



# Bridging Empirics and Theory: Unveiling Asymptotic Universality Across Gaussian and Gaussian Mixture Inputs in Deep Learning

Jaeyong Bae, Hawoong Jeong<sup>†</sup> Department of Physics, KAIST

## I. Introduction

PhysRevX.10.041044

### Previous Findings: Dynamics of Hidden Manifold Teacher-Student Model

The core property  $\mathcal{C}$  is based on the latent space (manifold) dimension  $D$ .

In real world, this core property becomes wrapped, resulting to  $X$ , on the real space dimension  $N$ .

$$\mathcal{C} \sim \mathcal{N}(0, I)$$

#### Gaussian Equivalence Property

In the Limit of  $N \rightarrow \infty, D \rightarrow \infty$ , the preactivations of the teacher and student model  $\lambda = XW^T/\sqrt{N}, \nu = C\tilde{W}^T/\sqrt{D}$  conform jointly Gaussian variables. **Statistics involving  $\{\lambda, \nu\}$  are entirely represented by their mean and covariances**

#### SGD to ODE + Change of Basis

Dynamics are defined by statistics of preactivation distribution

“The Dynamics of Student Model are fully theoretically tractable”

Q: The Simple Gaussian  $\mathcal{C} \sim \mathcal{N}(0, I)$  setting is enough?

How dynamics changes as the inherent distribution deviate from simple Gaussian to Gaussian mixture?

#### Ex) $\nu_k$ dynamics

$$W_{k,i} := W_{k,i} - \frac{\eta}{\sqrt{N}} v_k (\hat{y} - y) g'(\lambda_k) f(U_i)$$

$$v_k := v_k - \frac{\eta}{N} g(\lambda_k) (\hat{y} - y)$$

#### Scaled SGD

$$\Delta t \rightarrow dt$$

$$\frac{dv_k}{dt} = \eta \left[ \sum_n^M \tilde{v}_n I_2(k, m) - \sum_j^K v_j I_2(k, j) \right]$$

Expectation of functions

$$I_2(k, m) = \mathbb{E}[g(\lambda_k) \tilde{g}(\nu_m)]$$

$$I_2(k, j) = \mathbb{E}[g(\lambda_k) g(\lambda_j)]$$

Tractable under GEP;

$\therefore \{\lambda, \nu\}$  Statistics are tractable

## II. Setting

### Gaussian Mixture Setting

-  $m$  component Gaussian mixture with moderating parameter  $\alpha$

$$r = r_i \sim \mathcal{N}(\mu_i, I), \quad \mu_i \sim \mathcal{U}[-\alpha, \alpha], \quad \text{with } p_i, \sum p_i = 1$$

Empirical investigations across a spectrum ( $m, \alpha$ ) of Gaussian mixture

### Dynamic Metrics

- Covariance of preactivations  $Q_{k,\ell}, R_{k,m}$
- generalization error  $\epsilon_g$  and 2<sup>nd</sup> layer weight  $\nu_k$

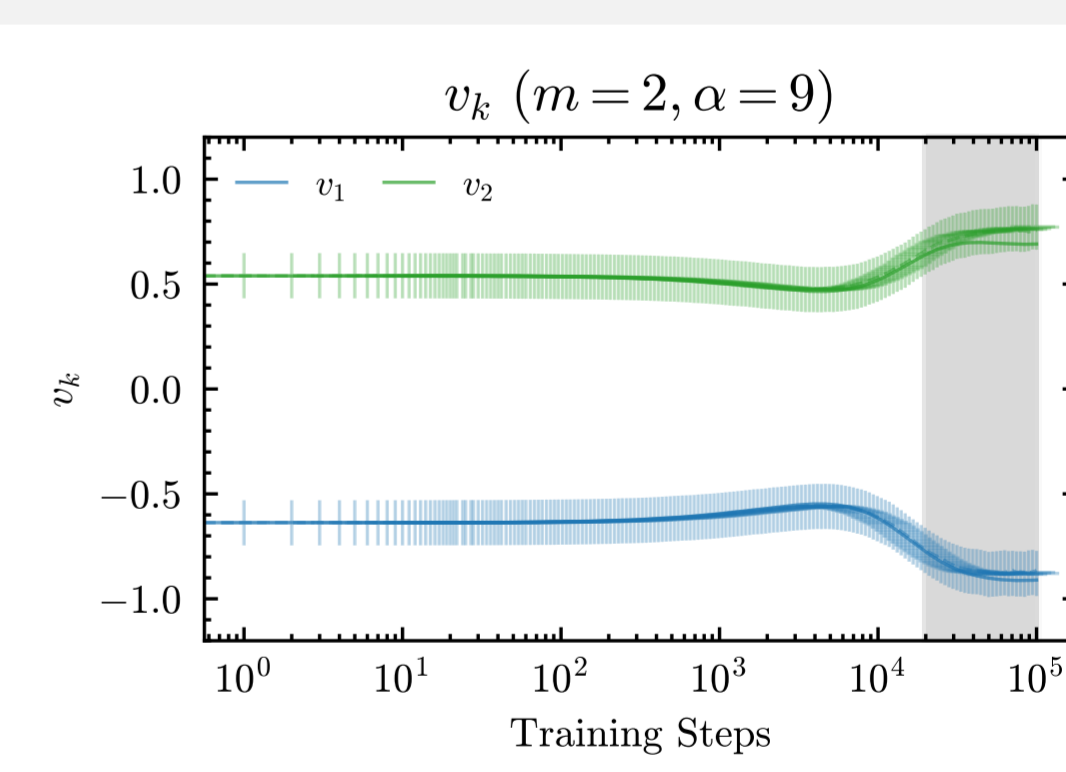
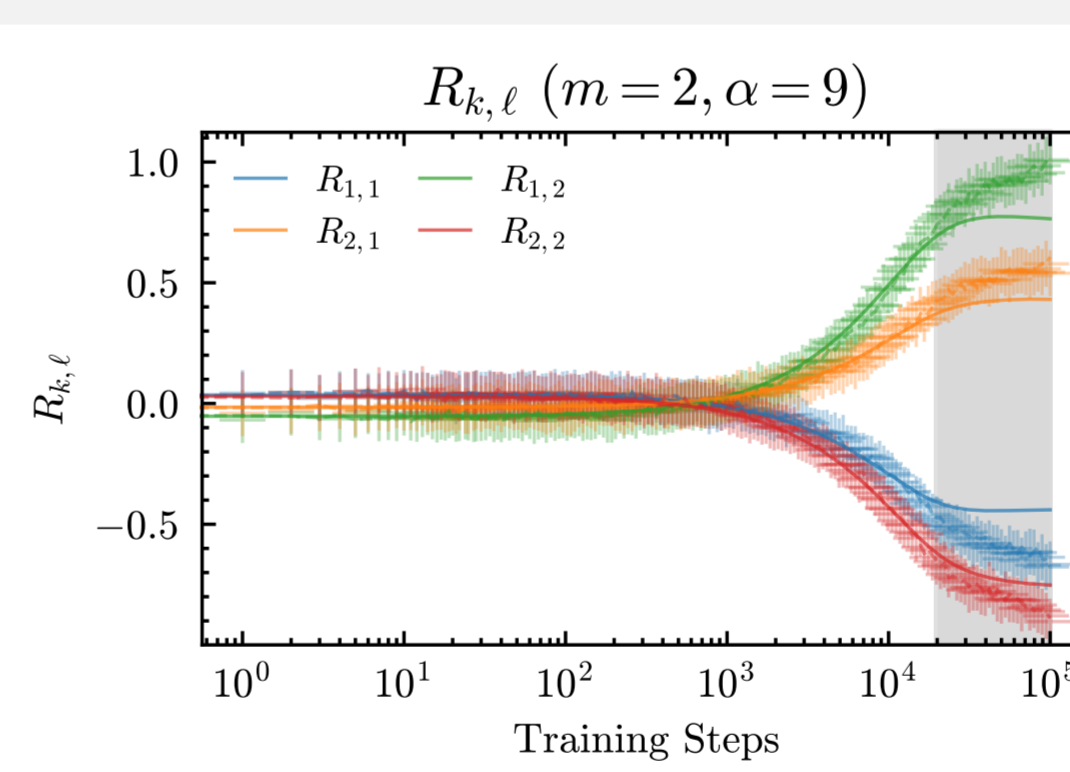
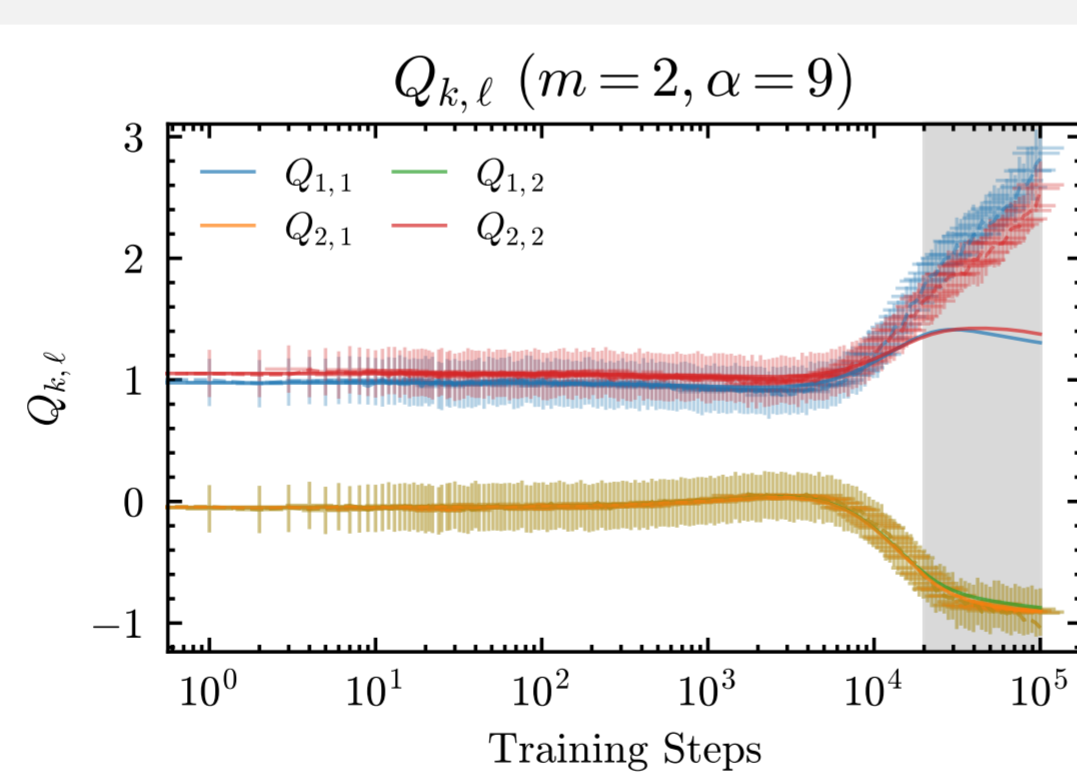
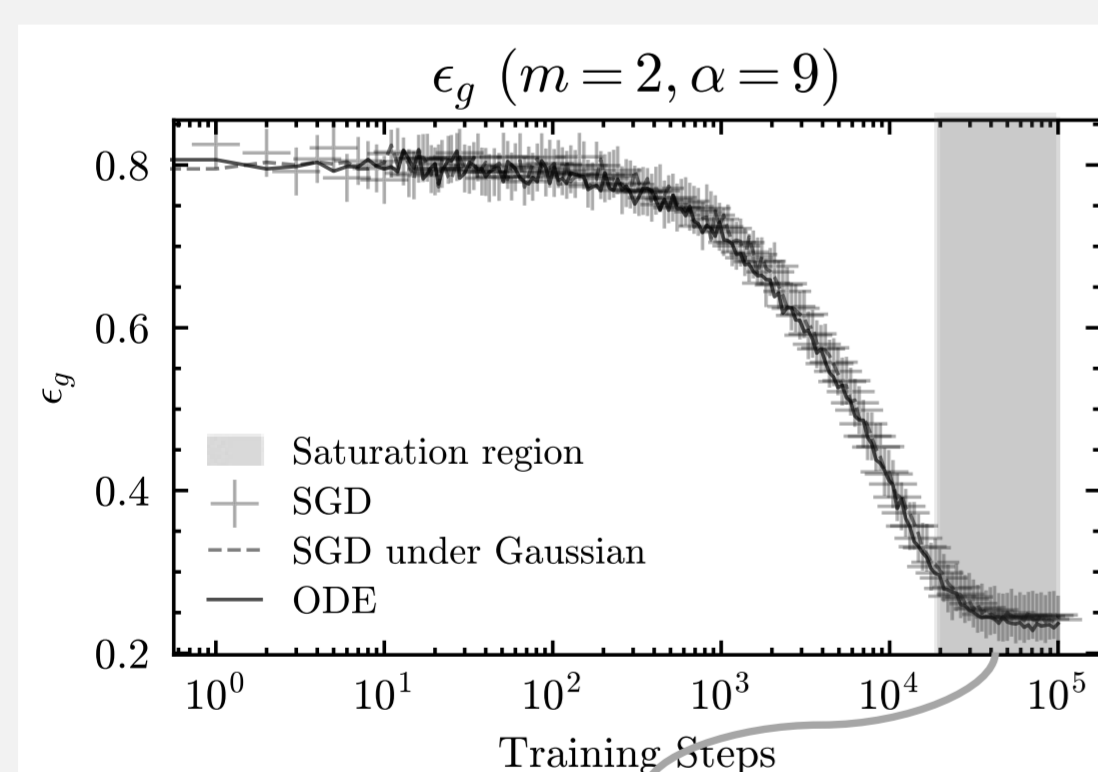
## III. Results & Discussion

### Dynamics across various Gaussian mixture

- With **non-standardized** Gaussian mixture, the theoretical prediction completely **collapse** for **all spectrum of mixture**
- With **standardized** Gaussian mixture, **mixture spectrum dose not effect**

### Dynamics under **standardized** Gaussian mixture

$$\text{Sample-wise } C \sim \mathcal{P}, \quad C := (C - \mathbb{E}[C]) / \sqrt{\mathbb{E}[(C - \mathbb{E}[C])^2]}$$



Saturation region; Randomness has a significant effect

- Theoretical prediction under the Simple Gaussian **surprisingly aligned** with Mixture dynamics

### Convergence in Standardized Gaussian Mixture

**Key quantity: “expectation value of functions” in dynamics**

#### Dominance of Moments in the Expectation value of Functions

[Definition] If  $\mathcal{P}$  shares identical cumulants with  $\mathcal{D}$  up to order 2, represent  $\mathcal{P}$  as  $\mathcal{P}_{\mathcal{D}2}$

[Lemma] If a function  $f$  has the property of **erasing the influence of high-order cumulants**, then the expectation value of the function tends towards the same value.

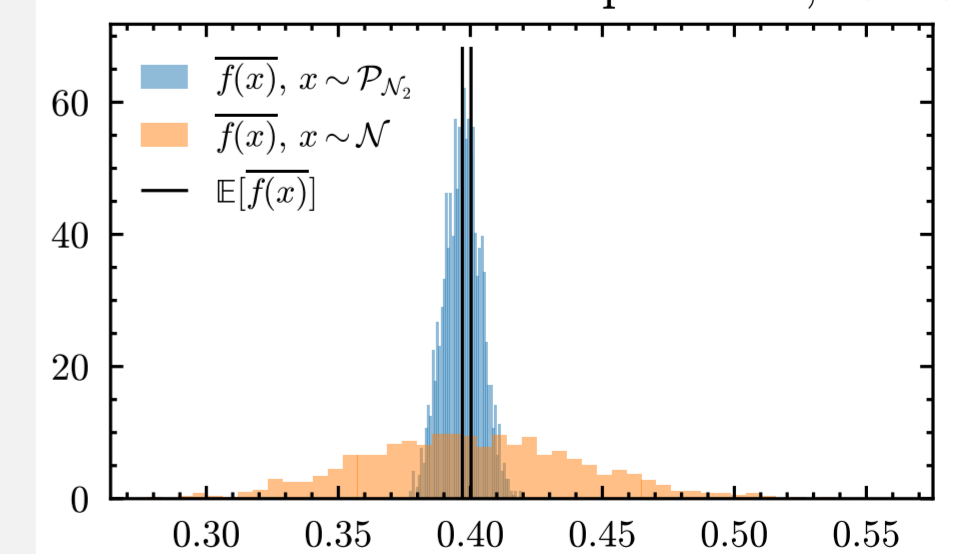
$$\mathbb{E}_{x \sim \mathcal{P}_{\mathcal{D}2}}[f(x)] \rightarrow \mathbb{E}_{x \sim \mathcal{D}}[f(x)]$$

[Proof] Taylor expansion of expectation + Function property vanish high order derivation

$$\mathbb{E}_{x \sim \mathcal{P}_{\mathcal{D}2}}[f(x)]_{x \in I} = \mathbb{E}[f(\mu) + \dots + f^{(n)}(\mu) \frac{(x - \mu)^n}{n!} + \dots]_{x \in I} \approx f(\mu) + f''(\mu) \mathbb{E}\left[\frac{(x - \mu)^2}{2!}\right] \rightarrow \mathbb{E}_{x \sim \mathcal{D}}[f(x)]_{x \in I}$$

$$\mathbb{E}_{x \sim \mathcal{P}_{\mathcal{D}2}}[f(x)] \rightarrow \mathbb{E}_{x \sim \mathcal{D}}[f(x)]$$

Distribution of the sample mean, ReLU



Empirical results, with  $f \equiv \text{ReLU}$

[Lemma] ReLU activation function **well erasing** influence of high-order cumulants

[Proof] In piecewise view, ReLU function is zero in 2<sup>nd</sup> Derivative

## IV. Summary

**Specific function + standardized**  $\Rightarrow$  Make Dynamics dominantly up to 2<sup>nd</sup> Cumulants.

Even the mixture, the **dynamics also tend to approximate same.**

### Acknowledgement

This study was supported by the Basic Science Research Program through the National Research Foundation of Korea (NRF Grant No. 2022R1A2B5B02001752).