



# TRAINING LLMS: SCALING AND EFFICIENT FINE-TUNING

Lauri Seppäläinen



# OUTLINE

Part 1: Kaplan & al. *Scaling Laws for Neural Language Models* (2020)

Part 2: Han & al. *Parameter-Efficient Fine-Tuning for Large Models: A Comprehensive Survey* (2024)



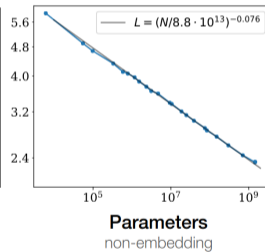
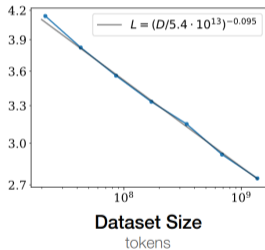
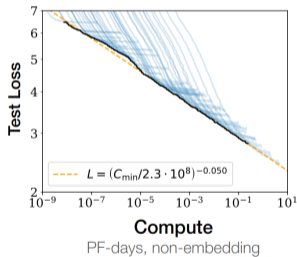
# SO YOU WANT TO TRAIN AN LLM FROM SCRATCH

- What architecture should you use?
- How big should your model be?
- How much data do you need?
- How long should you train for?





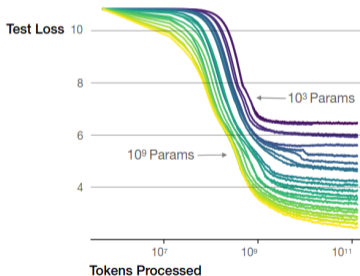
# POWER LAWS



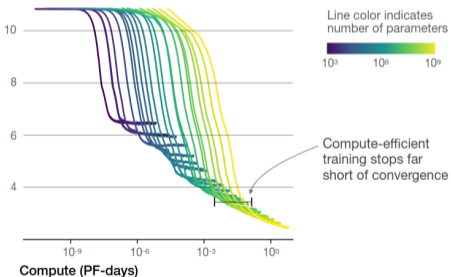


# MODEL SIZE

Larger models require **fewer samples** to reach the same performance

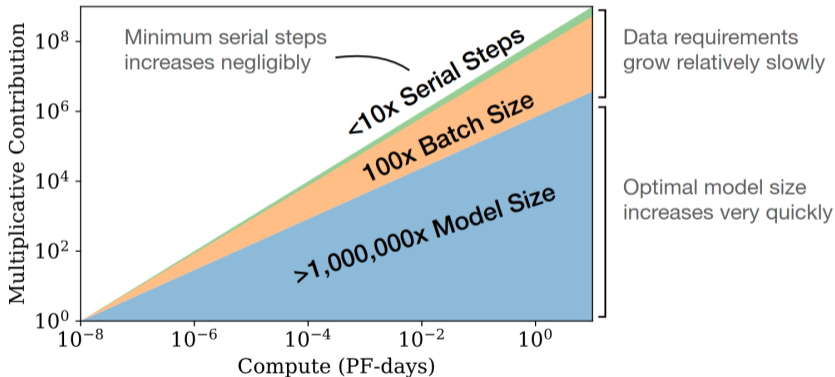


The optimal model size grows smoothly with the loss target and compute budget





# OPTIMAL COMPUTE ALLOCATION





# OTHER CONSIDERATIONS

- Model architecture and optimization hyperparameters seem to have minimal impact; model scale much more important
- Training until convergence is inefficient
- Overfitting can result both from training batch size and parameter count



# SO YOU WANT TO TRAIN AN LLM FROM SCRATCH

- What architecture should you use?
  - ✓ *Matters far less than model scale.*
- How big should your model be?
  - ✓ *Large models reach the same level of performance with fewer optimisation steps and data points.*
- How much data do you need?
  - ✓ *Surprisingly little; when model size increases 8x, dataset size should only increase 5x.*
- How long should you train for?
  - ✓ *Should stop well before convergence.*





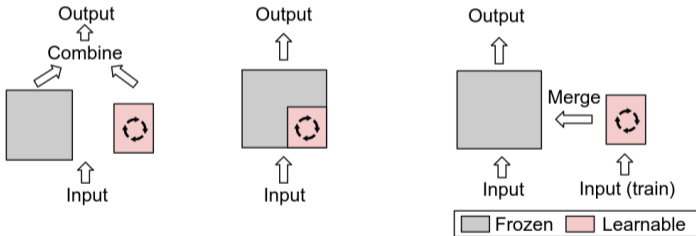
# FINE-TUNING - WHY?

- LLMs are huge (billions of params)
- Much of training is spent on learning basic semantic connections between tokens
- More efficient to take a (large) pre-trained model and **fine-tune** it for the specific downstream task
- Parameter Efficient Fine-Tuning (PEFT) = fine-tuning while minimizing the number of additional parameters or computational costs involved



# PEFT STRATEGIES

(a) Additive PEFT (b) Selective PEFT (c) Reparameterization PEFT

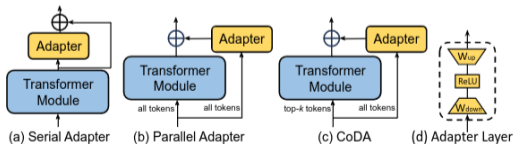


+ hybrid PEFT



# ADDITIVE PEFT

- **General idea:** freeze the entire pretrained model and **add** components to fine-tune
- **Adapters:** add small adapter layers to transformer blocks
- **Soft prompt:** fine-tune prompts with a learnable prefix
- **Prefix tuning:** add learnable prefixes to key and value vectors in all transformer layers





# SELECTIVE PEFT

- **General idea:** freeze most of the parameters; **select** a subset of parameters to fine-tune
- Often implemented as (learnable) binary “freeze/unfreeze” mask
- Imposing structure to the learnable selection may increase computational efficiency

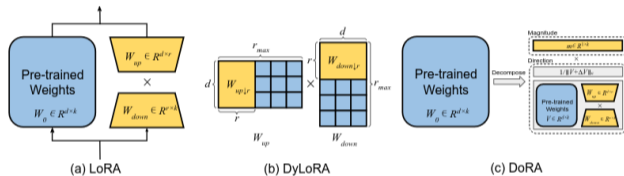


# REPARAMETERIZATION PEFT

- **Observation:** LLMs often have a low-dimensional intrinsic space
- **General idea:** train a small parallel network which projects the input to a low-dimensional space, effectively **reparameterizing** the entire pre-trained parameter space
- Most widely-recognized technique is LoRA



# LOW-RANK ADAPTATION (LoRA)



- **Basic implementation:** train a parallel network which bottlenecks the input to low-dimensional space (with **rank**  $r$ )
- During inference, as efficient as the original model
- More sophisticated variants choose the rank in a flexible fashion instead adhering to single value



# EFFICIENT PEFT

## Main considerations

- Processing latency
- Peak memory overhead

## Some approaches

- KV Caching
- Pruning
- Quantization