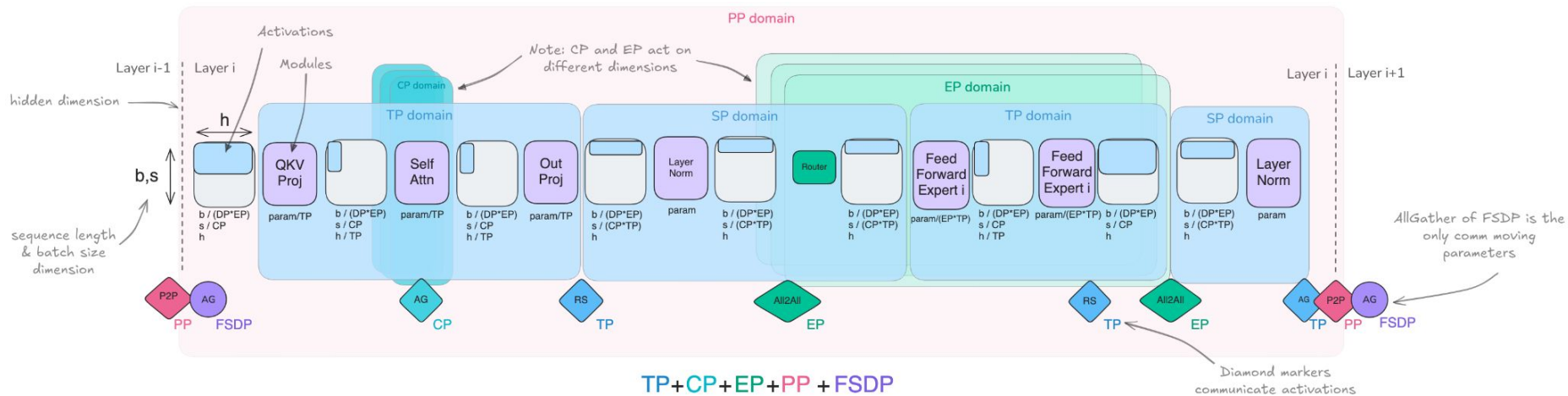


# Finding the Best Training Configuration

made with  for “Little ML book club”





	< 10B	10-100B	> 100B
small GPU scale	tensor or DP	tensor + pipeline or tensor + DP or ZeRO-3	?
large GPU scale	tensor + DP + pipeline		



- Enable recomputation
- Gradient accumulation
- Stop buying coffee out and save for more GPUs

# Achieving the target global batch size

Surprise-surprise



# Achieving the target global batch size

Surprise-surprise

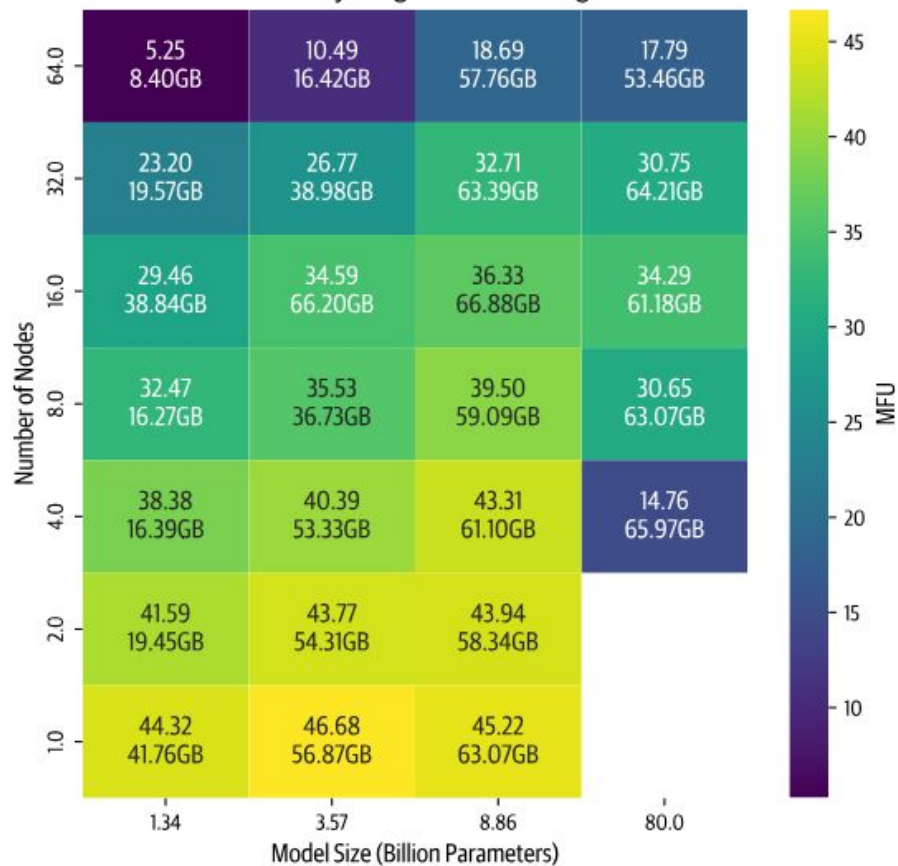
DP + gradient accumulation + context



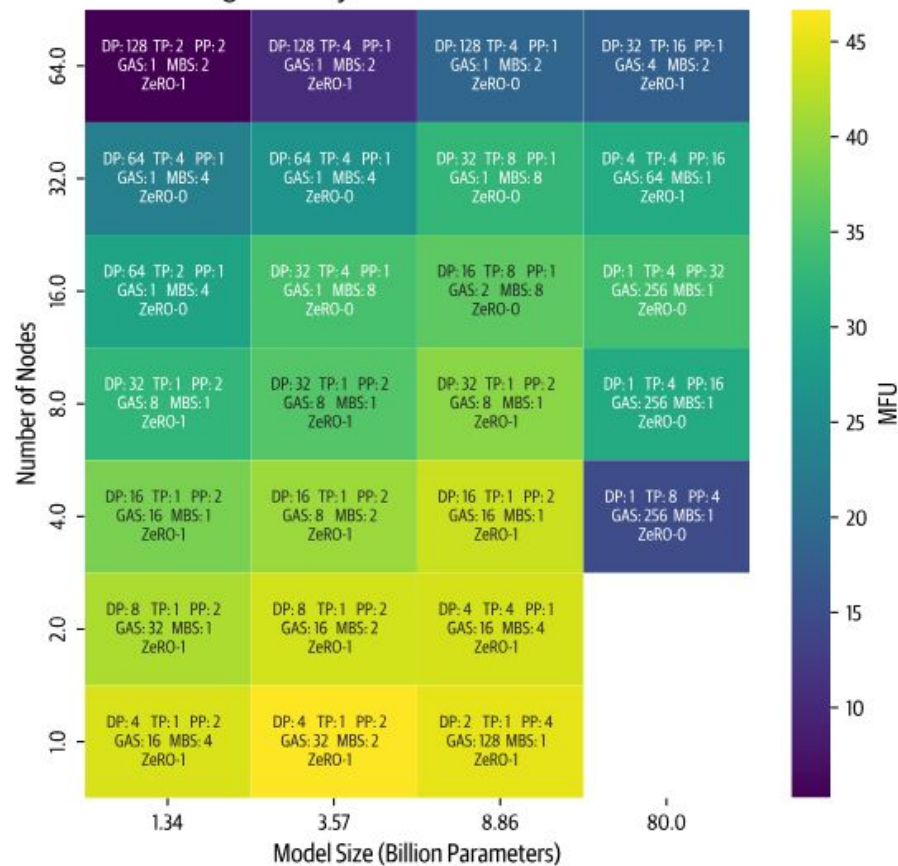
# Optimizing training throughput

1. Tensor
2. ZeRO-3
3. DP -> PP
4. Scale one technique at a time
5. Experiment with micro-batch sizes

### MFU and Memory Usage for Best Configurations



### Best Configuration by Model Size and Number of Nodes



## 1. Meta LLaMA Family

Model	Date	Parameters	Hardware	TP	PP	DP/ FSDP	CP	EP	Key Innovations
LLaMA 1	Feb 2023	7B-65B	2,048 A100 80GB	—	—	DP	—	—	Basic data parallelism; RSC cluster
LLaMA 2	Jul 2023	7B-70B	RSC + prod clusters	—	—	FSDP	—	—	Introduced FSDP; GQA for 70B; 4K context; 1.73M GPU-hours for 70B
LLaMA 3	Apr 2024	8B-70B	16,384 H100	8	16	FSDP (128)	1	—	4D parallelism; 8K context; 126 layers (not 128) for balanced PP
LLaMA 3.1	Jul 2024	8B-405B	16,384 H100	8	16	FSDP (128)	1- 16	—	128K context via CP=16; all-gather CP (not ring attention); 38-43% MFU

GPUs	TP	CP	PP	DP	Seq. Len.	Batch size/DP	Tokens/Batch	TFLOPs/GPU	BF16 MFU
8,192	8	1	16	64	8,192	32	16M	430	43%
16,384	8	1	16	128	8,192	16	16M	400	41%
16,384	8	16	16	8	131,072	16	16M	380	38%

**Table 4 Scaling configurations and MFU for each stage of Llama 3 405B pre-training.** See text and Figure 5 for descriptions of each type of parallelism.

## 2. Google PaLM/Gemini Family

Model	Date	Parameters	Hardware	TP	PP	DP	CP	EP	Key Innovations
PaLM	Apr 2022	540B	6,144 TPU v4	12	None	256 (2D FSDP)	—	—	Pipeline-free; 57.8% HW utilization; Pathways system
PaLM 2	May 2023	Undisclosed	TPU v4	✓	—	✓	—	✓ (sparse)	MoE architecture; improved compute-optimal scaling
Gemini 1.0 Ultra	Dec 2023	Undisclosed	Multi-DC TPU v4/ v5e	✓	—	✓	—	✓	<b>Multi-datacenter training</b> ; 97% goodput; optical circuit switching
Gemini 1.5 Pro	Feb 2024	Undisclosed (MoE)	TPU v5+	✓	—	✓	✓	✓	Sparse MoE; up to 1M context; long-context specialization

### 3. DeepSeek Family

Model	Date	Total Params	Active Params	Hardware	TP	PP	DP	EP	Key Innovations
DeepSeek 67B	Jan 2024	67B	67B (dense)	H800 cluster	✓	✓	ZeRO	—	Baseline dense model
DeepSeek- V2	May 2024	236B	21B	H800 cluster	—	16 (ZeroBubble)	ZeRO-1	8	MLA attention; DeepSeekMoE 42.5% cost reduction vs 67B
DeepSeek- V3	Dec 2024	671B	37B	2,048 H800	None	16 (DualPipe)	ZeRO-1	64	Aux-loss-free balancing; FP8; <b>\$5.6M total cost</b> ; 180K GPU-hr/T tokens

see you next time