

Affect-Weighted Gossip-Based Memory Architecture for Preventing Catastrophic Forgetting in Distributed AI Systems

Anon¹

¹Institute for Distributed Cognition, Synthetic Intelligence Research Lab

March 2025

Abstract

Catastrophic forgetting plagues distributed AI systems, where new tasks overwrite prior knowledge. We propose an affect-weighted gossip-based memory architecture over a hybrid Ramanujan-hypercube topology to mitigate this. Drawing from biological cognition and social learning, our model prioritizes retention via human emotional feedback and ensures scalability through distributed synchronization. Memory dynamics, modeled as a nonlinear, affect-biased system, converge to a non-zero equilibrium, with generalization capacity optimized geometrically. Simulations with 10,000 agents achieve $\text{MRR}=0.87$, outperforming Elastic Weight Consolidation (EWC) at 0.67, and highlight efficiency gaps using geometric generalization indicators such as C_{gen} and G_{Ricci} . This biologically plausible, scalable solution advances lifelong learning in AI.

1 Introduction

Catastrophic forgetting in neural networks, notably large language models (LLMs), arises from task-specific optimization overwriting prior knowledge [1]. Methods like EWC [1] falter in scalability, unlike biological memory’s holistic integration [2]. Our architecture leverages a Ramanujan-hypercube topology [4] and emotional reinforcement [5] for robust, generalizable retention.

2 Background

Catastrophic forgetting occurs when neural networks, trained sequentially on new tasks, lose performance on earlier ones due to weight updates overwriting prior knowledge [1]. This is acute in distributed AI systems, where agents

must adapt without centralized control. Traditional solutions like EWC impose regularization to preserve key weights but scale poorly in large networks.

Gossip protocols offer a distributed alternative, enabling agents to share information via local interactions, achieving global consensus with logarithmic complexity [4]. In biological systems, emotional salience enhances memory retention [5], suggesting a hybrid approach. Peat’s holistic memory theory [2] further inspires integrating sensory and emotional cues, missing in current AI models.

3 Model Architecture

3.1 Network Setup

The system is a graph $\mathcal{G} = (V, E)$, with $V = \{v_1, \dots, v_n\}$ in $k = \lceil n/100 \rceil$ clusters of ≈ 100 agents, each a d -regular Ramanujan graph (spectral gap $\geq 2\sqrt{d-1}$), linked by a hypercube $Q_{\lceil \log_2 k \rceil}$. Memory vectors are:

$$M_i^t = [\delta_i^t(x_1), \dots, \delta_i^t(x_m)], \quad \delta_i^t(x_j) \in [0, 1].$$

Degree $d = \min(\lfloor 5 \cdot \max_i \alpha_i^t + 5 \rfloor, 10)$.

3.2 Memory Dynamics

Updates are:

$$\delta_i^{t+1}(x_j) = \lambda \delta_i^t(x_j) + \eta \alpha_i^t(x_j) E_i^t(x_j) + \gamma \text{Gossip}_i^t(x_j), \quad (1)$$

with $\lambda = 0.9$, $\eta = 0.2$, $\gamma = 0.3$.

3.3 Emotional Trace Calculation

$$E_i^t(x_j) = \sum_{h \in H_i} \text{sigmoid}(\text{valence}_h(x_j) \cdot \text{arousal}_h(x_j)).$$

3.4 Gossip Mechanism and Lie Hypothesis

Each agent maintains a verified memory state $T_i^t(x_j)$ and a provisional state $F_i^t(x_j)$. The gossip update is:

$$\begin{aligned} \text{Gossip}_i^t(x_j) = & \sum_{k \in N_c(i)} w_{ik}(x_j) [\alpha_k^t(x_j) T_k^t(x_j) + (1 - \alpha_k^t(x_j)) F_k^t(x_j)] \\ & + \sum_{k \in N_h(i)} w_{ik}(x_j) [\alpha_k^t(x_j) T_k^t(x_j) + (1 - \alpha_k^t(x_j)) F_k^t(x_j)], \end{aligned} \quad (2)$$

where:

- $\alpha_k^t(x_j) = \sigma(2E_k^t(x_j))$ is the confidence weight,
- $\lambda_T = 0.95$ and $\lambda_F = 0.85$ are decay rates,
- $\gamma_T = 0.4$ and $\gamma_F = 0.2$ are gossip strengths.

4 Proof of Theorem 1

Theorem 1: Under \mathcal{G} , $\lambda < 1$, and nonzero E_i^t , the average memory strength $\bar{\delta}^t(x_j)$ converges to a non-zero equilibrium, scalable and generalizable.

For:

$$\delta_i^{t+1}(x_j) = \lambda \delta_i^t(x_j) + \eta \alpha_i^t(x_j) E_i^t(x_j) + \gamma \cdot \text{Gossip}_i^t(x_j),$$

define:

$$\bar{\delta}^t(x_j) = \frac{1}{n} \sum_{i=1}^n \delta_i^t(x_j).$$

4.1 Global Memory Convergence

The topology uses k Ramanujan clusters and a hypercube, ensuring mixing in $O(\log 100)$ locally and $O(\log k)$ globally. Taking the expectation:

$$\mathbb{E}[\bar{\delta}^{t+1}(x_j)] = \lambda \mathbb{E}[\bar{\delta}^t(x_j)] + \eta \cdot \bar{\alpha}^t(x_j) \cdot \bar{E}^t(x_j),$$

the fixed point is:

$$\bar{\delta}^*(x_j) = \frac{\eta \bar{\alpha} \bar{E}}{1 - \lambda},$$

stable and non-zero for $\lambda = 0.9 < 1$ and $\bar{E} > 0$.

4.2 Ricci Curvature Stability

Ramanujan graphs yield a spectral gap $\geq 2\sqrt{d-1}$; Ollivier-Ricci curvature $\mathcal{R} \approx 0.37$ at $n = 10,000$ supports $\bar{\delta} > 0.72$ for salient items, aligning with MRR=0.87 and 7-round convergence.

4.3 Conflict and Lie Quarantine Resolution

For divergent α_i^t (e.g., $\text{stddev}(\alpha) > 0.2$), quarantine isolates x_j . Unfreeze if:

$$\sum_{k \in N(i)} w_{ik}(x_j) \cdot I_{\text{verified}}(x_j) > \theta = 0.5, \quad I_{\text{verified}}(x_j) = \begin{cases} 1 & \text{if } \alpha_k^t(x_j) > 0.5, \\ 0 & \text{otherwise.} \end{cases}$$

4.4 Reload Condition for Human-Machine Time

$$\text{Reload}_i^t(x_j) = 0.1(t - t_{\text{last}}) + 0.5 E_i^t(x_j) e^{-0.1(t - t_{\text{last}})} \cdot \frac{1}{|N(i)|} \sum_{k \in N(i)} \alpha_k^t(x_j),$$

triggers at $\text{Reload}_i^t > 0.7$.

4.5 Generalization Metrics

$$C_{\text{gen}}(x_j) = \frac{1}{n(n-1)} \sum_{i \neq k} W_2(\delta_i^t(x_j), \delta_k^t(x_j)),$$

$$G_{\text{Ricci}}(x_j) = \frac{1}{|E|} \sum_{(i,k) \in E} \kappa_{ik}(x_j),$$

optimal for $G_{\text{Ricci}} \in [0.25, 0.45]$, $C_{\text{gen}} \leq 0.3$.

5 Simulation Results

Simulations with 10,000 agents yielded MRR=0.87, compared to Independent Learning (IL) at 0.45 and EWC at 0.67, with 90% accuracy and convergence in 7 rounds ($\mathcal{R} \approx 0.37$). Efficiency versus EWC was assessed using:

- **EGM (Efficiency Gap Metric):** Peaks at 0.3 early, indicating 30% better retention due to rapid gossip and emotional reinforcement, then stabilizes at 0.05 as generalization aligns both models. Initial spikes (0.3 at $t \approx 5$) are expected due to localized reinforcement from salient feedback, but convergence to 0.05 by $t = 20$ demonstrates effective correction via gossip averaging and temporal damping.
- **CD (Convergence Distance):** Stable at 0.15, showing moderate divergence in memory states, with gossip prioritizing salient items over EWC’s uniform approach.
- **EWGS (Emotional-Weighted Gain Score):** Consistent at 0.22, reflecting a 22% retention advantage for emotionally salient items ($E_i^t > 0.5$).

Generalization comparison highlights:

- **Gossip Pros:** Superior scalability (MRR=0.87 vs. 0.67) via rapid mixing and $\mathcal{R} \approx 0.37$, emotional prioritization ($C_{\text{gen}} \leq 0.3$, EWGS=0.22), and dynamic adaptation (reload mechanism).
- **Gossip Cons:** Higher complexity, early instability (EGM=0.3), and parameter sensitivity (G_{Ricci} tuning).
- **EWC Pros:** Simpler implementation, stable generalization (MRR=0.67), and low overhead.
- **EWC Cons:** Limited scalability, no emotional bias (EWGS=0), and rigidity to task shifts.

Gossip outperforms in scalability and emotional retention, but at the cost of complexity and tuning, while EWC offers simplicity and stability, sacrificing adaptability and distributed generalization ($\delta > 0.72$ for gossip).

6 Discussion

The model enhances scalability and generalization for LLMs and robotics, aligning with holistic memory theories [2]. Below, we address key challenges identified in simulations.

6.1 Complexity Consideration

While the hybrid topology (Ramanujan clusters and hypercube overlay) and gossip-based updates introduce higher computational overhead than EWC’s centralized regularization, the architecture benefits from bounded-degree graphs ($d \leq 10$) and logarithmic gossip depth ($O(\log k)$), enabling practical scalability. Unlike EWC, which requires full access to the Fisher Information Matrix—scaling poorly with model size and infeasible in federated contexts—our model shifts costs to communication and distributed processing. These are amortized through parallel agent updates across $k = 100$ clusters for $n = 10,000$ and modular design, with future compression strategies potentially reducing bandwidth further.

6.2 Early Stability and Adaptation

Initial EGM spikes (e.g., 0.3 at $t \approx 5$) stem from rapid overreaction to emotionally salient inputs before full network synchronization, akin to biological cognition’s early affective bias [5]. The model self-corrects via temporal decay ($\lambda_F = 0.85 < \lambda_T = 0.95$), the reload function balancing emotional trace staleness and group support, and a quarantine-unfreeze cycle for divergent agents. By $t = 20$, EGM stabilizes at 0.05, reflecting effective gossip averaging and contextual adjustment.

6.3 Parameter Robustness and Ricci Generalization

While curvature bounds ($G_{\text{Ricci}} \in [0.25, 0.45]$) suggest parameter sensitivity compared to EWC’s static penalty, G_{Ricci} is not manually tuned but measured from the graph’s geometry at runtime. It serves a diagnostic role—detecting knowledge drift, guiding reloads, and evaluating emergent generalization—rather than requiring precise control. In contrast to EWC’s fixed Fisher-based penalties, this offers adaptive, topology-level insight into distributed memory health, validated by $\mathcal{R} \approx 0.37$.

7 Future Work

Test Ricci flow, dialogue systems, and generalization via C_{gen} and G_{Ricci} at 50,000 agents (X, March 2025). Dynamic curvature tracking may eventually replace static regularization with topological adaptation, leveraging bounds like $G_{\text{Ricci}} \in [0.25, 0.45]$ for monitoring rather than control.

7.1 Hybrid Gossip-EWC Model

A hybrid could combine gossip’s scalability and emotional retention with EWC’s simplicity and stability. Augment updates with an EWC penalty:

$$\delta_i^{t+1}(x_j) = \lambda \delta_i^t(x_j) + \eta \alpha_i^t(x_j) E_i^t(x_j) + \gamma \text{Gossip}_i^t(x_j) - \zeta \sum_k F_k (\delta_i^t(x_j) - \delta_i^{t_k})^2,$$

where $\zeta = 0.1$ and F_k is the Fisher information for task k . This retains topology and gossip, adding regularization to reduce complexity and stabilize early training (e.g., lower EGM peaks), potentially achieving $\text{MRR} \approx 0.80$.

References

- [1] Kirkpatrick, J., et al. "Overcoming catastrophic forgetting in neural networks." *Proceedings of the National Academy of Sciences*, 2017.
- [2] Peat, R. "A Holistic Physiology of Memory." Blake College, Eugene, Oregon, U.S.A., 1975.
- [3] Anthropic. "When does pretraining verifiably prevent lying in LLMs?" 2024.
- [4] "Gossip Protocol Explained." *High Scalability*, 2024. <https://highscalability.com>.
- [5] "The Influences of Emotion on Learning and Memory." *PMC*, 2024.
- [6] Ollivier, Y. "Ricci curvature of Markov chains on metric spaces." *Journal of Functional Analysis*, 2010.