

Lattice: Large AI Model Marketplace on Solana

TLDR: A Solana-powered marketplace where consumer GPU owners earn by processing specialized portions of large AI models, overcoming the VRAM limitations of individual devices.

The Problem: Large Models Don't Fit on Consumer Hardware

- **Hard Technical Limit:** Modern AI models require 24-80GB VRAM
- **Consumer Reality:** Most GPUs have only 8-16GB VRAM
- **Current Options:**
 - Purchase expensive specialized hardware (\$5K-\$10K+ per unit)
 - Rent cloud GPU instances (\$2-\$8+ per hour)
 - Use someone's API (untrusted, no verification)

This creates a significant barrier to AI democratization.

Introducing Lattice: The AI Compute Marketplace

Lattice creates a two-sided marketplace on Solana:

- **Buyers:** AI developers who need compute for large models
- **Sellers:** GPU owners contributing compute power

The correctness of computations is ensured through our Proof of Sampling Protocol (PoSP) with only 8% verification overhead - dramatically more efficient than traditional verification methods that require 100%+ redundancy.

Result: Consumer GPUs can collectively run models too large for any single one... and it's cost-effective.

Why Lattice is Special: Key Technical Advantages

Built on EigenTensor

1. **Tensor-Centric Computation:** Universal format for any ML workload
 - Compatible with popular frameworks (PyTorch, TensorFlow models)
 - Entire models represented as computational graphs
 - TinyGrad compatible - familiar API for ML developers
2. **Memory Safety:** No arbitrary code execution - only memory-safe tensor operations allowed
 - Prevents security exploits or physical hardware damage
 - Enables trustless computation with mathematical guarantees

3. **GPU Agnosticism:** Works on any consumer GPU regardless of manufacturer

- NVIDIA, AMD, Intel all supported
- Nodes compile optimized code for their specific hardware

4. **Automatic Model Splitting:** VRAM no longer limits model size

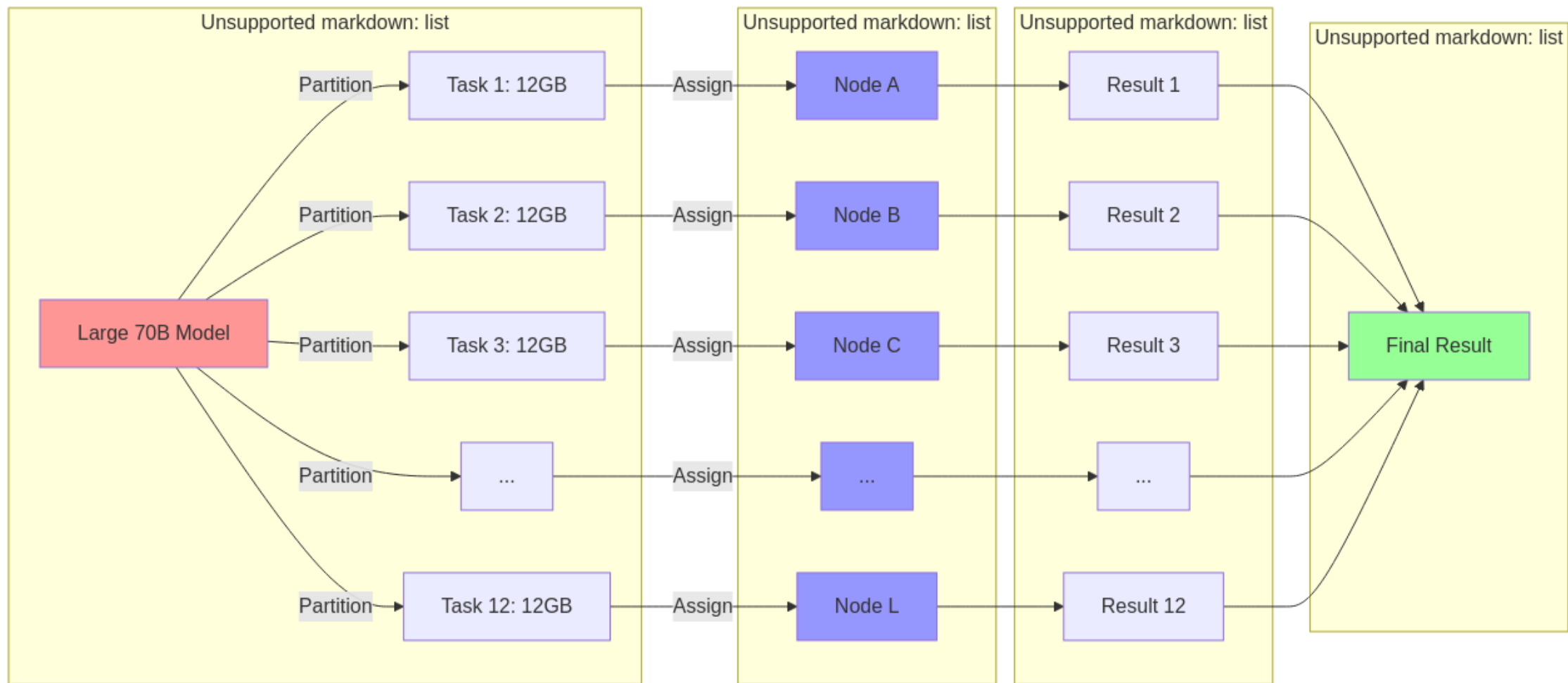
- Partitioning algorithm finds optimal split points
- Memory requirements distributed across multiple nodes
- Solves the #1 bottleneck in AI democratization
- Currently TODO, not a difficult task

Uses Proof of Sampling Protocol (PoSP)

1. **Efficient Verification:** Only 8% of work needs verification
 - Proof of Sampling Protocol vs traditional 100%+ overhead
 - Economic incentives make honesty more profitable than cheating

No other platform combines all these advantages in a working marketplace.

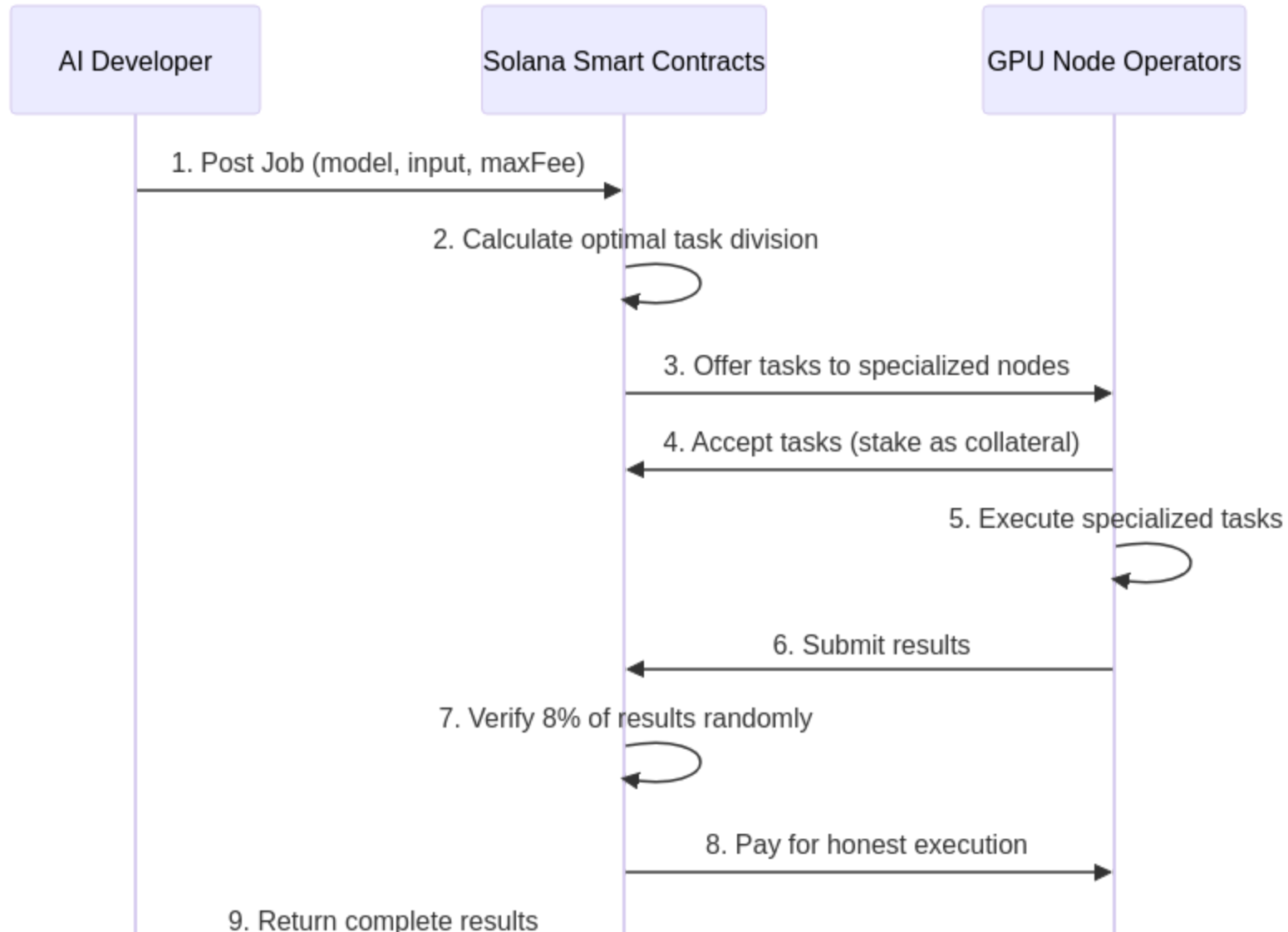
How It Works: The Big Picture



Technical Deep Dive: EigenTensor Integration

Lattice transforms this memory-safe computation into a Solana marketplace.

Solana-Powered Marketplace Flow



Market Dynamics & Incentives

For GPU Owners (Sellers):

- **Specialization:** Choose specific model components to specialize in
- **Requirements:** Stake SOL as security deposit (slashed if dishonest)
- **Optimization:** Run multiple adjacent tasks for higher earnings

For AI Developers (Buyers):

Very cheap infrastructure, leverages underutilized compute.

Comparison: What Makes Lattice Unique

Feature	Lattice	Other Decentralized Platforms
Large Model Support	✓ Any size via partitioning	✗ Limited by node VRAM
Memory Safety	✓ Only tensor operations	✓ Limited operations
GPU Agnosticism	✓ Works with any GPU hardware	✗ Often limited to specific GPUs
Model Splitting	✓ Automatic partitioning	✗ Not supported
Verification Overhead	✓ Just 8% (via PoSP)	~ Higher overhead
Marketplace Model	✓ Fully open on Solana	✓ Less efficient

Concrete Example: Running LLaMA-70B on Lattice

Traditional Approach:

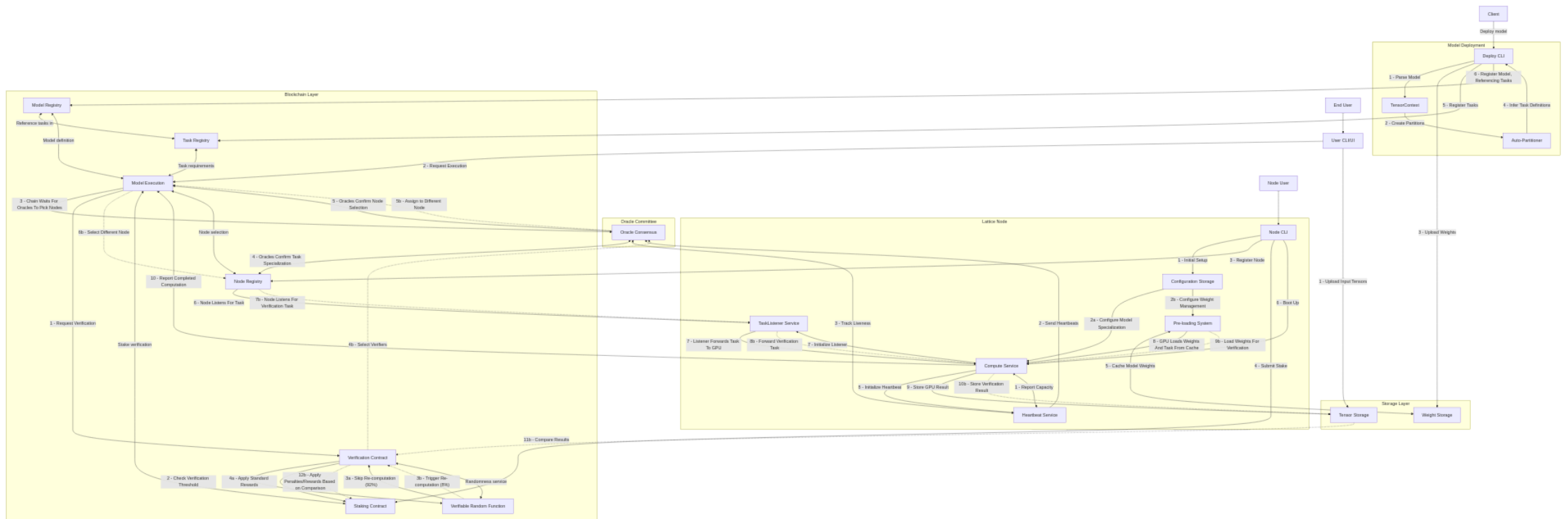
- Requires expensive A100 GPUs or high-cost cloud instances
- VRAM constraints force costly hardware choices

Lattice Approach:

- Automatically partitions the LLaMA-70B model into 12 tasks (~12GB each)
- Distributes tasks among consumer GPUs (e.g., RTX 3060+)
- **Cost:** ~\$0.50 per inference vs. \$2+ on cloud platforms

This makes state-of-the-art models accessible to everyone.

Technical Architecture Overview



Hackathon Implementation: Prototype Highlights

We've built a functional prototype demonstrating:

- EigenTensor integration for memory-safe tensor operations
- Developed a partitioning engine to break models into VRAM-friendly tasks
- Deployed Solana contracts for marketplace coordination
- Implemented Proof of Sampling (8% verification) for efficient validation
- Built an economic model for staking and payments

Demo: Partitioned MNIST for distributed execution across consumer GPUs

Join Our Marketplace

For AI Developers:

- Run large models without upfront hardware investments
- Pay only for the compute used
- Use a simple API for seamless integration

For GPU Owners:

- Earn SOL by contributing idle GPU power
- Specialize in high-demand model components
- Enter a secure, low-barrier marketplace

Thank You

Lattice: The Solana-powered Marketplace for Distributed AI Computation

[Created for a 2025 AI + Web3 Hackathon]