



COMPUTING AND
INFORMATION TECHNOLOGY
Research Computing and Data

LLM fine-tuning

How can we maximize computational efficiency?

Batch size considerations

Batch size is the number of training samples processed simultaneously during training.

Batch size affects how much memory is used during training, **and** affects training outcomes.

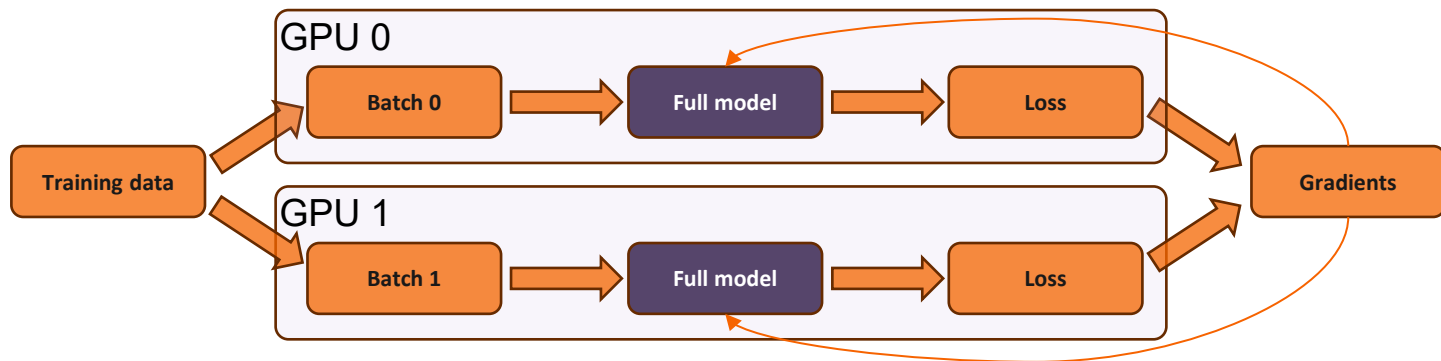
Batch size	Advantages	Disadvantages
Large	<ul style="list-style-type: none">• Faster training per epoch• More stable gradients• Better GPU utilization	<ul style="list-style-type: none">• Risk of poor generalization• Can lead to training instability• High memory usage
Small	<ul style="list-style-type: none">• More frequent updates and faster convergence• Better generalization• Lower memory usage	<ul style="list-style-type: none">• Noisier gradients and training instability• Slower per-epoch training• May require gradient accumulation to match large batch performance



Data parallelism

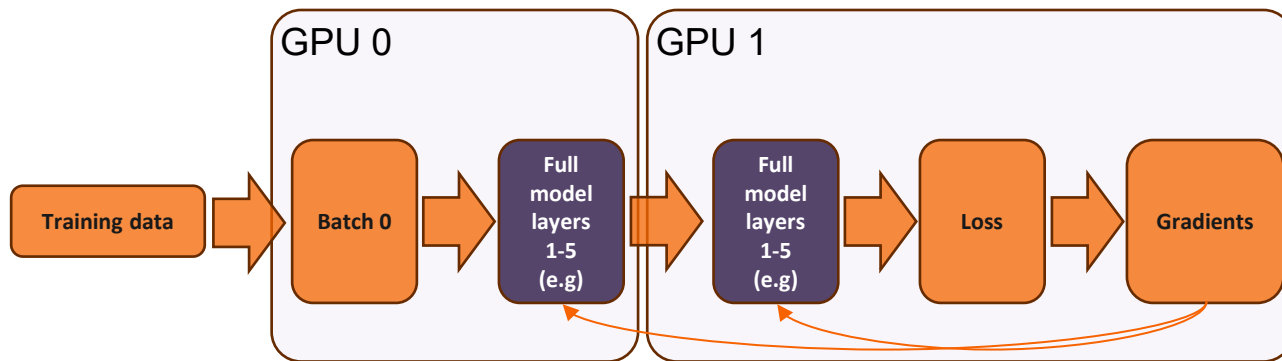
There are many ways to use **multiple GPUs** during training, with different goals and outcomes.

Data parallelism puts a copy of the model on **each** GPU, speeding training.



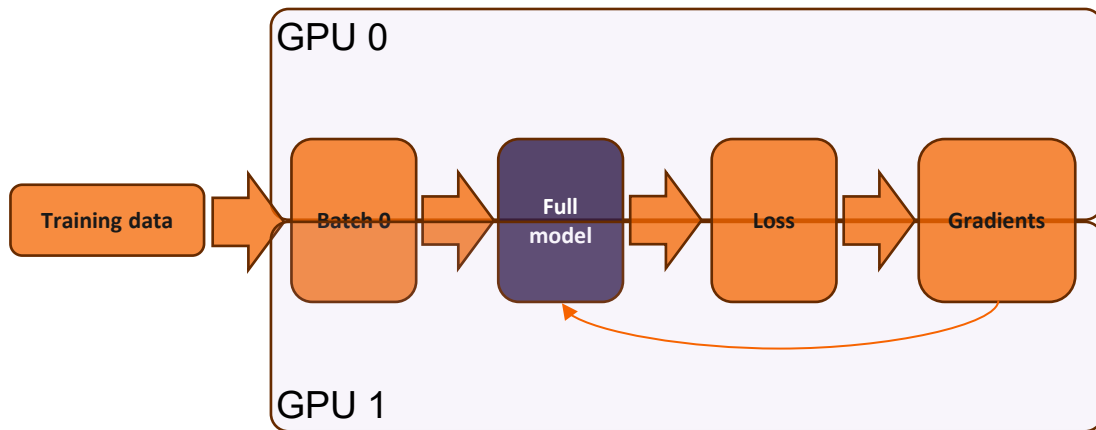
Model parallelism

Model parallelism puts a some of the layers of the model on each GPU, allowing the use of a model that is too large to fit on a single GPU.



Tensor parallelism

Tensor parallelism splits individual layers of the model across GPUs, allowing the use of a model with individual layers so large they can't fit on a single GPU.



Comparison of parallelism strategies

Parallelism type	What it splits	Best for	Pros	Cons
Data parallelism	Dataset (batch split across GPUs)	<ul style="list-style-type: none">• Large datasets• Small to medium models	<ul style="list-style-type: none">• Easy to implement• Works with most models	<ul style="list-style-type: none">• Requires full model copy on each GPU• Some overhead
Model parallelism	Model (some layers on each GPU)	Large models that don't fit on one GPU	<ul style="list-style-type: none">• Reduces memory load per GPU• Usually easy to implement	<ul style="list-style-type: none">• Slower due to inter-GPU communication• Can be tough to implement
Tensor parallelism	Individual layers (weights split across models)	Extremely large models with layers too big for 1 GPU	<ul style="list-style-type: none">• Enables massive model scaling• More memory efficient than model parallelism	<ul style="list-style-type: none">• High communication overhead• Very difficult to implement, requires specialized libraries

