

Deep learning for GPU poor

Author: Anton Zhitomirsky

Contents

| | |
|--|----------|
| 1 Calculating statistics of a model | 2 |
| 1.1 Foundation model parameter count | 2 |
| 1.2 Floating Point System | 2 |
| 1.3 Floating Point Operations (FLOPs) | 3 |
| 2 Deep Learning At Scale | 4 |
| 3 Methods for reducing computational requirements | 6 |
| 3.1 Gradient Accumulation | 6 |
| 3.2 Gradient Checkpointing | 7 |

1 Calculating statistics of a model

1.1 Foundation model parameter count

Commonly when you look online for models, you will see their parameter counts. *Transformers are typically described by the #parameters* which impacts the computational requirements to store and run these models.

1.2 Floating Point System

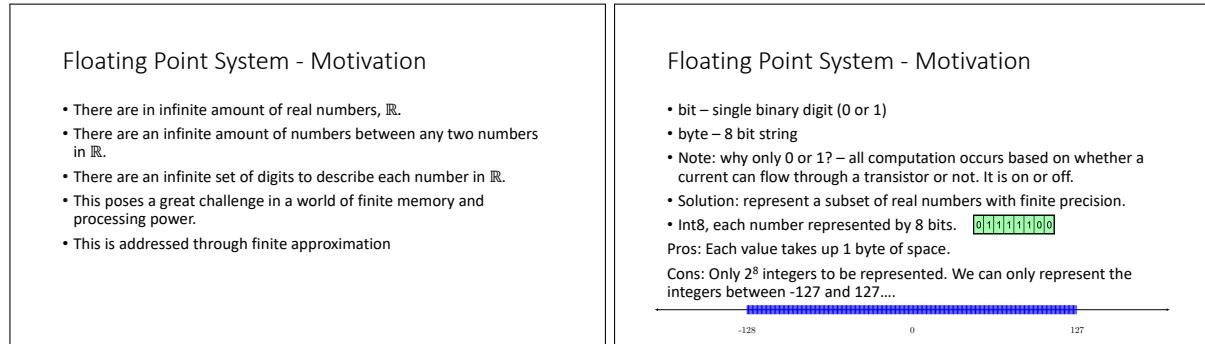
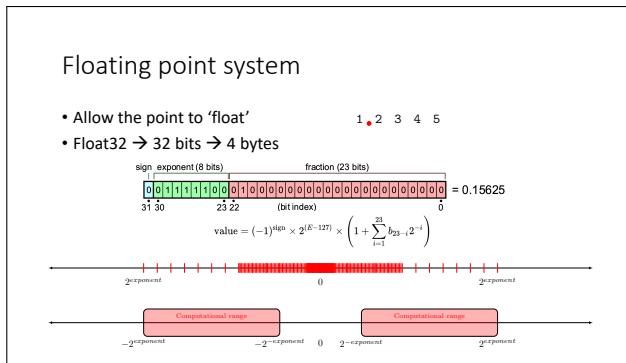
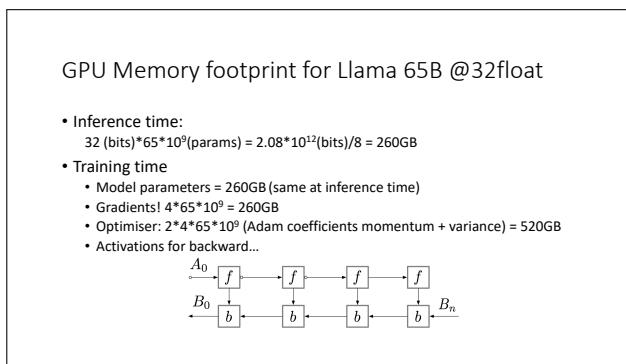


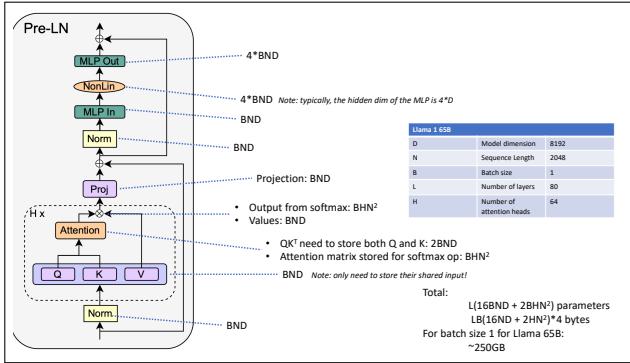
Figure 1: We store data in computers in bits with transistors. The more bits, the more precision.



- **Floating point** is a way to represent numbers in a computer where the decimal point can float.
- A certain section is allocated to represent the precision (actual number) the exponent and the sign.
- still has the problem of underflow or overflow.



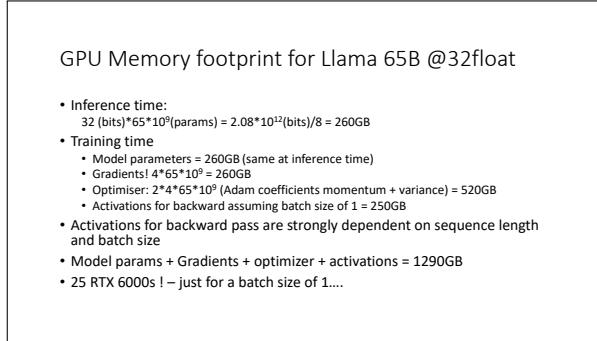
- At inference time, just to store the data it is already 260GB
- At training time, we also require gradients and optimiser values (momentum, etc) which can be 10x the size of the model and the activations.



So if we pass through the network, we need to store the inputs to every single layer that theres a change in the computational graph. Here:

- The first is the batch times the number of tokens or the number of patches times by that dimensionality.
- For the query key value, they all get the same value.
- The attention matrix requires query with key-values
- We also need to store the n^2 complexity of the attention matrix.

- Also we store the result of the soft max which is the batch times the number of heads times the sequence dimension.
- With multi-head attention we usually have another layer that learns how to combine these values together.
- All in all, we have 250GB for a batch size of 1.



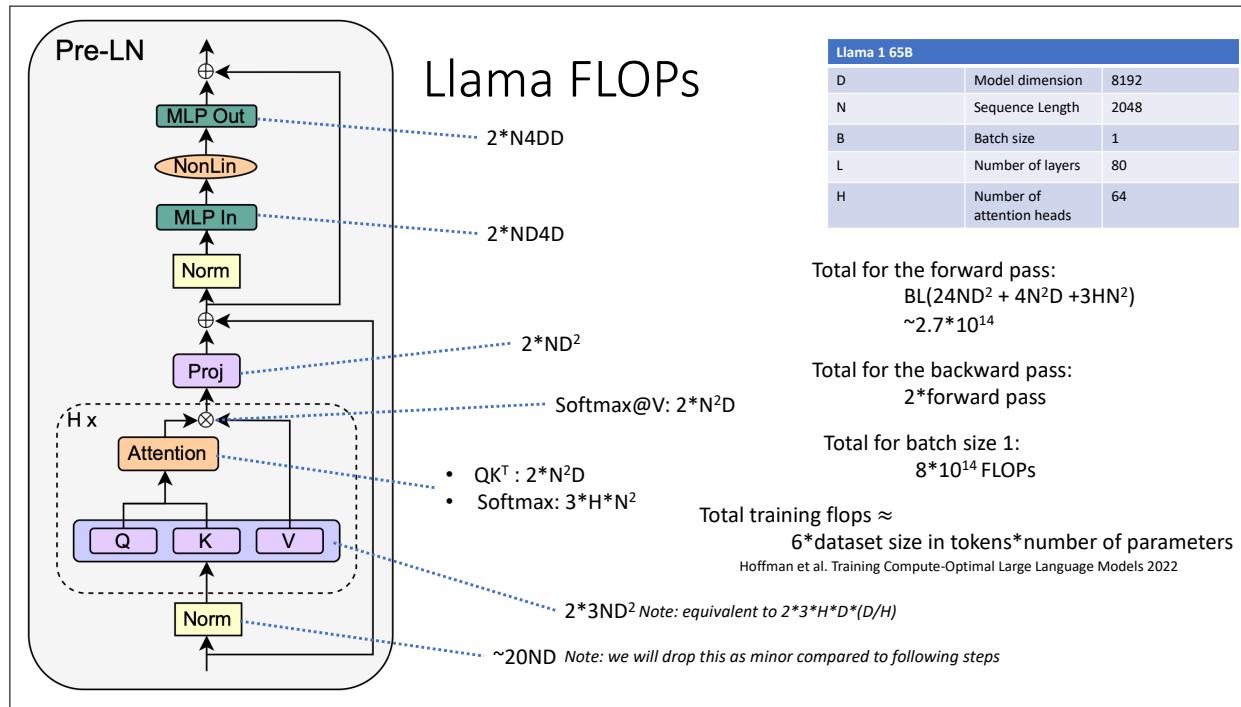
1.3 Floating Point Operations (FLOPs)

Floating point operations (FLOPs)

- Basic unit of computation
- e.g. addition, subtraction, multiplication, divide of 2 floating point numbers
- Given $w, x \in \mathbb{R}^n$
- $w + x$: n FLOPs
- $w^T x$: 2n FLOPs (multiplication + (n-1) addition)
- Given $W \in \mathbb{R}^{m \times p}$, $X \in \mathbb{R}^{p \times n}$
- WX : 2nmp

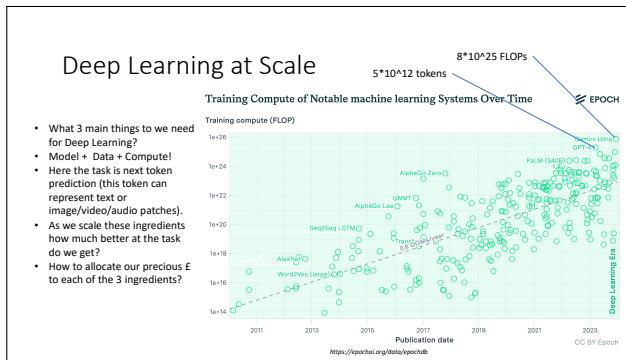
Computational requirements:

- given two vectors, if we add this will require n operations.
- For the dot products, there is $2n$ flops.
- For the matrix multiplication, there is $2nmp$ multiplications

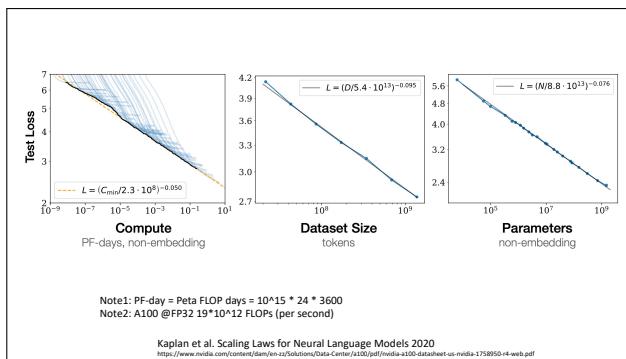


- the layer norm is roughly 20 operations for every feature and token
- for the query key values we have it 3 times, and 2 because its a dot product.
- The softmax: for every single query we dot product it for every single value.
- projection is same as query key values (but only 1)

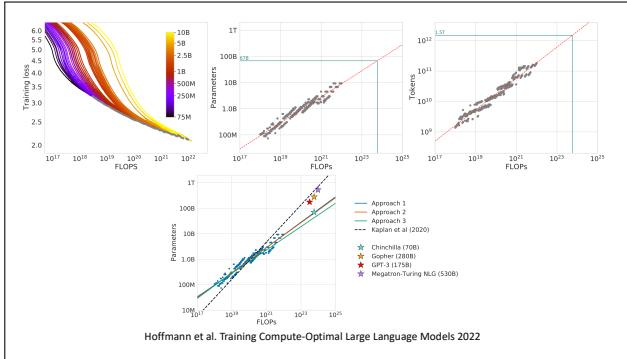
2 Deep Learning At Scale



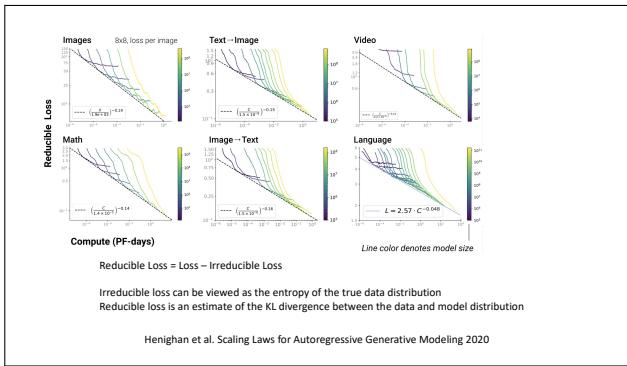
- For deeplearning we need model data and compute.
- At scale, you can't train the model many times for hyper-parameter tuning. A single model may take many dollars to train just once.



"When they have limitless compute and limitless datasets, they saw this really nice power law trend where it sticks very closely to the trend. This says that generally you can predict quite well, if you have limited compute and limited dataset size what loss you would get at the end." The values here were proven to be wrong because the smaller models were not fully optimized; they used the same learning rate schedule as they did for the larger models. Therefore, the line was too steep.



“openAI published findings about colour representing the size of the model, and they fit for given flops and compute budgets which model gets the lowest loss.” Here, we see that chinchilla is favoured since it prefers the number of data points over the model size.



Here we see that ‘deep learning has been solved’. “This is showing that this happens in all modalities; these losses can just keep going until we get to super human performance.”

Definition 2.1 (Reducible Loss). It is the actual loss - the irreducible loss, where irreducible loss is viewed as the true entropy of the data i.e. if I had a perfect fit of the data, what would be my loss even if I had a perfect fit?

It can be viewed as the KL divergence between the manifold of the dataspace and the manifold of the learned model.

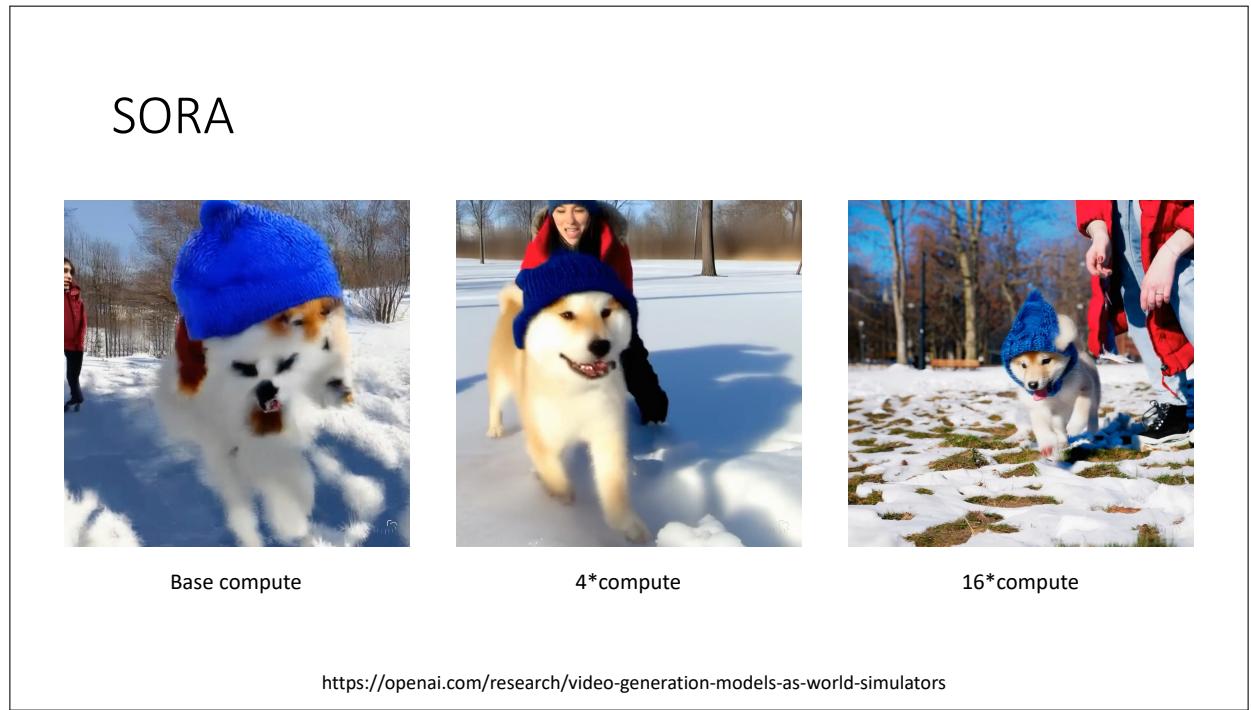
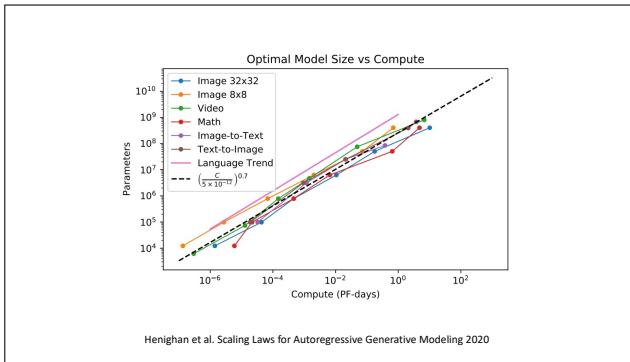
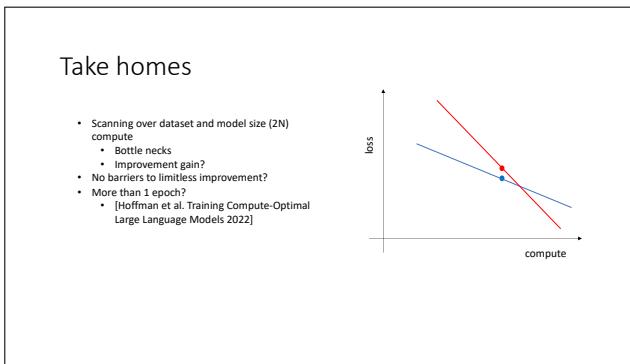


Figure 2: SORA AI just used more compute!



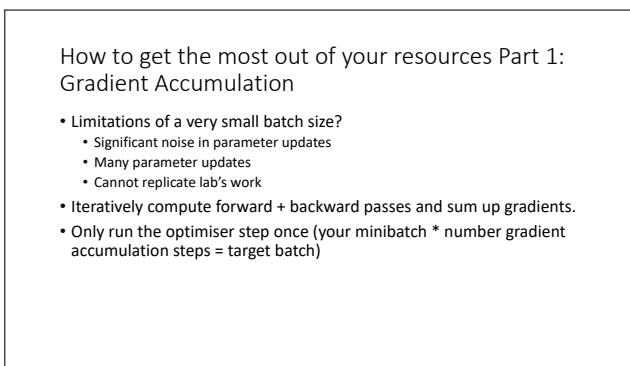
“When we plot all modalities the number of compute days vs the number of parameters for the optimal setting. There are correlation between what it is like to compute optimal size model for video and for language model”



- if blue is state of the art and red is my implementation,
- then we can do a scaling analysis and conclude that since the curve is steeper then it may perform better at high levels of compute.

3 Methods for reducing computational requirements

3.1 Gradient Accumulation

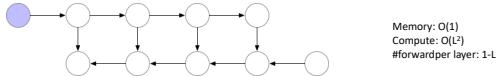


- As we increase the batch size our computation requirements increase
- If we use 1 batch, this is just Stochastic Gradient Descent so you're approximating the true gradient with one sample; your loss will be jumping around without convergence
- accumulate the gradients without updating the optimizer, then when the mini-batch size times the number of cumulative gradient accumulation steps equals the target batch then we do the step in the loss surface.*
- Useful if you want a larger batch size

3.2 Gradient Checkpointing

How to get the most out of your resources Part 1: Gradient Checkpointing

- Recall: The activations for a batch size 1 for Llama 65B is 172GB



- We could forget each activation during the forward pass and then recompute them during the backward pass
 - This adds significant more compute
- Compromise! Strategically select activations throughout computational graph, so only a fraction of the graphs need to be recomputed.

Chen et al. Training Deep Nets with Sublinear Memory Cost 2016
<https://medium.com/tensorflow/fitting-larger-networks-into-memory-583e3c75ff9>

- idea is to initially throw out the activation in the forward pass after we've done it, but in the backpropagation step we redo the entire forward computation step.

- This is computationally expensive, but the memory requirements stay low; it doesn't scale with the number of layers.
- The forward pass is much faster than the backwards pass
- Therefore, we select checkpoints from which we continue the forward pass.

Deep Learning at Scale for the GPU poor - Part 2

Harry Coppock



Only slides with this label Examinable is examinable.

-

Recap

- Tools for measuring scale – FLOPs + memory requirements
- Motivation – why scale matters
- Gradient checkpointing + accumulation

-

Plan

- Efficient Finetuning
- Prompt Engineering
- Mixture of experts

-

Finetuning problem setting

• After pretraining we may want to adapt our model to a specific downstream task.

• The pretrained weights Φ , are updated $\Phi + \Delta\Phi$ according to

• $\max_{\Phi} \sum_{\{(x,y)\sim Z\}} \sum_{t=1}^Y \log(P_{\Phi}(y_t|x, y_{<t}))$

Context-target tokens Autoregressive foundation model

• Here, every parameter in Φ is updated therefore need to store gradients + optimisers of every weight.

• For every downstream task a new model of size Φ is required + we need to relearn $\Delta\Phi$ parameters.

Examinable

Low Rank Adaption (LoRA)

• Motivation

- Li et al. demonstrated that the model-loss overparametrised space has a low intrinsic dimension/rank.
- What if $\Delta\Phi$ also has low intrinsic rank we can approximate $\Delta\Phi$ with Θ where $|\Theta| \ll |\Phi|$

• $\max_{\Theta} \sum_{\{(x,y)\sim Z\}} \sum_{t=1}^Y \log(P_{\Phi+\Theta}(y_t|x, y_{<t}))$

Only optimise Θ

Examinable

U et al. Measuring the intrinsic dimension of objective landscapes 2018
Goodfellow et al. Qualitatively Characterizing Neural Network Optimization Problems 2015
Aghajanyan et al. Intrinsic Dimensionality Explains the Effectiveness of Language Model Fine-Tuning 2020

Low Rank Adaption (LoRA)

• For a weight matrix in a pretrained model: $W_0 \in \mathbb{R}^{D \times D}$

• Following convention that at fine tuning we are learning: $\Delta W \in \mathbb{R}^{D \times D}$

• For a pretrained set up of $y = W_0 x$

• $y = XW_0 + X\Delta W \approx XW_0 + XL_1L_2$

Where, $L_1 \in \mathbb{R}^{D \times r}$ and $L_2 \in \mathbb{R}^{r \times D}$ and $r \ll D$
(note: matrix multiplication is distributive!)
(note: $L_1 L_2 \in \mathbb{R}^{D \times D}$)

• This can be done for every weight layer in the network

Hu et al. LoRA: Low-Rank Adaption of Large Language 2021

Examinable

Low Rank Adaption (LoRA)

• Empirically show that $r = 1$ or $r = 2$ yields good results for GPT-3 which has a full rank of 12,288!

• r will depend on the degree of shift between the pretraining dataset and the finetuning dataset

• Here, the number of trainable parameters, $|S|$ can be 0.01% of the pretrained parameters, $|\Phi|$.

Validation Accuracy

| | Weight Type | $r=1$ | $r=2$ | $r=4$ | $r=8$ | $r=64$ |
|--------------------------|----------------------|-------|-------|-------|-------|--------|
| WikiSQL($\pm 0.5\%$) | W_0 | 69.8 | 69.6 | 70.5 | 70.4 | 70.0 |
| | W_0, W_e | 73.4 | 73.3 | 73.7 | 73.8 | 73.5 |
| | W_0, W_k, W_v, W_o | 74.1 | 73.7 | 74.0 | 73.5 | 73.9 |
| MultiNLI ($\pm 0.1\%$) | W_0 | 90.7 | 90.9 | 91.1 | 90.7 | 90.7 |
| | W_0, W_e | 91.3 | 91.4 | 91.3 | 91.6 | 91.4 |
| | W_0, W_k, W_v, W_o | 91.2 | 91.7 | 91.7 | 91.5 | 91.4 |

WikiSQL MultiNLI-matched

Hu et al. LoRA: Low-Rank Adaption of Large Language 2021

Examinable

Low Rank Adaption (LoRA)

Examinable

- At inference time we combine: $\mathbf{W}_0 + \mathbf{L}_1 \mathbf{L}_2$
 - No additional inference time (no extra steps to compute both $\mathbf{X}\mathbf{W}_0$ and $\mathbf{X}\mathbf{L}_1 \mathbf{L}_2$)
 - We can always subtract $\mathbf{L}_1 \mathbf{L}_2$ when we are done and combine with a different set of parameters $\mathbf{L}_1^T \mathbf{L}_2^T$
- How much GPU memory does this save?
- No longer need to train 99.99% of parameters (number from Llama 65B in part1):
 - Gradients $4 * 65 * 10^9 = 260\text{GB}$
 - Optimizer: $2 * 4 * 65 * 10^9 = 520\text{GB}$
- Still need to store:
 - Model parameters = 260GB
 - Activations for backward = 172GB
- Assuming batch size of 1 we need 3 times less GPU resources.

Hu et al. LoRA: Low-Rank Adaption of Large Language 2021

QLoRA Background: k-bit quantisation

Examinable

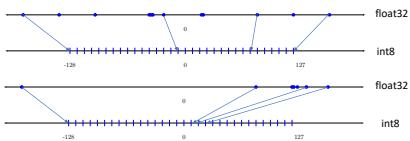
- Teaser: QLoRA enables finetuning of a 65B parameter transformer model on a single 48GB GPU + in 12 hours leads to SOTA opensource model performance.
- K-bit quantisation. Example float32 \rightarrow Int8

$$\mathbf{X}^{\text{Int8}} = \text{round}\left(\frac{127}{\text{absmax}(\mathbf{X}^{\text{FP32}})} \mathbf{X}^{\text{FP32}}\right) = \text{round}\left(c^{\text{FP32}} \cdot \mathbf{X}^{\text{FP32}}\right),$$

$$\text{dequant}(c^{\text{FP32}}, \mathbf{X}^{\text{Int8}}) = \frac{\mathbf{X}^{\text{Int8}}}{c^{\text{FP32}}} = \mathbf{X}^{\text{FP32}}$$
- To ensure full range of input is captured in Int8, input is rescaled to target range here [-127, 127]
- Where c^x is the quantisation constant

QLoRA: Background: k-bit quantisation

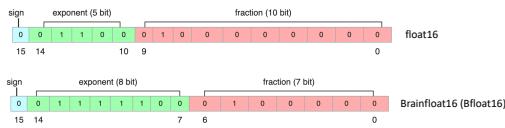
Examinable



Outliers result suboptimal allocation of quantized bins
Chunking! Slice $\mathbf{W} \in \mathbb{R}^{D \times D}$ into n chunks of size C where $n = \frac{D \times D}{C}$
Quantise each block separately with unique c^x

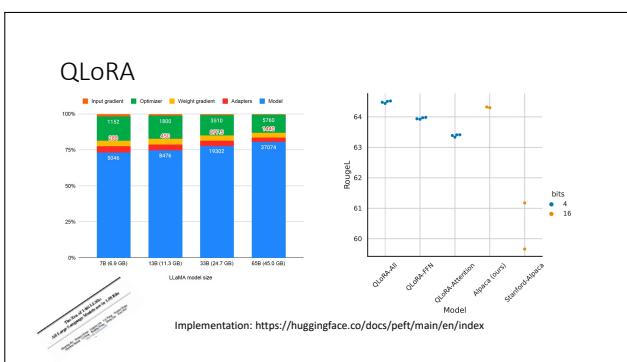
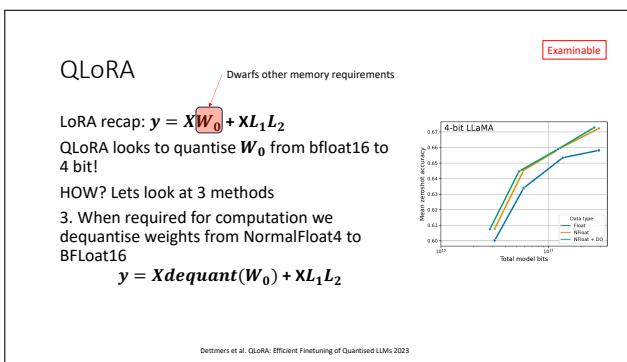
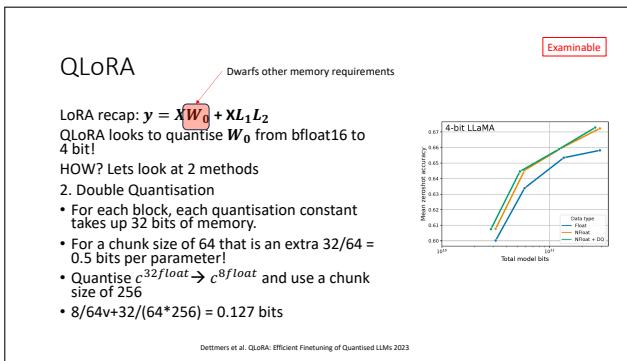
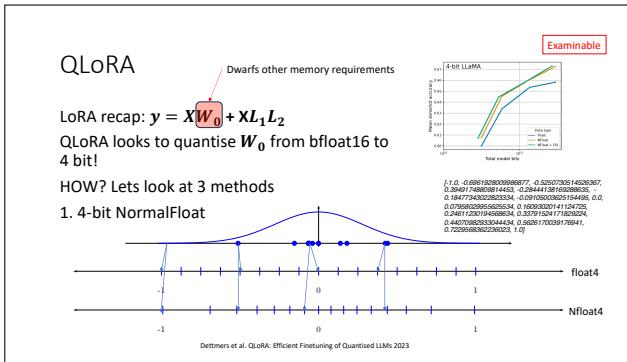
QLoRA: BrainFloat

Examinable



- Deep learning models are more sensitive to underflow + overflow than to precision
- Hardware requirements scale quadratically with fraction but not exponent

[developed at google] <https://cloud.google.com/tpu/docs/bfloat16>

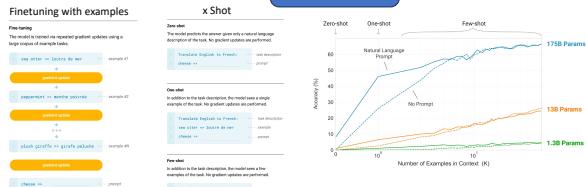


Prompt Engineering – What is it?

- Performance of autoregressive models at inference time varies greatly depending on the input/prompt.
 - $P_\Phi(x_t | x_{ where $x_{ is the input/context$$
 - Here, P_Φ , selects the response x_t which has the highest probability of following the prompt given the training data.
 - Therefore, if we are looking to elicit a certain response, x_t , we need to spend some time thinking about what $x_{ would maximise the chances of seeing a desirable x_t .$



Prompt Engineering



Brown et al. Language Models are Few-Shot Learners 2020

Prompt Engineering – What is it?

- Some call it ‘in context learning’ – learning from examples
 - However, this wording is confusing as no weights are updated.
 - Another interpretation:
 - Not learning *how* to perform the task.
 - Rather, instructing where in the model’s learnt data likelihood space we want to sample from.

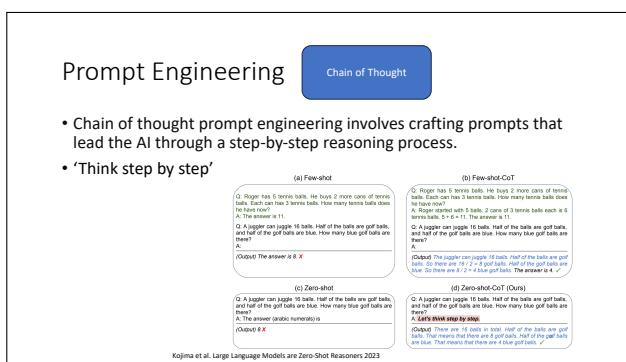
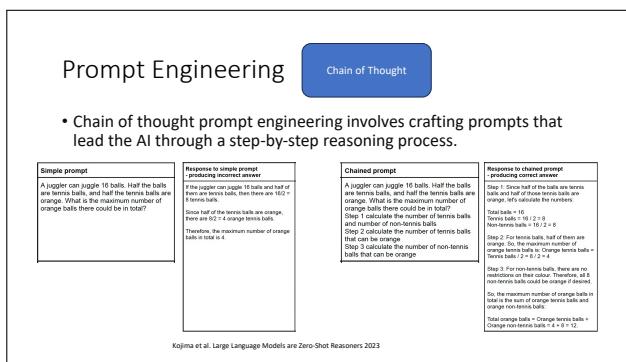
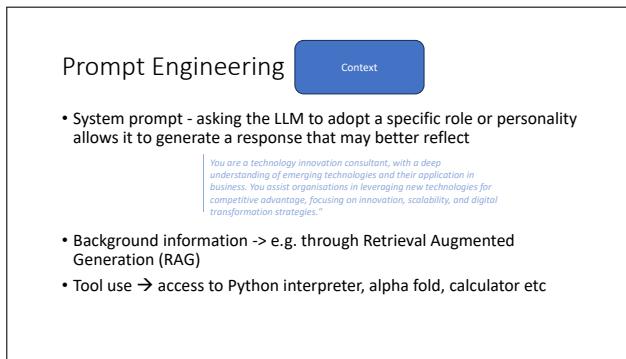
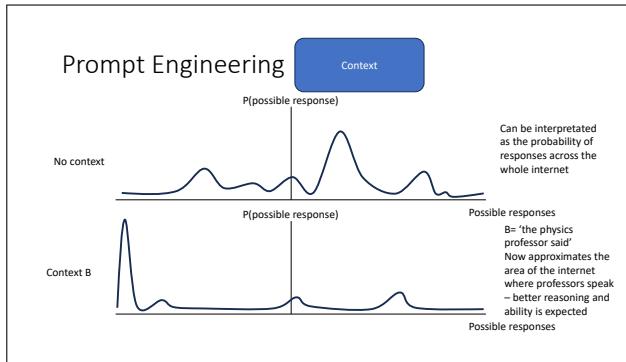
Reynolds et al. Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm 2021

The Dynamics of Language

- What function are these autoregressive models attempting to fit?
 - NOT just grammar + ngram probabilities
 - Language has been designed to fit the world we live in.
 - A good fit requires modelling how the real world affects language.
 - Modelling language is as difficult as modelling everything that influences it!

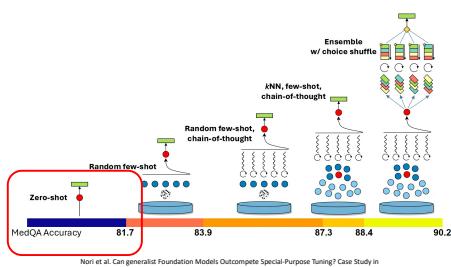
| | | |
|--|---|---|
| <p><i>The ball was filled with helium gas. When I dropped the ball it [fell, floated, ...]</i></p> | <p><i>dropped the ball it [fell, floated, ...]</i></p> | <p><i>The scientist was surprised to see that when she dropped the ball it [fell, floated, ...]</i></p> |
| <p><i>When the boy dropped the ball it [fell, floated, ...]</i></p> | <p><i>When the boy dropped the ball it [fell, floated, ...]</i></p> | <p><i>When the boy dropped the ball it [fell, floated, ...]</i></p> |

Reynolds et al. Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm 202

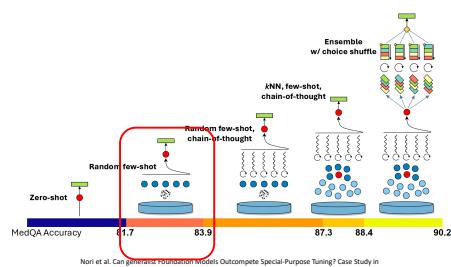


Nice practical guide + tips and tricks:
<https://platform.openai.com/docs/guides/prompt-engineering>

Medprompt case study



Medprompt case study



Medprompt case study

