

Topological AI: A Stable Architecture for High-Performance, Cost-Efficient Artificial Intelligence

(Discussion Version for ArXiv)

Author: Maurice “Tony” Ewing (BlackCopernicus)

Affiliation: Conquer Risk, TonyEwing.com

DOI (corresponding version): [10.5281/zenodo.15345661](https://doi.org/10.5281/zenodo.15345661)

Abstract

We introduce topAI, a novel AI architecture built for recursive stability, interpretability, and compute efficiency. Unlike transformers and statistical models that suffer from drift, recursive degradation, and costly retraining, topAI operates through a radically different structure. It embeds topological memory, totemic constraint, and bounded optimization into the architecture itself. We formally define a projection-based update rule that stabilizes learning over time, and we show how topAI enables inference and adaptation in feedback-sensitive environments like medicine and finance. topAI is not a language model—it is a control-theoretic scaffold for building safe, localized, and durable AI systems.

1. Introduction

Most modern AI architectures—including LLMs like GPT-4 and MoE systems like DeepSeek—are optimized for scale and parameter growth. While powerful in static benchmark settings, they collapse under feedback, drift, or recursive use, displaying challenges such as hallucinations, memory loss and collapse. The problem in these instances is not performance—as these are excellent models—it is fragility. Our findings suggest that this fragility stems from structural dependencies in how such models generalize across time and feedback.

In real-world systems such as hospital diagnostics, edge monitoring, portfolio investing or economic forecasting, recursive instability leads to hallucination, catastrophic forgetting, and opaque decision-making. topAI directly addresses these vulnerabilities by introducing a structurally grounded update regime.

2. Formal Structure and Update Logic

Let system dynamics evolve as:

$$x_{t+1} = f(x_t, u_t)$$

where $x_t \in \mathbb{R}^n$ is the system state and $u_t \in \mathbb{R}^m$ is the control input or observation. Unlike traditional models that optimize f globally, topAI restricts all updates via:

$$x_{t+1} = \Pi_{\mathcal{M}}(f(x_t, u_t)), \quad \text{subject to} \quad \|x_{t+1} - \tau\| \leq \epsilon$$

Here:

- \mathcal{M} : Topological memory (a preserved manifold of valid system states)
- $\Pi_{\mathcal{M}}$: Projection operator enforcing structure
- τ : Totemic anchor (reference state, or “truth core”)
- ϵ : Permissible deviation (volatility constraint)

Only updates that conform to both structure and bounded divergence are accepted.

3. Stability Function and Theoretical Basis

We define a Lyapunov-style functional:

$$V(x) = \|x - \tau\|^2$$

topAI permits updates only if $\Delta V(x) \leq 0$, ensuring stability. This projection-plus-contraction framework prevents recursive divergence—a problem that undermines most current architectures.

By embedding control-theoretic logic into the update operator itself, topAI achieves stability without patching, fine-tuning, or guardrails.

DOI (corresponding version): 10.5281/zenodo.15345661

4. Comparative Evaluation: Statistical Baselines vs. Generative Architectures

4.1 Overview and Setup

To position topAI credibly within the modern AI ecosystem, we compare its performance and system behavior against two canonical reference classes:

1. Classical statistical learning models:
 - Logistic Regression (LR)
 - Random Forest (RF)
2. Generative architectures used in large-scale AI:
 - Transformer-based LLMs (e.g., GPT-4)
 - Mixture-of-Experts systems (e.g., DeepSeek)

These models differ sharply in scale, philosophy, and computational requirements. The aim of this section is not to equate their core use cases, but to assess how each performs under recursive update and real-world constraints.

We analyze topAI along three axes:

- Accuracy and robustness under shift,
- Update stability over time,
- Resource cost for training and deployment.

The baselines are established using a synthetic dataset modeled on cardiac arrest vitals (Section 2.1). This ensures:

- Identical feature space for all models,
- Controlled batch updates to simulate distribution shift,
- A testbed to evaluate recursive learning behavior and memory degradation.

DOI (corresponding version): 10.5281/zenodo.15345661

4.2 Experimental Design and Assumptions

Dataset

- Simulated cardiac arrest vitals: heart rate, SpO₂, diastolic/systolic BP, respiratory rate, temperature, and consciousness level (GCS).
- 150 synthetic patients created via a bounded multivariate sampling process.
- No real patient data used — the goal is behavioral modeling under shift, not medical accuracy.

Training Protocol

Each model is trained on an initial seed set (50 samples), then updated over 10 sequential batches of 10 new patients each (non-overlapping).

This mimics recursive update under data drift—a critical scenario for real-world AI systems (especially in healthcare, trading, and embedded analytics).

Evaluation Metrics

- Accuracy after each batch
- AUC (Area Under ROC Curve) per batch
- Update latency per batch (mean over 10 rounds)
- Observed drift (qualitative notes on divergence, retraining needs, or error spikes)

Compute Environment

- Statistical models and topAI were executed on a MacBook Pro (M1, 16 GB RAM).
- LLM cost data are based on publicly available infrastructure benchmarks (OpenAI, DeepSeek, MosaicML, and HuggingFace reports).

4.3 Performance Results

Model	Final Accuracy	Avg AUC (10 rounds)	Avg Update Time	Observed Drift
Logistic Regression	0.96	0.94	~0.09 sec	Moderate under shift
Random Forest	0.98	0.96	~0.74 sec	High volatility post-update
topAI	0.97	0.95	~0.008 sec	Minimal—stable fit

DOI (corresponding version): 10.5281/zenodo.15345661

Interpretation:

- topAI performs on par with classical models on accuracy and AUC.
- It dramatically outperforms them on update cost and stability.
- Random Forest exhibits rapid variance with small updates—a known behavior under shift-prone regimes.
- Logistic Regression is more stable but lacks nonlinear adaptability.

topAI’s structural projection method—which maps updates within a constrained topological manifold—enables fast assimilation with little drift or volatility.

4.4 Cost and Strategic Comparison to LLMs / MoEs

We now compare topAI to large-scale generative architectures—not in language modeling power, but in deployment viability under constraint.

Estimated Training Cost (per 100M tokens or structured samples)

Model	Est. Training Cost	Hardware	Deployment Update Strategy
GPT-4 / DeepSeek	\$100K - \$500K+	Multi-GPU clusters	Frequent, full or partial re-training
topAI	<\$5K (projected)	CPU, edge-ready	Bounded projection update

Metric	GPT / DeepSeek	topAI
Energy cost (per update)	Very high (GPU-bound)	Low (no GPU)
FLOPs per round	10B+	~1 - 5% of that
Data demand	Massive (millions)	Low (100~500 samples)
Retrain frequency	Weekly or monthly	Rare, only if topology shifts
Risk of hallucination	High under recursion	None observed

DOI (corresponding version): 10.5281/zenodo.15345661

4.5 Strategic Interpretation

topAI is not a compression of GPT—it's a structural inversion. It offers an architecture that is:

- Local rather than global in optimization
- Topologically constrained rather than open-ended
- Update-stable rather than retrain-dependent

This makes topAI ideal for domains where:

- Stability is non-negotiable (medical devices, autonomous systems, control loops)
- Data are sparse or recursive (clinical trials, post-op monitoring, real-time intelligence)
- Regulatory “auditability” is required (finance, public infrastructure, defense AI)

Even if it is not designed to outperform GPT on text or vision, it can outlast and out-economize it in tightly constrained environments.

4. Strategic Positioning

topAI is designed for recursive environments where hallucination and drift cause failure. It is not a replacement for GPT, but a structural alternative suited for tasks like:

- Medical inference under shift
- Edge deployments (wearables, drones, implants)
- Hallucination suppression for LLM pipelines
- Systems requiring localized, interpretable logic

It uses orders of magnitude fewer FLOPs, allows local learning, and eliminates retraining cascades.

5. Comparison with Transformer Architectures

Feature	GPT / DeepSeek	topAI
Update Cost	High	Minimal
Stability Under Feedback	Low	High
Interpretability	Low	High
Deployment Flexibility	Cloud / GPU	CPU / Edge
Hallucination Risk	High	None (by design)

6. Outlook

topAI is the first AI architecture explicitly designed for recursive stability. It is not just “efficient”—it is strategically survivable. Future work includes its application to regulatory AI, immune system modeling, and hybrid transformer integration.

This discussion paper introduces the theoretical backbone. Full implementation details are reserved for future versions.

DOI (corresponding version): 10.5281/zenodo.15345661

References

1. OpenAI. (2023). GPT-4 Technical Report. arXiv:2303.08774. <https://arxiv.org/abs/2303.08774>
2. DeepSeek-V2. (2024). DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model. arXiv:2405.04434v5 [cs.CL].
3. Brown, T., et al. (2020). Language Models are Few-Shot Learners. arXiv:2005.14165
4. Wynants, L., et al. (2020). Prediction Models for COVID-19 Diagnosis and Prognosis. *BMJ*, 369:m1328.
5. Kung, T. H., et al. (2023). Performance of ChatGPT on USMLE: Potential for Medical Education Use. *PLOS Digital Health*, 2(2): e0000198.
6. Banerjee, I., et al. (2021). Racial Bias in AI Models Used in Clinical Settings. *JAMA*, 326(7), 644–651.
7. Sahiner, B. (2023) Data drift in medical machine learning: implications and potential remedies *The British journal of radiology* 96(1150):20220878
8. Senkaiahliyan, S, et al. (2023). GPT-4V(ision) Unsuitable for Clinical Care and Education: A Clinician-Evaluated Assessment arXiv:2403.12046