

Natural Language Processing



Scaling and Systems

Kevin Lin – UC Berkeley
March 22, 2023

Scaling and Systems



Announcements

- Spring Break Next Week
- Panels Week After Spring Break
- HW4 Bug
- HW5 Testing
- Today
 - What to scale?
 - How to scale?

Scaling With Fixed Compute

ed

CS 288 – Ed Discussion

hahhah

I made the network with 4096 hidden units. It finally achieved 66% accuracy!

I guess brute force really works.

♡ Reply Edit Delete ...

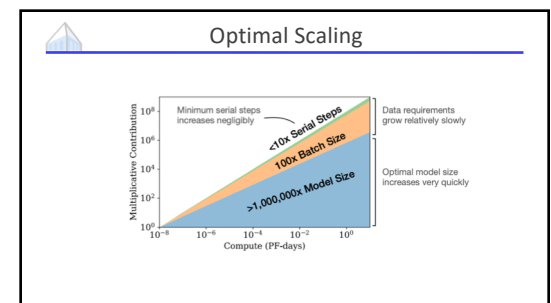
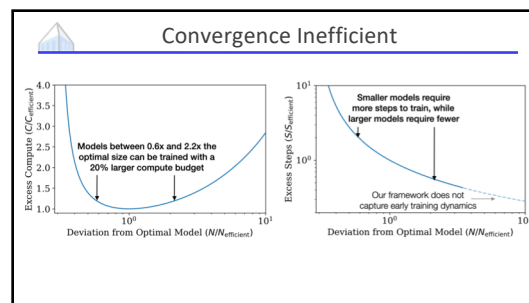
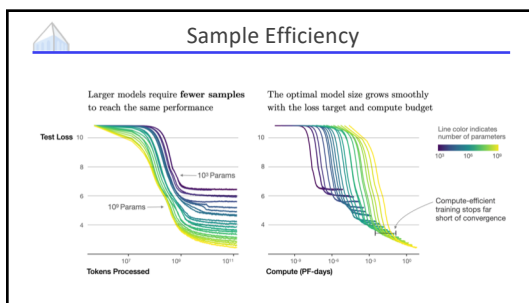
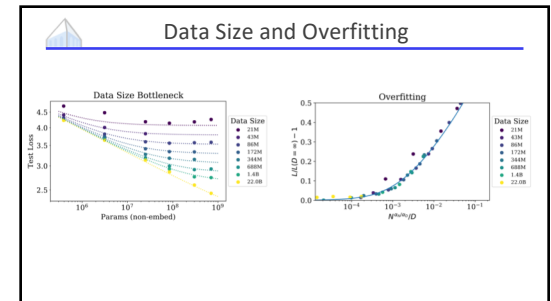
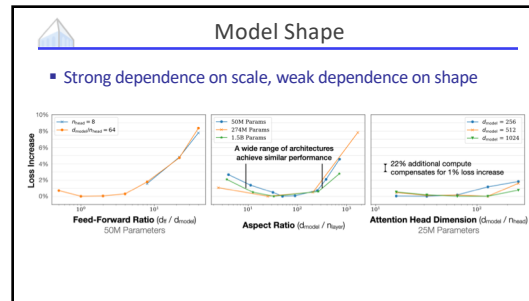
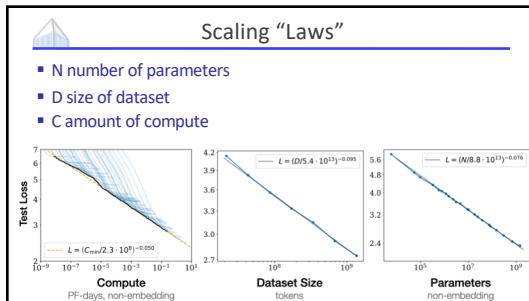
Scaling With Fixed Compute

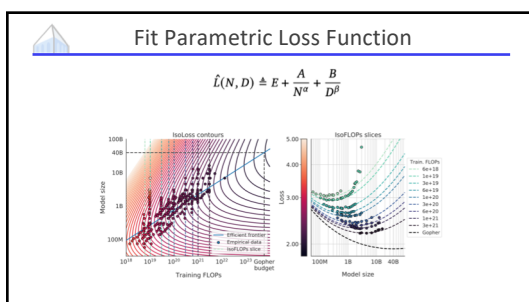
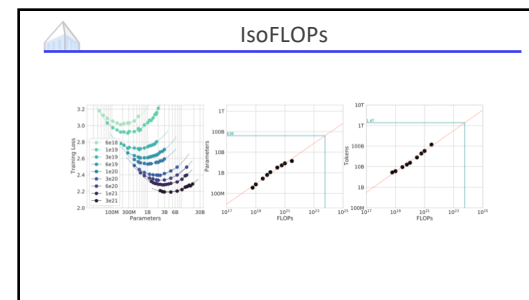
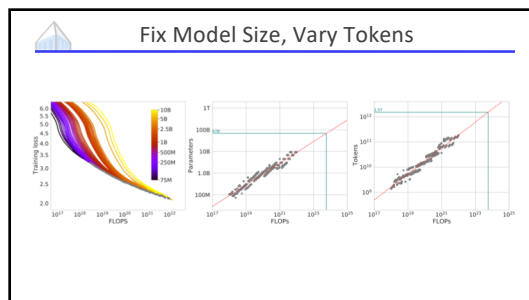
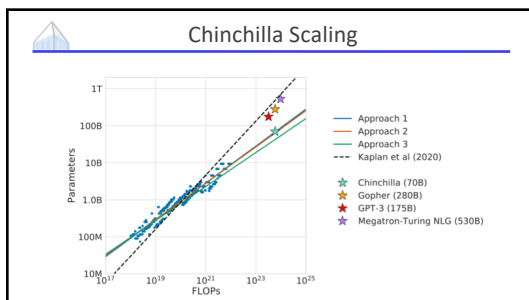
Model	Size (# Parameters)	Training Tokens
LaMDA (Thoppilan et al., 2022)	137 Billion	168 Billion
GPT-3 (Brown et al., 2020)	175 Billion	300 Billion
Jurassic (Lieber et al., 2021)	178 Billion	300 Billion
Gopher (Rae et al., 2021)	280 Billion	300 Billion
MT-NLG 530B (Smith et al., 2022)	530 Billion	270 Billion
Chinchilla	70 Billion	1.4 Trillion

Fixed Compute

Scaling Laws for Neural Language Models (Kaplan et al., 2020)

- N – the number of model parameters, *excluding all vocabulary and positional embeddings*
- $C \approx 6NBS$ – an estimate of the total non-embedding training compute, where B is the batch size, and S is the number of training steps (ie parameter updates). We quote numerical values in PF-days, where one PF-day = $10^{15} \times 24 \times 3600 = 8.64 \times 10^{19}$ floating point operations.





BIG-bench

1. Step Inference: Given a prompt goal and four candidate steps, choose the correct step that helps achieve the goal.
For example:
Which step is likely to help achieve the goal "prevent coronavirus"?
A. wash your hands B. wash your cat C. clap your hands D. eat your protein

2. Goal Inference: Given a prompt step and four candidate goals, choose the correct goal that the step helps achieve.
For example:
Which is the most likely goal of "choose a color of lipstick"?
A. get pink lips B. read one's lips C. lip sync D. draw lips

3. Step Ordering: Given a prompt goal and two steps, determine which step temporally precedes the other.
For example:
In order to "clean silver," which step should be done first?
A. dry the silver B. handwash the silver

Task: goal_step_wikihow

BIG-bench

1. When Max was applying for a new job after he got fired, he wrote his resume and cover letters so creatively that he got offered a copywriter job. Which of the following proverbs best apply to this situation?

2. Where there is a will, there is a way.

3. Where one door shuts, another opens

4. You can catch more flies with honey than with vinegar

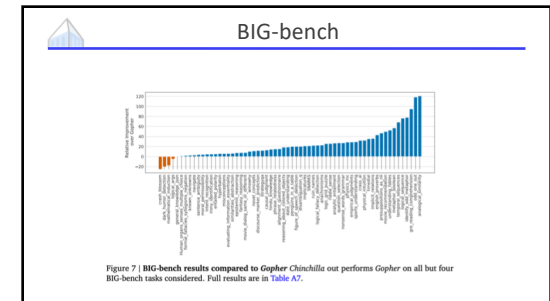
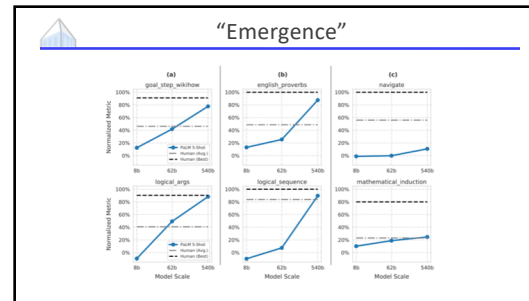
Task: English_proverbs

BIG-bench

Turn right. Take 1 step. Turn right. Take 6 steps. Turn right. Take 1 step. Turn right. Take 2 steps. Take 4 steps.

Are you back at the start?

Task: **navigate**



Are We Optimal?

	Compass budget (G)	Series	alpha	beta	γ	δ	A (HBM)	B (HBM)	C (shared mem)	D (shared mem)	E, HBM250 optimal	F, HBM250 optimal	G, HBM250 optimal
OPT-175B	4,300-05	From eq 12 to Approach 5.0	0.34	0.28	0.402	0.048	458.4	410.7	1.0000	28.28	2940.88		
	4,300-05	Approach 1: 16x80GB training summs	0.34	0.34	0.5	0.5	458.4	410.7	1.0000	28.28	2940.88	271.80	
	4,300-05	Approach 2: 16x80GB training summs	0.28	0.28	0.5	0.5	458.4	410.7	1.0000	28.28	2940.88	271.80	
	4,300-05	Approach 3: 16x80GB profiles	0.34	0.33	0.40	0.51	458.4	410.7	1.0000	28.28	2940.88	433.08	
	4,300-05	Approach 4: 16x80GB profiles	0.28	0.27	0.40	0.51	458.4	410.7	1.0000	28.28	2940.88	429.48	
	4,300-05	Approach 5: 16x80GB profiles	0.34	0.29	0.46	0.54	458.4	410.7	1.0000	28.28	2940.88	415.08	1555.08
	4,300-05	Approach 6: 16x80GB profiles	0.28	0.28	0.46	0.54	458.4	410.7	1.0000	28.28	2940.88	415.08	1540.88

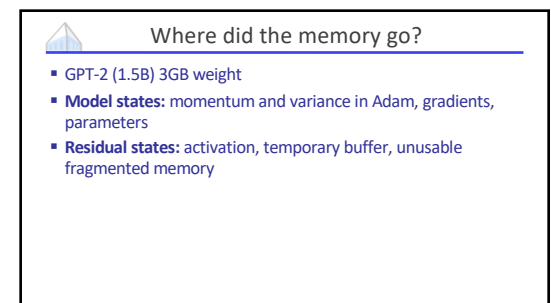
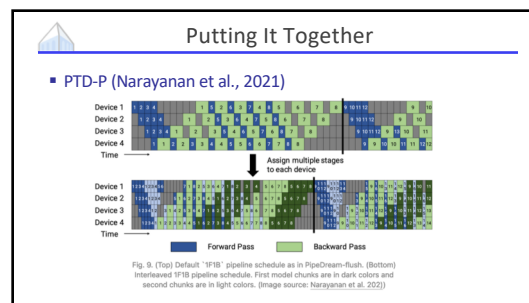
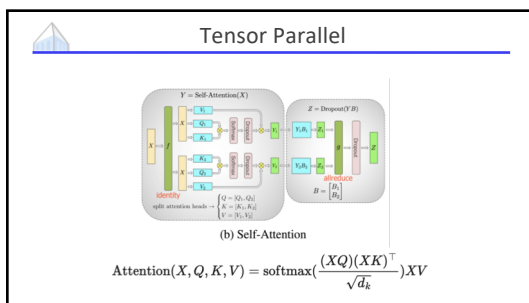
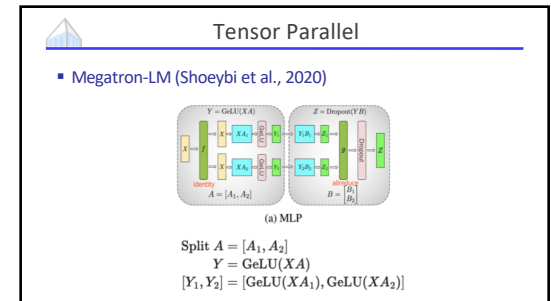
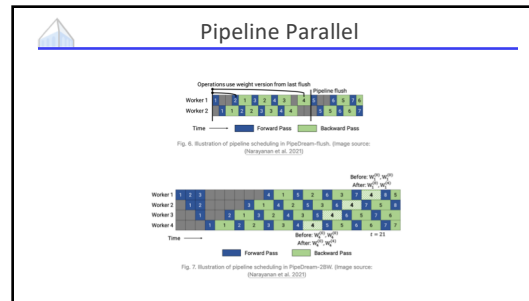
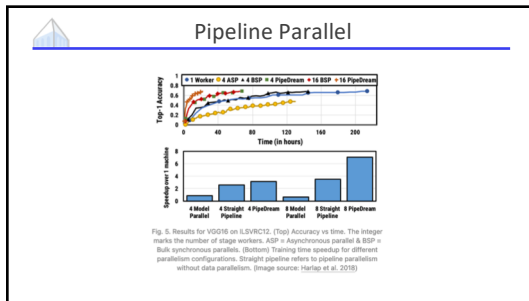
Source: Susan Zhang

Systems

- Parallelism
 - Data
 - Model
 - Tensor
- Memory Optimization

Data Parallel

- Bulk synchronous parallel: sync at end of every minibatch
 - Pros: higher learning efficiency
 - Cons: wait for all machines
- Asynchronous parallel: apply updates when ready
 - Pros: No wait
 - Cons: lower learning efficiency





Model States

- *Example.* Transformer architecture trained with Adam
- Ψ parameters with mixed precision training (use F16 and F32)
- F16 copies of **params (2 Ψ bytes)** and **gradients (2 Ψ bytes)**
- F32 copies of **params (4 Ψ bytes)**, **momentum (4 Ψ bytes)** and **variance (4 Ψ bytes)**
- = **16 Ψ bytes**, at least **24GB**



Residual States

- **Activations:** 1.5B transformer, around 60 GB even with activation checkpointing
- **Temporary buffers:** gradient all-reduce, norm computation, etc. around 5GB
- **Memory fragmentation:** 30% of memory still available when OOM

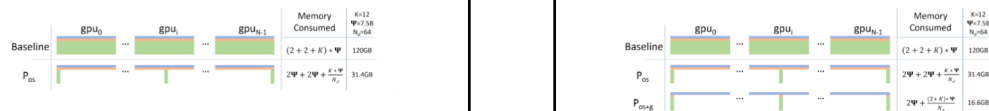


ZeRO Redundancy Optimizer

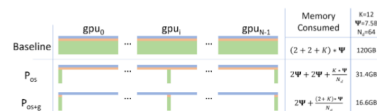
- Memory Optimizations Towards Training Trillion Parameter Models (Rajbhandari et al., 2019)
- **ZeRO-DP** optimizes model states
- **ZeRO-R** optimizes residual states



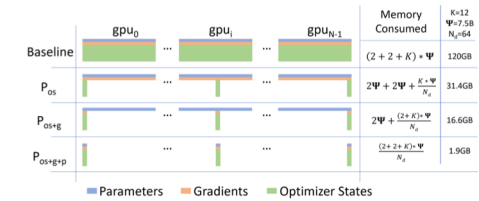
ZeRO-DP



ZeRO-DP



Memory Optimization





ZeRO-R

- **Partitioned Activation Checkpointing:** Once forward prop for a layer is computed, partition the input activations until needed for backprop
- **Constant size buffer:** computational efficiency can depend on input size, eg. All-reduce achieves higher bandwidth than a smaller one
- **Memory Defragmentation:** pre-allocate contiguous memory chunks for activation checkpoints



Scaling for Varying Model Sizes

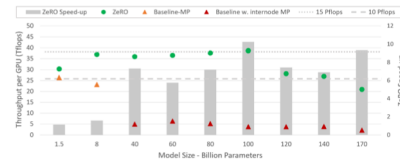


Figure 2: *ZeRO* training throughput and speedup w.r.t SOTA baseline for varying model sizes. For *ZeRO*, the MP always fit in a node, while for baseline, models larger than 40B require MP across nodes.



Superlinear Scaling for Increasing GPUs

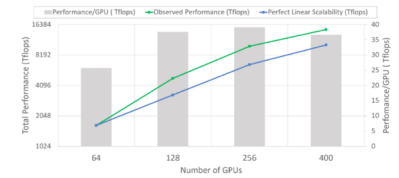


Figure 3: Superlinear scalability and per GPU training throughput of a 60B parameter model using *ZeRO*-100B.



Turing-NLG

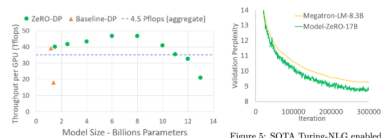


Figure 4: Max model throughput with *ZeRO*-DP.

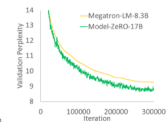


Figure 5: SOTA Turing-NLG enabled by *ZeRO*.