
DOCUMENT

Tensor Parallelism for Multi-GPU Transformer Inference

Source: temp_input_17d68ee69aa24e07b8e372ff8a3aa9da.md

Contents

Estimated reading time: 16 min

- Introduction
 - 1. Introduction
 - 2. The Need for Tensor Parallelism
 - 2.1 The Memory Wall Problem with Large Models
 - 2.2 Comparing Parallelism Strategies
 - 3. How Tensor Parallelism Works
 - 3.1 Column-Parallel Matrix Multiplication
 - 3.2 Row-Parallel Matrix Multiplication
 - 4. Applying TP in Transformer Layers
 - 4.1 Tensor Parallelism in the Attention Layer
 - 4.2 Tensor Parallelism in the Feed-Forward Network
 - 4.3 Communication and Practical Constraints
 - 5. Implementing TP with HuggingFace Transformers
 - 5.1 A Simple Code Implementation
 - 5.2 Launching with torchrun
 - 5.3 Supported Models and Partitioning Strategies
 - 6. Hands-On Guide: Tensor Parallelism on RunPod
 - 6.1 Setting Up a Multi-GPU Pod
 - 6.2 Environment Setup and Inference Script
- **Update and install dependencies**
- **Verify GPU setup**
- **Check NCCL (the communication backend)**
- **runpod_tp_inference.py**
 - 6.3 Launching the Inference Job

- **Launch on 4 GPUs**

- 6.4 Using vLLM with Tensor Parallelism

- **Install vLLM**

- **Run with tensor parallelism**

- **In your RunPod serverless handler**

- 7. Hands-On Guide: Tensor Parallelism on Lambda Cloud

- 7.1 Launching a Multi-GPU Instance

- 7.2 Creating and Launching the Inference Script

- **lambda_tp_inference.py**

- **For a single node with 4 GPUs**

- **For 8 GPUs**

- 7.3 Multi-Node Setup for Large-Scale Inference

- 8. Performance Benchmarks and Observations

- 8.1 Llama-3-70B Inference Throughput

- 8.2 Key Performance Scaling Observations

- 9. Limitations and Advanced Strategies

- 9.1 What Tensor Parallelism Does Not Solve

- 9.2 Combining with Pipeline Parallelism

- 10. Key Takeaways

- Key Takeaways

Introduction

Running large language models, such as a 70-billion parameter model, on a single GPU is often infeasible due to memory limitations. Even a high-end H100 GPU with 80GB of VRAM cannot accommodate a model like Llama-2-70B in full precision. Tensor Parallelism (TP) addresses this challenge by splitting a model's weight matrices across multiple GPUs, enabling the execution of models that would otherwise be too large. This guide provides a hands-on approach to understanding and implementing tensor-parallel inference, covering the underlying theory before diving into practical code examples for running inference on cloud platforms like RunPod and Lambda Cloud.

1. Introduction

This guide focuses on the practical application of Tensor Parallelism for large model inference. We will cover the necessary theoretical concepts to understand its mechanics and then proceed with implementation details. By the conclusion, you will have functional scripts for executing tensor-parallel inference on multi-GPU environments.

2. The Need for Tensor Parallelism

The primary driver for using parallelism strategies is the immense size of modern large language models, which creates a significant memory bottleneck.

2.1 The Memory Wall Problem with Large Models

Modern LLMs require substantial GPU memory, often exceeding the capacity of a single device. The memory footprint grows directly with the number of parameters and the precision used (e.g., FP16 or FP32).

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

As the table illustrates, a single A100 GPU with 80GB of VRAM can barely hold the Llama-3-70B model in FP16, and this does not account for the additional memory required for the KV cache, activations, and batch size overhead. To run larger models, it is essential to distribute the model across multiple GPUs.

2.2 Comparing Parallelism Strategies

Several strategies exist for distributing computational workloads across multiple GPUs, each with distinct characteristics regarding what is split, memory usage, and communication patterns.

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

Tensor Parallelism is the preferred strategy when the model is too large to fit on a single GPU, when fast interconnects like NVLink or InfiniBand are available, and when the primary goal is to minimize inference latency.

3. How Tensor Parallelism Works

The fundamental concept behind Tensor Parallelism is that matrix multiplications, which are core operations in Transformers, can be parallelized by splitting the matrices themselves across different processing units.

3.1 Column-Parallel Matrix Multiplication

Consider the computation $\mathbf{y} = \mathbf{x}\mathbf{w}$, where \mathbf{x} is the input and \mathbf{w} is a weight matrix. In a column-parallel approach, the weight matrix \mathbf{w} is split into column blocks. Each GPU computes its portion of the output by multiplying the full input \mathbf{x} with its assigned column slice of \mathbf{w} . The resulting outputs \mathbf{y}_i are also sharded by columns. This step can be performed independently on each GPU without any inter-GPU communication.

3.2 Row-Parallel Matrix Multiplication

Now, if the input \mathbf{x} is already sharded by columns, as in $\mathbf{x} = [\mathbf{x}_1 \mid \mathbf{x}_2 \mid \dots \mid \mathbf{x}_n]$, and the weight matrix \mathbf{w} is split into corresponding row blocks, each GPU can compute a partial result. To obtain the final output \mathbf{y} , the partial results from all GPUs must be summed together. This is achieved through an **all-reduce** communication operation, which is a key step requiring GPU-to-GPU data transfer.

4. Applying TP in Transformer Layers

Tensor Parallelism can be systematically applied to the main components of a Transformer layer: the attention mechanism and the feed-forward network.

4.1 Tensor Parallelism in the Attention Layer

The attention mechanism involves three projection matrices for queries (w_Q), keys (w_K), and values (w_V), along with an output projection matrix (w_O). TP is applied in three steps.

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

The w_Q , w_K , and w_V weight matrices are split by columns. This is equivalent to distributing the attention heads across the available GPUs. For instance, with 32 attention heads and 4 GPUs, each GPU would handle 8 heads. Each GPU computes its slice independently.

```

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]
```

This stage requires no communication.

Step 2: Local Attention Computation

Since attention heads operate independently, each GPU can compute the attention mechanism for its assigned heads locally without communicating with other GPUs.

```

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

```

Step 3: Output Projection (Row-Parallel)

The output projection matrix W_O is split by rows. Each GPU multiplies its local attention output by its slice of W_O to get a partial result. An all-reduce operation is then performed to sum these partial results into the final output.

```

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)

```

This process results in one all-reduce operation per attention layer.

4.2 Tensor Parallelism in the Feed-Forward Network

The feed-forward network (FFN) in a Transformer typically consists of two linear layers. TP is applied to these layers as follows:

First Linear Layer (Column-Parallel)

The first linear layer's weight matrix, W_1 , is split by columns. Each GPU computes a partial hidden state independently.

```

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

```

Second Linear Layer (Row-Parallel)

The second linear layer's weight matrix, W_2 , is split by rows. Each GPU computes a partial output using its hidden state slice. These partial outputs are then summed across all GPUs using an all-reduce operation.

```

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)

```

This adds a second all-reduce operation, bringing the total to two all-reduce operations per complete Transformer layer.

4.3 Communication and Practical Constraints

In total, a standard Transformer layer requires **two all-reduce operations** when using Tensor Parallelism. However, there are practical constraints to consider:

1. The TP size (number of GPUs) cannot exceed the number of attention heads.
2. The number of attention heads must be divisible by the TP size to ensure an equal workload on each GPU.
3. The FFN hidden dimension must also be divisible by the TP size.

For example, the Llama-3-70B model has 64 attention heads. Therefore, valid TP sizes are 1, 2, 4, 8, 16, 32, and 64.

5. Implementing TP with HuggingFace Transformers

The HuggingFace `transformers` library offers built-in support for Tensor Parallelism, simplifying its implementation for many models to just a single parameter.

5.1 A Simple Code Implementation

For supported models, enabling Tensor Parallelism can be done by setting the `tp_plan` argument to "auto" when loading the model. This instructs the library to automatically partition the model's weights across the available GPUs.

```

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-Instruct")
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 ■
18 ■ print(tokenizer.decode(outputs[0], skip_special_tokens=True))

```

5.2 Launching with torchrun

A script using Tensor Parallelism cannot be run directly with `python`. Instead, it must be launched using `torchrun`, which spawns multiple processes and manages the distributed environment. Each process is assigned to a separate GPU.

To run the script on 4 GPUs:

```
1 ■ torchrun --nproc-per-node 4 tp_inference.py
```

5.3 Supported Models and Partitioning Strategies

As of late 2025, HuggingFace provides native TP support for a growing list of models, including Llama, Mistral, Mixtral, Qwen, and Gemma. You can check if a model is supported by inspecting its configuration for a `_tp_plan` attribute.

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

Internally, HuggingFace uses several partitioning strategies:

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

For advanced use cases, you can define a custom `tp_plan` dictionary to specify how each layer should be partitioned.

```

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

```

6. Hands-On Guide: Tensor Parallelism on RunPod

RunPod provides on-demand multi-GPU instances suitable for running large models with Tensor Parallelism. This section outlines the steps to run Llama-3-70B on a RunPod pod.

6.1 Setting Up a Multi-GPU Pod

First, deploy a new pod on RunPod:

1. Navigate to **RunPod** → **Pods** → **Deploy**.
2. Select a suitable PyTorch template, such as
runpod/pytorch:2.1.0-py3.10-cuda11.8.0.
3. Choose a multi-GPU configuration, like 4x A100 80GB for Llama-3-70B or 8x H100 for larger models.

Critical: It is essential to select instances with **NVLink** interconnects (e.g., SXM

Tensor variants like A100-SXM or H100-SXM) rather than PCIe. NVLink offers 600-900 GB/s

Tensor Parallelism for Large Language Models | Page 5, 2026

of inter-GPU bandwidth, whereas PCIe is limited to approximately 64 GB/s. The high communication cost of all-reduce operations makes NVLink crucial for TP performance; without it, the communication overhead becomes a severe bottleneck and negates performance gains.

6.2 Environment Setup and Inference Script

Once connected to the pod via SSH, set up the environment by installing the necessary libraries.

```
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()}\";GPU c
```

The expected output should confirm CUDA availability and list the correct number of GPUs. Next, create the inference script:

```

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        with torch.no_grad():
30            outputs = model.generate(
31                **inputs,
32                max_new_tokens=200,
33                do_sample=True,
34                temperature=0.7,
35                top_p=0.9,
36            )
37
38        if rank == 0:
39            response = tokenizer.decode(outputs[0], skip_special_tokens=True)
40            print(f"\n{'='*50}")
41            print(f"Prompt: {prompt}")
42            print(f"Response: {response}")
43            print(f"{'='*50}\n")
44
45    if __name__ == "__main__":
46        main()

```

6.3 Launching the Inference Job

To run the script, first set your HuggingFace token as an environment variable, then use `torchrun` to launch the job on all available GPUs.

```

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

```

6.4 Using vLLM with Tensor Parallelism

For production-level inference, vLLM is often a more performant option. To use vLLM with tensor parallelism, you can run its OpenAI-compatible API server:

```

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

```

Alternatively, RunPod's serverless platform can automatically handle vLLM and TP configuration through a handler.

```

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 └ }

```

7. Hands-On Guide: Tensor Parallelism on Lambda Cloud

Lambda Cloud provides high-performance GPU instances, including configurations with up to 8x H100 GPUs, making it another excellent choice for large model inference.

7.1 Launching a Multi-GPU Instance

To start, launch a multi-GPU instance from the Lambda Cloud dashboard:

1. Navigate to **Lambda Cloud** → **Instances** → **Launch**.
2. Select an appropriate instance type, such as `gpu_8x_h100_sxm5` or `gpu_4x_a100_80gb_sxm4`.

Critical: As with RunPod, always choose **SXM** variants (e.g., `sxm4`, `sxm5`) over PCIe. The SXM designation confirms that the GPUs are interconnected with high-bandwidth NVLink (600-900 GB/s), which is essential for efficient tensor parallelism. PCIe-based instances are limited to the CPU's PCIe lane bandwidth (~64 GB/s), creating a bottleneck that severely degrades performance.

7.2 Creating and Launching the Inference Script

After SSHing into the instance and completing the environment setup, create the following script, which includes simple benchmarking to measure performance.

```

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     rank = int(os.environ.get("RANK", 0))
10    world_size = int(os.environ.get("WORLD_SIZE", 1))
11
12    if rank == 0:
13        print(f"Loading {model_id} with TP across {world_size} GPUs...")
14
15    start_time = time.time()
16    model = AutoModelForCausalLM.from_pretrained(
17        model_id,
18        torch_dtype=torch.bfloat16,
19        tp_plan="auto"
20    )
21    if rank == 0:
22        load_time = time.time() - start_time
23        print(f"Model loaded in {load_time:.2f}s")
24
25    tokenizer = AutoTokenizer.from_pretrained(model_id)
26    tokenizer.pad_token = tokenizer.eos_token
27
28    # Benchmark inference
29    prompt = "Write a short poem about distributed computing"
30    inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
31
32    # Warmup
33    with torch.no_grad():
34        _ = model.generate(**inputs, max_new_tokens=10)
35
36    # Timed generation
37    torch.cuda.synchronize()
38    start = time.time()
39    with torch.no_grad():
40        outputs = model.generate(
41            **inputs,
42            max_new_tokens=100,
43            do_sample=False, # Greedy for reproducibility
44        )
45    torch.cuda.synchronize()
46    gen_time = time.time() - start

```

```

47 ■
48 ■     if rank == 0:
49 ■         response = tokenizer.decode(outputs[0], skip_special_tokens=True)
50 ■         tokens_generated = outputs.shape[1] - inputs['input_ids'].shape[1]
51 ■         tokens_per_sec = tokens_generated / gen_time
52 ■
53 ■         print(f'\"Prompt: {prompt}\"')
54 ■         print(f'\"Response: {response}\"')
55 ■         print(f'\"--- Performance ---\"')
56 ■         print(f'\"Tokens generated: {tokens_generated}\"')
57 ■         print(f'\"Time: {gen_time:.2f}s\"')
58 ■         print(f'\"Throughput: {tokens_per_sec:.1f} tokens/sec\"')
59 ■
60 ■ if __name__ == "__main__":
61 ■     main()

```

Launch the script using `torchrun`, specifying the number of GPUs on the node.

```

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

```

7.3 Multi-Node Setup for Large-Scale Inference

For models that require more than 8 GPUs, you can extend Tensor Parallelism across multiple nodes using `torchrun`'s multi-node capabilities. This requires a high-bandwidth interconnect like InfiniBand between nodes to be effective.

On the master node (Node 0):

```

1 ■ torchrun \
2 ■     --nproc-per-node 8 \
3 ■     --nnodes 2 \
4 ■     --node-rank 0 \
5 ■     --master-addr <master-ip> \
6 ■     --master-port 29500 \
7 ■     lambda_tp_inference.py

```

On the worker node (Node 1):

```

1 └─ torchrun \
2   └─ --nproc-per-node 8 \
3   └─ --nnodes 2 \
4   └─ --node-rank 1 \
5   └─ --master-addr &lt;master-ip&gt; \
6   └─ --master-port 29500 \
7   └─ lambda_tp_inference.py

```

This configuration creates a 16-GPU setup with tensor parallelism. However, be aware that cross-node communication overhead can significantly impact performance if the interconnect is not sufficiently fast.

8. Performance Benchmarks and Observations

Performance with Tensor Parallelism scales with the number of GPUs, but not perfectly due to communication overhead.

8.1 Llama-3-70B Inference Throughput

The following table shows approximate inference throughput for Llama-3-70B on different H100 GPU configurations.

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

8.2 Key Performance Scaling Observations

Based on these benchmarks, several key patterns emerge:

- 1. Memory scales linearly:** Using four times the number of GPUs reduces the memory requirement per GPU by approximately four times.
- 2. Throughput scales sub-linearly:** Performance gains diminish as the TP size increases because the communication overhead from `all-reduce` operations becomes more significant.
- 3. There is a sweet spot:** For many models, the optimal balance between performance and cost is often found with a TP size of 4 to 8 GPUs. Beyond this, communication can start to dominate computation time.

9. Limitations and Advanced Strategies

While powerful, Tensor Parallelism has inherent limitations and is often combined with other strategies for very large-scale systems.

9.1 What Tensor Parallelism Does Not Solve

Tensor Parallelism has three main limitations:

- 1. Scalability is Capped by Attention Heads:** The TP size cannot be larger than the model's number of attention heads. For a model with 64 heads, you cannot use more than 64 GPUs for TP.
- 2. Communication Overhead Across Nodes:** TP requires frequent communication. This is fast within a node using NVLink (900 GB/s), but becomes a bottleneck across nodes connected by InfiniBand (~400 GB/s) or Ethernet (~100 Gbps).

3. Activation Memory: TP primarily reduces the memory required for model weights, but it does not reduce the memory needed for activations. For very long sequences, activation memory can still become a bottleneck.

9.2 Combining with Pipeline Parallelism

For extremely large models (400B+ parameters), a hybrid approach combining Tensor Parallelism (TP) and Pipeline Parallelism (PP) is common. The general rule of thumb is to use TP within a single node to leverage fast NVLink interconnects and use PP across different nodes.

For example, to run a model on 32 GPUs, you could use an 8-way TP within each of four nodes, and a 4-way PP across the nodes:

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

This hybrid strategy effectively scales to a large number of GPUs by minimizing slow cross-node communication.

10. Key Takeaways

Tensor parallelism is an essential technique for running large language models that do not fit on a single GPU. The core principles involve splitting the model's weight matrices across multiple GPUs—using a column-wise split for projection layers and a row-wise split for output layers—and using all-reduce operations to aggregate partial results. A standard Transformer layer requires two such communication steps. For optimal performance, it is best to confine Tensor Parallelism within a single, NVLink-connected node. The HuggingFace `transformers` library simplifies

implementation with its `tp_plan="auto"` feature. For production environments, frameworks like vLLM provide highly optimized TP implementations, while for training, FSDP (Fully Sharded Data Parallel) offers an alternative that combines elements of both tensor and data parallelism.

Key Takeaways

Tensor Parallelism is the go-to technique for running models that don't fit on a single GPU. The key ideas are splitting weight matrices across GPUs (column-wise for projections, row-wise for outputs) and using all-reduce to aggregate partial results, with two such operations per transformer layer. For best performance, TP should be kept within a single node with fast interconnects like NVLink. For ease of use, HuggingFace's `tp_plan="auto"` is a straightforward solution, while specialized frameworks like vLLM are recommended for production inference.