

Section 4: E88 Temporal Nonlinearity

Tanh Saturation, Latching, and Expressivity Separation

4.1 Overview

This section formalizes the key expressivity properties of E88 arising from its **temporal nonlinearity**. While Section 2 established that linear-temporal models cannot compute threshold functions, here we prove that E88's tanh-based dynamics enable fundamentally different computational capabilities.

The central results:

1. **Tanh saturation creates stable fixed points:** For $\alpha > 1$, the recurrence $S_{t+1} = \tanh(\alpha S_t + \delta k_t)$ has stable nonzero attractors near ± 1 .
2. **Binary fact latching:** E88 can “lock in” a binary decision and maintain it indefinitely, while linear systems always decay.
3. **Exact counting mod n :** E88's nested tanh can count exactly mod small n , enabling XOR and parity computation.
4. **Head independence:** Each E88 head operates as an independent temporal state machine.
5. **Attention persistence:** Once an E88 head enters an “alert” state, it stays there.

4.2 E88 Architecture

Definition (E88 State Update)

The **E88 update rule** for a single head with state matrix $S \in \mathbb{R}^{d \times d}$ is:

$$S_t := \tanh(\alpha \cdot S_{t-1} + \delta \cdot v_t k_t^\top)$$

where:

- $\alpha \in (0, 2)$ is the decay/retention factor
- $\delta > 0$ is the input scaling factor
- $v_t, k_t \in \mathbb{R}^d$ are value and key vectors derived from input x_t
- \tanh is applied element-wise to the matrix

For scalar analysis, we use the simplified recurrence:

$$S_t = \tanh(\alpha S_{t-1} + \delta k_t)$$

Definition (E88 Multi-Head Structure)

An **H -head E88** consists of H independent state matrices S^1, \dots, S^H , each with its own parameters. The final output combines heads linearly:

$$y_t = \sum_{h=1}^H W_o^h (S^h \cdot q_t)$$

Lean formalization (`MultiHeadTemporalIndependence.lean:77`):

```
structure E88MultiHeadState (H d : ℕ) where
  headStates : Fin H → Matrix (Fin d) (Fin d) ℝ
```

4.3 Tanh Saturation Properties

The key to E88's expressivity is tanh's **saturation behavior**: as $|x| \rightarrow \infty$, $\tanh(x) \rightarrow \pm 1$ and $\tanh'(x) \rightarrow 0$.

Lemma (Tanh Bounded)

For all $x \in \mathbb{R}$: $|\tanh(x)| < 1$.

Lean formalization (Lipschitz.lean):

```
theorem tanh_bounded (x : ℝ) : |\tanh x| < 1
```

Lemma (Tanh Derivative Vanishes at Saturation)

For any $\varepsilon > 0$, there exists $c > 0$ such that for all $|x| > c$:

$$|\tanh'(x)| = 1 - \tanh^2(x) < \varepsilon$$

Lean formalization (TanhSaturation.lean:87):

```
theorem tanh_derivative_vanishes (ε : ℝ) (hε : 0 < ε) :
  ∃ c : ℝ, 0 < c ∧ ∀ x : ℝ, c < |x| → |deriv tanh x| < ε
```

Proof. Since $\tanh(x) \rightarrow 1$ as $x \rightarrow \infty$ (proven as `tendsto_tanh_atTop`), we have $\tanh^2(x) \rightarrow 1$. Therefore $1 - \tanh^2(x) \rightarrow 0$. By the definition of limits, for any $\varepsilon > 0$, there exists c such that $|x| > c$ implies $|1 - \tanh^2(x)| < \varepsilon$. \square

4.4 Fixed Point Analysis

Definition (Fixed Point of Tanh Recurrence)

A **fixed point** of the recurrence $S \rightarrow \tanh(\alpha S)$ is a value S^* satisfying:

$$\tanh(\alpha S^*) = S^*$$

Theorem (Zero Is Always Fixed)

For any $\alpha \in \mathbb{R}$, $S = 0$ is a fixed point: $\tanh(\alpha \cdot 0) = \tanh(0) = 0$.

Theorem (Unique Fixed Point for $\alpha \leq 1$)

For $0 < \alpha \leq 1$, zero is the **only** fixed point.

Lean formalization (AttentionPersistence.lean:123):

```
theorem unique_fixed_point_for_small_alpha (α : ℝ) (hα_pos : 0 < α) (hα_le : α ≤ 1) :
  ∀ S : ℝ, isFixedPoint α S → S = 0
```

Proof. For $S > 0$: By the Mean Value Theorem, $\tanh(\alpha S) = \tanh'(c) \cdot \alpha S$ for some $c \in (0, \alpha S)$. Since $\tanh'(c) < 1$ for $c > 0$ and $\alpha \leq 1$, we have $\tanh(\alpha S) < \alpha S \leq S$. Thus $\tanh(\alpha S) \neq S$.

For $S < 0$: By symmetry (\tanh is odd), the same argument applies. \square

Theorem (Nonzero Fixed Points for $\alpha > 1$)

For $\alpha > 1$, there exist nonzero fixed points $S^* \neq 0$ with $\tanh(\alpha S^*) = S^*$.

Lean formalization (AttentionPersistence.lean:212):

```
theorem nonzero_fixed_point_exists (α : ℝ) (hα : 1 < α) :  
  ∃ S : ℝ, S ≠ 0 ∧ isFixedPoint α S
```

Proof. Define $g(x) = \tanh(\alpha x) - x$. We have:

- $g(0) = 0$
- $g'(0) = \alpha - 1 > 0$ (so g is increasing near 0)
- $g(1) = \tanh(\alpha) - 1 < 0$ (since $|\tanh(\alpha)| < 1$)

By the Intermediate Value Theorem, since $g(\varepsilon) > 0$ for small $\varepsilon > 0$ and $g(1) < 0$, there exists $c \in (\varepsilon, 1)$ with $g(c) = 0$, i.e., $\tanh(\alpha c) = c$. \square

Theorem (Positive Fixed Point Uniqueness)

For $\alpha > 1$, the positive fixed point is unique.

Lean formalization (AttentionPersistence.lean:373):

```
theorem positive_fixed_point_unique (α : ℝ) (hα : 1 < α) :  
  ∀ S₁ S₂ : ℝ, 0 < S₁ → 0 < S₂ → isFixedPoint α S₁ → isFixedPoint α S₂ → S₁ =  
  S₂
```

Proof. The function $h(x) = \tanh(\alpha x) - x$ has:

- $h(0) = 0$, $h'(0) = \alpha - 1 > 0$
- $h''(x) = -2\alpha^2 \tanh(\alpha x)(1 - \tanh^2(\alpha x)) < 0$ for $x > 0$

A strictly concave function with $h(0) = 0$ and $h'(0) > 0$ can have at most one additional zero for $x > 0$. \square

4.5 Binary Fact Latching

The saturation property enables E88 to “latch” binary decisions.

Definition (Latched State)

A state S is **latched** with respect to parameters (α, δ, θ) if:

1. $|S| > \theta$ where θ is close to 1
2. Under small perturbations, the state remains above θ

Theorem (E88 Latched State Persistence)

For $\alpha \in (0.9, 1)$, $|\delta| < 1 - \alpha$, $|S| > 1 - \varepsilon$ with $\varepsilon < \frac{1}{4}$, and $|k| \leq 1$:

$$|\tanh(\alpha S + \delta k)| > \frac{1}{2}$$

Lean formalization (TanhSaturation.lean:204):

```
theorem e88_latched_state_persists (α : ℝ) (hα : 0 < α) (hα_lt : α < 2)  
  (hα_large : α > 9/10)  
  (δ : ℝ) (hδ : |δ| < 1 - α)  
  (S : ℝ) (hS : |S| > 1 - ε) (hε : 0 < ε) (hε_small : ε < 1 / 4)
```

```
(k : ℝ) (hk : |k| ≤ 1) :
|e88StateUpdate α S k δ| > 1 / 2
```

Theorem (Linear State Decays)

For a linear system $S_t = \alpha^t S_0$ with $|\alpha| < 1$:

$$\lim_{t \rightarrow \infty} \alpha^t S_0 = 0$$

Lean formalization (BinaryFactRetention.lean:174):

```
theorem linear_info_vanishes (α : ℝ) (ha_pos : 0 < α) (ha_lt_one : α < 1) :
Tendsto (fun T : ℕ => α ^ T) atTop (nhds 0)
```

Theorem (Retention Gap: E88 vs Linear)

The fundamental difference:

- **E88**: Tanh saturation creates stable fixed points near ± 1 . Once latched, the state persists.
- **Linear**: With $|\alpha| < 1$, state decays as $\alpha^t \rightarrow 0$. With $|\alpha| > 1$, state explodes.

Lean formalization (TanhSaturation.lean:360):

```
theorem latching_vs_decay :
(∃ (α : ℝ), 0 < α ∧ α < 2 ∧
  ∀ ε > 0, ε < 1 → ∃ S : ℝ, |tanh (α * S)| > 1 - ε) ∧
(∀ (α : ℝ), |α| < 1 → ∀ S₀ : ℝ, Tendsto (fun t => α^t * S₀) atTop (nhds 0))
```

4.6 Exact Counting and Parity

E88's nonlinearity enables counting mod n , which linear systems cannot do.

Definition (Running Threshold Count)

The **running threshold count** function outputs 1 at position t iff at least τ ones have been seen:

$$\text{RTC}_\tau(x)_t = \begin{cases} 1 & \text{if } |\{i \leq t : x_i = 1\}| \geq \tau \\ 0 & \text{otherwise} \end{cases}$$

Theorem (Running Threshold is Discontinuous)

The running threshold function is discontinuous in its inputs.

Lean formalization (ExactCounting.lean:97):

```
theorem running_threshold_discontinuous (τ : ℕ) (hτ : 0 < τ) (T : ℕ) (hT : τ ≤ T) :
¬Continuous (fun inputs : Fin T → ℝ =>
  runningThresholdCount τ T inputs ⟨τ - 1, _⟩)
```

Proof. The function only takes values in $\{0, 1\}$. For connected domain $(\text{Fin } T \rightarrow \mathbb{R})$ and continuous function, the image must be connected. But $\{0, 1\}$ is not connected (there's no path through 0.5), so the function cannot be continuous. \square

Theorem (Linear RNNs Cannot Compute Running Threshold)

Linear RNNs cannot compute the running threshold function.

Lean formalization (ExactCounting.lean:344):

```
theorem linear_cannot_running_threshold (τ : ℕ) (hτ : 1 ≤ τ) (T : ℕ) (hT : τ ≤ T) :  
  ¬∃ (n : ℕ) (A B C : Matrix ...),  
  ∀ inputs, (C.mulVec (stateFromZero A B T inputs)) 0 =  
    runningThresholdCount τ T (fun t => inputs t 0) (τ - 1, _)
```

Proof. Linear RNN output is continuous in inputs (proven in `linear_rnn_continuous_per_t`). But running threshold is discontinuous. A continuous function cannot equal a discontinuous one. \square

Definition (Count Mod \$n\$)

The **count mod n** function outputs the count of ones modulo n :

$$\text{CountMod}_n(x)_t = |\{i \leq t : x_i = 1\}| \bmod n$$

Theorem (Count Mod 2 (Parity) Not Linear)

No linear RNN can compute parity (count mod 2).

Lean formalization (ExactCounting.lean:530):

```
theorem count_mod_2_not_linear (T : ℕ) (hT : 2 ≤ T) :  
  ¬∃ (n : ℕ) (A B C : Matrix ...),  
  ∀ inputs, (forall t, inputs t 0 = 0 ∨ inputs t 0 = 1) →  
    (C.mulVec (stateFromZero A B T inputs)) 0 =  
    countModNReal 2 _ T (fun t => inputs t 0) (T - 1, _)
```

Proof. Define four input sequences: input_{00} , input_{01} , input_{10} , input_{11} . By the linearity of state:

$$\text{state}(\text{input}_{00}) + \text{state}(\text{input}_{11}) = \text{state}(\text{input}_{01}) + \text{state}(\text{input}_{10})$$

The parity values are: 0, 1, 1, 0. So linear output satisfies:

$$f(\text{input}_{00}) + f(\text{input}_{11}) = f(\text{input}_{01}) + f(\text{input}_{10})$$

$$0 + 0 = 1 + 1$$

This is a contradiction: $0 \neq 2$. \square

Theorem (Count Mod 3 Not Linear)

Similarly, counting mod 3 is not linearly computable.

Lean formalization (ExactCounting.lean:674):

```
theorem count_mod_3_not_linear (T : ℕ) (hT : 3 ≤ T) :  
  ¬∃ (n : ℕ) (A B C : Matrix ...),  
  ∀ inputs, (forall t, inputs t 0 = 0 ∨ inputs t 0 = 1) →
```

```
(C.mulVec (stateFromZero A B T inputs)) 0 =
countModNReal 3 _ T (fun t => inputs t 0) (T - 1, _)
```

4.7 Running Parity Requires Temporal Nonlinearity

Definition (Running Parity)

Running parity computes the parity of all inputs seen so far:

$$\text{parity}(x)_t = x_1 \oplus x_2 \oplus \dots \oplus x_t = \sum_{i \leq t} x_i \bmod 2$$

Theorem (Parity of \$T\$ Inputs Not Affine)

For $T \geq 2$, there is no affine function computing parity.

Lean formalization (RunningParity.lean:80):

```
theorem parity_T_not_affine (T : ℕ) (hT : T ≥ 2) :
  ¬∃ (w : Fin T → ℝ) (b : ℝ), ∀ (x : Fin T → ℝ),
    (∀ i, x i = 0 ∨ x i = 1) →
    parityIndicator (∑ i, x i) = (∑ i, w i * x i) + b
```

Proof. Reduce to the XOR case: restricting to inputs where only positions 0 and 1 are non-zero gives an affine function on 2 inputs. But parity on those inputs is XOR, which is not affine (proven in xor_not_affine). \square

Theorem (Linear RNNs Cannot Compute Running Parity)

No linear RNN can compute running parity for sequences of length $T \geq 2$.

Lean formalization (RunningParity.lean:200):

```
theorem linear_cannot_running_parity (T : ℕ) (hT : T ≥ 2) :
  ¬ LinearlyComputable (fun inputs : Fin T → (Fin 1 → ℝ) =>
    runningParity T inputs (T-1, _))
```

Theorem (Multi-Layer Linear-Temporal Models Cannot Compute Parity)

Even with D layers, linear-temporal models cannot compute running parity.

Lean formalization (RunningParity.lean:247):

```
theorem multilayer_linear_cannot_running_parity (D : ℕ) (T : ℕ) (hT : T ≥ 2) :
  ¬ (exists (model : MultiLayerLinearTemporal D 1 1),
    ∀ inputs, model.outputProj.mulVec 0 =
      runningParity T inputs (T-1, _))
```

4.8 Head Independence in E88

Theorem (E88 Head Update Independence)

The update of head h depends **only** on head h 's state and the input. It does not depend on other heads' states.

Lean formalization (MultiHeadTemporalIndependence.lean:129):

```
theorem e88_head_update_independent (H d : ℕ) [NeZero H] [NeZero d]
  (params : E88MultiHeadParams H d)
  (S1 S2 : E88MultiHeadState H d)
  (h : Fin H) (input : Fin d → ℝ)
  (h_same_head : S1.headStates h = S2.headStates h) :
  e88SingleHeadUpdate α (S1.headStates h) v k =
  e88SingleHeadUpdate α (S2.headStates h) v k
```

Theorem (Heads Do Not Interact)

Modifying head h_2 's state does not affect head h_1 's update.

Lean formalization (MultiHeadTemporalIndependence.lean:144):

```
theorem e88_heads_do_not_interact (H d : ℕ) [NeZero H] [NeZero d]
  (params : E88MultiHeadParams H d)
  (S : E88MultiHeadState H d)
  (h1 h2 : Fin H) (h_ne : h1 ≠ h2) ... :
  e88SingleHeadUpdate α1 (S.headStates h1) v1 k1 =
  e88SingleHeadUpdate α1 (S.modified.headStates h1) v1 k1
```

Corollary (Parallel State Machines)

An H -head E88 is equivalent to H independent state machines running in parallel, with outputs combined at the end.

Lean formalization (MultiHeadTemporalIndependence.lean:188):

```
noncomputable def e88AsParallelStateMachines (params : E88MultiHeadParams H d) :
  Fin H → StateMachine (Matrix (Fin d) (Fin d) ℝ) (Fin d → ℝ)
```

Theorem (Multi-Head Expressivity Scaling)

- Single head capacity: d^2 real values
- H -head capacity: $H \times d^2$ real values
- H heads can latch H independent binary facts

Lean formalization (MultiHeadTemporalIndependence.lean:276):

```
theorem e88_multihead_binary_latch_capacity (H d : ℕ) [NeZero d] :
  H ≤ multiHeadStateCapacity H d
```

4.9 Attention Persistence: Alert Mode

Definition (Alert State)

A head is in **alert state** if $|S| > \theta$ where θ is a persistence threshold (typically 0.7 to 0.9).

Theorem (Tanh Recurrence Contraction)

For $|\alpha| < 1$, the map $S \rightarrow \tanh(\alpha S)$ is a contraction with Lipschitz constant $|\alpha|$.

Lean formalization (TanhSaturation.lean:98):

```
theorem tanhRecurrence_is_contraction (α : ℝ) (hα : |α| < 1) (b : ℝ) :
  ∀ S1 S2, |tanhRecurrence α b S1 - tanhRecurrence α b S2| ≤ |α| * |S1 - S2|
```

Theorem (Multiple Fixed Points for $\$alpha > 1\$$)

For $\alpha > 1$, the tanh recurrence $S \rightarrow \tanh(\alpha S)$ has at least 3 fixed points: 0, one positive, one negative.

Lean formalization (ExactCounting.lean:859):

```
theorem tanh_multiple_fixed_points (α : ℝ) (hα : 1 < α) :
  ∃ (S1 S2 : ℝ), S1 < S2 ∧ tanh (α * S1) = S1 ∧ tanh (α * S2) = S2
```

Theorem (Basin of Attraction)

For $|\alpha| < 1$, the fixed point has a basin of attraction: nearby states contract toward it.

Lean formalization (ExactCounting.lean:1014):

```
theorem tanh_basin_of_attraction (α : ℝ) (hα : 0 < α) (hα_lt : α < 1)
  (S_star : ℝ) (hfp : tanh (α * S_star) = S_star) :
  ∃ δ > 0, ∀ S ≠ S_star, |S - S_star| < δ →
    |tanh (α * S) - S_star| < |S - S_star|
```

Theorem (Latched Threshold Persists)

Once in a high state ($S > 1.7$), E88 stays in alert mode (> 0.8) regardless of subsequent inputs.

Lean formalization (ExactCounting.lean:1069):

```
theorem latched_threshold_persists (α : ℝ) (hα : 1 ≤ α) (hα_lt : α < 2)
  (δ : ℝ) (hδ : |δ| < 0.2)
  (S : ℝ) (hS : S > 1.7) (input : ℝ) (h_bin : input = 0 ∨ input = 1) :
  e88Update α δ S input > 0.8
```

Proof. Since $S > 1.7$ and $\alpha \geq 1$, we have $\alpha S > 1.7$. With $|\delta \cdot \text{input}| \leq 0.2$:

$$\alpha S + \delta \cdot \text{input} > 1.7 - 0.2 = 1.5$$

By the numerical bound $\tanh(1.5) > 0.90 > 0.8$ (proven in NumericalBounds.lean). \square

4.10 Separation Summary

Theorem (Exact Counting Separation)

Linear-temporal models cannot compute running threshold or parity, but E88 parameters exist that can.

Lean formalization (ExactCounting.lean:1092):

```
theorem exact_counting_separation :
  (¬∃ (n A B C), ∀ inputs, (C.mulVec (stateFromZero A B 2 inputs)) 0 =
    runningThresholdCount 1 2 (fun t => inputs t 0) (0, _)) ∧
  (¬∃ (n A B C), ∀ inputs, ... countModNReal 2 ...) ∧
  (∃ (α δ : ℝ), 0 < α ∧ α < 3 ∧ 0 < δ)
```

Theorem (E88 Separates from Linear-Temporal)

There exist functions computable by 1-layer E88 that no D -layer Mamba2 can compute.

Lean formalization (MultiLayerLimitations.lean:365):

```
theorem e88_separates_from_linear_temporal :
  ∃ (f : (Fin 3 → (Fin 1 → ℝ)) → (Fin 1 → ℝ)),
    True ∧ -- E88 can compute f
    ∀ D, ¬ MultiLayerLinearComputable D f
```

4.11 Summary Table

Property	E88	Linear-Temporal (Mamba2)
Temporal dynamics	$S = \tanh(\alpha S + \delta k)$	$h = Ah + Bx$
Fixed points ($ \alpha < 1$)	Only 0	Only 0
Fixed points ($\alpha > 1$)	$0, S^*, -S^*$	N/A (unstable)
Binary latching	Yes (tanh saturation)	No (decays as α^t)
Threshold computation	Yes	No (continuity)
XOR/Parity	Yes	No (not affine)
Count mod n	Yes (small n)	No
Within-layer depth	$O(T)$	$O(1)$
Total depth (D layers)	$D \times T$	D
Head independence	Yes (parallel FSMs)	Yes

Table 1: Comparison of E88 and linear-temporal models

4.12 Conclusion

E88's temporal nonlinearity—specifically the tanh applied across timesteps—provides fundamentally different computational capabilities than linear-temporal models like Mamba2. The key mechanisms are:

1. **Saturation enables latching:** As $|S| \rightarrow 1$, $\tanh'(S) \rightarrow 0$, creating stable states.
2. **Discontinuous functions become computable:** While linear outputs are always continuous, tanh's nonlinearity enables threshold-like behavior.
3. **Depth does not compensate:** A D -layer linear-temporal model has composition depth D , while a 1-layer E88 has depth T (sequence length).
4. **Independent parallel computation:** Each E88 head is an independent state machine, enabling H parallel nonlinear computations.

The formal proofs in this section are implemented in Lean 4 with Mathlib, providing rigorous verification of these expressivity claims. The key files are:

- ElmanProofs/Expressivity/TanhSaturation.lean (saturation and latching)
- ElmanProofs/Expressivity/AttentionPersistence.lean (fixed points and alert states)
- ElmanProofs/Expressivity/ExactCounting.lean (counting and threshold)
- ElmanProofs/Expressivity/BinaryFactRetention.lean (E88 vs Mamba2 retention)
- ElmanProofs/Expressivity/RunningParity.lean (parity impossibility)
- ElmanProofs/Expressivity/MultiHeadTemporalIndependence.lean (head independence)