

JAXMg Basics

Introduction to sharding and multi-GPU linear solvers

1 Introduction

JAXMg wraps NVIDIA's **cuSolverMg** (multi-GPU solver library) for JAX, enabling distributed linear algebra operations across multiple GPUs. This tutorial covers the essential concepts for using JAXMg effectively.

Function	Solves	Complexity
potrs	$Ax = b$ via Cholesky (A symmetric positive definite)	$O(N^3/3)$
potri	A^{-1} via Cholesky	$O(N^3)$
syevd	Eigendecomposition of symmetric matrix	$O(N^3)$

Table 1: JAXMg distributed linear algebra functions

2 Core Concepts: JAX Sharding

2.1 The Mesh

A **mesh** is a logical n -dimensional grid of devices with named axes. Think of it as organizing your GPUs into a coordinate system that JAX uses to distribute data and computation.

```
1 import jax
2 from jax.sharding import Mesh, NamedSharding, PartitionSpec as P
3
4 # 1D Mesh: 4 GPUs in a line, axis named "x"
5 devices = jax.devices("gpu") # [gpu:0, gpu:1, gpu:2, gpu:3]
6 mesh_1d = jax.make_mesh((4,), ("x",))
7
8 # 2D Mesh: 4 GPUs as 2x2 grid, axes "S" and "T"
9 import numpy as np
10 mesh_2d = Mesh(np.array(devices).reshape(2, 2), ("S", "T"))
```

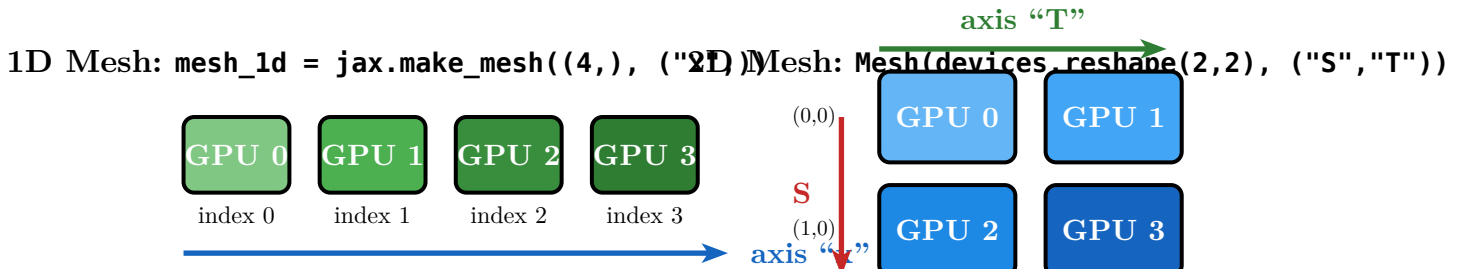


Figure 1: Device mesh configurations. Left: 1D mesh with 4 GPUs along axis “x”. Right: 2D mesh with 4 GPUs arranged in a 2x2 grid with axes “S” (rows) and “T” (columns).

2.2 PartitionSpec: How to Shard

`PartitionSpec` (aliased as `P`) tells JAX which mesh axis to shard each array dimension along. `None` means that dimension is replicated (not sharded).

PartitionSpec	Meaning (for 2D array of shape (N, M))
<code>P("x", None)</code>	Shard dim 0 (rows) across “x”, replicate dim 1 (columns)
<code>P(None, "x")</code>	Replicate dim 0 (rows), shard dim 1 (columns) across “x”
<code>P("S", "T")</code>	Shard rows across “S”, columns across “T”
<code>P(None, None)</code>	Fully replicate — each device has a complete copy
<code>P(("S", "T"), None)</code>	Flatten $S \times T$ mesh into 1D, shard rows across all devices

Table 2: `PartitionSpec` examples for a 2D array

2.3 Placing Data on Devices

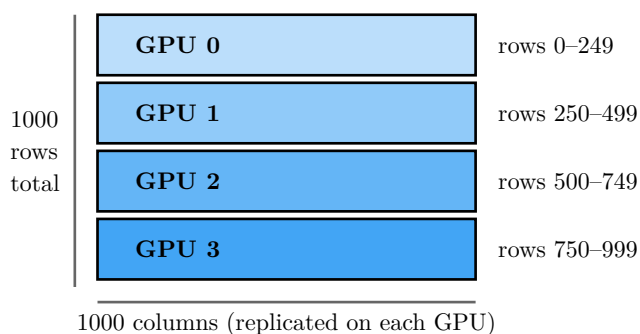
```

1 import jax.numpy as jnp
2
3 N = 1000
4 A = jnp.eye(N, dtype=jnp.float64) # Matrix: 1000 × 1000
5 b = jnp.ones((N, 1), dtype=jnp.float64) # Vector: 1000 × 1
6
7 # Create mesh and shard data
8 mesh = jax.make_mesh((4,), ("x",))
9 A_sharded = jax.device_put(A, NamedSharding(mesh, P("x", None))) # Row-sharded
10 b_replicated = jax.device_put(b, NamedSharding(mesh, P(None, None))) # Replicated

```

Matrix A (1000×1000)

Sharding: `P("x", None)`



Vector b (1000×1)

Sharding: `P(None, None)`

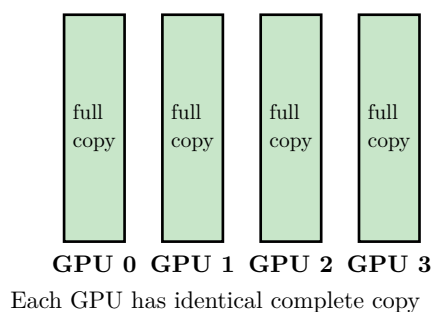


Figure 2: Data distribution across 4 GPUs. Left: Matrix A is row-sharded — each GPU owns 250 rows but all 1000 columns. Right: Vector b is replicated — each GPU has a complete copy.

3 Using `potrs`: Distributed Cholesky Solve

3.1 Basic Usage

```

1 from jaxmg import potrs
2
3 # Problem: Solve Ax = b where A is symmetric positive definite

```

```

4 N = 1000
5 A = jnp.eye(N, dtype=jnp.float64) * 2 # Diagonal matrix
6 b = jnp.ones((N, 1), dtype=jnp.float64)
7
8 # Setup mesh and sharding
9 mesh = jax.make_mesh((4,), ("x",))
10 A = jax.device_put(A, NamedSharding(mesh, P("x", None)))
11 b = jax.device_put(b, NamedSharding(mesh, P(None, None)))
12
13 # Solve!
14 x = potrs(A, b, T_A=128, mesh=mesh, in_specs=(P("x", None), P(None, None)))
15 # x ≈ [0.5, 0.5, 0.5, ...]

```

3.2 Key Parameters

Parameter	Description
T_A	Tile size — block width for distributed algorithm. Larger = faster but more memory. Recommended: 128–2048. Auto-pads if shard size not divisible.
mesh	JAX Mesh object defining device topology
in_specs	Tuple of PartitionSpecs: (P("axis", None), P(None, None)) — matrix row-sharded, RHS replicated
return_status	If True, returns (solution, status_code) tuple
pad	If True (default), auto-pad shards to be divisible by T_A

Table 3: `potrs` function parameters

3.3 Internal Workflow

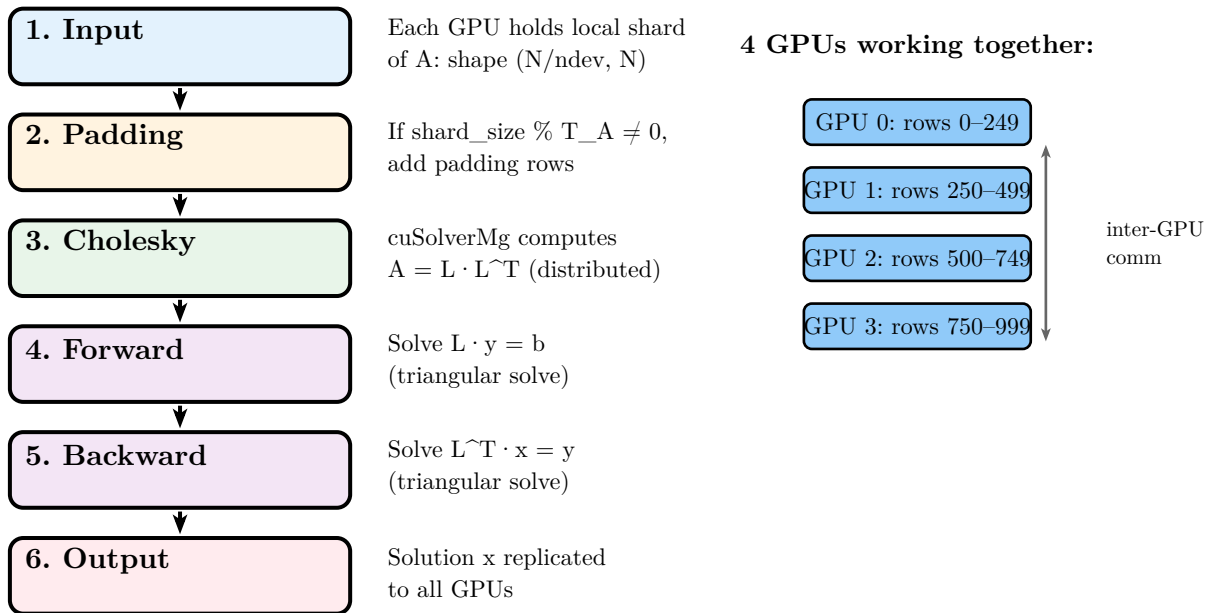


Figure 3: Distributed Cholesky solve workflow. The matrix is distributed across GPUs; cuSolverMg handles the distributed factorization and triangular solves with automatic inter-GPU communication.

3.4 Tile Size Selection

The tile size T_A controls the block granularity of the distributed algorithm:

```
1 N = 10000
2 ndev = 4
3 shard_size = N // ndev # 2500 rows per GPU
4
5 # Good choices: T_A divides shard_size evenly (no padding needed)
6 T_A = 250 # 2500 / 250 = 10 tiles ✓
7 T_A = 500 # 2500 / 500 = 5 tiles ✓
8 T_A = 2048 # 2500 / 2048 = 1.22 → padding added automatically ⚠
```

💡 Tips for choosing T_A :

- Use $T_A \geq 128$ for good performance (small tiles = slow)
- Larger $T_A \rightarrow$ fewer communication rounds, but more memory per tile
- Best: choose T_A that divides $N // \text{ndev}$ evenly to avoid padding overhead
- Common good choices: 256, 512, 1024, 2048

4 The shard_map Pattern

4.1 What is shard_map?

`jax.shard_map` lets you write code that runs **per-shard** on each device, with explicit collective operations for cross-device communication:

```
1 from functools import partial
2
3 @partial(jax.shard_map,
4         mesh=mesh,
5         in_specs=P("x", None), # Input: row-sharded
6         out_specs=P(None))      # Output: replicated
7 def my_distributed_fn(local_A):
8     # This code runs on EACH GPU independently
9     # local_A shape: (N/ndev, M) – just this GPU's shard
10
11     my_idx = jax.lax.axis_index("x") # Which GPU am I? (0, 1, 2, or 3)
12
13     # Collective: sum across all GPUs
14     global_sum = jax.lax.psum(local_sum, axis_name="x")
15
16     # Collective: gather all shards
17     full_A = jax.lax.all_gather(local_A, axis_name="x", axis=0, tiled=True)
18
19     return result
```

4.2 Common Collective Operations

Operation	Description	Output Shape
<code>psum(x, "axis")</code>	Sum values across all devices on axis	Same shape, replicated
<code>pmean(x, "axis")</code>	Mean across devices	Same shape, replicated
<code>all_gather(x, "axis")</code>	Concatenate shards from all devices	Larger (combined)
<code>axis_index("axis")</code>	Get this device's index on axis	Scalar: 0 to ndev-1
<code>pbroadcast(x, "axis")</code>	Broadcast from source device	Same on all devices

Table 4: JAX collective operations for `shard_map`

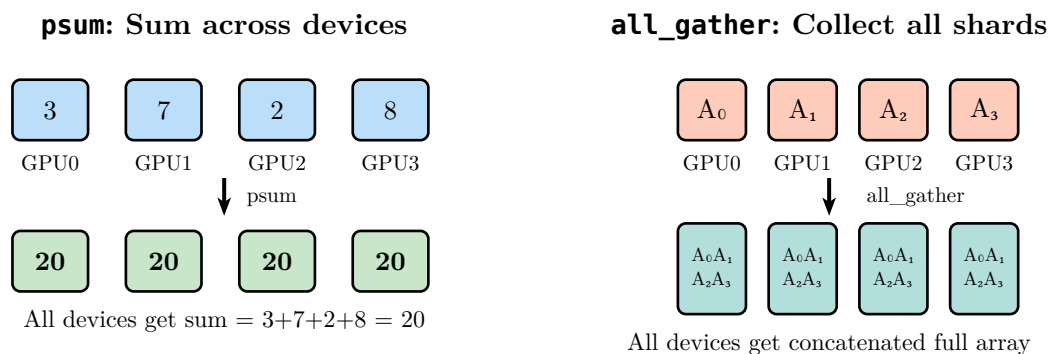


Figure 4: Collective operations illustrated. Left: `psum` sums values and replicates result. Right: `all_gather` collects all shards into complete array on each device.

4.3 Using potrs_shardmap_ctx

For advanced use inside your own `shard_map` (when you need custom logic before/after the solve):

```
1  from jaxmg import potrs_shardmap_ctx
2
3  def custom_solve(local_A, b):
4      """Runs inside shard_map context – local_A is just this GPU's shard"""
5      local_A = local_A * some_scaling_factor # Custom pre-processing
6      x, status = potrs_shardmap_ctx(local_A, b, T_A=256) # No extra shard_map
7      return x, status
8
9  result = jax.shard_map(
10     partial(custom_solve, b=b_replicated),
11     mesh=mesh,
12     in_specs=P("x", None),
13     out_specs=(P(None, None), P(None)),
14     check_vma=False, # Required for FFI calls
15 )(A_sharded)
```

5 Application: MinSR with JAXMg

5.1 The MinSR Algorithm

In Variational Monte Carlo, **minSR** (minimum-step Stochastic Reconfiguration) computes parameter updates efficiently when $N_s \ll N_p$ (fewer samples than parameters):

$$\delta\theta = \tau \cdot X^\dagger (XX^\dagger + \lambda I)^{-1} \cdot E^{\text{loc}}$$

$X \in \mathbb{R}^{N_s \times N_p}$	Centered Jacobian (samples \times parameters)
$T = XX^\dagger \in \mathbb{R}^{N_s \times N_s}$	Gram matrix — much smaller than $N_p \times N_p$!
λ	Regularization (diagonal shift, typically 10^{-4} to 10^{-2})
E^{loc}	Centered local energy vector

5.2 Multi-Node Mesh Setup

For clusters with multiple nodes, use a 2D mesh where one axis is within-node and another is across-nodes:

```

1 def get_device_grid():
2     """Organize devices: rows = local GPUs, columns = nodes"""
3     by_proc = {}
4     for d in jax.devices():
5         by_proc.setdefault(d.process_index, []).append(d)
6     hosts = sorted(by_proc)
7     return np.array([[by_proc[h][x] for h in hosts]
8                     for x in range(jax.local_device_count())]).T
9
10 def create_2d_mesh():
11     dev_grid = get_device_grid() # shape: (n_nodes, n_gpus_per_node)
12     return Mesh(dev_grid, ("S", "T"))
13     # "S" = across nodes, "T" = within node

```

Multi-Node Setup: 2 Nodes \times 4 GPUs = 8 Total

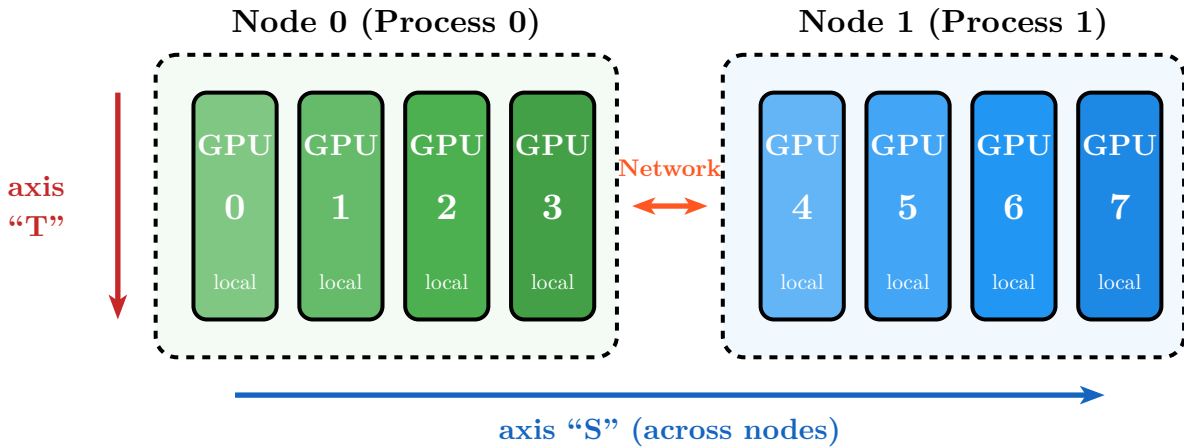


Figure 5: 2D mesh for multi-node cluster. Axis “S” spans across nodes (inter-node communication via network). Axis “T” spans GPUs within each node (fast NVLink/PCIe). Total 8 devices in a (2, 4) grid.

5.3 MinSR Data Flow

MinSR Computation Pipeline

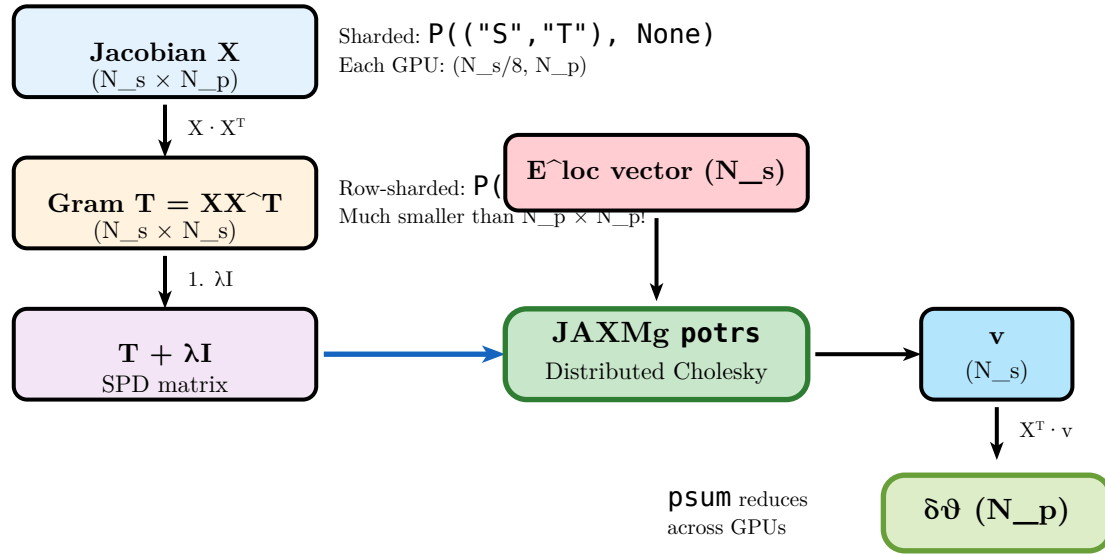


Figure 6: MinSR data flow. The Jacobian X is sharded by samples. The Gram matrix $T=XX^T$ is computed distributedly. JAXMg's **potrs** solves $(T+\lambda I)v = E^{\text{loc}}$. Final update $\delta\theta = X^T v$ uses distributed matrix-vector product with **psum** reduction.

5.4 Key Code Snippets

Sharding the Jacobian:

```
1 0_LT = jax.lax.with_sharding_constraint(
2     0_LT, NamedSharding(mesh_2d, P(("S", "T"), None))
3 ) # Samples sharded across all 8 GPUs, params replicated
```

Distributed Gram matrix (memory-efficient streaming):

```
1 def streamed_gram(0_L, chunk_size):
2     """Compute T = X @ X.T without materializing full intermediate"""
3     def step(acc, chunk):
4         full_chunk = jax.lax.all_gather(chunk, ('S', 'T'), tiled=True)
5         return acc + chunk @ full_chunk.T, None
6     return jax.lax.scan(step, zeros, 0_L_chunked)[0]
```

Distributed Cholesky solve:

```
1 v = potrs(T_reg, E_loc[:, None], T_A=2048, mesh=mesh_2d,
2           in_specs=(P("T", None), P(None, None)))
```

Final update with reduction:

```
1 def update(0_LT, v):
2     return jax.lax.psum(v @ 0_LT, axis_name=("S", "T"))
3 delta_theta = jax.shard_map(update, mesh_2d,
4     in_specs=(P(("S", "T"), None), P(("S", "T"))), out_specs=P(None))(0_LT, v)
```


6 Complete Example

```
1 import jax
2 import jax.numpy as jnp
3 from jax.sharding import Mesh, NamedSharding, PartitionSpec as P
4 from jaxmg import potrs
5
6 # Setup
7 jax.config.update("jax_enable_x64", True)
8 devices = jax.devices("gpu")
9 ndev = len(devices)
10 mesh = jax.make_mesh((ndev,), ("x",))
11
12 # Create Gram matrix (like in minSR: T = X @ X.T)
13 N_samples, N_params = 4000, 100000
14 key = jax.random.PRNGKey(0)
15 X = jax.random.normal(key, (N_samples, N_params)) / jnp.sqrt(N_samples)
16 T = X @ X.T # (4000, 4000) — much smaller than (100000, 100000)!
17
18 # Add regularization
19 T_reg = T + 1e-4 * jnp.eye(N_samples)
20
21 # RHS (centered local energies)
22 b = jax.random.normal(key, (N_samples, 1))
23
24 # Shard for distributed solve
25 T_sharded = jax.device_put(T_reg, NamedSharding(mesh, P("x", None)))
26 b_replicated = jax.device_put(b, NamedSharding(mesh, P(None, None)))
27
28 # Distributed Cholesky solve: (T + λI) v = b
29 v = potrs(T_sharded, b_replicated, T_A=256, mesh=mesh,
30           in_specs=(P("x", None), P(None, None)))
31
32 # Final minSR update: δθ = X.T @ v
33 delta_theta = X.T @ v # (100000, 1)
34 print(f"Update shape: {delta_theta.shape}")
```

7 Quick Reference

Aspect	Rule
Matrix sharding	Always $P(\langle \text{axis} \rangle, \text{None})$ — rows sharded, columns replicated
RHS vector	Always $P(\text{None}, \text{None})$ — fully replicated on all devices
Tile size T_A	Use 128–2048; choose to divide $N // \text{ndev}$ evenly if possible
Multi-node mesh	2D mesh: $(\text{"across_nodes"}, \text{"within_node"})$ axes
Memory efficiency	Use streaming Gram computation for large N_{params}
Output	Solution is always replicated to all devices

Table 5: JAXMg quick reference

✓ Summary

JAXMg makes NVIDIA's multi-GPU linear algebra “just work” with JAX:

1. **Create a mesh** — organize your GPUs into a logical grid
2. **Shard your data** — use `PartitionSpec` to distribute arrays
3. **Call the solver** — `potrs`, `potri`, or `syevd` handle everything

The library manages distributed Cholesky factorization, inter-GPU communication, padding, and memory automatically. You write high-level JAX code; JAXMg handles the CUDA complexity.