

E79: Coupled Memory-Modulation Matrix System

A Mathematical Analysis of Hierarchical Delta Rules

Formal verification in Lean 4 with Mathlib

1 Introduction

E79 represents the culmination of 79 architectural experiments in recurrent neural network design. Its key innovation is **coupled delta rules**: two $n \times n$ matrix states where the second learns to predict the residuals of the first.

This document provides:

1. Complete mathematical specification of E79
2. Analysis of how M modulates S
3. Jacobian and gradient flow analysis
4. Generalizations to K-level hierarchies
5. Conditions for simplification

2 Mathematical Specification

2.1 State Definition

E79 maintains two matrix states:

$$S \in \mathbb{R}^{n \times n} \quad \text{Content Memory (primary associative storage)} \tag{1}$$

$$M \in \mathbb{R}^{n \times n} \quad \text{Modulation Memory (controls S's gating)} \tag{2}$$

Total state: $2n^2$ real values. For $n = 32$, this is 2048 elements.

2.2 Input Vectors

At each timestep, E79 receives:

- $k \in \mathbb{R}^n$: Key vector for content addressing
- $v \in \mathbb{R}^n$: Value to store
- $q \in \mathbb{R}^n$: Query for output
- $m \in \mathbb{R}^n$: Modulation key for M addressing

2.3 The E79 Update Rule (Actual Implementation)

Input: State (S, M) , vectors (k, v, q, m) , biases (b_S, b_M)

Step 1: Normalize keys

$$\hat{k} = \frac{k}{\|k\|_2}, \quad \hat{m} = \frac{m}{\|m\|_2} \quad (3)$$

Step 2: M controls S's decay gates ($M \rightarrow S$ coupling)

$$g_{\text{row}}^S = \sigma(M\hat{k} + b_S) \quad (\text{row decay from } M) \quad (4)$$

$$g_{\text{col}}^S = \sigma(M^\top \hat{k} + b_S) \quad (\text{col decay from } M) \quad (5)$$

Step 3: S delta rule update with M-controlled gating

$$\delta_S = v - S\hat{k} \quad (6)$$

$$S' = \underbrace{(g_{\text{row}}^S g_{\text{col}\{\cdot\}^\top})}_{M\text{-controlled decay}} \odot S + \delta_S \hat{k}^\top \quad (7)$$

Step 4: S controls M's decay gates ($S \rightarrow M$ coupling)

$$g_{\text{row}}^M = \sigma(S\hat{m} + b_M) \quad (8)$$

$$g_{\text{col}}^M = \sigma(S^\top \hat{m} + b_M) \quad (9)$$

Step 5: M delta rule update (M predicts S's changes)

$$\delta_M = \delta_S - M\hat{m} \quad (10)$$

$$M' = \underbrace{(g_{\text{row}}^M g_{\text{col}\{\cdot\}^\top})}_{M\text{-predicted changes}} \odot M + \delta_M \hat{m}^\top \quad (11)$$

Step 6: Output with self-gating

$$o = (S'q) \odot \text{silu}(S'q) \quad (12)$$

Return: New state (S', M') , output o

Algorithm 1: E79 Forward Pass - Mutual Gating Control

2.4 Explicit Matrix Form

The key insight is the **factorized gating**:

$$S' = \underbrace{(g_{\text{row}}^S g_{\text{col}\{\cdot\}^\top})}_{M\text{-controlled decay}} \odot S + \delta_S \hat{k}^\top \quad (13)$$

Where the decay gate is an outer product of M's outputs:

$$g_{\text{row}}^S g_{\text{col}\{\cdot\}^\top} = \sigma(M\hat{k})\sigma(M^\top \hat{k})^\top \quad (14)$$

This means **M directly influences S's update** and thus the output. Similarly:

$$M' = \underbrace{(g_{\text{row}}^M g_{\text{col}\{\cdot\}^\top})}_{S\text{-predicted changes}} \odot M + \delta_M \hat{m}^\top \quad (15)$$

With S controlling M's decay gates.

3 How M Modulates S

3.1 The Coupling Mechanism (Actual Implementation)

Unlike the simplified description in E79_RESULTS.md, the **actual** E79 implements **mutual gating control**:

Key Insight: M directly controls S's decay gates, and S controls M's decay gates. This creates a bidirectional dynamical coupling where each memory controls what the other forgets.

3.1.1 M → S Coupling (M controls S's forgetting)

The decay factors for S come from M:

$$\mathbf{g}_{\text{row}}^S = \sigma(M\hat{\mathbf{k}} + \mathbf{b}_S) \in (0, 1)^n \quad (16)$$

$$\mathbf{g}_{\text{col}}^S = \sigma(M^\top \hat{\mathbf{k}} + \mathbf{b}_S) \in (0, 1)^n \quad (17)$$

The S update becomes:

$$S'_{ij} = g_{\text{row},i}^S \cdot g_{\text{col},j}^S \cdot S_{ij} + (\delta_S)_i \hat{k}_j \quad (18)$$

M controls what S retains. When $M\hat{\mathbf{k}}$ is large and positive, $\mathbf{g}_{\text{row}}^S \rightarrow 1$ and S preserves its rows. When negative, S forgets.

3.1.2 S → M Coupling (S controls M's forgetting)

Symmetrically, S controls M's decay:

$$\mathbf{g}_{\text{row}}^M = \sigma(S\hat{\mathbf{m}} + \mathbf{b}_M) \quad (19)$$

$$\mathbf{g}_{\text{col}}^M = \sigma(S^\top \hat{\mathbf{m}} + \mathbf{b}_M) \quad (20)$$

S controls what M retains. This creates a feedback loop where the memories regulate each other.

3.2 Gradient Flow Through M

Theorem (M Gets Gradients Through S). M influences the output through the gating path:

$$\text{Loss} \rightarrow \mathbf{o} \rightarrow S' \rightarrow \mathbf{g}_{\text{row}}^S, \mathbf{g}_{\text{col}}^S \rightarrow M \quad (21)$$

Specifically:

$$\frac{\partial \mathcal{L}}{\partial M} = \frac{\partial \mathcal{L}}{\partial S'} \cdot \frac{\partial S'}{\partial g^S} \cdot \frac{\partial g^S}{\partial M} \quad (22)$$

Proof. From Equation 7: $S' = (\mathbf{g}_{\text{row}}^S \mathbf{g}_{\text{col}}^S)^\top \odot S + \delta_S \hat{\mathbf{k}}^\top$

The gradient of S' with respect to $\mathbf{g}_{\text{row}}^S$ is:

$$\frac{\partial S'_{ij}}{\partial g_{\text{row},i}^S} = g_{\text{col},j}^S \cdot S_{ij} \quad (23)$$

And $\mathbf{g}_{\text{row}}^S = \sigma(M\hat{\mathbf{k}} + \mathbf{b}_S)$, so:

$$\frac{\partial g_{\text{row},i}^S}{\partial M_{il}} = \sigma'(\dots) \cdot \hat{k}_l \quad (24)$$

Composing these gives a non-zero gradient path from Loss to M. \square

3.3 Interpretation: Mutual Control Dynamical System

The E79 coupling creates a **self-organizing** memory system:

Aspect	Mechanism
$M \rightarrow S$	M decides what S should forget based on current key
$S \rightarrow M$	S decides what M should forget based on modulation key
S delta	Standard delta rule with M-modulated decay
M delta	Learns S's prediction errors for meta-learning

This is analogous to:

- **Neural gating** (LSTM): Forget gates control information flow
- **Attention** (Transformers): Context-dependent routing
- **Neuromodulation** (Biology): One system modulates another's plasticity

3.4 Why Mutual Control Helps

Proposition (Adaptive Forgetting). With M-controlled gating, S can learn **input-dependent forgetting**:

- For familiar keys: M outputs high gates \rightarrow S preserves old information
- For novel keys: M outputs low gates \rightarrow S makes room for new content

This is impossible with fixed decay α_S .

Proposition (Meta-Learning Through Coupling). M can learn to recognize **when** S should update strongly vs. weakly:

- Systematic input patterns \rightarrow M learns predictable gating
- Noisy inputs \rightarrow M learns to gate conservatively

This is a form of “learning to learn” for associative memory.

4 Jacobian Analysis

4.1 State Space Jacobian

The full E79 state is $z = \text{vec}([S; M]) \in \mathbb{R}^{2n^2}$.

Theorem (Lower-Triangular Jacobian). The Jacobian of the E79 update has block lower-triangular structure:

$$\frac{\partial z'}{\partial z} = \begin{pmatrix} J_S & 0 \\ J_{MS} & J_M \end{pmatrix} \quad (25)$$

where:

- $J_S = \frac{\partial S'}{\partial S}$: How S affects S'
- $J_M = \frac{\partial M'}{\partial M}$: How M affects M'
- $J_{MS} = \frac{\partial M'}{\partial S}$: How S affects M' (the coupling!)

- $\frac{\partial \mathbf{S}'}{\partial \mathbf{M}} = \mathbf{0}$: \mathbf{M} does not affect \mathbf{S}' directly

Proof. From Equation 7, \mathbf{S}' depends only on \mathbf{S} , not \mathbf{M} :

$$\mathbf{S}' = \alpha_S \mathbf{S} + (\mathbf{v} - \mathbf{S}\hat{\mathbf{k}})\hat{\mathbf{k}}^\top \quad (26)$$

Therefore $\frac{\partial \mathbf{S}'}{\partial \mathbf{M}} = \mathbf{0}$.

From Equation 11, \mathbf{M}' depends on both \mathbf{S} and \mathbf{M} :

$$\mathbf{M}' = \alpha_M \mathbf{M} + \left(\underbrace{\mathbf{v} - \mathbf{S}\hat{\mathbf{k}}}_{\delta_S} - \mathbf{M}\hat{\mathbf{m}} \right) \hat{\mathbf{m}}^\top \quad (27)$$

The dependence on \mathbf{S} comes through δ_S :

$$\frac{\partial \mathbf{M}'}{\partial \mathbf{S}} = \frac{\partial}{\partial \mathbf{S}} [(\mathbf{v} - \mathbf{S}\hat{\mathbf{k}} - \mathbf{M}\hat{\mathbf{m}})\hat{\mathbf{m}}^\top] = -\hat{\mathbf{k}}\hat{\mathbf{m}}^\top \quad (28)$$

(in the appropriate tensor form). \square

4.2 Individual Block Jacobians

4.2.1 Content Memory Jacobian \mathbf{J}_S

For the delta rule update with decay:

$$\mathbf{S}' = \alpha_S \mathbf{S} + (\mathbf{v} - \mathbf{S}\hat{\mathbf{k}})\hat{\mathbf{k}}^\top = \alpha_S \mathbf{S} + \mathbf{v}\hat{\mathbf{k}}^\top - \mathbf{S}\hat{\mathbf{k}}\hat{\mathbf{k}}^\top \quad (29)$$

Taking the derivative with respect to \mathbf{S} :

$$\mathbf{J}_S = \alpha_S \mathbf{I} - \hat{\mathbf{k}}\hat{\mathbf{k}}^\top \otimes \mathbf{I}_n \quad (30)$$

In terms of action on a perturbation $\delta \mathbf{S}$:

$$\delta \mathbf{S}' = \alpha_S \delta \mathbf{S} - (\delta \mathbf{S}\hat{\mathbf{k}})\hat{\mathbf{k}}^\top \quad (31)$$

Theorem (S Jacobian Spectral Properties). With $\|\hat{\mathbf{k}}\|_2 = 1$, the Jacobian \mathbf{J}_S (as a linear map on $n \times n$ matrices) has eigenvalues:

- α_S with multiplicity $n^2 - n$ (eigenvectors: matrices with $\hat{\mathbf{k}}$ in null space)
- $\alpha_S - 1$ with multiplicity n (eigenvectors: outer products with $\hat{\mathbf{k}}$)

For stability, we need $|\alpha_S| \leq 1$ and $|\alpha_S - 1| \leq 1$, giving $\alpha_S \in [0, 1]$.

4.2.2 Modulation Memory Jacobian \mathbf{J}_M

Similarly:

$$\mathbf{J}_M = \alpha_M \mathbf{I} - \hat{\mathbf{m}}\hat{\mathbf{m}}^\top \otimes \mathbf{I}_n \quad (32)$$

Same spectral structure with $\hat{\mathbf{m}}$ replacing $\hat{\mathbf{k}}$.

4.3 Gradient Flow Through the Coupling

The coupling term \mathbf{J}_{MS} enables **gradient sharing**:

$$\frac{\partial \mathcal{L}}{\partial S} = \frac{\partial \mathcal{L}}{\partial S'} J_S + \frac{\partial \mathcal{L}}{\partial M'} J_{MS} \quad (33)$$

The second term means: **M's gradient signal flows back to S.**

Corollary. When training E79 end-to-end, S receives gradients from:

1. Direct output path: $\mathbf{o} \rightarrow S'$
2. Indirect coupling path: $\mathbf{o} \rightarrow M' \rightarrow S$ (through δ_S)

This coupling allows M to “tell” S about systematic errors.

5 Exact Retrieval and Capacity

5.1 Single Write Exact Retrieval

Theorem (Exact Retrieval). If $S = \mathbf{0}$ (empty memory) and we write (\mathbf{v}, \mathbf{k}) with $\|\mathbf{k}\| = 1$, then:

$$S' \hat{\mathbf{k}} = \mathbf{v} \quad (34)$$

Proof.

$$\begin{aligned} S' \hat{\mathbf{k}} &= [\mathbf{0} + (\mathbf{v} - \mathbf{0} \cdot \hat{\mathbf{k}}) \hat{\mathbf{k}}^\top] \hat{\mathbf{k}} \\ &= (\mathbf{v}) \hat{\mathbf{k}}^\top \hat{\mathbf{k}} \\ &= \mathbf{v}(\hat{\mathbf{k}}^\top \hat{\mathbf{k}}) \\ &= \mathbf{v} \cdot 1 = \mathbf{v} \end{aligned} \quad (35)$$

5.2 Orthogonal Keys Preserve Information

Theorem (Selective Update). If $\hat{\mathbf{k}}_1 \perp \hat{\mathbf{k}}_2$ (orthogonal keys), then writing $(\mathbf{v}_2, \mathbf{k}_2)$ does not affect retrieval with \mathbf{k}_1 :

$$S' \hat{\mathbf{k}}_1 = S \hat{\mathbf{k}}_1 \quad (36)$$

Proof. The update adds $(\mathbf{v}_2 - S \hat{\mathbf{k}}_2) \hat{\mathbf{k}}_2^\top$. Applying to $\hat{\mathbf{k}}_1$:

$$[(\mathbf{v}_2 - S \hat{\mathbf{k}}_2) \hat{\mathbf{k}}_2^\top] \hat{\mathbf{k}}_1 = (\mathbf{v}_2 - S \hat{\mathbf{k}}_2)(\hat{\mathbf{k}}_2^\top \hat{\mathbf{k}}_1) = \mathbf{0} \quad (37)$$

since $\hat{\mathbf{k}}_2^\top \hat{\mathbf{k}}_1 = 0$.

5.3 Capacity Analysis

Proposition (E79 Capacity). With orthonormal keys $\{\hat{\mathbf{k}}_i\}_{i=1}^n$ for S and $\{\hat{\mathbf{m}}_j\}_{j=1}^n$ for M:

- S can store n independent (value, key) pairs: n^2 real values
- M can store n independent (residual, modulation-key) pairs: n^2 real values
- Total: $2n^2$ real values (equal to state size)

This is **optimal** capacity utilization.

6 Generalizations

6.1 K-Level Hierarchies

E79 is the $K = 2$ case of a general construction:

Definition (K-Level Coupled Memory). For $K \geq 1$, define matrices $\mathbf{M}_0, \mathbf{M}_1, \dots, \mathbf{M}_{K-1} \in \mathbb{R}^{n \times n}$ with:

$$\mathbf{r}_0 = \mathbf{v} - \mathbf{M}_0 \hat{\mathbf{k}}_0 \quad (\text{Level 0 residual}) \quad (38)$$

$$\mathbf{r}_i = \mathbf{r}_{i-1} - \mathbf{M}_i \hat{\mathbf{k}}_i \quad \text{for } i = 1, \dots, K-1 \quad (39)$$

$$\mathbf{M}'_i = \alpha_i \mathbf{M}_i + \mathbf{r}_i \hat{\mathbf{k}}_i^\top \quad (\text{Each level's update}) \quad (40)$$

K	Description
1	Standard delta rule (E74). Single matrix S .
2	E79. Content memory S + Modulation memory M .
3	Triple hierarchy. $S + M + N$ where N predicts M 's residuals.
K	Chain of K residual predictors.

6.2 Diminishing Returns Conjecture

Conjecture. For a fixed compute budget, there exists an optimal K^* such that:

- $K < K^*$: Adding levels improves performance
- $K > K^*$: Additional levels have negligible benefit

K^* depends on:

1. **Task complexity**: How much structure exists in residuals
2. **Training time**: Deeper hierarchies need more convergence time
3. **State budget**: Each level costs n^2 parameters

The benchmark showing $n = 32$ optimal for 10-minute training suggests E79 is near optimal for that regime.

7 Simplifications

7.1 Tied Keys: $m = k$

Theorem (Tied Keys Reduction). If $m = k$ (same key for both levels), then E79 reduces to a single delta rule on the combined matrix $S + M$.

Proof. With $\hat{\mathbf{m}} = \hat{\mathbf{k}}$:

$$\delta_M = \mathbf{v} - S\hat{\mathbf{k}} - M\hat{\mathbf{k}} = \mathbf{v} - (S + M)\hat{\mathbf{k}} \quad (41)$$

Combined update:

$$S' + M' = \alpha_S S + \alpha_M M + (\mathbf{v} - S\hat{\mathbf{k}})\hat{\mathbf{k}}^\top + (\mathbf{v} - S\hat{\mathbf{k}} - M\hat{\mathbf{k}})\hat{\mathbf{k}}^\top \quad (42)$$

If $\alpha_S = \alpha_M = \alpha$:

$$S' + M' = \alpha(S + M) + [2v - 2S\hat{k} - M\hat{k}] \hat{k}^\top \quad (43)$$

This is **not** exactly a single delta rule, but the key insight is: tied keys limit M's ability to organize independently.

7.2 Zero M Decay: $\alpha_M = 0$

Proposition (Instantaneous Modulation). With $\alpha_M = 0$:

$$M' = \delta_M \hat{m}^\top \quad (44)$$

M becomes an “instantaneous” residual predictor with no memory of past residuals. This is useful when residuals have no temporal structure.

7.3 No Modulation: $M = 0$

Proposition (Reduction to E74). Setting $M = 0$ and $\alpha_M = 0$ recovers E74 (single delta rule with self-gating):

$$S' = \alpha_S S + (v - S\hat{k}) \hat{k}^\top \quad (45)$$

$$o = (S'q) \odot \text{silu}(S'q) \quad (46)$$

8 Empirical Results Summary

From the benchmark (100M params, 10-minute training):

Model	Loss	tok/s	State
Mamba2	1.27	78.7K	SSM (parallel)
E79 n=32	1.51	31.5K	$2 \times 32^2 = 2048$
E1 (gated)	1.53	45.5K	vector
E42 (linear)	1.59	137K	vector
FLA-GDN	1.99	18.7K	matrix

Key observations:

- E79 beats E1 (1.51 vs 1.53): modulation helps
- n=32 optimal for 10-min training (larger n under-converged)
- 40% of Mamba2 throughput despite sequential scan

9 Conclusions

9.1 What E79 Teaches Us

1. **Hierarchical error correction works:** Even one level of residual prediction (M on S) provides measurable benefit.
2. **Separate addressing enables specialization:** k for content, m for error patterns.
3. **Gradient coupling is essential:** M receives loss-relevant gradients through δ_S , not just residual reconstruction loss.
4. **Training time vs capacity tradeoff:** Larger state needs more training to converge.

9.2 Open Questions

1. What is optimal K for K-level hierarchies?
2. Can we learn the coupling adaptively?
3. How does E79 scale beyond 100M parameters?
4. Can parallel scan be applied to coupled matrices?