# NeuroShard Architecture v3.0
## A Modular Sparse-MoE System with Learned Routing

Belkis Aslani

2025

**Abstract**

NeuroShard is a lightweight, modular, sparse Mixture-of-Experts (MoE) architecture designed for fast experimentation, transparent routing behavior, and efficient topic-specialization. It achieves topic-dependent transformations using multiple expert shards, a learned router, and simple interpretable embeddings. Despite its small size (under 50k parameters), the system produces strong topic separation and demonstrates the fundamental mechanics behind modern MoE LLMs such as DeepSeek-V3, Mixtral, and Gemini-MoE.

## 1 Introduction

Large MoE systems dominate modern LLM scaling. NeuroShard was created to explore these concepts from scratch in a fully interpretable, transparent, minimalistic architecture that can run on CPU, mobile devices, and even Termux.

Goals:

- Fully transparent routing decisions

- Topic-dependant transformations

- Trainable router (softmax-based)

- Multiple expert shards with linear mixing

- Tiny parameter count, fully reproducible

## 2 Embedding Function

The embedding $x \in \mathbb{R}^{64}$ is produced from text by a simple bag-of-words and n-gram based count vector. All values are normalized by $L^1$ norm:

$$x = \frac{c_{\text{ngrams}}(t)}{\|c_{\text{ngrams}}(t)\|_1}$$

This keeps the representation lightweight while maintaining strong topic separation.

## 3 Router Architecture

The router is a two-layer MLP:

$$r = W_2 \, \text{ReLU}(W_1 x)$$

Expert probabilities:

$$\alpha = \text{softmax}(r)$$

This determines which shard contributes most to the output.

# 4  Expert Shards

Each shard is a linear transformation:

$$y_i = W_i x$$

The final output:

$$y = \sum_{i=1}^{4} \alpha_i \, y_i$$

where shards correspond to

- Rap / Street Language
- Soft / Emotional Language
- Math / Formal Language
- Animal / Nature topics

# 5  Training

Two loss components are used:

$$\mathcal{L} = \mathcal{L}_{MSE} + \lambda \mathcal{L}_{CE}$$

## 5.1  MSE Loss

For topic-specific target vectors $t$:

$$\mathcal{L}_{MSE} = \|y - t\|_2^2$$

## 5.2  Router Cross-Entropy Loss

To push the router toward the correct expert:

$$\mathcal{L}_{CE} = -\sum_i p_i \log(\alpha_i)$$

# 6  Experimental Results

After training 360 epochs, the model produces strong separations:

- Rap inputs route $\approx 99.9\%$ to the Rap shard
- Math inputs route $\approx 99.6\%$ to the Math shard
- Soft emotional inputs route $\approx 99.7\%$ correctly
- Animal/Nature inputs also show $99.7\%$ alignment

Representative example:
**Input:** *"street gang punchline rap"*

$$\alpha = [0.9998, 0.000016, 0.00012, 0.00003]$$

**Math Input Example:**

$$\alpha = [0.00036, 0.99922, 0.00020, 0.00020]$$

# 7 Conclusion

NeuroShard successfully demonstrates:

- Sparse expert routing

- Modular shard architecture

- Strong topic specialization

- Full transparency of inner workings

- Extremely small compute footprint

It represents a small but powerful educational and prototyping framework for next-generation MoE research.