

LLM LE4 VT2025

Pre-training and Scaling

Fredrik Heintz

Dept. of Computer Science

Linköping University

fredrik.heintz@liu.se

@FredrikHeintz

Outline:

- Pre-training
- Parallel Training
- Scaling

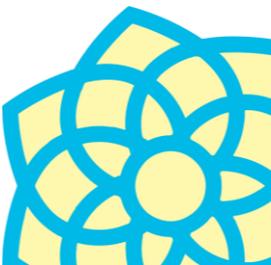
Language modelling

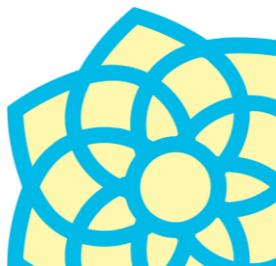
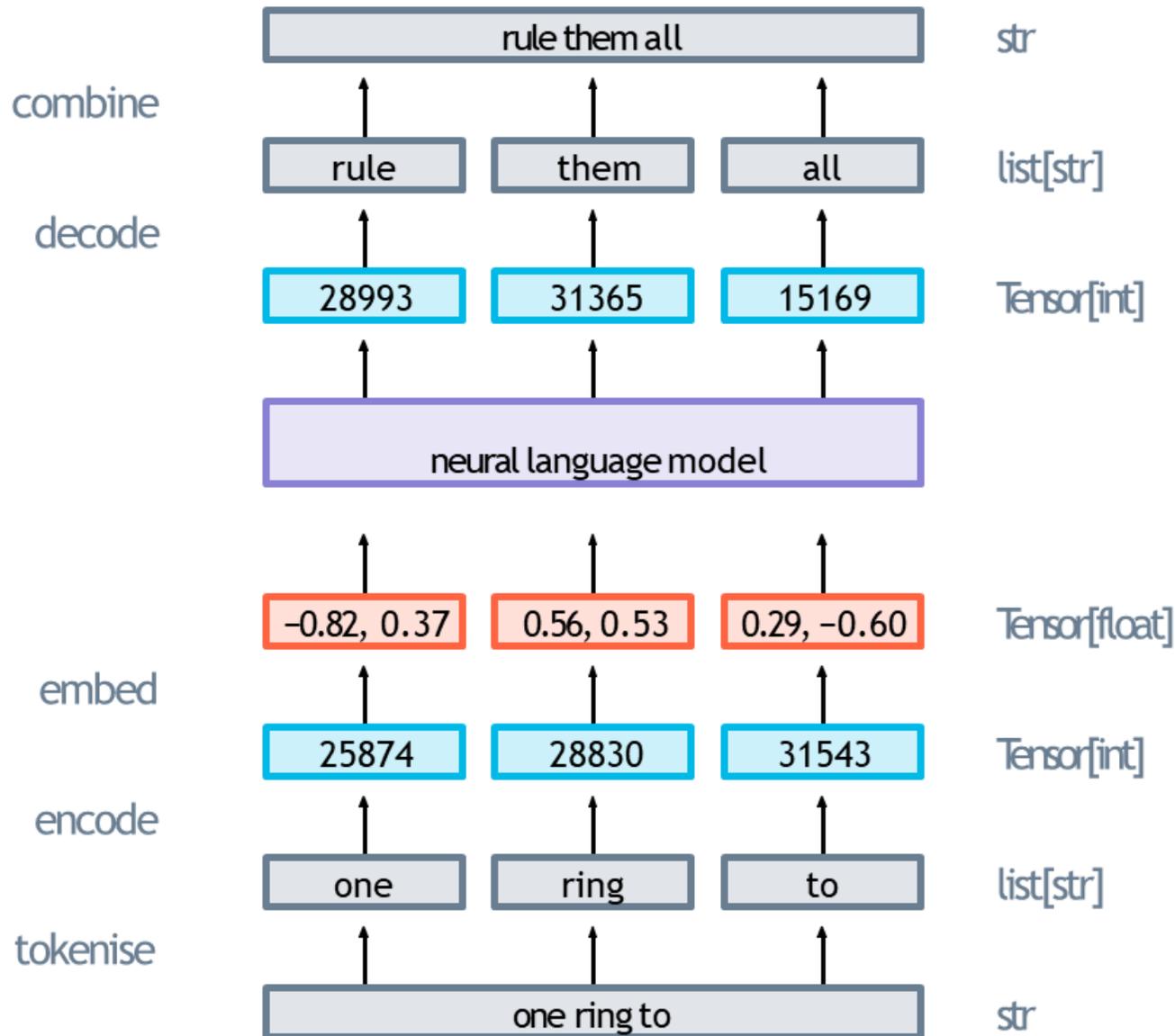
- **Language modelling** is the task of predicting which word comes next in a sequence of words.
- More formally, given a sequence of words w_1, \dots, w_t we want to know the probability of the next word, w_{t+1} :

$$P(w_{t+1} | w_1, \dots, w_t)$$

- We are assuming that w_{t+1} comes from a finite vocabulary V .

language models = classifiers

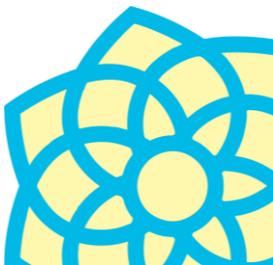
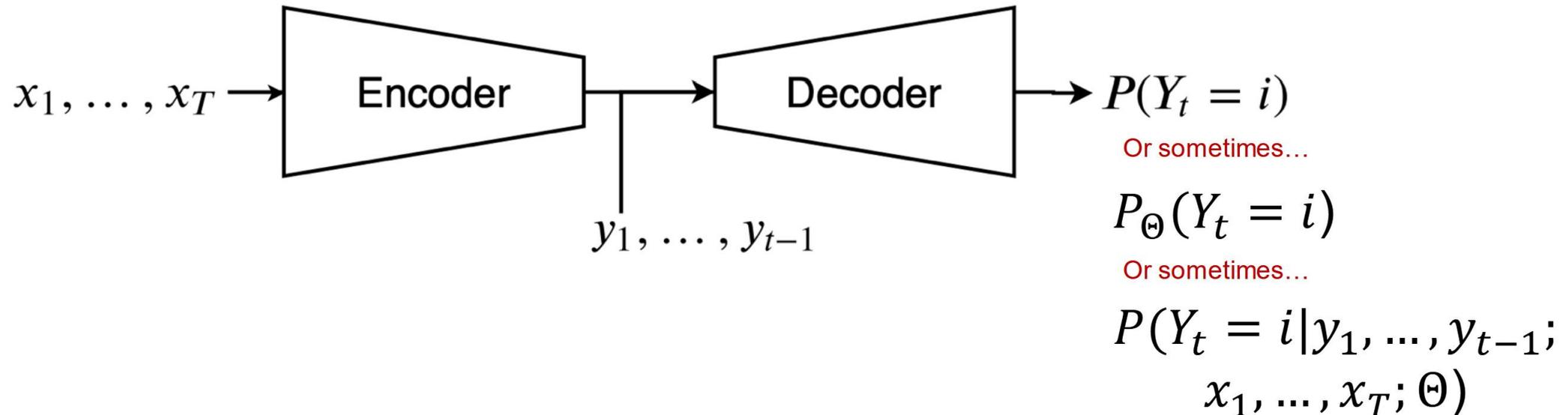




Neural Language Models

Input sequence: x_1, \dots, x_T

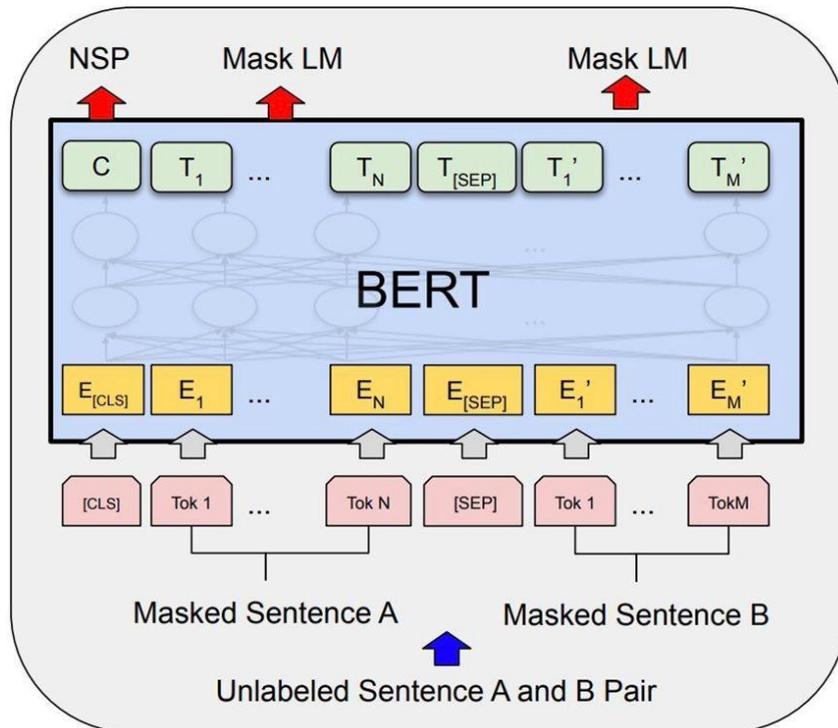
Target sequence: y_1, \dots, y_T



Language models: Broad Sense

- ❖ Decoder-only models (GPT-x models)
- ❖ Encoder-only models (BERT, RoBERTa, ELECTRA)
- ❖ Encoder-decoder models (T5, BART)

The latter two usually involve a different **pre-training** objective.



"translate English to German: That is good."

"cola sentence: The course is jumping well."

"stsb sentence1: The rhino grazed on the grass. sentence2: A rhino is grazing in a field."

"summarize: state authorities dispatched emergency crews tuesday to survey the damage after an onslaught of severe weather in mississippi..."

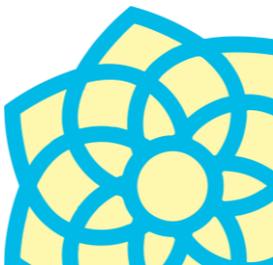
T5

"Das ist gut."

"not acceptable"

"3.8"

"six people hospitalized after a storm in attala county."



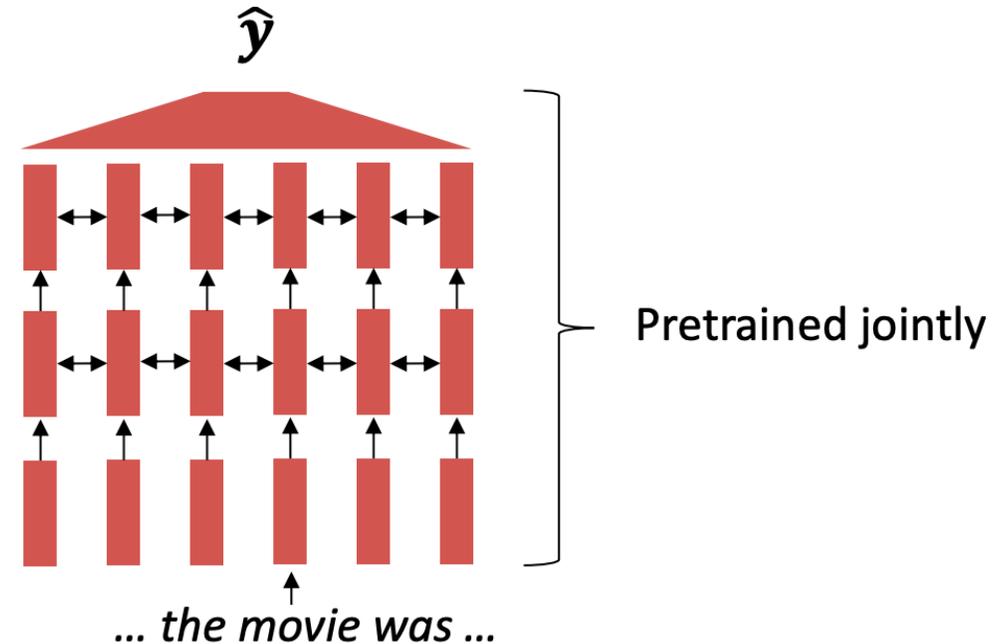
Why Pretraining?

In modern NLP:

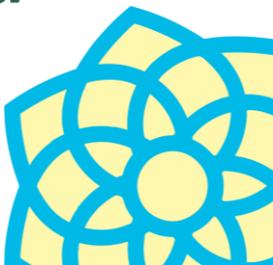
- All (or almost all) parameters in NLP networks are initialized via **pre-training**.
- Pretraining methods hide parts of the input from the model, and then train the model to reconstruct those parts.

This has been exceptionally effective at building strong:

- **representations of language**
- **parameter initializations** for strong NLP models.
- **probability distributions** over language that we can sample from

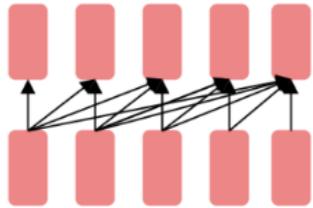


[This model has learned how to represent entire sentences through pretraining]



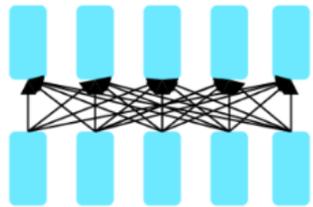
Pretraining for three types of architectures

The neural architecture influences the type of pretraining, and natural use cases.



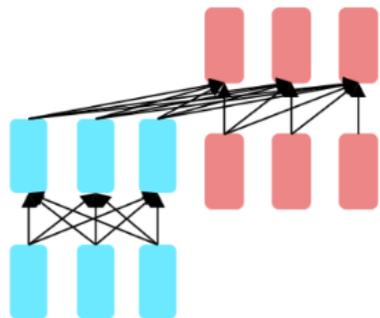
Decoders

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words
- **Examples:** GPT-2, GPT-3, LaMDA



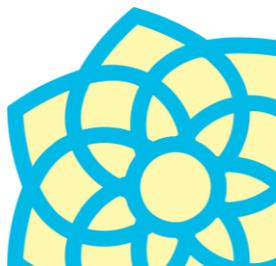
Encoders

- Gets bidirectional context – can condition on future!
- Wait, how do we pretrain them?
- **Examples:** BERT and its many variants, e.g. RoBERTa



**Encoder-
Decoders**

- Good parts of decoders and encoders?
- What's the best way to pretrain them?
- **Examples:** Transformer, T5, Meena

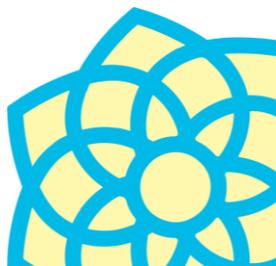


What do we want from our learning objective.

The goals of **pre-training** (the first stage of training) are to get the language model to:

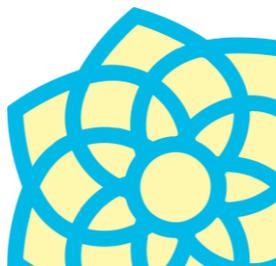
- learn the structure of natural language
- learn humans' understanding of the world (as encoded in the training data).

We want a learning objective that facilitates these goals.



Objectives for pre-training

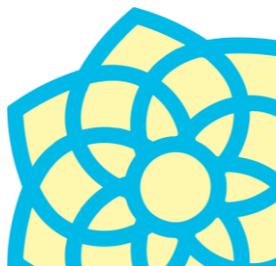
- Predict a suffix given a prefix.
 - Input: I took my dog, Fido, to the
 - Target: park for his walk.
- Masked language modeling
 - Input: I took <x> to <y> his walk.
 - Target: <x> my dog, Fido, <y> the park for
- Full language modeling: predict next token given previous tokens
 - Target: I took my dog, Fido, to the park for his walk.



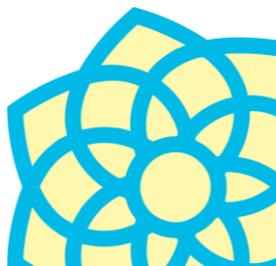
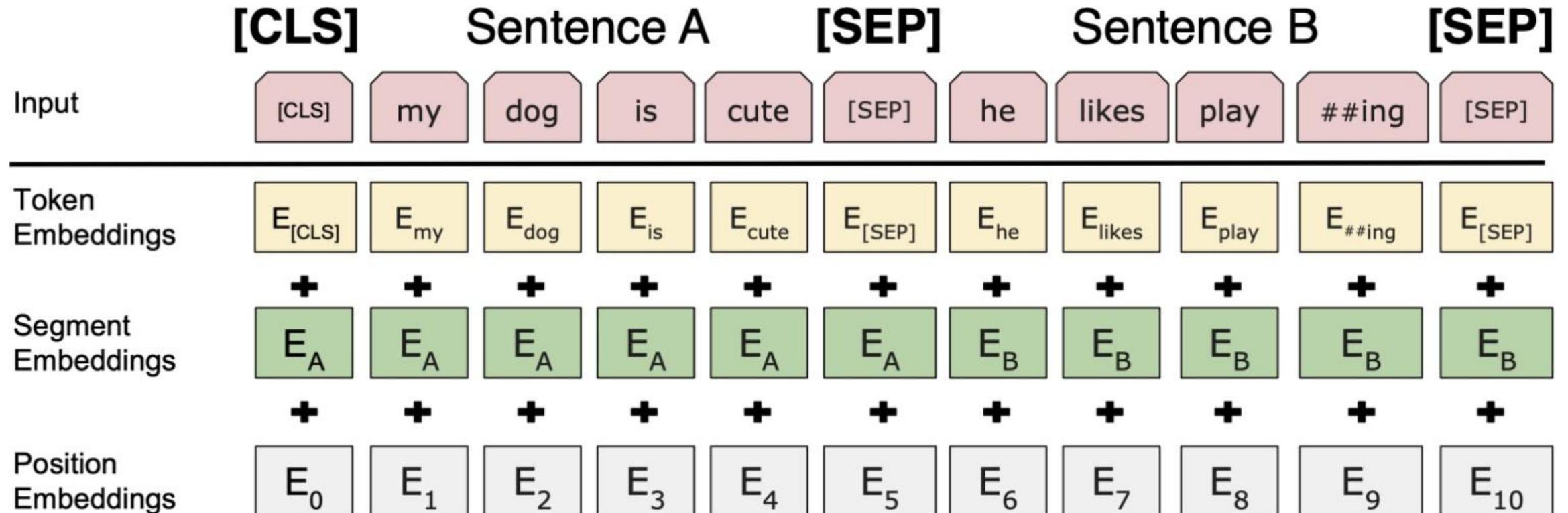
Pretraining Encoders

BERT [[Devlin et al, NAACL 2019](#)]

- **Fully bidirectional transformer encoder**
 - BERTbase: 12 layers, hidden size=768, 12 att'n heads (110M parameters)
 - BERTlarge: 24 layers, hidden size=1024, 16 att'n heads (340M parameters)
- **Input:** sum of token, positional, segment embeddings
 - **Segment embeddings** (A and B): is this token part of sentence A (before SEP) or sentence B (after SEP)?
- **[CLS]** and **[SEP]** tokens: added during pre-training
- **Pre-training tasks:**
 - Masked language modeling
 - Next sentence prediction



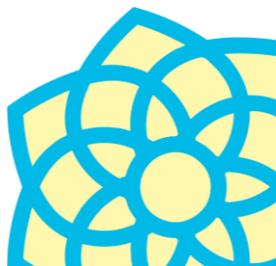
BERT Input



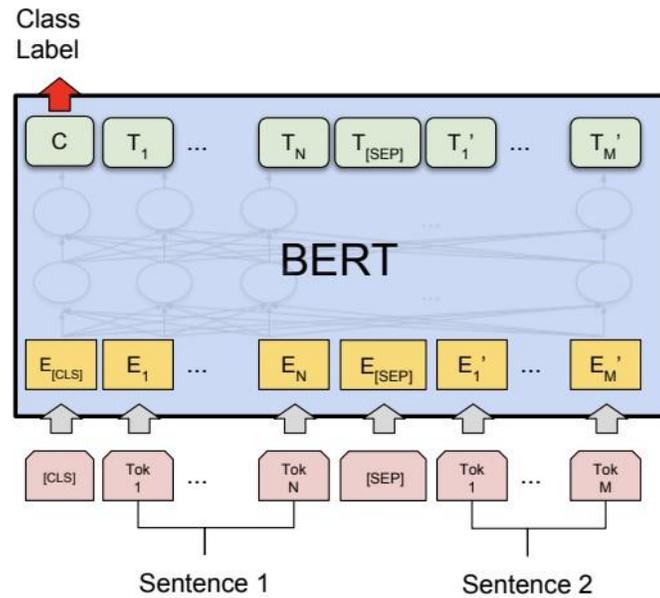
BERT Pre-training Tasks

BERT is jointly pre-trained on two tasks:

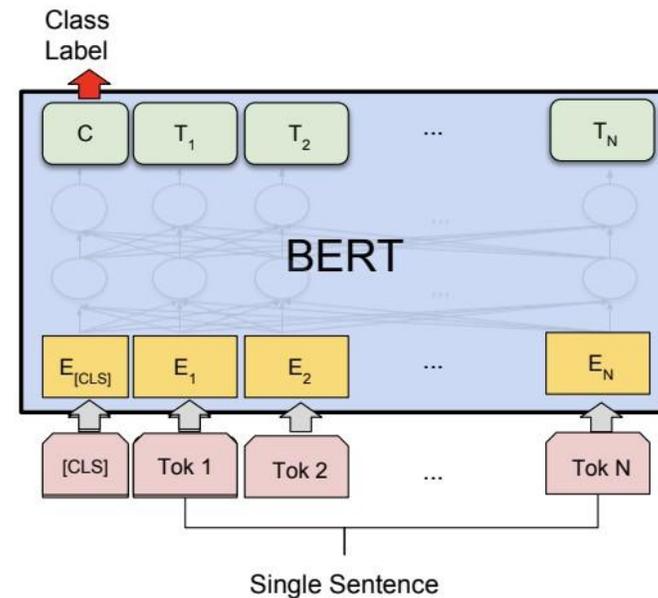
- **Next-sentence prediction:** [based on CLS token]
 - Does sentence B follow sentence A in a real document?
- **Mask language modeling:**
 - **15%** of tokens are randomly chosen as masking tokens
 - 10% of the time, a masking token remains unchanged
 - 10% of the time, a masking token is replaced by a random token
 - **80% of the time**, a masking token is replaced by [MASK], and the **output layer has to predict the original token**



Using BERT for Classification

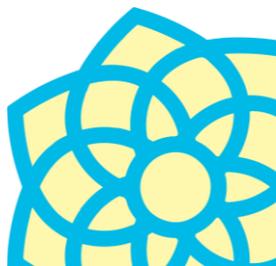


Sentence Pair
Classification

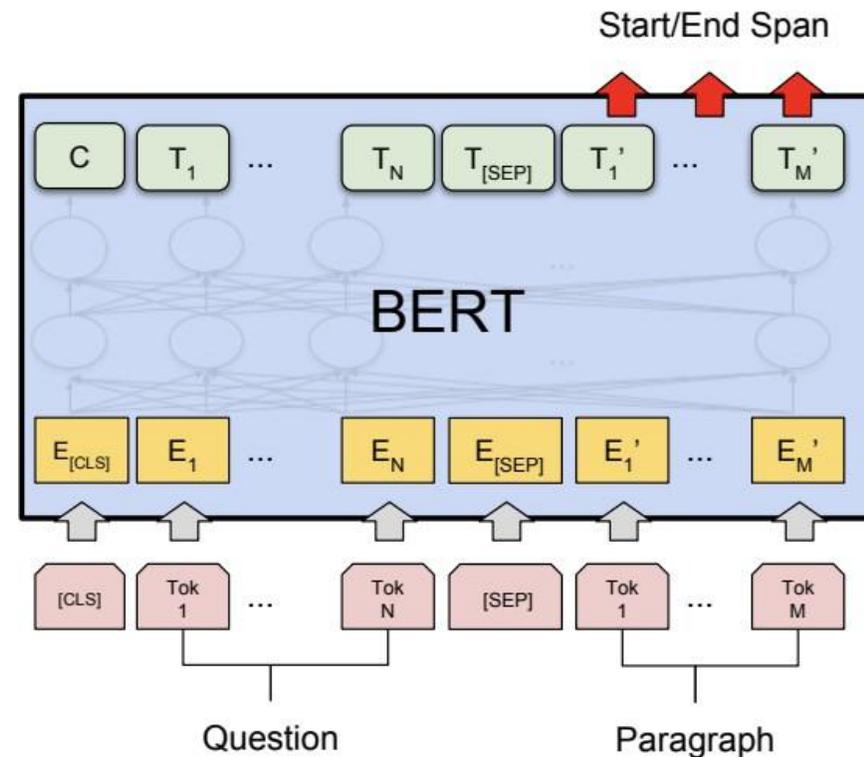


Single Sentence
Classification

Add a **softmax classifier** on final layer of [CLS] token

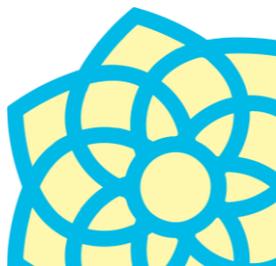


Using BERT for Question-Answering



Input: [CLS] question [SEP] answer passage [SEP]

Learn to predict a START and an END token on answer tokens



Examples of language models pretraining objectives



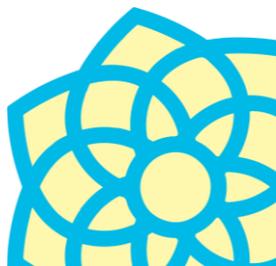
Guess the next word in the sentence (GPT)



Guess some masked words in the sentence (BERT)

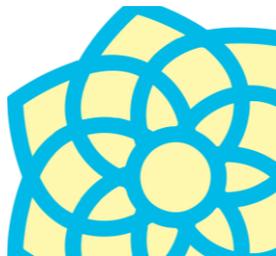
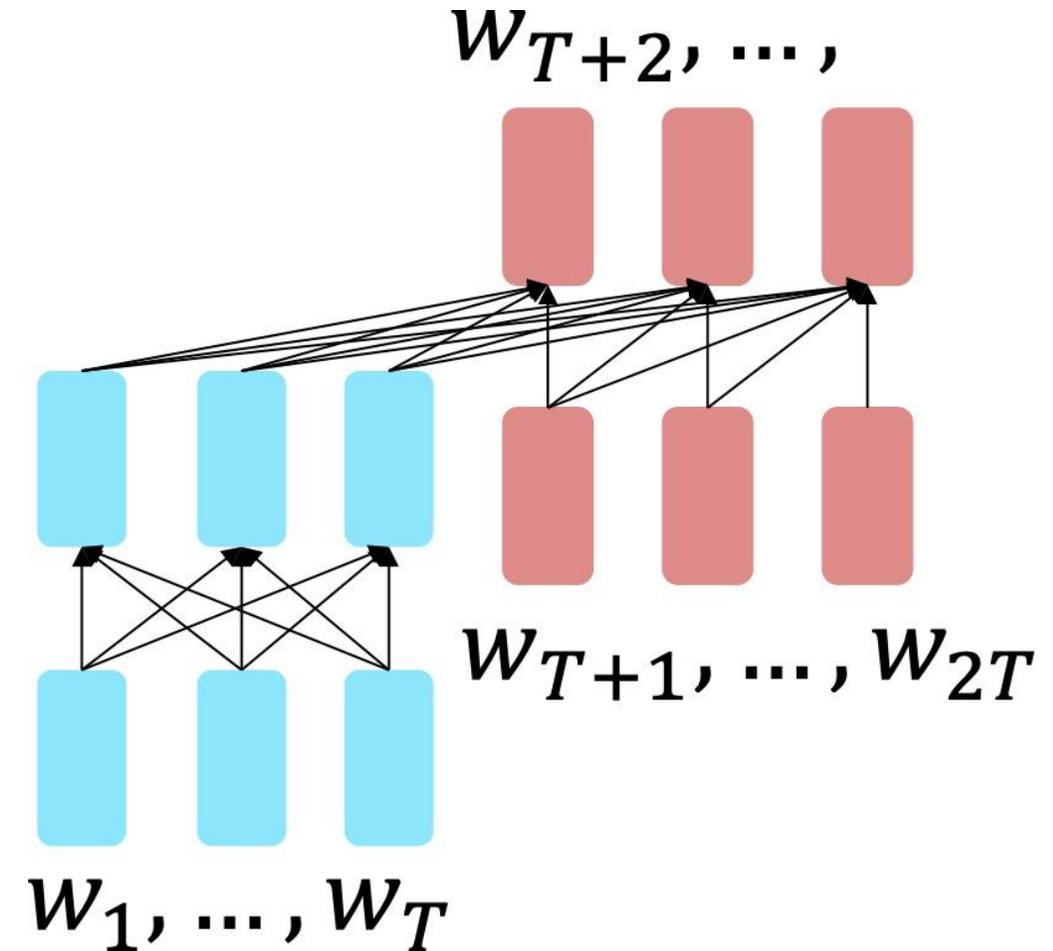
Why not encoder-based LLMs?

1. **I cannot generate (it discriminates):** It can only work for classification (discrimination) tasks, it is not easy to generate something new.
1. **Its objective is not scalable:** Its self-supervised tasks (masked language model) are just too simple for LLMs, and increasing model size does not improve performance too much.



Pretraining Encoder-Decoders

The **encoder** portion benefits from bidirectional context; the **decoder** portion is used to train the whole model through language modeling.



Pretraining Encoder-Decoders: Span Corruption

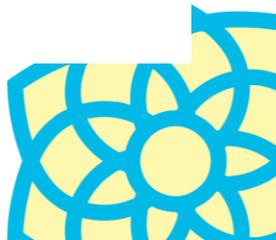
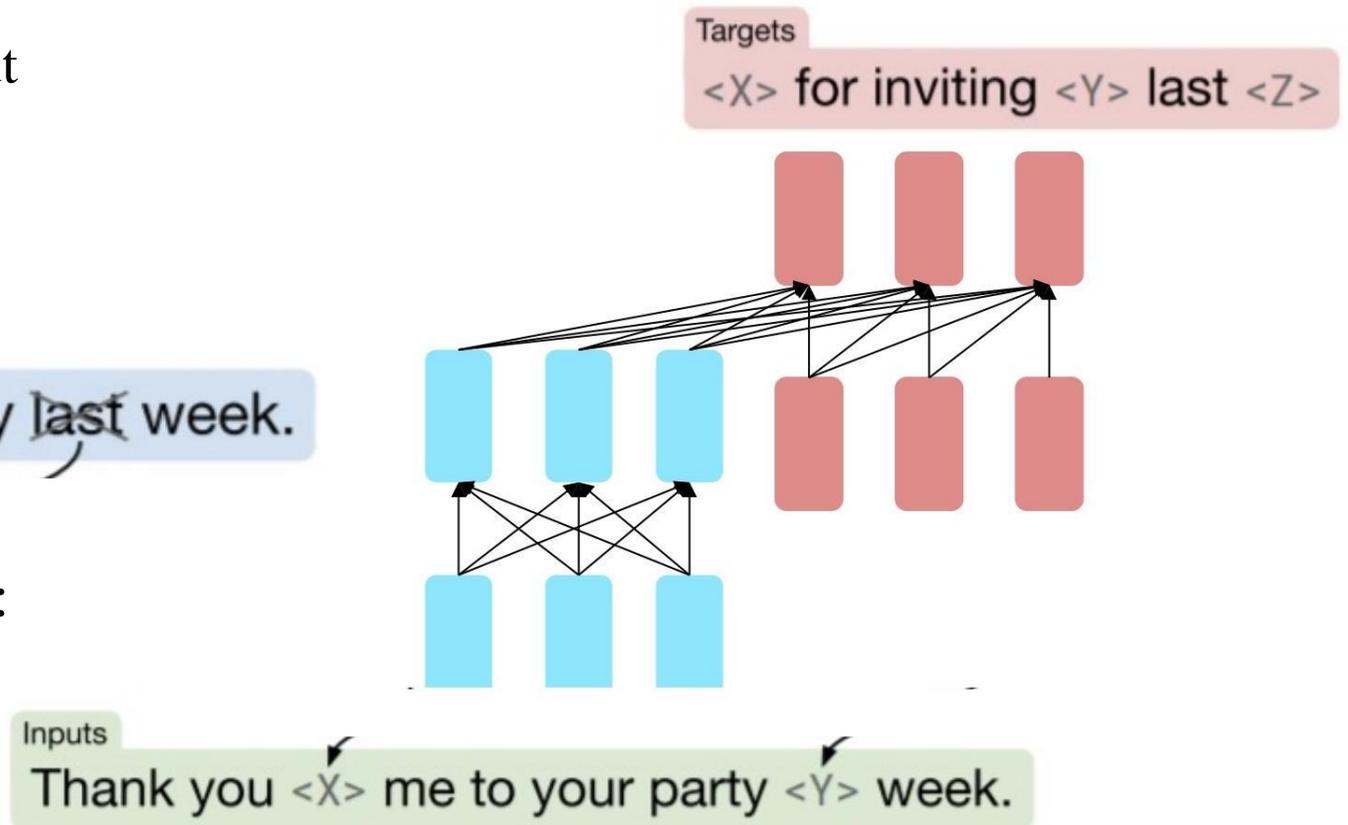
What [[Raffel et al., 2018](#)] found to work best was span corruption. Their model: T5.

Replace different-length spans from the input with unique placeholders; decode out the spans that were removed!

Original text

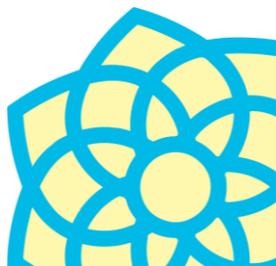
Thank you ~~for~~ ~~inviting~~ me to your party ~~last~~ week.

This is implemented in text preprocessing:
it's still an objective that looks like **language modeling** at the decoder side.



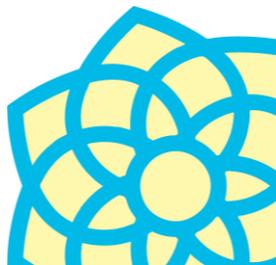
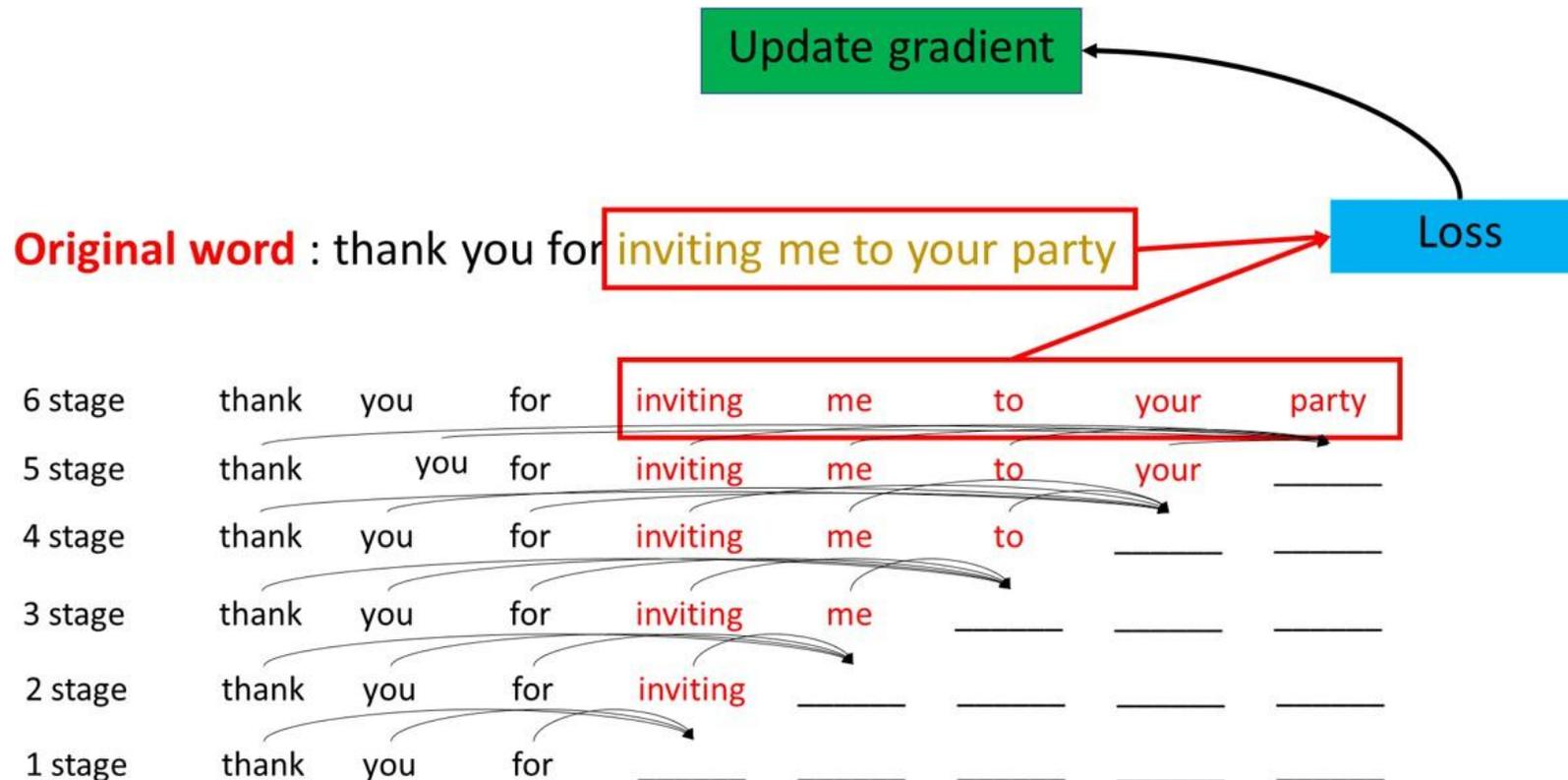
Why not Encoder-Decoder LLMs?

1. Decoder could work also as a seq-2-seq task, its protocol is much easier
2. When performing multi-turn generation, it is not easy to cache previous values.



Pretraining Decoders

It's natural to pretrain decoders as language models and then use them as generators, finetuning their $p_{\theta}(w_t|w_{1:t-1})!$



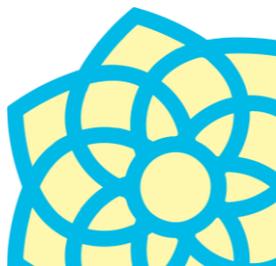
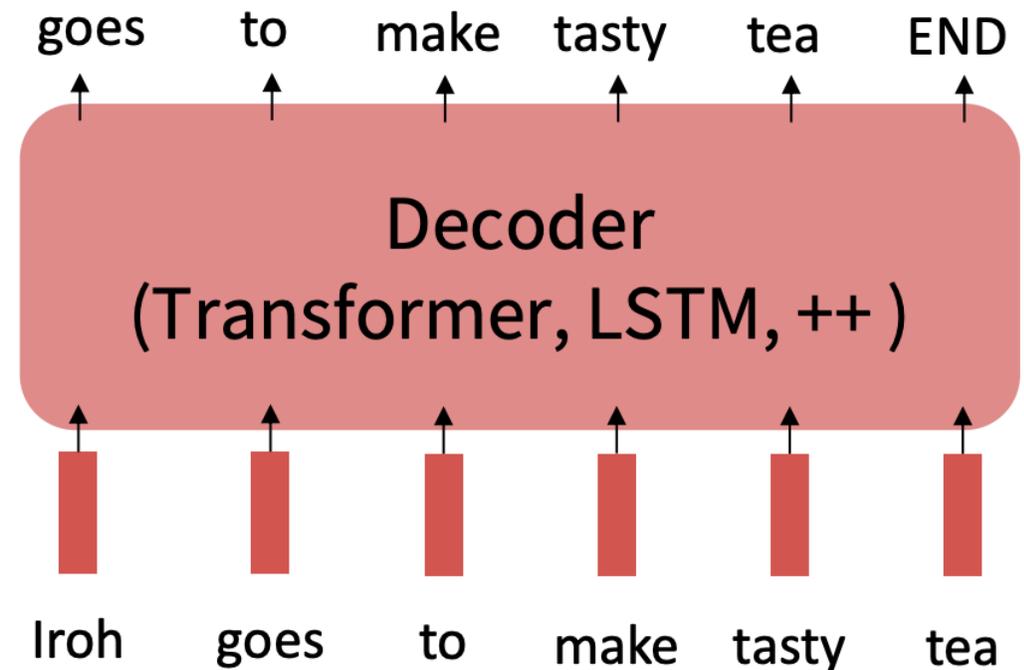
Pretraining through language modeling

Recall the language modeling task:

- Model the probability distribution over words given their past contexts.
- There's lots of data for this! (In English.)

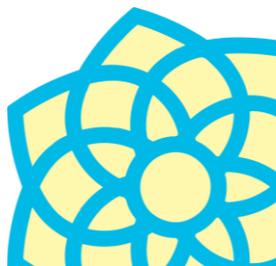
Pretraining through language modeling:

- Train a neural network to perform language modeling on a large amount of text.
- Save the network parameters.



Common roadmap for LLMs

- Phase 1: pre-training
 - Learn **general** world knowledge, ability, etc.
- Phase 2: Supervised finetuning
 - Tailor to **tasks** (**unlock** some abilities)
- Phase 3: RLHF
 - Tailor to **humans**
 - *Even you could teach ChatGPT to do something*



Parallel Training

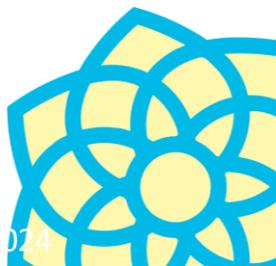
Parallel Training: Overview

As scale grows, training with one GPU is not enough

- There are many ways to improve efficiency on single-GPU training
 - Checkpointing: moving part of the operations to CPU memory
 - Quantizing different part of the optimization to reduce GPU memory cost
- Eventually more FLOPs are needed

Different setups of parallel training:

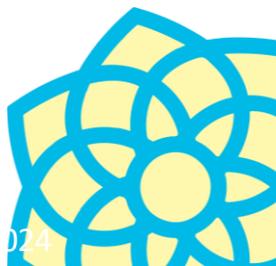
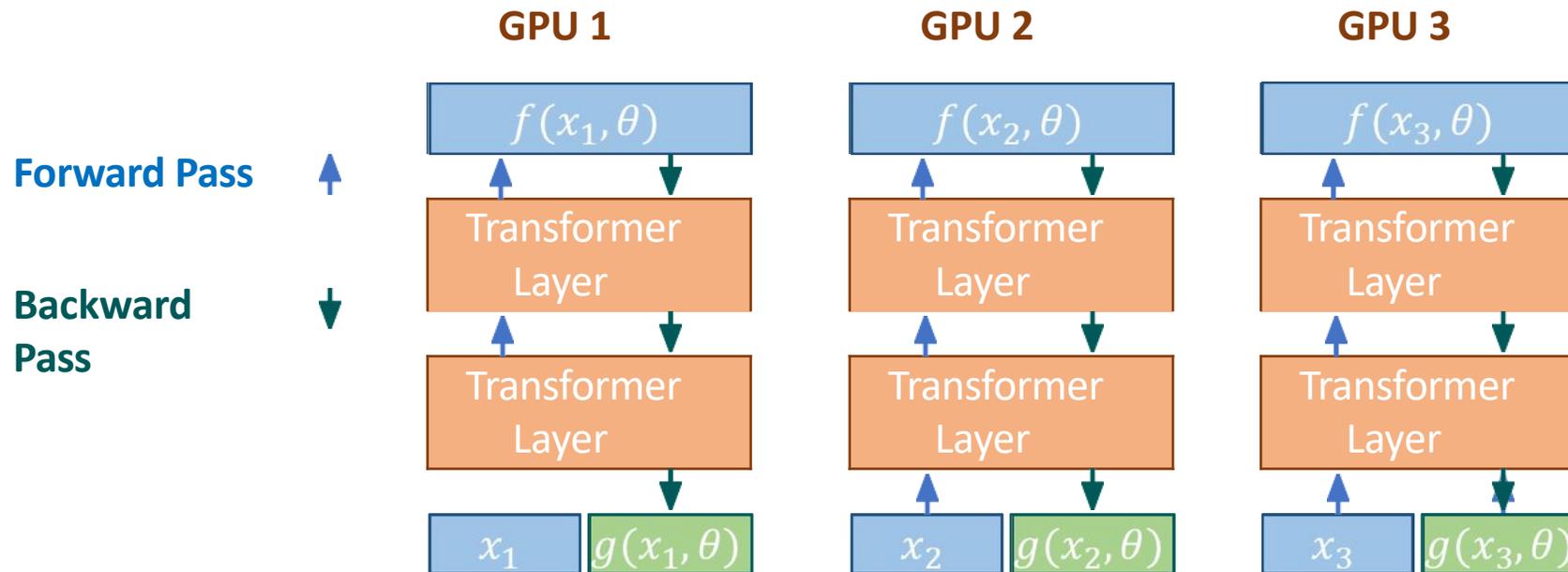
- When model training can fit into single-GPU
 - Data parallelism
- When model training cannot fit into single-GPU
 - Model parallelism: pipeline or tensor



Parallel Training: Data Parallelism

Split training data batch into different GPUs

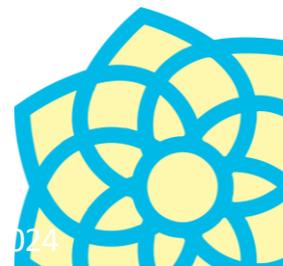
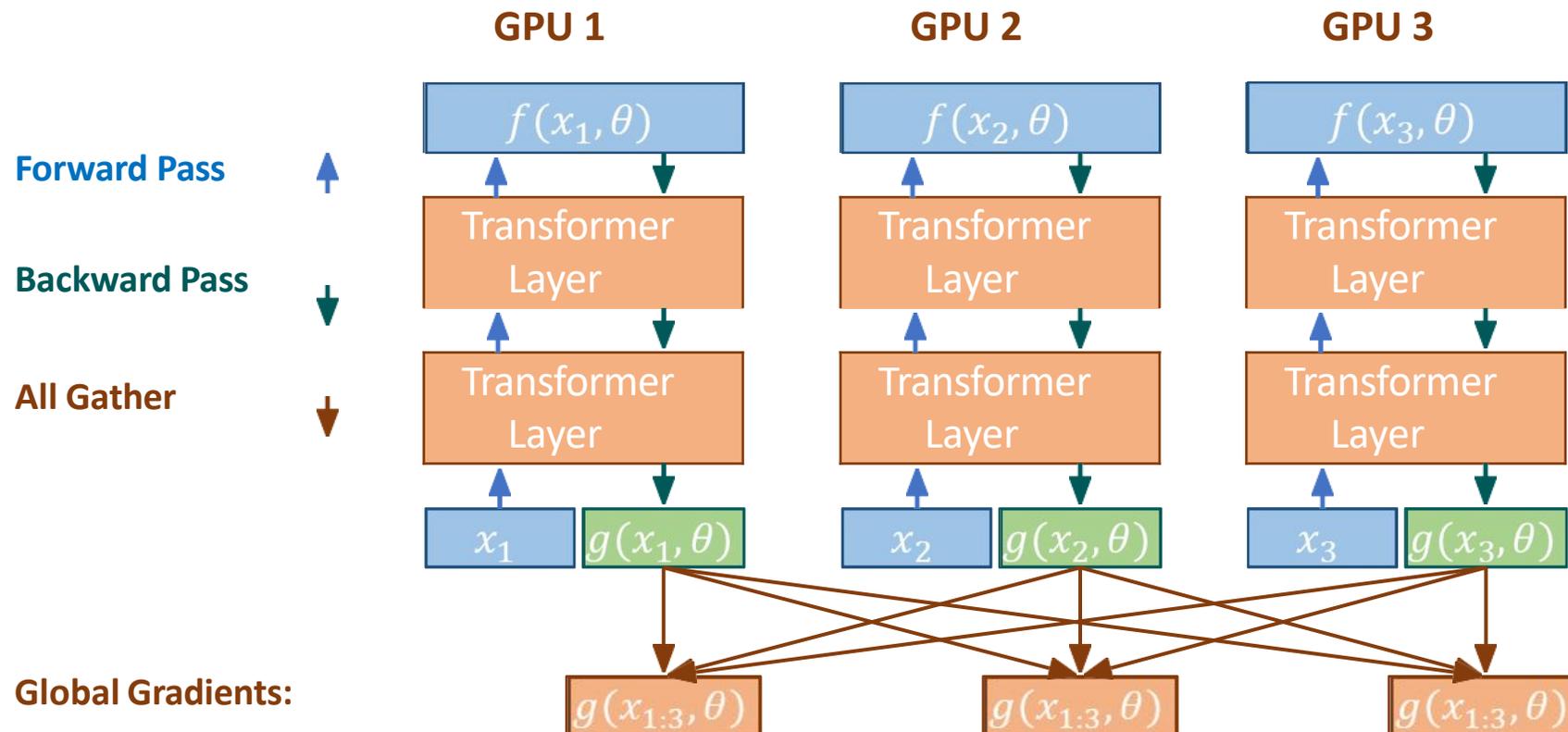
- Each GPU maintains its own copy of model and optimizer
- Each GPU gets a different local data batch, calculates its gradients



Parallel Training: Data Parallelism

Split training data batch into different GPUs

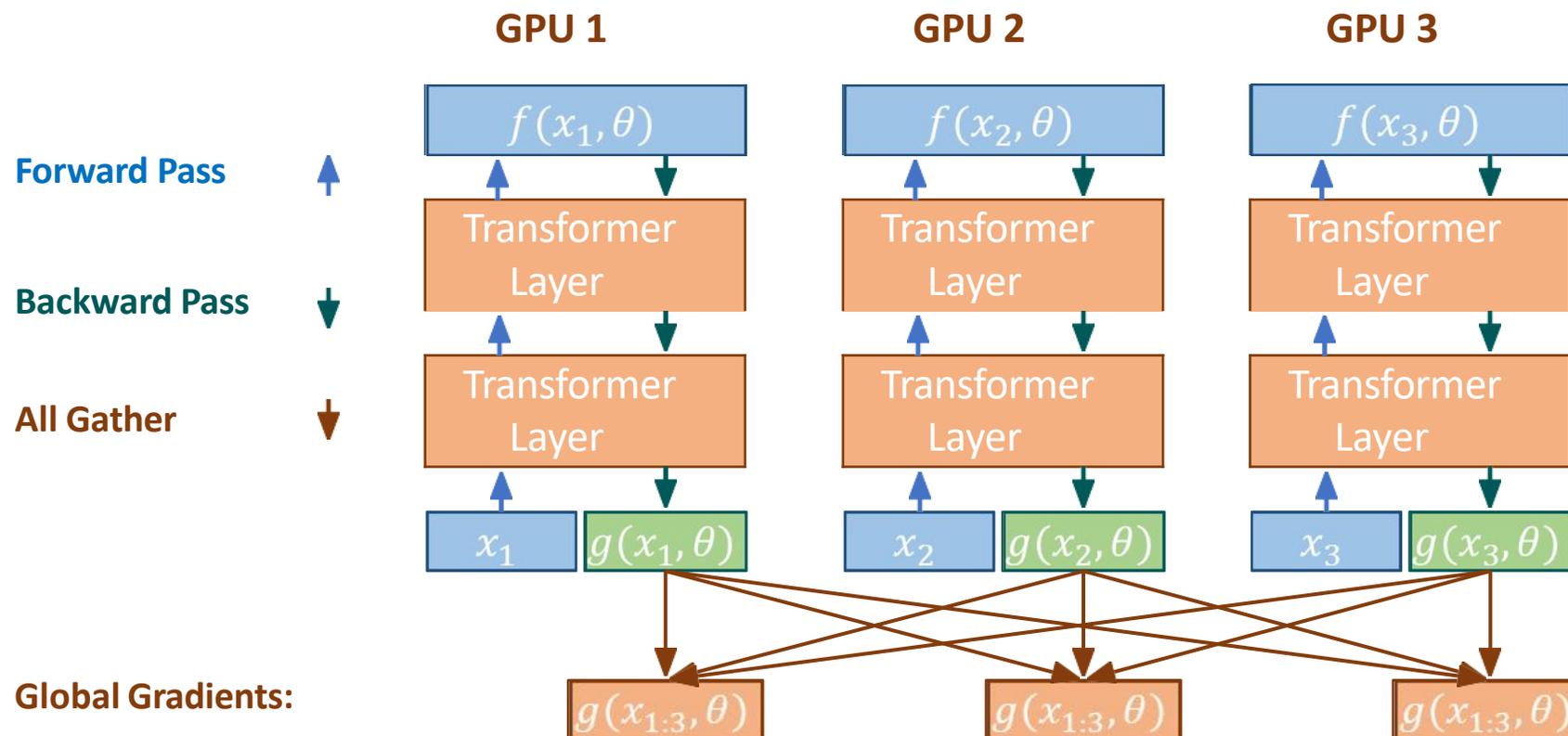
- Each GPU maintains its own copy of model and optimizer
- Each GPU gets a different local data batch, calculates its gradients
- Gather local gradients together to each GPU for global updates



Parallel Training: Data Parallelism

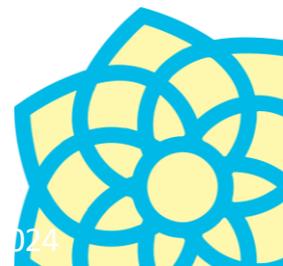
Split training data batch into different GPUs

- Each GPU maintains its own copy of model and optimizer
- Each GPU gets a different local data batch, calculates its gradients
- Gather local gradients together to each GPU for global updates



Communication:

- The full gradient tensor between every pair of GPUs, at each training batch.
- Not an issue between GPUs in the same machine or machines with infinity band
- Will need work around without fast cross-GPU connection



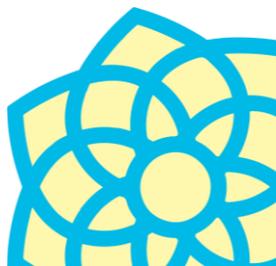
Parallel Training: Model Parallelism

LLM size grew quickly and passed the limit of single GPU memory

	Cost of 10B Model
Parameter Bytes	20GB
Gradient Bytes	20GB
Optimizer State: 1st Order Momentum	20GB
Optimizer State: 2nd Order Momentum	20GB
Total Per Model Instance	80GB

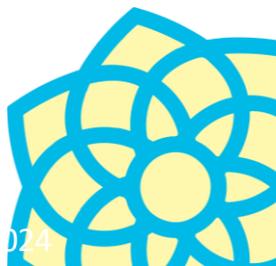
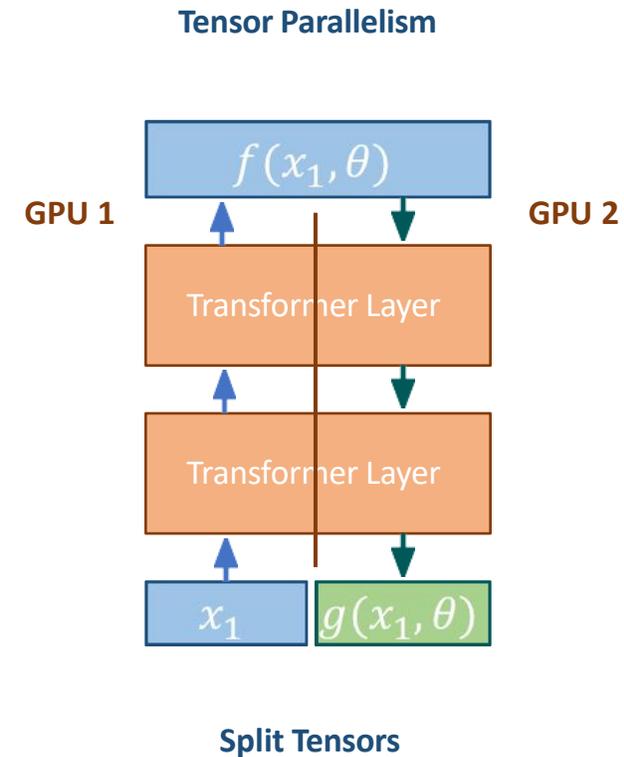
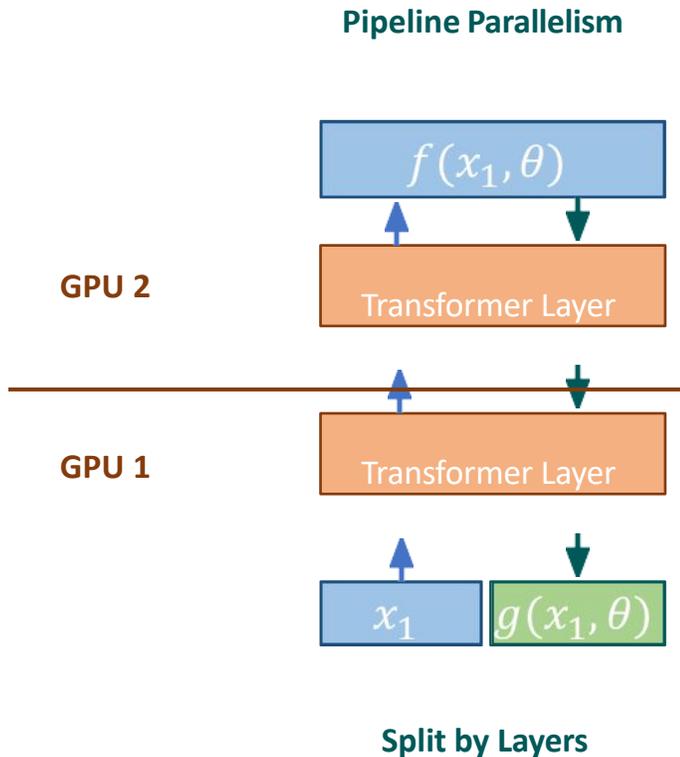
Memory Consumption of Training Solely with **BF16** (Ideal case) of a model sized Ψ

Solution: Split network parameters (thus their gradients and corresponding optimizer states) to different GPUs



Parallel Training: Model Parallelism

Two ways of splitting network parameters



Parallel Training: Pipeline Parallelism

Split network by layers, aligning devices by layer order to a pipeline, and pass data through devices [7]

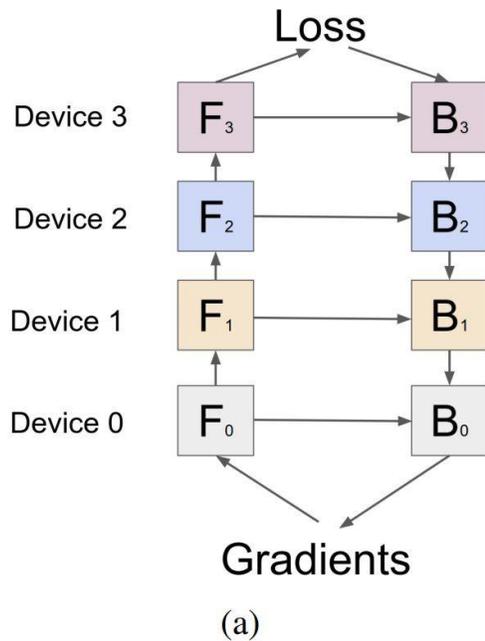
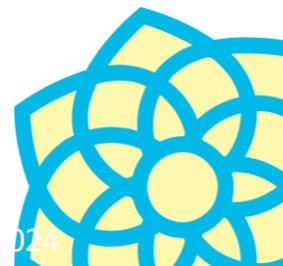
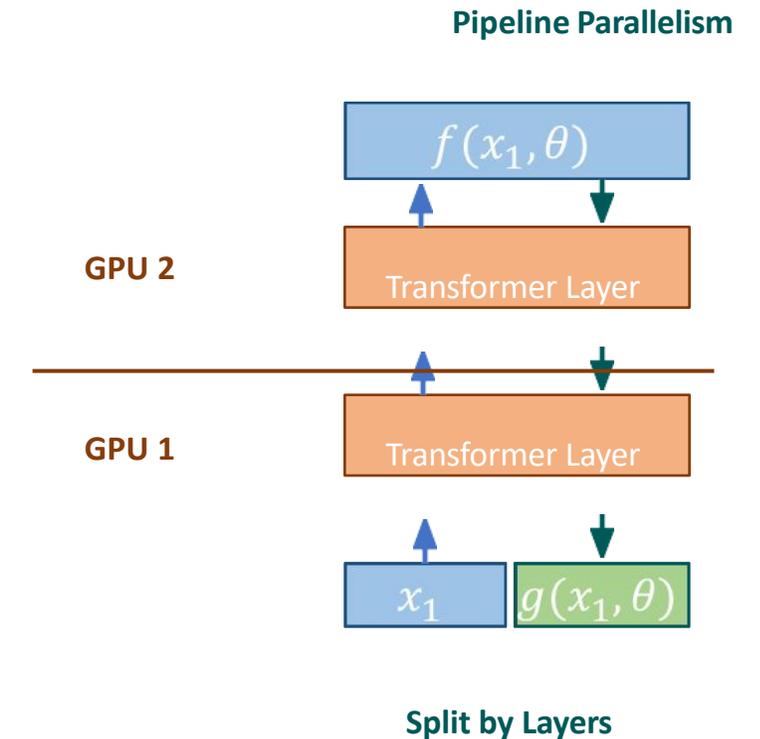


Illustration of Pipeline Parallelism [7]



Parallel Training: Pipeline Parallelism

Split network by layers, aligning devices by layer order to a pipeline, and pass data through devices [7]

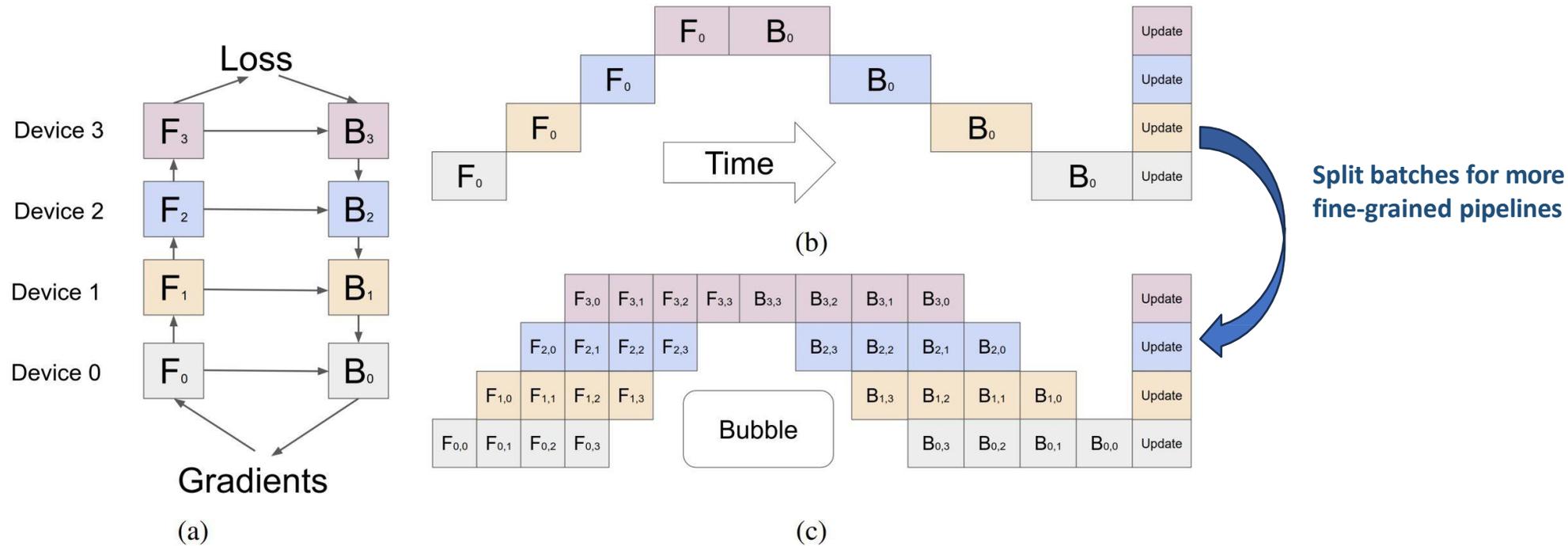
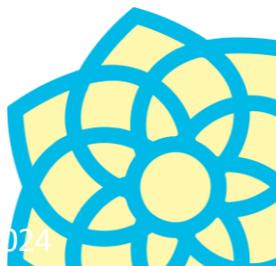
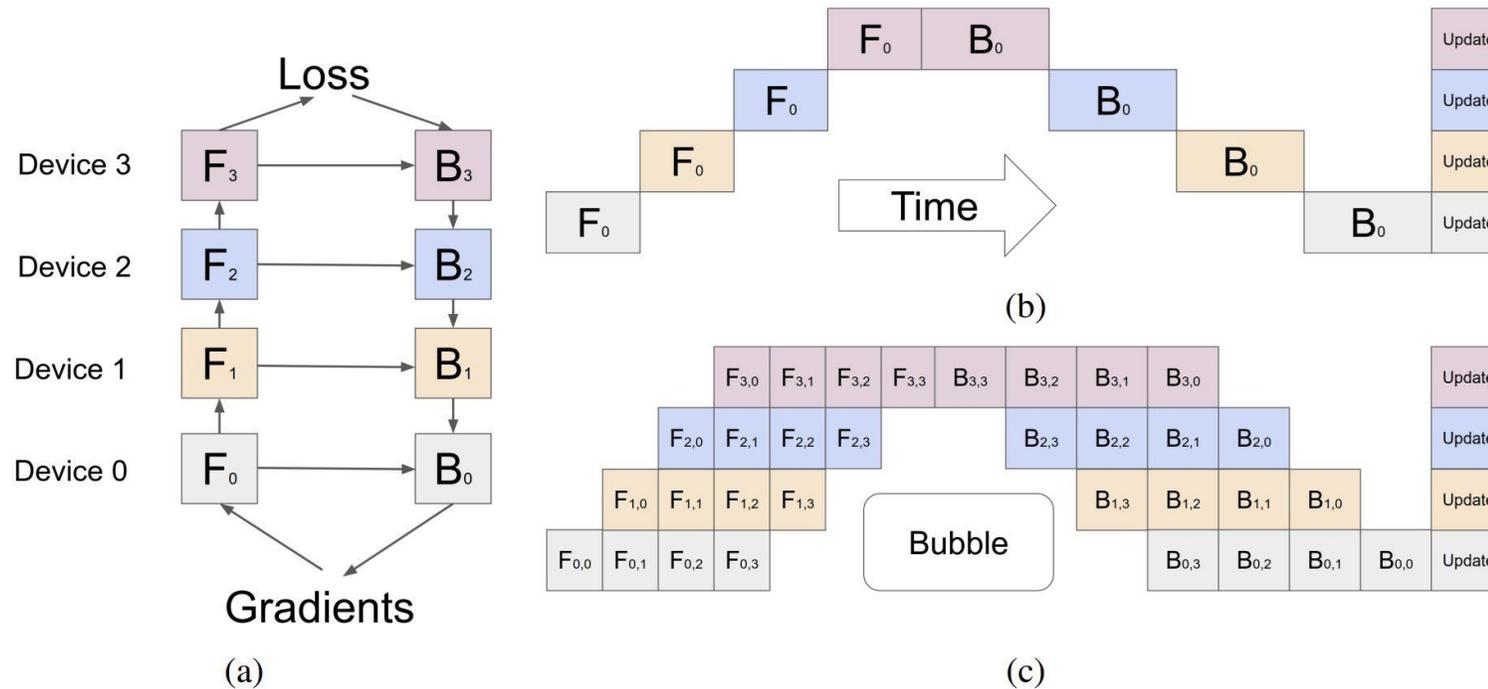


Illustration of Pipeline Parallelism [7]



Parallel Training: Pipeline Parallelism

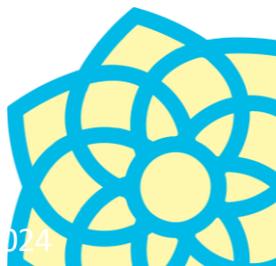
Split network by layers, aligning devices by layer order to a pipeline, and pass data through devices [7]



Communication:

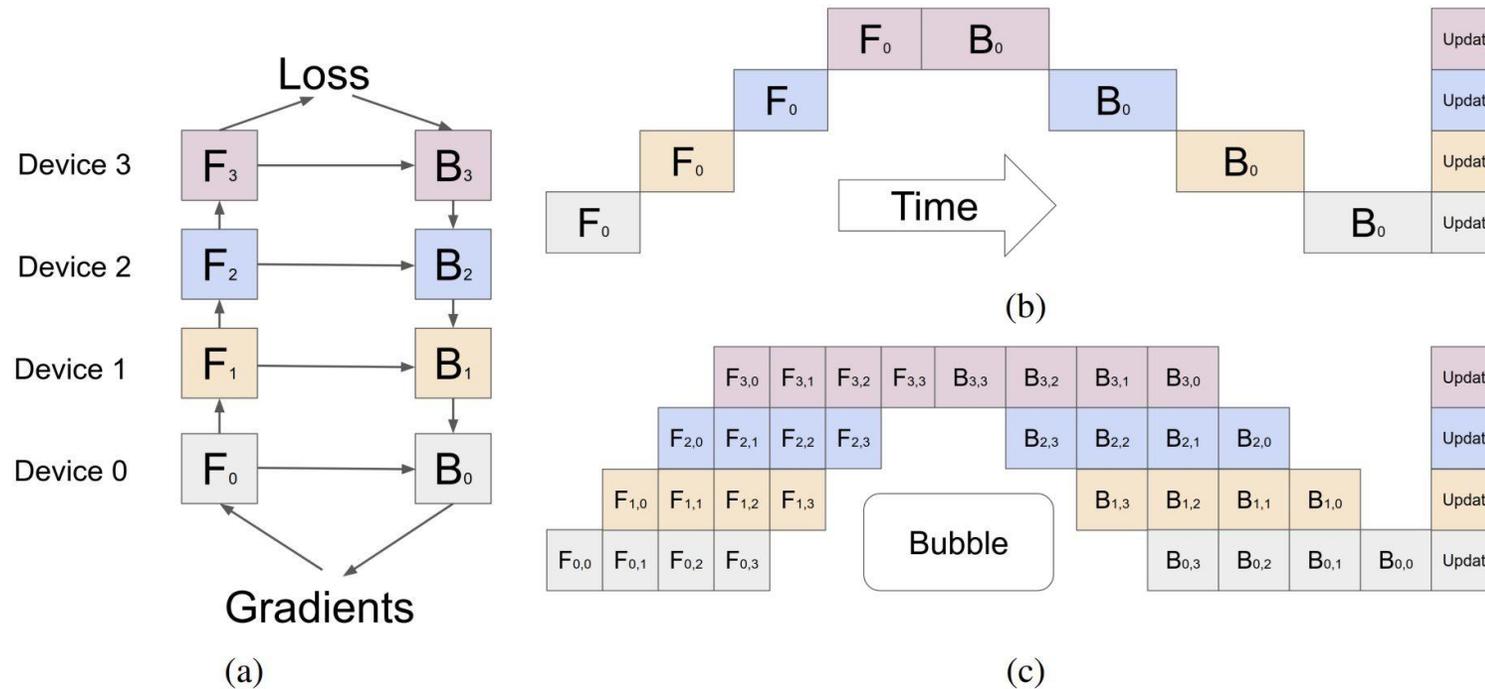
- Activations between nearby devices in forward pass
- Partial gradients between nearby devices in backward

Illustration of Pipeline Parallelism [7]



Parallel Training: Pipeline Parallelism

Split network by layers, aligning devices by layer order to a pipeline, and pass data through devices [7]



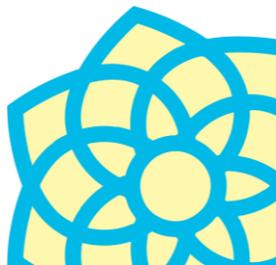
Communication:

- Activations between nearby devices in forward pass
- Partial gradients between nearby devices in backward

Illustration of Pipeline Parallelism [7]

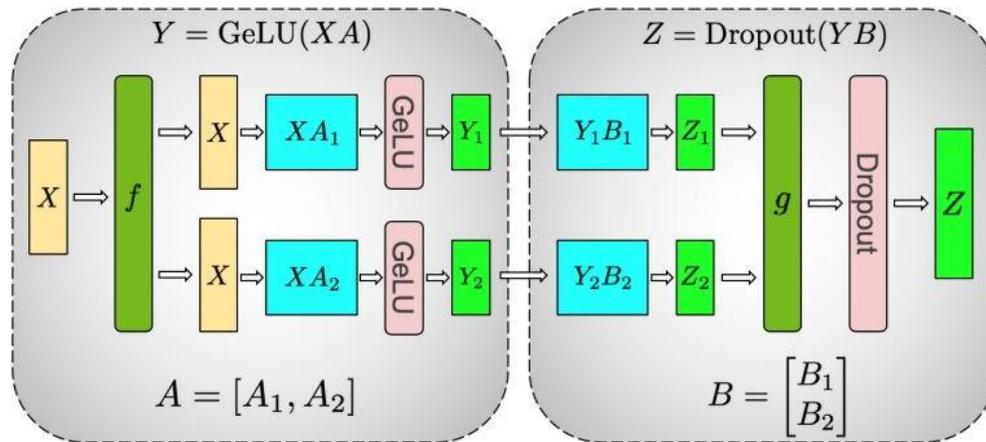
Pros: Conceptually simple and not coupled with network architectures. All networks have multiple layers.

Cons: Waste of compute in the Bubble. Bubble gets bigger with more devices and bigger batches.

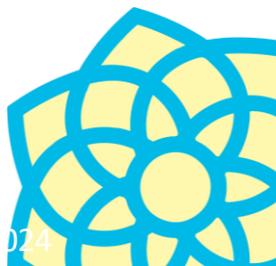


Parallel Training: Tensor Parallelism

Split the parameter tensors of network layers into different devices for parallel matrix operations

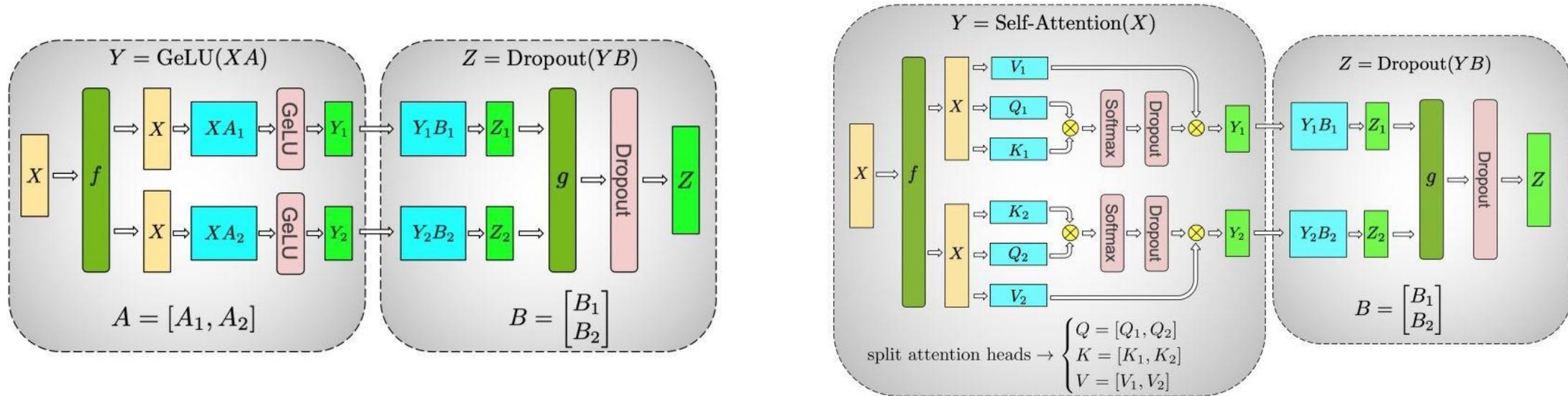


Tensor Parallelism of MLP blocks and Self-attention Blocks [8]

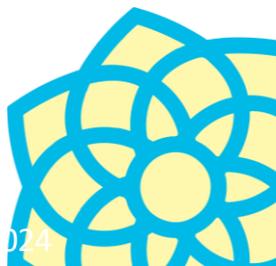


Parallel Training: Tensor Parallelism

Split the parameter tensors of network layers into different devices for parallel matrix operations

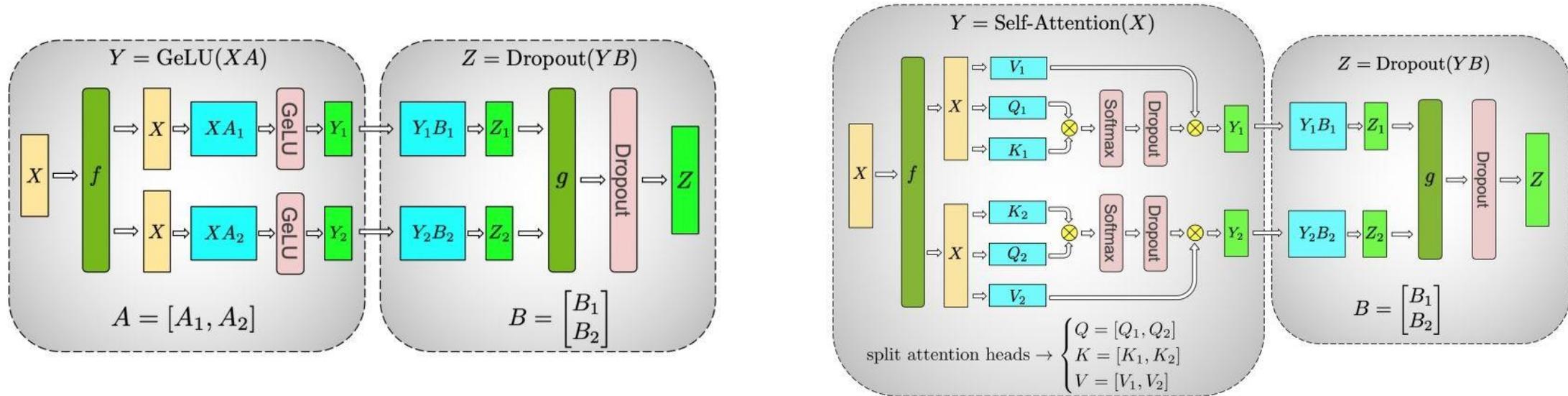


Tensor Parallelism of MLP blocks and Self-attention Blocks [8]



Parallel Training: Tensor Parallelism

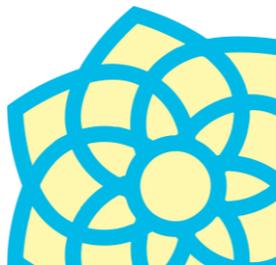
Split the parameter tensors of network layers into different devices for parallel matrix operations



Tensor Parallelism of MLP blocks and Self-attention Blocks [8]

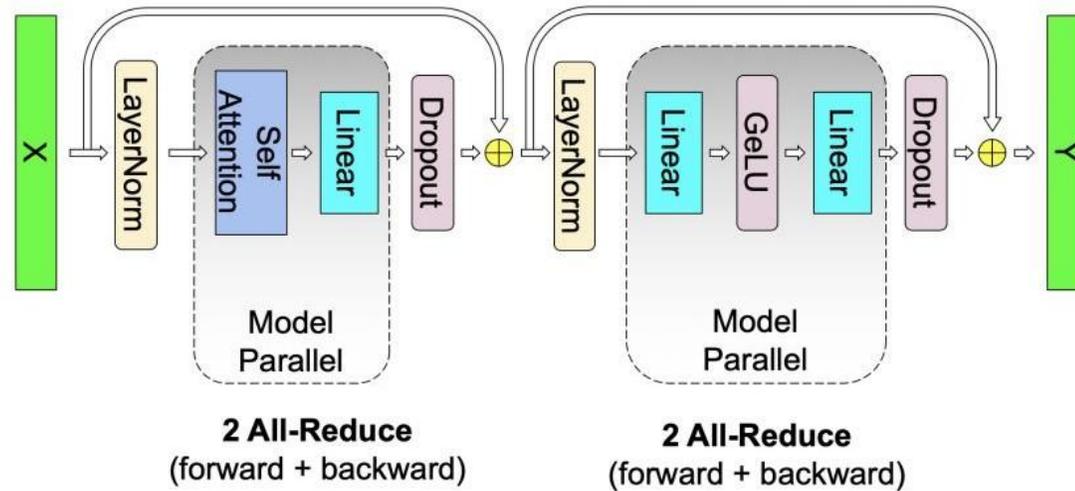
Pros: No bubble

Cons: Different blocks are better split differently, lots of customizations



Parallel Training: Tensor Parallelism

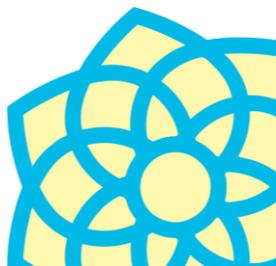
Split the parameter tensors of network layers into different devices for parallel matrix operations



Communication of Tensor Parallelism for a Transformer Layer [8]

Communication:

- All-gather of partial activations and gradients for each split tensor

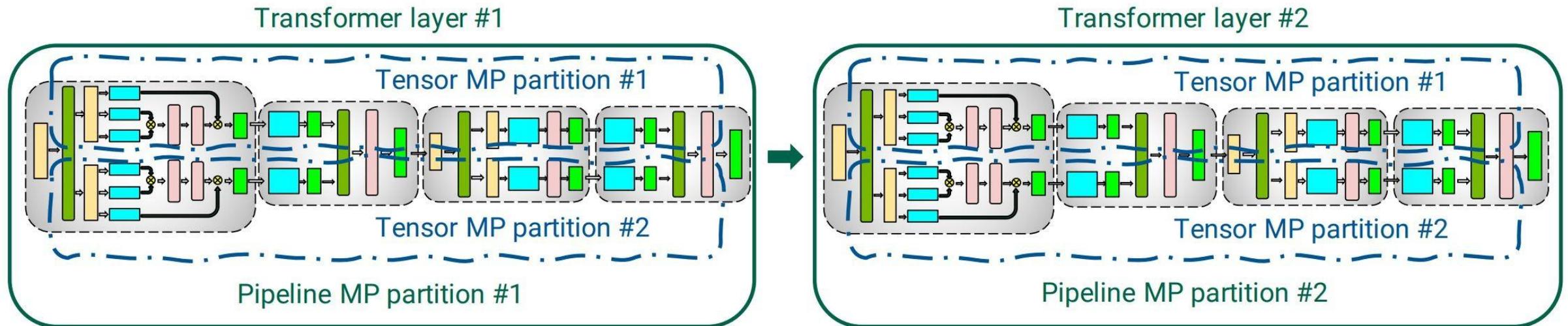


Parallel Training: Combining Different Approaches

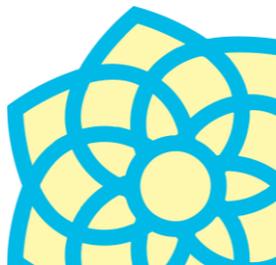
Often data parallelism and model parallelism are used together.

- No need not to use data parallelism

Pipeline Parallelism and Tensor Parallelism can also be used together.



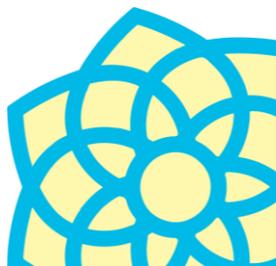
Combination of Tensor Parallelism and Pipeline Parallelism [9]



Scaling

Scaling laws in language modelling

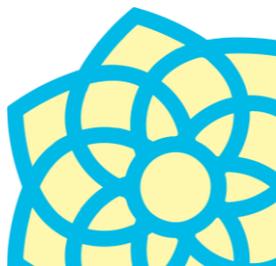
- **Scaling laws** describe how model performance improves as we increase key factors such as model size and training data size.
- Empirical results suggest that performance improvements obey a power law: performance increases, but at a diminishing rate.
 - cf. Heap's law
- Scaling laws can help developers answer many practically relevant questions about resource allocation.



Scaling Factors

Many factors in configuring a scaled up pretraining run for Transformer Decoder + Autoregressive LM

- Model size (parameter counts)
- Pretraining dataset size
- Pretraining compute (FLOPs or TPU/GPU hours)
- Network shape (Parameters allocations)
- Effective batch size
- Learning rate & learning rate scheduler
- Context length



Scaling Factors

Many factors in configuring a scaled up pretraining run for Transformer Decoder + Autoregressive LM

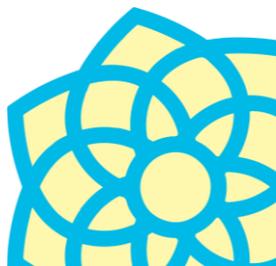
- Model size (parameter counts)
- Pretraining dataset size
- Pretraining compute (FLOPs or TPU/GPU hours)
- Network shape (Parameters allocations)
- Effective batch size
- Learning rate & learning rate scheduler
- Context length

Main factors to study

Hyper-parameters with rule of thumb

1. Batch size determined by GPU memory
2. Try biggest LR before blowing up

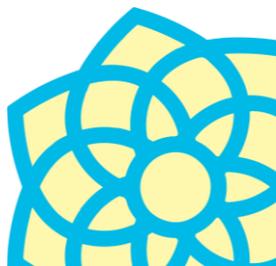
Balance complexity and downstream needs



Scaling Law Study: Setup

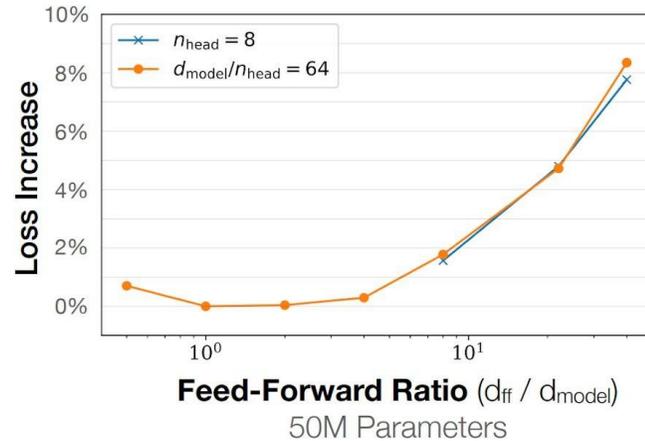
Empirically study the relationship between various factors to language model performances [4]

- Model: GPT-style, auto-regressive loss, maximum 1.5 billion non-embedding parameters
- Pretraining data: WebText2, harvest from Reddit out links, at max 23 billion tokens
- Metric: language modeling loss on testing data



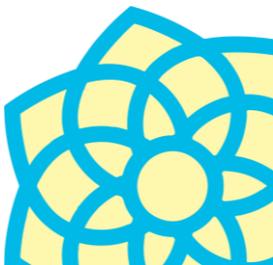
Scaling Law Study: Observations

Network shape (allocation of parameters at different parts) does not matter as much



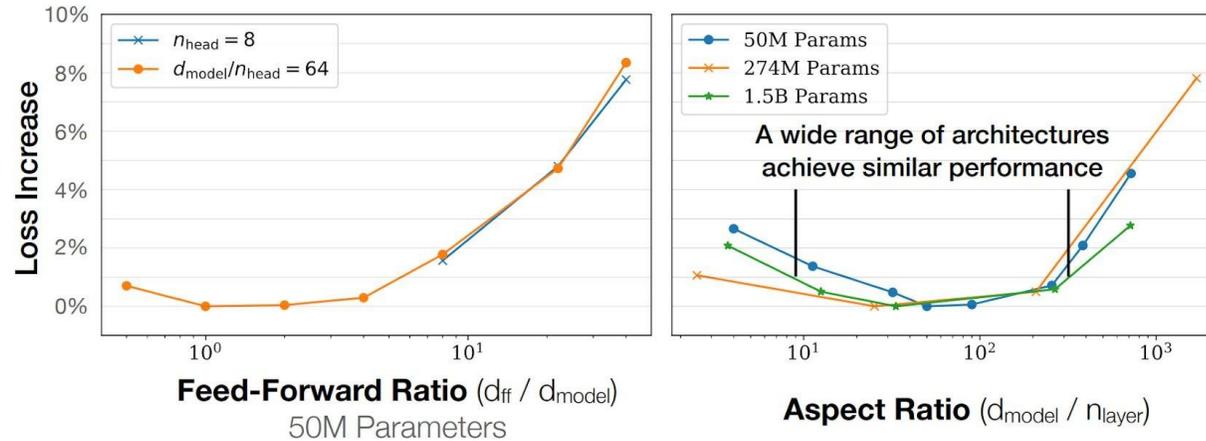
Language model loss changes with different network shape configurations [4].

- As long as the network shape is in a general sweet range, it does not impact performance much



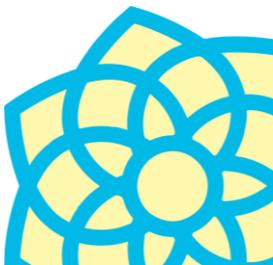
Scaling Law Study: Observations

Network shape (allocation of parameters at different parts) does not matter as much



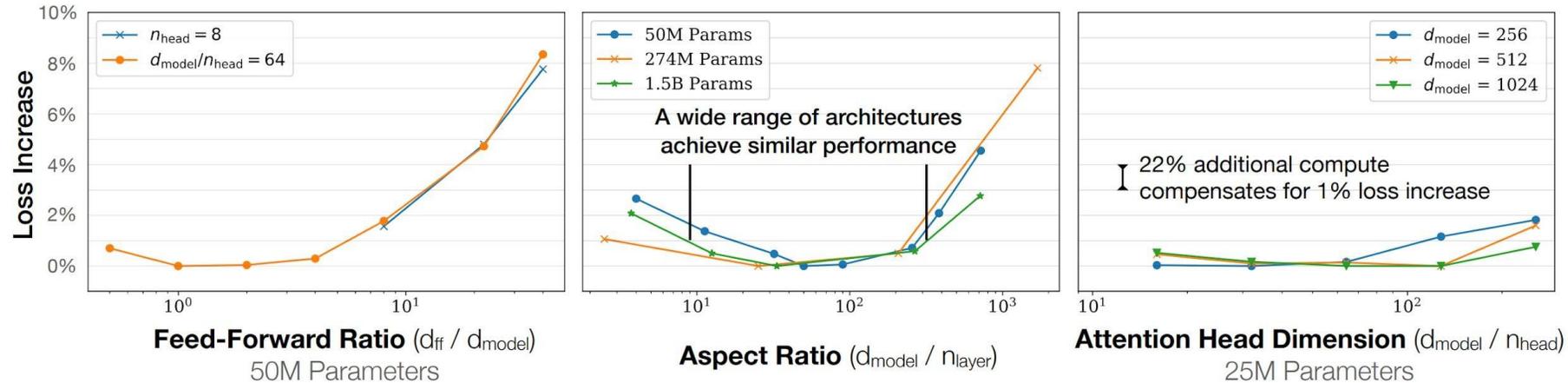
Language model loss changes with different network shape configurations [4].

- As long as the network shape is in a general sweet range, it does not impact performance much



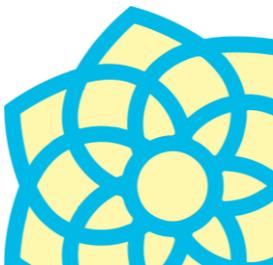
Scaling Law Study: Observations

Network shape (allocation of parameters at different parts) does not matter as much



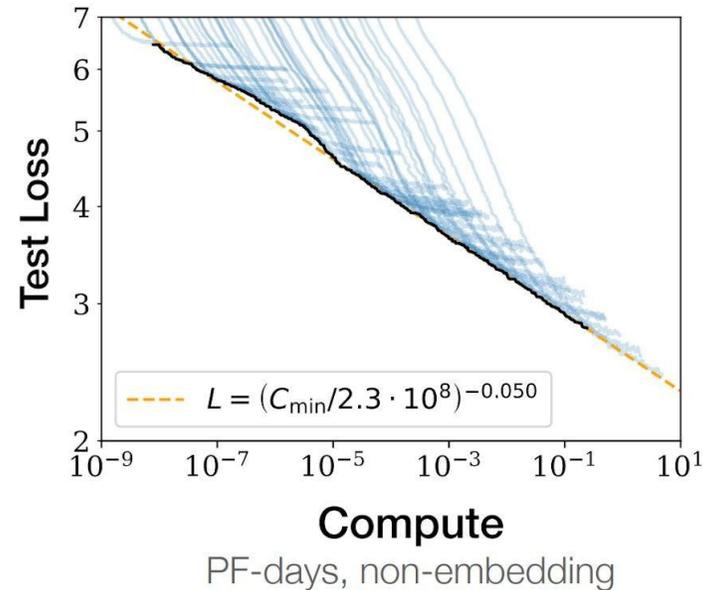
Language model loss changes with different network shape configurations [4].

- As long as the network shape is in a general sweet range, it does not impact performance much

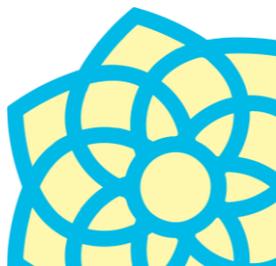


Scaling Law Study: Observations

A clear mapping from compute, data size, and parameter counts to testing loss

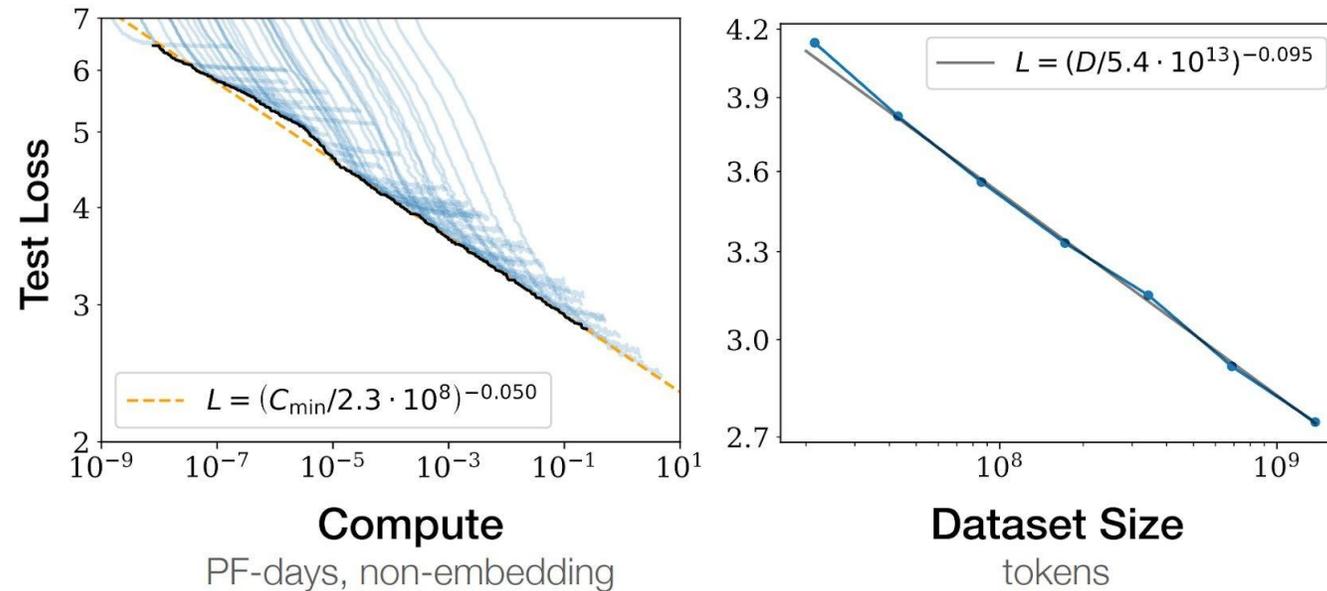


Mapping from compute (Peta-Flops days), data size, and model parameters to language modeling loss on testing data [4].

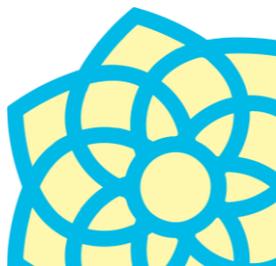


Scaling Law Study: Observations

A clear mapping from compute, data size, and parameter counts to testing loss

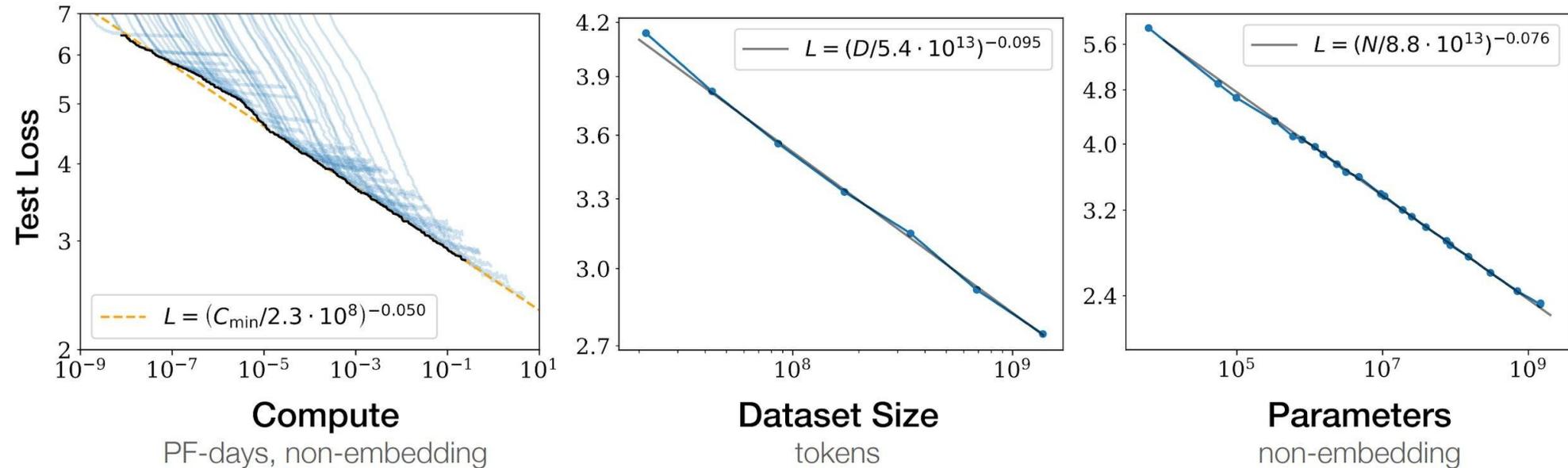


Mapping from compute (Peta-Flops days), data size, and model parameters to language modeling loss on testing data [4].

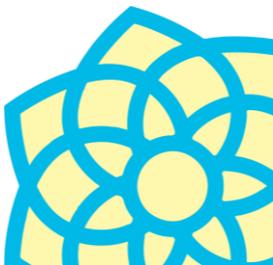


Scaling Law Study: Observations

A clear mapping from compute, data size, and parameter counts to testing loss

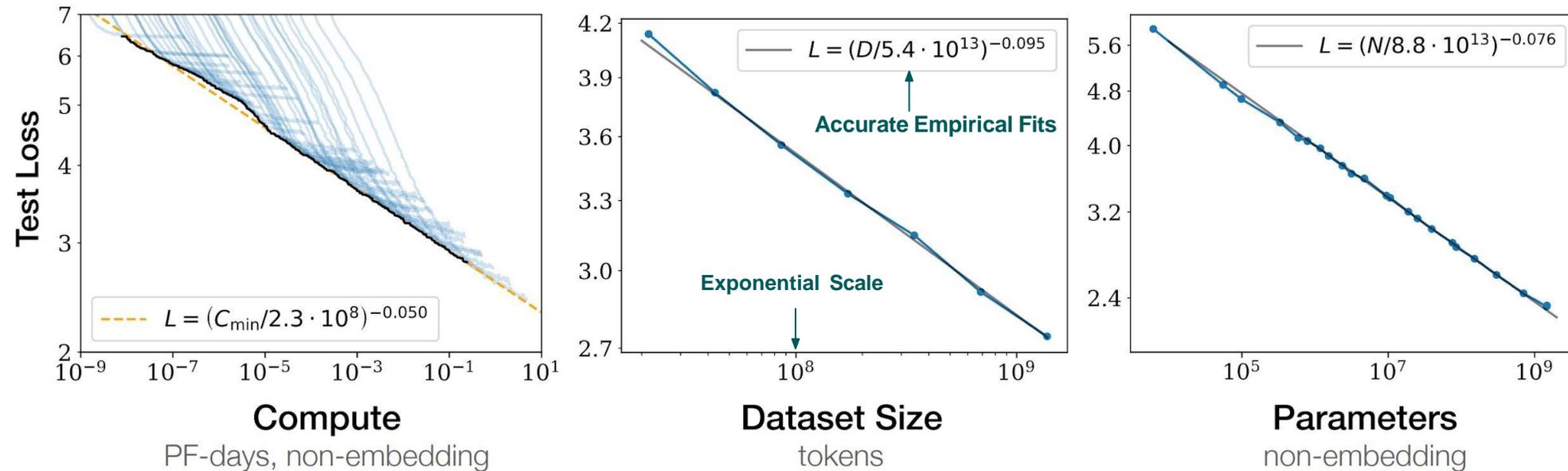


Mapping from compute (Peta-Flops days), data size, and model parameters to language modeling loss on testing data [4].



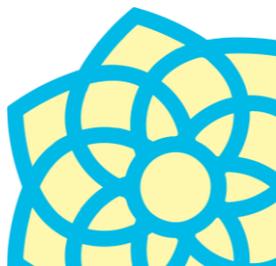
Scaling Law Study: Observations

A clear mapping from compute, data size, and parameter counts to testing loss



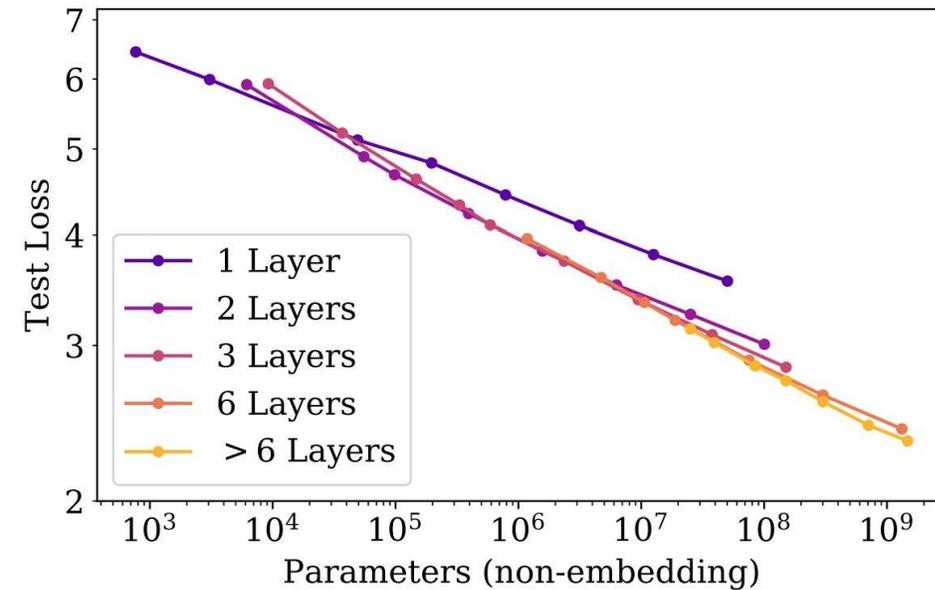
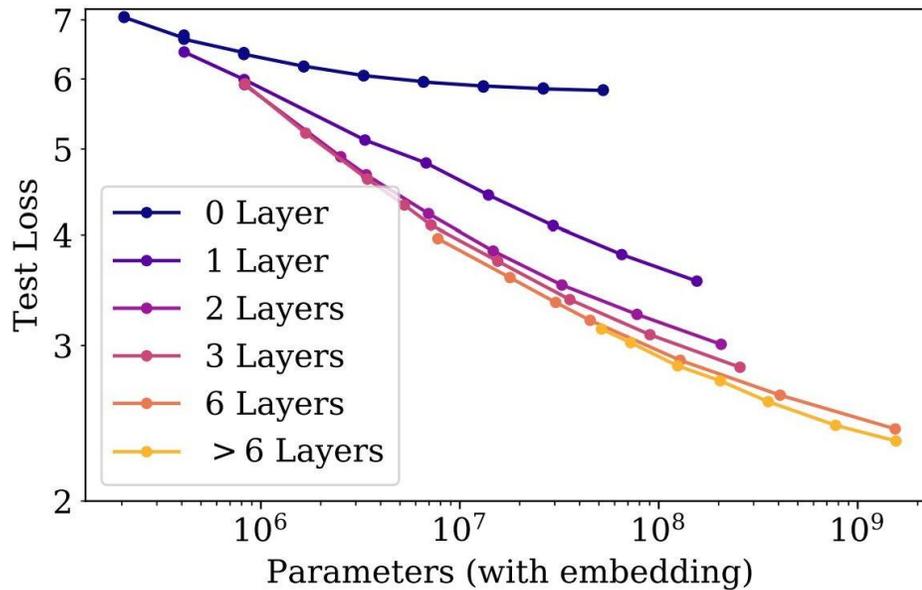
Mapping from compute (Peta-Flops days), data size, and model parameters to language modeling loss on testing data [4].

- Linear increasement of language modeling accuracy requires exponential scaling
- Three factors need to scale jointly to reach target model performance improvements

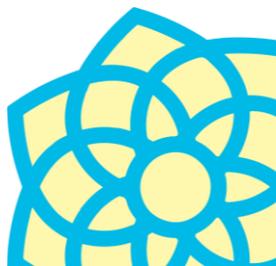


Scaling Law Study: Observations

Network parameters matter more than embedding parameters



Scaling law with network parameter counts include (left) and exclude (right) embeddings [4].



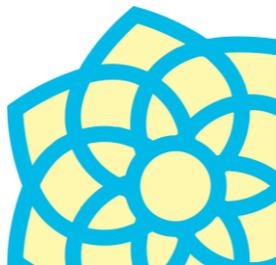
What Configurations to Scale Up

Goal: Given a computing budget and a candidate language model, select the optimal scaling up configurations

- E.g., One million H100 hours, pretrain the best LLaMA style LLM
- Configurations to choose: Model size (# of parameters) and pretraining data size (# of tokens)

A common question when scaling up

- Computing budget is the biggest constraint
- No room for exploration at target scale
- Only one scaled up pretraining run allowed, both budget-wise and time-wise.



What Configurations to Scale Up

Goal: Given a computing budget and a candidate language model, select the optimal scaling up configurations

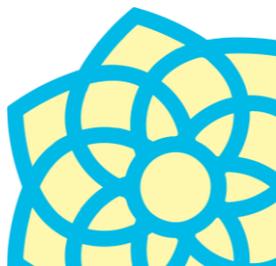
- E.g., One million H100 hours, pretrain the best LLaMA style LLM
- Configurations to choose: Model size (# of parameters) and pretraining data size (# of tokens)

A common question when scaling up

- Computing budget is the biggest constraint
- No room for exploration at target scale
- Only one scaled up pretraining run allowed, both budget-wise and time-wise.

Solution: Scaling law

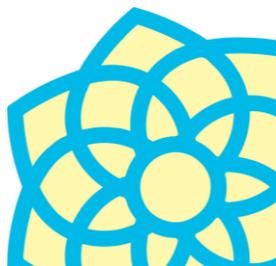
- Use many experiments at small scale to establish the scaling law
- Use scaling Law to predict best configuration at target compute



Scaling Configuration: Empirical Scaling Law

Empirical Approach #1: Fix model size and varying pretraining tokens [7]

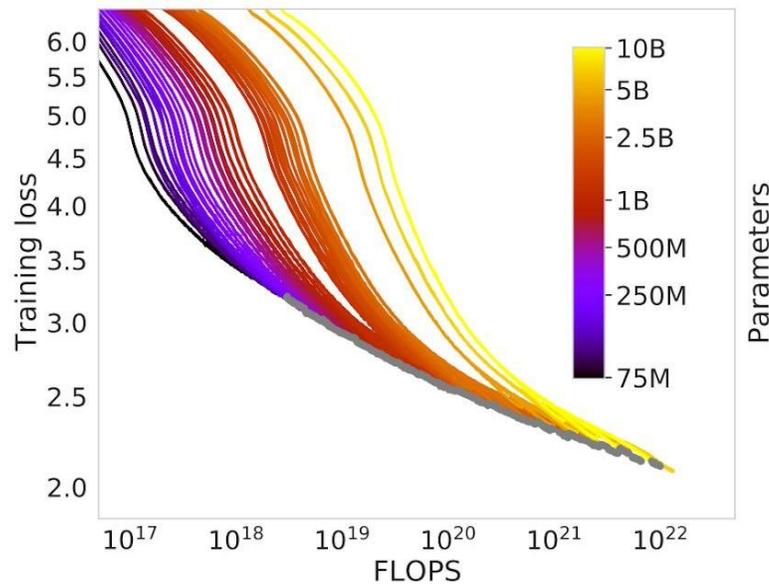
1. Pretrain different sized models to near converge and track loss
2. Record best (model size, data size) at each FLOP.
3. Estimate the scaling law



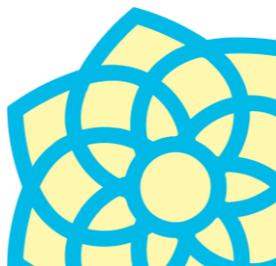
Scaling Configuration: Empirical Scaling Law

Empirical Approach #1: Fix model size and varying pretraining tokens [7]

1. Pretrain different sized models to near converge and track loss
2. Record best (model size, data size) at each FLOP.
3. Estimate the scaling law



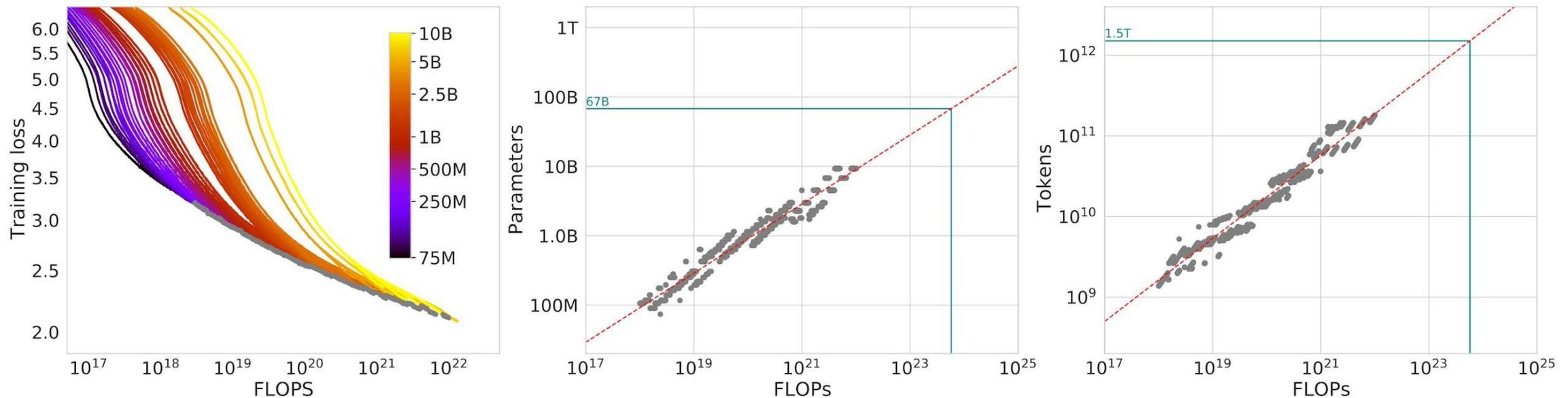
Pretraining loss of varying model (left), and the identified optimal parameters (mid) and tokens (right) at different FLOPS [6]



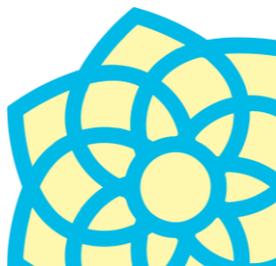
Scaling Configuration: Empirical Scaling Law

Empirical Approach #1: Fix model size and varying pretraining tokens [7]

1. Pretrain different sized models to near converge and track loss
2. Record best (model size, data size) at each FLOP.
3. Estimate the scaling law



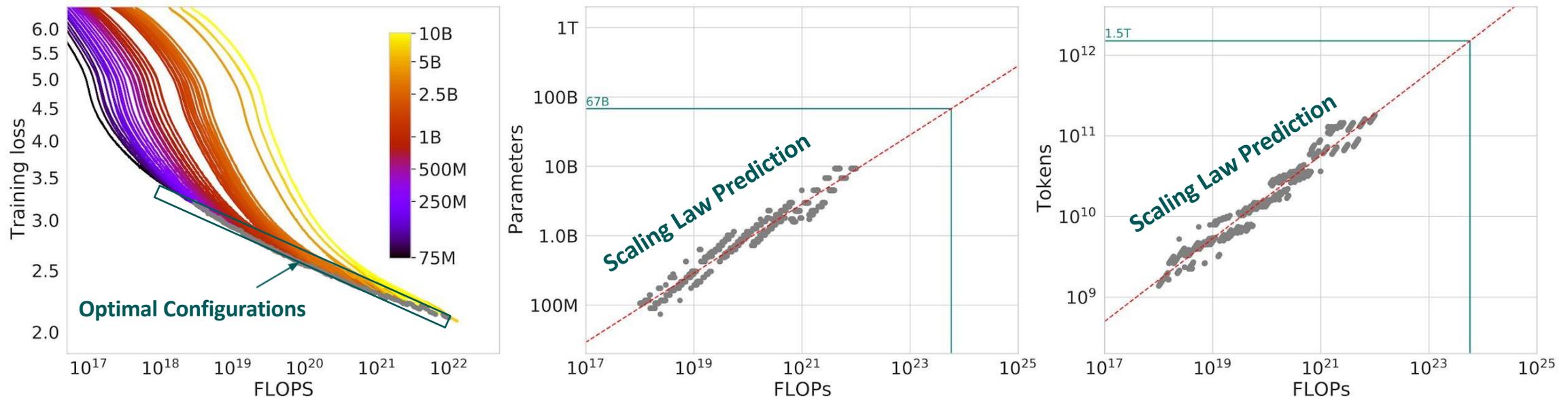
Pretraining loss of varying model (left), and the identified optimal parameters (mid) and tokens (right) at different FLOPs [6]



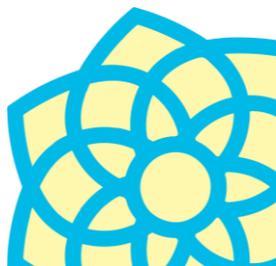
Scaling Configuration: Empirical Scaling Law

Empirical Approach #1: Fix model size and varying pretraining tokens [7]

1. Pretrain different sized models to near converge and track loss
2. Record best (model size, data size) at each FLOP.
3. Estimate the scaling law



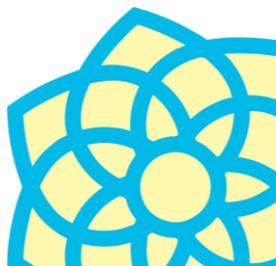
Pretraining loss of varying model (left), and the identified optimal parameters (mid) and tokens (right) at different FLOPs [6]



Scaling Configuration: Empirical Scaling Law

Empirical Approach #2: Fix total FLOPs, pretrain different sized models [7]

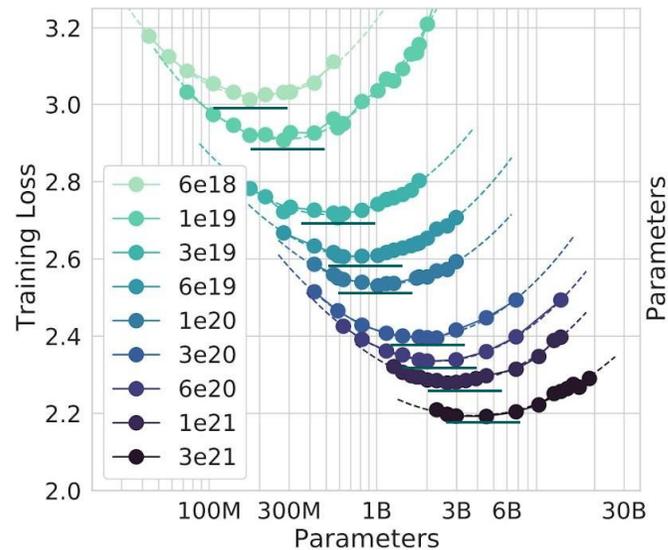
1. Pretrain to the # of tokens using total FLOPs and track final loss
2. Track best configurations and vary the total FLOPs and rerun #1
3. Estimate the scaling law



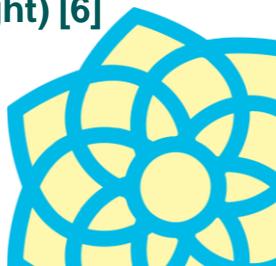
Scaling Configuration: Empirical Scaling Law

Empirical Approach #2: Fix total FLOPs, pretrain different sized models [7]

1. Pretrain to the # of tokens using total FLOPs and track final loss
2. Track best configurations and vary the total FLOPs and rerun #1
3. Estimate the scaling law



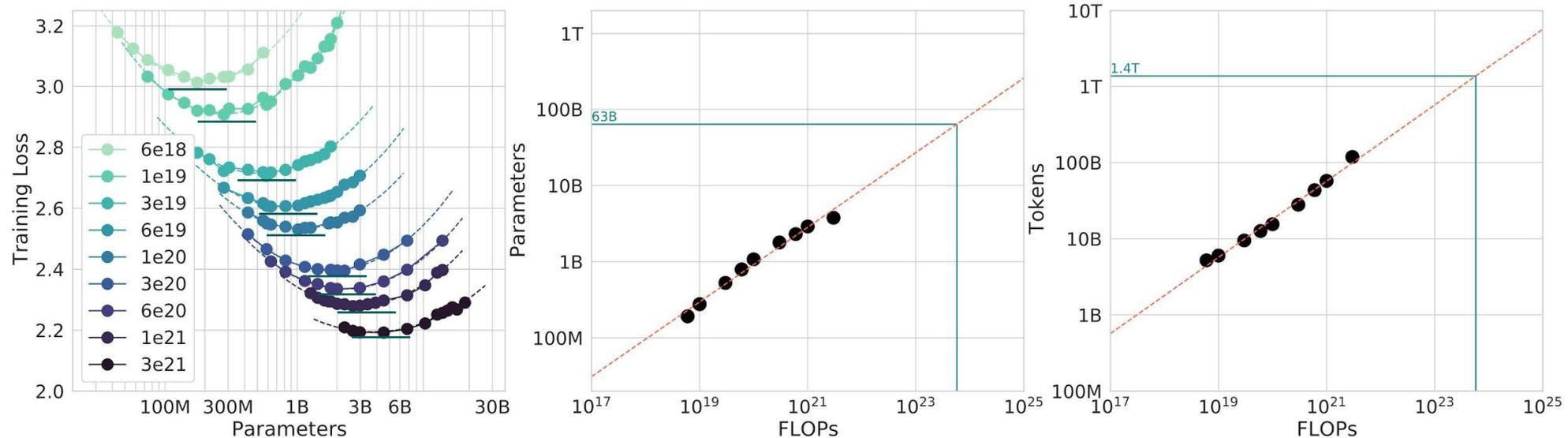
Pretraining loss of varying model sizes at varying FLOPs (left), and the identified optimal parameters (mid) and tokens (right) [6]



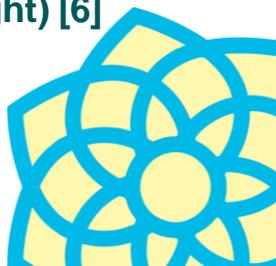
Scaling Configuration: Empirical Scaling Law

Empirical Approach #2: Fix total FLOPs, pretrain different sized models [7]

1. Pretrain to the # of tokens using total FLOPs and track final loss
2. Track best configurations and vary the total FLOPs and rerun #1
3. Estimate the scaling law



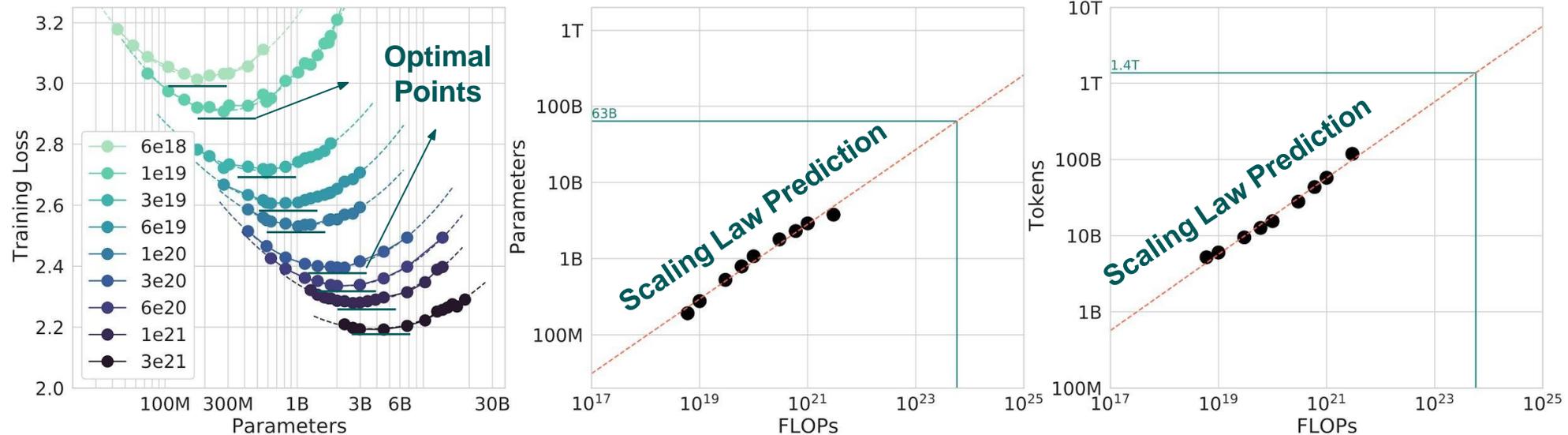
Pretraining loss of varying model sizes at varying FLOPs (left), and the identified optimal parameters (mid) and tokens (right) [6]



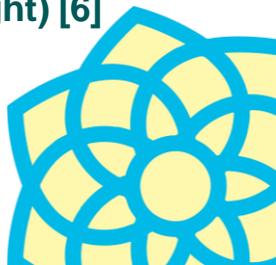
Scaling Configuration: Empirical Scaling Law

Empirical Approach #2: Fix total FLOPs, pretrain different sized models [7]

1. Pretrain to the # of tokens using total FLOPs and track final loss
2. Track best configurations and vary the total FLOPs and rerun #1
3. Estimate the scaling law

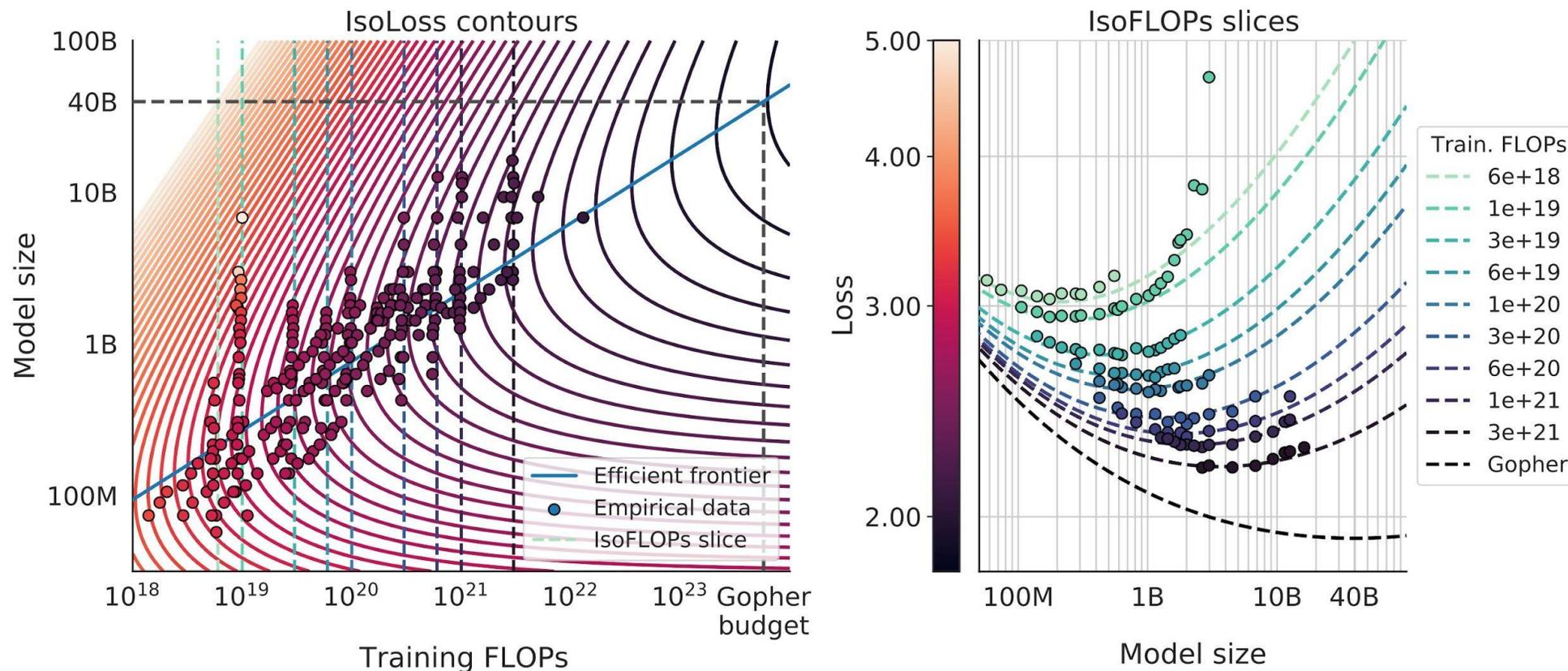


Pretraining loss of varying model sizes at varying FLOPs (left), and the identified optimal parameters (mid) and tokens (right) [6]

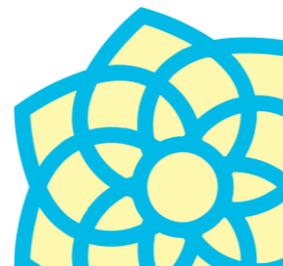


Scaling Configuration: Empirical Scaling Law

Empirical Approach #3: Using data points collected from previous two approaches and fix a parametric functions



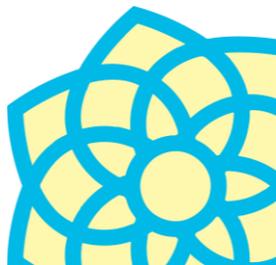
Fitted parametric function of (model size, FLOPs) \rightarrow Loss using data from approach one (left) and two (right) [6]



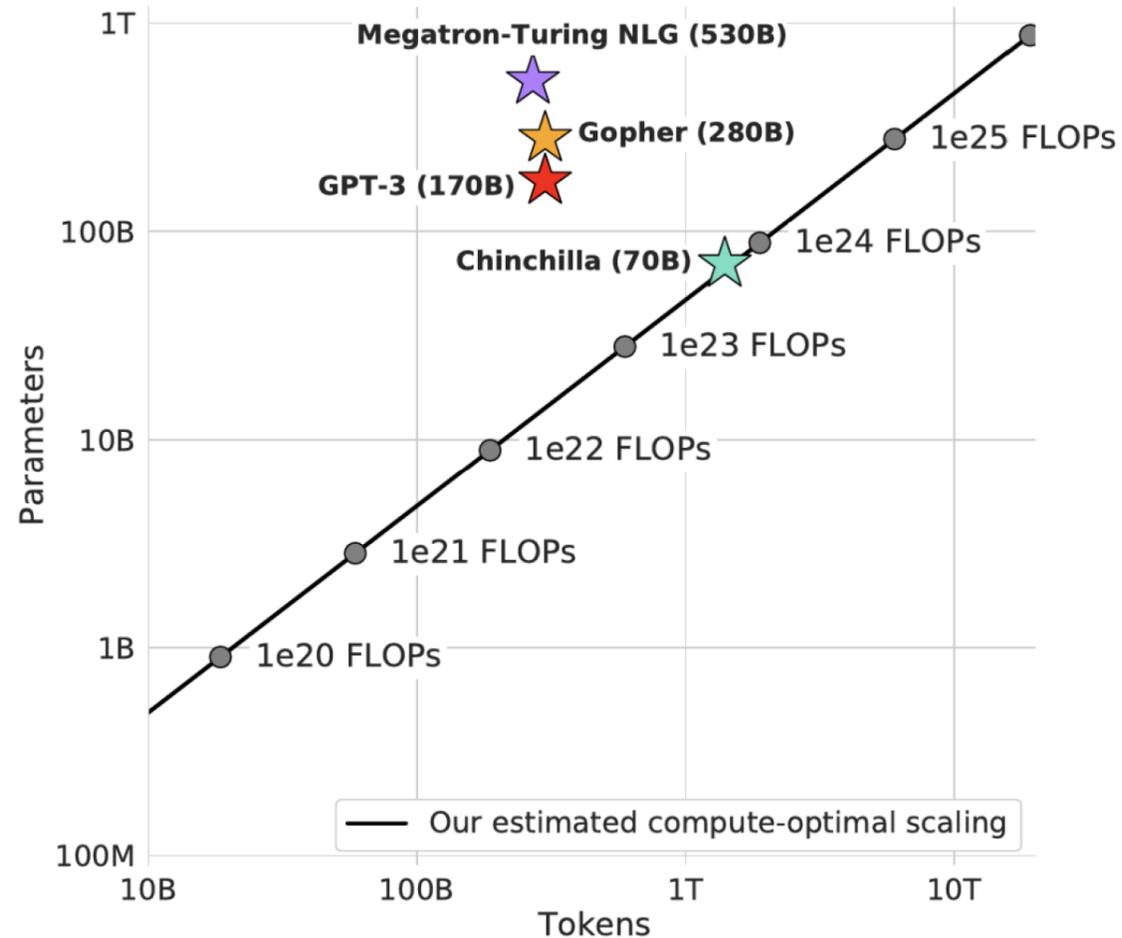
Scaling Configuration: Estimated Optimal Configurations

Parameters	FLOPs	FLOPs (in <i>Gopher</i> unit)	Tokens
400 Million	1.92e+19	1/29,968	8.0 Billion
1 Billion	1.21e+20	1/4,761	20.2 Billion
10 Billion	1.23e+22	1/46	205.1 Billion
67 Billion	5.76e+23	1	1.5 Trillion
175 Billion	3.85e+24	6.7	3.7 Trillion
280 Billion	9.90e+24	17.2	5.9 Trillion
520 Billion	3.43e+25	59.5	11.0 Trillion
1 Trillion	1.27e+26	221.3	21.2 Trillion
10 Trillion	1.30e+28	22515.9	216.2 Trillion

Examples of estimated scaling configurations at different model sizes [3].



Large language models can be too large

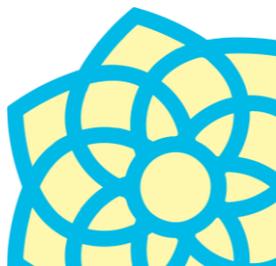


Recent models and its training tokens:

LlaMA-1: 1-1.4 T tokens

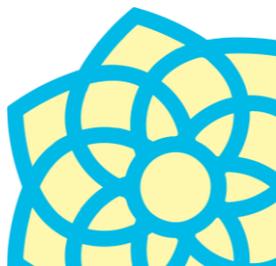
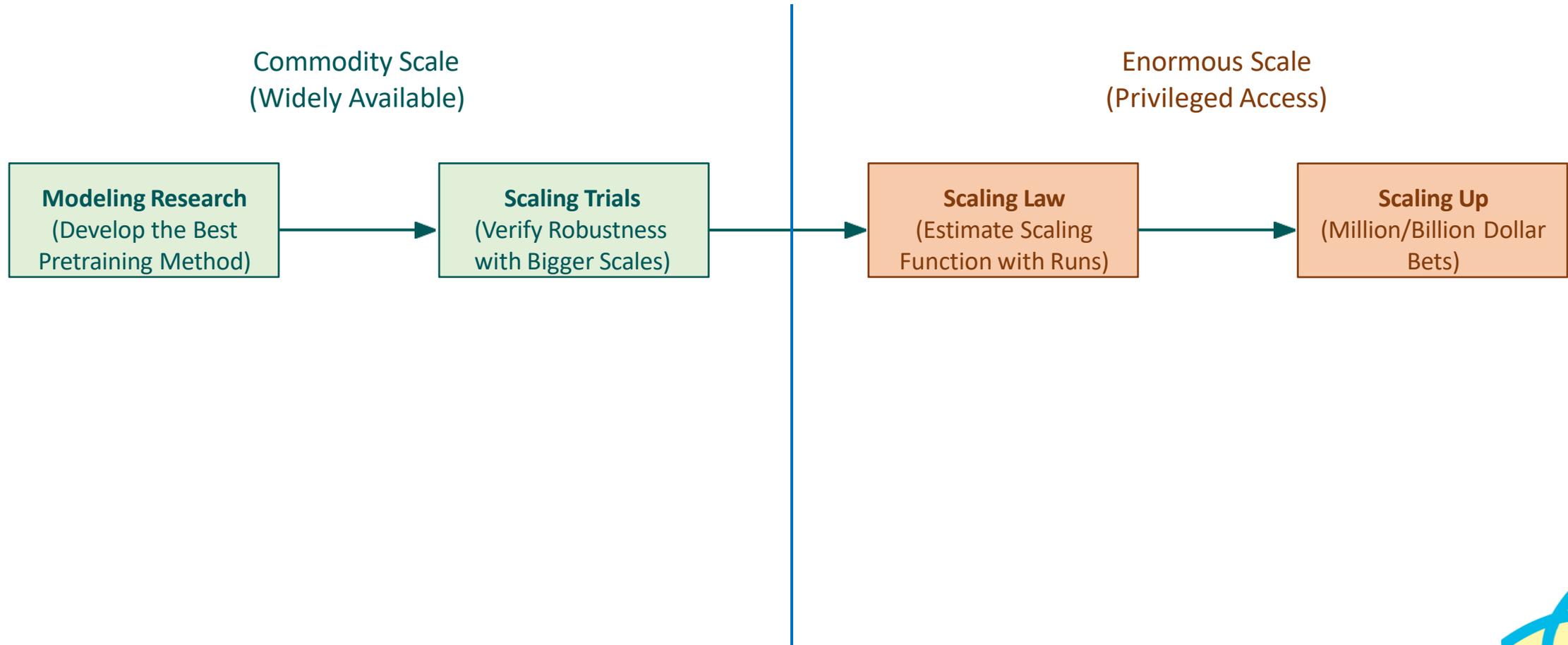
LlaMA-2: 2T tokens

Mistral-7B: much more...



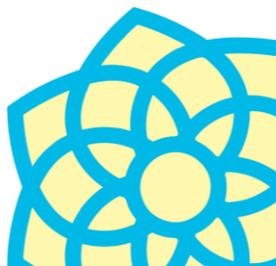
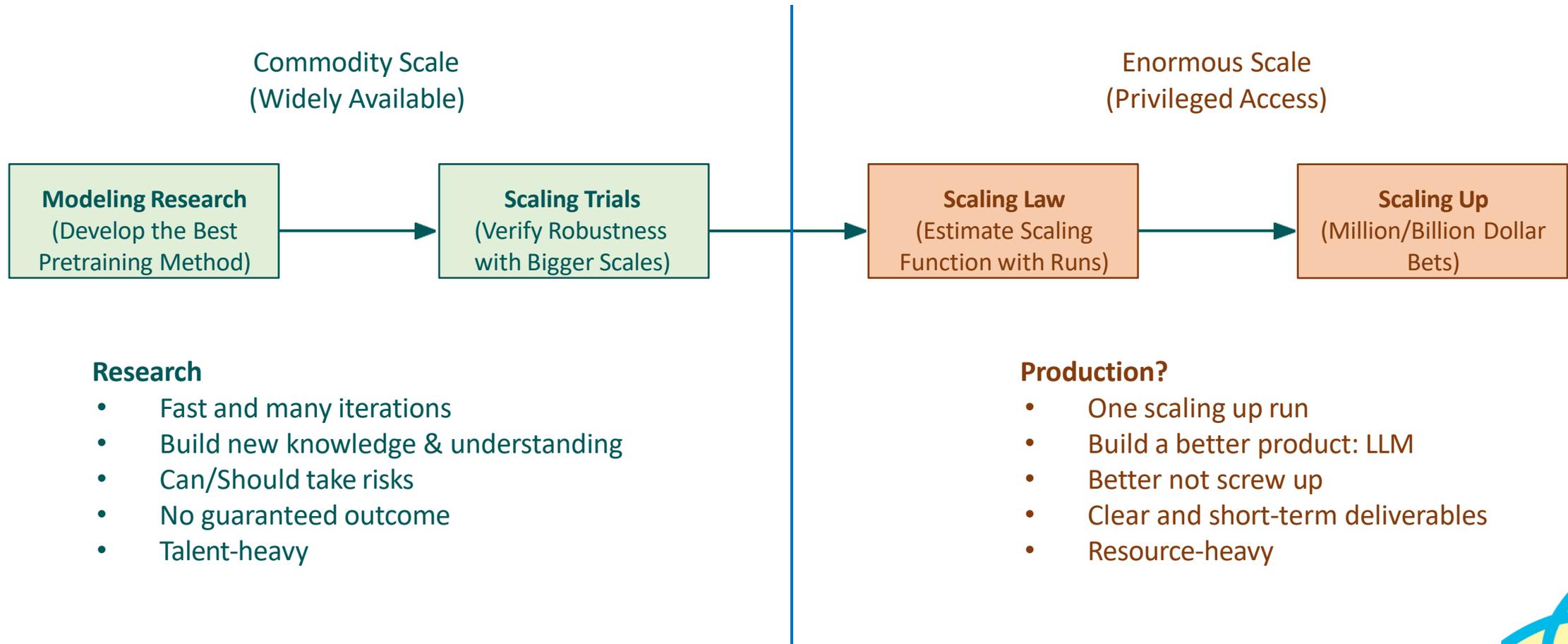
Scaling Up Pipeline

The current development pipeline of scaling up LLM pretraining, e.g., used by GPT-4, PaLM-2, and many more

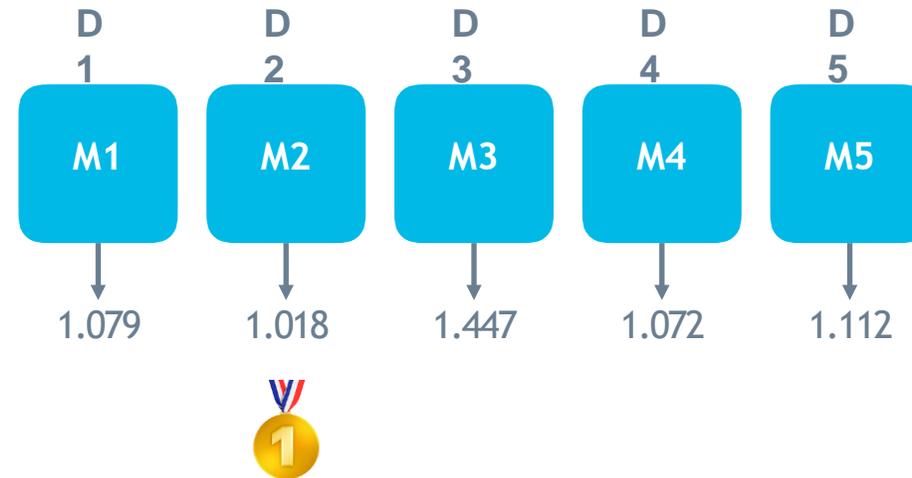


Scaling Up Pipeline

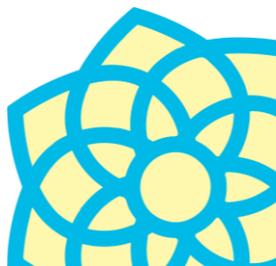
The current development pipeline of scaling up LLM pretraining, e.g., used by GPT-4, PaLM-2, and many more



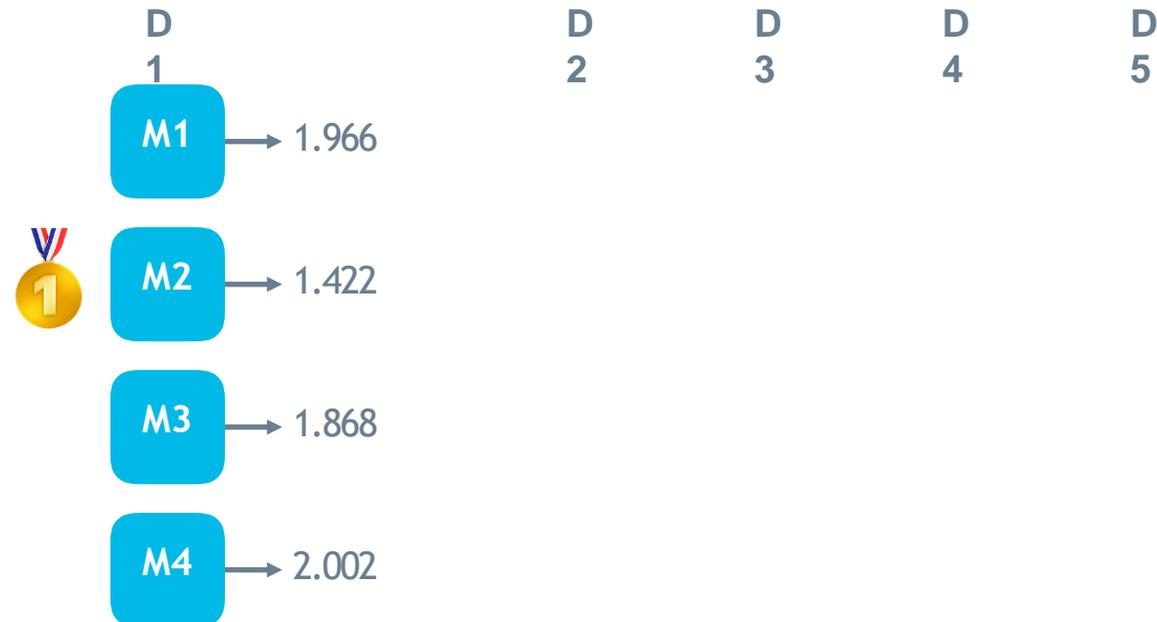
Putting scaling laws into practice



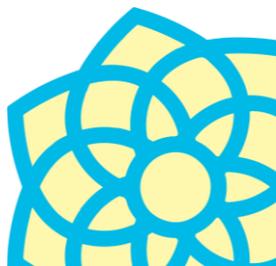
old paradigm: train a few models, select the best one



Putting scaling laws into practice

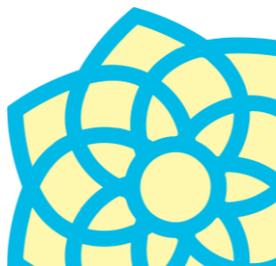


new paradigm: train many small models, up-scale the best one



Scaling Law: Summary

- Why Scaling Up
 - Predictable benefits in nearly all scenarios
- Which Language Model to Scale Up
 - Benefits of decoder models
- What Factors Matter in Scaling
 - Strong mapping from compute, model size, and pretraining data size to language model performances
- What Configurations to Scale Up
 - Establish scaling law with small scale explorations, scaling up based on scaling law predictions
- Capabilities Emerged from Scaling Up
 - Lots of unknowns and challenges!



- **LAIM LE4 VT2025:**
Pre-training
Parallel Training
Scaling

www.ida.liu.se/~frehe08/llm