



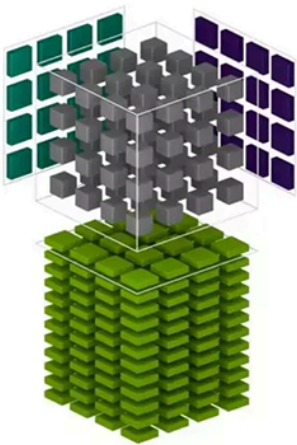
Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT

Tensor Parallelism in Transformers: A Hands-On Guide for Multi-GPU Inference

Posted Dec 5, 2025 • Updated Dec 9, 2025



Tensor Parallelism

$$D = \begin{matrix} \text{FP16 or FP32} \\ \begin{matrix} \begin{matrix} A_{0,0} & A_{0,1} & A_{0,2} \\ A_{1,0} & A_{1,1} & A_{1,2} \\ A_{2,0} & A_{2,1} & A_{2,2} \\ A_{3,0} & A_{3,1} & A_{3,2} \end{matrix} & \begin{matrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{matrix} & \begin{matrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{matrix} \end{matrix} \end{matrix} + \begin{matrix} \text{FP16} \\ \begin{matrix} A_{0,0} & A_{0,1} & A_{0,2} \\ A_{1,0} & A_{1,1} & A_{1,2} \\ A_{2,0} & A_{2,1} & A_{2,2} \\ A_{3,0} & A_{3,1} & A_{3,2} \end{matrix} \end{matrix}$$

By Gian Paolo Santopaolo

15 min read

Contents >

Running a 70B parameter model on a single GPU? Not happening. Even the beefiest H100 with 80GB of VRAM can’t hold Llama-2-70B in full precision. This is where **Tensor Parallelism (TP)** comes in — it splits the model’s weight matrices across multiple GPUs so you can run models that would otherwise be impossible.

This guide is hands-on. We’ll cover the theory just enough to understand what’s happening, then dive straight into code. By the end, you’ll have working scripts for running tensor-parallel inference on **RunPod** and **Lambda Cloud**.

Why Tensor Parallelism? The Memory Wall Problem

Modern LLMs are massive. Here’s a quick reality check:

Model	Parameters	FP16 Memory	FP32 Memory
Llama-3-8B	8B	~16 GB	~32 GB
Llama-3-70B	70B	~140 GB	~280 GB
Llama-3-405B	405B	~810 GB	~1.6 TB

A single A100 (80GB) can barely fit Llama-3-70B in FP16 — and that’s before accounting for KV cache, activations, and batch size overhead. For anything larger, you need to split the model across GPUs.

The Parallelism Zoo



Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT

There are several ways to distribute work across GPUs:



Strategy	What's Split	Memory per GPU	Communication
Data Parallelism	Data batches	Full model on each GPU	Gradient sync after
Pipeline Parallelism	Layers	Subset of layers	Activations between
Tensor Parallelism	Weight matrices	Slice of every layer	All-reduce within e

When to use Tensor Parallelism:

- Model doesn't fit on a single GPU
- You have fast interconnects (NVLink, InfiniBand)
- You want to minimize latency for inference

How Tensor Parallelism Works

The core insight is simple: **matrix multiplications can be parallelized by splitting the matrices.**

Column-Parallel Matrix Multiplication

Suppose you need to compute $Y = XW$ where X is your input and W is a weight matrix. If you split W into column blocks:

$$W = [W_1 \mid W_2 \mid \dots \mid W_n]$$

Then each GPU computes its slice independently:

$$Y_i = X \cdot W_i$$

The outputs are naturally sharded by columns — no communication needed yet.

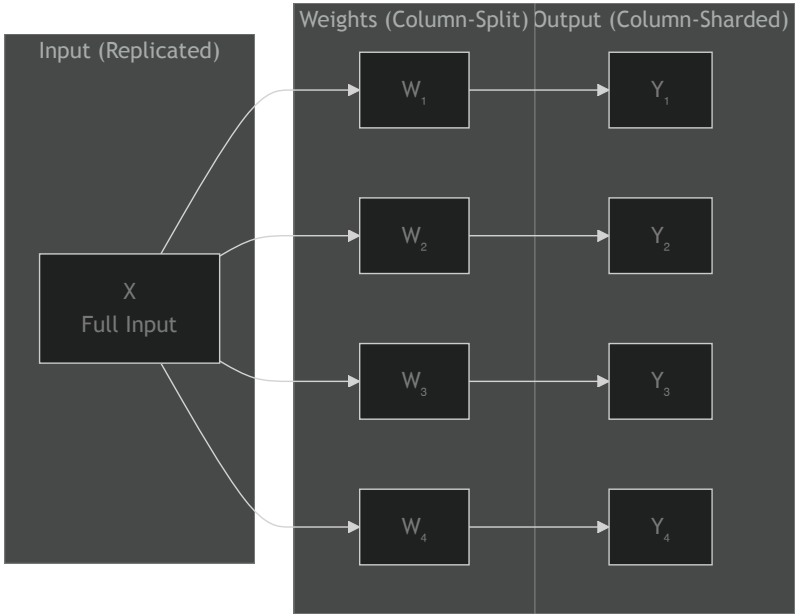




Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT



Row-Parallel Matrix Multiplication

Now suppose you have column-sharded input $X = [X_1 \mid X_2 \mid \dots \mid X_n]$ and you split W into matching row blocks:

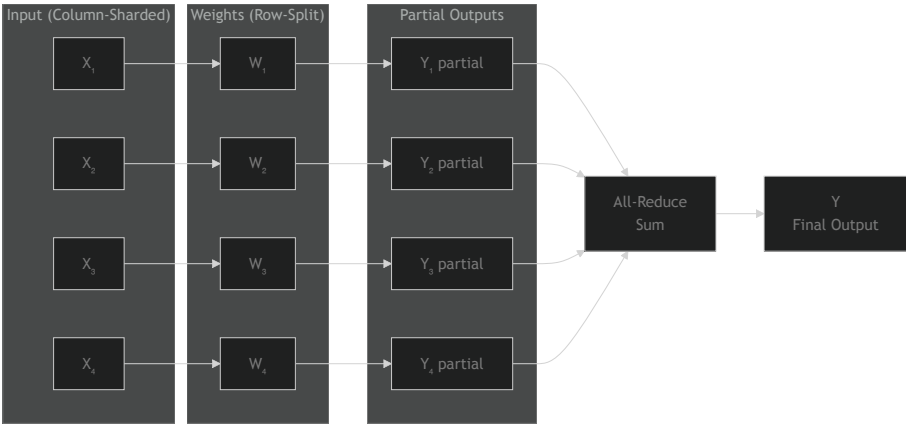
$$W = \begin{bmatrix} W_1 \\ W_2 \\ \vdots \\ W_n \end{bmatrix}$$

Each GPU computes a partial result:

$$Y_i = X_i \cdot W_i$$

Then you **sum across GPUs** (all-reduce) to get the final output:

$$Y = \sum_i Y_i$$



This is the key operation that requires GPU-to-GPU communication.

TP in Transformer Layers

Now let's see how these primitives apply to actual Transformer components.







Attention Layer





Gian Paolo Santopaolo

*Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact*

-  HOME
-  POSTS
-  PROJECTS
-  ABOUT

The attention mechanism has three projection matrices: W_Q , W_K , W_V (queries, keys, values) and an output projection W_O .

Step 1: Split Q, K, V Projections (Column-Parallel)

Each GPU gets a subset of attention heads. If you have 32 heads and 4 GPUs, each GPU handles 8 heads.

```

1 GPU i computes (column slices of weight matrices):
2   Q_i = X × W_Q[all_rows, columns_for_heads_i]
3   K_i = X × W_K[all_rows, columns_for_heads_i]
4   V_i = X × W_V[all_rows, columns_for_heads_i]
```

No communication needed — each GPU works independently.

Step 2: Local Attention Computation

Since attention heads are independent, each GPU computes attention for its heads locally:

```

1 GPU i computes attention locally:
2   attn_i = softmax(Q_i × K_i^T / √d_k) × V_i
```

Still no communication.

Step 3: Output Projection (Row-Parallel)

The output projection W_O is split by rows. Each GPU multiplies its attention output by its slice of W_O , then we all-reduce:

```

1 GPU i computes partial output (row slice of W_O):
2   partial_i = attn_i × W_O[rows_for_gpu_i, all_cols]
3
4 All GPUs synchronize:
5   output = AllReduce(partial_0 + partial_1 + ... + partial_n)
```

One all-reduce per attention layer.

Feed-Forward Network (FFN)

The FFN typically has two linear layers with an activation in between:

$$\text{FFN}(x) = \text{GELU}(xW_1)W_2$$

First Linear (Column-Parallel):

```

1 GPU i computes:
2   hidden_i = GELU(x × W1[all_rows, cols_for_gpu_i])
```

Second Linear (Row-Parallel):

```





1 GPU i computes:
```





Gian Paolo Santopaolo

*Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact*

-  HOME
-  POSTS
-  PROJECTS
-  ABOUT

```
1 GPU i computes:
2   partial_i = hidden_i × W2[rows_for_gpu_i, all_cols]
3
4 All GPUs synchronize:
5   output = AllReduce(sum of all partial_i)
```

One all-reduce per FFN layer.

The Full Picture



Total communication per layer: 2 all-reduce operations.

Constraints

TP comes with a few practical constraints:

1. **TP size ≤ number of attention heads** — you can't split a single head across GPUs
2. **Heads must be divisible by TP size** — each GPU needs an equal share
3. **FFN hidden dimension must be divisible by TP size**

For Llama-3-70B with 64 heads, valid TP sizes are: 1, 2, 4, 8, 16, 32, 64.

Tensor Parallelism with HuggingFace Transformers

The good news: HuggingFace Transformers now has built-in TP support. For supported models, it's a one-liner.

The 3-Line Solution

```
1 # tp_inference.py
2 import torch
3 from transformers import AutoModelForCausalLM, AutoTokenizer
4
5 model = AutoModelForCausalLM.from_pretrained(
6     "meta-llama/Meta-Llama-3-8B-Instruct",
7     torch_dtype=torch.bfloat16,
8     tp_plan="auto" # <-- This enables tensor parallelism
9 )
10
11 tokenizer = AutoTokenizer.from_pretrained("meta-llama/Meta-Llama-3-8B-Ins
12
13 prompt = "Explain tensor parallelism in one paragraph:"
14 inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
15
16 outputs = model.generate(**inputs, max_new_tokens=100)
17 print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

Launching with torchrun

You can't just run `python tp_inference.py`. You need to launch it with `torchrun` to spawn multiple processes:

```
</> Shell
```



Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT

```
1 # Run on 4 GPUs
2 torchrun --nproc-per-node 4 tp_inference.py
```

Each process gets assigned to one GPU, and PyTorch’s distributed runtime handles the communication.

Supported Models

As of late 2025, HuggingFace supports TP for:

- Llama (all versions)
- Mistral
- Mixtral
- Qwen
- Gemma
- And more...

Check the model’s config for `_tp_plan` to see if it’s supported:

```
</> Python
1 from transformers import AutoConfig
2 config = AutoConfig.from_pretrained("meta-llama/Meta-Llama-3-8B-Instruct")
3 print(config._tp_plan) # Shows the default TP plan
```

Partitioning Strategies

Under the hood, HuggingFace uses these strategies:

Strategy	Description
colwise	Column-parallel (for Q, K, V projections)
rowwise	Row-parallel (for output projections)
sequence_parallel	For LayerNorm, Dropout
replicate	Keep full copy on each GPU

You can define a custom `tp_plan` if needed:

```
</> Python
1 tp_plan = {
2     "model.layers.*.self_attn.q_proj": "colwise",
3     "model.layers.*.self_attn.k_proj": "colwise",
4     "model.layers.*.self_attn.v_proj": "colwise",
5     "model.layers.*.self_attn.o_proj": "rowwise",
6     "model.layers.*.mlp.gate_proj": "colwise",
7     "model.layers.*.mlp.up_proj": "colwise",
8     "model.layers.*.mlp.down_proj": "rowwise",
9 }
10
11 model = AutoModelForCausalLM.from_pretrained(
12     "meta-llama/Meta-Llama-3-8B-Instruct",
13     torch_dtype=torch.bfloat16,
14     tp_plan=tp_plan
15 )
```





Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT

Hands-On: Running TP on RunPod

RunPod offers on-demand GPU pods with multi-GPU configurations. Let’s run Llama-3-70B with tensor parallelism.

Step 1: Spin Up a Multi-GPU Pod

- 1. Go to RunPod → Pods → Deploy
- 2. Select a template with PyTorch (e.g., runpod/pytorch:2.1.0-py3.10-cuda11.8.0)
- 3. Choose a multi-GPU configuration:
 - 4× A100 80GB for Llama-3-70B
 - 8× H100 for larger models or faster inference


Critical: Select instances with **NVLink** interconnect (e.g., SXM variants like A100-SXM or H100-SXM), not PCIe. NVLink provides 600-900 GB/s bandwidth between GPUs, while PCIe is limited to ~64 GB/s. Without NVLink, the all-reduce operations in tensor parallelism become a severe bottleneck, negating most of the performance gains.

Configure Deployment

Pod name

selective_silver_bug

Pod Template

 Runpod Pytorch 2.8.0
runpod/pytorch:1.0.2-cu1281-torch280-ubuntu2404

⌕

✎ Edit

Change Template

GPU count

1

2

3

4

5

6

7

8

Instance pricing

On-Demand

Non-Interruptible

\$5.96/hr

Pay as you go, with costs based on actual usage time.

3 Month Savings Plan

Save \$1,641.60

\$5.20/hr

\$11,232.00

Reserve a GPU for three months at a discounted hourly cost.

6 Month Savings Plan

Save \$3,842.96

\$5.08/hr

\$22,184.36

Reserve a GPU for six months at a discounted hourly cost.

1 Year Savings Plan

Save \$9,460.80

\$4.88/hr

\$42,748.80

Reserve a GPU for one year at a discounted hourly cost.

Spot

Interruptible

\$3.80/hr

Pay much less for an interruptible instance.

Looking for multi-year savings plans? [Contact Sales](#)

☐ Encrypt volume

☒ SSH terminal access

☒ Start Jupyter notebook

Selecting a 4xA100 pod on RunPod — look for SXM variants with NVLink

Step 2: Environment Setup

SSH into your pod and set up the environment:





Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT



```
1 # Update and install dependencies
2 pip install --upgrade transformers accelerate torch
3
4 # Verify GPU setup
5 nvidia-smi
6
7 # Check NCCL (the communication backend)
8 python -c "import torch; print(f'CUDA available: {torch.cuda.is_available()})"
```

Expected output:

● ●

</> Plaintext

CUDA available: True
GPU count: 4

Step 3: Create the Inference Script

● ●





</> Python





Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

-  HOME
-  POSTS
-  PROJECTS
-  ABOUT

```

1  # runpod_tp_inference.py
2  import os
3  import torch
4  from transformers import AutoModelForCausalLM, AutoTokenizer
5
6  def main():
7      model_id = "meta-llama/Meta-Llama-3-70B-Instruct"
8
9      # Load model with tensor parallelism
10     model = AutoModelForCausalLM.from_pretrained(
11         model_id,
12         torch_dtype=torch.bfloat16,
13         tp_plan="auto"
14     )
15
16     tokenizer = AutoTokenizer.from_pretrained(model_id)
17     tokenizer.pad_token = tokenizer.eos_token
18
19     # Only rank 0 should print
20     rank = int(os.environ.get("RANK", 0))
21
22     prompts = [
23         "What is tensor parallelism?",
24         "Explain the difference between data and model parallelism.",
25     ]
26
27     for prompt in prompts:
28         inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
29
30         with torch.no_grad():
31             outputs = model.generate(
32                 **inputs,
33                 max_new_tokens=200,
34                 do_sample=True,
35                 temperature=0.7,
36                 top_p=0.9,
37             )
38
39         if rank == 0:
40             response = tokenizer.decode(outputs[0], skip_special_tokens=True)
41             print(f"\n{'='*50}")
42             print(f"Prompt: {prompt}")
43             print(f"Response: {response}")
44             print(f"{'='*50}\n")
45
46 if __name__ == "__main__":
47     main()

```

Step 4: Launch with torchrun

```

1  # Set your HuggingFace token for gated models
2  export HF_TOKEN="your_token_here"
3
4  # Launch on 4 GPUs
5  torchrun --nproc-per-node 4 runpod_tp_inference.py

```

Using vLLM with Tensor Parallelism on RunPod

For production inference, [vLLM](#) is often faster. RunPod has native vLLM support:

```

1  # Shell

```





Gian Paolo Santopaolo

*Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact*

- HOME
- POSTS
- PROJECTS
- ABOUT

```
1 # Install vLLM
2 pip install vllm
3
4 # Run with tensor parallelism
5 python -m vllm.entrypoints.openai.api_server \
6     --model meta-llama/Meta-Llama-3-70B-Instruct \
7     --tensor-parallel-size 4 \
8     --dtype bfloat16 \
9     --port 8000
```

Or use RunPod's serverless vLLM workers which handle TP automatically:

```
</> Python
1 # In your RunPod serverless handler
2 handler_config = {
3     "model_name": "meta-llama/Meta-Llama-3-70B-Instruct",
4     "tensor_parallel_size": 4,
5     "dtype": "bfloat16",
6 }
```

Hands-On: Running TP on Lambda Cloud

[Lambda Cloud](#) offers GPU instances with up to 8× H100s. The setup is similar but with some Lambda-specific details.

Step 1: Launch a Multi-GPU Instance

- Go to Lambda Cloud → Instances → Launch
- Select instance type:
 - gpu_8x_h100_sxm5** (8× H100 80GB) — best for large models
 - gpu_4x_a100_80gb_sxm4** (4× A100 80GB) — good for 70B models

! Critical: Always choose **SXM** variants (e.g., `sxm4`, `sxm5`) over PCIe. The “SXM” designation indicates GPUs connected via **NVLink** with 600-900 GB/s inter-GPU bandwidth. PCIe-based instances share bandwidth through the CPU's PCIe lanes (~64 GB/s), creating a communication bottleneck that cripples tensor parallelism performance.





Gian Paolo Santopaolo

*Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact*

- HOME
- POSTS
- PROJECTS
- ABOUT

LAUNCH INSTANCE

Instance type

Region

Base image

Filesystem

Security

Select instance type

1x GH200 (96 GB)

ARM64 + H100

\$1.49 / hr

(64 vCPUs, 432 GiB RAM, 4 TiB SSD)

8x B200 (180 GB SXM6)

New

\$39.92 / hr

(208 vCPUs, 2900 GiB RAM, 22 TiB SSD)

8x H100 (80 GB SXM5)

\$23.92 / hr

(208 vCPUs, 1800 GiB RAM, 22 TiB SSD)

4x H100 (80 GB SXM5)

\$12.36 / hr

(104 vCPUs, 900 GiB RAM, 11 TiB SSD)

2x H100 (80 GB SXM5)

\$6.38 / hr

(52 vCPUs, 450 GiB RAM, 5.5 TiB SSD)

1x H100 (80 GB SXM5)

\$3.29 / hr

(26 vCPUs, 225 GiB RAM, 2.8 TiB SSD)

1x H100 (80 GB PCIe)

\$2.49 / hr

(26 vCPUs, 200 GiB RAM, 1 TiB SSD)

8x A100 (80 GB SXM4)

\$14.32 / hr

(240 vCPUs, 1800 GiB RAM, 20 TiB SSD)

Selecting a multi-GPU instance on Lambda Cloud — SXM variants have NVLink

Step 2: SSH and Setup

```
1 # SSH into your instance
2 ssh ubuntu@<your-instance-ip>
3
4 # Lambda instances come with PyTorch pre-installed
5 # Just update transformers
6 pip install --upgrade transformers accelerate
7
8 # Verify setup
9 python -c "import torch; print(f'GPUs: {torch.cuda.device_count()}')"
```

Step 3: Create the Inference Script





```
</> Python
```





Gian Paolo Santopaolo

*Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact*

-  HOME
-  POSTS
-  PROJECTS
-  ABOUT



```

1  # lambda_tp_inference.py
2  import os
3  import time
4  import torch
5  from transformers import AutoModelForCausalLM, AutoTokenizer
6
7  def main():
8      model_id = "meta-llama/Meta-Llama-3-70B-Instruct"
9
10     rank = int(os.environ.get("RANK", 0))
11     world_size = int(os.environ.get("WORLD_SIZE", 1))
12
13     if rank == 0:
14         print(f"Loading {model_id} with TP across {world_size} GPUs...")
15         start_time = time.time()
16
17     model = AutoModelForCausalLM.from_pretrained(
18         model_id,
19         torch_dtype=torch.bfloat16,
20         tp_plan="auto"
21     )
22
23     if rank == 0:
24         load_time = time.time() - start_time
25         print(f"Model loaded in {load_time:.2f}s")
26
27     tokenizer = AutoTokenizer.from_pretrained(model_id)
28     tokenizer.pad_token = tokenizer.eos_token
29
30     # Benchmark inference
31     prompt = "Write a short poem about distributed computing:"
32     inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
33
34     # Warmup
35     with torch.no_grad():
36         _ = model.generate(**inputs, max_new_tokens=10)
37
38     # Timed generation
39     torch.cuda.synchronize()
40     start = time.time()
41
42     with torch.no_grad():
43         outputs = model.generate(
44             **inputs,
45             max_new_tokens=100,
46             do_sample=False, # Greedy for reproducibility
47         )
48
49     torch.cuda.synchronize()
50     gen_time = time.time() - start
51
52     if rank == 0:
53         response = tokenizer.decode(outputs[0], skip_special_tokens=True)
54         tokens_generated = outputs.shape[1] - inputs["input_ids"].shape[1]
55         tokens_per_sec = tokens_generated / gen_time
56
57         print(f"\nPrompt: {prompt}")
58         print(f"Response: {response}")
59         print(f"\n--- Performance ---")
60         print(f"Tokens generated: {tokens_generated}")
61         print(f"Time: {gen_time:.2f}s")
62         print(f"Throughput: {tokens_per_sec:.1f} tokens/sec")
63
64     if __name__ == "__main__":
65         main()

```





Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT

Step 4: Launch with torchrun

```
</> Shell
1 # For a single node with 4 GPUs
2 torchrun --nproc-per-node 4 lambda_tp_inference.py
3
4 # For 8 GPUs
5 torchrun --nproc-per-node 8 lambda_tp_inference.py
```

Multi-Node Setup on Lambda Cloud

If you need more than 8 GPUs, you can run across multiple nodes. Lambda instances support this via `torchrun` :

```
</> Shell
1 # On Node 0 (master)
2 torchrun \
3     --nproc-per-node 8 \
4     --nnodes 2 \
5     --node-rank 0 \
6     --master-addr <master-ip> \
7     --master-port 29500 \
8     lambda_tp_inference.py
9
10 # On Node 1 (worker)
11 torchrun \
12     --nproc-per-node 8 \
13     --nnodes 2 \
14     --node-rank 1 \
15     --master-addr <master-ip> \
16     --master-port 29500 \
17     lambda_tp_inference.py
```

This gives you 16 GPUs with tensor parallelism across nodes.

Warning: Cross-node TP requires high-bandwidth interconnects (InfiniBand). Without it, communication overhead can kill performance.

Performance Benchmarks

Here’s what you can expect with tensor parallelism on different configurations:

Llama-3-70B Inference Throughput

Configuration	TP Size	Tokens/sec	Memory/GPU
1× H100 80GB	1	OOM	—
2× H100 80GB	2	~45	~38 GB
4× H100 80GB	4	~85	~20 GB
8× H100 80GB	8	~140	~12 GB

Key Observations

- 1. **Memory scales linearly** — 4× GPUs = ~4× less memory per GPU





Gian Paolo Santopaolo

*Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact*

- HOME
- POSTS
- PROJECTS
- ABOUT

2. **Throughput scales sub-linearly** — communication overhead increases with TP size
3. **Sweet spot is often 4-8 GPUs** — beyond that, communication dominates

What TP Doesn't Solve

Tensor parallelism is powerful, but it has limitations:

1. Scalability is Capped by Attention Heads

If your model has 64 attention heads, TP size can't exceed 64. In practice, you want TP size much smaller than head count to maintain efficiency.

2. Communication Overhead Across Nodes

TP requires frequent all-reduce operations (2 per layer). Within a node with NVLink (900 GB/s), this is fast. Across nodes with InfiniBand (~400 GB/s) or worse, Ethernet (~100 Gbps), it becomes a bottleneck.

Rule of thumb: Keep TP within a single node. Use Pipeline Parallelism (PP) across nodes.

3. Doesn't Help with Activation Memory

TP reduces weight memory but not activation memory. For very long sequences, you may still need gradient checkpointing or other techniques.

When to Combine with Pipeline Parallelism

For truly massive models (400B+), combine TP and PP:

```
Node 0: Layers 0-19 (TP=8 within node)
Node 1: Layers 20-39 (TP=8 within node)
Node 2: Layers 40-59 (TP=8 within node)
Node 3: Layers 60-79 (TP=8 within node)
```

This gives you 32 GPUs total: 8-way TP × 4-way PP.

Practical Takeaways

Decision Tree: Which Parallelism Strategy?

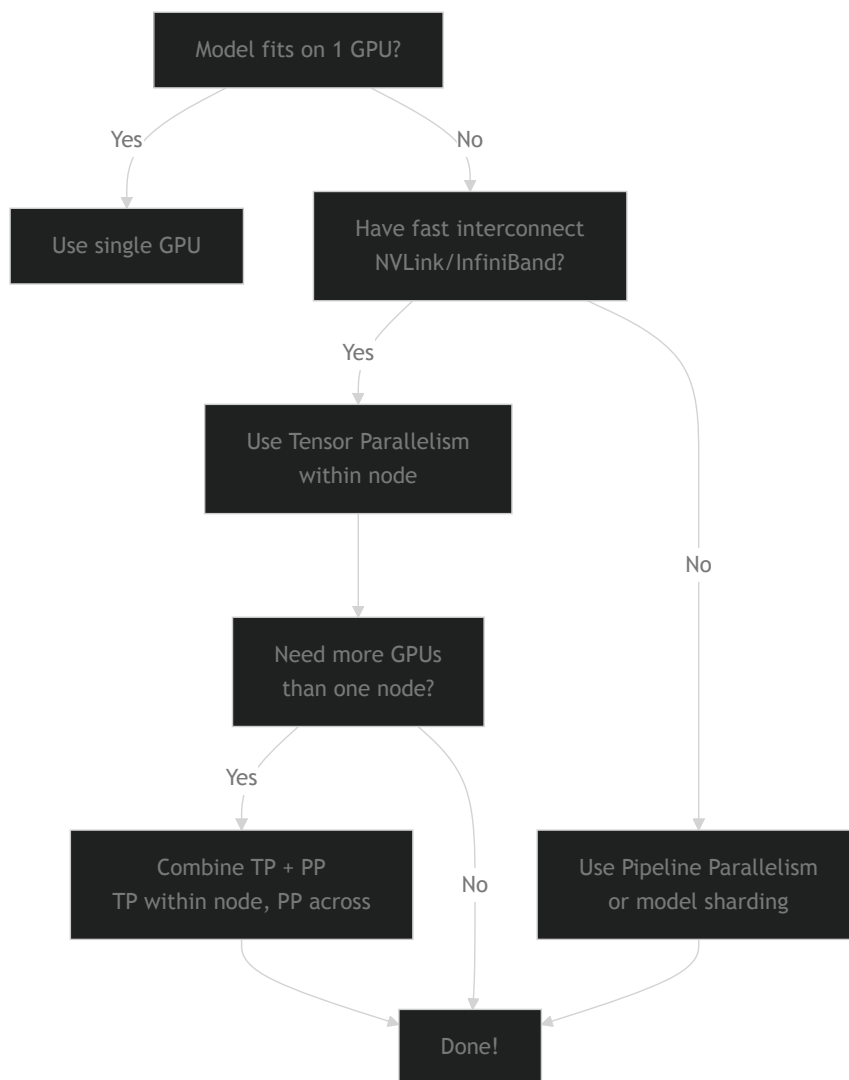




Gian Paolo Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT



Quick Reference Commands

```
1 # HuggingFace Transformers with TP
2 torchrun --nproc-per-node 4 inference.py
3
4 # vLLM with TP
5 python -m vllm.entrypoints.openai.api_server \
6     --model meta-llama/Meta-Llama-3-70B-Instruct \
7     --tensor-parallel-size 4
8
9 # Check GPU topology (important for TP performance)
10 nvidia-smi topo -m
11
12 # Monitor GPU usage during inference
13 watch -n 0.5 nvidia-smi
```

Key Constraints Checklist

Before deploying with TP, verify:

- ☐ TP size \leq number of attention heads
- ☐ Attention heads divisible by TP size
- ☐ FFN hidden dim divisible by TP size
- ☐ All GPUs have NVLink or fast interconnect
- ☐ Using `torchrun` or equivalent launcher





Gian Paolo Santopaolo

*Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact*

- HOME
- POSTS
- PROJECTS
- ABOUT

Wrapping Up

Tensor parallelism is the go-to technique for running models that don't fit on a single GPU. The key ideas:

1. **Split weight matrices** across GPUs (column-wise for projections, row-wise for outputs)
2. **All-reduce** to aggregate partial results (2× per transformer layer)
3. **Keep TP within a node** for best performance
4. Use `tp_plan="auto"` in HuggingFace for the easy path

For production inference, consider vLLM which has highly optimized TP implementations. For training, look into FSDP (Fully Sharded Data Parallel) which combines aspects of TP and data parallelism.

References

- [HuggingFace Distributed Inference Documentation](#)
- [Megatron-LM: Training Multi-Billion Parameter Language Models](#)
- [RunPod Multi-GPU Training Guide](#)
- [Lambda Labs Multi-Node PyTorch Guide](#)
- [vLLM Documentation](#)
- [PyTorch Distributed Overview](#)

[Tensor Parallelism](#), [Tensor](#), [Transformers](#), [Multi-GPU](#), [PyTorch](#), [RunPod](#), [Lambda Cloud](#)

This post is licensed under CC BY 4.0 by the author.

Share:

Further Reading

Dec 2, 2025

What Is a Tensor? A Practical Guide for AI Engineers

When dealing with deep learning, we quickly encounter the word tensor. But what exactly i...

Dec 3, 2025

Why GPUs Love Tensors: Understanding Tensor Cores and...

Understanding why modern AI is so fast requires understanding the hardware that ...

Dec 18, 2025

PyTorch Training Loop: A Production-Ready Template for ...

You can train almost anything in PyTorch with seven lines of code. The problem? Those seve...



OLDER

NEWER

Why GPUs Love Tensors: Understanding
Tensor Cores and AI Acceleration

Choosing the Right Loss Function for your
ML Problem



Gian Paolo
Santopaolo

Hands-On ML Engineering
Meets Strategic Vision—
Bridging Code and Cloud to
Drive Enterprise AI for Real-
World Impact

- HOME
- POSTS
- PROJECTS
- ABOUT

© 2025 Gian Paolo Santopaolo. Some rights reserved.

Opinions expressed in this website are my own.

