

# **Noncommutative probability for deep learning**

**Statistics Seminar  
Department of Mathematics  
The University of Mississippi**

Cong Zhou October 14, 2021

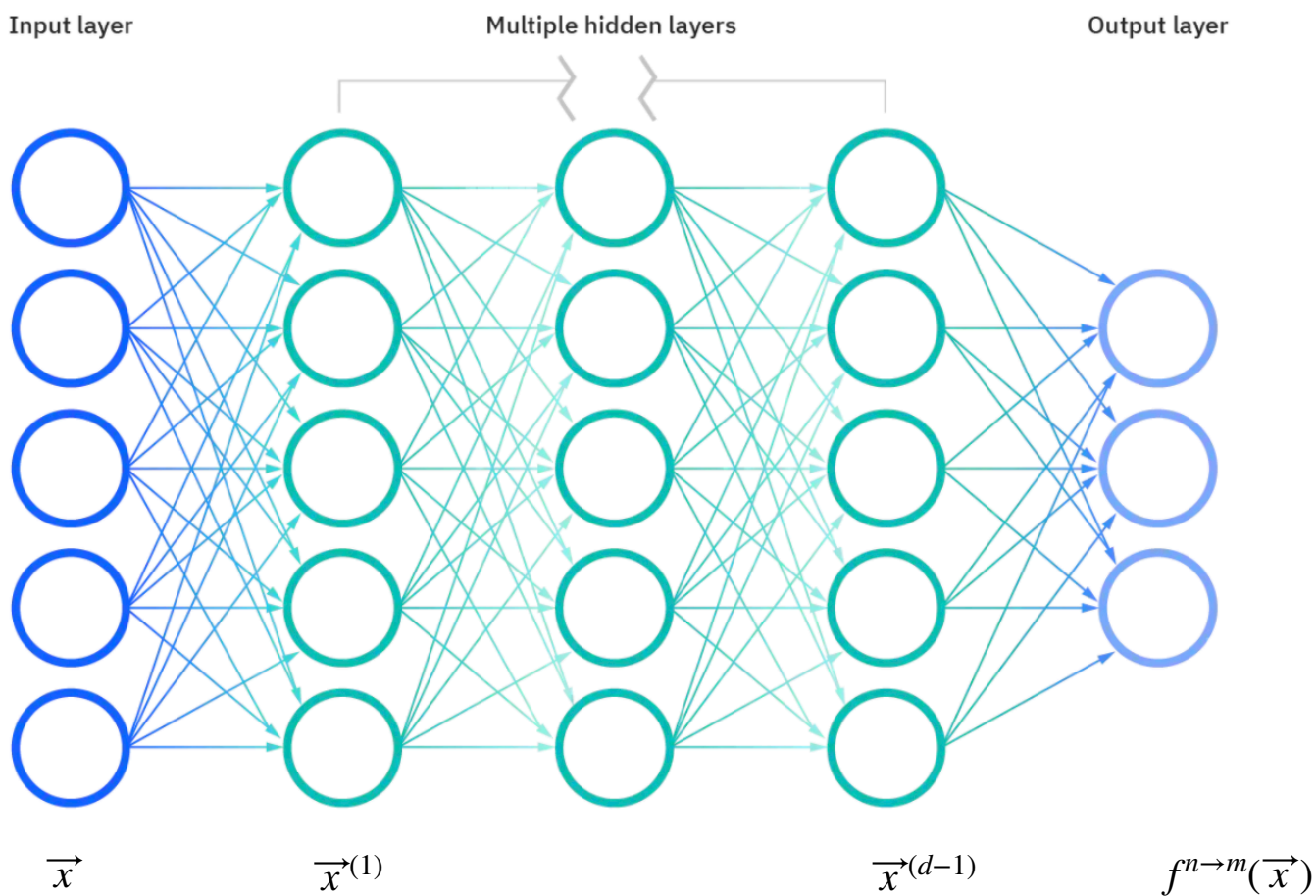
# Outline

- 1. Background and motivation**
- 2. Connection to random matrices**
- 3. Essentials of random matrix and free probability theory**
- 4. Condition for fast training**
- 5. Conclusions**

# Background and motivation

# What is deep learning?

## Deep neural network



$$\vec{x}^{(0)} = \vec{x}$$

$$Z^{(1)} = W^{(1)}\vec{x}^{(0)} + b^{(1)}$$

$$\vec{x}^{(1)} = \phi(Z^{(1)})$$

⋮

⋮

⋮

$$\vec{x}^{(d)} = \phi(Z^{(d)})$$

$$Z^{(d+1)} = W^{(d+1)}\vec{x}^{(d)} + b^{(d+1)}$$

$$f^{n \rightarrow m}(\vec{x}) = Z^{(d+1)}$$

Weights	$W^{(i)}$
Bias	$b^{(i)}$
Parameter vector	$\theta = (W^{(i)}, b^{(i)})$
Activation function	$\phi$
Pre-activation	$Z^{(i)}$
Post-activation	$\vec{x}^{(i)}$
Input data	$\vec{x} = \vec{x}^{(0)}$
Output data	$f^{n \rightarrow m}(\vec{x}) = Z^{(d+1)}$
Targets	$\vec{y}$
Loss function (MSE)	$\text{Error}(W, b) = \frac{1}{2m} \ f^{n \rightarrow m}(\vec{x}) - \vec{y}\ ^2$

$$\begin{aligned}
 \vec{x}^{(0)} &= \vec{x} \\
 Z^{(1)} &= W^{(1)} \vec{x}^{(0)} + b^{(1)} \\
 \vec{x}^{(1)} &= \phi(Z^{(1)}) \\
 &\vdots \\
 &\vdots \\
 &\vdots \\
 \vec{x}^{(d)} &= \phi(Z^{(d)}) \\
 Z^{(d+1)} &= W^{(d+1)} \vec{x}^{(d)} + b^{(d+1)} \\
 f^{n \rightarrow m}(\vec{x}) &= Z^{(d+1)}
 \end{aligned}$$

How to find the parameters?

### Supervised learning Problem

Given examples  $\{\vec{x}_\varepsilon\}_{\varepsilon \in \text{Examples}} \subset \mathbb{R}^n$ , and labels  $\{\vec{y}_\varepsilon\}_{\varepsilon \in \text{Examples}} \subset \mathbb{R}^m$ . How to find parameters  $\underline{W}^{(i)}$  and  $\vec{b}^{(i)}$  so that  $f^{n \rightarrow m}$  minimizes the prediction error:

$$\text{Error}(W, b) = \sum_{\varepsilon \in \text{Examples}} \|f^{n \rightarrow m}(\vec{x}_\varepsilon) - \vec{y}_\varepsilon\|^2$$

### Supervised Learning: Solution Idea

0. **Invent** the architecture: depth  $d$  and layer width  $n_1, \dots, n_{d-1}$ . Set  $n_0 = n, n_d = m$ .

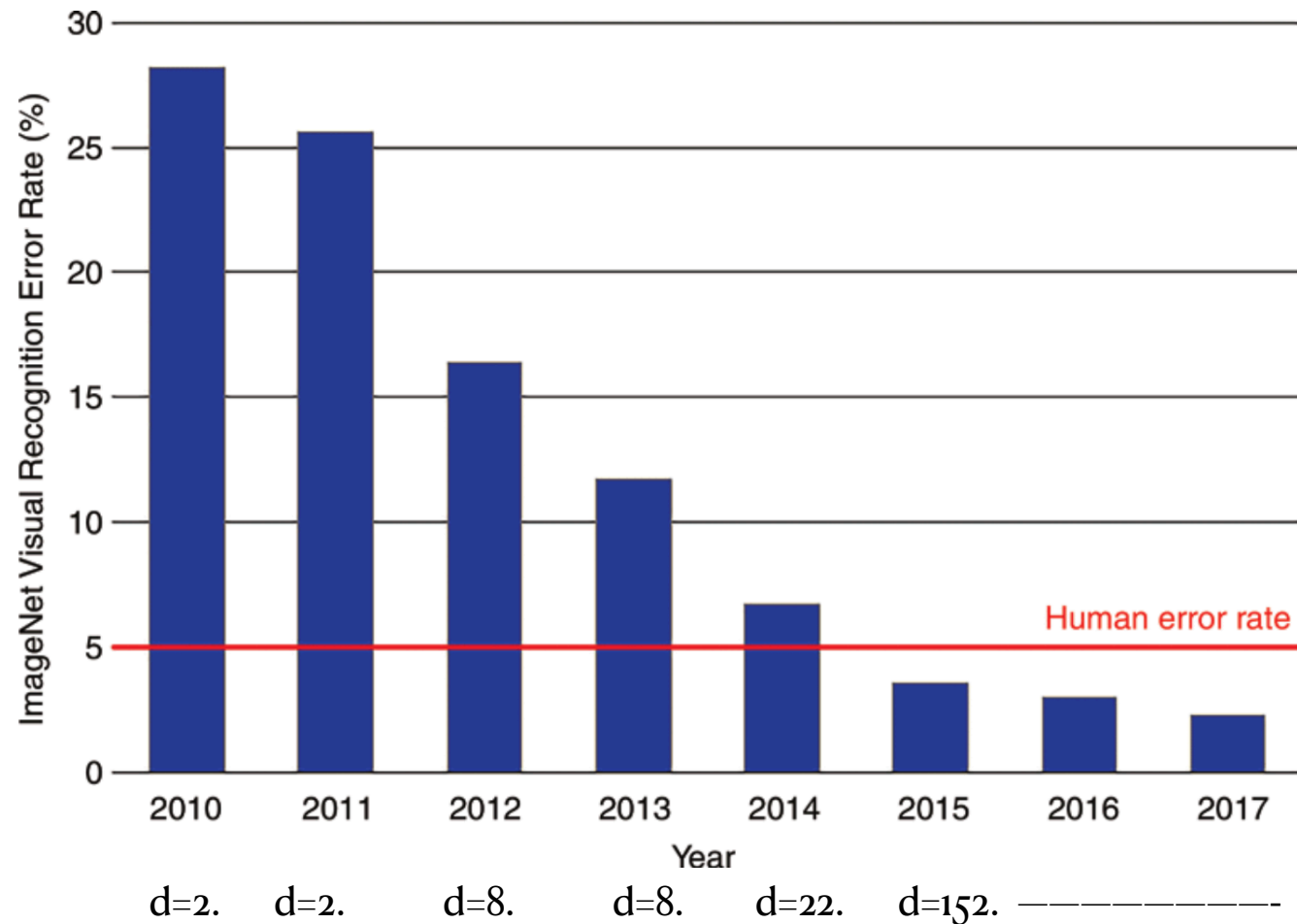
1. **Initialization**: Pick parameters  $\underline{W}, \vec{b}$  at **random**.

2. **Modify** parameters to **shrink**  $\text{Error}(W, b)$  by **gradient descent**.

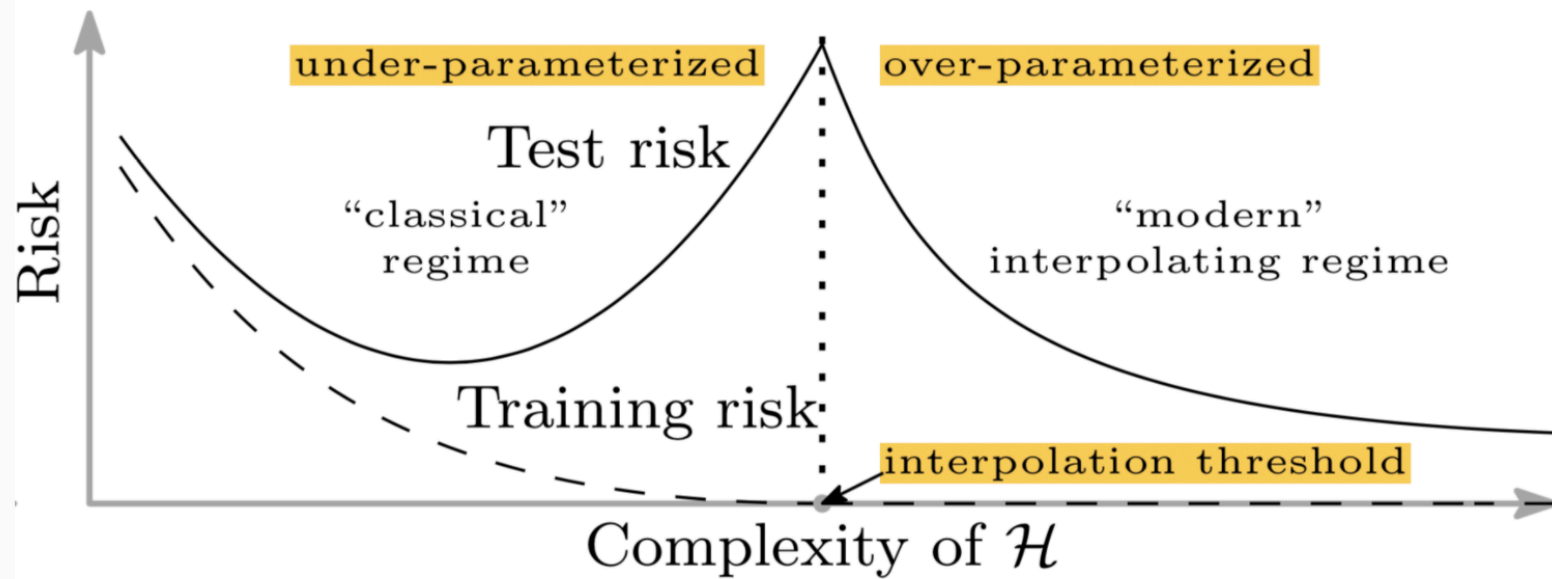
$$\text{new } W_{j,k}^{(i)} := \text{old } W_{j,k}^{(i)} - \partial_{W_{j,k}^{(i)}} \text{Error}(W, b)$$

3. Repeat step 2 many time. **Hope** that the error is now small.

# ImageNet Large Scale Visual Recognition Challenge Results



## Double descent phenomenon





**Empirical Issue:**  $d$  large has **vanishing and exploding gradients**.

**Random initialization:**  $\partial_{W_{j,k}^{(i)}} \text{Error}(W, b)$  is very large or small.

**Today's story is about what happens on the initialization.**

Connection to random matrices

## Part 2: Products of random matrices

- Connect to neural network
- Limit theories for products of random matrices.

**Definition.** The **Input-Output Jacobian** matrix  $\underline{J} = \text{Jac}\{f^{n_0 \rightarrow n_d}\}$  is the  $n_d \times n_0$  matrix

$$J_{ij}(\vec{x}) := \partial_j f_i^{n_0 \rightarrow n_d}(\vec{x})$$

*Remark:*  $\partial_{W_{j,k}^{(i)}} \text{Error}(W, b)$  can be written in terms of  $\underline{J}$ .

$J$  when  $\phi(x) = \max\{x, 0\}$

Recall  $f^{n_{i-1} \rightarrow n_i}(\vec{x}) := \phi(\underline{W}^{(i)}\vec{x} + \vec{b}^{(i)})$  and want to compute

$$\underline{J} = \text{Jac}(f^{n_{d-1} \rightarrow n_d} \circ f^{n_{d-2} \rightarrow n_{d-1}} \circ \dots \circ f^{n_1 \rightarrow n_0})$$

Since  $\phi(x) = \max\{x, 0\}$ ,  $\phi'(x) = 1\{x > 0\}$ , then the gradient of each layer is

$$\text{Jac}(f^{n_{i-1} \rightarrow n_i}) = \text{Diag}\left(1\{\underline{W}^{(i)}\vec{x} + \vec{b}^{(i)} > 0\}\right) \underline{W}^{(i)}.$$

Assume all **random matrices** are **symmetric**:

$$\text{Diag}\left(1\{\underline{W}^{(i)}\vec{x} + \vec{b}^{(i)} > 0\}\right) \stackrel{d}{=} \text{Diag}\left(\vec{X}^{(i)}\right)$$

where  $\vec{X}^{(i)} \in \mathbb{R}^n$  has i.i.d entries  $X_j^{(i)} \sim \text{Bernoulli}\left(\frac{1}{2}\right)$ .

By the chain rule, we shall expect

$$\underline{J} \stackrel{?}{=} \underbrace{\text{Diag}\left(\vec{X}^{(d)}\right) W^{(d)} \dots \text{Diag}\left(\vec{X}^{(1)}\right) W^{(1)}}_{\underline{M}}$$

Could show  $\underline{M} \stackrel{d}{=} \underline{J}$  up to conjugation by random  $\pm 1$  Bernoulli's.

# Essentials of random matrix and free probability theory



Dan-Virgil Voiculescu



Hari Bercovici

## Spectral density

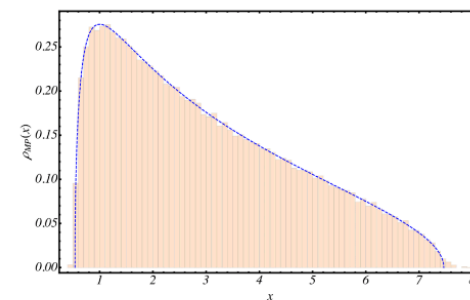
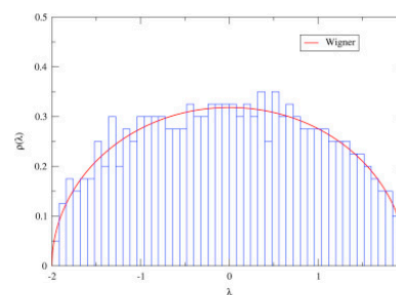
Let  $(A_n)_{n=1}^{\infty}$  be a sequence of  $n \times n$  matrix.

The empirical spectral density of  $\frac{1}{\sqrt{n}}A_n$  is:

$$\mu_{\frac{1}{\sqrt{n}}A_n}(\lambda) = \frac{1}{n} \sum_{j=1}^n \delta(\lambda - \lambda_j(A_n)).$$

For a sequence of matrices with increasing size,  $(A_n)_{n=1}^{\infty}$ , the limiting spectral density is:

$$\mu(\lambda) = \lim_{n \rightarrow \infty} \mu_{\frac{1}{\sqrt{n}}A_n}(\lambda).$$





For  $z \in \mathbb{C} \setminus \text{supp}(\mu)$  the Cauchy transform is defined as:

$$G(z) = \sum_{n=1}^{\infty} \frac{m_n}{z^{n+1}}, \quad m_n = \int_{\text{supp}(\mu)} t^n d\mu(t)$$

The moment generating function,  $M(z)$ , defined by a functional equation,

$$M(z) = \sum_{n=1}^{\infty} m_n z^n$$

The spectral density can be recovered from  $G$  using the Stieltjes inversion formula,

$$\mu(\lambda) = -\frac{1}{\pi} \lim_{\varepsilon \rightarrow 0^+} \Im G(\lambda + i\varepsilon).$$

The Cauchy transform  $G$  relates to the  $R$ -transform,  $R$ , defined by the functional equation,

$$1 + R(G(z)) = zG(z).$$

The moment series  $M$  and  $R$ -transform has the following relationship:

$$R(z(1 + M(z))) = M(z).$$

and the  $S$ -transform,

$$S(zG(z) - 1) = \frac{G(z)}{zG(z) - 1}$$

If  $a$  and  $b$  are **free independent**, then the distribution of sum  $a + b$  can be computed using the  $R$ -transform:

$$\begin{array}{ccc}
 \mu_a(\lambda) \rightarrow G_a(z) \rightarrow R_a & \searrow & \\
 & & R_a + R_b = R_{a+b} \rightarrow G_{a+b}(z) \rightarrow \mu_{a+b}(\lambda) \\
 & \nearrow & \\
 \mu_b(\lambda) \rightarrow G_b(z) \rightarrow R_b & & 
 \end{array}$$

S-transform:

$$\begin{array}{ccc}
 \mu_a(\lambda) \rightarrow G_a(z) \rightarrow S_a & \searrow & \\
 & & S_a S_b = S_{ab} \rightarrow G_{ab}(z) \rightarrow \mu_{ab}(\lambda) \\
 & \nearrow & \\
 \mu_b(\lambda) \rightarrow G_b(z) \rightarrow S_b & & 
 \end{array}$$

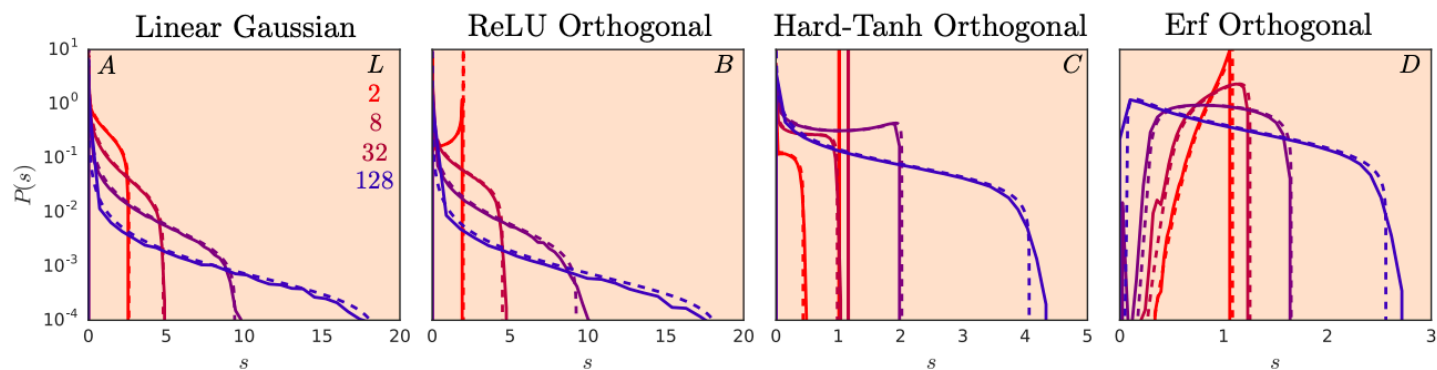
Classical independence	Free independence
$X, Y$ are independent if one has $\mathbf{E}f(X)g(Y) = 0$ whenever $f$ and $g$ are such that $\mathbf{E}f(X) = 0$ $\mathbf{E}g(Y) = 0$	$X, Y$ are freely independent if $\tau(f_1(X)g_1(Y) \dots f_k(X)g_k(Y)) = 0$ whenever $f_i$ and $g_i$ are such that $\tau(f_i(X)) = 0$ $\tau(g_i(Y)) = 0$
	$\ast$ -free is $\{X, X^*\}$ and $\{Y, Y^*\}$ are free.

Define the  $n_d \times n_0$  product random matrix  $\underline{M}$

$$\underline{M} = \text{Diag} \left( \vec{X}^{(d)} \right) W^{(d)} \dots \text{Diag} \left( \vec{X}^{(1)} \right) W^{(1)}$$

where  $\vec{X}^{(i)} \in \mathbb{R}^n$  has i.i.d entries  $X_j^{(i)} \sim \text{Bernoulli} \left( \frac{1}{2} \right)$ .

$W$  and  $\text{Diag} \left( \vec{X} \right)$  are asymptotically free [Yang '19].



Condition for fast training

## Conditions for trainability

**The back-propagated signals (error signals) should not explode/vanish with depth.**

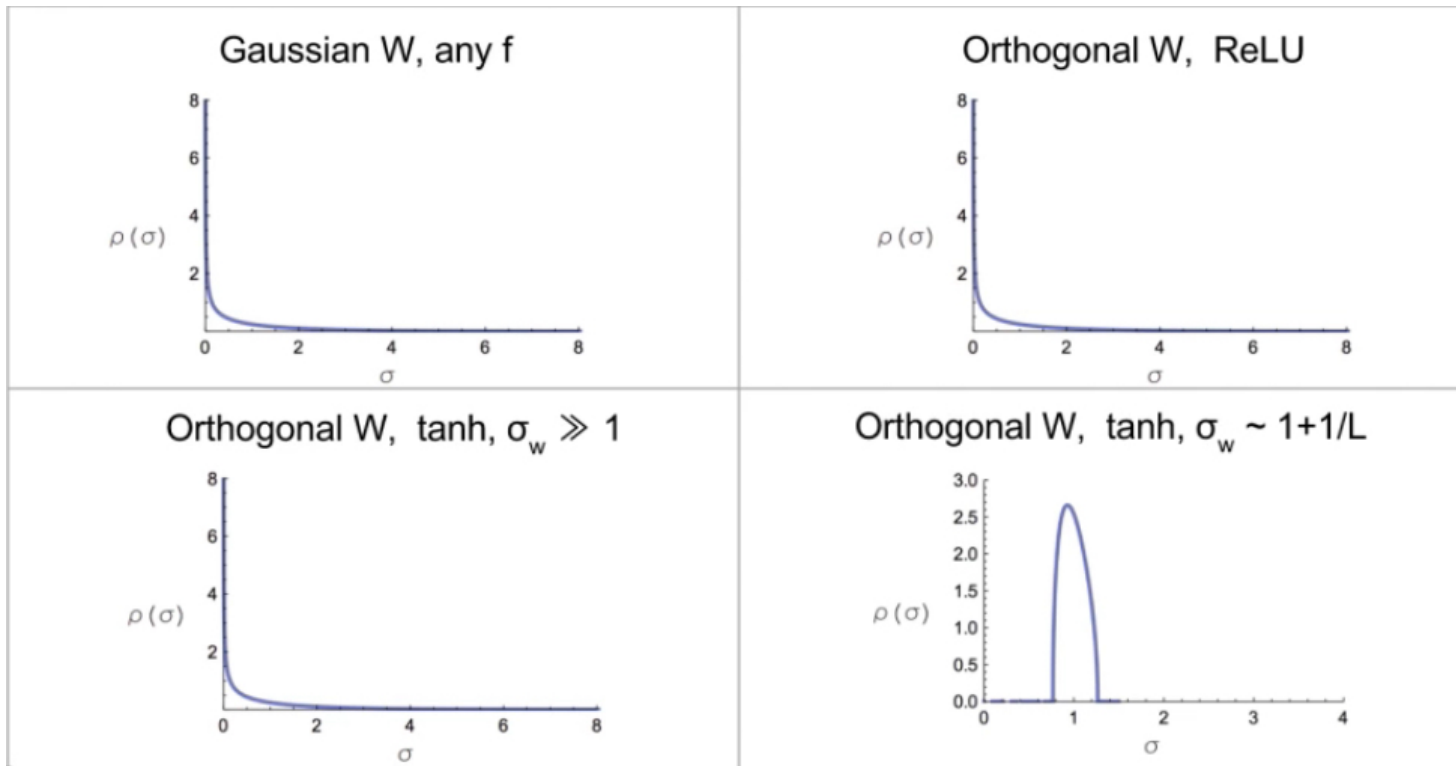
Along the critical line, the gradient norms are well-behaved on average:

$$\delta^1 = \epsilon^T J, \quad \epsilon = \delta^d \text{ (error signal)}$$

$$\underline{J} = \text{Diag} \left( \overrightarrow{X}^{(d)} \right) W^{(d)} \dots \text{Diag} \left( \overrightarrow{X}^{(1)} \right) W^{(1)}$$
$$\mathbb{E}_{\epsilon} [\|\delta^1\|^2] = \mathbb{E}_{\epsilon} [\epsilon^T J J^T \epsilon] = \text{tr}[J J^T] = \sum_{i=1}^n \sigma_i^2$$
$$\frac{1}{n} \sum_{i=1}^n \sigma_i^2 = 1 \Leftarrow \text{trainable}$$

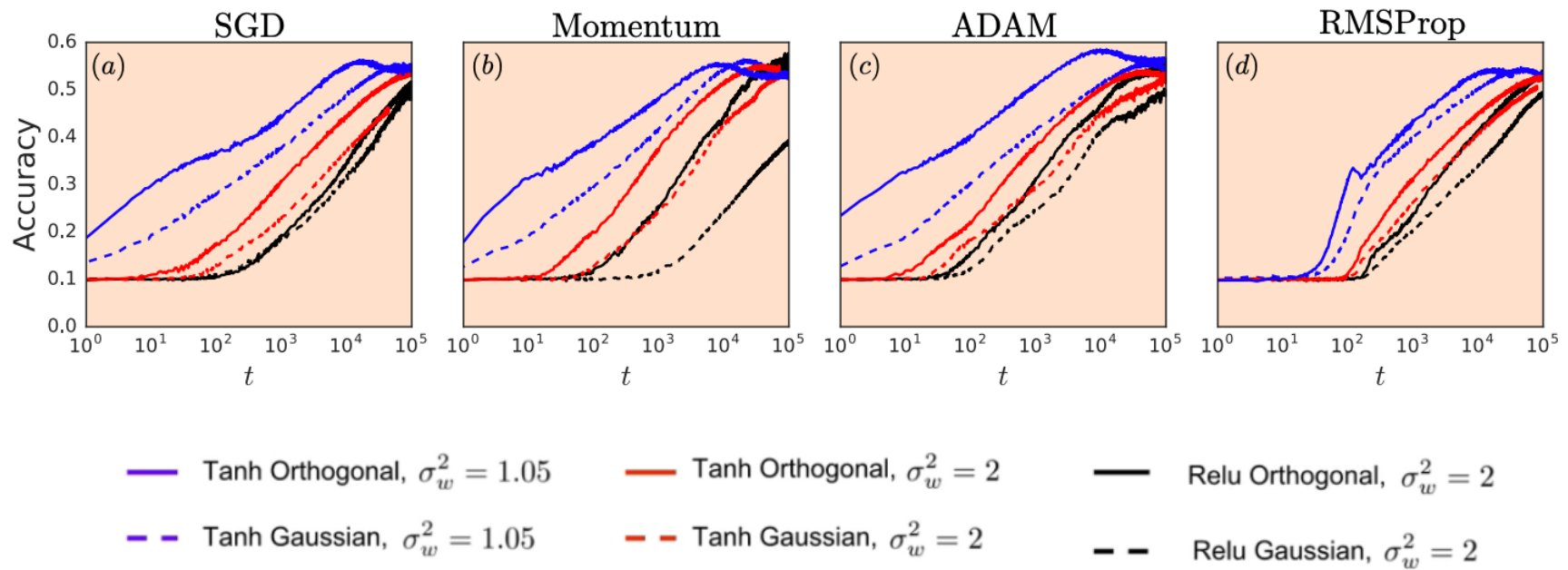
$$\underline{J} = \text{Diag} \left( \overrightarrow{X}^{(d)} \right) W^{(d)} \dots \text{Diag} \left( \overrightarrow{X}^{(1)} \right) W^{(1)}$$

Jacobian spectra for large depth





## Implications for training speed



# Conclusions

Conclusion:

Free probability and random matrix theory provides powerful tools for studying deep learning:

- Jacobian (showed how to equilibrate the distribution of singular values of the input-output Jacobian for faster training).
- Hessian

Thank you!