

Recurrent Self-Attention Dynamics: An Energy-Agnostic Perspective from Jacobians

Akiyoshi Tomihari & Ryo Karakida

Artificial Intelligence Research Center, AIST, Japan

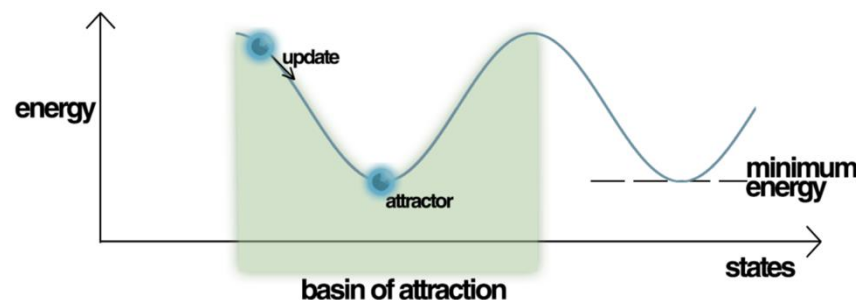


Theoretical Analysis of Self-attention

- Self-attention (SA) is a core mechanism in Transformers
 - Transformers are the backbone of modern language models

$$\text{SA}_h(\mathbf{X}) := \text{softmax}(\beta \mathbf{X} \mathbf{W}_h^Q \mathbf{W}_h^{K\top} \mathbf{X}^\top) \mathbf{X} \mathbf{W}_h^V$$

- Theoretical analysis of SA is challenging
 - Mainly due to the Softmax function
- A major approach involves energy-based analysis
 - SA can be viewed as minimizing an implicit or explicit energy function



Energy-based analysis of SA

- SA can be viewed as minimizing an energy function [Geshkovski et al., 2023, 2024, Bruno et al., 2025]
- Often formulated using continuous equations and a particle-based interpretation of tokens

$$\dot{\mathbf{X}} = \text{Proj}_{\mathbf{X}} (\text{softmax}(\beta \mathbf{X} \mathbf{W}^Q \mathbf{W}^{K\top} \mathbf{X}^\top) \mathbf{X} \mathbf{W}^V)$$

- To define an energy function, constraints are typically assumed
 - Weight constraints and single-head constraint

$$\mathbf{W}^Q \mathbf{W}^{K\top} = \mathbf{W}^V = \mathbf{W}^{V\top} \text{ or } \mathbf{W}^Q = \mathbf{W}^K = \mathbf{W}^V = \mathbf{I}_D$$

- However, these constraints deviate from practical scenarios

Our approach

1. Relax constraints in energy-based formulations
 - remove symmetry constraints on weights
 - extend beyond single-head SA
 2. Apply Jacobian-based analysis to SA
 - Captures linear stability without requiring an energy function
 - Compute Lyapunov exponents from Jacobians
- **Our study reveals Jacobians and Lyapunov exponents as fundamental tools for realistic SA.**

Revisit Energy-based analysis

What is a more realistic setting for energy guarantees?

- Remove symmetry constraints on weights (W^Q, W^K)

Proposition 4.1. *Consider the continuous-time dynamics for single-head SA equipped with projection (3). The energy function*

$$E_{\text{single}}(\mathbf{X}) = - \sum_{i,j} \exp \left(\beta \mathbf{X}_{[i,:]}^\top \mathbf{W}^Q \mathbf{W}^{K\top} \mathbf{X}_{[j,:]} \right) \quad (10)$$

is monotonically decreasing as $dE_{\text{single}}(\mathbf{X})/dt \leq 0$ under the condition:

$$\mathbf{W}^V = (\mathbf{W}^{Q\top} \mathbf{W}^K + \mathbf{W}^Q \mathbf{W}^{K\top})/2. \quad (11)$$

- Extend beyond single-head self-attention

Proposition 4.2. *Consider the continuous-time dynamics for multi-head SA without projection: $d\mathbf{X}/dt = \sum_{h=1}^H SA_h(\mathbf{X})$. An energy function*

$$E_{\text{multi}}(\mathbf{X}) = - \sum_h \sum_{i,j} \exp \left(\beta \mathbf{X}_{[i,:]}^\top \mathbf{W}_h^Q \mathbf{W}_h^{K\top} \mathbf{X}_{[j,:]} \right) \quad (12)$$

is monotonically decreasing as $dE_{\text{multi}}(\mathbf{X})/dt \leq 0$ under the condition

$$\mathbf{W}_h^V = (\mathbf{W}_h^{Q\top} \mathbf{W}_h^K + \mathbf{W}_h^Q \mathbf{W}_h^{K\top})/2, \quad \mathbf{W}_h^Q \mathbf{W}_h^{K\top} = \mathbf{U}_{1,h} \mathbf{U}_{2,h}^\top, \quad (13)$$

where $\mathbf{U}_{1(2),h} \in \mathbb{R}^{D \times D/(2H)}$ ($h \in [1, H]$) satisfies the orthogonality condition $\mathbf{U}_{k,h}^\top \mathbf{U}_{k',h'} = \delta_{hh'} \delta_{kk'} \mathbf{I}_{D/(2H)}$.

Apply Jacobian-based analysis to SA

- More general than energy-based analysis:
 - captures **linear stability**
 - detects **non-stationary dynamics** more easily
- **Normalization** plays a crucial role in discrete systems
 - no counterpart in continuous-time dynamics

Proposition 5.1. *Suppose that, in the update of ItrSA (9), the input to the normalization layer satisfies $\|\mathbf{X}_{[i,:]} + \eta\Delta\mathbf{X}_{[i,:]}\| \geq R$ for all $i \in [1, S]$. Then, the spectral norm of the Jacobian satisfies*

$$\left\| \frac{\partial \text{RMSNorm}(\mathbf{X} + \eta\Delta\mathbf{X})}{\partial \mathbf{X}} \right\|_2 \leq \frac{\max_j(|\gamma_j|)}{R} (1 + |\eta| \|\mathbf{J}_{MSA}(\mathbf{X})\|_2), \quad (14)$$

where $\mathbf{J}_{MSA}(\mathbf{X}) := \partial \text{MSA}(\mathbf{X}) / \partial \mathbf{X}$ denotes the Jacobian of MSA.

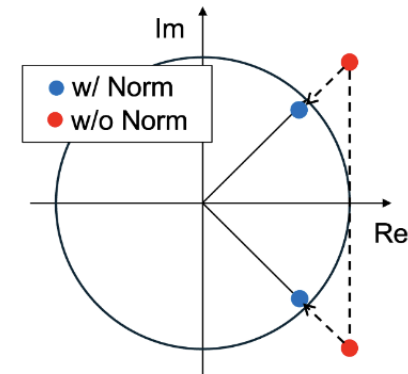
- Normalization operators suppress the Jacobian eigenvalues

Example: oscillatory components

$$\begin{aligned}
 \dot{x} &= \Omega x, & J(x) &= \Omega & (\text{continuous}) \\
 x^{(t+1)} &= (I_D + \eta \Omega) x^{(t)}, & J(x^{(t)}) &= I_D + \eta \Omega & (\text{discrete w/o Norm}) \\
 x^{(t+1)} &= \Pi((I_D + \eta \Omega) x^{(t)}), & J(x^{(t)}) &= \left(I_D - \frac{yy^\top}{\|y\|^2} \right) \frac{I_D + \eta \Omega}{\|y\|} & (\text{discrete w/ Norm})
 \end{aligned}$$

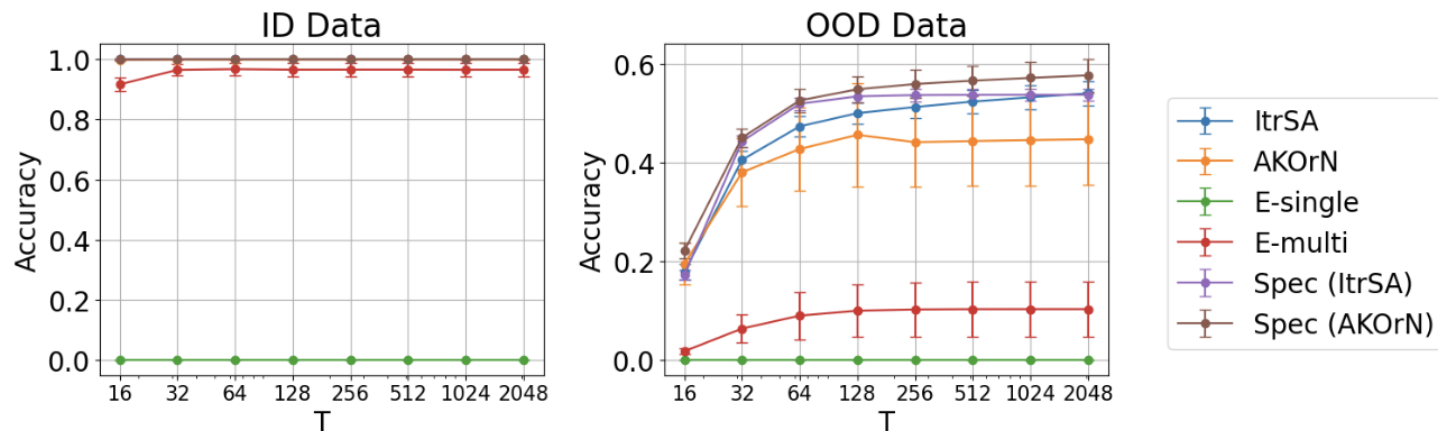
Ω : anti-symmetric matrix

- The continuous systems can exhibit stability
- Discrete counterparts require normalization for stability



(b) Effect of normalization on eigenvalues in oscillatory case

Application: Regularization



- Apply energy-based and Jacobian spectral regularization
 - Regularization based on each approach
- Energy-based regularization performs worse than the original
- Jacobian spectral regularization outperforms both

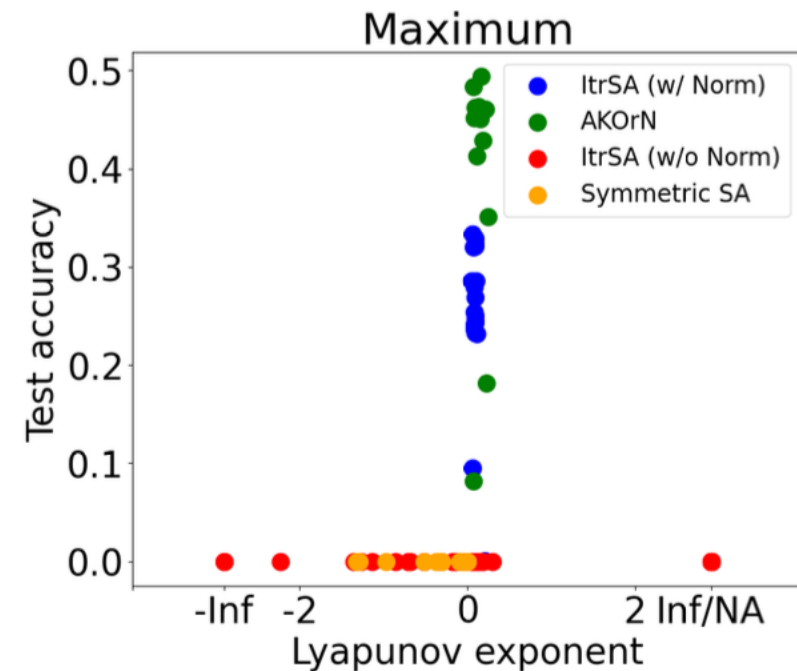
Lyapunov Exponent and Criticality

- Measures local divergence/convergence rate in dynamics
 - Positive \rightarrow instability, Negative \rightarrow convergence, Zero \rightarrow criticality

Normalization

\rightarrow exponents toward zero & higher acc.

- Successful models:
 - max exponent $\approx +0.1$
 - \rightarrow near-critical, from the chaotic side



Conclusion

- Relax constraints in energy-based formulations
- Extend Jacobian-based analysis to SA
- Empirically, strong SA models have Lyapunov exponent ≈ 0
- Jacobians and Lyapunov exponent emerge as fundamental tools for realistic SA architectures

Paper link:

