

# Hyper-SET: Designing Transformers via Hyperspherical Energy Minimization

**Yunzhe Hu**, Difan Zou, and Dong Xu

The University of Hong Kong



SCHOOL OF  
**COMPUTING &  
DATA SCIENCE**  
The University of Hong Kong

# 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 believe 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 symmetric attention to skip connections—emerges naturally from a single, principled objective:

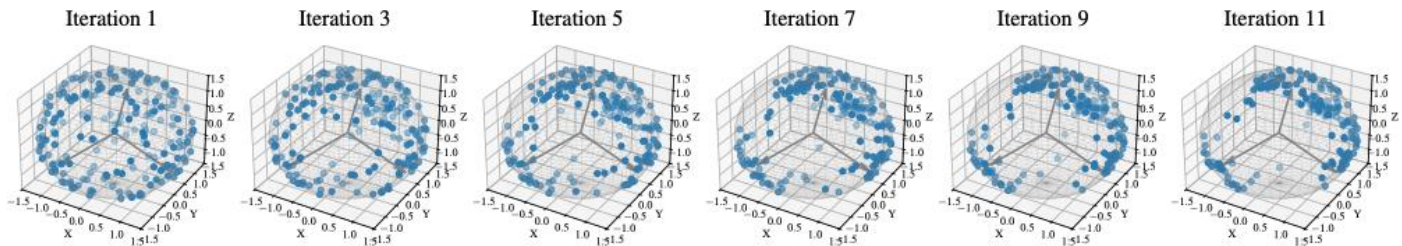
maximum likelihood estimation on the hypersphere.

# Conceptualization

We conjecture that token dynamics should satisfy two complementary properties:

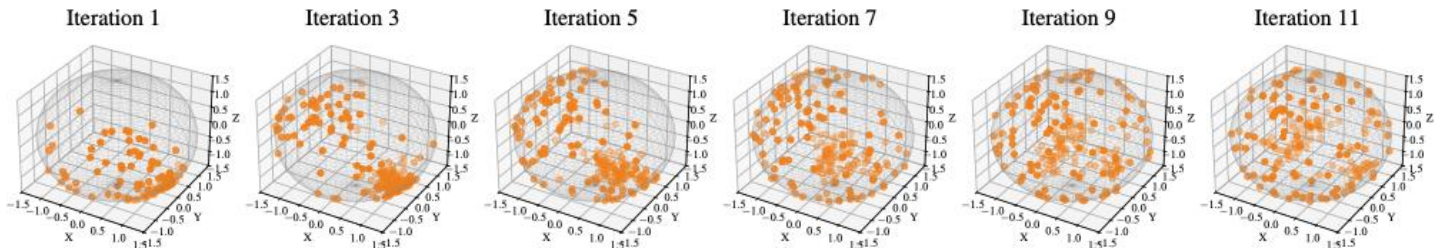
- **Semantic Alignment** aligns tokens with learned semantic directions to compress uninformative redundancy.

High-dim.  
Space




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

Low-dim.  
Subspace




# Theoretical Justification

1. Joint maximum likelihood on the hypersphere (under mild assumptions):


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

2. Two complementary objectives: **Isotropic Gaussian supported on hypersphere**


$$\min_{\mathbf{x}} \sum_{h=1}^H \underbrace{D_{\text{KL}}(p(\mathbf{z}) \| p_{\phi}(\mathbf{z}^h | \mathbf{x}))}_{\text{uniformity}} - \underbrace{\log p_{\theta}(\mathbf{x})}_{\text{alignment}}$$

↑

3. Quantified into optimizable, modified Hopfield energy functions:

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

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

# Hyper-SET Architecture

## Symmetric self-attention:

$$\dot{\mathbf{X}} = -\nabla_{\mathbf{X}} E_{\text{ATTN}} \quad \text{ODE discretization}$$
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \alpha_t \sum_{h=1}^H \left( \mathbf{W}_h \mathbf{Z}_{\text{RMS},t}^h \underbrace{\text{softmax}([\mathbf{Q}\mathbf{K}]_{\text{RMS},t})}_{\text{column-wise}} + \mathbf{W}_h \mathbf{Z}_{\text{RMS},t}^h \underbrace{\text{softmax}([\mathbf{Q}\mathbf{K}]_{\text{RMS},t})}_{\text{row-wise}} \right)$$

## Symmetric feedforward:

$$\dot{\mathbf{X}} = -\nabla_{\mathbf{X}} E_{\text{FF}} = \mathbf{D} \text{ReLU}(\mathbf{D}^\top \mathbf{X}) \quad \text{ODE discretization}$$
$$\mathbf{X}_{t+1} = \mathbf{X}_t + \gamma_t \mathbf{D} \text{ReLU}(\text{RMSNorm}(\mathbf{D}^\top \mathbf{X}_t))$$

where **RMSNorm** emerges as spherical constraints

$$\mathbf{Z}_{\text{RMS}}^h = \text{RMSNorm}(\mathbf{Z}^h) = \text{RMSNorm}(\mathbf{W}_h^\top \mathbf{X})$$

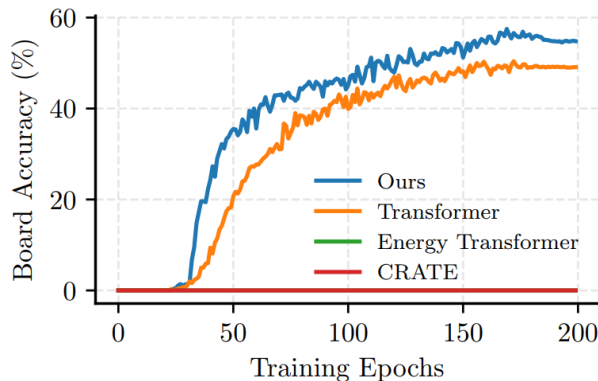
$$[\mathbf{Q}\mathbf{K}]_{\text{RMS},t} = \beta(\mathbf{Z}_{\text{RMS},t}^h)^\top (\mathbf{Z}_{\text{RMS},t}^h)$$



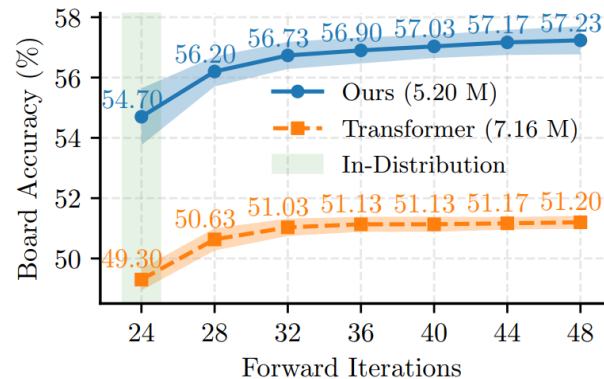
# Practical Competitiveness, Validation & Extensibility

## Sudoku Reasoning

- Train / Test



(a) Sudoku training dynamics.



(b) Sudoku test-time extrapolation.

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.

## Image Classification

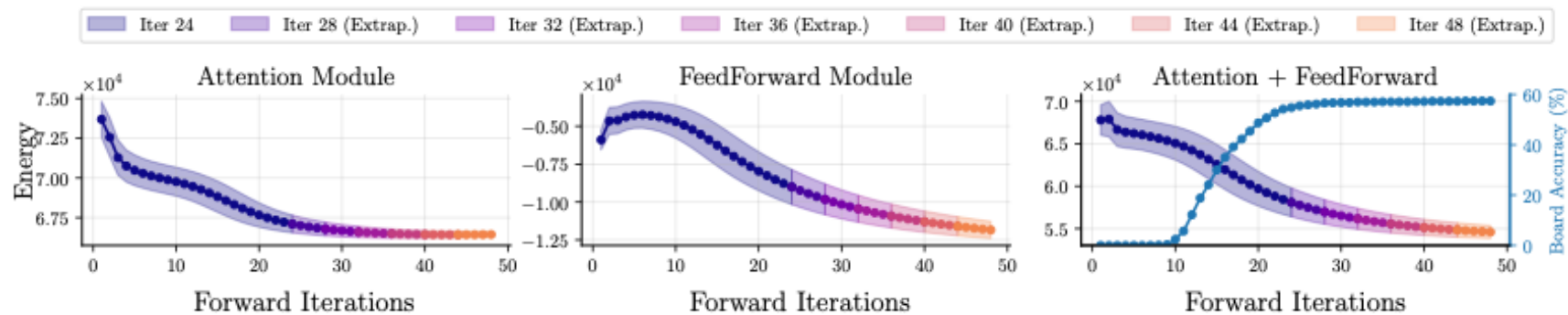
- ImageNet-1K

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
HYPER-SET (Ours)	512	2.39	<b>90.11</b>	63.41	<b>70.16</b>	62.76
HYPER-SET (Ours)	640	3.40	89.96	<b>64.60</b>	69.31	66.21

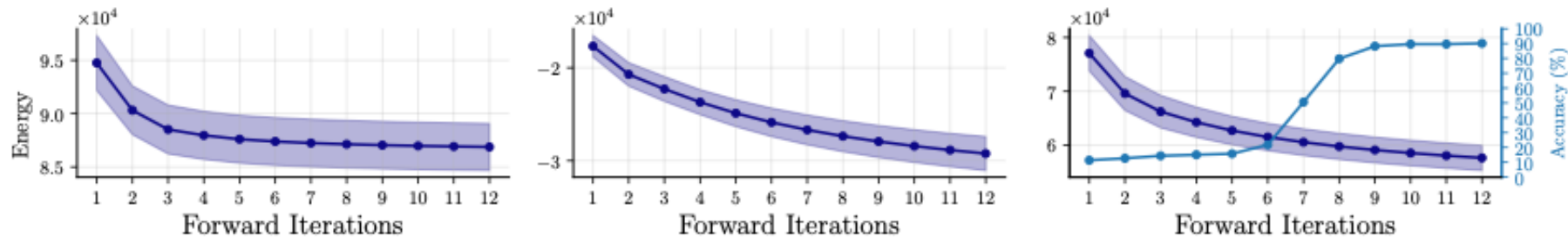
# Practical Competitiveness, Validation & Extensibility

## Energy Descent *positively* correlates with increase in accuracy

- Sudoku



- Image Classification



# Practical Competitiveness, Validation & Extensibility

**Extensibility:** we are **NOT** re-interpreting existing Transformers!! but rather providing a more general design framework.

By defining various energy functions under our formalism, we can design novel components beyond what's widely known, like **linear attention** and **gated feedforward**.

Operator	$f(x)$	$K(\mathbf{x}, \mathbf{y})$	$E_{\text{ATTN}}$	$-\nabla_{\mathbf{X}} E_{\text{ATTN}}$
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	$\frac{\beta^{-1}}{2} x$	$\frac{1}{2} (\beta \Phi(\mathbf{x})^\top \Phi(\mathbf{y}))^2$	$\frac{1}{4} \sum_{h=1}^H \sum_{i,j=1}^N (\beta \Phi(\mathbf{W}_h^\top \mathbf{x}_i)^\top \Phi(\mathbf{W}_h^\top \mathbf{x}_j))^2 \beta^{-1}$	$\sum_{h=1}^H \mathbf{W}_h \Phi'(\mathbf{W}_h^\top \mathbf{X}) \odot (\beta \Phi(\mathbf{W}_h^\top \mathbf{X}) \Phi(\mathbf{W}_h^\top \mathbf{X})^\top \Phi(\mathbf{W}_h^\top \mathbf{X}))$

Operator	$g(x)$	$h(x)$	$E_{\text{FF}}$	$-\nabla_{\mathbf{X}} E_{\text{FF}}$
ReLU FF (Default)	$x$	$\frac{1}{2} \text{ReLU}^2(x)$	Eq. 5	Eq. 10
Softmax FF	$\log(x)$	$\exp(x)$	$-\sum_{i=1}^N \log\left(\sum_{m=1}^M \exp(\mathbf{d}_m^\top \mathbf{x}_i)\right)$	$\underbrace{D}_{\text{column-wise}} \text{softmax}(D^\top \mathbf{X})$
Gated FF	$\frac{1}{2} x^2$	$\Phi(x)$	$-\frac{1}{2} \sum_{i=1}^N \left(\sum_{m=1}^M \Phi(\mathbf{d}_m^\top \mathbf{x}_i)\right)^2$	$\underbrace{D \Phi(D^\top \mathbf{X})}_{\text{column sum}} \odot \Phi'(D^\top \mathbf{X})$

# Summary & What's next

## Hyper-SET

- 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

## Looking ahead:

- How to extend it to autoregressive modeling?
- Can it connect to flow matching to discover new (layer-wise) training strategies?



Project Page



Full Paper

Project Page: [hyper-set.github.io](https://hyper-set.github.io)

Github Repo: <https://github.com/huyunzhe/hyper-set>