# **Efficiency & Advanced Techniques**

#### Quantization, MoE, linear attention, LoRA, RAG, reasoning

Zdeněk Kasner, Michal Novák



Charles University Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated

## Today's learning outcomes

After today's class, you should be able to:

- **Describe the main bottlenecks for deploying and scaling up** the Transformer architecture.
- **Describe the ideas behind efficient algorithms** such as parameter quantization, (Q)LoRA, linear attention, and mixture of experts.
- Implement a simple retrieval-augmented generation pipeline.
- **Use LLM reasoning** to achieve better outputs for complex problems.



## Transformer: algorithmic complexity

Variables: hidden state dimension (D), sequence length (N), vocabulary (V).





📦 = model parameters, 🧮 = computed values

# LLM bottlenecks: Computational throughput

# Time complexity: $N \cdot D^2 + N^2 \cdot D + N \cdot V \cdot D$ $\begin{bmatrix} long context \\ N^2 \gg D^2 \end{bmatrix}$ $\begin{bmatrix} long context \\ only once \end{bmatrix}$

Computing **multi-head attention** is the main limitation when scaling to longer input contexts.

**Solution**: Efficient attention, mixture of experts

par.	min max	
N	512	128,000+
D	768	4,096+
V	30,000	100,000+
blocks	8	128+

Typical hyperparameter values

#### LLM bottlenecks: Disk space

Space complexity:  $V \cdot D + D^2$ 

(model parameters 📦)

MLPs typically take up majority of model parameters as they are repeated in each layer, unlike the embedding matrix.

stored only once

Solution: Parameter quantization, (Q)LoRA

par.	min	тах
N	512	128,000+
D	768	4,096+
V	30,000	100,000+
blocks	8	128+

Typical hyperparameter values

#### **LLM bottlenecks: Memory**

gets more important with growing input context size

# **Space complexity:** $N^2 + N \cdot D + N \cdot V$

(activations 🧮)

For loading a model to GPU memory, we need to account for both model parameters ( $\rightarrow$ fixed) and for computed activations ( $\rightarrow$ depend on the input context size).

Solution: Parameter quantization, (Q)LoRA

par.	min	тах
N	512	128,000+
D	768	4,096+
V	30,000	100,000+
blocks	8	128+

Typical hyperparameter values

# LLM bottlenecks: Training vs. inference

In terms of resource requirements, **pre-training > finetuning > inference**.

#### Dataset size:

- Finetuning: up to 100s GB.
- Pre-training: 10T tokens  $\rightarrow$  **50TB.**

#### **Memory requirements:**

- Inference: 2 bytes per parameter  $\rightarrow$  14 GB for a 7B model.
- Training: 2 bytes per parameter + 2 bytes per gradient + 12 bytes per optimizer weight (Adam) → 112 GB for a 7B model.

(Note that we typically want to run inference on consumer hardware.)

# **Efficient algorithms: overview**

Algorithm	What it targets	Training memory	Training speed	Inference memory	Inference speed
Parameter quantization	Reducing parameter size	-	-	<i>\ \</i>	1
Mixture of Experts (MoE)	Reducing active parameter count	-	<b>J J</b>	-	1
Linear attention	Faster attention using linear algebra	✓	1	1	1
FlashAttention	Faster attention using HW optimizations	<b>√</b> √	1	1	<b>\</b>
(Q)LoRA	Efficient model fine-tuning	1	1	-	-

# Model parameters are floating point numbers. How do we store them? And can we store them more efficiently?

- **float64** native Python (typically not used on GPUs)
- float32 baseline precision
- **float16** half precision
- For inference, float16 typically causes only little performance degradation.
- For training, float16 can lead to vanishing gradients → mixed-precision training.



#### bfloat16

- Developed by Google Brain
   (→"brain float").
- 16-bit float optimized for ML models.
- Greater dynamic range than float16 (→supports outlier weights) at the cost of lower precision.



Can we go further?

**Bob:** Let's use **int8**! It is an 8-bit numerical format. We'll save twice as much space compared to float16.



Alice: But we cannot just round the number, can we? Floats have a huge dynamic range. How do we squeeze them into (-127, 127)?

**Bob:** We don't. We only use the range of parameters that we have in our model weights.

Alice: The highest weight in my model can be still larger than 127.

**Bob:** But it will not be that large. Now we can squeeze them!



Alice: I am not convinced. What if some of my weights are outliers? Wouldn't this map most of the weights to zero (or near zero)?



👷 **Bob**: Ahem, right... Let's just clip the outliers?



Alice: Ok, and how will you pre-compute the range for activations? These are dynamic.







Alice: I give up. Just go on and quantize your model to 1-bit if it works for you.

👷 **Bob:** Can I? Awesome!

#### BitNet: Scaling 1-bit Transformers for Large Language Models

Hongyu Wang<sup>\*†‡</sup> Shuming Ma<sup>\*†</sup> Li Dong<sup>†</sup> Shaohan Huang<sup>†</sup> Huaijie Wang<sup>§</sup> Lingxiao Ma<sup>†</sup> Fan Yang<sup>†</sup> Ruiping Wang<sup>‡</sup> Yi Wu<sup>§</sup> Furu Wei<sup>†</sup><sup>◊</sup> <sup>†</sup> Microsoft Research <sup>‡</sup> University of Chinese Academy of Sciences <sup>§</sup> Tsinghua University https://aka.ms/GeneralAI

Quantization in practice:

#### • <u>GPTQ</u>

- Quantization method that uses a few extra tricks to go robustly beyond
   8-bit quantization (6-bit, 4-bit, 2-bit).
- <u>GGUF</u>
  - File format that performs block-wise quantization  $\rightarrow$  enables offloading parts of the model to CPU.

Models in Ollama are quantized by default to 4-bit.

#### **Mixture of experts (MoE)**

"Experts" = multiple feed-forward networks in each MLP layer of the Transformer.



## **Mixture of experts (MoE)**

We can choose a different

#### expert at each **layer**...

# ...and a different set of experts for each **token**.



class MoeLayer(nn.Module):
 def \_\_init\_\_(self, experts: List[nn.Module],
 super().\_\_init\_\_()
 assert len(experts) > 0
 self.experts = nn.ModuleList(experts)
 self.gate = gate
 self.args = moe\_args

# **Mixture of experts (MoE)**

#### Why it's a good idea?

- Individual experts can specialize to solving certain kinds of problems.
- Faster training for the same number of total parameters (→ we backpropagate only through the selected experts).
- **Faster inference** (although we still need to load the full model into memory).

 $8x7B \rightarrow 8$  standalone models, approx. equivalent to a 47B dense model (Not 56B, as other parts of the model than MLPs are shared.)



#### **Linear attention**

Let's revisit the  $O(N^2)$  attention. Can we do better?

```
\operatorname{Attention}(Q, K, V) = \operatorname{Softmax}(QK^T)V
```



#### **Linear attention**

We can use a cheaper way of computing similarity:



This also allows us to re-arrange matrix multiplications:

$$\begin{array}{c} Q \cdot K \in \\ R^{N \times N} \end{array} \quad \begin{array}{c} \hline \phi(Q)\phi(K)^T \\ \hline \sum_{i=1}^L \phi(Q)\phi(K_i)^T \end{array} V = \frac{\phi(Q)\phi(K)^T V)}{\phi(Q)\sum_{i=1}^L \phi(K_i)^T} \quad \begin{array}{c} K \cdot V \in \\ R^{D \times D} \end{array}$$

#### The resulting operations are $O(\mathbf{N} \cdot \mathbf{D}^2)$ , which is linear in sequence length.

Wo-hoo! 🎉



 r/MachineLearning · 2 yr. ago currentscurrents Top 1% Poster
 Why aren't we all using linear transformers?

There's a bunch of them - Linformer, Longformer, Performer, Nystromformer, Big Bird, etc etc. Plus a bunch more that have similar goals but don't necessarily aim for linear complexity, like memory-augmented transformers.

As far as I know, none of them have really seen much use. Even for image problems, which have very long input sizes, people are using regular transformers with tokenization schemes.

- Am I wrong? Are they actually good, or are at least some of them better than regular transformers?
- If not, what's wrong with them? Do they have lower accuracy? Are they slower to train?

...

Approximations such as linear attention can lead to degraded performance.

Meanwhile, hardware optimizations of the full attention mechanism can go a long way.



#### **FlashAttention**

#### FlashAttention re-arranges the operations to use

the underlying hardware more efficiently.

It is:

- Fast: 2-3x faster than baselines.
- Memory-efficient: linear in sequence length.
- **Exact**: uses no approximations.

It is now commonly implemented in LLM frameworks.



Memory Hierarchy with Bandwidth & Memory Size

## Low-rank adaptation (LoRA)

Alice: Remember the 112 GB of GPU memory I would need to finetune a 7B model? Where should I get that?

👷 **Bob**: Buy a new GPU cluster!

👰 Alice: ...

👷 **Bob**: Or use efficient finetuning methods?

#### Key ideas of LoRA:

- Instead of updating model weights directly, we can keep the Δ (the "diff") in a separate matrix.
- The  $\Delta$  matrix has a low rank  $\rightarrow$  it can be approximated by two smaller matrices.

#### Low-rank adaptation (LoRA)



For inference, LoRA weights are applied to the model  $\rightarrow$  no latency issues.

# QLoRA: finetuning quantized models

Can we apply LoRA to quantized models for even higher efficiency?

Yes: **QLoRA** does that with (1) careful parameter quantization and (2) optimizing CPU-GPU memory transfers.



## LoRA and QLoRA in practice

All things being equal, full finetuning > LoRA > QLoRA in terms of model performance.

 $\rightarrow$  The goal is to make the most out of the memory you have at your disposal.

#### **Starting points:**

- <u>https://huggingface.co/blog/4bit-transformer</u> <u>s-bitsandbytes</u>
- <u>https://www.kaggle.com/code/neerajmohan/</u> <u>finetuning-large-language-models-using-qlora</u>

Method	Bits	~ memory required for finetuning a 7B model
Full finetuning	32	120GB
Full finetuning	16	60GB
LoRA	16	16GB
QLoRA	8	10GB
QLoRA	4	6GB
QLoRA	2	4GB

https://github.com/hiyouga/LLaMA-Factory?tab =readme-ov-file#hardware-requirement

The knowledge saved in model parameters is lossy and incomplete, and the models tend to hallucinate when they don't know the answer. What can we do about it?

#### **Retrieval-augmented generation (RAG):**

- Retrieve additional documents related to user prompt (websites/ files / database entries / code / ...).
- Concatenate the most relevant documents with the prompt.
- Let the LLM generate the answer.





#### **Pros:**

- Allows searching private knowledge bases.
- Can reduce the amount of hallucinated facts.
- Makes use of long input context windows.

#### Cons:

- The model tries to account for all the documents, even the irrelevant ones → the answer may be worse than without RAG.
- Can give a false sense of "groundedness".





# <u>Wei et al., (2022)</u>: "The models struggle with math problems. What if we showed them how to do intermediate reasoning steps?"



#### <u>Kojima et al (2022)</u>: "Finding good examples of intermediate reasoning steps can be tricky. Do we need to do that at all?"



#### (c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 X

#### (d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

#### A: Let's think step by step.

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

Chain-of-thought prompting:

- ...is nowadays a standard prompting method.
- ...generally increases performance on problems requiring multi-step reasoning.
- ...does not require more than appending (a variant of) "Think step-by-step" to the prompt.



image source

Why does it work?

- The "emergent" ability can be learned from pretraining on code data (<u>Ma et al., 2024</u>; <u>Puerto et</u> <u>al., 2024</u>).
- The model may be just using extended inference time to perform more computation (<u>Pfau et al.</u>, <u>2024</u>).

Prompt: "How many of the first 6 digits of e are >5?"





"7 digits are greater than 5"

## **Test-time scaling**

CoT prompting has shown a new path: test-time scaling.

The simplest example of

test-time scaling:

- Generate multiple
   CoT paths for the
   given problem.
- Use majority voting to select the final answer.



#### **Test-time scaling**

More advanced extensions involve tree-like search techniques and a verifying the solution with a verifier model (e.g., code interpreter).



# Large Reasoning Models (LRMs)

**Goal:** instilling the reasoning process into the behavior of the LLM.

→ The model is trained on reasoning traces.

The traces include problem decomposition, intermediate steps, and back-tracking.

Where to get the reasoning traces?



# Examples of LRMs: OpenA1 o1, o3, o4 DeepSeek-R1 Gemini 2.5 Claude 3.7 Thinking

## **Training techniques for LRMs**

#### **Option #1: Pure reasoning-oriented RL**

- 1. Get a base LLM (no instruction tuning).
- 2. Run the model on datasets for code generation and math problem solving.
- 3. Perform RL training with rewards for accuracy (using a verifier) and format (rule-based).
- Natural emergence of test-time computation.
- Early unstable cold-start phase.
- Missing general capabilities.

Use to make DeepSeek-R1-Zero.





# **Training techniques for LRMs**

#### Option #2. Supervised fine-tuning (SFT), followed by RL.

- Get a pure RL-tuned LRM, generate outputs <u>A</u> and post-process them by human annotators.
- 2. Fine-tune an instruction-tuned LLM on <u>A</u>  $\rightarrow$  generate outputs <u>B</u>.

- addressing unstable cold start.

- 3. Fine-tune a base LLM on a mixture of <u>B</u> and non-reasoning training data.
  - addressing missing general capabilities
- 4. Continue training with reasoning-oriented as well as human preference RL.

The process was used to make DeepSeek-R1 (671B parameters); likely also used in OpenAI-o1.



# **Training techniques for LRMs**

#### **Option #3: Pure supervised fine-tuning**.

- Fine-tuning a small base LLM (e.g., Llama or Qwen models) on the combination of reasoning and non-reasoning outputs.
- Cheaper to run but still expensive to train (800k SFT samples)
- <u>NovaSky's Sky-T1</u> used only 17k SFT samples (\$450) to perform on par with o1

Used to make smaller ("distilled") DeepSeek models.



#### NPFL140 Large Language Models

#### **Efficiency & Advanced Techniques**

- LLMs can be implemented more efficiently:
  - **Lower memory footprint**: parameter quantization, (Q)LoRA
  - **Faster inference**: MoE, linear attention, FlashAttention
- Post-training performance of LLMs can improved at test time:
  - RAG for **extending the model knowledge**
  - CoT and reasoning for **tackling complex problems**

#### https://ufal.cz/courses/npfl140

## **Additional links**

- <u>Mixture of Experts Explained</u> HuggingFace blogpost
- <u>Understanding Reasoning LLMs</u> blogpost.
- <u>awesome-o1</u> list of resources related to OpenAI O1.
- <u>Scaling Test-time Compute</u> HuggingFace blogpost.