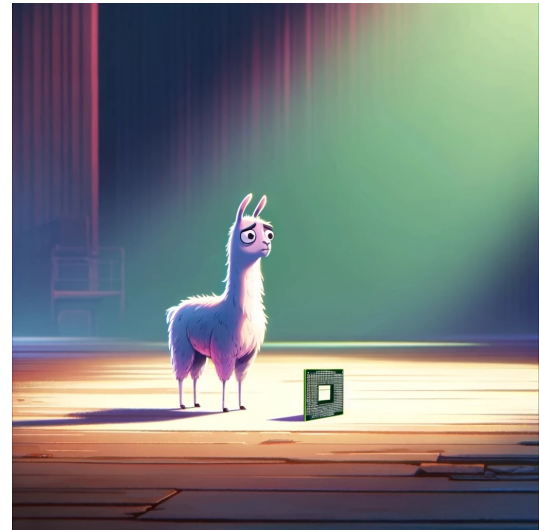


Deep Learning at Scale for the GPU poor - Part 1

Harry Coppock



Plan

- Tools for measuring scale – FLOPs + memory requirements
- Motivation – why scale matters
- How to get the most out of your resources
- Finetuning vs Prompt Engineering
- Foundation model safety/evals

Foundation model parameter count

- Transformers are typically described by the #parameters.
- But how does this translate to the computational requirements to store and run these models?

Model Name	# Parameters (Billion)
GPT3	175
Gopher	280
Llama 1 + 2	6-70
Mixtral	8*7
PaLM	540
GPT4	1700 (rumour)

Side note: A Fruit Fly has 100 million synapses

Side note: The Human brain has 100 trillion synapses

Floating Point System - Motivation

- There are an infinite amount of real numbers, \mathbb{R} .
- There are an infinite amount of numbers between any two numbers in \mathbb{R} .
- There are an infinite set of digits to describe each number in \mathbb{R} .
- This poses a great challenge in a world of finite memory and processing power.
- This is addressed through finite approximation

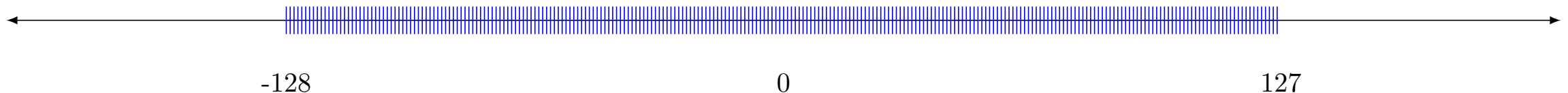
Floating Point System - Motivation

- bit – single binary digit (0 or 1)
- byte – 8 bit string
- Note: why only 0 or 1? – all computation occurs based on whether a current can flow through a transistor or not. It is on or off.
- Solution: represent a subset of real numbers with finite precision.
- Int8, each number represented by 8 bits.

0	1	1	1	1	1	0	0
---	---	---	---	---	---	---	---

Pros: Each value takes up 1 byte of space.

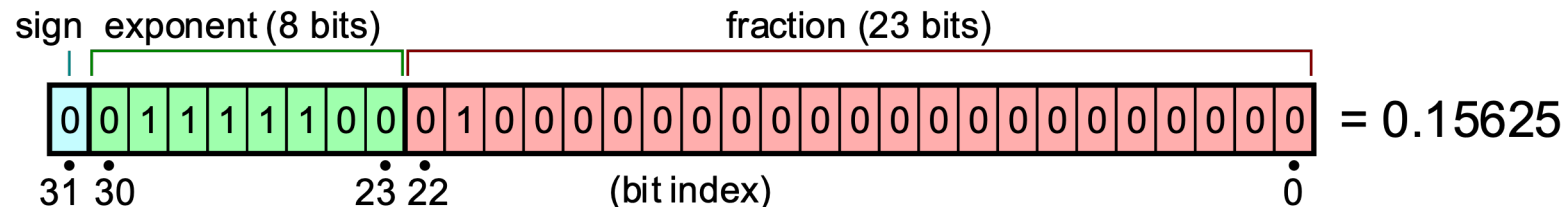
Cons: Only 2^8 integers to be represented. We can only represent the integers between -127 and 127....



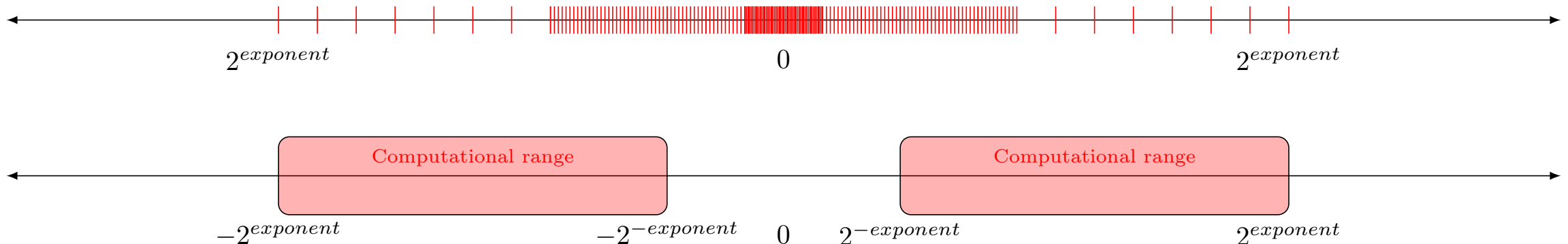
Floating point system

- Allow the point to 'float'
- Float32 → 32 bits → 4 bytes

1 . 2 3 4 5



$$\text{value} = (-1)^{\text{sign}} \times 2^{(E-127)} \times \left(1 + \sum_{i=1}^{23} b_{23-i} 2^{-i} \right)$$



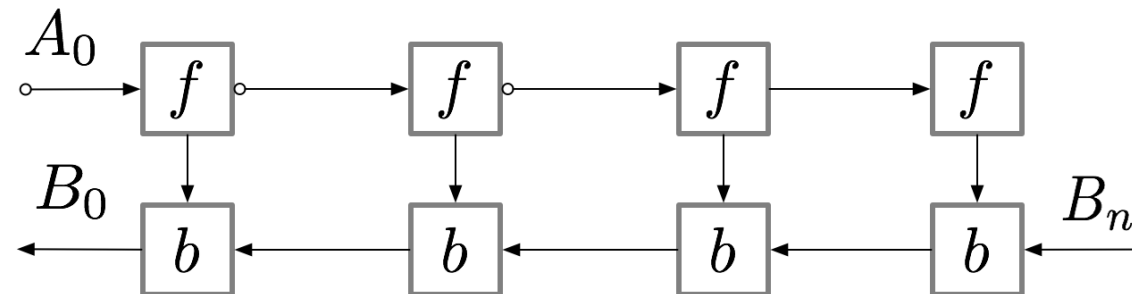
GPU Memory footprint for Llama 65B @32float

- Inference time:

$$32 \text{ (bits)} * 65 * 10^9 \text{ (params)} = 2.08 * 10^{12} \text{ (bits)} / 8 = 260 \text{ GB}$$

- Training time

- Model parameters = 260GB (same at inference time)
- Gradients! $4 * 65 * 10^9 = 260 \text{ GB}$
- Optimiser: $2 * 4 * 65 * 10^9$ (Adam coefficients momentum + variance) = 520GB
- Activations for backward...



GPU Memory footprint for Llama 65B @32float

- Inference time:
 $32 \text{ (bits)} * 65 * 10^9 \text{ (params)} = 2.08 * 10^{12} \text{ (bits)} / 8 = 260\text{GB}$
- Training time
 - Model parameters = 260GB (same at inference time)
 - Gradients! $4 * 65 * 10^9 = 260\text{GB}$
 - Optimiser: $2 * 4 * 65 * 10^9$ (Adam coefficients momentum + variance) = 520GB
 - Activations for backward assuming batch size of 1 = 250GB
- Activations for backward pass are strongly dependent on sequence length and batch size
- Model params + Gradients + optimizer + activations = 1290GB
- 25 RTX 6000s ! – just for a batch size of 1....

Floating point operations (FLOPs)

- Basic unit of computation
- e.g. addition, subtraction, multiplication, divide of 2 floating point numbers

Examples

Given $\mathbf{w}, \mathbf{x} \in \mathbb{R}^n$

- $\mathbf{w} + \mathbf{x}$:
n FLOPs
- $\mathbf{w}^T \mathbf{x}$:
2n FLOPs (multiplication + (n-1) addition)

Given $\mathbf{W} \in \mathbb{R}^{m \times p}$, $\mathbf{X} \in \mathbb{R}^{p \times n}$

- \mathbf{WX} :
2nmp

Pre-LN

Llama FLOPs

Llama 1 65B		
D	Model dimension	8192
N	Sequence Length	2048
B	Batch size	1
L	Number of layers	80
H	Number of attention heads	64

Total for the forward pass:

$$BL(24ND^2 + 4N^2D + 3HN^2)$$

$$\sim 2.7 * 10^{14}$$

Total for the backward pass:

2*forward pass

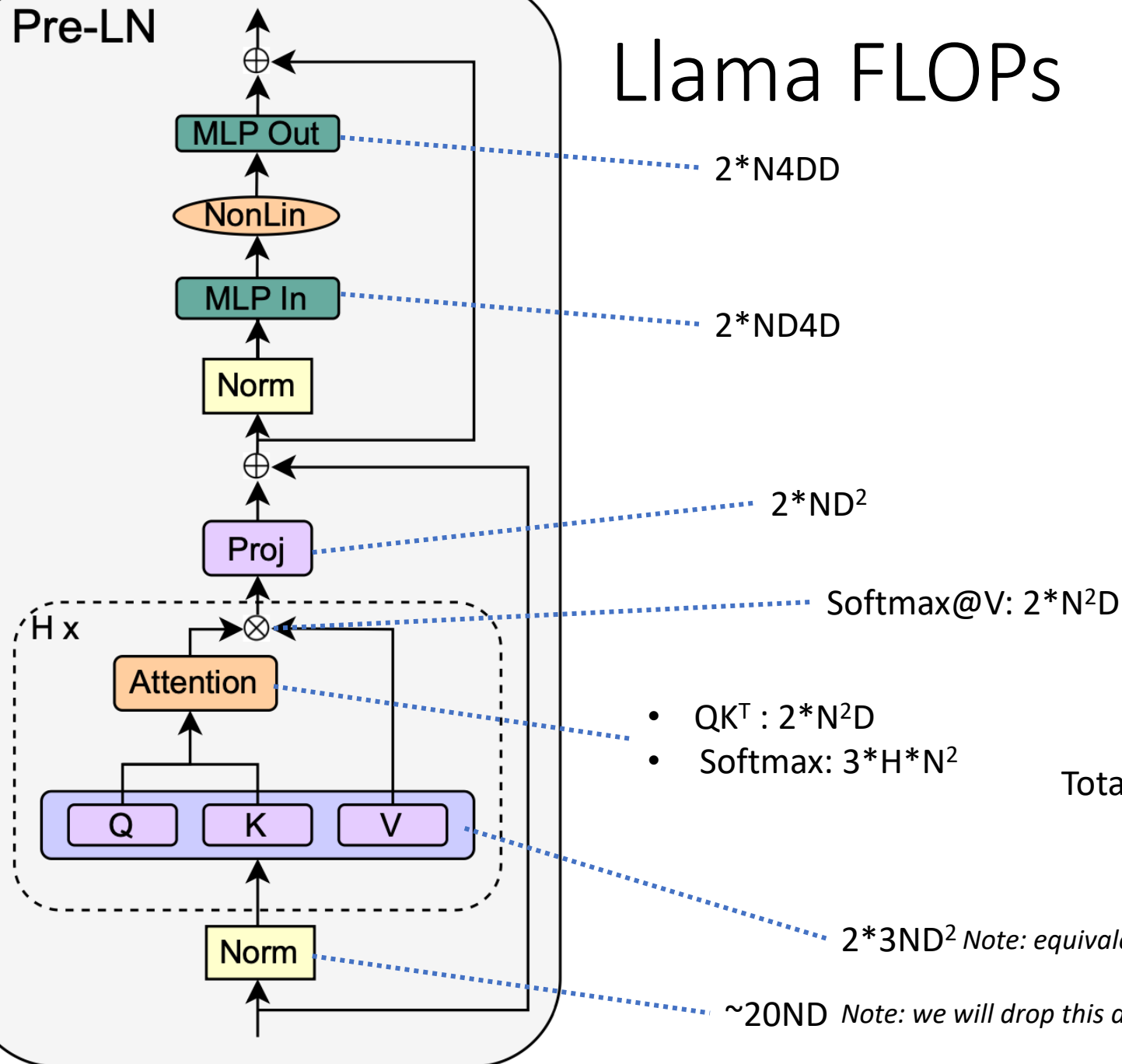
Total for batch size 1:

$$8 * 10^{14} \text{ FLOPs}$$

Total training flops \approx

$$6 * \text{dataset size in tokens} * \text{number of parameters}$$

Hoffman et al. Training Compute-Optimal Large Language Models 2022

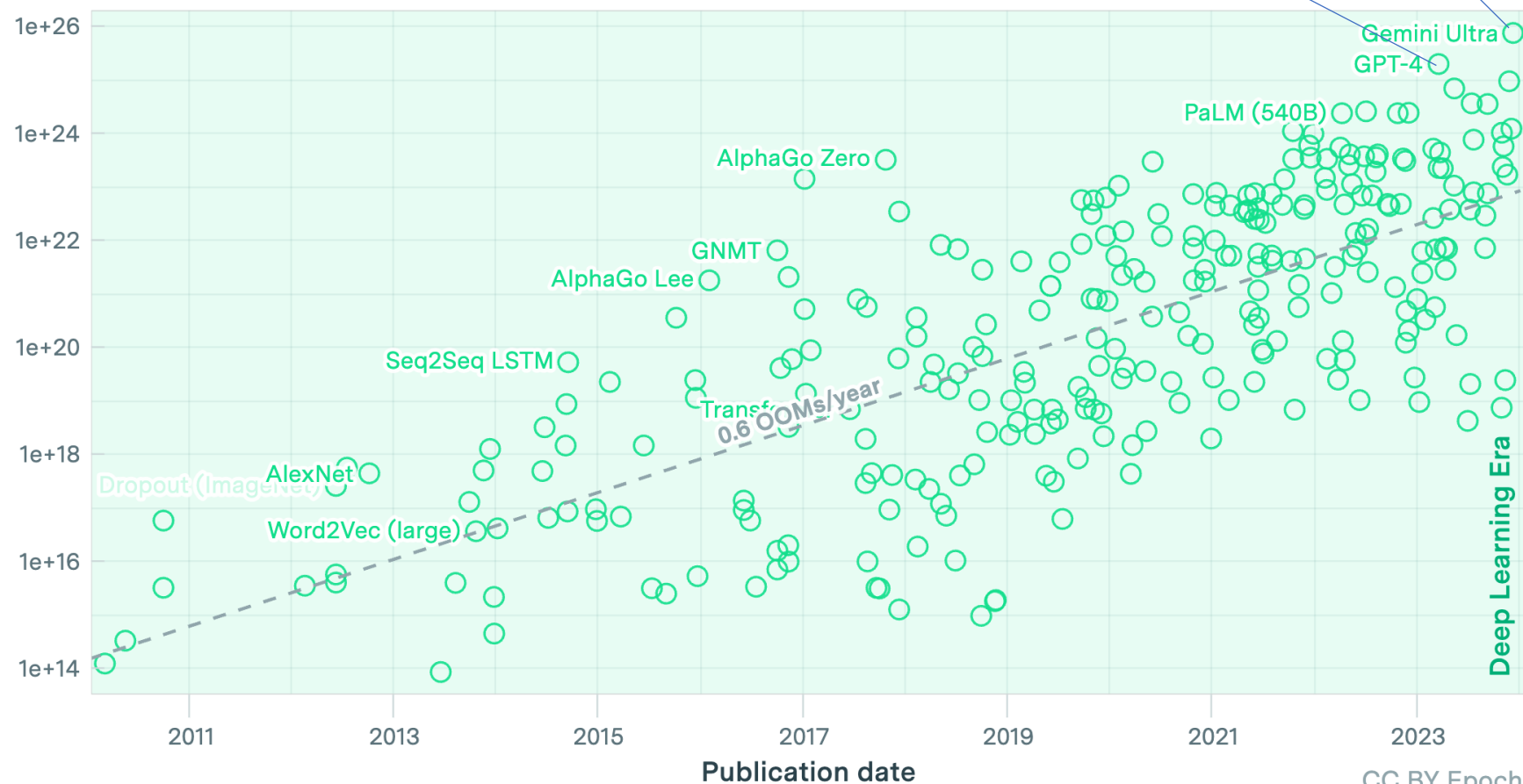


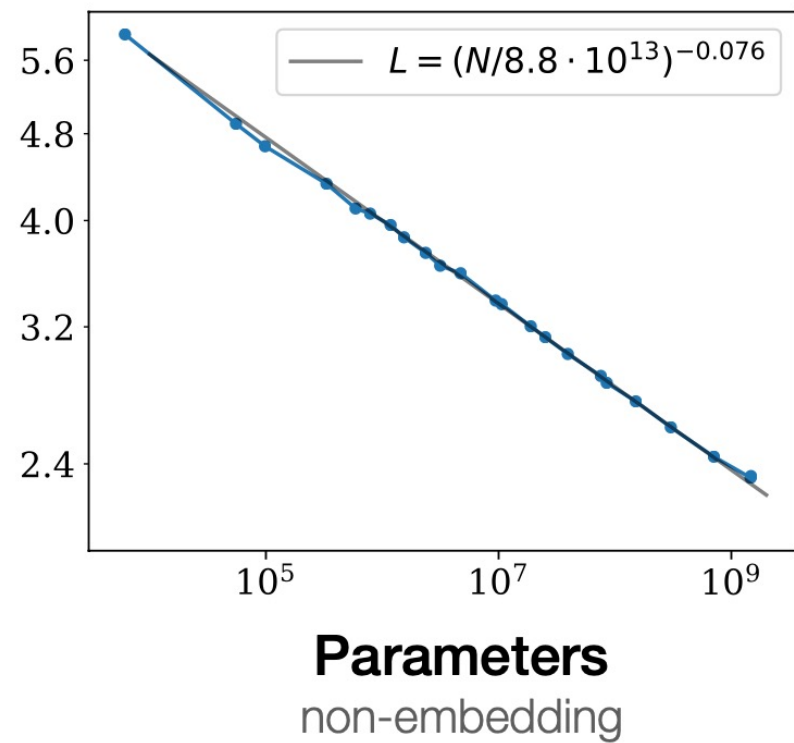
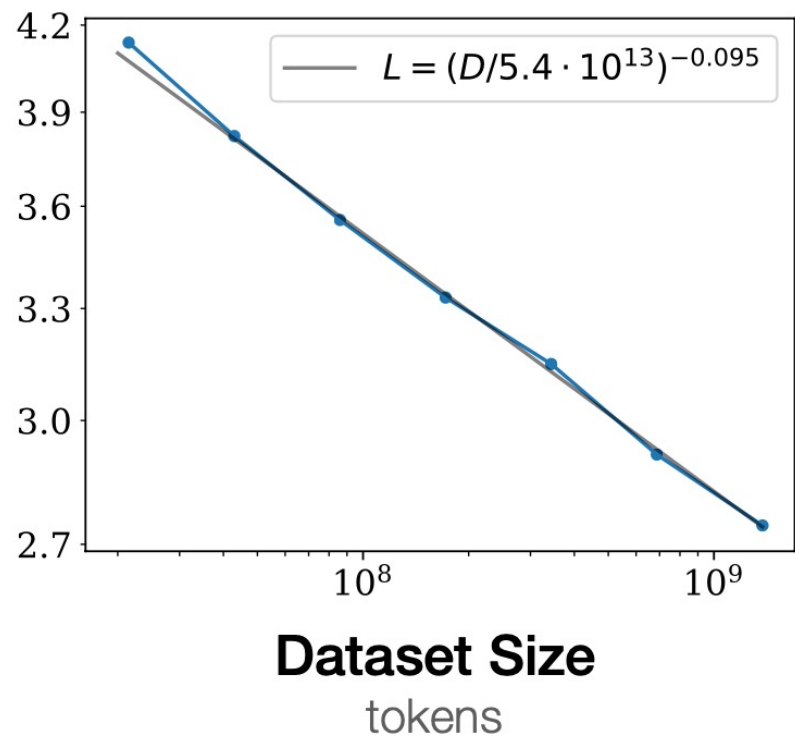
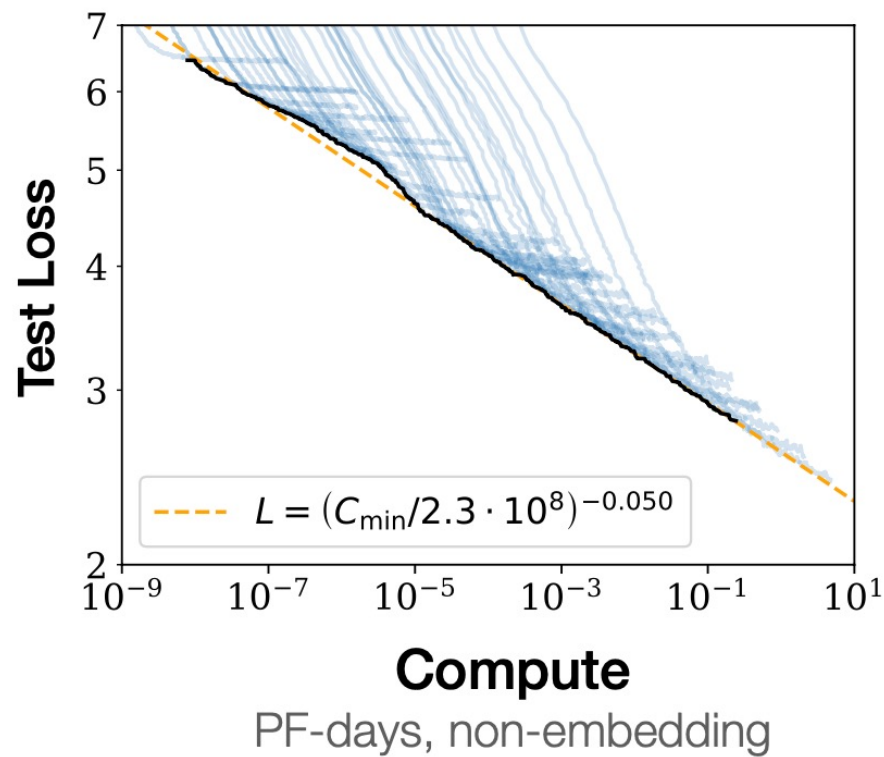
Deep Learning at Scale

- What 3 main things do we need for Deep Learning?
- Model + Data + Compute!
- Here the task is next token prediction (this token can represent text or image/video/audio patches).
- As we scale these ingredients how much better at the task do we get?
- How to allocate our precious £ to each of the 3 ingredients?

Training Compute of Notable machine learning Systems Over Time

Training compute (FLOP)



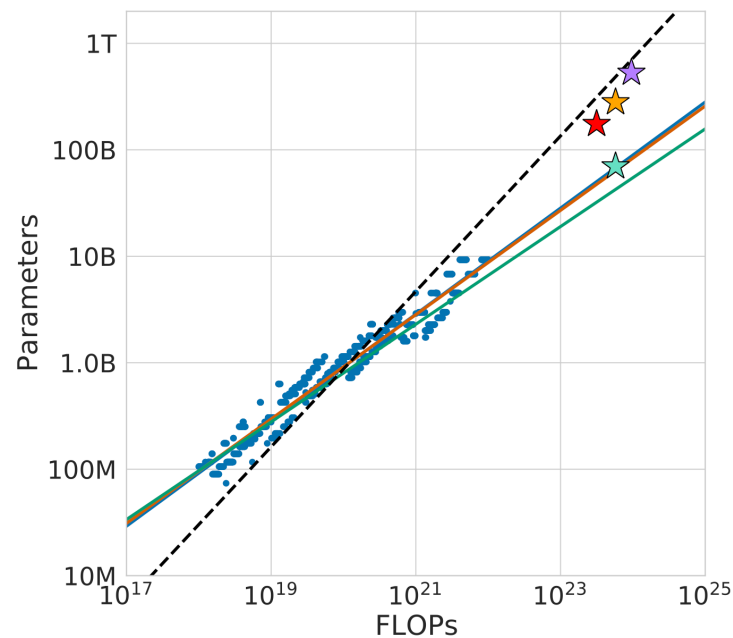
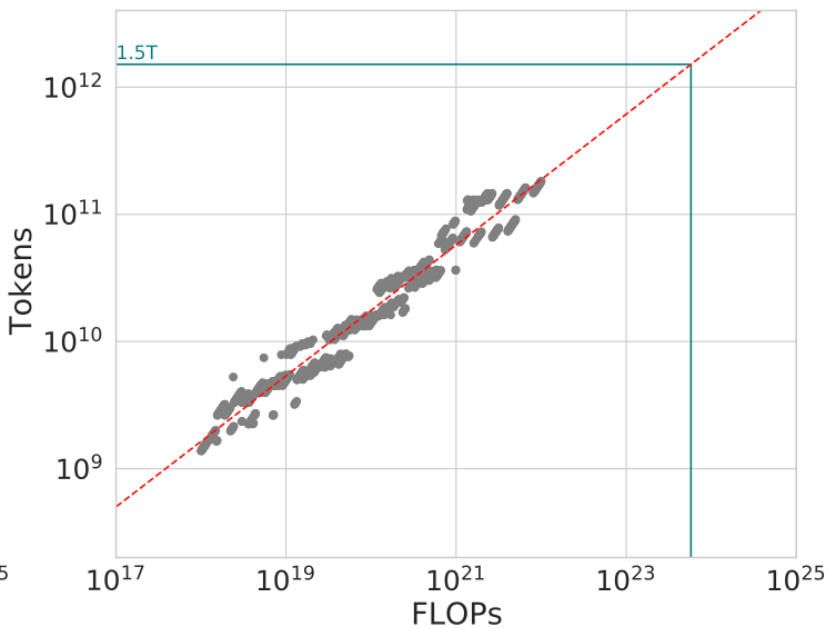
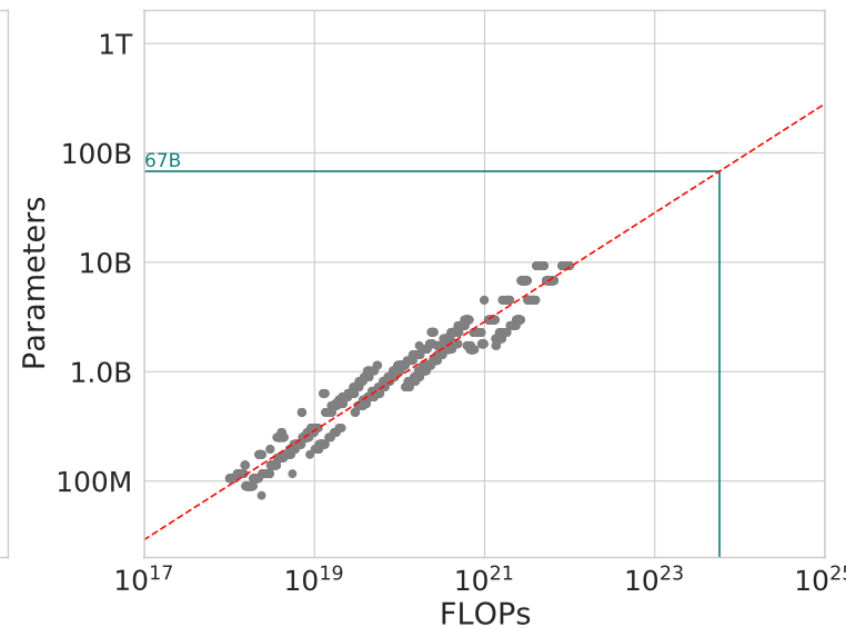
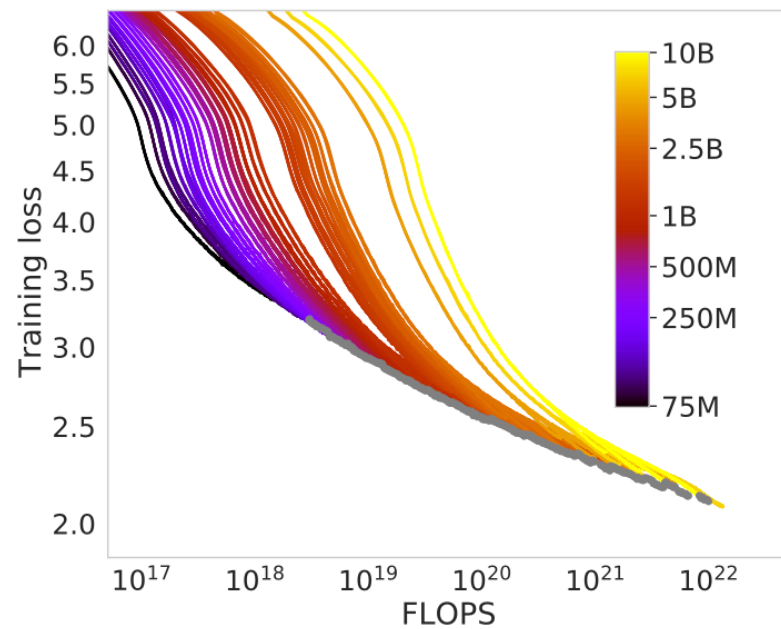


Note1: PF-day = Peta FLOP days = $10^{15} \cdot 24 \cdot 3600$

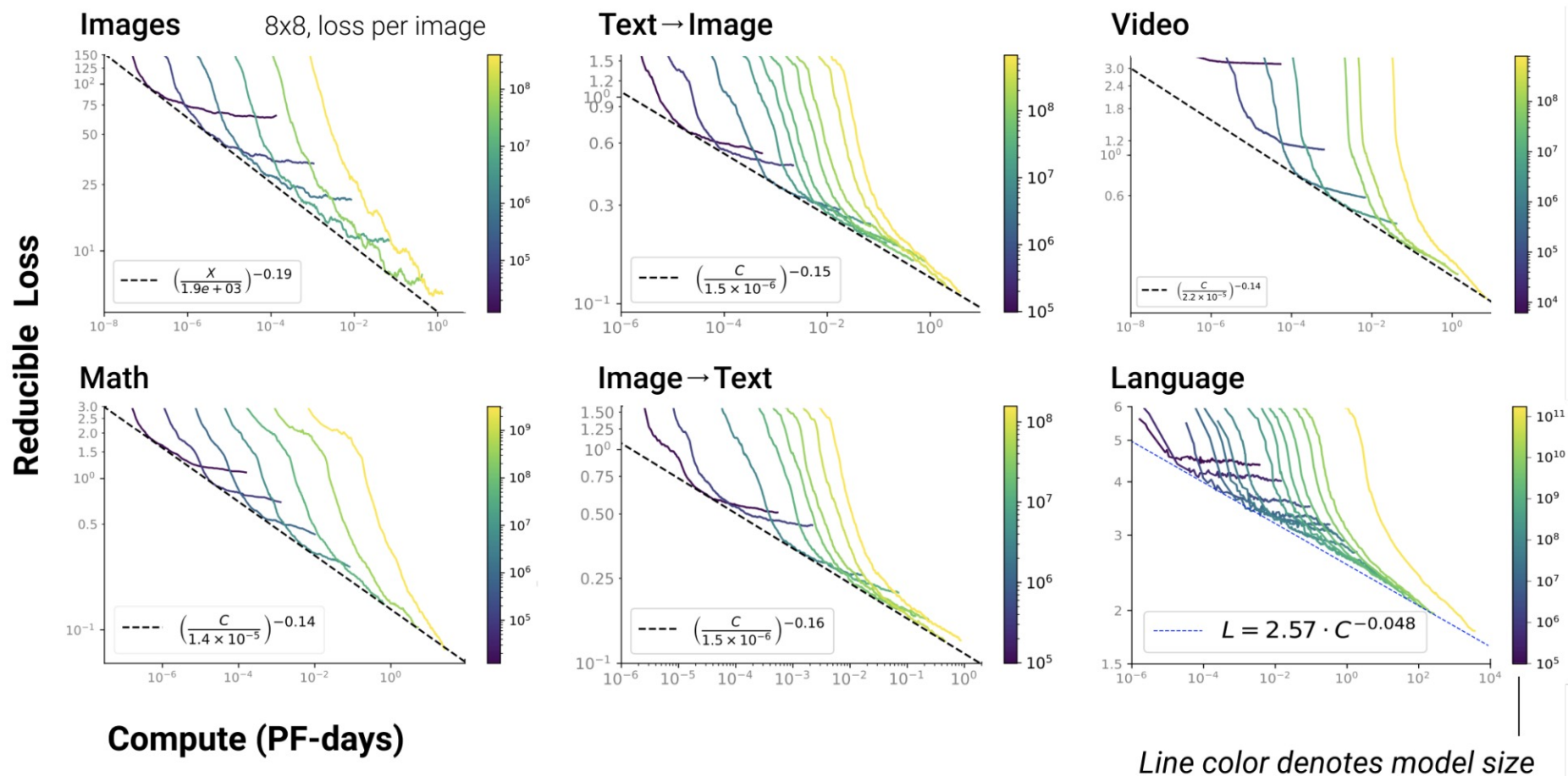
Note2: A100 @FP32 $19 \cdot 10^{12}$ FLOPs (per second)

Kaplan et al. Scaling Laws for Neural Language Models 2020

<https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-us-nvidia-1758950-r4-web.pdf>



- Approach 1
- Approach 2
- Approach 3
- - - Kaplan et al (2020)
- ★ Chinchilla (70B)
- ★ Gopher (280B)
- ★ GPT-3 (175B)
- ★ Megatron-Turing NLG (530B)



Reducible Loss = Loss – Irreducible Loss

Irreducible loss can be viewed as the entropy of the true data distribution

Reducible loss is an estimate of the KL divergence between the data and model distribution

SORA



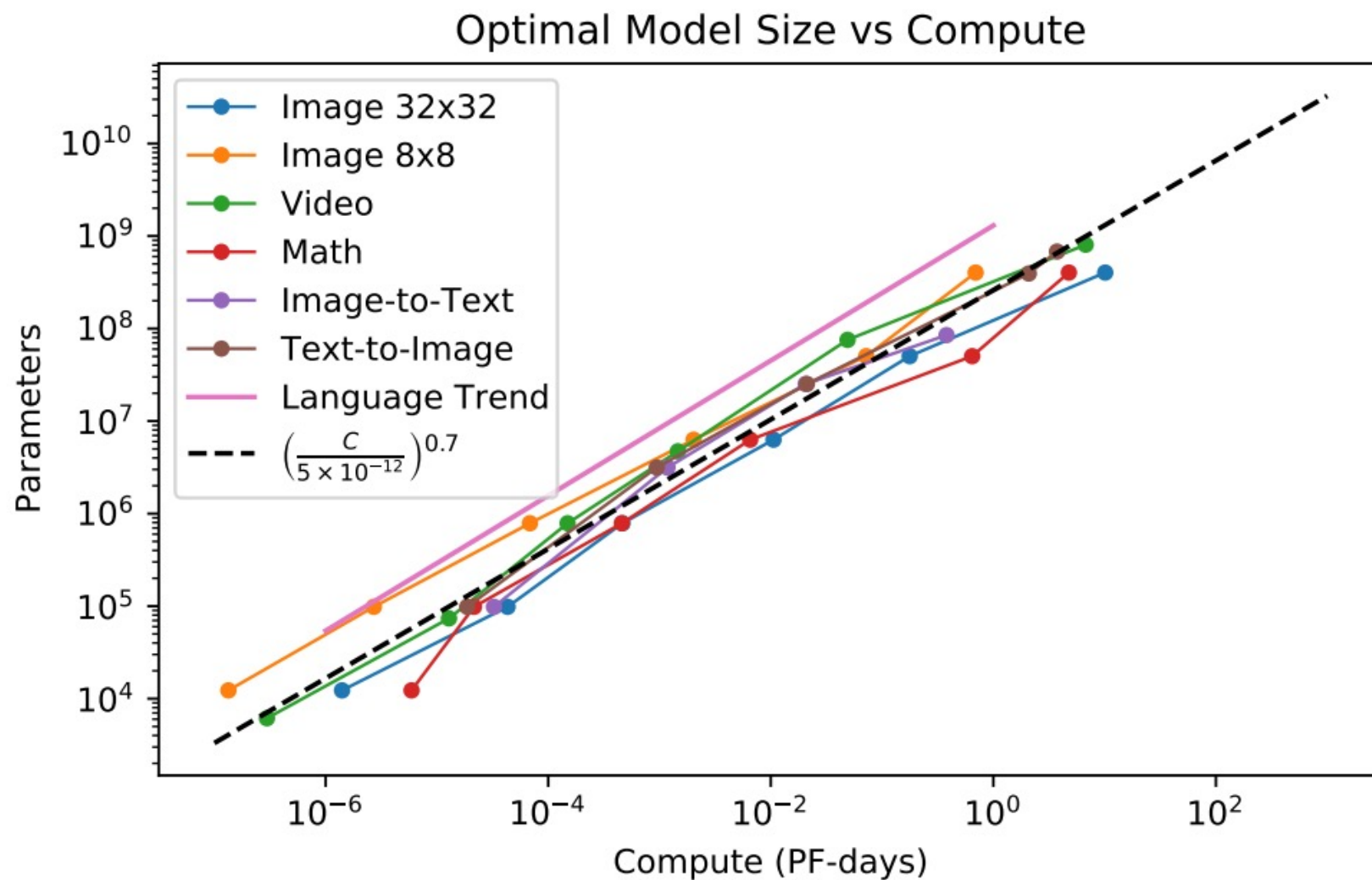
Base compute



4*compute

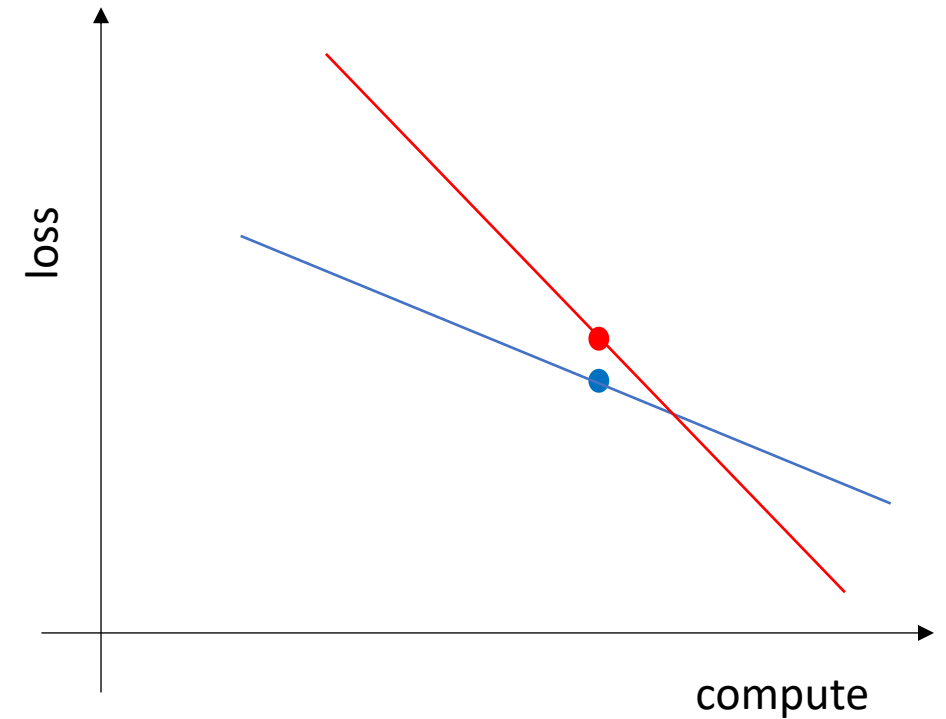


16*compute



Take homes

- Scanning over dataset and model size ($2N$) compute
 - Bottle necks
 - Improvement gain?
- No barriers to limitless improvement?
- More than 1 epoch?
 - [Hoffman et al. Training Compute-Optimal Large Language Models 2022]

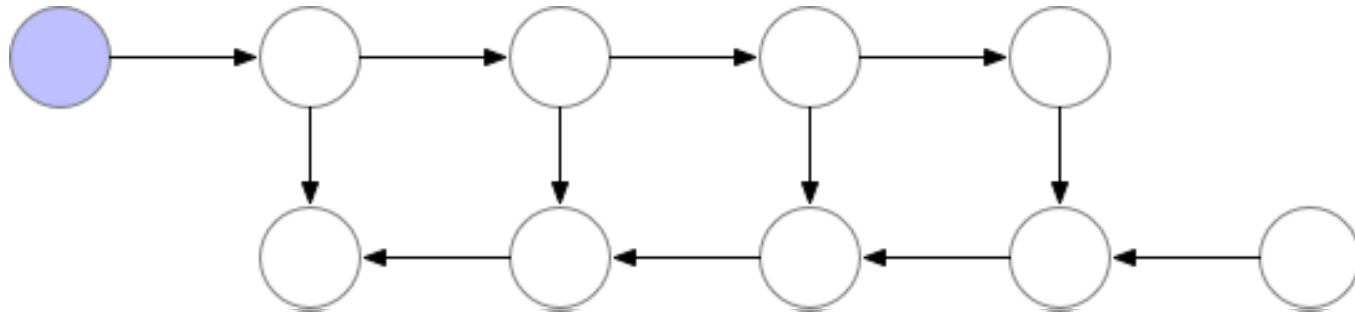


How to get the most out of your resources Part 1: Gradient Accumulation

- Limitations of a very small batch size?
 - Significant noise in parameter updates
 - Many parameter updates
 - Cannot replicate lab's work
- Iteratively compute forward + backward passes and sum up gradients.
- Only run the optimiser step once (your minibatch * number gradient accumulation steps = target batch)

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



Memory: $O(1)$
Compute: $O(L^2)$
#forwardper layer: $1-L$

- 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.

Next week

- Finetuning at Scale
 - Prompt Engineering
 - Foundation model Evals!
-
- Note on examinable content