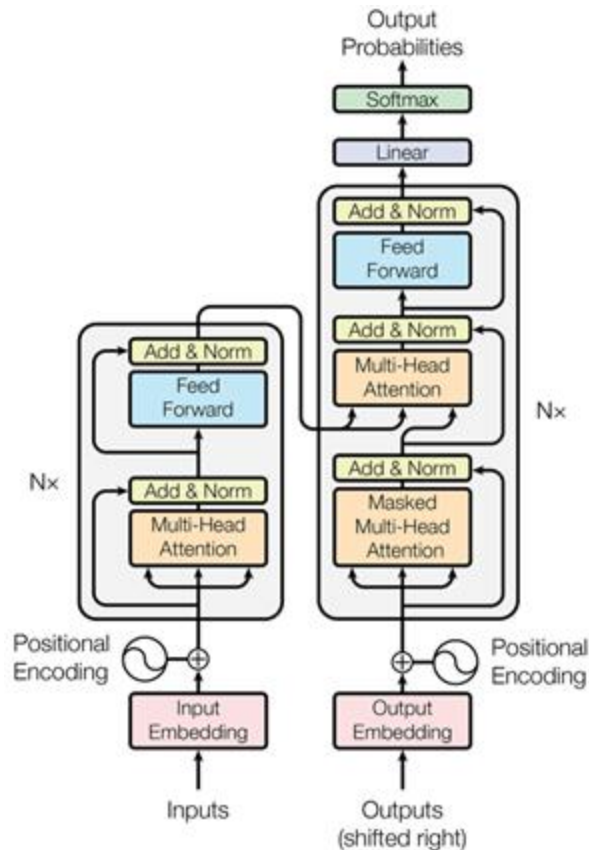


Transformers

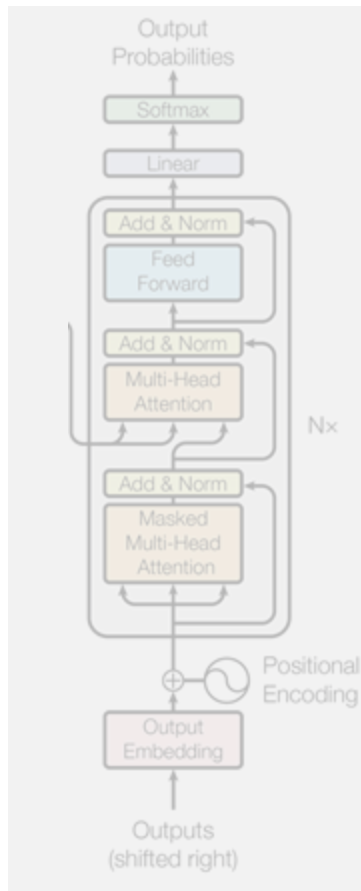
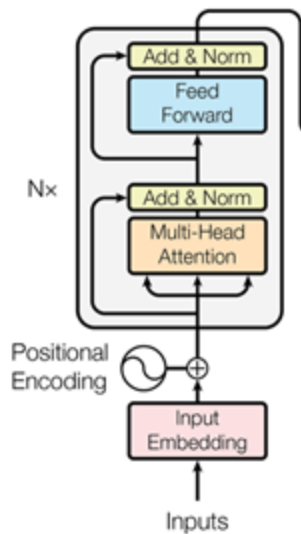
What is the **importance** of every element in the input sequence with respect to **itself**?

Every word directly connected to every other word

More parallelization!



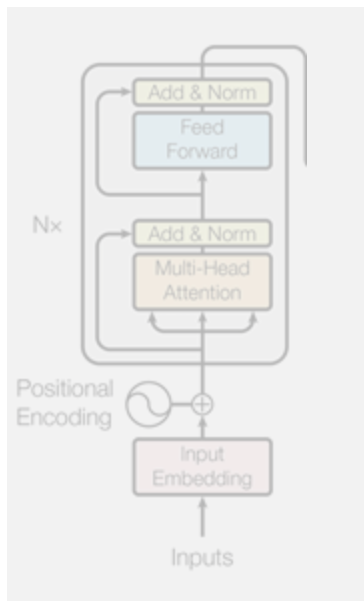
Dissection of the Transformer



Encoder

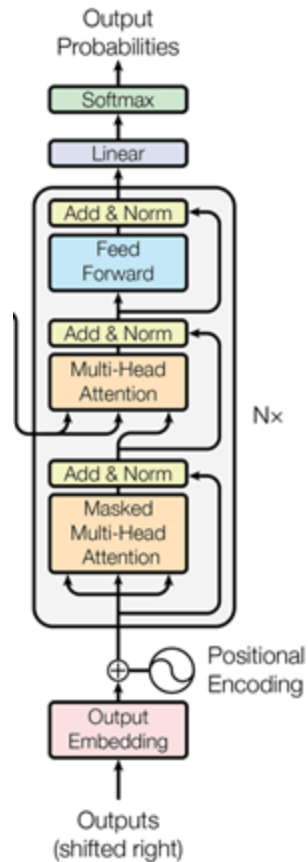
Understanding the input

Dissection of the Transformer

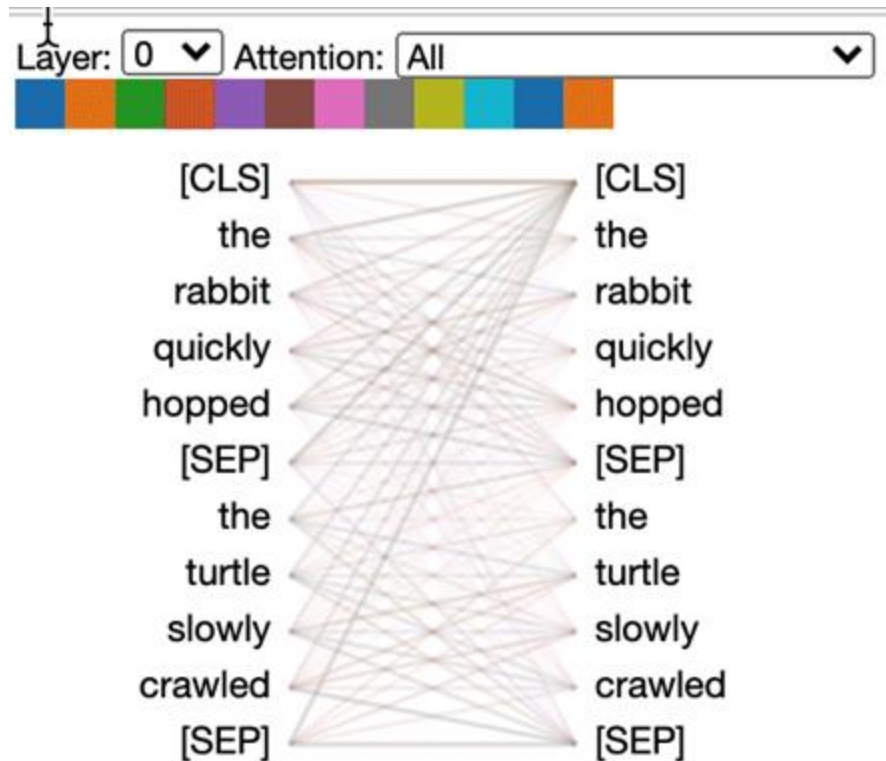
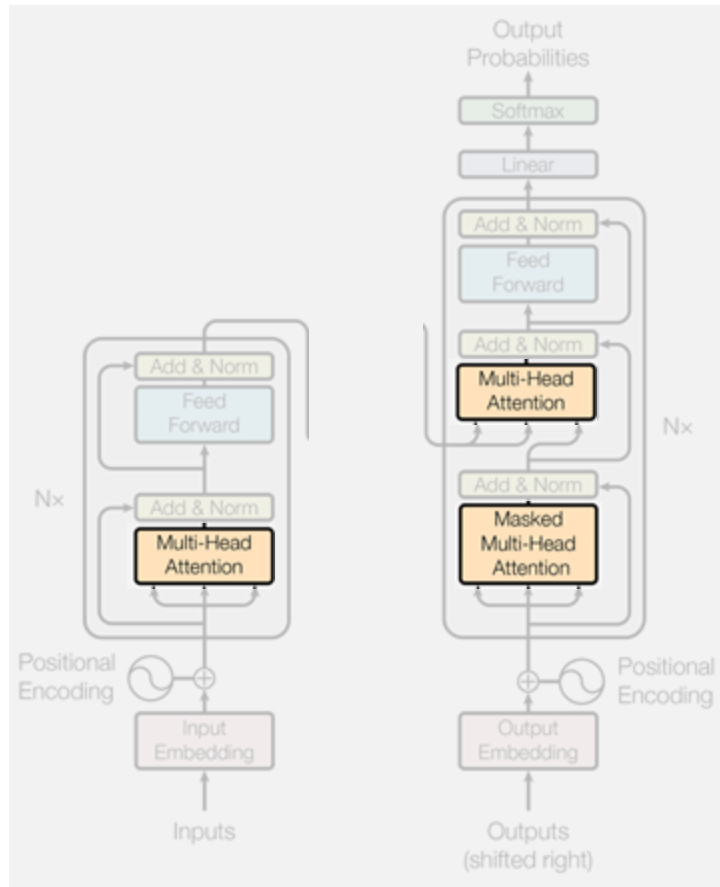


Decoder

Generating the output



Dissection of the Transformer



Self Attention

Input

Embedding

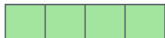
Queries

Keys

Values

Score

Thinking

x_1 

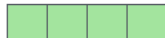
q_1 

k_1 

v_1 

$q_1 \cdot k_1 = 112$

Machines

x_2 

q_2 

k_2 

v_2 

$q_1 \cdot k_2 = 96$

Self Attention

Input

Embedding

Queries

Keys

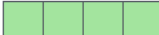
Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Thinking

x_1 

q_1 

k_1 

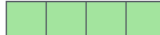
v_1 

$$q_1 \cdot k_1 = 112$$

14

0.88

Machines

x_2 

q_2 

k_2 

v_2 

$$q_1 \cdot k_2 = 96$$

12

0.12

Self Attention

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax

X

Value

Sum

Thinking

x_1 

q_1 

k_1 

v_1 

$q_1 \cdot k_1 = 112$

14

0.88

v_1 

z_1 

Machines

x_2 

q_2 

k_2 

v_2 

$q_1 \cdot k_2 = 96$

12

0.12

v_2 

z_2 

Self Attention

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$
$$= \begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Self Attention

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix}$$

$$\frac{e^{x_i}}{\sum_{j=0}^K e^{x_j}}$$

Self Attention

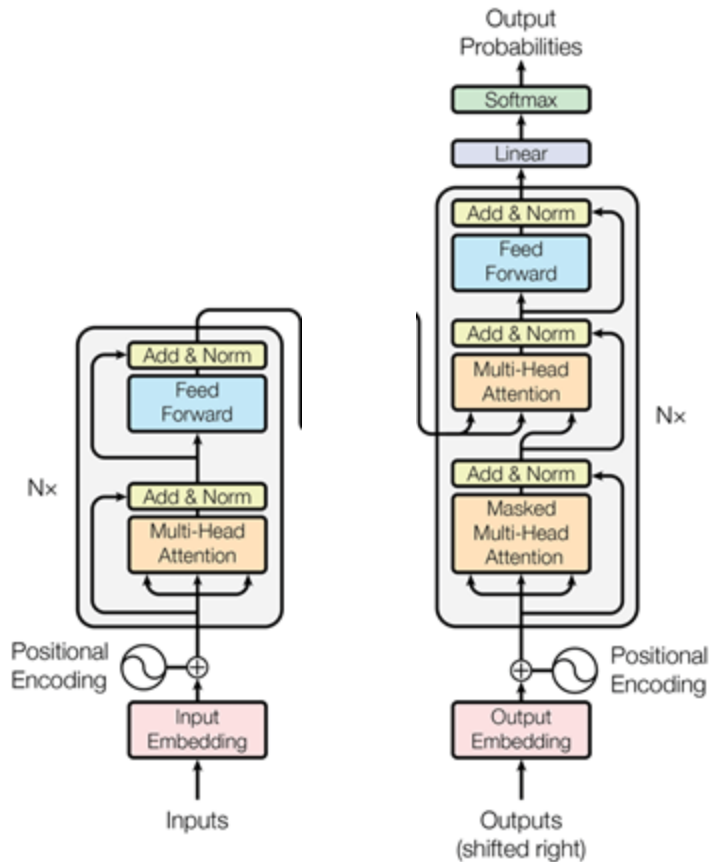
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix}$$

$$\frac{e^{x_i}}{\sum_{j=0}^K e^{x_j}}$$

$$\begin{bmatrix} 0.02 \\ 0.90 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$

Dissection of the Transformer



Layer Normalization

Skip Connections

Positional Encoding

Byte Pair Encodings

Layer Stacking

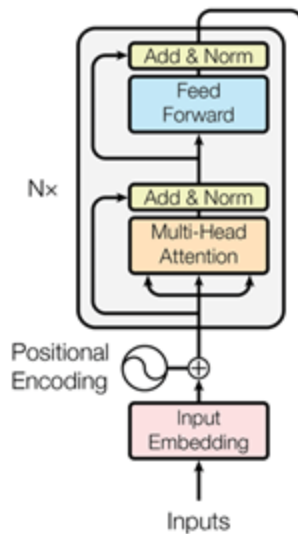
Cross Attention

Model types

Encoder

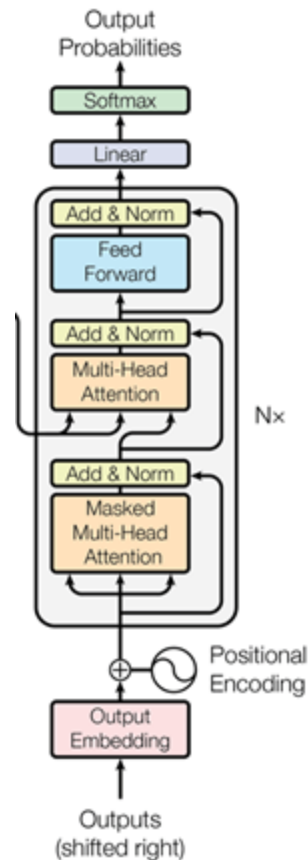
BERT

Sequence classification



Encoder-decoder: T5

Sequence-to-sequence e.g.
machine translation

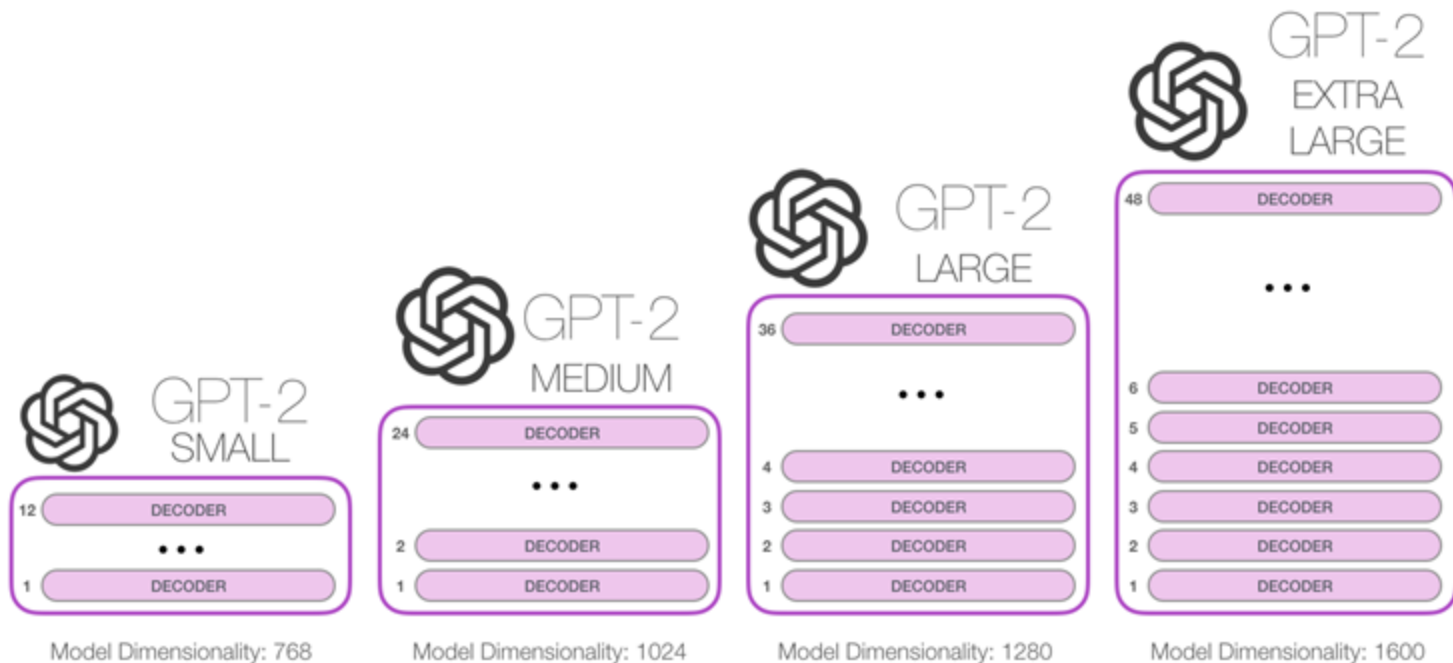


Decoder

GPT

Text generation

Generative Pre-Trained Transformers



From language models to chatbots

GPT: language model

ChatGPT: *aligned* language model

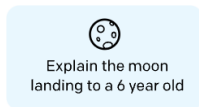
Alignment: what type of response is desirable? How can we make model behaviour be *aligned* with human norms and values?

Alignment and instruction finetuning

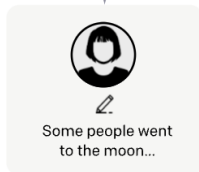
Step 1

**Collect demonstration data,
and train a supervised policy.**

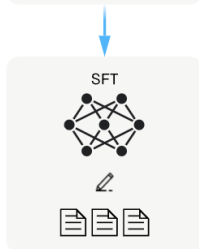
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.



This data is used
to fine-tune GPT-3
with supervised
learning.

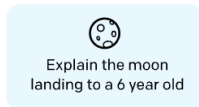


Alignment and instruction finetuning

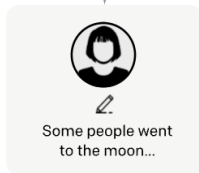
Step 1

**Collect demonstration data,
and train a supervised policy.**

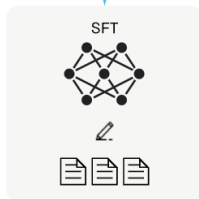
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.



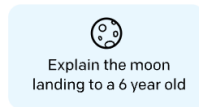
This data is used
to fine-tune GPT-3
with supervised
learning.



Step 2

**Collect comparison data,
and train a reward model.**

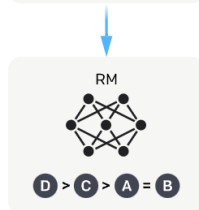
A prompt and
several model
outputs are
sampled.



A labeler
ranks the
outputs from
best to worst.



This data is used
to train our
reward model.

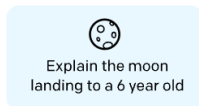


Alignment and instruction finetuning

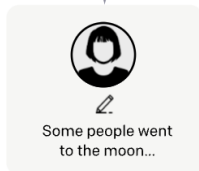
Step 1

Collect demonstration data, and train a supervised policy.

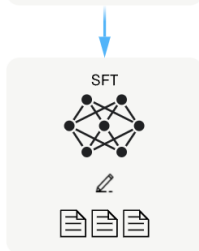
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



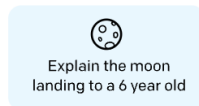
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

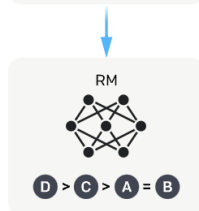
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



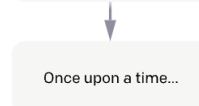
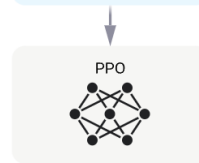
Step 3

Optimize a policy against the reward model using reinforcement learning.

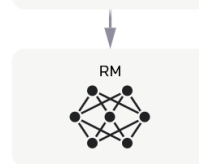
A new prompt is sampled from the dataset.



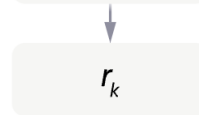
The policy generates an output.

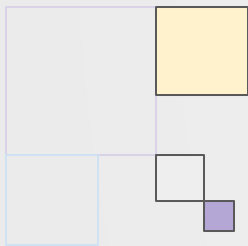


The reward model calculates a reward for the output.

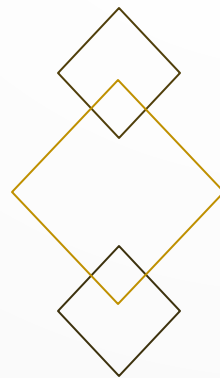


The reward is used to update the policy using PPO.

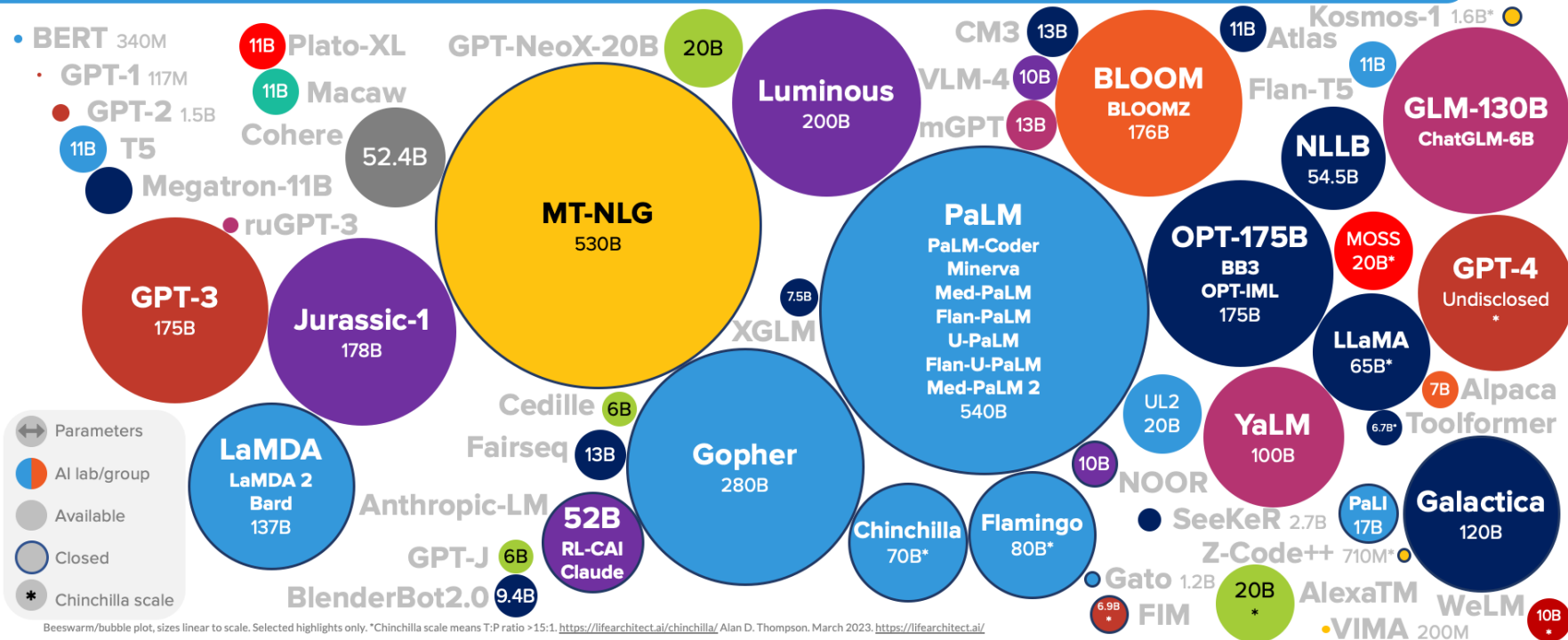




Scaling Transformers



Cambrian Explosion of LMs

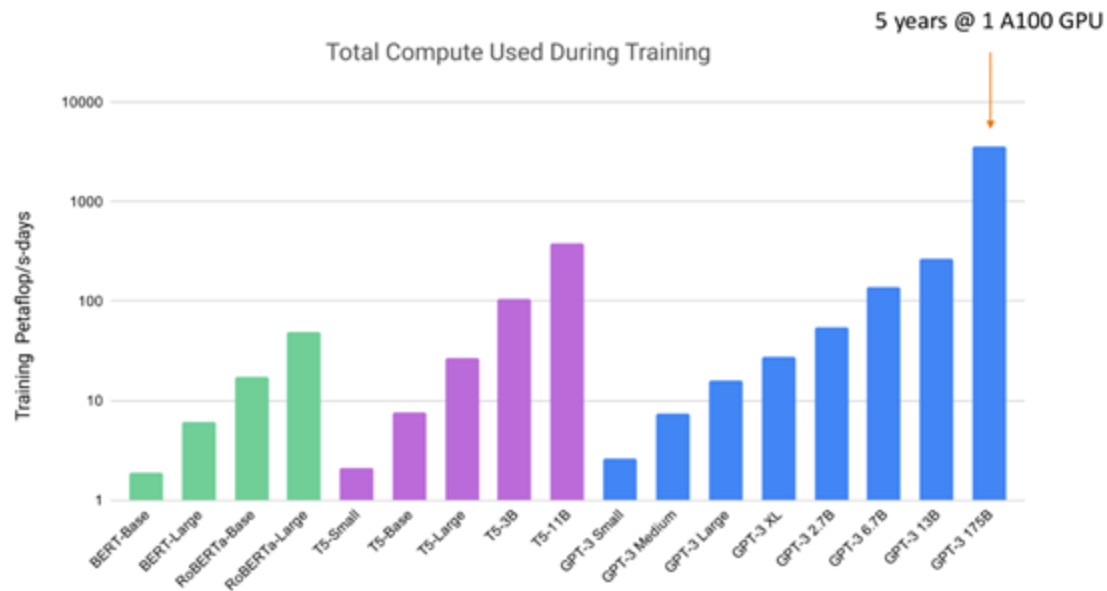


In comparison: the “big” model in original attention paper had only 213M parameters

Transformer Explosion

GPT-3 Trained on ~350 Billion tokens

GPT-3 has 175B parameters!



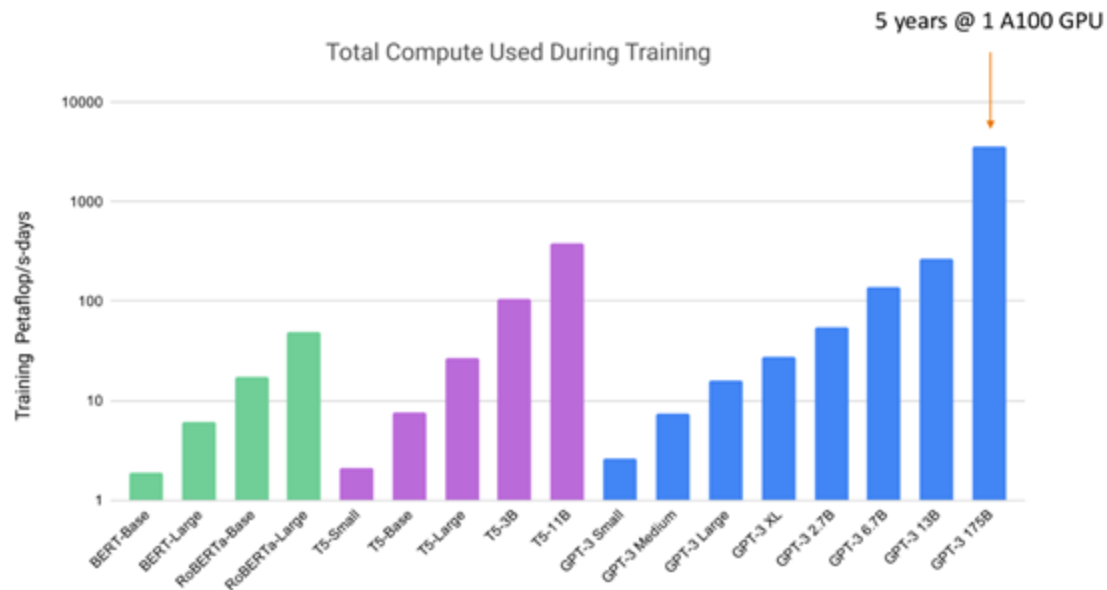
Transformer Explosion

GPT-3 Trained on ~350 Billion tokens

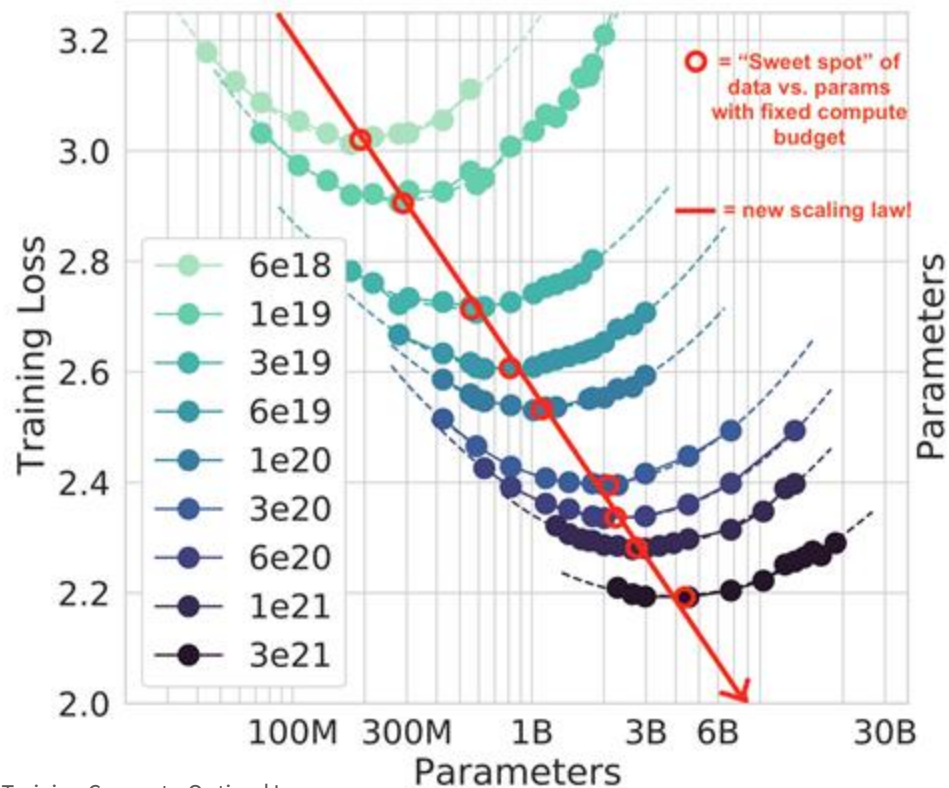
GPT-3 has 175B parameters!

GPT4: in the range of
1.8T parameters, trained on >1T tokens.

3 months at 20.000 A100s...

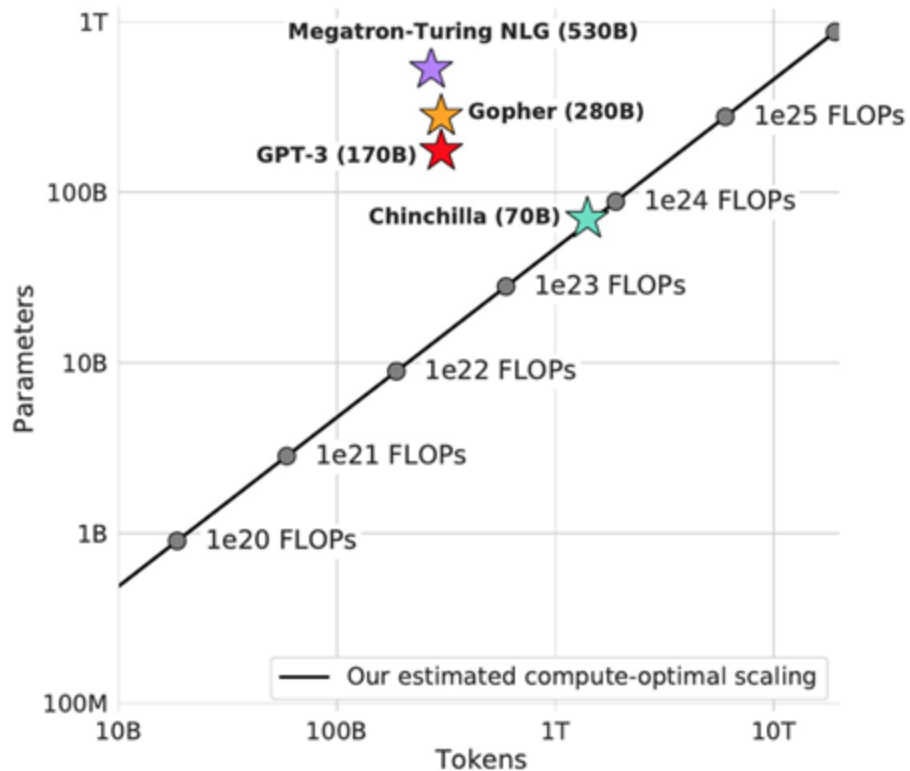


Chinchilla Scaling Laws

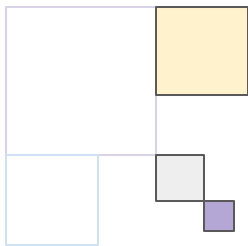


Hoffmann, Jordan et al. "Training Compute-Optimal Large Language Models." *ArXiv abs/2203.15556* (2022): n. pag.

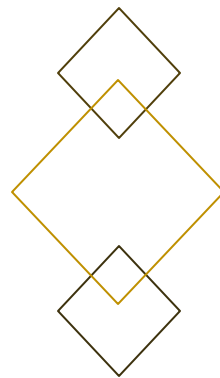
Chinchilla Scaling Laws



Hoffmann, Jordan et al. "Training Compute-Optimal Large Language Models." *ArXiv abs/2203.15556* (2022): n. pag.

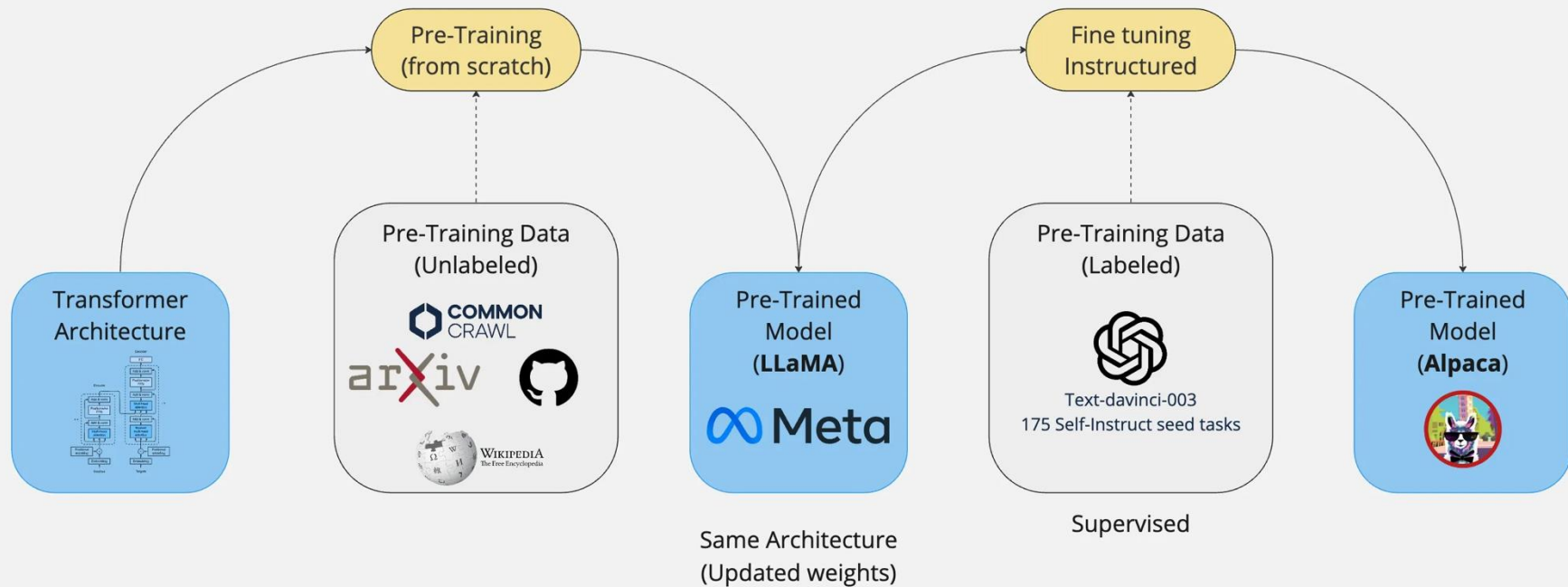


Finetuning LLMs at scale



High Performance Machine Learning Group

From training to finetuning



From training to finetuning

Current era of LLMS: few parties with enough data and compute to train large models

What can you do without requiring 1000 H100 GPUs? 😊

From training to finetuning

Current era of LLMS: few parties with enough data and compute to train large models

What can you do without requiring 1000 H100 GPUs? 😊

No training

Use foundation models as they are:

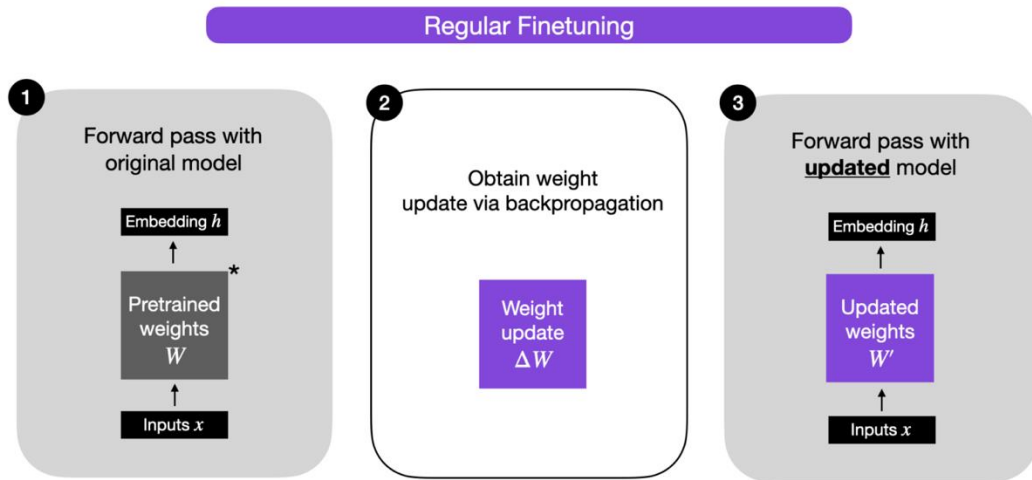
- In-context learning
- Retrieval augmented generation

Finetuning

Transfer learning: finetune model to adapt to specific task

Efficient parameter optimization

Often not necessary to change **all** parameters

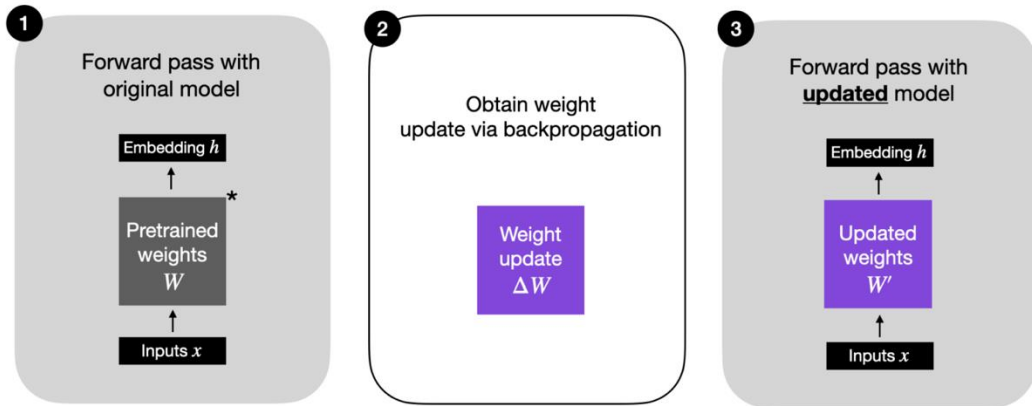


* The pretrained model could be any LLM, e.g., an encoder-style LLM (like BERT) or a generative decoder-style LLM (like GPT)

Efficient parameter optimization

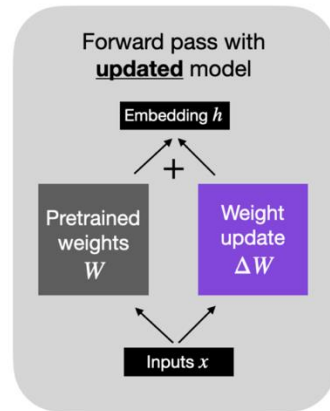
Often not necessary to change **all** parameters

Regular Finetuning



* The pretrained model could be any LLM, e.g., an encoder-style LLM (like BERT) or a generative decoder-style LLM (like GPT)

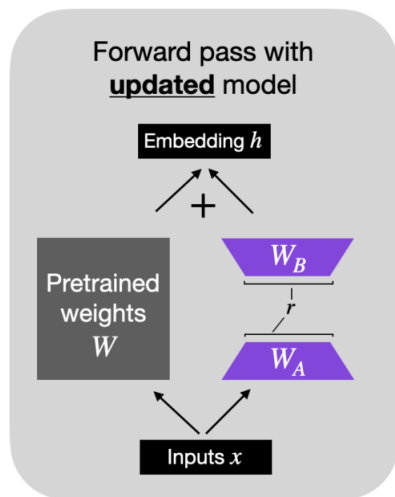
Alternative formulation (regular finetuning)



Efficient parameter optimization

Often not necessary to change **all** parameters

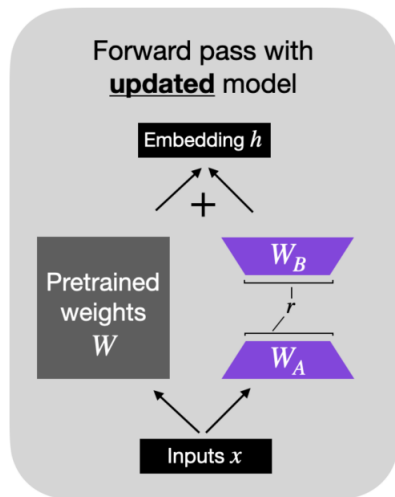
LoRA weights, W_A and W_B , represent ΔW



Efficient parameter optimization

Often not necessary to change **all** parameters

LoRA weights, W_A and W_B , represent ΔW



Reduces parameters in
vram by about 4x!

Example:

LLama 7B in FP:
4 (bytes) x 7 = 28GB
(inference)

8 (bytes) x 7 = 56GB
(training)

This can go down to 14GB
with PEFT!

Considerations



Social bias



Legal limitations



Trust and reliability

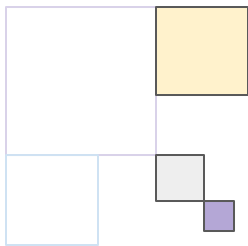
Confident talking parrot



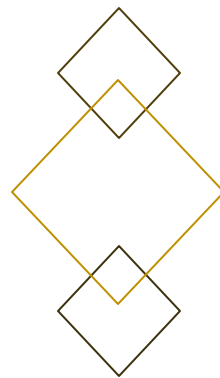
Disinformation



Environmental impact



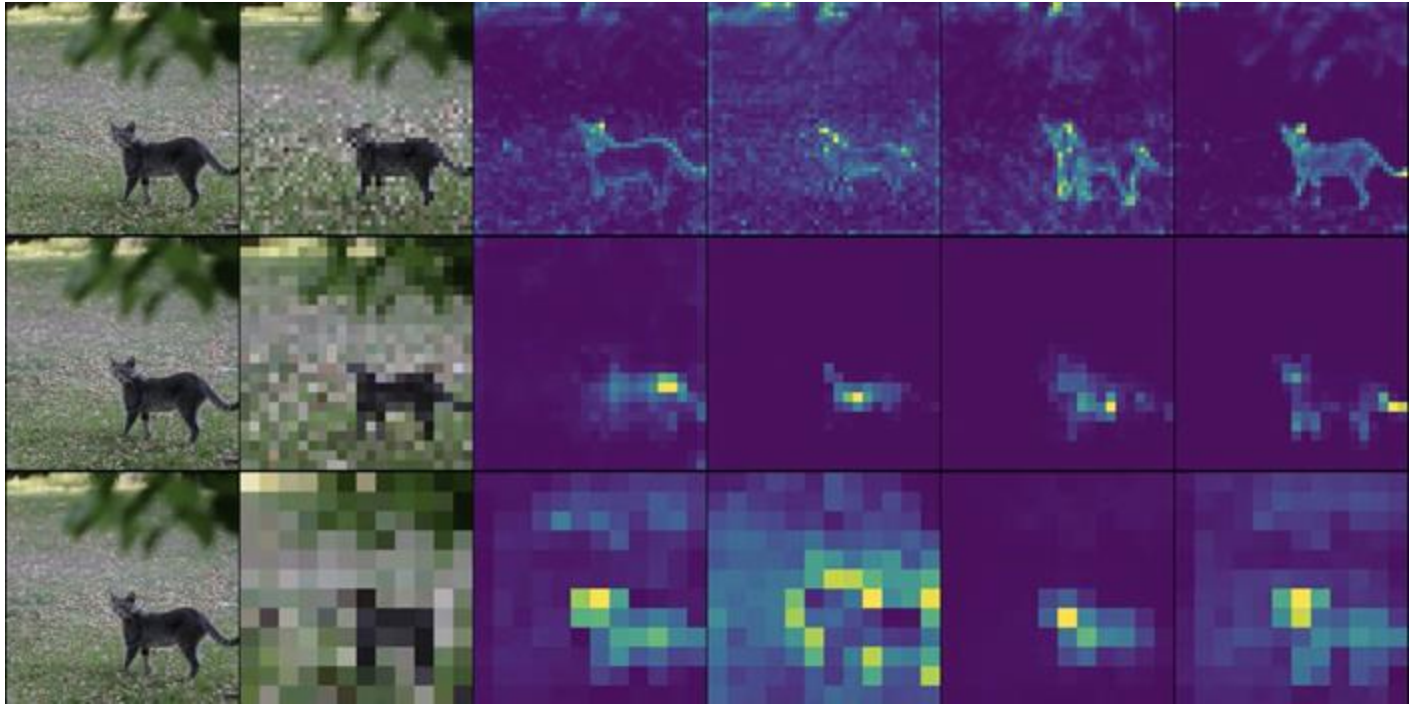
Transformers for other modalities



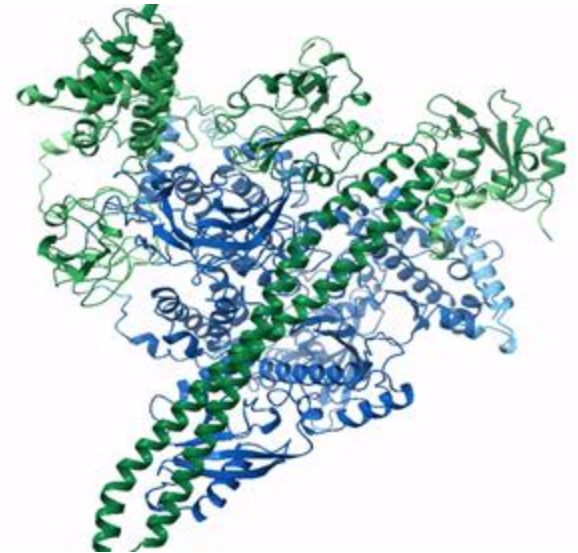
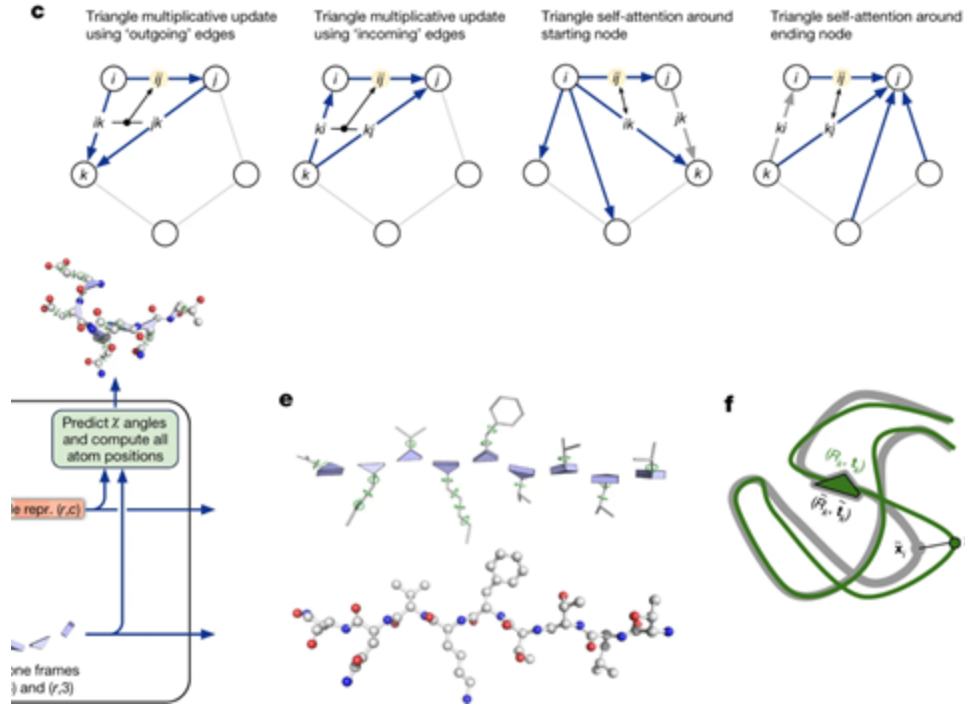
High Performance Machine Learning Group

SURF

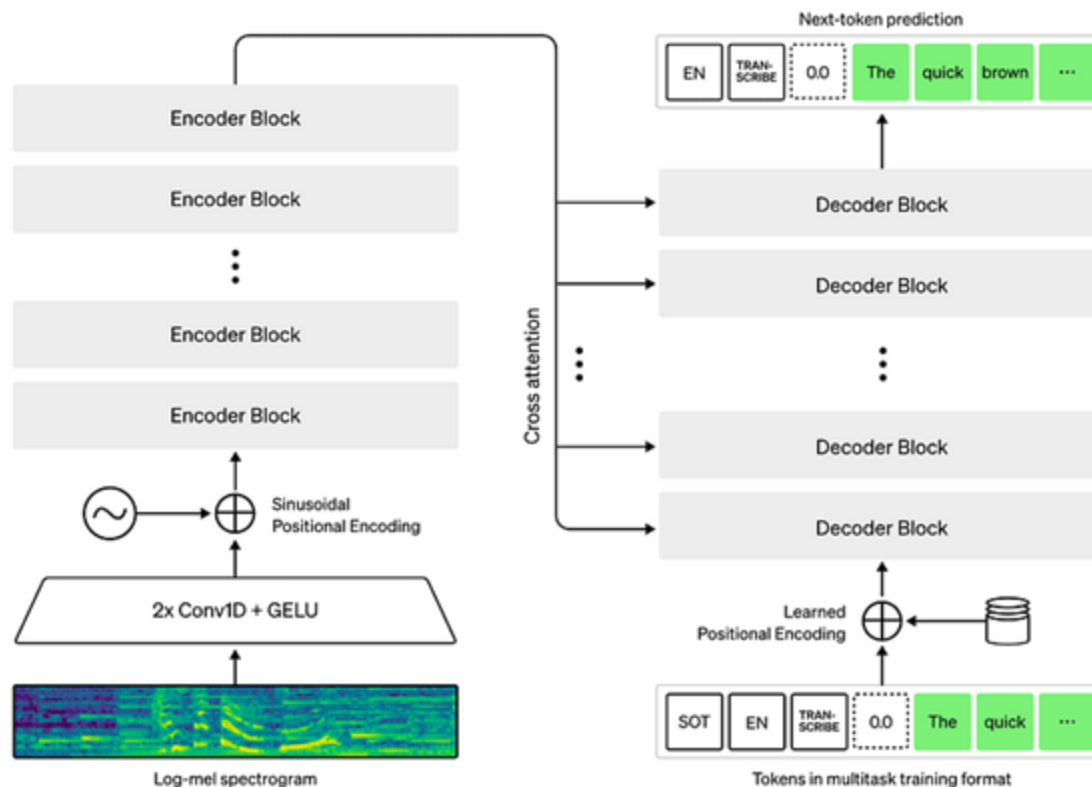
Transformers for other modalities: ViT



Transformers for other modalities: EvoFormer



Transformers for other modalities: Whisper



Want to learn more?

Use existing (efficient!) frameworks



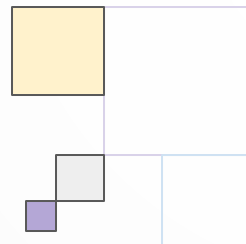
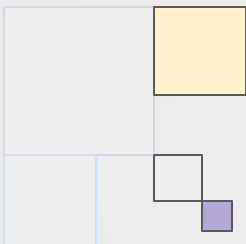
Hugging Face



unsloth

Make use of SURF guides for efficient
finetuning:

<https://servicedesk.surf.nl/wiki/display/WIKI/LLM+finetuning+on+Snellius>



Fin

Introduction Series

