# Generative AI - Text to Text - Part 2



## Ahmed Elbagoury

# Introduction – Ahmed Elbagoury

Senior ML Engineer, Google.

ML Research + Teaching

Applied ML

Teaching Applied ML

# Optimize Your Experience

✅ Ask questions in the class

✅ Don't ignore the math (Whys and Hows)

✅ Fortify understanding by reading papers, walking through code implementations and trying out notebooks!

✅ Enjoy the subject. Study deeply, seek understanding, read papers, practice problems, ask questions

✅ Don't spend time in self-doubt. Unlike real life tortoise & hare race, slow & steady literally wins the career race

# Here's the plan

**Part 1** ➡ **How LLM can be used in applications** ✅

**Part 2** ➡ **Decoder details and Optimizations** ⇐ **Today**

# Today's Agenda

**1** GPT Models

**2** Sampling Strategies

**3** LLM Optimizations

# We talked about encoder-decoder models

# We talked about



BERT

7

# What is missing?

{ik} | INTERVIEW KICKSTART

# Decoder Models



Transformer          GPT

# Quiz

**Which Component should be dropped from the decoder in GPT <u>(decoder only)</u> models?**
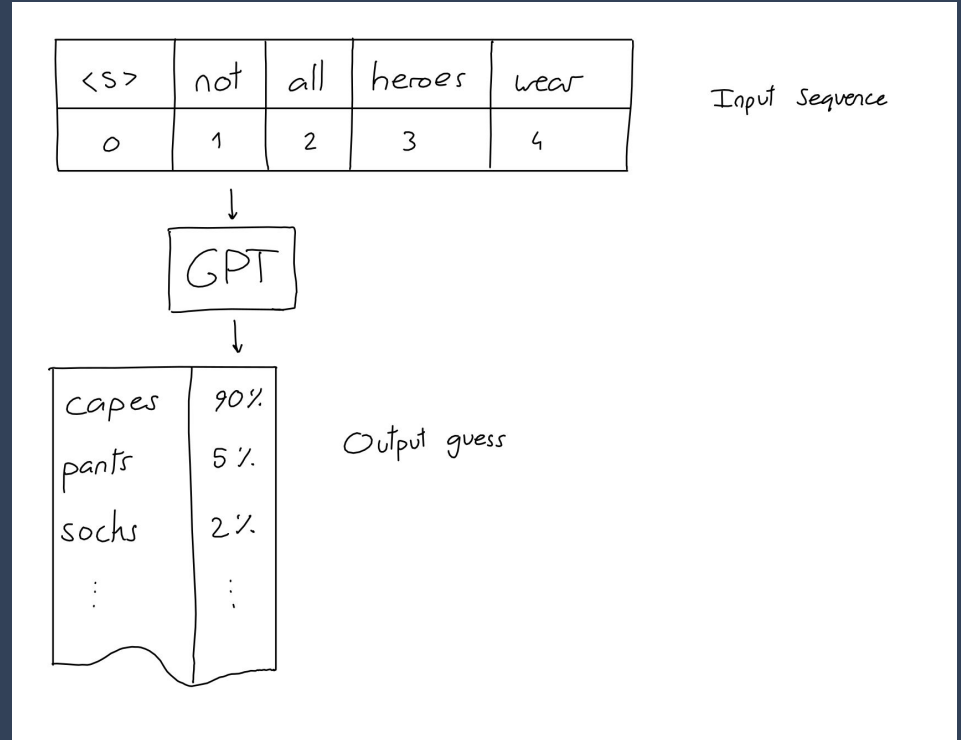
# Quiz

**Which Component should be dropped from the decoder in GPT <u>(decoder only)</u> models?**

⇒ Encoder-decoder attention should be dropped since there is no encoder!

{ik} | INTERVIEW KICKSTART

# Decoder Models – Text Prediction

"Not all heroes wear capes" ⟹ 5

| <s> | not | all | heroes | wear |
|-----|-----|-----|--------|------|
| 0 | 1 | 2 | 3 | 4 |

Input Sequence

↓

GPT

↓

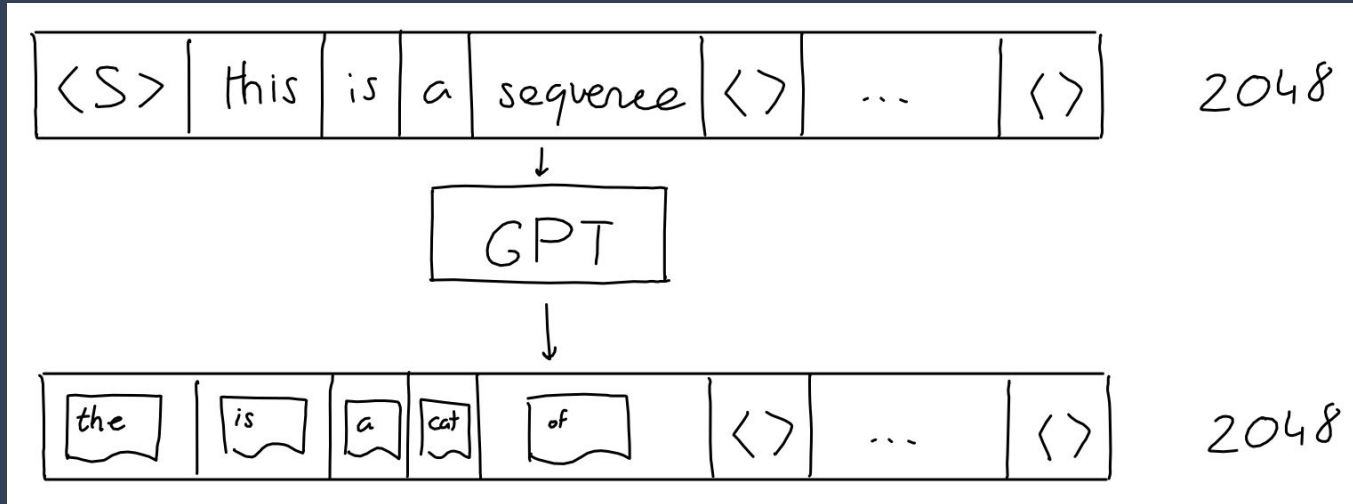| capes | 90% |
|-------|-----|
| pants | 5% |
| socks | 2% |
| ⋮ | ⋮ |

Output guess

# Decoder Models – Text Prediction

How to get multiple words and for different lengths?!

# Decoder Models – Text Prediction

How to get multiple words and for different lengths?!



When generating text, we typically only look at the guess for the last word of the sequence.

# Decoder Models – Task Classifier

After pre-training

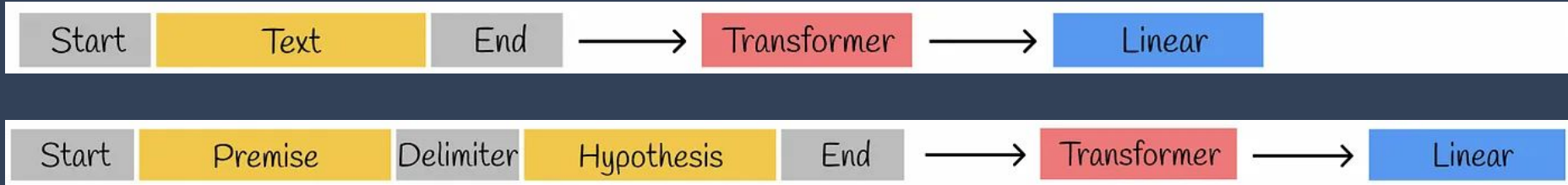- GPT can capture linguistic knowledge of input sequences ✅
- However, to make it better perform on downstream tasks ❌

⇒ it needs to be fine-tuned on a supervised problem.

# Decoder Models – Task Classifier

Start | Text | End ⟶ Transformer ⟶ Linear

# Decoder Models – Task Classifier

Start | Text | End → Transformer → Linear

Start | Premise | Delimiter | Hypothesis | End → Transformer → Linear
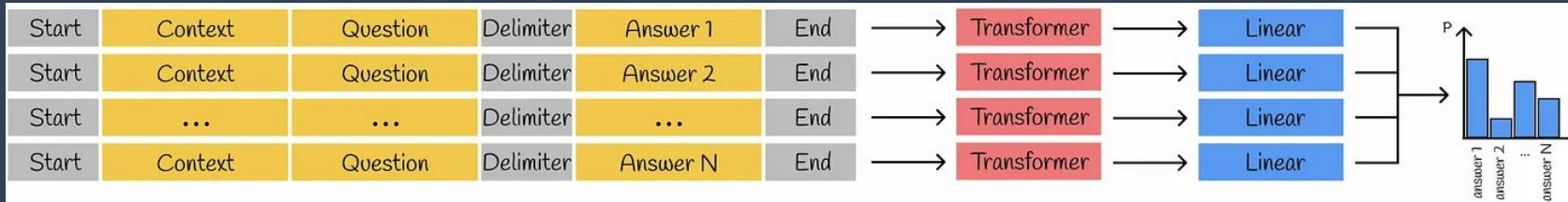
# Decoder Models – Task Classifier

# Decoder Models – Task Classifier

# Why Decoder Models aka Causal Decoder/Models?

- Cost of training for Causal Decoder (CD) is cheaper
- CD works better for In-Context learning
  $\Rightarrow$ It has a more straightforward effect for CD

# Since Decoding is very important

# Let's zoom in to it

# Today's Agenda

**1** GPT Models

**2** Sampling Strategies

**3** LLM Optimizations

# Decoding Details

# Decoding Details

# What is the Optimal Solution

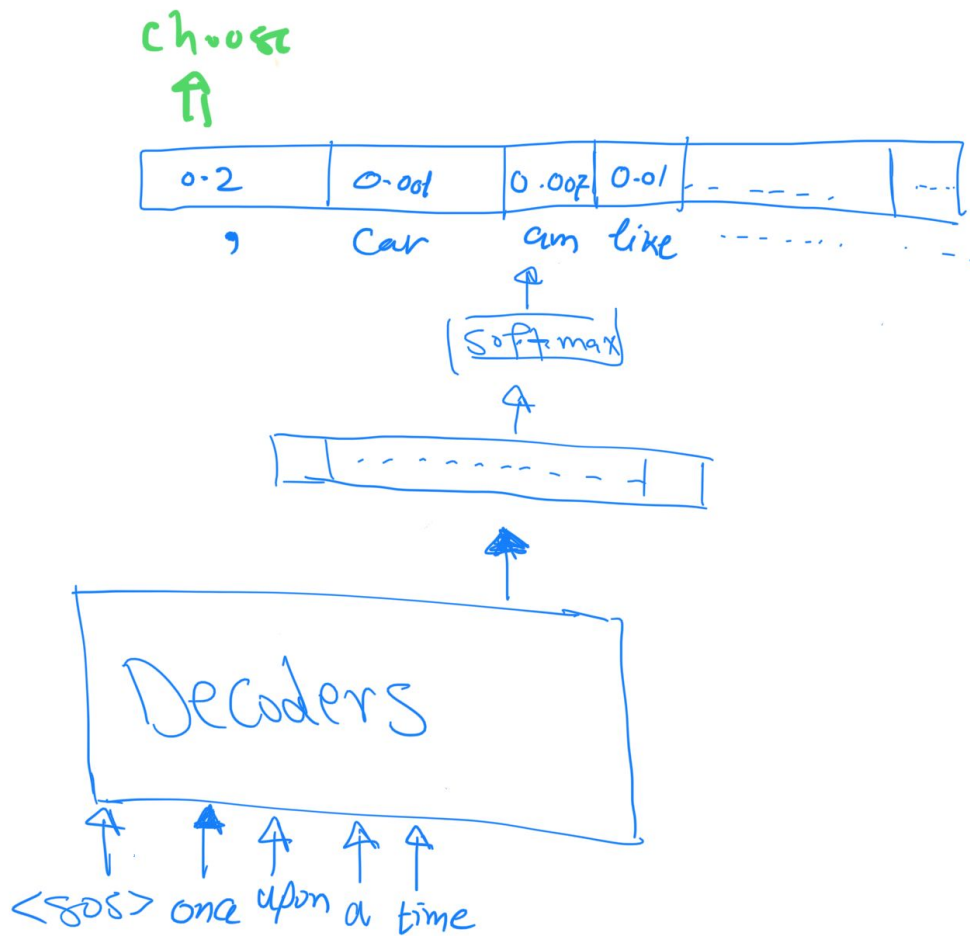- We want to maximize the likelihood of the output sequence
  - This is the multiplication of conditional probabilities
  - Or minimizing the summation of log likelihood

# What is the Optimal Solution

- We want to maximize the likelihood of the output sequence
    - This is the multiplication of conditional probabilities
    - Or minimizing the summation of log likelihood

What are the problems with the previous approach? (greedy)

{ik} | INTERVIEW KICKSTART

total = 4·2

Student ___1.9___ •__3__ EOS

0·3

I am a

total = 1·2

0·7 player __0·2__ eats ___0·3___ EOS

0·1
xyz

0·2
abc

28

# Decoding Details: Beam Search

- Keep your options open (to some extent)
- Just in case your local choices are not optimal
- It requires more memory and computational power to keep track of the beam

beam = ??

beam = 3

# Decoding Details: What if we got the answer wrong

Help me write a congratulation email to a new colleague ⟹ **Decoder Model** ⟹ I hope this message finds you well. I wanted to take a moment to personally welcome you to ...

Let's say I want to see other drafts!

# Decoding Details: What if we got the answer wrong

Help me write a congratulation email to a new colleague → **Decoder Model** → I hope this message finds you well. I wanted to take a moment to personally welcome you to …

Help me write a congratulation email to a new colleague → **Decoder Model** → We're thrilled to have you on board and are excited about the great things you'll bring to the team

33

# Decoding Details: What if we got the answer wrong

- Random (aka stochastic) sampling to the rescue

**Sorry more terminology!**
**So that you know the lingo!**

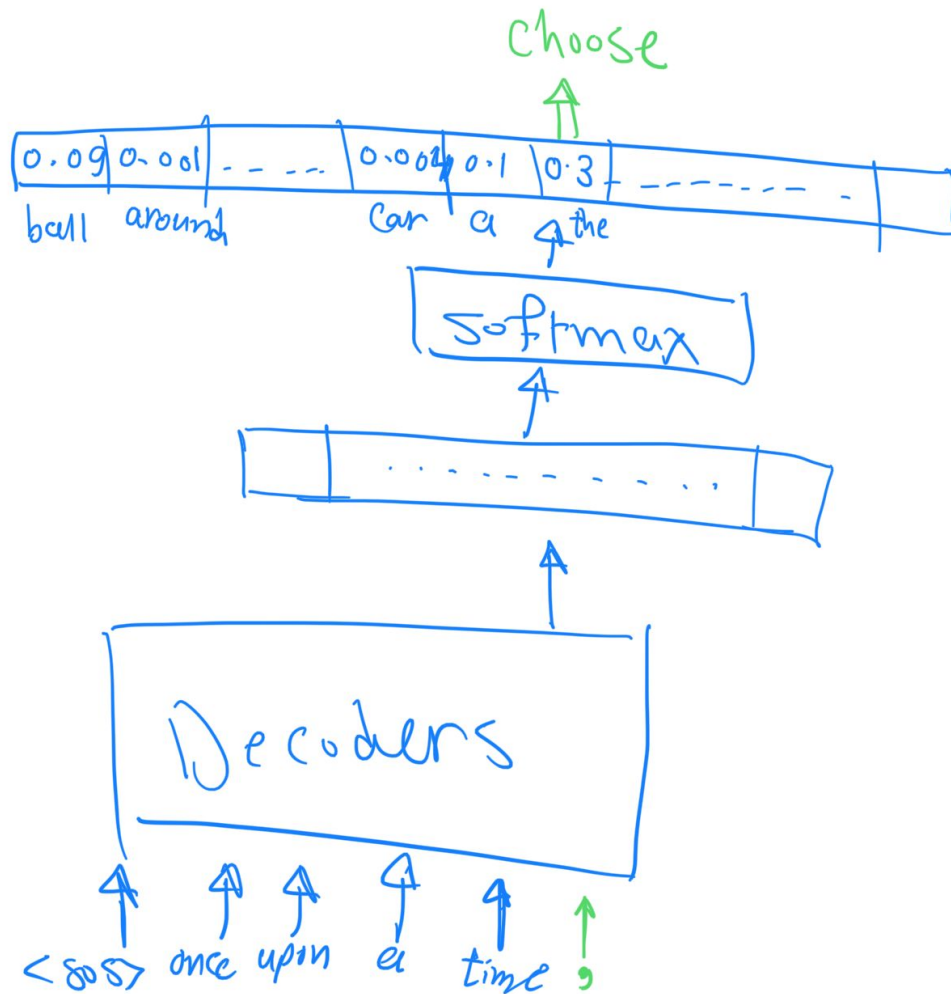# Decoding Details: What if we got the answer wrong

- Random (aka stochastic) sampling to the rescue

**Sorry more terminology!
So that you know the lingo!**

- **It just means**

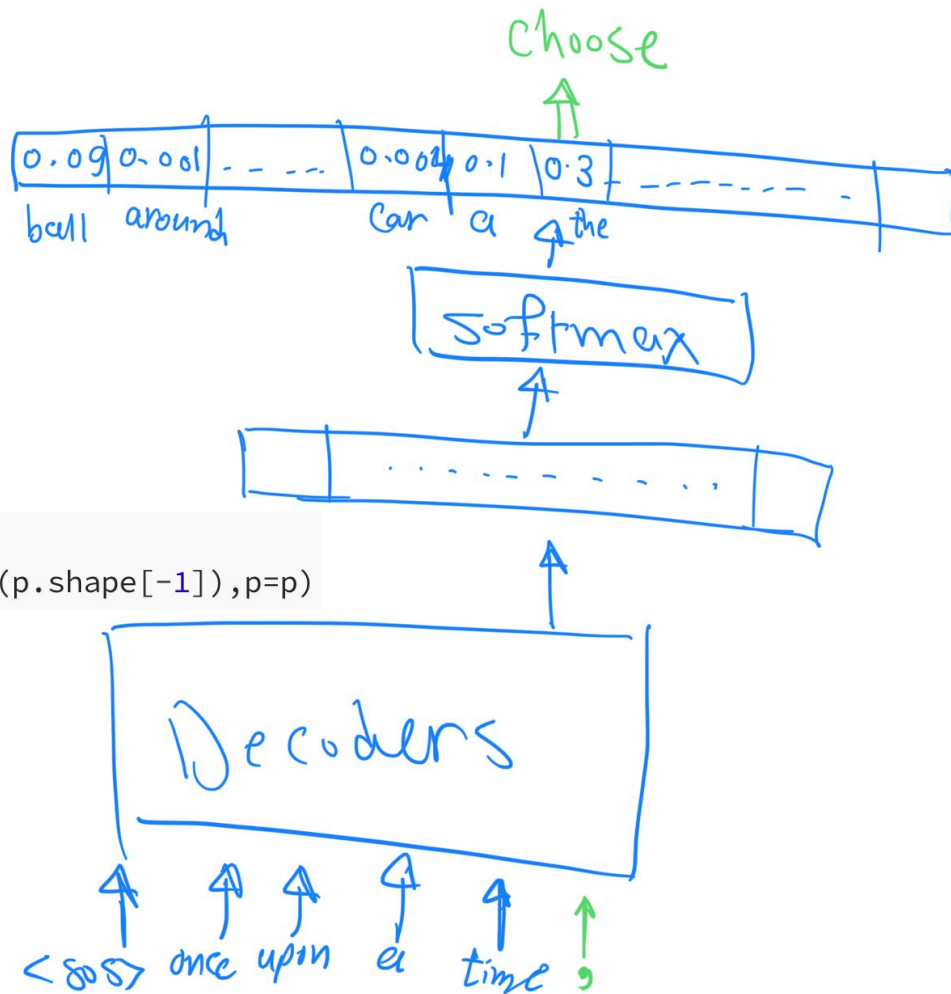  - Choose randomly according to the probability scores

# Random Sampling

# Random Sampling



```
def sample(p):
    return np.random.choice(np.arange(p.shape[-1]),p=p)
```

# Decoding Details: Top-K

# Decoding Details: Top-p aka Nucleus

# Temperature



SOFTMAX WITHOUT TEMPERATURE (T=1)

$$\frac{e^{z_i}}{\sum_j e^{z_i}}$$

SOFTMAX WITH TEMPERATURE

$$\frac{e^{z_i/T}}{\sum_j e^{z_i/T}}$$

Temperature: 0.7

LESS ENTROPY → INCREASE IN ENTROPY WITH INCREASE IN T → MORE ENTROPY

{ik} | INTERVIEW KICKSTART

# Demo Time

https://colab.research.google.com/drive/1CejK4UaV0O6L3c4jnKqtWENp9YFboNlg

# Today's Agenda

**1** GPT Models

**2** Sampling Strategies

**3** LLM Optimizations

# Optimization Strategies for LLM

Larger Models Mean More Capacity
- To handle variety of tasks ⇒ Models need to be trained on large amount of data
- To avoid underfitting ⇒ Models need to be larger!

**Number of Parameters**

(in Millions)

- 65M — (2017)
- 117M — GPT1 (2018)
- 1500 — GPT2 (2019)
- 175000 M — GPT3 (2020)
- 280000 M — GOPHER (2021)
- 540000 M — PALM (2022)

# Larger Models Mean More Capacity

- To handle variety of tasks
  - Models need to be trained on  large amount of data
- To avoid underfitting
  - Models need to be larger!

# Implications of Increasing Model Size

# Implications of Increasing Model Size

- **More hardware resources**
  - On-device use cases?
- Higher Latency
- Worse Carbon footprint

# Do we have to use these huge models?

- A line of research motivated by the challenges of training and productionizing LLMs
- It aims to
  - Improving training and/or inference latency
  - Reducing model sizes with little to no impact on quality

# Do we have to use these huge models?

- A line of research motivated by the challenges of training and productionizing LLMs
- It aims to
  - Improving training and/or inference latency
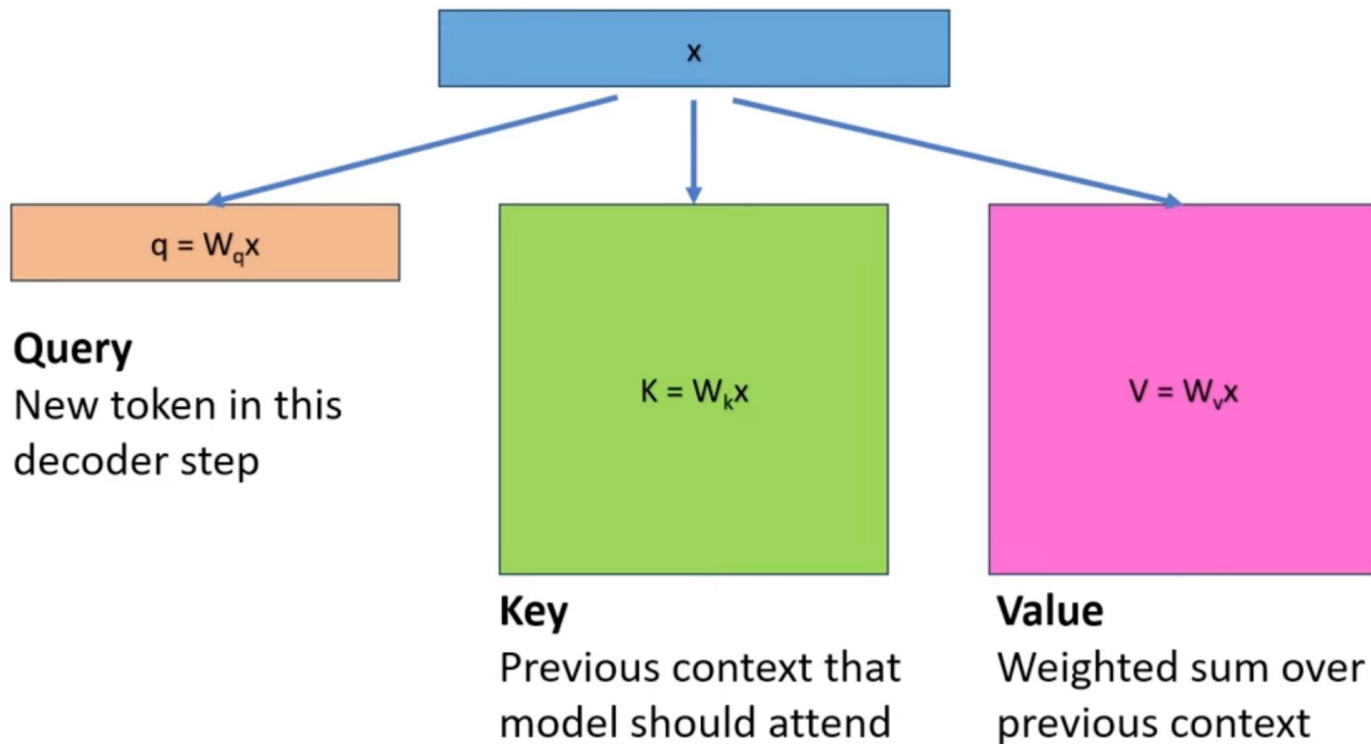  - Reducing model sizes with little to no impact on quality

Major techniques:

- KV Caching
- LoRA Tuning (More generally PEFT)
- Distillation
- Model Pruning
- Quantization

# KV Caching

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Photo from

# KV Caching



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

x

$q = W_q x$

**Query**
New token in this decoder step

$K = W_k x$

**Key**
Previous context that model should attend

$V = W_v x$

**Value**
Weighted sum over previous context

51

Photo from

{ik} | INTERVIEW KICKSTART

# KV Caching

It was a cold windy morning when I stepped outside, feeling a chill

$$q = W_q x$$

$$K = W_k x$$

$$V = W_v x$$

# KV Caching



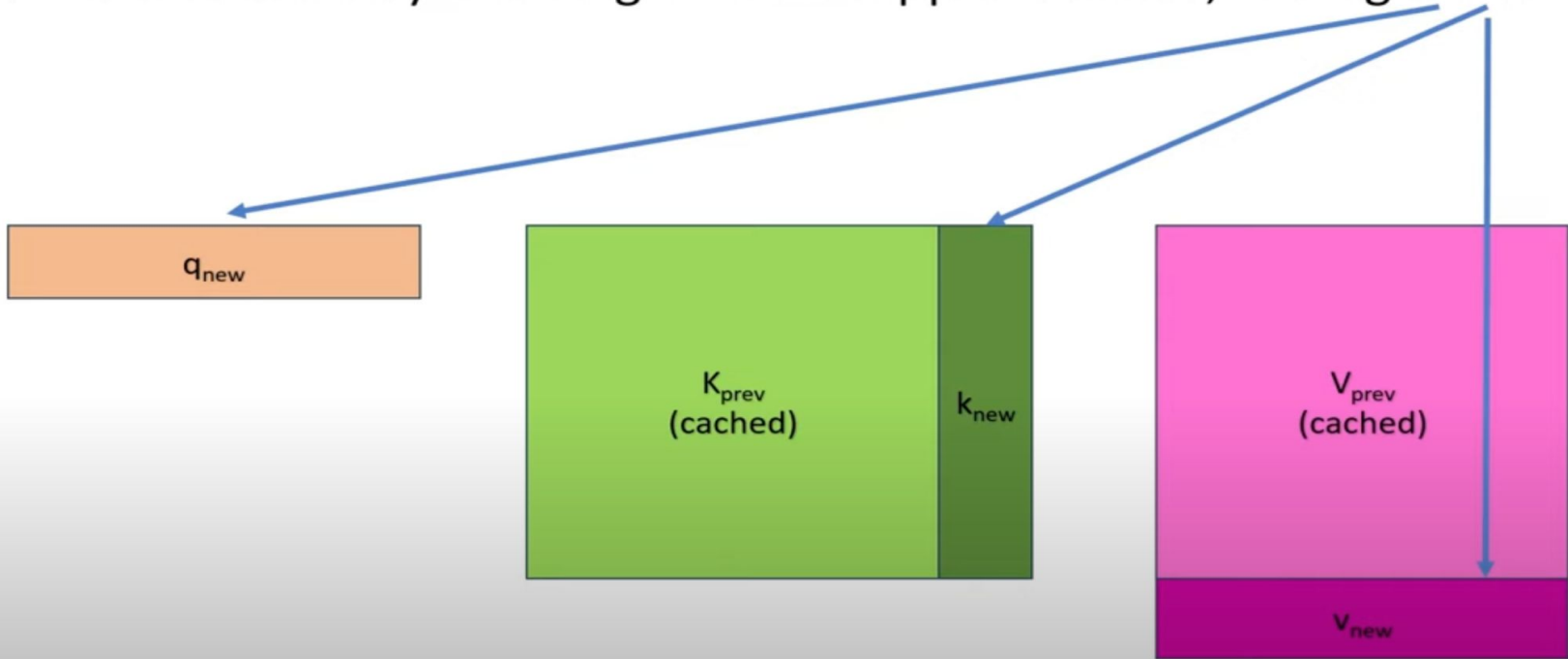It was a cold windy morning when I stepped outside, feeling a chill
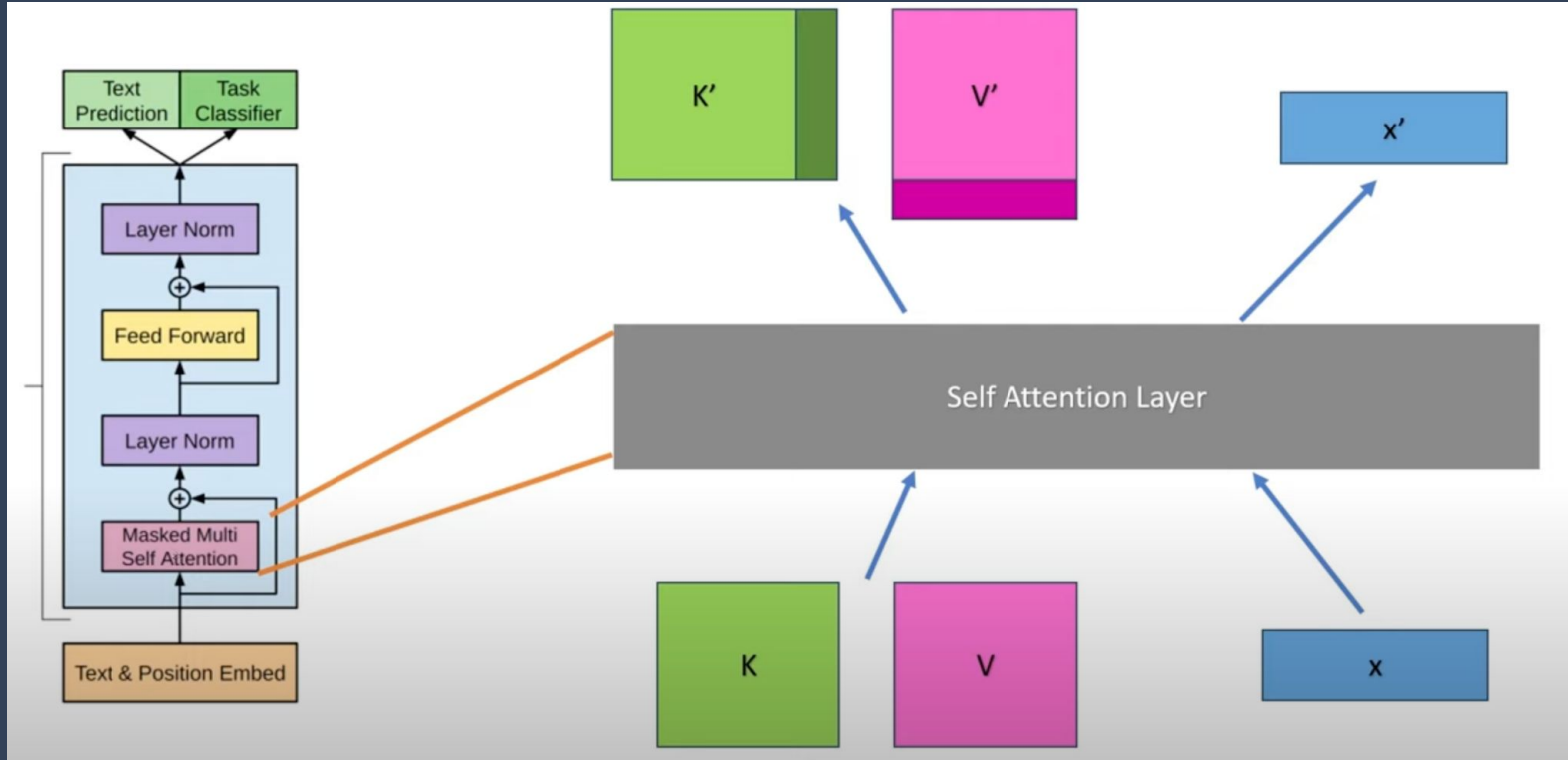
$q = W_q x$

$K = W_k x$

$V = W_v x$

# KV Caching



It was a cold windy morning when I stepped outside, feeling a chill

$q_{new}$

$K_{prev}$ (cached)    $k_{new}$

$V_{prev}$ (cached)

$v_{new}$

# KV Caching

# Demo Time

https://colab.research.google.com/drive/12ioUtylE5BuWTNjDHdecNQ8Xnb1dRISW

# Quiz

**Which part will run faster: Red or Blue? and why?**
**Assuming that red is the prompt**



It was a cold windy morning when I stepped outside, feeling a chill crawl
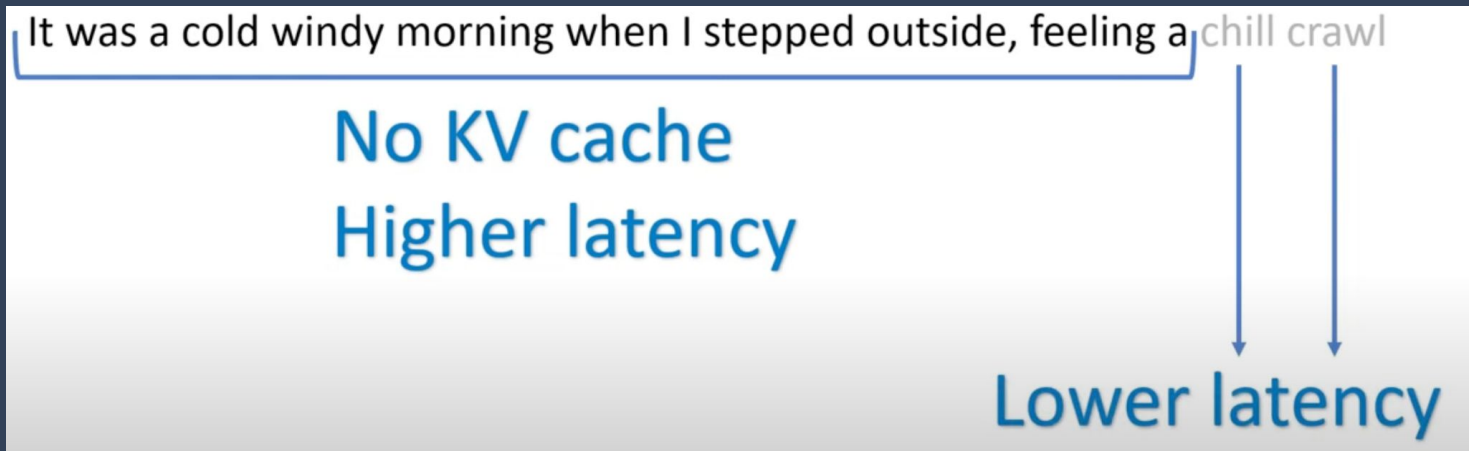
# Quiz

**Which part will run faster: Red or Blue? and why?**
**Assuming that red is the prompt**



It was a cold windy morning when I stepped outside, feeling a chill crawl

It was a cold windy morning when I stepped outside, feeling a chill crawl

No KV cache
Higher latency

Lower latency

# One more Quiz

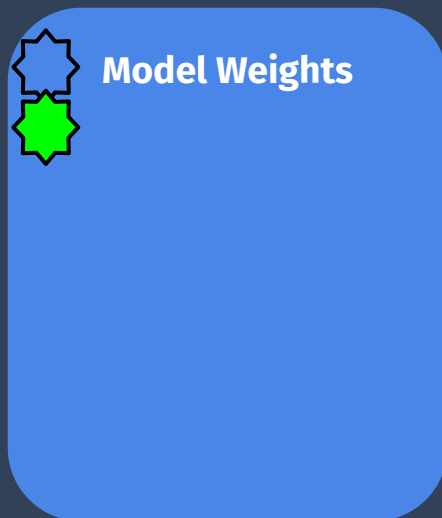**KV Caching improves <u>inference</u> latency at the expense of**

A) Memory Usage during inference
B) Training Time
C) A & B
D) Memory Usage during training

# Quiz

**KV Caching improves <u>inference</u> latency at the expense of**

A) Memory Usage during inference
B) Training Time
C) A & B
D) Memory Usage during training

# LoRA

- Low-Rank Adaptation of Large Language Models
- Is a type of what is known as Parameter-Efficient tuning (PEFT)
- The goal is fine-tuning a model
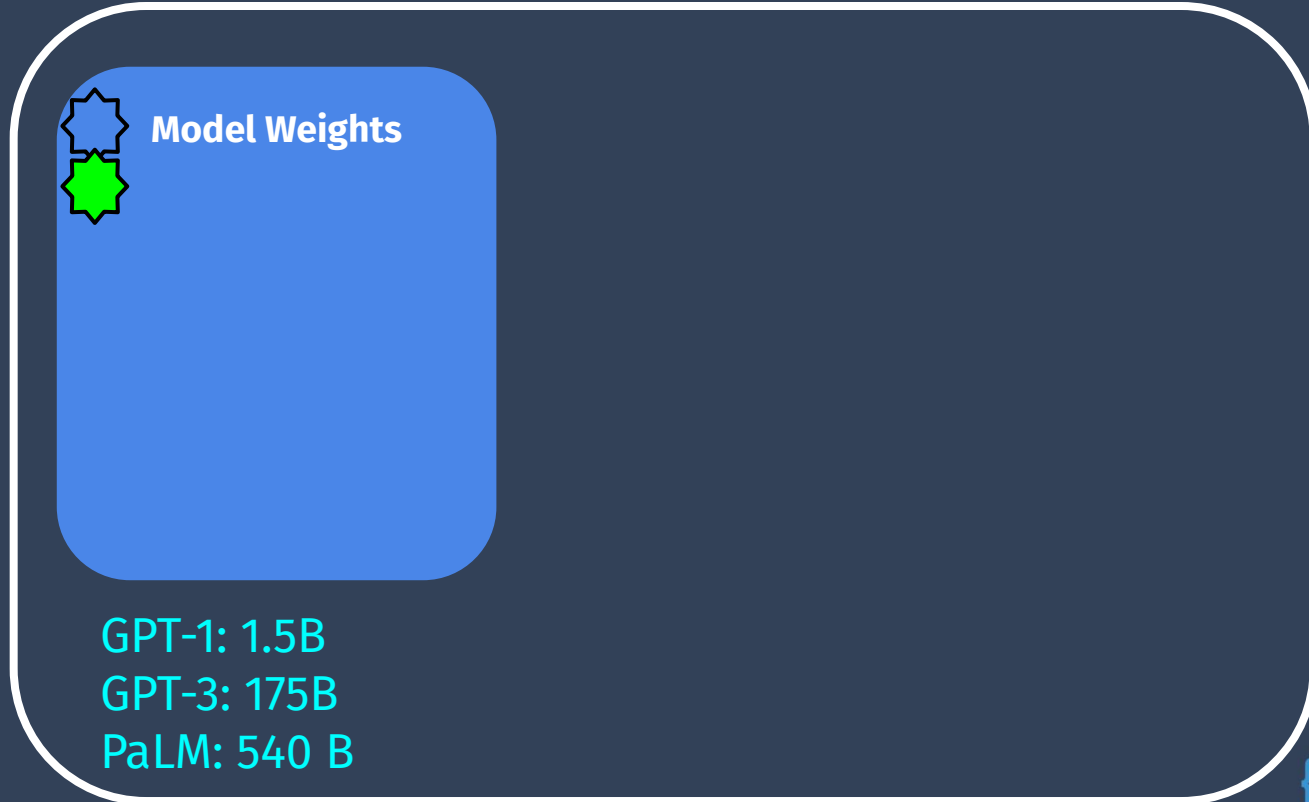  - By updating smaller number of parameters

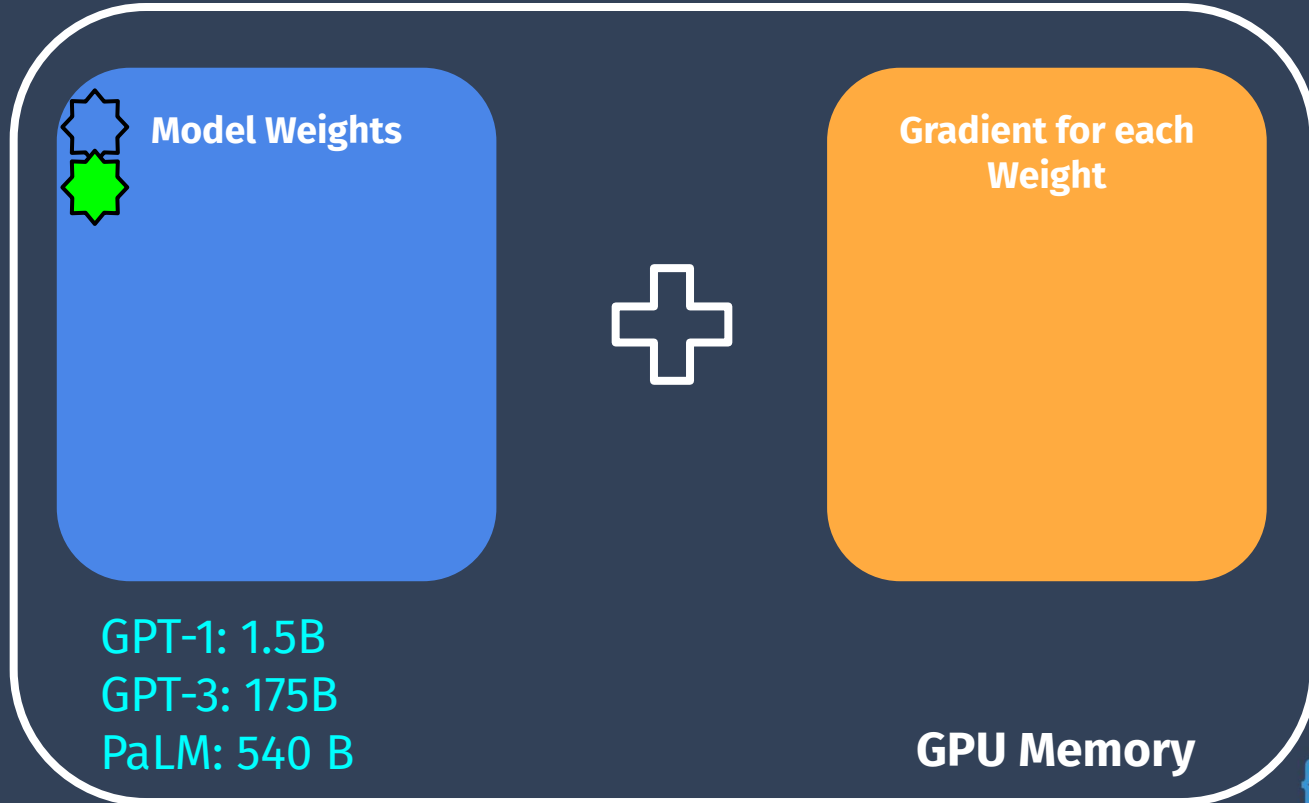# Why LoRA?

**Model Weights**

GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

# Regular Fine-Tuning

During
fine-tuning ⟹
we need to
load this in
Memory

**Model Weights**

GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

{ik} | INTERVIEW KICKSTART

# Regular Fine-Tuning

During
fine-tuning ⇒
we need to
load this in
Memory

**Model Weights**

➕

**Gradient for each Weight**

GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

**GPU Memory**

65

{ik} | INTERVIEW KICKSTART

# Regular Fine-Tuning

During fine-tuning ⇒ we need to load this in Memory

**Model Weights**

$$\begin{pmatrix} w11 & w12 & w13 \\ w21 & w22 & w23 \\ w31 & w32 & w33 \end{pmatrix}$$
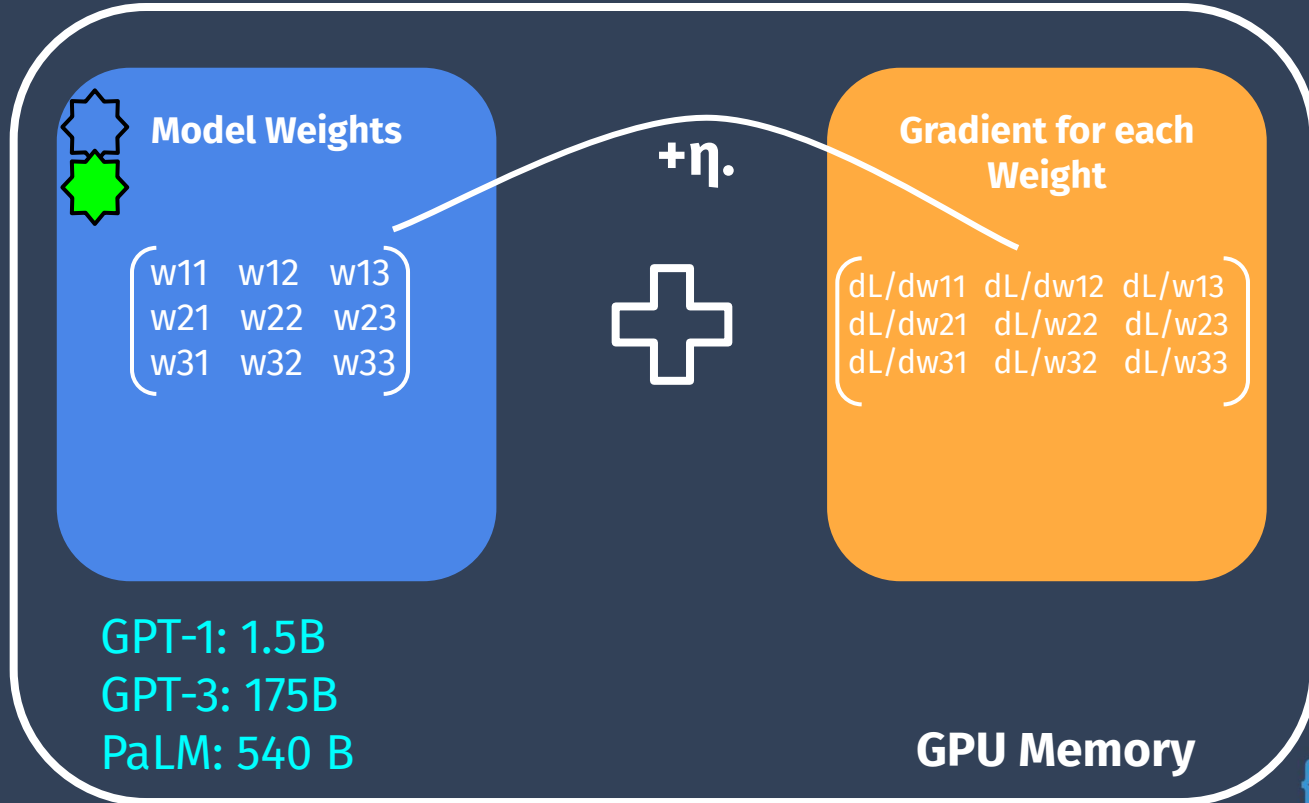
GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

**+**

**Gradient for each Weight**

$$\begin{pmatrix} dL/dw11 & dL/dw12 & dL/w13 \\ dL/dw21 & dL/w22 & dL/w23 \\ dL/dw31 & dL/w32 & dL/w33 \end{pmatrix}$$

**GPU Memory**

{ik} | INTERVIEW KICKSTART

# Regular Fine-Tuning

During fine-tuning we need to load ⇒

**Model Weights**

$$\begin{pmatrix} w11 & w12 & w13 \\ w21 & w22 & w23 \\ w31 & w32 & w33 \end{pmatrix}$$

**+η.**

**+**

**Gradient for each Weight**

$$\begin{pmatrix} dL/dw11 & dL/dw12 & dL/w13 \\ dL/dw21 & dL/w22 & dL/w23 \\ dL/dw31 & dL/w32 & dL/w33 \end{pmatrix}$$
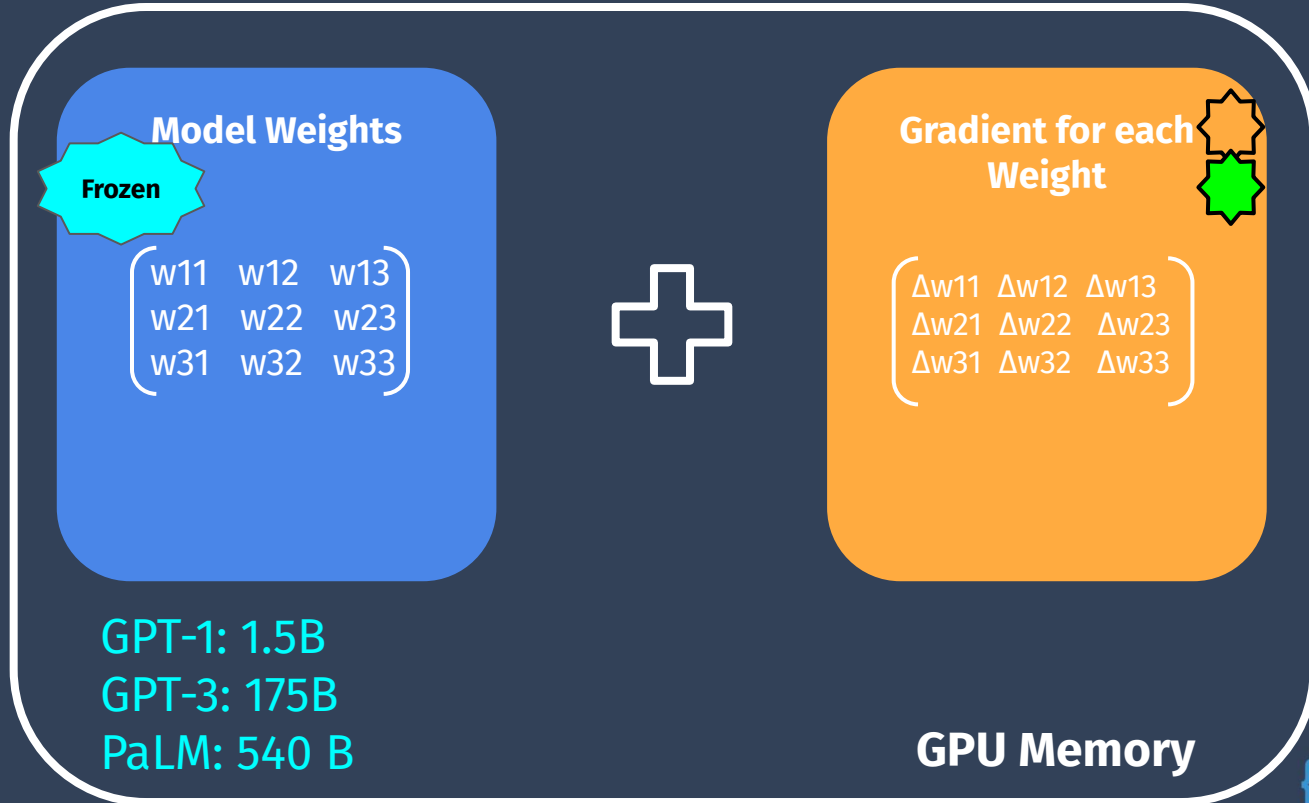
GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

**GPU Memory**

{ik} | INTERVIEW KICKSTART

# LoRA

**Model Weights**

Frozen

$$\begin{bmatrix} w11 & w12 & w13 \\ w21 & w22 & w23 \\ w31 & w32 & w33 \end{bmatrix}$$

**+**

**Gradient for each Weight**

$$\begin{bmatrix} \Delta w11 & \Delta w12 & \Delta w13 \\ \Delta w21 & \Delta w22 & \Delta w23 \\ \Delta w31 & \Delta w32 & \Delta w33 \end{bmatrix}$$

GPT-1: 1.5B

GPT-3: 175B

PaLM: 540 B

**GPU Memory**

{ik} | INTERVIEW KICKSTART

# LoRA

**Model Weights**

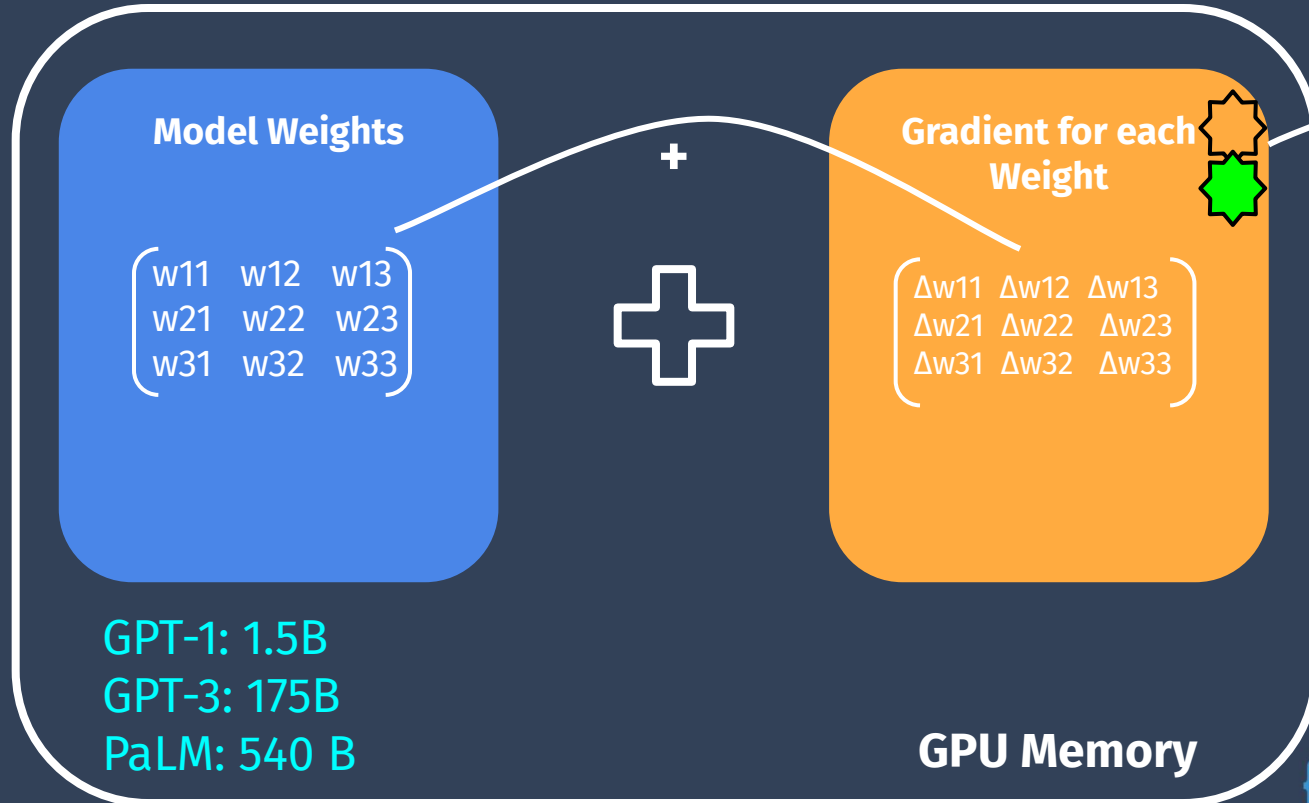$$\begin{pmatrix} w11 & w12 & w13 \\ w21 & w22 & w23 \\ w31 & w32 & w33 \end{pmatrix}$$

**+**

**Gradient for each Weight**

$$\begin{pmatrix} \Delta w11 & \Delta w12 & \Delta w13 \\ \Delta w21 & \Delta w22 & \Delta w23 \\ \Delta w31 & \Delta w32 & \Delta w33 \end{pmatrix}$$

GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

**GPU Memory**

{ik} | INTERVIEW KICKSTART

# LoRA

**Does not it this double the number of parameters as well??**

GPU Memory

**Model Weights**

$$\begin{pmatrix} w11 & w12 & w13 \\ w21 & w22 & w23 \\ w31 & w32 & w33 \end{pmatrix}$$

**+**

**Gradient for each Weight**

$$\begin{pmatrix} \Delta w11 & \Delta w12 & \Delta w13 \\ \Delta w21 & \Delta w22 & \Delta w23 \\ \Delta w31 & \Delta w32 & \Delta w33 \end{pmatrix}$$

GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

# LoRA

**Model Weights**

$$\begin{bmatrix} w11+\Delta w11 & w12+\Delta w12 & w13+\Delta w13 \\ w21+\Delta w21 & w22+\Delta w22 & w23+\Delta w23 \\ w31+\Delta w31 & w32+\Delta w32 & w33+\Delta w33 \end{bmatrix}$$

$=$

**Model Weights**

$$\begin{bmatrix} w11 & w12 & w13 \\ w21 & w22 & w23 \\ w31 & w32 & w33 \end{bmatrix}$$

$+$

**Gradient for each Weight**

$$\begin{bmatrix} \Delta w11 & \Delta w12 & \Delta w13 \\ \Delta w21 & \Delta w22 & \Delta w23 \\ \Delta w31 & \Delta w32 & \Delta w33 \end{bmatrix}$$

**GPU Memory**

GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

# LoRA

**Model Weights**

$$\begin{pmatrix} w11 & w12 & w13 \\ w21 & w22 & w23 \\ w31 & w32 & w33 \end{pmatrix}$$

**+**

**Gradient for each Weight**

$$\begin{pmatrix} \Delta w11 & \Delta w12 & \Delta w13 \\ \Delta w21 & \Delta w22 & \Delta w23 \\ \Delta w31 & \Delta w32 & \Delta w33 \end{pmatrix}$$

GPT-1: 1.5B
GPT-3: 175B
PaLM: 540 B

**GPU Memory**

72

{ik} | INTERVIEW KICKSTART

# LoRA Tuning

$$\Delta W = \begin{bmatrix} \Delta \omega_{11} & \Delta \omega_{12} & \Delta \omega_{13} \\ \Delta \omega_{21} & \Delta \omega_{22} & \Delta \omega_{23} \\ \Delta \omega_{31} & \Delta \omega_{32} & \Delta \omega_{33} \end{bmatrix} = \underset{(d \times r)\,(r \times k)}{B\,A} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} \cdot (a_1 \; a_2 \; a_3)$$

$3 \leftarrow d \times k \rightarrow 3$

$(3 \times 1) \qquad (1 \times 3)$

or

$$= \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix}$$

$(3 \times 2) \qquad (2 \times 3)$

# LoRA Tuning

$$\Delta W = \begin{bmatrix} \Delta \omega_{11} & \Delta \omega_{12} & \Delta \omega_{13} \\ \Delta \omega_{21} & \Delta \omega_{22} & \Delta \omega_{23} \\ \Delta \omega_{31} & \Delta \omega_{32} & \Delta \omega_{33} \end{bmatrix} = \underset{(d \times r)\,(r \times k)}{B\,A} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} \cdot (a_1 \; a_2 \; a_3)$$

$3 \leftarrow d \times k \rightarrow 3$

$(3 \times 1) \qquad (1 \times 3)$

$B = 0$

$r$

$A = N(0, \sigma^2)$

or

$$= \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix}$$

$(3 \times 2) \qquad (2 \times 3)$

# LoRA Tuning

# Quiz

**LoRA Tuning**

A) Improves Training Memory footprint
B) Improves Inference Memory footprint
C) A & B

# Quiz

**LoRA Tuning**

A) Improves Training Memory footprint
B) Improves Inference Memory footprint
C) A & B

Knowledge Distillation

# Knowledge Distillation

# Response-based knowledge

# Distillation

# Distillation

# Feature-based knowledge

# Demo Time

https://colab.research.google.com/drive/1MH2kIPtcGzu0g2VkQ8lcCHFVeR6nAls7

# Distillation in LLM

- DistilBERT
  - The student model imitate the teacher model
  - Has about half the total number of parameters of BERT base and retains 95% of BERT's performances on the language understanding benchmark

TinyBert

# Distillation: Practical Considerations

- The student model will be good at the distillation task **<u>only</u>**
- You still need unlabeled data
  - In some enterprise setups you may not be allowed to use customer data
- Model refresh

# Quiz

**Distillation Improves:**

A) Training latency
B) Inference latency
C) A & B

# Quiz

**Distillation Improves:**
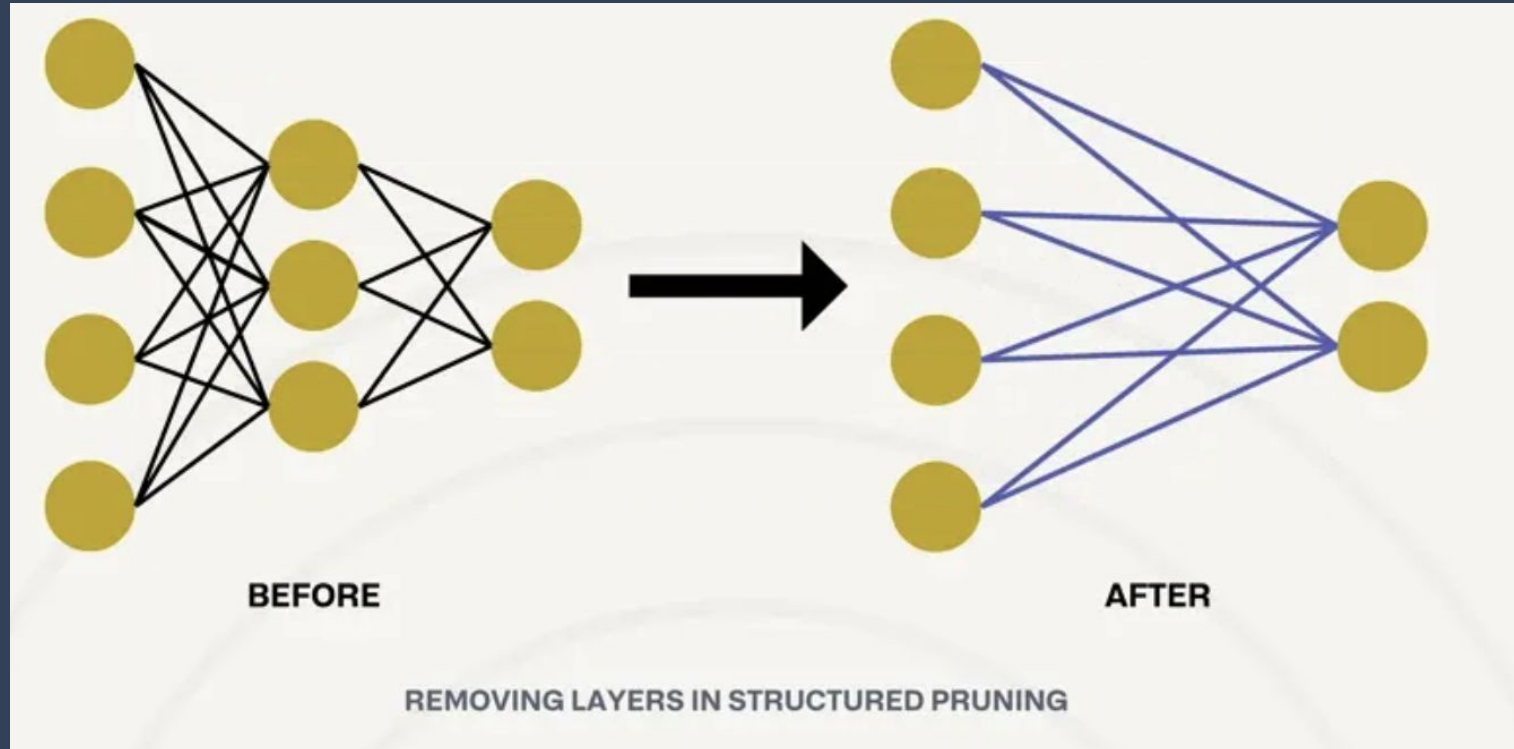
A) Training latency
B) Inference latency
C) A & B

# Model Pruning

- The reasoning is that:
  "Neural networks have an excess of parameters needed to generalize well and make accurate predictions"
- Then, we should drop these extra parameters

# Model Pruning: Unstructured



ZEROING WEIGHTS IN UNSTRUCTURED PRUNING

# Model Pruning: Structured



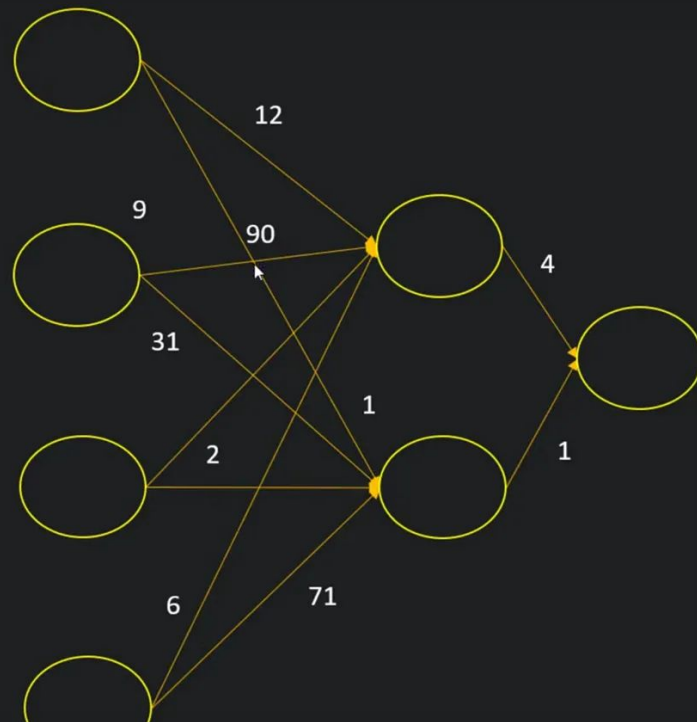REMOVING LAYERS IN STRUCTURED PRUNING

# Model Pruning

- **At training time**
  - ✅✅ More efficient as the model is trained with sparsity at training time
  - ❌❌ Can be more complex
- **Post training time**
  - ✅✅ Simpler to implement and adjust per use-case
  - ❌❌ May require fine-tuning

# Model Quantization

# Quiz

**Model Pruning and Quantization Improve**

A) Training latency
B) Inference latency
C) A & B

# Quiz

**Model Pruning and Quantization Improve**

A) Training latency
B) Inference latency
C) A & B

# Is that it?

No, there are many many other optimizations
- Speculative decoding
- Linear attention
- Model parallelism
- …

# Summary and recap

- We went through the decoder details
- Sampling strategies
- Optimizations
    - KV Caching
    - LoRA
    - Distillation
    - Pruning
    - Quantizations

# Questions & Discussions

{ik} | INTERVIEW KICKSTART

# References

- https://dugas.ch/artificial_curiosity/GPT_architecture.html
- https://towardsdatascience.com/large-language-models-gpt-1-generative-pre-trained-transformer-7b895f296d3b
- https://medium.com/@yumo-bai/why-are-most-llms-decoder-only-590c903e4789
- https://medium.com/@javaid.nabi/all-you-need-to-know-about-llm-text-generation-03b138e0ed19
- https://mlabonne.github.io/blog/posts/2023-06-07-Decoding_strategies.html
- https://www.appypie.com/blog/hardware-requirements-for-llm-training#Hardware
- https://sujayskumar.com/2020/04/24/reducing-model-size/
- https://snorkel.ai/llm-distillation-demystified-a-complete-guide/
- https://neptune.ai/blog/knowledge-distillation
- https://medium.com/huggingface/distilbert-8cf3380435b5
- https://www.datature.io/blog/a-comprehensive-guide-to-neural-network-model-pruning
- https://www.youtube.com/watch?v=80bIUggRJf4
- https://medium.com/@joaolages/kv-caching-explained-276520203249
- https://huggingface.co/docs/peft/main/en/conceptual_guides/lora
- https://www.youtube.com/watch?v=KEv-F5UkhxU
- https://www.youtube.com/watch?v=KEv-F5UkhxU

{ik} | INTERVIEW KICKSTART