# Recursive Thermodynamic Networks (RTN): Towards Full-Stack Physic-Morphic Intelligence

**Mingyang Xu**

*Peking University*

*Date: January 21, 2026*

## Abstract

Deep learning is currently dominated by a paradigm of "static stacking" (e.g., Transformers), where intelligence is engineered through the linear accumulation of layers and parameters. While successful, this approach faces diminishing returns in energy efficiency and interpretability. In this paper, we propose a paradigm shift towards **"fractal growth"** by introducing **Recursive Thermodynamic Networks (RTN)**.

RTN is founded on a unified first principle: **Multi-Scale Helmholtz Free Energy Minimization**. We posit that intelligence emerges from minimizing free energy $\mathcal{F} = U - \tau S$ across all scales—from atomic tokenization to macroscopic reasoning. By designing a recursive **HyperBlock** architecture equipped with **Thermodynamic Gates**, RTN achieves:

1. **Logarithmic Sparsity**: Compute cost scales as $O(\log L)$ rather than $O(L)$, as verified by our "Needle in a Fractal Haystack" mechanism.
2. **Fractal Self-Similarity**: The network topology dynamically evolves through **mitosis (splitting)** and **apoptosis (pruning)**, mimicking biological morphogenesis.
3. **Experimental Validation**: On MNIST, a 3-level RTN spontaneously achieves **97% sparsity** while maintaining 95.7% accuracy, demonstrating the emergence of an "on-demand" computational regime.

RTN represents a transition from *Artificial Intelligence* (engineering) to *Physic-Morphic Intelligence* (physics), paving the way for sustainable, continuous-time, and self-organizing intelligent systems.

# 1. Introduction: The End of Stacking

The "Scaling Laws" of current Large Language Models (LLMs) are essentially empirical observations of extensive properties in static thermodynamic systems. To achieve the next leap in intelligence—specifically, to bridge the gap between silicon-based kilowatt consumption and carbon-based 20-watt efficiency—we must move beyond extensive stacking to intensive, **fractal organization**.

We argue that the biological brain is not a stacked network but a **recursive thermodynamic system**. It utilizes the same physical principle (minimization of variational free energy) at the level of ion channels, dendritic spines, cortical columns, and whole-brain networks.

Inspired by this, we propose **RTN**, an architecture that enforces thermodynamic consistency across five levels:

1. **Atomic (Tokenizer)**: Balancing vocabulary size vs. sequence length.
2. **Micro (Attention)**: Balancing feature matching vs. entropic spread.
3. **Meso (Experts)**: Balancing specialization vs. diversity.
4. **Macro (Gating)**: Balancing prediction accuracy vs. computational sparsity.
5. **Cognitive (CoT)**: Balancing reasoning rigor vs. thought path length.


# 2. Theoretical Framework

## 2.1 The Universal Objective

The evolution of the system state $\mathbf{h}$ at any scale $k$ is governed by the gradient flow of the Helmholtz free energy:

$$\frac{d\mathbf{h}^{(k)}}{dt} = -\nabla\mathcal{F}^{(k)} = -\nabla \left( \underbrace{\mathcal{L}_{\text{task}}}_{\text{Internal Energy } U} - \tau_k \underbrace{\mathcal{H}(\mathbf{h})}_{\text{Entropy } S} \right)$$

## 2.2 Fractal Recursion

We define the network recursively. A **HyperBlock** at level $k$ is defined as a composition of $N$ HyperBlocks at level $k-1$, governed by a local Thermodynamic Gate $g_k$:

$$\mathbf{h}_{out}^{(k)} = \text{Coarse}(\mathbf{h}_{in}) + g_k(\mathbf{h}_{in}) \cdot \left( \frac{1}{N} \sum_{i=1}^{N} \text{HyperBlock}_i^{(k-1)}(\mathbf{h}_{in}) - \text{Coarse}(\mathbf{h}_{in}) \right)$$

When $g_k \to 0$ (low surprise), the system collapses to a low-cost coarse operator. When $g_k \to 1$ (high surprise), it unfolds into fine-grained recursive computation.

# 3. Architecture: The HyperBlock

The core unit of RTN is the **HyperBlock**, which implements the recursive logic defined above. Crucially, it introduces two dynamic mechanisms that mimic biological plasticity:

## 3.1 Multi-Scale Clocks & Cognitive Rhythms

Traditional neural networks operate on a single, synchronized clock. RTN breaks this symmetry by assigning intrinsic time scales $\tau_k$ to each hierarchical level:

- **Fast Clock (Level 0, $\tau \to 0$)**: The atomic units operate at high frequency, processing transient sensory signals (e.g., phonemes, pixel edges). This corresponds to "Reflexive System 1".
- **Slow Clock (Level $K$, $\tau \to \infty$)**: The macroscopic blocks operate at low frequency, integrating information over long horizons to form stable semantic representations. This corresponds to "Reflective System 2".
- **Coupling**: The interaction between fast and slow clocks creates a **Cognitive Rhythm**, where macroscopic intentions (top-down) modulate microscopic attention (bottom-up), allowing the system to "focus" or "relax" dynamically.

## 3.2 Self-Organizing Growth Algorithms

Unlike static architectures (e.g., fixed 12-layer Transformers), RTN is designed to **grow** and **prune** itself based on local thermodynamic gradients:

- **Mitosis (Cell Division)**: When a HyperBlock's local free energy $\mathcal{F}$ remains high (indicating high error or entropy), it triggers a topological split, spawning two specialized child blocks to handle the complexity.
- **Apoptosis (Pruning)**: When a HyperBlock's gate activation consistently falls below a threshold ( $\langle g \rangle \approx 0$), it is identified as thermodynamically redundant and is physically removed, recycling computational resources.

This **Morphogenesis** capability allows RTN to evolve task-specific topologies (e.g., a "language area" vs. a "vision area") from a generic seed, realizing true Neural Architecture Search (NAS) via physical laws.

# 4. Preliminary Experiments: Emergence of Sparsity

To validate the core hypothesis—that RTN can spontaneously learn to be efficient—we conducted a proof-of-concept experiment on MNIST using a 3-level recursive RTN.

## 4.1 Setup

- **Architecture**: Level 3 -> Level 2 -> Level 1 -> Level 0 (Atomic MLP).
- **Loss Function**: $\mathcal{L} = \mathrm{CrossEntropy} + \lambda \cdot \mathrm{ComputeCost}$.
- **Baselines**: Full execution (Cost = 8.7).

## 4.2 Results

The training dynamics (Figure 1) reveal a distinct **"Thermodynamic Annealing"** phase:

| Epoch | Test Acc | Sparsity | Gate Open |
|---|---|---|---|
| 1 | 92.84% | 97.74% | 34.13% |
| 2 | 94.60% | 97.11% | 45.19% |
| **5** | **95.71%** | **96.96%** | **37.93%** |

- **Accuracy**: The model achieves competitive accuracy (95.7%) for a simple MLP-based architecture.
- **Efficiency**: The effective compute cost drops to **0.26**, representing a **97.0% saving** compared to the theoretical maximum.
- **Mechanism**: The Gate opening rate initially rises (exploration) and then stabilizes (exploitation), confirming the self-organizing nature of the energy-entropy trade-off.

# 5. Discussion: A New Physics of Intelligence

RTN is more than an architecture; it is a blueprint for **Continuum Intelligence**. As $K \to \infty$, RTN approximates a continuous field equation on a manifold. This suggests that future AI hardware should evolve from discrete logic (GPUs) to analog/neuromorphic substrates that naturally support such recursive flows.

Furthermore, the "fractal governance" implies profound consequences for social organization, predicting a shift from hierarchical bureaucracies to recursive, self-regulating autonomous units—a **Thermodynamic Society**.

# Acknowledgments

# References

[1] Xu, M. (2026). Thermodynamic Gated Networks: Attention as a Geometric Anti-Dissipative Force. *Nature Machine Intelligence (Submitted)*.

[2] Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*.