

# Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures

## ICML 2020

MEA. Seddik<sup>12</sup>, C.Louart<sup>13</sup>, M. Tamaazousti<sup>1</sup>, R. Couillet<sup>23</sup>

<sup>1</sup><http://melaseddik.github.io/>

<sup>1</sup>CEA List, France

<sup>2</sup>CentraleSupélec, L2S, France

<sup>3</sup>GIPSA Lab Grenoble-Alpes University, France

June 8, 2020



CentraleSupélec

# Abstract

## Context:

- ▶ Study of large **Gram** matrices of **concentrated** data.

## Motivation:

- ▶ **Gram** matrices are at the core of various ML algorithms.
- ▶ RMT predicts their performances under **Gaussian** assumptions on the data.
- ▶ **BUT Real data** are **unlikely close** to **Gaussian** vectors.

## Results:

- ▶ **GAN data** ( $\approx$  **Real data**) fall within the class of **Concentrated** vectors.
- ▶ **Universality result:**  
*Only first and second order statistics of **Concentrated** data matter to describe the behavior of **Gram** matrices.*

# Notion of Concentrated Vectors

## Definition (Concentrated Vectors)

Given a normed space  $(E, \|\cdot\|_E)$  and  $q \in \mathbb{R}$ , a random vector  $\mathbf{Z} \in E$  is  $q$ -exponentially concentrated if for any 1-Lipschitz<sup>1</sup> function  $\mathcal{F} : E \rightarrow \mathbb{R}$ , there exists  $C, c > 0$  such that

$$\forall t > 0, \mathbb{P}\{|\mathcal{F}(\mathbf{Z}) - \mathbb{E}\mathcal{F}(\mathbf{Z})| \geq t\} \leq Ce^{-(t/c)^q} \xrightarrow{\text{denoted}} \boxed{\mathbf{Z} \in \mathcal{E}_q(c)}$$

If  $c$  independent of  $\dim(E)$ , we denote  $\boxed{\mathbf{Z} \in \mathcal{E}_q(1)}$

Concentrated vectors enjoy:

**(P1)** If  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  then  $\mathbf{X} \in \mathcal{E}_2(1)$

“Gaussian vectors are concentrated vectors”

**(P2)** If  $\mathbf{X} \in \mathcal{E}_q(1)$  and  $\mathcal{G}$  is a  $\lambda_{\mathcal{G}}$ -Lipschitz map, then  $\mathcal{G}(\mathbf{X}) \in \mathcal{E}_q(\lambda_{\mathcal{G}})$

“Concentrated vectors are stable through Lipschitz maps”

---

<sup>1</sup>Reminder:  $\mathcal{F} : E \rightarrow F$  is  $\lambda_{\mathcal{F}}$ -Lipschitz if  $\forall (\mathbf{x}, \mathbf{y}) \in E^2 : \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{y})\|_F \leq \lambda_{\mathcal{F}} \|\mathbf{x} - \mathbf{y}\|_E$ .

# Why Concentrated Vectors?



Figure: Images artificially generated using the BigGAN model [Brock et al, ICLR'19].

$$\text{Real Data} \approx \text{GAN Data} = \underbrace{\mathcal{F}_L \circ \mathcal{F}_{L-1} \circ \cdots \circ \mathcal{F}_1}_{\mathcal{G}}(\text{Gaussian})$$

where the  $\mathcal{F}_i$ 's correspond to Fully Connected layers, Convolutional layers, Sub-sampling, Pooling and activation functions, residual connections or Batch Normalisation.

⇒ The  $\mathcal{F}_i$ 's are essentially *Lipschitz* operations.

## Why Concentrated Vectors?

- ▶ Fully Connected Layers and Convolutional Layers are affine operations:

$$\mathcal{F}_i(\mathbf{x}) = \mathbf{W}_i \mathbf{x} + \mathbf{b}_i,$$

and  $\|\mathcal{F}_i\|_{lip} = \sup_{\mathbf{u} \neq 0} \frac{\|\mathbf{W}_i \mathbf{u}\|_p}{\|\mathbf{u}\|_p}$ , for any  $p$ -norm.

- ▶ Pooling Layers and Activation Functions: Are 1-Lipschitz operations with respect to any  $p$ -norm (e.g., ReLU and Max-pooling).
- ▶ Residual Connections:  $\mathcal{F}_i(\mathbf{x}) = \mathbf{x} + \mathcal{F}_i^{(l)} \circ \dots \circ \mathcal{F}_i^{(1)}(\mathbf{x})$   
where the  $\mathcal{F}_i^{(j)}$ 's are Lipschitz operations, thus  $\mathcal{F}_i$  is a Lipschitz operation with Lipschitz constant bounded by  $1 + \prod_{j=1}^l \|\mathcal{F}_i^{(j)}\|_{lip}$ .
- ▶ ...

By:

- (P1) If  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  then  $\mathbf{X} \in \mathcal{E}_2(1)$
- (P2) If  $\mathbf{X} \in \mathcal{E}_q(1)$  and  $\mathcal{G}$  is a  $\lambda_{\mathcal{G}}$ -Lipschitz map, then  $\mathcal{G}(\mathbf{X}) \in \mathcal{E}_q(\lambda_{\mathcal{G}})$

⇒ GAN data are concentrated vectors by design.

Remark: Still we need to control  $\lambda_{\mathcal{G}}$ .

## Control of $\lambda_{\mathcal{G}}$ with Spectral Normalization

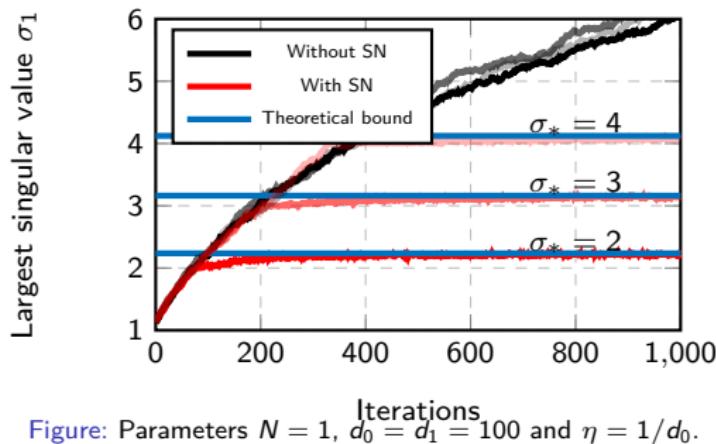
Let  $\sigma_* > 0$  and  $\mathcal{G}$  be a neural network composed of  $N$  affine layers, each one of input dimension  $d_{i-1}$  and output dimension  $d_i$  for  $i \in [N]$ , with 1-Lipschitz activation functions. Consider the following dynamics with learning rate  $\eta$ :

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \mathbf{E}, \text{ with } \mathbf{E}_{i,j} \sim \mathcal{N}(0, 1)$$

$$\mathbf{W} \leftarrow \mathbf{W} - \max(0, \sigma_1(\mathbf{W}) - \sigma_*) \mathbf{u}_1(\mathbf{W}) \mathbf{v}_1(\mathbf{W})^\top.$$

The Lipschitz constant of  $\mathcal{G}$  is bounded at convergence with high probability as:

$$\lambda_{\mathcal{G}} \leq \prod_{i=1}^N \left( \varepsilon + \sqrt{\sigma_*^2 + \eta^2 d_i d_{i-1}} \right).$$



# Model & Assumptions

(A1) Data matrix (distributed in  $k$  classes  $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_k$ ):

$$\mathbf{X} = \left[ \underbrace{\mathbf{x}_1, \dots, \mathbf{x}_{n_1}}_{\in \mathcal{E}_{q_1}(1)}, \underbrace{\mathbf{x}_{n_1+1}, \dots, \mathbf{x}_{n_2}}_{\in \mathcal{E}_{q_2}(1)}, \dots, \underbrace{\mathbf{x}_{n-n_k+1}, \dots, \mathbf{x}_n}_{\in \mathcal{E}_{q_k}(1)} \right] \in \mathbb{R}^{p \times n}$$

Model statistics:  $\mu_\ell = \mathbb{E}_{\mathbf{x}_i \in \mathcal{C}_\ell} [\mathbf{x}_i], \quad \mathbf{C}_\ell = \mathbb{E}_{\mathbf{x}_i \in \mathcal{C}_\ell} [\mathbf{x}_i \mathbf{x}_i^\top]$

(A2) Growth rate assumptions: As  $p \rightarrow \infty$ ,

1.  $p/n \rightarrow c \in (0, \infty)$ .
2. The number of classes  $k$  is bounded.
3. For any  $\ell \in [k]$ ,  $\|\mu_\ell\| = \mathcal{O}(\sqrt{p})$ .

Gram matrix and its resolvent:

$$\boxed{\mathbf{G} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}, \quad \mathbf{Q}(z) = (\mathbf{G} + z \mathbf{I}_n)^{-1}}$$

$$m_L(z) = \frac{1}{n} \text{tr}(\mathbf{Q}(-z)), \quad \mathbf{U} \mathbf{U}^\top = \frac{-1}{2\pi i} \oint_{\gamma} \mathbf{Q}(-z) dz$$

## Main Result

### Theorem

Under Assumptions **(A1)** and **(A2)**, we have  $\mathbf{Q}(z) \in \mathcal{E}_q(p^{-\frac{1}{2}})$ . Furthermore,

$$\left\| \mathbb{E}[\mathbf{Q}(z)] - \tilde{\mathbf{Q}}(z) \right\| = \mathcal{O} \left( \sqrt{\frac{\log p}{p}} \right) \text{ where } \tilde{\mathbf{Q}}(z) = \frac{1}{z} \Lambda(z) + \frac{1}{pz} \mathbf{J} \Omega(z) \mathbf{J}^T$$

with  $\Lambda(z) = \text{diag} \left\{ \frac{1_{n_\ell}}{1 + \delta_\ell(z)} \right\}_{\ell=1}^k$  and  $\Omega(z) = \text{diag}\{\mu_\ell^T \tilde{\mathbf{R}}(z) \mu_\ell\}_{\ell=1}^k$

$$\tilde{\mathbf{R}}(z) = \left( \frac{1}{k} \sum_{\ell=1}^k \frac{\mathbf{C}_\ell}{1 + \delta_\ell(z)} + z \mathbf{I}_p \right)^{-1}$$

with  $\delta(z) = [\delta_1(z), \dots, \delta_k(z)]$  is the unique fixed point of the system of equations

$$\delta_\ell(z) = \text{tr} \left( \mathbf{C}_\ell \left( \frac{1}{k} \sum_{j=1}^k \frac{\mathbf{C}_j}{1 + \delta_j(z)} + z \mathbf{I}_p \right)^{-1} \right) \text{ for each } \ell \in [k].$$

# Main Result

## Theorem

Under Assumptions **(A1)** and **(A2)**, we have  $\mathbf{Q}(z) \in \mathcal{E}_q(p^{-\frac{1}{2}})$ . Furthermore,

$$\left\| \mathbb{E}[\mathbf{Q}(z)] - \tilde{\mathbf{Q}}(z) \right\| = \mathcal{O} \left( \sqrt{\frac{\log p}{p}} \right) \text{ where } \tilde{\mathbf{Q}}(z) = \frac{1}{z} \Lambda(z) + \frac{1}{p z} \mathbf{J} \Omega(z) \mathbf{J}^\top$$

with  $\Lambda(z) = \text{diag} \left\{ \frac{\mathbf{1}_{n_\ell}}{1 + \delta_\ell(z)} \right\}_{\ell=1}^k$  and  $\Omega(z) = \text{diag} \{ \boldsymbol{\mu}_\ell^\top \tilde{\mathbf{R}}(z) \boldsymbol{\mu}_\ell \}_{\ell=1}^k$

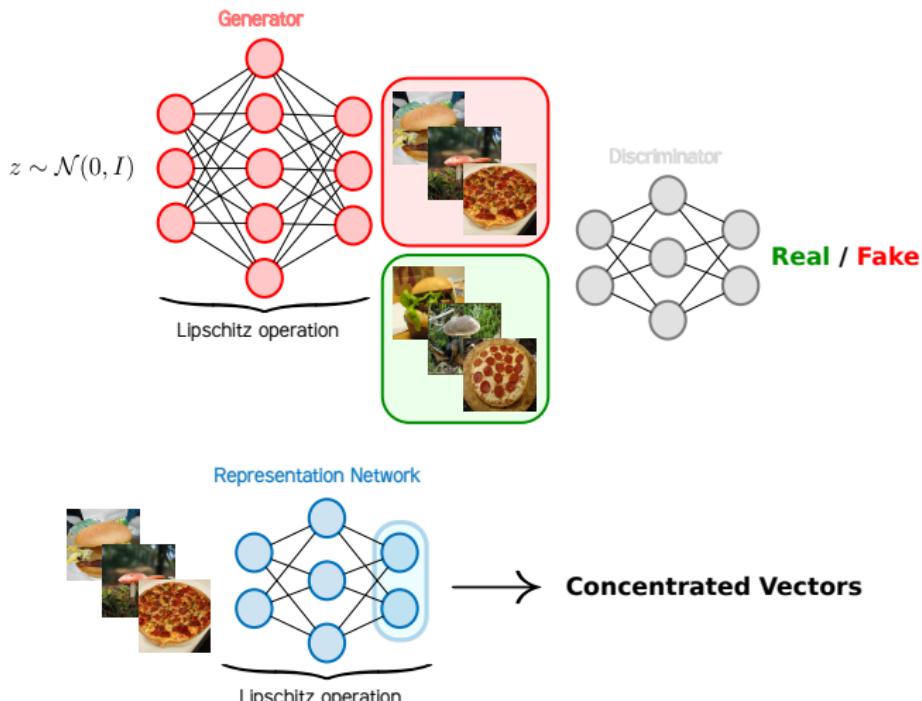
$$\tilde{\mathbf{R}}(z) = \left( \frac{1}{k} \sum_{\ell=1}^k \frac{\mathbf{C}_\ell}{1 + \delta_\ell(z)} + z \mathbf{I}_p \right)^{-1}$$

with  $\delta(z) = [\delta_1(z), \dots, \delta_k(z)]$  is the unique fixed point of the system of equations

$$\delta_\ell(z) = \text{tr} \left( \mathbf{C}_\ell \left( \frac{1}{k} \sum_{j=1}^k \frac{\mathbf{C}_j}{1 + \delta_j(z)} + z \mathbf{I}_p \right)^{-1} \right) \text{ for each } \ell \in [k].$$

**Key Observation:** Only **first** and **second** order statistics matter!

## Application to CNN Representations of GAN Images



- ▶ CNN representations correspond to the **penultimate** layer.
- ▶ Popular architectures considered in practice are: **Resnet**, **VGG**, **Densenet**.

# Application to CNN Representations of GAN Images

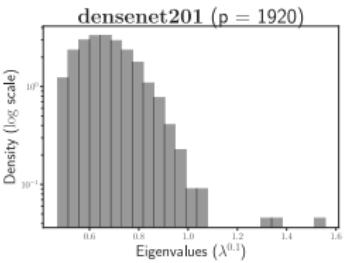
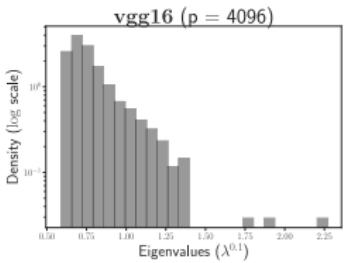
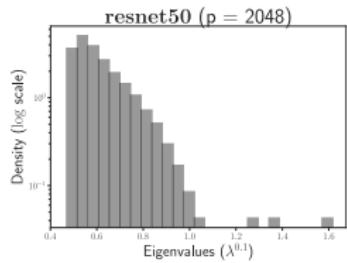
GAN Images



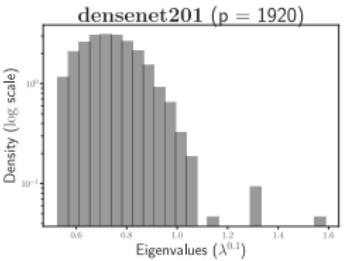
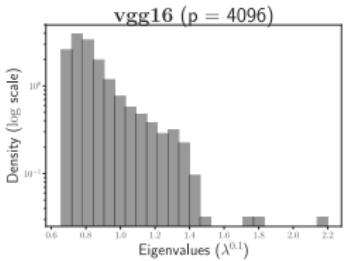
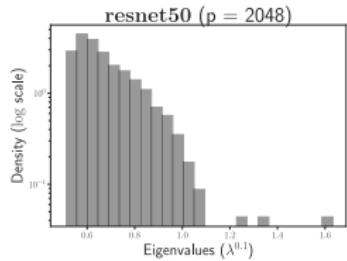
Figure:  $k = 3$  classes,  $n = 3000$  images.

# Application to CNN Representations of GAN Images

## GAN Images

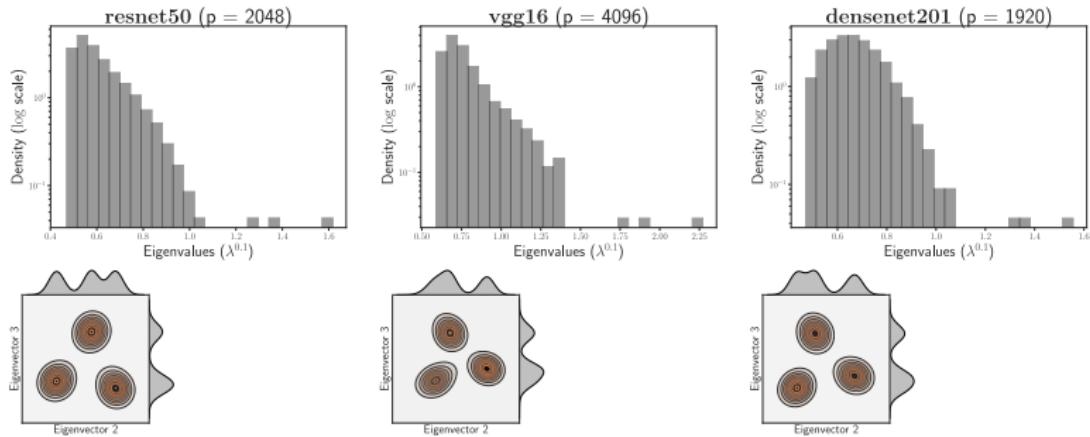


## Real Images

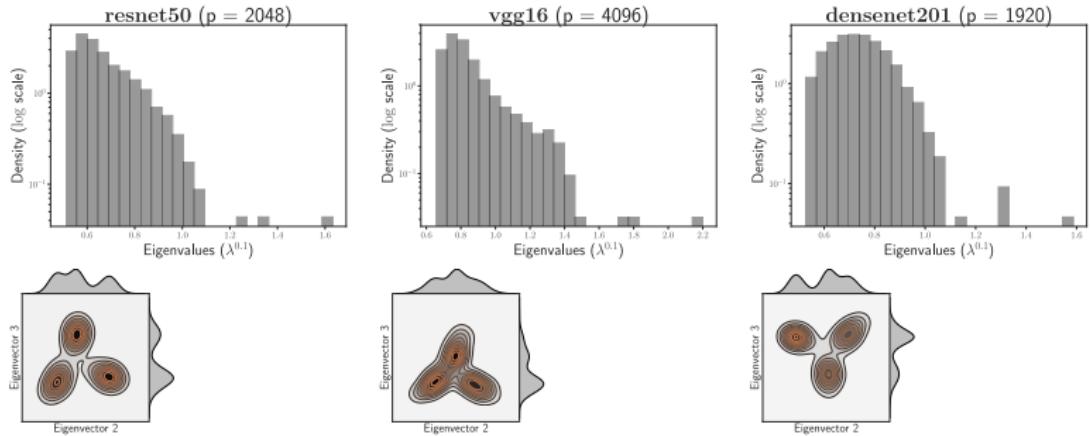


# Application to CNN Representations of GAN Images

## GAN Images

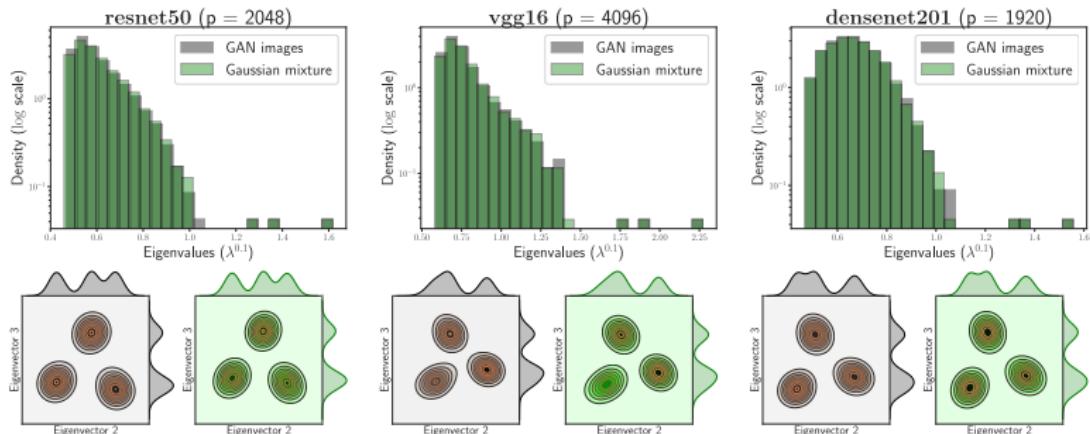


## Real Images

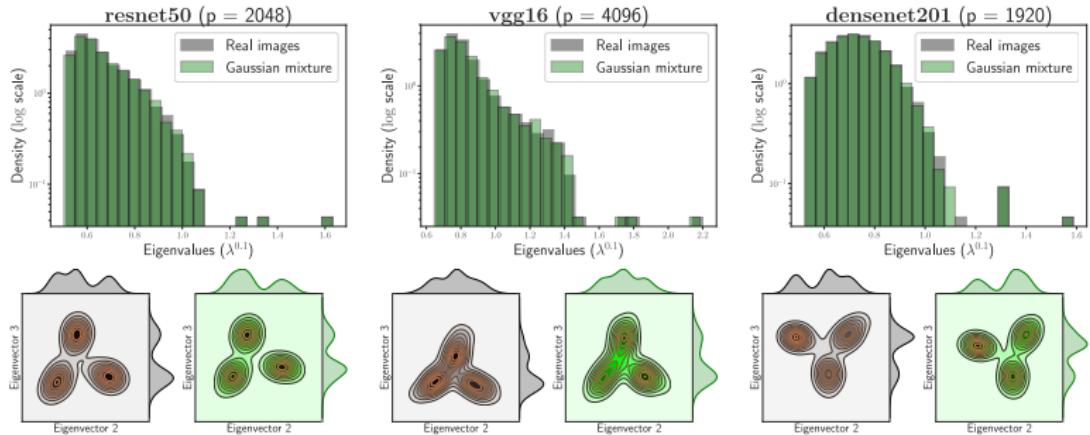


# Application to CNN Representations of GAN Images

## GAN Images

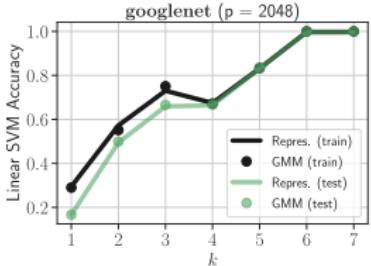
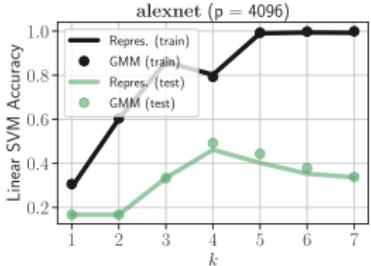
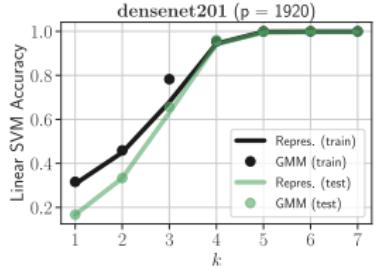
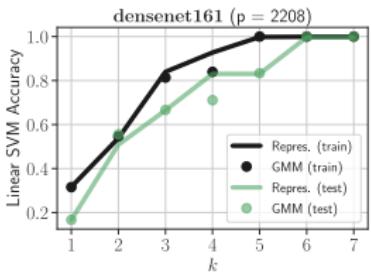
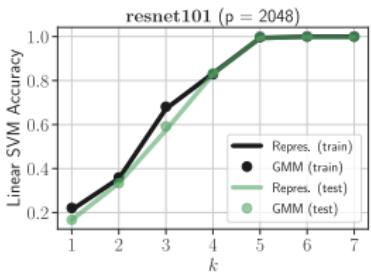
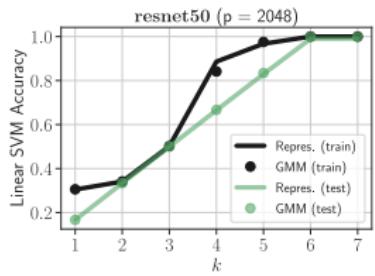
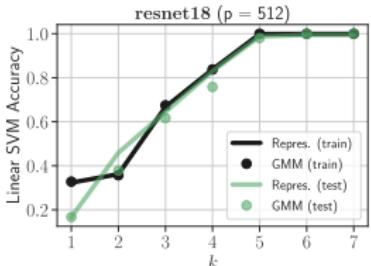
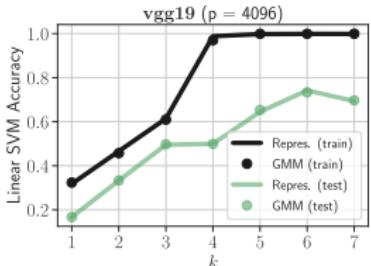
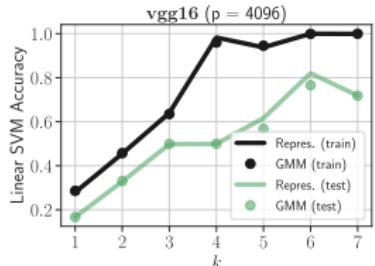


## Real Images



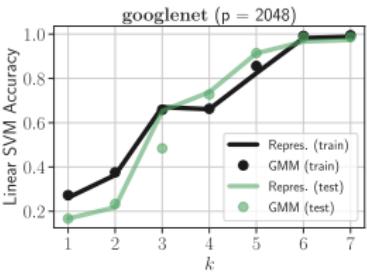
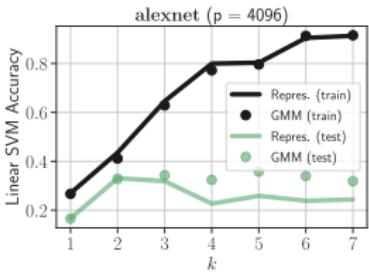
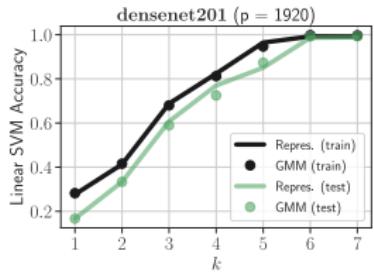
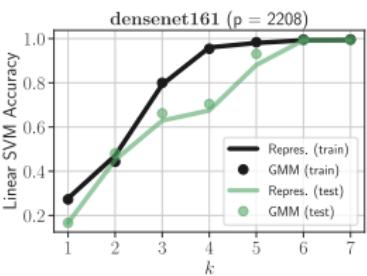
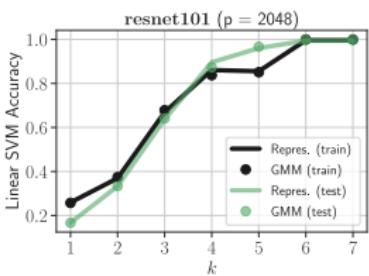
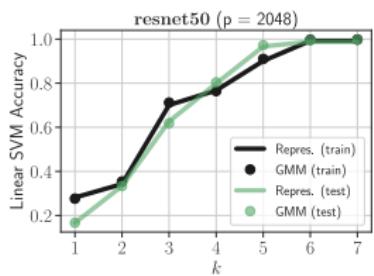
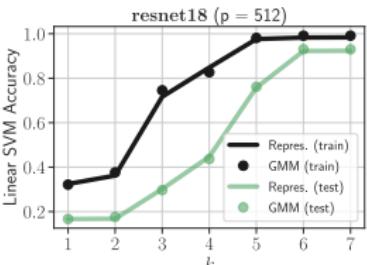
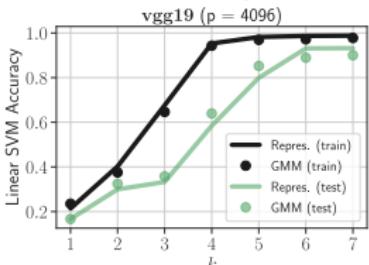
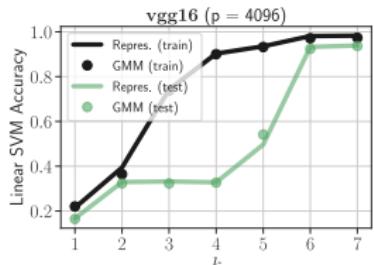
# Performance of a linear SVM classifier

## GAN Images



# Performance of a linear SVM classifier

## Real Images



## Take away messages

- ▶ Concentrated Vectors seem appropriate for realistic data modelling.
- ▶ Universality of linear classifiers regardless of the data distribution.
- ▶ RMT can anticipate the performances of standard classifiers for DL representations of GAN images.
- ▶ Universality supports the Gaussianity assumption on the data representations as considered in the literature, e.g., the FID metric

$$d^2((\mu, \mathbf{C}), (\mu_w, \mathbf{C}_w)) = \|\mu - \mu_w\|^2 + \text{tr} \left( \mathbf{C} + \mathbf{C}_w - 2(\mathbf{C}\mathbf{C}_w)^{\frac{1}{2}} \right).$$