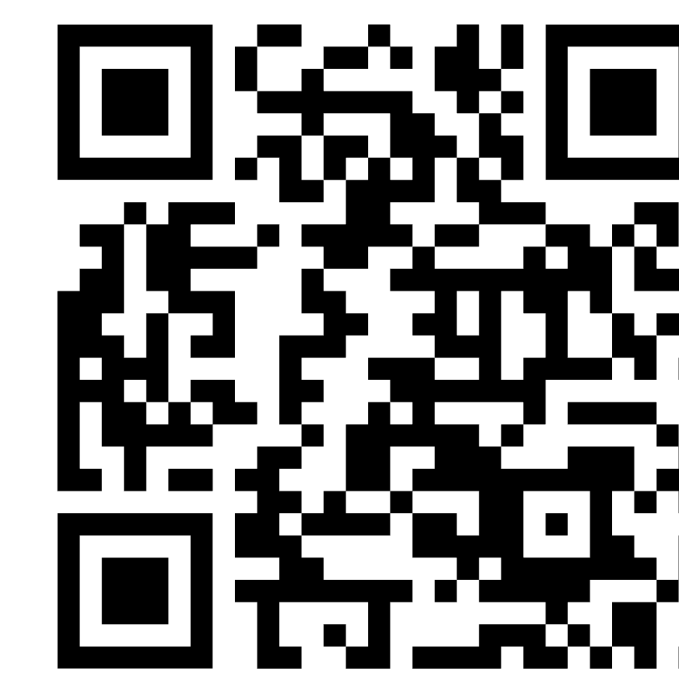


# Hyper-SET: Designing Transformers via Hyperspherical Energy Minimization

Yunzhe Hu<sup>1</sup>, Difan Zou<sup>1,2</sup>, and Dong Xu<sup>1</sup>

<sup>1</sup>The University of Hong Kong, <sup>2</sup>Shenzhen Loop Area Institute



Project Page w/ Code



Full Paper

## Problem

Transformers are great, but...

- **Engineered from the bottom up.** Their architecture remains largely heuristics-driven—key components are arranged by trial and error.
- **Mysteriously redundant.** Evidence that representations are similar in the middle layers of LLMs suggests a convergent layer functionality.
- **Mostly interpreted post hoc.** Current tools to interpret their inner workings (e.g., SAEs, circuit analysis) are mostly with hindsight—hard to break the performance ceiling.

## What we do

A top-down approach by asking:

Can we find or design a function prior that induces a model interpretable by construction?

**Our response:** We think the answer is **YES**, at least for a family of Transformers.

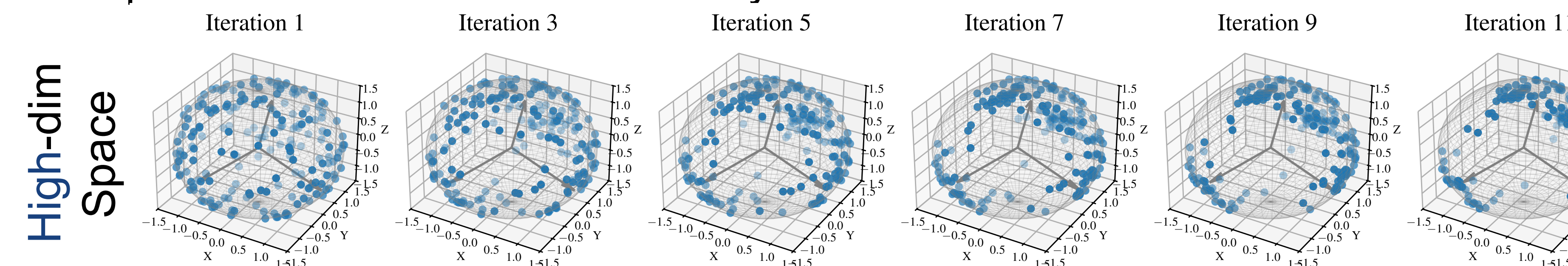
**Main result:** We introduce *Hyper-SET*, an **intrinsically interpretable** Transformer where every core component—from self-attention to skip connections—emerges naturally from a single, principled objective:

maximum likelihood estimation  
on the hypersphere

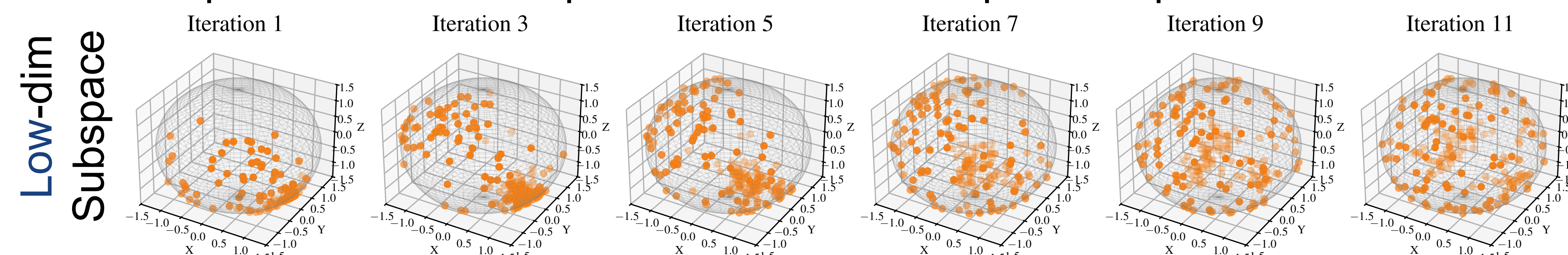
## Conceptualization

We begin by conceptualizing what effective representations should look like. We posit that token dynamics should satisfy two complementary properties:

- **Semantic Alignment** aligns tokens with learned semantic directions to compress uninformative redundancy.



- **Distributional Uniformity** prevents representation collapse and ensures tokens spread out as isotropic Gaussian on the sphere to preserve volume.



## Objective-driven Architecture

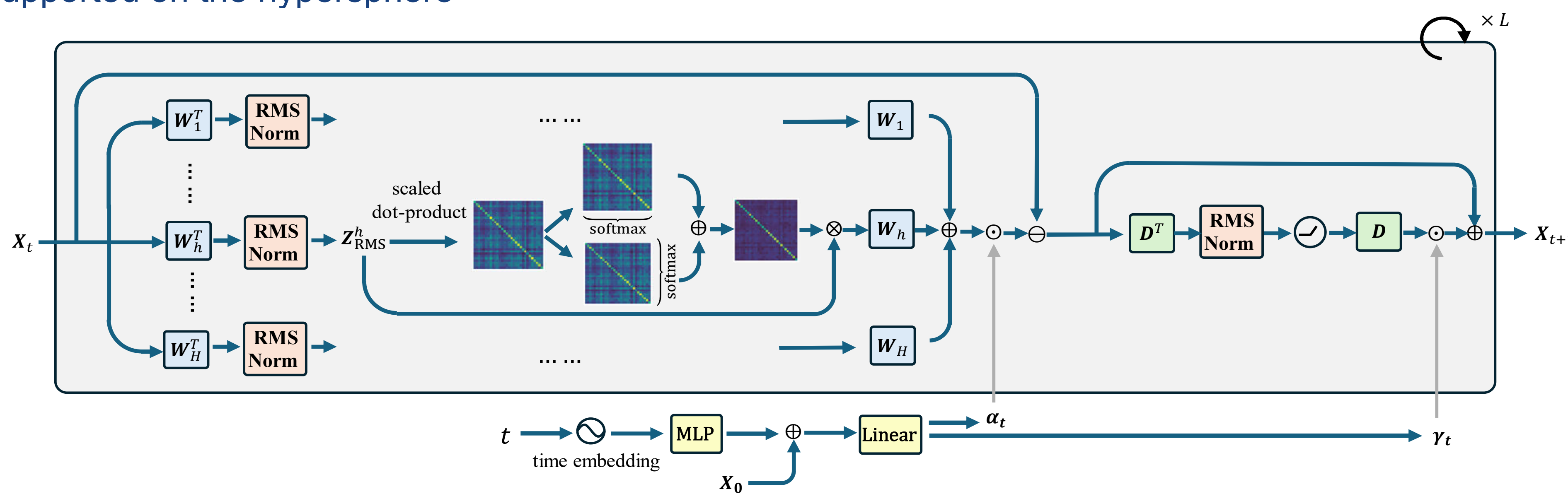
$$\max_{\mathbf{x}} \mathbb{E}_{(z^1, \dots, z^H) \sim p(z^1, \dots, z^H)} [\log p(\mathbf{x}, z^1, \dots, z^H; \theta, \phi)]$$

$$\min_{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}} E(\mathbf{X}; \mathbf{W}, \mathbf{D}) = E_{\text{ATTN}} + E_{\text{FF}},$$

$$\text{subject to } \|\mathbf{W}_h^\top \mathbf{x}_i\|_2 = \sqrt{p}, \quad \|\mathbf{D}^\top \mathbf{x}_i\|_2 = \sqrt{M}, \quad i = 1, \dots, N.$$

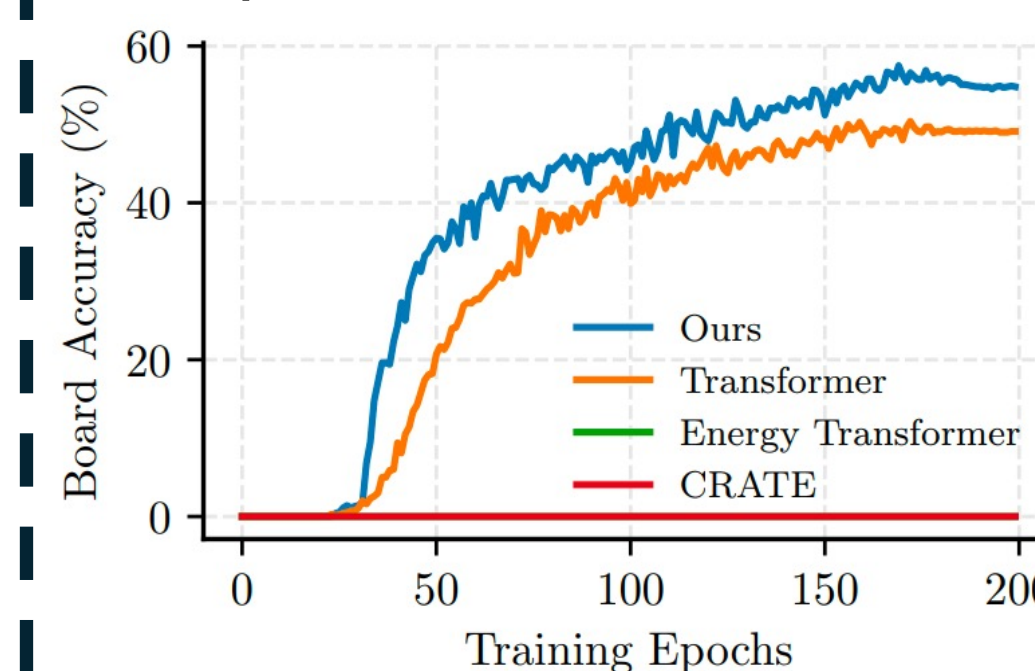
Quantified into optimizable, modified Hopfield energy functions

**Isotropic Gaussian**  
supported on the hypersphere

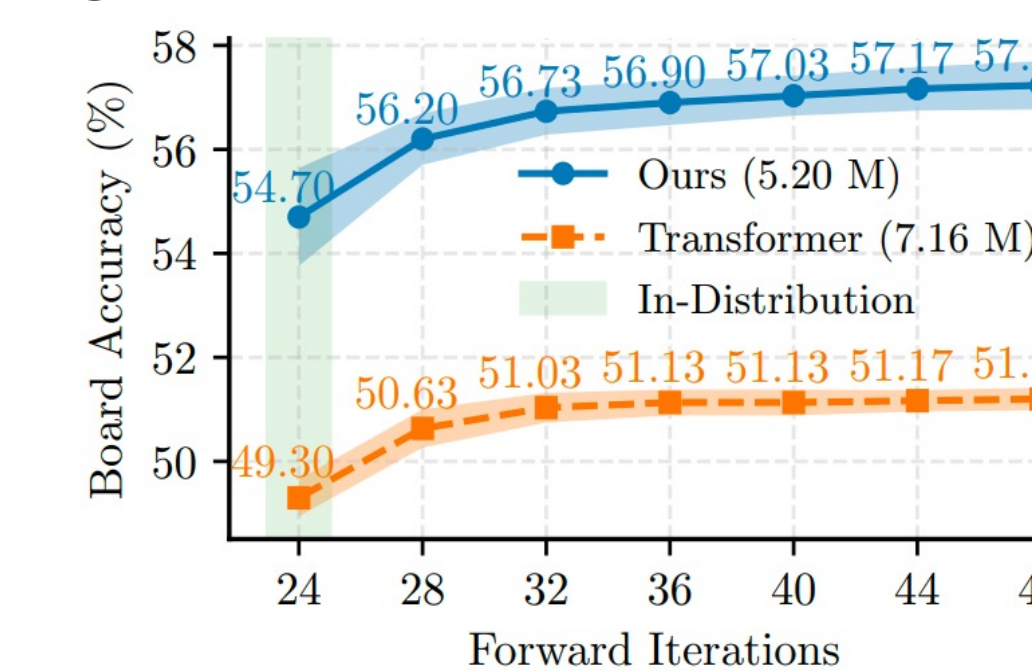


## Empirical Results

### a) Sudoku Reasoning



(a) Sudoku training dynamics.



(b) Sudoku test-time extrapolation.

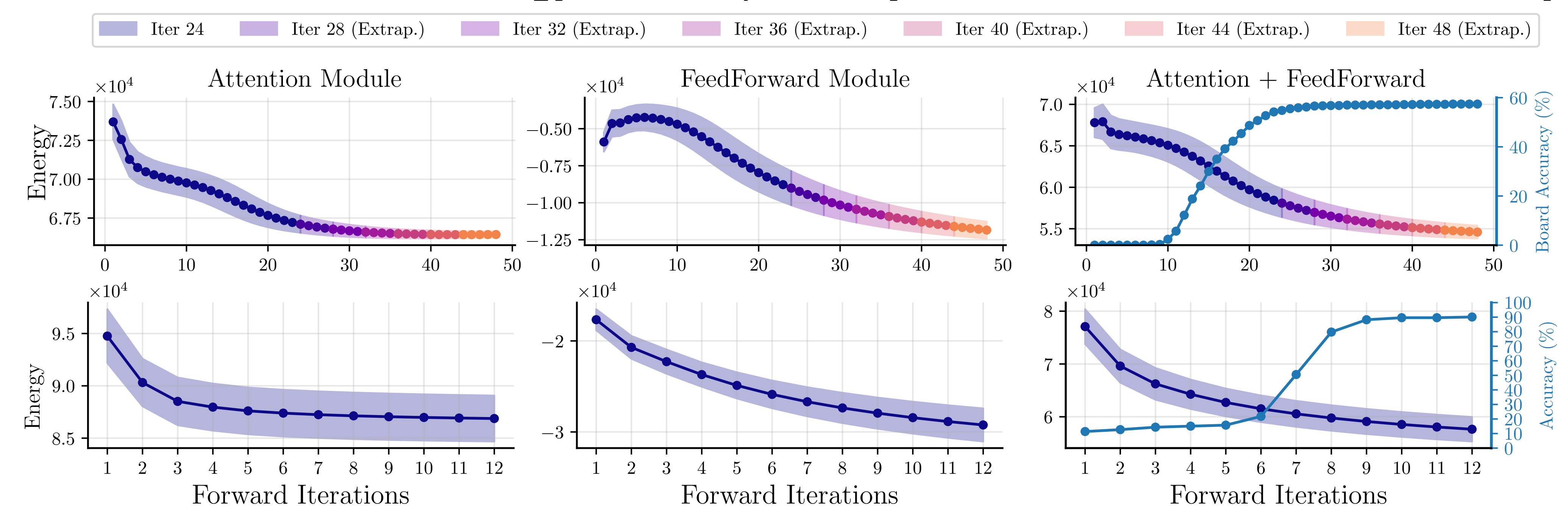
### b) Image Classification

Table 1: Top-1 accuracy (%) for image classification with single-layer recurrent-depth models. Parameters are measured on ImageNet-1K. All models are trained from scratch on the listed datasets.

| Model                                    | Width $d$ | # Params (M) | Dataset      |              |              |              |
|--|-----------|--------------|--------------|--------------|--------------|--------------|
|  |           |              | CIFAR-10     | CIFAR-100    | IN-100       | IN-1K        |
| Transformer                              | 384       | 2.38         | 89.90        | 61.89        | 69.44        | <b>66.94</b> |
| CRATE-T (Hu et al., 2024c)               | 896       | 3.04         | 87.54        | 60.23        | 68.16        | 57.89        |
| CRATE (Yu et al., 2023)                  | 768       | 3.00         | 84.81        | 58.22        | 68.52        | 57.00        |
| Energy Transformer (Hoover et al., 2024) | 512       | 2.39         | 76.39        | 50.60        | 36.68        | 34.24        |
| <b>HYPER-SET (Ours)</b>                  | 512       | 2.39         | <b>90.11</b> | 63.41        | <b>70.16</b> | 62.76        |
| <b>HYPER-SET (Ours)</b>                  | 640       | 3.40         | 89.96        | <b>64.60</b> | 69.31        | 66.21        |

### c) Energy Evolution

Energy Descent *positively* correlates with increase in accuracy



## Extensibility

We are **NOT** re-interpreting existing Transformers! but instead providing a general design framework. Designing energy functions unlocks novel variants (e.g., **linear attention**).

| Operator             | $f(x)$                   | $K(x, y)$  | $E_{\text{ATTN}}$   | $-\nabla_{\mathbf{X}} E_{\text{ATTN}}$  |
|----------------------|--------------------------|--|---|---|
| Bi-Softmax (Default) | $\beta^{-1} \log(x)$     | $\exp(\beta \mathbf{x}^\top \mathbf{y})$                       | Eq. 3   | Eq. 7   |
| Sigmoid Attention    | $\frac{\beta^{-1}}{2} x$ | $\sigma(\beta \mathbf{x}^\top \mathbf{y})$                     | $\frac{1}{2} \sum_{h=1}^H \sum_{i,j=1}^N \sigma(\beta (\mathbf{W}_h^\top \mathbf{x}_i)^\top \mathbf{W}_h^\top \mathbf{x}_j) \beta^{-1}$       | $\sum_{h=1}^H \mathbf{W}_h \mathbf{W}_h^\top \mathbf{X} \sigma(1 - \sigma) (\beta (\mathbf{W}_h^\top \mathbf{X})^\top \mathbf{W}_h^\top \mathbf{X})$  |
| Linear Attention     | $\frac{\beta^{-1}}{2} x$ | $\frac{1}{2} (\beta \Phi(\mathbf{x})^\top \Phi(\mathbf{y}))^2$ | $\frac{1}{4} \sum_{h=1}^H \sum_{i,j=1}^N (\beta \Phi(\mathbf{W}_h^\top \mathbf{x}_i)^\top \Phi(\mathbf{W}_h^\top \mathbf{x}_j))^2 \beta^{-1}$ | $\sum_{h=1}^H \mathbf{W}_h \Phi'(\mathbf{W}_h^\top \mathbf{X}) \odot (\beta \Phi(\mathbf{W}_h^\top \mathbf{X}) \Phi(\mathbf{W}_h^\top \mathbf{X})^\top \Phi(\mathbf{W}_h^\top \mathbf{X}))$ |

## Summary

- Frames representation learning as joint maximum likelihood estimation on hyperspheres
- Bridges the gap between energy-based learning and practical Transformer design
- Provides a general design principle that unlocks novel variants in core Transformer block