

Geometric Routing and Continual Learning

Rethinking How Language Models Evolve



Opaque Decisions



Catastrophic Forgetting



Hallucinations
& Uncertainty

A New Approach to Smarter, Adaptive AI

★ Geometric Routing and Continual Learning: Rethinking How Language Models Evolve

Despite impressive gains in scale, modern language models still struggle with three fundamental problems: opaque decision-making, catastrophic forgetting, and hallucinations under uncertainty. These issues are not incidental; they are structural consequences of how models are trained and routed today. I'll outline an alternative architecture that addresses these failures by making routing, learning, and uncertainty explicit.

Motivation and Industry Context

The machine learning industry has struggled with a recurring set of systemic problems as models have grown larger, more complex, and more deeply integrated into real-world systems. While no single architecture can resolve all challenges in machine learning, the **primary focus of GRCLM (Geometric Routing and Continual Learning Model)** is to address three foundational issues that repeatedly limit the reliability, safety, and evolvability of modern language models.

At a high level, GRCLM is designed to:

1. **Increase transparency and accuracy** by replacing opaque learned gating mechanisms with explicit geometric routing.
2. **Enable continual learning** in large language models through geometry-driven specialization rather than destructive retraining.
3. **Provide uncertainty clarity** by explicitly identifying low-confidence and unknown inputs instead of producing false certainty.

These goals reflect long-standing industry pain points and motivate the architectural choices behind GRCLM.

Core Industry Problems Restated

Opaque expert selection in Mixture-of-Experts systems

Traditional MoE architectures rely on learned gating networks that are difficult to interpret, debug, and audit. GRCLM replaces these black-box gates with **deterministic geometric routing** based on similarity, distance, entropy, and support, making expert selection transparent, explainable, and stable.

Catastrophic forgetting in continually evolving language models

Conventional fine-tuning approaches overwrite existing knowledge when new data is introduced. GRCLM enables **true continual learning** by isolating new knowledge into specialists and using geometric signals to determine when experts should be added, reused, merged, or retired, without retraining the entire model.

Uncontrolled hallucinations and false confidence

Standard language models lack a principled way to represent uncertainty and often hallucinate when operating outside their training distribution. GRCLM provides **explicit uncertainty signaling** by detecting when inputs fall outside known domain geometry, allowing the system to respond with “I don’t know,” defer action, or collect evidence rather than fabricate answers.

Secondary Problems Addressed by GRCLM

While not the sole intention of the architecture, the proposed solution naturally encompasses and offers plausible resolutions to several additional industry challenges:

- Opaque expert selection
- Catastrophic forgetting
- Hallucinations under uncertainty
- Costly full retraining
- Model drift
- Lack of provenance
- Cold-start expert instability
- Unknown-domain blindness
- ML-SDLC mismatch
- Unsafe scaling of knowledge

These issues emerge as downstream consequences of monolithic model design and opaque learning dynamics, and GRCLM addresses them through structural rather than heuristic changes.

Description of Each Issue

Opaque expert selection

Learned gating functions obscure why specific experts are chosen, limiting trust and debuggability.

Catastrophic forgetting

New training data overwrites previously learned representations, erasing prior capabilities.

Hallucinations under uncertainty

Models produce confident but incorrect outputs when encountering unfamiliar inputs.

Costly full retraining

Updating knowledge typically requires retraining massive models, incurring high computational and operational costs.

Model drift

Repeated fine-tuning causes gradual, uncontrolled changes in model behavior over time.

Lack of provenance

Most models cannot explain which internal components contributed to an output or why.

Cold-start expert instability

Newly added experts perform poorly until sufficient training stabilizes learned gates.

Unknown-domain blindness

Models lack mechanisms to detect when inputs fall outside all known domains.

ML–SDLC mismatch

Machine learning systems do not align well with standard software engineering practices such as versioning, modularity, and rollback.

Unsafe scaling of knowledge

Increasing model capacity often increases hallucination risk rather than reliability.

Introduction to GRCLM

GRCLM is a modular architecture that decomposes a language model into a **stable generalist trunk** and a set of **specialist heads**, coordinated through a **geometric routing mechanism**. Instead of learning routing decisions implicitly, GRCLM represents each domain geometrically using centroids, covariance estimates, entropy, and support statistics.

This geometric representation enables the system to reason explicitly about similarity, uncertainty, and novelty, forming the basis for controlled specialization and continual learning.

How GRCLM Resolves Each Problem

- **Opaque expert selection** is resolved through explicit geometric routing signals that can be inspected and audited.
- **Catastrophic forgetting** is avoided by isolating new knowledge into specialists without modifying the frozen trunk.
- **Hallucinations under uncertainty** are reduced by detecting out-of-distribution geometry and enabling abstention or deferral.
- **Costly full retraining** is mitigated by retraining only affected specialists rather than the entire model.
- **Model drift** is constrained by limiting updates to localized components.
- **Lack of provenance** is addressed by tracking which specialists contribute to each response.
- **Cold-start instability** is reduced because geometric routing works immediately without requiring learned gates.
- **Unknown-domain blindness** is resolved by explicit out-of-domain detection using geometric thresholds.
- **ML–SDLC mismatch** is addressed by treating specialists as versioned, swappable components aligned with software lifecycles.
- **Unsafe scaling of knowledge** is mitigated by scaling expertise through selective activation rather than monolithic expansion.

Summary

GRCLM is not merely an optimization of Mixture-of-Experts architectures, but a structural rethinking of how knowledge is represented, routed, and evolved in language models. By grounding routing and learning decisions in geometry, GRCLM introduces interpretability, stability, and uncertainty awareness into systems that have traditionally relied on opaque and brittle mechanisms.

This makes GRCLM particularly well-suited for long-lived, safety-critical, and continually evolving AI systems.