

A Transparent MoE Architecture For Sequence Level Routing: Synthesis of Predictable Specialization

1 Executive Summary and Theoretical Foundation of Predictable Specialization

1.1 Project Context: The Specialization Crisis in Overlapping Domains

The efficacy of Mixture-of-Experts (MoE) models fundamentally rests on the ability of experts to specialize. Conventional practice relies on **token-level routing (R_T)**, a mechanism known primarily for cultivating specialization based on fine-grained syntactic or lexical features.¹ This strategy proves robust for highly divergent tasks, such as digital image recognition (MNIST, 96% accuracy observed for one expert) or code/mathematics tasks (where specialization concentrates clearly in a designated expert).¹

However, this architecture encounters a significant challenge—often termed the **“specialization crisis”**—when dealing with semantically overlapping domains, such as business, psychology, and history. Initial experiments confirm that in these generalist, overlapping fields, the specialization effect vanishes, leading to an undesirable **“collapse”** where a single, general expert dominates the routing outcome.¹ This phenomenon is hypothesized to occur because token-level routing tends to centralize similar underlying semantic and syntactic structures into one expert, a separation that collapses when input domains share high linguistic and contextual overlap. The limitation of existing sequential specialization approaches, such as Branch-Train-MiX (BTX), further emphasizes this problem, as BTX, which uses a pretraining step followed by R_T fine-tuning, also experiences this collapse in general domains.¹

1.2 The Sequence-Level Hypothesis and Necessary Trade-offs

The core hypothesis under investigation proposes that enforcing routing based on the entire input sequence or sentence (\mathbf{R}_S) may overcome these limitations by promoting context-dependent, global specialization. Sequence-level routing compels the router to consider holistic contextual information, which preliminary studies suggest can justify the existence of weak expert specialization related to topics.¹ This explicit focus on the sequence context aims to draw clearer domain boundaries during training, mitigating the centralizing tendencies observed with token-level systems.

However, pursuing specialization via \mathbf{R}_S requires navigating a necessary performance trade-off. Empirical evaluations comparing routing paradigms demonstrate that \mathbf{R}_T configurations, particularly the Layer-wise Top-2 standard prevalent in modern architecture, consistently achieve superior validation perplexity (PPL), often surpassing dense model baselines that match the MoE’s total parameter count.¹ In contrast, Layer-wise \mathbf{R}_S (Top-2) achieves performance only comparable to a dense model counterpart matched for active parameter count.¹ This outcome implies that the architectural goal must prioritize **“architectural transparency and predictable specialization”**—essential for achieving the secondary objective of alignment steering—over raw PPL efficiency.

2 The Architectural Conflict: Routing Granularity and Specialization Dynamics

2.1 Token-Level Routing (R_T) Dynamics: Syntactic Bias and Capacity

Token-level routing, despite its performance efficiency, exhibits a strong bias toward syntactic and lexical features, often failing to capture high-level topical specialization. Mixtral 8x7B, a contemporary high-capacity MoE model (47B total parameters, 13B active, Top-2 routing), confirms this pattern. Analysis of Mixtral’s routing reveals that expert selection is heavily structured around syntax, with routing paths specializing in elements like punctuation, Python keywords (e.g., ‘self’), or code indentation, but lacking

obvious patterns based on high-level topic or domain.¹ This architectural inclination explains why R_T struggles when domains are semantically intertwined, as observed in the initial experiments.

A crucial counterpoint to this limitation is demonstrated by models leveraging extreme granularity. The OLMOE-1B-7B model, utilizing highly fine-grained routing (64 total experts with 8 active per layer), successfully achieved high degrees of domain and vocabulary specialization.¹ This demonstrates that specialization is achievable with R_T by dramatically increasing combinatorial capacity, rather than necessarily changing the routing unit. However, the resulting specialization remains intrinsically token-dependent, potentially yielding noise or inconsistency unsuitable for reliable behavioral steering applications.

For performance optimization under R_T , empirical studies indicate that increasing the total number of experts (N) provides greater benefit, supporting the routing of tokens within a sequence to different specialized experts (high divergence).¹

2.2 Sequence-Level Routing (R_S) for Topic Enforcement

The rationale for shifting to sequence-level routing centers on compelling the model to select experts based on the comprehensive context of the input, thereby enforcing explicit domain boundaries. Empirical work confirms that R_S demands an increased number of activated experts (K) to perform effectively, highlighting that it relies on aggregating contributions from multiple specialized pathways to capture the overall meaning of the sequence.¹ This aggregation is necessary to compensate for the computational loss incurred by making a single routing decision for the entire context, rather than localized decisions for every token.

The theoretical foundation for domain enforcement through R_S is strongly supported by the **Domain Expert Mixture (DEMIX)** architecture.¹ DEMIX implements sequence-level routing by utilizing explicit domain metadata (document provenance) as a non-learned conditioning variable during training.¹ This structural enforcement of domain segregation yields a modular language model capable of adaptation without catastrophic forgetting, validating that hard-coded or explicitly conditioned R_S successfully creates specialized, removable knowledge containers.¹ To achieve competitive outcomes using learned R_S , the architectural design must prioritize increasing the number of activated experts (K) to maximize contextual processing depth, as empirical data shows that R_S suffers significant performance degradation when constrained to $K = 1$.¹

Layer-wise routing mechanisms are crucial for maintaining efficiency and architectural flexibility, irrespective of whether routing occurs at the token or sequence level. Studies confirm that layer-wise routing consistently outperforms architectures employing only a single, global router (where a routing decision made at the first layer applies to all subsequent layers).¹ Global routing creates an information bottleneck, preventing the model from adapting its expert selection as the input representation evolves and abstracts through deeper layers. Layer-wise R_S is particularly necessary for transparent architectural analysis, allowing researchers to study how the perceived topic or context shifts through the depth of the transformer stack.

3 Strategies for Enforcing Predictable Specialization

The central difficulty in achieving predictable specialization is ensuring that the specialized expertise gained during initial training remains durable and localized, resisting the centralizing pressures encountered when training on generalist, overlapping data.

3.1 Advanced Initialization: Moving Beyond BTX

The initial BTX-style approach, involving asynchronous domain pretraining followed by fusion, proved insufficient in preventing specialization collapse in overlapping domains.¹ To overcome the convergence slowdown and lack of specialization observed in traditional upcycling (where all experts are initialized as replicas of a dense FFN)¹, an advanced initialization method is required.

Drop-Upcycling offers a refined solution through **Diversity Re-initialization**.¹ This method strategically injects divergence into the replicated expert weights. It involves determining an intensity ratio (r , typically $r \approx 0.5$ is optimal), randomly sampling common indices along the intermediate dimension of the FFN, dropping those selected weights (column-wise or row-wise), and re-initializing them statistically using the mean and standard deviation of the original weights.¹

This statistical intervention is critical because it forces independent learning pathways from the onset. The final model is a hybrid, retaining original pre-trained expertise (approximately 50% of the knowledge) while introducing enough structured noise (approximately 50% re-initialized weights) to compel experts to specialize uniquely, thereby preventing the symmetrical attraction of all tokens to the generalist expert.¹ The research project must adopt this Diversity Re-initialization methodology, applying it to its pre-trained BTX experts during the MoE fusion stage to engineer the required divergence in the semantic space.

3.2 Expert Granularity and Combinatorics

Modern high-performance MoE models often utilize extreme granularity. The OLMOE system, for instance, employs 64 experts with 8 active per token.¹ This level of granularity significantly increases the available routing space, providing billions of possible expert combinations per layer (**4,426,165,368** combinations for 64 choose 8).¹ For R_S , which depends on activating multiple experts ($K > 1$) to capture nuanced context, this combinatorial richness is indispensable for enabling the model to distinguish between complex, overlapping inputs (e.g., separating "business history" context from "political history" context).

Regarding architecture simplification, research suggests caution: experiments with explicit shared experts—designed to centralize common knowledge—were found to be ineffective in the OLMOE framework, performing slightly worse than purely routed architectures.¹ This outcome supports the preference for maximizing the combinatorial richness of the routing mechanism, relying solely on learned routing rather than pre-defining parameter usage via fixed, shared experts.

4 MoE Stability and Mitigating Catastrophic Drift

The strategy of combining initial specialization training with subsequent generalist MoE fine-tuning (continual learning) inherently raises the risk of catastrophic forgetting (CF), where the initial domain expertise dissipates or "drifts" as weights are optimized for broader tasks.²

4.1 Regularization for Utilization and Stability

The MoE training regime commonly incorporates auxiliary losses to address stability and expert utilization.

- **Router Z-Loss (\mathcal{L}_{RZ}):** This auxiliary loss penalizes excessive magnitude in the router logits, which is known to cause numerical instability, particularly in low-precision training environments.¹ Empirical evidence from OLMOE confirms that incorporating \mathcal{L}_{RZ} (with a typical weight $\beta \approx 0.001$) enhances overall training stability and performance, making its adoption standard practice as it does not conflict with the goal of specialization.¹
- **Load Balancing Loss (LBL):** LBL is designed to prevent expert collapse or the emergence of "dead experts" by penalizing unequal utilization.¹ While LBL (typically $\alpha \approx 0.01$) is crucial for utility, promoting balanced load distribution¹, its benefit to pure performance metrics (PPL) is questionable, as some studies found LBL to be negligible for validation perplexity.¹ Furthermore, a strongly enforced LBL can be counterproductive to the specialization goal, potentially forcing domain experts to process tokens outside their learned domain simply to satisfy the mandated utilization targets.

4.2 Mitigating Catastrophic Forgetting via Weight Constraints

To address the inherent conflict between specialization retention and mandated utilization (LBL), an alternative regularization method targeting weight stability is necessary. The architectural proposal includes a recommendation for regularization via **weight distance constraints** (as discussed by Ludziejewski et al.).⁴

Mechanisms employing weight distance constraints (such as variations of Elastic Weight Consolidation) directly address CF by penalizing large deviations in the weights (θ) away from their previously learned optimal state (θ^*).² This form of regularization provides an ideal lever for the project, as it specifically mitigates the catastrophic drift of the pre-trained domain expert weights during the generalist MoE fine-tuning stage. By stabilizing the core specialized knowledge, it ensures that the pre-trained

expertise remains identifiable and durable, reinforcing the "path of least resistance" for routing the intended domain input. This approach offers a more direct control over preserving specialization than using a strong LBL, which primarily enforces uniform usage.

5 Architectural Transparency for Alignment and Safety Steering

The concentration of misaligned behaviors within identifiable experts, a secondary goal of the project, is a critical application of the achieved architectural transparency.

5.1 Behavioral Specialization and Localization

Recent advancements confirm the viability of using expert modules as levers for behavioral steering. The **SteerMoE** framework demonstrates that highly specialized experts can be detected based on their distinct activation patterns when processing inputs designed to elicit contrasting behaviors (e.g., faithful response versus unfaithful response, or safe output versus adversarial attack).⁶ By actively controlling these behavior-linked experts (selective deactivation or amplification) during inference, SteerMoE successfully achieved quantifiable control over alignment properties, increasing safety by up to 20% and enabling severe adversarial attacks that drop safety by 100%.⁶

This confirms that the levers for alignment steering exist within the modular architecture. The purpose of the R_S architecture is to elevate this implicit, often noisy, behavioral specialization into an explicit, predictable, topic-based specialization. If a sequence-level routing architecture successfully routes all inputs concerning "psychological advice" to a designated expert, then any misaligned outputs related to that specific topic are geographically localized within that expert, making the subsequent behavioral control highly targeted and reliable.

5.2 Metrics for Transparency and Predictability

Achieving architectural transparency necessitates a stringent measurement regime that validates the predictability and durability of the routing decisions. Key metrics derived from modern MoE analysis are essential for this purpose:

- **Domain Specialization:** This core metric quantifies the proportion of tokens from a specific domain D that are routed to a particular expert E , directly addressing the primary research goal and confirming the existence and strength of predictable domain boundaries.¹
- **Vocabulary Specialization:** This metric links domain routing to low-level features by calculating the proportion of tokens with a specific token ID (vocabulary element) routed to an expert E .¹ It ensures that the sequence-level routing decision is coherently supported by underlying lexical features, reinforcing the distinction between specialized domains.
- **Router Saturation:** This measures the stability of the routing mechanism, quantifying how early in the training process an expert's selection probability locks onto the final learned selection probability for a given input.¹ High saturation confirms that the R_S mechanism quickly and reliably establishes the predictable domain containers, crucial for the reliability required for subsequent steering.

The implementation path for leveraging this transparency involves three primary steps: (1) deploying the highly specialized R_S architecture with Drop-Upcycling initialization; (2) training until high Domain Specialization metrics confirm the predictable domain containers; and (3) subsequently utilizing SteerMoE-like techniques to detect and isolate misaligned behaviors that manifest within these already localized topic containers, enabling precise inference-time suppression.⁶

6 Conclusion and Implementation Recommendations

6.1 Key Research Answers and Confirmation

The comprehensive architectural synthesis provides answers to the key research questions posed by the project:

- **Will sequence-level routing combined with pretraining strengthen domain expertise?**
Yes. Sequence-level routing compels the router to utilize global contextual information, justifying the observed existence of weak topic specialization.¹ The success of explicit domain-conditioned routing (DEMIX) further confirms that this architecture successfully enforces specialization boundaries.¹ However, performance depends heavily on enabling a high number of activated experts (\mathbf{K}) to process the comprehensive context.¹
- **How well does pretraining predict expert routing in large MoEs?** Pretraining alone (BTX-style initial experts) is insufficient and fails in overlapping semantic domains due to the centralization effects of downstream generalist training.¹ Predictable routing requires an architectural intervention—specifically the **Diversity Re-initialization mechanism** inherent in Drop-Upcycling—to statistically break the initial parameter symmetry and engineer specialized pathways prior to full MoE fine-tuning.¹
- **How do outcomes differ between narrow vs. broad tasks?** Narrow, easily distinguishable tasks (e.g., math, code) lead to predictable and highly localized specialization even under token-level routing (confirmed by initial BTX results).¹ In contrast, broad or overlapping tasks necessitate the use of R_S and stringent regularization (such as weight distance constraints) to prevent the inevitable collapse into the dominating generalist expert.⁴

6.2 Final Recommendations

Based on this analysis, the development of a transparent MoE architecture for predictable sequence-level routing should adopt the following specific strategies:

- **Routing Strategy and Configuration:** Implement **Layer-wise Sequence-level Routing (\mathbf{R}_S)**. The design should prioritize increasing the number of activated experts ($\mathbf{K} > \mathbf{1}$) to maximize the capacity for capturing complex contextual relationships, thereby offsetting the inherent performance efficiency advantage of token-level systems.¹
- **Initialization Strategy:** Adopt a hybrid training method: initial parallel domain training (BTX) followed by fusing those experts into the MoE layer using **Drop-Upcycling’s Diversity Re-initialization** principle (targeting an optimal re-initialization ratio $r \approx 0.5$).¹ This ensures both domain knowledge transfer and the necessary architectural divergence in overlapping fields.
- **Regularization and Stability:** Incorporate **Router Z-Loss ($\beta \approx 0.001$)** universally to maintain training stability.¹ Critically, substitute an aggressive Load Balancing Loss with **weight distance constraints** (e.g., regularization penalty on weight drift).⁴ This targeted regulation mitigates catastrophic forgetting of the pre-trained domain expertise without forcing experts into unnatural utilization patterns that compromise specialization.
- **Architectural Validation:** Employ rigorous quantitative metrics, specifically **Domain Specialization** and **Router Saturation**, throughout training to verify that the R_S mechanism successfully locks onto durable and predictable domain boundaries, thereby confirming the architectural transparency required for successful steering applications.¹

7 Literature Review: MoE Specialization and Routing Architectures

Table 1: Comparative Literature Review on MoE Specialization and Routing Architectures

Paper/Method	Routing Unit	Specialization Mechanism
Fan et al. (2024) ¹	Token/Sequence	Empirical tuning of N (experts) and K (accuracy)
Sukhbaatar et al. (BTX) (2024) ¹	Token (Top-2)	Parallel domain-specific pretraining followed by fine-tuning
Gururangan et al. (DEMIX) (2021) ¹	Sequence/Document	Explicit conditioning on domain metadata
Nakamura et al. (Drop-Upcycling) (2025) ¹	Token (Top-2)	Diversity Re-initialization ($r = 0.5$) via stratified sampling
Muennighoff et al. (OLMOE) (2025) ¹	Token (Dropless)	Fine-grained granularity (64/8 experts), trajectory-based routing
Fayyaz et al. (SteerMoE) (2025) ⁶	Token	Implicit behavioral specialization (safety, fairness)
Ludziejewski et al. (2025) ⁴	N/A (Scaling Laws Focus)	N/A