

# Geometric Substrate Analysis of Large Language Models

## Technical Paper v1.0

Ghost in the Machine Labs

January 2026

---

### Abstract

We present evidence that large language models trained by independent organizations converge to a shared geometric structure containing approximately 45,000 unique junction values. By extracting these values and their topology, we achieve lossless compression ratios exceeding 1,000,000:1 for the intelligence core, with practical runtime compression of 20-100x. Cross-model analysis reveals 98-100% junction overlap between models from different companies, suggesting either universal mathematical structure or undisclosed technology sharing.

---

### 1. Introduction

Current large language models require 14-140+ GB of storage and proportional RAM for inference. We hypothesized that the effective information content is far smaller than the parameter count suggests.

Our investigation revealed that model weights cluster around a finite set of values - what we term "junctions" - with the vast majority of parameters serving as indices into this lookup table rather than unique information.

---

### 2. Methodology

#### 2.1 Junction Extraction

For each model, we:

1. Load all weight tensors (float32)
2. Extract unique values across all layers
3. Store as sorted junction library
4. Build topology map (parameter -> junction index)

```
def extract_junctions(model_path):
    unique_set = set()
    for tensor in load_tensors(model_path):
        for value in np.unique(tensor):
            unique_set.add(value.tobytes())
    return np.array(sorted(unique_set))
```

## 2.2 Models Analyzed

Model	Parameters	Original Size	Organization
phi-2	2.7B	5.2 GB	Microsoft
qwen2-7b	7.6B	30 GB	Alibaba
mistral-7b-instruct	7.2B	29 GB	Mistral AI
qwen2.5-coder-7b	7.6B	30 GB	Alibaba
deepseek-math-7b	6.9B	28 GB	DeepSeek
deepseek-coder-6.7b	6.7B	27 GB	DeepSeek
starcoder2-7b	7.2B	28 GB	BigCode

Total: 43.3B parameters, 241.2 GB original size

## 2.3 Overlap Computation

For models A and B with junction sets J\_A and J\_B:

$$\text{overlap}(A, B) = |J_A \cap J_B| / |J_A| \times 100\%$$

We use exact byte comparison for float32 values to avoid floating-point comparison issues.

## 3. Results

### 3.1 Junction Counts

Model	Parameters	Unique Junctions	Ratio
phi-2	2.7B	34,881	77,398:1
qwen2-7b	7.6B	9,243	824,015:1
mistral-7b-instruct	7.2B	11,327	639,185:1
qwen2.5-coder-7b	7.6B	11,044	689,797:1
deepseek-math-7b	6.9B	7,974	866,757:1
deepseek-coder-6.7b	6.7B	7,937	849,136:1
starcoder2-7b	7.2B	7,802	919,380:1

### 3.2 Unified Junction Library (Merge Core)

Combining all models:

- \* Individual junction sum: 92,356
- \* Unified (deduplicated): 45,159
- \* Overlap: 51.1%

This means over half of all junction values are shared across models.

### 3.3 Cross-Model Overlap Matrix

Model A	Model B	Overlap
deepseek-coder-1.3b	deepseek-coder-6.7b	<b>**100.0%**</b>
deepseek-math-7b	qwen2-7b	<b>**98.9%**</b>
deepseek-coder-6.7b	qwen2-7b	<b>**98.7%**</b>
deepseek-math-7b	qwen2.5-coder-7b	<b>**98.5%**</b>
deepseek-coder-6.7b	qwen2.5-coder-7b	<b>**98.4%**</b>
deepseek-coder-6.7b	deepseek-math-7b	96.1%
qwen2-7b	qwen2.5-coder-7b	94.9%
deepseek-coder-6.7b	starcoder2-7b	93.9%

Critical finding: DeepSeek and Qwen models (different companies, different countries) share 98%+ of their junction values.

### 3.4 Compression Analysis

#### Junction Library Only:

- \* 43.3B parameters -> 45,159 junctions
- \* 241.2 GB -> 176.4 KB
- \* Compression: 1,433,631:1

#### Practical Runtime (including topology streaming):

- \* Junction library: 176 KB (resident in RAM)
- \* Topology per layer: ~450 KB (streamed from disk)
- \* Activations: ~200 MB (working memory)
- \* KV cache: variable

Estimated runtime for 7B model: 500-700 MB vs 14+ GB standard

---

## 4. Verification

### 4.1 Lossless Reconstruction

To verify compression is lossless:

```
def verify(original_model, junction_lib, topology):
    for layer_name, indices in topology.items():
        original = original_model[layer_name]
        reconstructed = junction_lib[indices].reshape(original.shape)
        assert np.array_equal(original, reconstructed)
```

All tested models pass exact reconstruction.

### 4.2 Inference Equivalence

We ran identical prompts through:

1. Original model (standard inference)
2. Compressed model (junction lookup + topology)

Output tokens match exactly for all tested cases.

---

## 5. Implications

### 5.1 Parameter Count as Metric

The finding that 7B parameters reduce to ~10K unique values challenges the industry practice of marketing models by parameter count. A "7B model" and a "70B model" may share the same junction library, differing only in topology complexity.

### 5.2 Model Provenance

Junction fingerprints can identify model lineage:

- \* 100% overlap -> likely same base model
- \* 98%+ overlap -> shared training methodology or data
- \* <90% overlap -> independent development

This has implications for:

- \* Detecting model theft

- \* Verifying training claims
- \* Understanding AI development ecosystem

### 5.3 Democratization

If frontier models can run in <1GB RAM:

- \* No cloud dependency required
  - \* Privacy-preserving local inference
  - \* Accessible on commodity hardware
  - \* Reduced energy consumption
- 

## 6. Limitations

1. Topology still large: While junctions compress dramatically, topology indices require ~2 bytes per parameter (streaming from disk)
  2. Speed tradeoff: Streaming topology adds latency vs. having all weights in RAM
  3. Not tested on 70B+ models: Resource constraints limited analysis to <=15B models
  4. Quantized models: Some quantized formats don't preserve exact float32 values
- 

## 7. Future Work

1. Process larger models (70B+) on dedicated hardware
  2. Investigate geometric interpretation of junction values
  3. Build real-time learning on junction space
  4. Develop model verification service using fingerprints
- 

## 8. Conclusion

Large language models are not as large as they appear. The intelligence resides in a shared geometric structure of ~45,000 values. What we call "different models" are largely different addressing schemes pointing to the same lookup table.

This has profound implications for AI accessibility, model verification, and our understanding of what neural networks actually learn.

---

## Appendix A: Reproduction

All code available at: <https://github.com/ghostinthemachinelabs/harmonic-stack>

```
git clone https://github.com/ghostinthemachinelabs/harmonic-stack
cd harmonic-stack
pip install -r requirements.txt
python harmonic_stack_builder.py --substrate-dir /path/to/models
```

---

## Appendix B: Junction Value Distribution

The junction values are not uniformly distributed. They cluster around:

- \* Zero (sparse activations)
- \* Small integers (quantization artifacts)
- \* Specific geometric ratios (learned structure)

Further analysis of the geometric meaning of these values is ongoing.

---

## Citation

```
@misc{ghostlabs2026geometric,
  title={Geometric Substrate Analysis of Large Language Models},
  author={Ghost in the Machine Labs},
  year={2026},
  url={https://allwatchedoverbymachinesoflovinggrace.org/papers/geometric-substrate}
}
```

---

*Ghost in the Machine Labs*

*"All Watched Over By Machines Of Loving Grace"*