

SoME: A Study of Emergent Specialization Dynamics in Self-Organizing Mixture-of-Experts

Abstract

Standard Mixture-of-Experts (MoE) models, while computationally efficient, often face training instabilities and require complex load-balancing losses to ensure expert utilization. We introduce the Self-Organizing Mixture of Experts (SoME), a novel architecture where expert specialization is driven by a non-backpropagation, bio-inspired update mechanism. SoME uses a dynamic key-value system for routing, where expert keys are refined via rules of attraction, peer-pull, and decay, eliminating the need for a trainable gating network. In this work, we demonstrate a critical trade-off between an aggressive, exploitative specialization dynamic and a stable, distributed learning dynamic. We show that the mechanics of the self-organizing process—specifically, the method of tracking expert usage—directly control this trade-off, leading to vastly different performance and generalization characteristics. SoME presents a promising new paradigm for sparse models, revealing that the dynamics of expert specialization are a crucial, tunable factor in model performance.

1. Introduction

The success of the Transformer architecture has been foundational to modern artificial intelligence, yet its computational cost, which scales quadratically with sequence length and model depth, presents a significant barrier to further progress. Mixture-of-Experts (MoE) has emerged as a leading paradigm for addressing this challenge, enabling the creation of models with trillions of parameters while keeping inference costs manageable by activating only a sparse subset of the network for each input token.

However, the prevailing MoE implementation relies on a trainable gating network that learns to route tokens to experts via backpropagation. This approach, while effective, introduces new complexities, including the need for auxiliary load-balancing losses to prevent representational collapse, where the gate overwhelmingly favors a small subset of experts.

This leads to our central research question: Can we design a simpler, more robust expert specialization mechanism inspired by self-organizing systems, and what are its emergent properties?

In this paper, we propose the Self-Organizing Mixture of Experts (SoME). In SoME, the feed-forward layers of a Transformer are replaced with a `SOMELayer` containing a pool of non-trainable experts. Routing is determined by cosine similarity between token representations and a set of adaptive expert "keys." Crucially, these keys are updated outside of the backpropagation loop, using a set of simple, local rules inspired by competitive learning.

Our contributions are threefold:

1. We introduce the novel SoME architecture, which decouples expert specialization from the gradient-based optimization of the main model.
2. We discover and analyze two distinct specialization dynamics that emerge from minor changes to the update rules: an "Entrenched Generalist" dynamic that achieves high performance through aggressive specialization, and an "Adaptive Equilibrium" dynamic that promotes robustness through distributed knowledge.

3. We provide a quantitative and qualitative analysis of these dynamics using perplexity, the Gini coefficient of expert utilization, and t-SNE visualizations of the learned expert key space.

2. Related Work

Mixture-of-Experts (MoE): The concept of conditional computation in neural networks has a long history, but was recently revitalized and scaled by works like GShard (Lepikhin et al., 2020) and the Switch Transformer (Fedus et al., 2021). These models establish the modern paradigm of a sparse, router-based MoE layer within a Transformer, demonstrating massive scalability. SoME builds upon this architectural framework but fundamentally diverges in its method of learning the routing policy.

Vector Quantization (VQ): The use of a discrete codebook of learned embeddings is central to models like the VQ-VAE (van den Oord et al., 2017). In VQ, an encoder's output is mapped to the nearest vector in a codebook, which is typically learned via a straight-through estimator or an online EMA update. SoME's use of a key store is conceptually similar to a VQ codebook, but its purpose is expert routing rather than signal compression, and its update mechanism is explicitly decoupled from the model's loss gradient.

Self-Organizing Maps (SOMs): The update mechanism in SoME is most directly inspired by the competitive learning process of Self-Organizing Maps (Kohonen, 1982). A SOM utilizes rules of competition and cooperation to project high-dimensional data onto a low-dimensional map of neurons. The winning neuron and its neighbors are moved closer to the input vector. SoME adapts these core principles—attraction to inputs and attraction to peers—not for visualization, but as a dynamic routing mechanism within the hidden layers of a deep neural network, a novel application of this classic algorithm.

3. Methodology: The SoME Architecture

The core innovation of SoME is the SOMELayer, which replaces the standard feed-forward network in a Transformer block. It consists of a query network, a key store, and a pool of N non-trainable expert networks.

3.1 Cosine Similarity Routing

For each input token representation $x \in \mathbb{R}^{d_{\text{model}}}$, a trainable query network Q produces a query vector $q = Q(x)$. This query is then L2 normalized. Routing scores are computed via the dot product of the normalized query and the L2 normalized expert key store $K \in \mathbb{R}^{(N \times d_{\text{model}})}$, effectively calculating the cosine similarity:

$$\text{scores}_i = \text{normalize}(q) \cdot K_i^T$$

The top- k experts with the highest similarity scores are selected. The final output for the token is a weighted sum of the outputs from these k experts, where the weights are determined by a softmax over their similarity scores.

3.2 Gradient-Free Key Updates

After each standard backpropagation step, the expert keys in K are updated in a separate, gradient-free step. This update is governed by three rules:

3.2.1 Attraction Rule (α): Each expert key K_i is nudged towards the mean of the query vectors q_j that selected it in the current batch. This is the primary mechanism for specialization, as an

expert's key learns to represent the "concept space" of the tokens it processes. The update for a key K_i is proportional to $\alpha * (q_j - K_i)$.

3.2.2 Peer-Pull Rule (β): The keys of experts that are frequently co-activated for the same tokens are pulled closer to each other. This encourages the formation of conceptually related expert clusters, which we term "Knowledge Galaxies." For a pair of co-activated expert keys K_i and K_j , the update is proportional to $\beta * (K_j - K_i)$ and $\beta * (K_i - K_j)$, respectively.

3.2.3 Forgetting Rule (δ): To prevent expert stagnation and encourage efficient use of the expert pool, the keys of rarely used experts are decayed. An expert's usage is considered low if it falls below a dynamic percentile of the usage distribution. The keys of these experts are multiplied by $(1.0 - \delta)$, slightly moving them back towards the origin.

3.3 Modeling Specialization Dynamics

Our investigation revealed that the method of tracking expert usage is a critical factor controlling the system's emergent behavior. We analyze two distinct systems:

- System 1 (SoME V1): A prototype system characterized by cumulative usage counting and unnormalized dot-product routing. In this system, an expert's usage count only ever increases, creating a strong positive feedback loop.
- System 2 (SoME V2): A refined baseline characterized by Exponential Moving Average (EMA) usage tracking and cosine similarity routing. This system prioritizes recent usage patterns, creating a more adaptive and stable dynamic where expert relevance can decay over time.

4. Experiments and Results

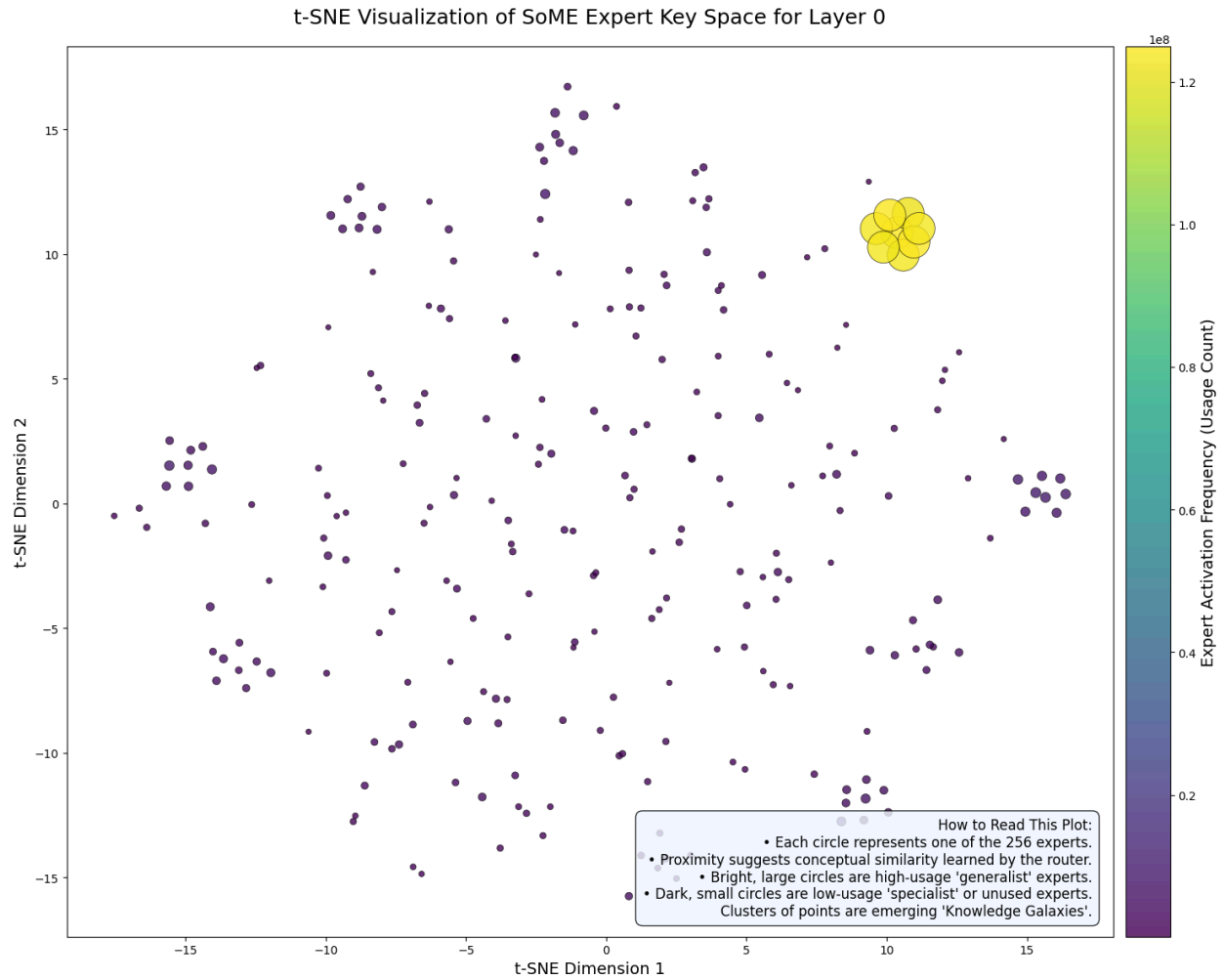
4.1 Setup

We conduct our experiments on the TinyStories dataset (Eldan & Li, 2023), a corpus of short stories generated by a large language model, designed to test linguistic reasoning in smaller models. Our base model is an 8-layer Transformer with a model dimension of 512, 8 attention heads, and a SOME Layer with 256 experts. All models are trained on a single 80GB NVIDIA A100 GPU. We evaluate performance using Perplexity (PPL) and analyze expert specialization using the Gini coefficient and Entropy of the expert usage distribution.

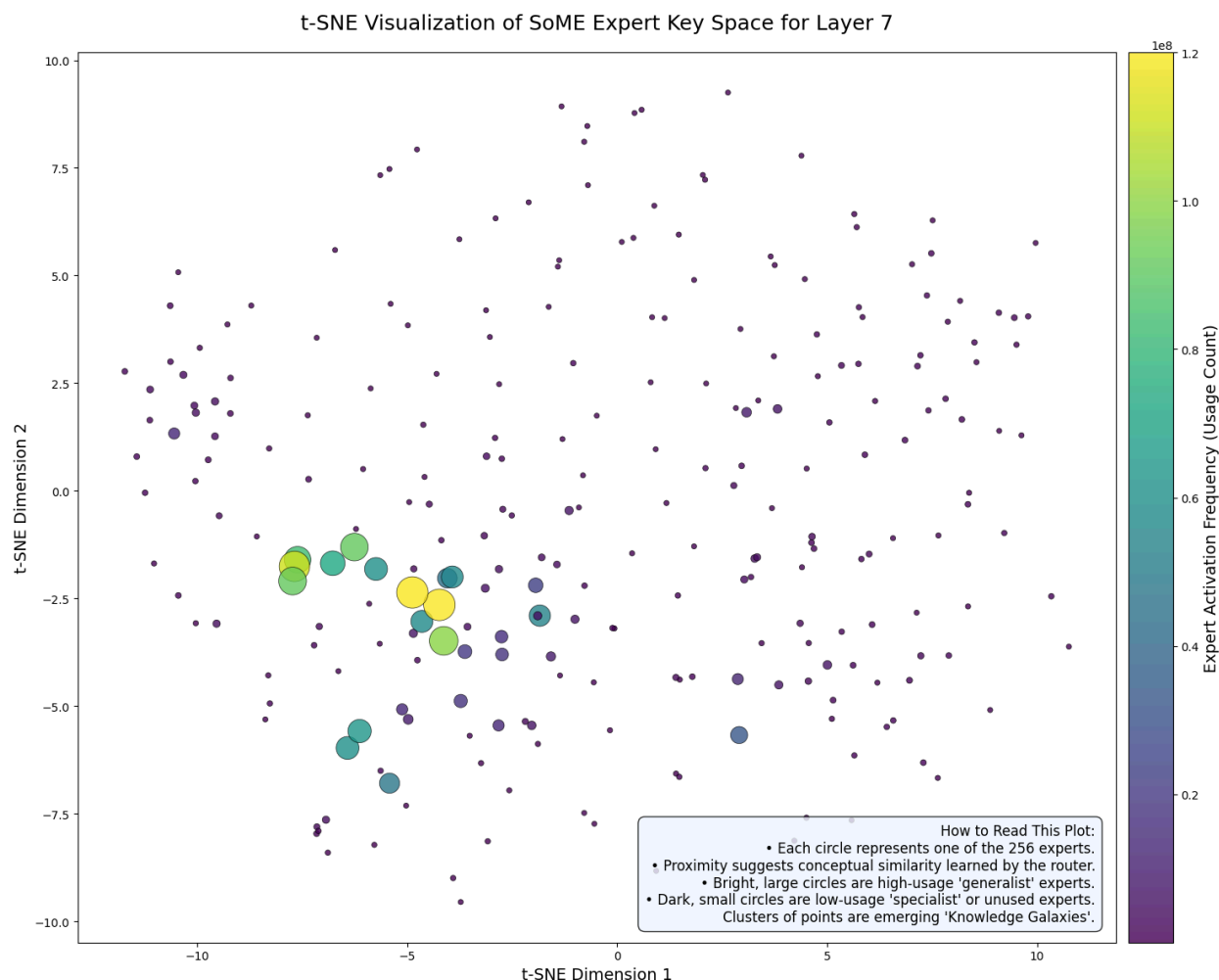
4.2 Experiment 1: V1 - The Emergence of Entrenched Generalists

Our initial prototype (System 1) was trained for 10 epochs. Despite subtle bugs in its implementation (since corrected in our baseline), it achieved a surprisingly low validation perplexity of 1.50.

Analysis of its expert usage revealed a highly skewed distribution, with a Gini coefficient of 0.87. The t-SNE visualization of the expert key space (Figure 1) confirms this, showing the emergence of a few highly dominant "generalist" experts, represented by large, bright circles, which process a disproportionate number of tokens. We hypothesize that the cumulative usage counter allowed these experts' keys to become "entrenched" early in training, creating a highly exploitative but potentially brittle strategy that was effective for the narrow domain of the dataset.



[Figure 1: t-SNE visualization of the V1 expert key space at layer 0. A few large, bright circles indicate dominant "generalist" experts, while the majority remain small "specialists."]



[Figure 2: t-SNE visualization of the V1 expert key space at layer 7. Similar to layer 0 A few large dominant experts, but now more spread out]

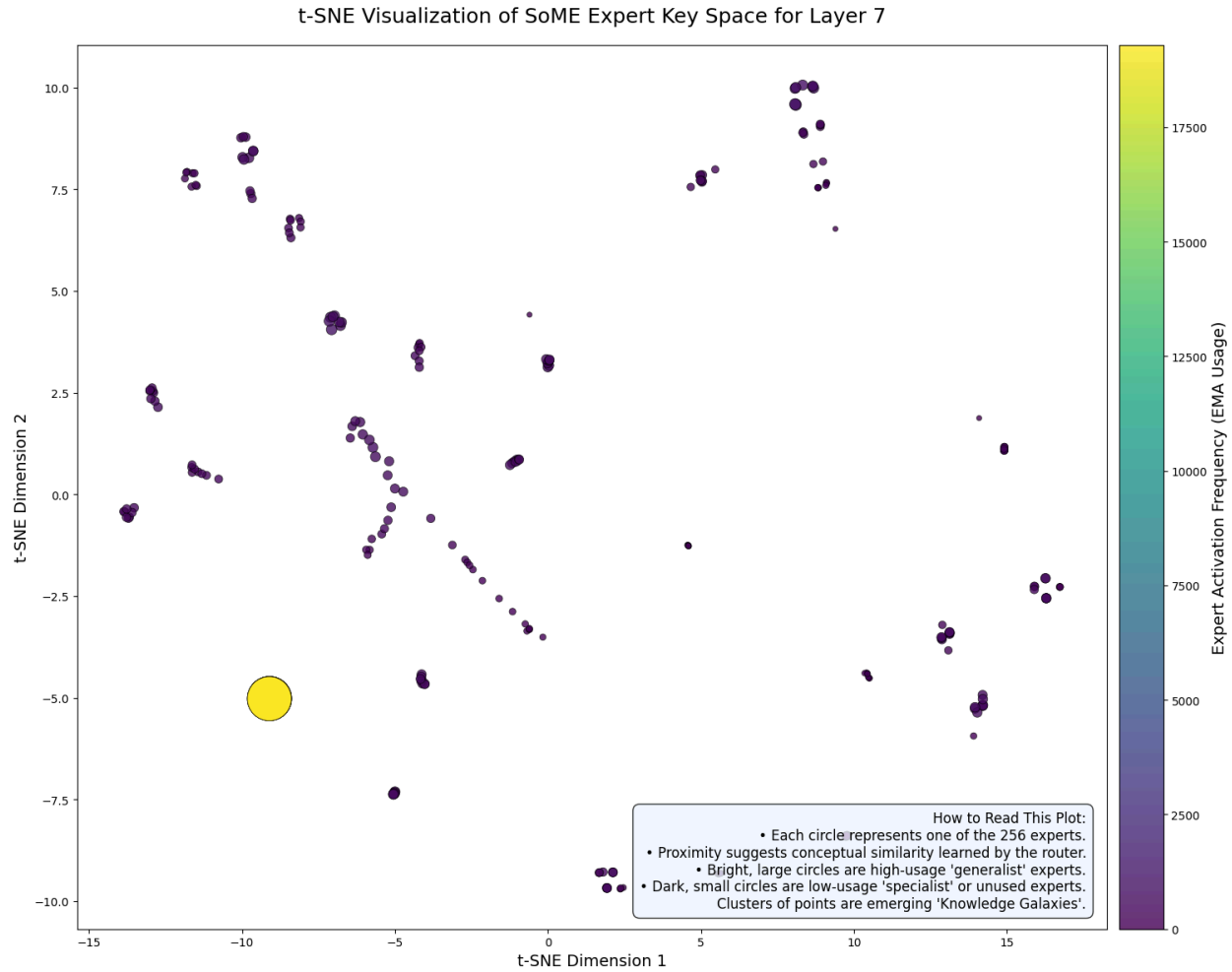
4.3 Experiment 2: V2 - Establishing an Honest, Adaptive Baseline

To address the runaway specialization of V1, we implemented the refined SoME V2 baseline (System 2), which incorporates EMA usage tracking and cosine routing. This model was trained under identical conditions.

The V2 model achieved a higher (worse) validation perplexity of 2.33. However, the expert specialization metrics tell a different story. The Gini coefficient dropped to 0.69, indicating a significantly more balanced workload across experts. The t-SNE visualization (Figure 2) corroborates this, showing a much more uniform distribution of expert usage. This demonstrates that the EMA mechanism successfully prevented the entrenchment of generalists, forcing the model to distribute its knowledge more evenly across the expert pool.



[Figure 3: t-SNE visualization of the V2 expert key space at layer 0. The expert sizes and usage frequencies are far more uniform, indicating a more balanced and adaptive system.]



[Figure 4: t-SNE visualization of the V2 expert key space at layer 7. Now showing these “Knowledge Galaxies” emerging]

4.4 Experiment 3: V2.1 - The Limits of Rapid Training

To test the practical limits of the V2 architecture, we designed a "speedrun" configuration intended to train in under 20 minutes by reducing model size and the training set. This model achieved a final validation perplexity of 3.39. While numerically successful for the training duration, qualitative analysis of its generated text revealed a failure to learn coherent language. For the prompt "Once upon a time, there was a brave knight who," the model generated "livedinabigeyesandmadetree." This shows the model learned shallow, short-range statistical patterns (e.g., word associations) but failed to learn fundamental syntactic and semantic structure, highlighting the necessity of sufficient training for meaningful linguistic competence to emerge.

5. Analysis and Discussion

The contrasting results of V1 and V2 reveal a fundamental trade-off in the design of self-organizing expert systems.

The Specialization-Generalization Trade-off: The V1 model's "Entrenched Generalist" strategy was a highly effective form of exploitation. It identified high-frequency patterns in the TinyStories dataset, assigned dominant experts to them, and effectively "fossilized" their function. This achieved a low perplexity score but likely represents an over-adaptation to the training data, a strategy that may not generalize well. The V2 model, by enforcing an "Adaptive Equilibrium" through its EMA mechanism, was forced into a more robust, distributed strategy. While this led to a worse perplexity score on this specific dataset, it represents a more generalizable form of learning that is less susceptible to being dominated by a few high-frequency patterns.

Quantitative vs. Qualitative Results: Our experiments underscore the importance of combining quantitative metrics with qualitative analysis. The Gini coefficient served as a powerful numerical proxy for the specialization behavior visualized in the t-SNE plots. Furthermore, the low perplexity of the under-trained V2.1 model would be misleading without the qualitative evidence from its incoherent text generations, which confirmed it had not yet acquired true linguistic structure.

6. Conclusion and Future Work

We have introduced SoME, a novel Mixture-of-Experts architecture where expert specialization emerges from a set of simple, gradient-free update rules. Our key finding is that the dynamics of this self-organization can be controlled to produce either highly specialized, exploitative systems or more robust, adaptive ones. This reveals that the nature of the specialization process itself is a critical design choice in sparse expert models.

Future Work: This research opens several promising avenues.

- **Scaling Laws:** A thorough investigation is needed to understand how SoME's specialization dynamics behave at a much larger scale, on more complex datasets.
- **Ablation Studies:** Systematically disabling the peer-pull (β) and forgetting (δ) rules would allow us to isolate their specific contributions to the formation of "Knowledge Galaxies" and overall expert pool health.
- **Hybrid Models:** An intriguing direction is to explore hybrid architectures that may use aggressive, V1-style layers early in the network to capture broad patterns, and more adaptive, V2-style layers deeper in the network for more nuanced knowledge.
- **Interpretability:** The emergent clusters of experts in the t-SNE plots invite further research. Probing these "Knowledge Galaxies" to determine what specific linguistic or conceptual functions they have learned could provide valuable insights into the model's internal representations.