

# Random dynamical systems in deep neural networks

Maximilian Engel  
KdVI (UvA) and FU Berlin

FoMI: A Seminar on the Foundations of Machine Intelligence  
April 29, 2026

# (Random) Dynamical Systems and Neural Networks

**Artificial Neural Network** as a function

$$\Phi : \mathbb{R}^D \times \mathbb{R}^d \rightarrow \mathbb{R}^q, \quad (x, y) \mapsto \Phi(x, y).$$

# (Random) Dynamical Systems and Neural Networks

**Artificial Neural Network** as a function

$$\Phi : \mathbb{R}^D \times \mathbb{R}^d \rightarrow \mathbb{R}^q, \quad (x, y) \mapsto \Phi(x, y).$$

- ▶ **Training data:**  $(y_1, z_1), \dots, (y_N, z_N) \in \mathbb{R}^d \times \mathbb{R}^q$

# (Random) Dynamical Systems and Neural Networks

**Artificial Neural Network** as a function

$$\Phi : \mathbb{R}^D \times \mathbb{R}^d \rightarrow \mathbb{R}^q, \quad (x, y) \mapsto \Phi(x, y).$$

- ▶ **Training data:**  $(y_1, z_1), \dots, (y_N, z_N) \in \mathbb{R}^d \times \mathbb{R}^q$
- ▶ **Training risk (loss):**  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N |\Phi(x, y_i) - z_i|^2$

# (Random) Dynamical Systems and Neural Networks

**Artificial Neural Network** as a function

$$\Phi : \mathbb{R}^D \times \mathbb{R}^d \rightarrow \mathbb{R}^q, \quad (x, y) \mapsto \Phi(x, y).$$

- ▶ **Training data:**  $(y_1, z_1), \dots, (y_N, z_N) \in \mathbb{R}^d \times \mathbb{R}^q$
- ▶ **Training risk (loss):**  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N |\Phi(x, y_i) - z_i|^2$

Two dynamical problem classes [CHEMNITZ/E./KUEHN/KUNTZ 2025+]:

# (Random) Dynamical Systems and Neural Networks

**Artificial Neural Network** as a function

$$\Phi : \mathbb{R}^D \times \mathbb{R}^d \rightarrow \mathbb{R}^q, \quad (x, y) \mapsto \Phi(x, y).$$

- ▶ **Training data:**  $(y_1, z_1), \dots, (y_N, z_N) \in \mathbb{R}^d \times \mathbb{R}^q$
- ▶ **Training risk (loss):**  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N |\Phi(x, y_i) - z_i|^2$

Two dynamical problem classes [CHEMNITZ/E./KUEHN/KUNTZ 2025+]:

- ▶ **A) Dynamics of the Network** (Learning Process)

$$\varphi : \mathbb{R}^D \rightarrow \mathbb{R}^D, \quad x \mapsto \varphi(x) = x - \eta \nabla \mathcal{L}(x) \quad (\text{GD})$$

$$\hat{\varphi} : \Omega \times \mathbb{R}^D \rightarrow \mathbb{R}^D, \quad x \mapsto \hat{\varphi}_\omega(x) = x - \eta \nabla \hat{\mathcal{L}}_\omega(x) \quad (\text{SGD})$$

# (Random) Dynamical Systems and Neural Networks

**Artificial Neural Network** as a function

$$\Phi : \mathbb{R}^D \times \mathbb{R}^d \rightarrow \mathbb{R}^q, \quad (x, y) \mapsto \Phi(x, y).$$

- ▶ **Training data:**  $(y_1, z_1), \dots, (y_N, z_N) \in \mathbb{R}^d \times \mathbb{R}^q$
- ▶ **Training risk (loss):**  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N |\Phi(x, y_i) - z_i|^2$

Two dynamical problem classes [CHEMNITZ/E./KUEHN/KUNTZ 2025+]:

- ▶ **A) Dynamics of the Network** (Learning Process)

$$\varphi : \mathbb{R}^D \rightarrow \mathbb{R}^D, \quad x \mapsto \varphi(x) = x - \eta \nabla \mathcal{L}(x) \quad (\text{GD})$$

$$\hat{\varphi} : \Omega \times \mathbb{R}^D \rightarrow \mathbb{R}^D, \quad x \mapsto \hat{\varphi}_\omega(x) = x - \eta \nabla \hat{\mathcal{L}}_\omega(x) \quad (\text{SGD})$$

- ▶ **B) Dynamics on the Network** (Information Propagation)

$$\Phi(x, \cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^q \quad \text{or} \quad \hat{\Phi}_\omega(x, \cdot) : \Omega \times \mathbb{R}^d \rightarrow \mathbb{R}^q. \quad (\text{Neural map})$$

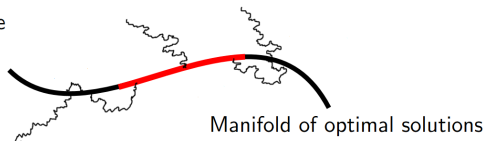
## Key questions for this talk

# Key questions for this talk

## A) Dynamics of the Network

- ▶ Characterization of parameter solutions learned via (S)GD?
- ▶ Related to **generalization** from training to new data sets?

Parameter space

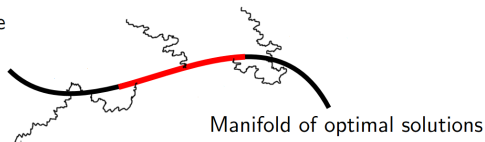


# Key questions for this talk

## A) Dynamics of the Network

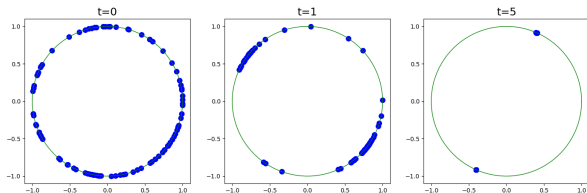
- ▶ Characterization of parameter solutions learned via (S)GD?
- ▶ Related to **generalization** from training to new data sets?

Parameter space



## B) Dynamics on the Network

- ▶ Mechanisms for clustering in **transformers** of LLMs?



## Random Dynamical Systems (RDS)

**Random dynamical system** (RDS)  $(\theta, \varphi)$  [ARNOLD 1998] consists of

## Random Dynamical Systems (RDS)

**Random dynamical system** (RDS)  $(\theta, \varphi)$  [ARNOLD 1998] consists of

- ▶ dynamical system  $(\theta_t)_{t \in \mathbb{T}}$  on probability space  $(\Omega, \mathcal{F}, \mathbb{P})$

## Random Dynamical Systems (RDS)

**Random dynamical system** (RDS)  $(\theta, \varphi)$  [ARNOLD 1998] consists of

- ▶ dynamical system  $(\theta_t)_{t \in \mathbb{T}}$  on probability space  $(\Omega, \mathcal{F}, \mathbb{P})$
- ▶ and *cocycle*  $(\varphi_\omega^t : X \rightarrow X)_{t \in \mathbb{T}}$  such that for all  $\omega \in \Omega$  and  $x \in X$

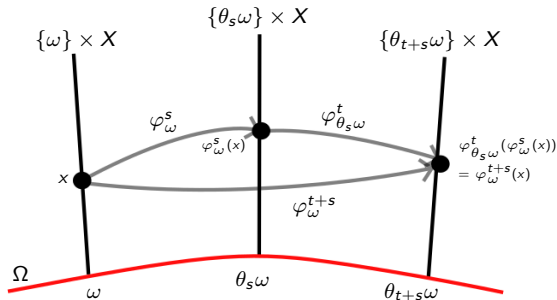
$$\varphi_\omega^{t+s}(x) = \varphi_{\theta_s \omega}^t(\varphi_\omega^s(x)).$$

# Random Dynamical Systems (RDS)

**Random dynamical system (RDS)**  $(\theta, \varphi)$  [ARNOLD 1998] consists of

- ▶ dynamical system  $(\theta_t)_{t \in \mathbb{T}}$  on probability space  $(\Omega, \mathcal{F}, \mathbb{P})$
- ▶ and *cocycle*  $(\varphi_t^t : X \rightarrow X)_{t \in \mathbb{T}}$  such that for all  $\omega \in \Omega$  and  $x \in X$

$$\varphi_\omega^{t+s}(x) = \varphi_{\theta_s \omega}^t(\varphi_\omega^s(x)).$$



## RDS in neural networks

**A) Stochastic Gradient Descent (SGD):**

### A) Stochastic Gradient Descent (SGD):

- ▶ Write  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$ , where  $\mathcal{L}_i(x) := |\Phi(x, y_i) - z_i|^2$

### A) Stochastic Gradient Descent (SGD):

- ▶ Write  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$ , where  $\mathcal{L}_i(x) := |\Phi(x, y_i) - z_i|^2$
- ▶ Then, for  $\omega = (\omega_1, \omega_2, \dots) \in \Omega$  uniform i.i.d. r.v.,  $(\theta\omega)_i = \omega_{i+1}$ ,

$$\varphi_{\omega}^{n+1}(x) = \varphi_{\omega}^n(x) - \eta \nabla \mathcal{L}_{\omega_{n+1}}(\varphi_{\omega}^n(x)).$$

### A) Stochastic Gradient Descent (SGD):

- ▶ Write  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$ , where  $\mathcal{L}_i(x) := |\Phi(x, y_i) - z_i|^2$
- ▶ Then, for  $\omega = (\omega_1, \omega_2, \dots) \in \Omega$  uniform i.i.d. r.v.,  $(\theta\omega)_i = \omega_{i+1}$ ,

$$\varphi_\omega^{n+1}(x) = \varphi_\omega^n(x) - \eta \nabla \mathcal{L}_{\omega_{n+1}}(\varphi_\omega^n(x)).$$

### B) Transformers as Stochastic Differential Equations (SDEs)

$$dX_t = F_0(X_t) dt + \sum_{j=1}^m F_j(X_t) \partial W_t^j, \quad X_0 = x \in \mathbb{S}^{n-1}. \quad (\text{SDE})$$

### A) Stochastic Gradient Descent (SGD):

- ▶ Write  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$ , where  $\mathcal{L}_i(x) := |\Phi(x, y_i) - z_i|^2$
- ▶ Then, for  $\omega = (\omega_1, \omega_2, \dots) \in \Omega$  uniform i.i.d. r.v.,  $(\theta\omega)_i = \omega_{i+1}$ ,

$$\varphi_\omega^{n+1}(x) = \varphi_\omega^n(x) - \eta \nabla \mathcal{L}_{\omega_{n+1}}(\varphi_\omega^n(x)).$$

### B) Transformers as Stochastic Differential Equations (SDEs)

$$dX_t = F_0(X_t) dt + \sum_{j=1}^m F_j(X_t) \partial W_t^j, \quad X_0 = x \in \mathbb{S}^{n-1}. \quad (\text{SDE})$$

- ▶  $\Omega = C_0(\mathbb{R}, \mathbb{R}^m)$  with Wiener measure  $\mathbb{P}$ ,  $\theta_t \omega(\cdot) := \omega(t + \cdot) - \omega(t)$

### A) Stochastic Gradient Descent (SGD):

- ▶ Write  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$ , where  $\mathcal{L}_i(x) := |\Phi(x, y_i) - z_i|^2$
- ▶ Then, for  $\omega = (\omega_1, \omega_2, \dots) \in \Omega$  uniform i.i.d. r.v.,  $(\theta\omega)_i = \omega_{i+1}$ ,

$$\varphi_\omega^{n+1}(x) = \varphi_\omega^n(x) - \eta \nabla \mathcal{L}_{\omega_{n+1}}(\varphi_\omega^n(x)).$$

### B) Transformers as Stochastic Differential Equations (SDEs)

$$dX_t = F_0(X_t) dt + \sum_{j=1}^m F_j(X_t) \partial W_t^j, \quad X_0 = x \in \mathbb{S}^{n-1}. \quad (\text{SDE})$$

- ▶  $\Omega = C_0(\mathbb{R}, \mathbb{R}^m)$  with Wiener measure  $\mathbb{P}$ ,  $\theta_t \omega(\cdot) := \omega(t + \cdot) - \omega(t)$
- ▶ The solution of (SDE) can be written as cocycle  $\varphi_\omega^t(x)$ .

### A) Stochastic Gradient Descent (SGD):

- ▶ Write  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$ , where  $\mathcal{L}_i(x) := |\Phi(x, y_i) - z_i|^2$
- ▶ Then, for  $\omega = (\omega_1, \omega_2, \dots) \in \Omega$  uniform i.i.d. r.v.,  $(\theta\omega)_i = \omega_{i+1}$ ,

$$\varphi_\omega^{n+1}(x) = \varphi_\omega^n(x) - \eta \nabla \mathcal{L}_{\omega_{n+1}}(\varphi_\omega^n(x)).$$

### B) Transformers as Stochastic Differential Equations (SDEs)

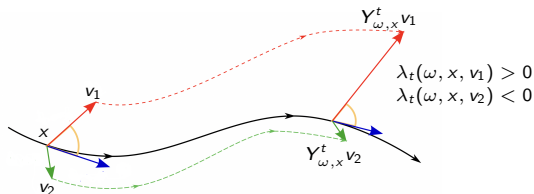
$$dX_t = F_0(X_t) dt + \sum_{j=1}^m F_j(X_t) \partial W_t^j, \quad X_0 = x \in \mathbb{S}^{n-1}. \quad (\text{SDE})$$

- ▶  $\Omega = C_0(\mathbb{R}, \mathbb{R}^m)$  with Wiener measure  $\mathbb{P}$ ,  $\theta_t \omega(\cdot) := \omega(t + \cdot) - \omega(t)$
- ▶ The solution of (SDE) can be written as cocycle  $\varphi_\omega^t(x)$ .

**Question for A) and B) :** What happens with  $\varphi_\omega^t(x)$  as  $t \rightarrow \infty$ ?

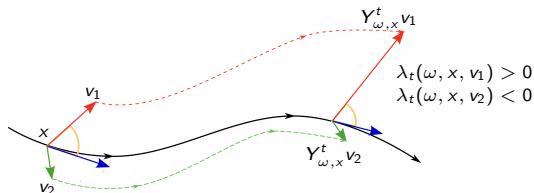
## Lyapunov exponents

- We denote  $Y_{\omega,x}^t := D_x \varphi_{\omega}^t(x)$  and set  $\lambda_t(\omega, x, v) = \frac{1}{t} \log \|Y_{\omega,x}^t v\|$ .



## Lyapunov exponents

- ▶ We denote  $Y_{\omega,x}^t := D_x \varphi_\omega^t(x)$  and set  $\lambda_t(\omega, x, v) = \frac{1}{t} \log \|Y_{\omega,x}^t v\|$ .



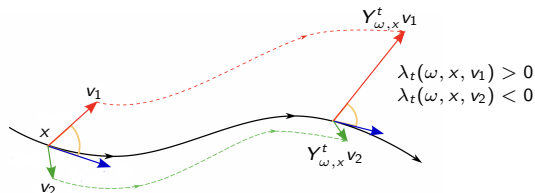
- ▶ For ergodic measure  $\mu$  of  $(\theta, \varphi)$

$$\lambda = \lim_{t \rightarrow \infty} \frac{1}{t} \log \|Y_{\omega,x}^t\|, \quad \text{for } \mu\text{-almost all } (\omega, x),$$

is the maximal **Lyapunov exponent**.

## Lyapunov exponents

- ▶ We denote  $Y_{\omega,x}^t := D_x \varphi_\omega^t(x)$  and set  $\lambda_t(\omega, x, v) = \frac{1}{t} \log \|Y_{\omega,x}^t v\|$ .



- ▶ For ergodic measure  $\mu$  of  $(\theta, \varphi)$

$$\lambda = \lim_{t \rightarrow \infty} \frac{1}{t} \log \|Y_{\omega,x}^t\|, \quad \text{for } \mu\text{-almost all } (\omega, x),$$

is the maximal **Lyapunov exponent**.

### Typical Paradigm:

- ▶  $\lambda < 0 \Rightarrow$  Synchronization/Accumulation
- ▶  $\lambda > 0 \Rightarrow$  Instability/Repulsion/Chaos

## A) Characterizing dynamical stability of stochastic gradient descent in overparameterized learning

[CHEMNITZ/E., JOURNAL OF MACHINE LEARNING RESEARCH 2025]

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

**Overparameterized**  $D > N$ :

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

- ▶ “classical” regime ( - 2010s)

**Overparameterized**  $D > N$ :

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$

**Overparameterized**  $D > N$ :

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$
- ▶ global minimum unique

**Overparameterized**  $D > N$ :

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$
- ▶ global minimum unique
- ▶ Generalization depends on network (mostly)

**Overparameterized**  $D > N$ :

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$
- ▶ global minimum unique
- ▶ Generalization depends on network (mostly)

**Overparameterized**  $D > N$ :

- ▶ “modern” regime (2010s - )

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$
- ▶ global minimum unique
- ▶ Generalization depends on network (mostly)

**Overparameterized**  $D > N$ :

- ▶ “modern” regime (2010s - )
- ▶  $(D - N)$ -dim. manifold of interp. sol.  $\{x : \mathcal{L}(x) = 0\}$

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

### Overdetermined $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$
- ▶ global minimum unique
- ▶ Generalization depends on network (mostly)

### Overparameterized $D > N$ :

- ▶ “modern” regime (2010s - )
- ▶  $(D - N)$ -dim. manifold of interp. sol.  $\{x : \mathcal{L}(x) = 0\}$
- ▶ global min. not unique

## Overdetermined vs. Overparameterized regime

$D$  = number of parameters,  $N$  = size of the training data set

### Overdetermined $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$
- ▶ global minimum unique
- ▶ Generalization depends on network (mostly)

### Overparameterized $D > N$ :

- ▶ “modern” regime (2010s - )
- ▶  $(D - N)$ -dim. manifold of interp. sol.  $\{x : \mathcal{L}(x) = 0\}$
- ▶ global min. not unique
- ▶ Generalization depends on network and **learning**

## Overdetermined vs. Overparameterized regime

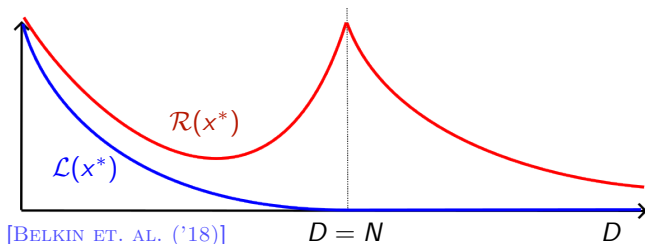
$D$  = number of parameters,  $N$  = size of the training data set

**Overdetermined**  $D < N$ :

- ▶ “classical” regime ( - 2010s)
- ▶ no interpolation solutions  
 $\{x : \mathcal{L}(x) = 0\} = \emptyset$
- ▶ global minimum unique
- ▶ Generalization depends on network (mostly)

**Overparameterized**  $D > N$ :

- ▶ “modern” regime (2010s - )
- ▶  $(D - N)$ -dim. manifold of interp. sol.  $\{x : \mathcal{L}(x) = 0\}$
- ▶ global min. not unique
- ▶ Generalization depends on network and **learning**



## Learning in the overparameterized regime

Loss:  $\mathcal{L}_i(x) := (\Phi(x, y_i) - z_i)^2$ ,  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$

Simplest algorithm: *Gradient descent (GD)*

## Learning in the overparameterized regime

Loss:  $\mathcal{L}_i(x) := (\Phi(x, y_i) - z_i)^2$ ,  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$

Simplest algorithm: *Gradient descent (GD)*

1. **Initialization:**  $X_0^{\text{GD}} \in \mathbb{R}^D$  random with dist  $\nu_{\text{init}}$

## Learning in the overparameterized regime

Loss:  $\mathcal{L}_i(x) := (\Phi(x, y_i) - z_i)^2$ ,  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$

Simplest algorithm: *Gradient descent (GD)*

1. **Initialization:**  $X_0^{\text{GD}} \in \mathbb{R}^D$  random with dist  $\nu_{\text{init}}$
2. **Update:**  $\eta > 0$  learning rate

$$X_{n+1}^{\text{GD}} = X_n^{\text{GD}} - \eta \nabla \mathcal{L}(X_n^{\text{GD}})$$

## Learning in the overparameterized regime

Loss:  $\mathcal{L}_i(x) := (\Phi(x, y_i) - z_i)^2$ ,  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$

Simplest algorithm: *Gradient descent (GD)*

1. **Initialization:**  $X_0^{\text{GD}} \in \mathbb{R}^D$  random with dist  $\nu_{\text{init}}$
2. **Update:**  $\eta > 0$  learning rate

$$X_{n+1}^{\text{GD}} = X_n^{\text{GD}} - \eta \nabla \mathcal{L}(X_n^{\text{GD}})$$

3. **Output:**  $X_{\text{lim}}^{\text{GD}} := \begin{cases} \lim_{n \rightarrow \infty} X_n^{\text{GD}}, & \text{if convergent} \\ \emptyset, & \text{else.} \end{cases}$

## Learning in the overparameterized regime

Loss:  $\mathcal{L}_i(x) := (\Phi(x, y_i) - z_i)^2$ ,  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$

Simplest algorithm: *Gradient descent (GD)*

Most common: Variants of *stochastic gradient descent (SGD)*

1. **Initialization:**  $X_0^{\text{GD}} \in \mathbb{R}^D$  random with dist  $\nu_{\text{init}}$
2. **Update:**  $\eta > 0$  learning rate

$$X_{n+1}^{\text{GD}} = X_n^{\text{GD}} - \eta \nabla \mathcal{L}(X_n^{\text{GD}})$$

3. **Output:**  $X_{\text{lim}}^{\text{GD}} := \begin{cases} \lim_{n \rightarrow \infty} X_n^{\text{GD}}, & \text{if convergent} \\ \emptyset, & \text{else.} \end{cases}$

## Learning in the overparameterized regime

Loss:  $\mathcal{L}_i(x) := (\Phi(x, y_i) - z_i)^2$ ,  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$

Simplest algorithm: *Gradient descent (GD)*

Most common: Variants of *stochastic gradient descent (SGD)*

1. **Initialization:**  $X_0^{\text{GD/SGD}} \in \mathbb{R}^D$  random with dist  $\nu_{\text{init}}$
2. **Update:**  $\eta > 0$  learning rate,  $\omega_n \in [N]$  uniform i.i.d.,

$$X_{n+1}^{\text{GD}} = X_n^{\text{GD}} - \eta \nabla \mathcal{L}(X_n^{\text{GD}})$$

$$X_{n+1}^{\text{SGD}} = X_n^{\text{SGD}} - \eta \nabla \mathcal{L}_{\omega_{n+1}}(X_n^{\text{SGD}})$$

3. **Output:**  $X_{\text{lim}}^{\text{GD/SGD}} := \begin{cases} \lim_{n \rightarrow \infty} X_n^{\text{GD/SGD}}, & \text{if convergent} \\ \emptyset, & \text{else.} \end{cases}$

## Learning in the overparameterized regime

Loss:  $\mathcal{L}_i(x) := (\Phi(x, y_i) - z_i)^2$ ,  $\mathcal{L}(x) := N^{-1} \sum_{i=1}^N \mathcal{L}_i(x)$

Simplest algorithm: *Gradient descent (GD)*

Most common: Variants of *stochastic gradient descent (SGD)*

1. **Initialization:**  $X_0^{\text{GD/SGD}} \in \mathbb{R}^D$  random with dist  $\nu_{\text{init}}$
2. **Update:**  $\eta > 0$  learning rate,  $\omega_n \in [N]$  uniform i.i.d.,

$$X_{n+1}^{\text{GD}} = X_n^{\text{GD}} - \eta \nabla \mathcal{L}(X_n^{\text{GD}})$$

$$X_{n+1}^{\text{SGD}} = X_n^{\text{SGD}} - \eta \nabla \mathcal{L}_{\omega_{n+1}}(X_n^{\text{SGD}})$$

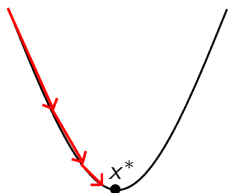
3. **Output:**  $X_{\text{lim}}^{\text{GD/SGD}} := \begin{cases} \lim_{n \rightarrow \infty} X_n^{\text{GD/SGD}}, & \text{if convergent} \\ \emptyset, & \text{else.} \end{cases}$

### Question:

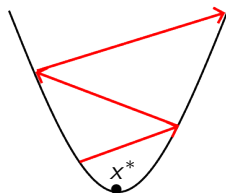
Do we converge and, if yes, whereto?

## Dynamical stability in gradient descent

$$x \mapsto x - \eta \nabla \mathcal{L}(x)$$

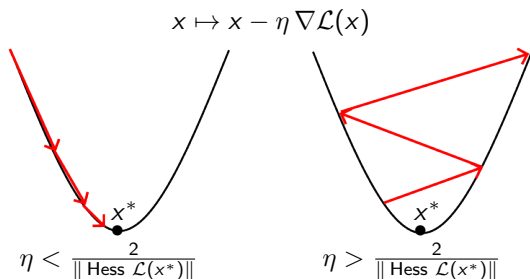


$$\eta < \frac{2}{\|\text{Hess } \mathcal{L}(x^*)\|}$$



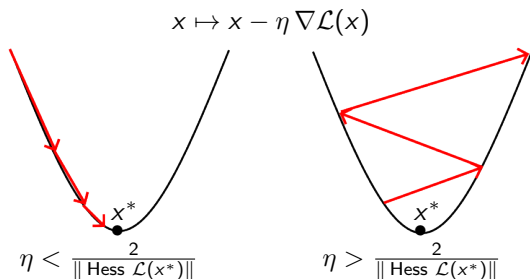
$$\eta > \frac{2}{\|\text{Hess } \mathcal{L}(x^*)\|}$$

## Dynamical stability in gradient descent



Overdetermined regime: Choose  $\eta \ll 1$ .

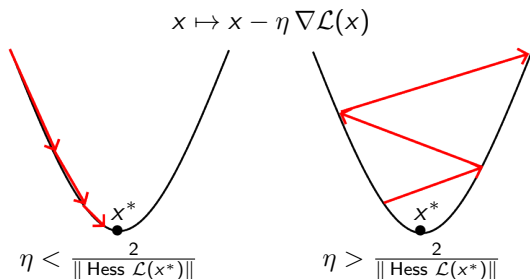
## Dynamical stability in gradient descent



Overdetermined regime: Choose  $\eta \ll 1$ .

Overparameterized regime:  $\|\text{Hess } \mathcal{L}(x^*)\| < 2\eta^{-1}$

## Dynamical stability in gradient descent

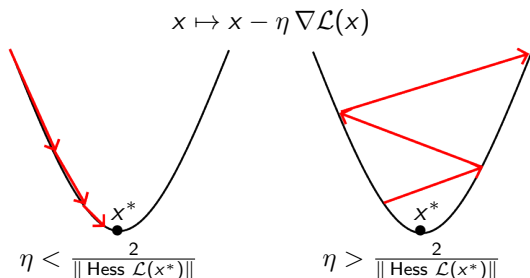


Overdetermined regime: Choose  $\eta \ll 1$ .

Overparameterized regime:  $\|\text{Hess } \mathcal{L}(x^*)\| < 2\eta^{-1}$

- ▶ Instability restricts the **effective hypothesis class**

## Dynamical stability in gradient descent

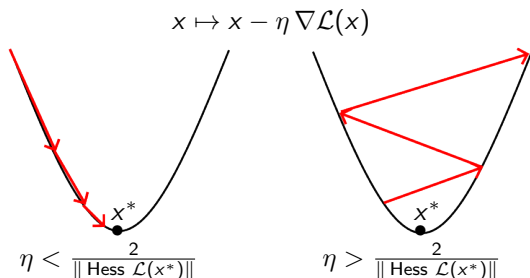


Overdetermined regime: Choose  $\eta \ll 1$ .

Overparameterized regime:  $\|\text{Hess } \mathcal{L}(x^*)\| < 2\eta^{-1}$

- ▶ Instability restricts the **effective hypothesis class**
- ▶  $\|\text{Hess } \mathcal{L}(x^*)\|$  small  $\rightsquigarrow$  Good generalization  
[HOCHREITER, SCHMIDTHUBER ('97)]

## Dynamical stability in gradient descent

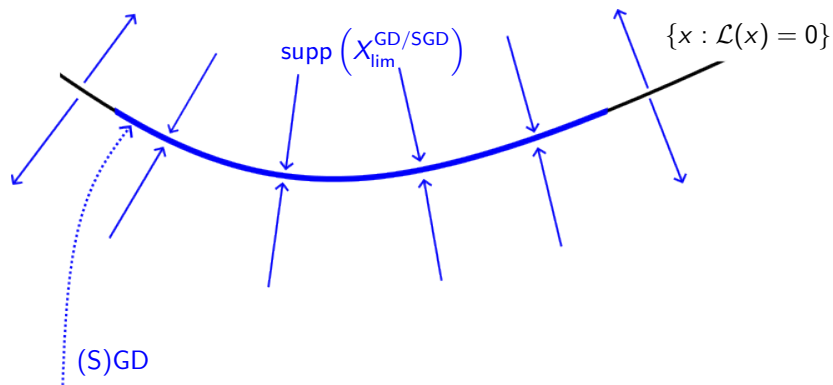


Overdetermined regime: Choose  $\eta \ll 1$ .

Overparameterized regime:  $\|\text{Hess } \mathcal{L}(x^*)\| < 2\eta^{-1}$

- ▶ Instability restricts the **effective hypothesis class**
- ▶  $\|\text{Hess } \mathcal{L}(x^*)\|$  small  $\rightsquigarrow$  Good generalization  
[HOCHREITER, SCHMIDTHUBER ('97)]
- ▶ “Edge of stability”:  $\|\text{Hess } \mathcal{L}(x^*)\| \approx 2\eta^{-1}$   
[WU ET AL. ('18)], [COHEN ET AL. ('20)]

## Dynamical stability in gradient descent



Goal: Characterize

$$\text{supp}(X_{\text{lim}}^{\text{GD/SGD}}) := \text{supp}(\text{Law}(X_{\text{lim}}^{\text{GD/SGD}})).$$

## Main Results - GD

Global min.:  $\mathcal{M} = \{x : \mathcal{L}(x) = 0\}$ , Update:  $\varphi(x) = x - \eta \nabla \mathcal{L}(x)$

## Main Results - GD

Global min.:  $\mathcal{M} = \{x : \mathcal{L}(x) = 0\}$ , Update:  $\varphi(x) = x - \eta \nabla \mathcal{L}(x)$

Assume

- ▶  $\mathcal{M}$  manifold [COOPER ('21)],
- ▶  $\nu_{\text{init}}$  equivalent to Lebesgue (usually Gaussian),
- ▶  $\varphi$  non-singular [CRAČIUN, GHOSHDASTIDAR ('24)].

## Main Results - GD

Global min.:  $\mathcal{M} = \{x : \mathcal{L}(x) = 0\}$ , Update:  $\varphi(x) = x - \eta \nabla \mathcal{L}(x)$

Assume

- ▶  $\mathcal{M}$  manifold [COOPER ('21)],
- ▶  $\nu_{\text{init}}$  equivalent to Lebesgue (usually Gaussian),
- ▶  $\varphi$  non-singular [CRAČIUN, GHOSHDASTIDAR ('24)].

For  $x^* \in \mathcal{M}$ , define

$$\mu(x^*) := \log \|D\varphi(x^*)|_{(T_{x^*}\mathcal{M})^\perp}\|,$$

such that  $\mu(x^*) < 0 \Leftrightarrow \eta < 2\|\text{Hess } \mathcal{L}(x^*)\|^{-1}$ .

## Main Results - GD

Global min.:  $\mathcal{M} = \{x : \mathcal{L}(x) = 0\}$ , Update:  $\varphi(x) = x - \eta \nabla \mathcal{L}(x)$

Assume

- ▶  $\mathcal{M}$  manifold [COOPER ('21)],
- ▶  $\nu_{\text{init}}$  equivalent to Lebesgue (usually Gaussian),
- ▶  $\varphi$  non-singular [CRAČIUN, GHOSHDASTIDAR ('24)].

For  $x^* \in \mathcal{M}$ , define

$$\mu(x^*) := \log \|D\varphi(x^*)|_{(T_{x^*}\mathcal{M})^\perp}\|,$$

such that  $\mu(x^*) < 0 \Leftrightarrow \eta < 2\|\text{Hess } \mathcal{L}(x^*)\|^{-1}$ .

### Theorem (Chemnitz/E. 2025)

Let  $x^* \in \mathcal{M}$ . Then

$$\begin{aligned}\mu(x^*) < 0 &\Rightarrow x^* \in \text{supp}(X_{\text{lim}}^{\text{GD}}), \\ \mu(x^*) > 0 &\Rightarrow x^* \notin \text{supp}(X_{\text{lim}}^{\text{GD}}).\end{aligned}$$

## Linearization - SGD

Consider the **random dynamical system** given by

$$\hat{\varphi}_i : \mathbb{R}^D \rightarrow \mathbb{R}^D, x \mapsto x - \eta \nabla \mathcal{L}_i(x).$$

Let  $\varphi_\omega^n = \hat{\varphi}_{\omega_n} \circ \dots \circ \hat{\varphi}_{\omega_1}$ , where  $\omega = (\omega_1, \omega_2, \dots)$  uniform i.i.d. r.v..

## Linearization - SGD

Consider the **random dynamical system** given by

$$\hat{\varphi}_i : \mathbb{R}^D \rightarrow \mathbb{R}^D, x \mapsto x - \eta \nabla \mathcal{L}_i(x).$$

Let  $\varphi_\omega^n = \hat{\varphi}_{\omega_n} \circ \dots \circ \hat{\varphi}_{\omega_1}$ , where  $\omega = (\omega_1, \omega_2, \dots)$  uniform i.i.d. r.v..

Note:  $x^* \in \mathcal{M} \Rightarrow \hat{\varphi}_i(x^*) = x^*$  for all  $1 \leq i \leq N$ .

## Linearization - SGD

Consider the **random dynamical system** given by

$$\hat{\varphi}_i : \mathbb{R}^D \rightarrow \mathbb{R}^D, x \mapsto x - \eta \nabla \mathcal{L}_i(x).$$

Let  $\varphi_\omega^n = \hat{\varphi}_{\omega_n} \circ \dots \circ \hat{\varphi}_{\omega_1}$ , where  $\omega = (\omega_1, \omega_2, \dots)$  uniform i.i.d. r.v..

Note:  $x^* \in \mathcal{M} \Rightarrow \hat{\varphi}_i(x^*) = x^*$  for all  $1 \leq i \leq N$ .

For  $x^* \in \mathcal{M}$ , let  $A_i(x^*) : (T_{x^*} \mathcal{M})^\perp \rightarrow (T_{x^*} \mathcal{M})^\perp$  be given by

$$A_i(x^*) := D\hat{\varphi}_i(x^*)|_{(T_{x^*} \mathcal{M})^\perp}.$$

## Linearization - SGD

Consider the **random dynamical system** given by

$$\hat{\varphi}_i : \mathbb{R}^D \rightarrow \mathbb{R}^D, \quad x \mapsto x - \eta \nabla \mathcal{L}_i(x).$$

Let  $\varphi_\omega^n = \hat{\varphi}_{\omega_n} \circ \dots \circ \hat{\varphi}_{\omega_1}$ , where  $\omega = (\omega_1, \omega_2, \dots)$  uniform i.i.d. r.v..

Note:  $x^* \in \mathcal{M} \Rightarrow \hat{\varphi}_i(x^*) = x^*$  for all  $1 \leq i \leq N$ .

For  $x^* \in \mathcal{M}$ , let  $A_i(x^*) : (T_{x^*} \mathcal{M})^\perp \rightarrow (T_{x^*} \mathcal{M})^\perp$  be given by

$$A_i(x^*) := D\hat{\varphi}_i(x^*)|_{(T_{x^*} \mathcal{M})^\perp}.$$

### Theorem (Furstenberg-Kesten)

*The Lyapunov exponent*

$$\lambda(x^*) := \lim_{n \rightarrow \infty} \frac{1}{n} \log \|A_{\omega_n}(x^*) \dots A_{\omega_1}(x^*)\|$$

*exists for almost all  $\omega$  and is essentially constant.*

# Main Results - SGD

## Assume

- ▶  $\mathcal{M}$  manifold [COOPER ('21)],
- ▶  $\nu_{\text{init}}$  equivalent to Lebesgue,
- ▶  $\hat{\varphi}_j$  non-singular [CRAČIUN, GHOSHDASTIDAR ('24)].

# Main Results - SGD

Assume

- ▶  $\mathcal{M}$  manifold [COOPER ('21)],
- ▶  $\nu_{\text{init}}$  equivalent to Lebesgue,
- ▶  $\hat{\varphi}_i$  non-singular [CRAČIUN, GHOSHDASTIDAR ('24)].

Theorem (Chemnitz/E. 2025)

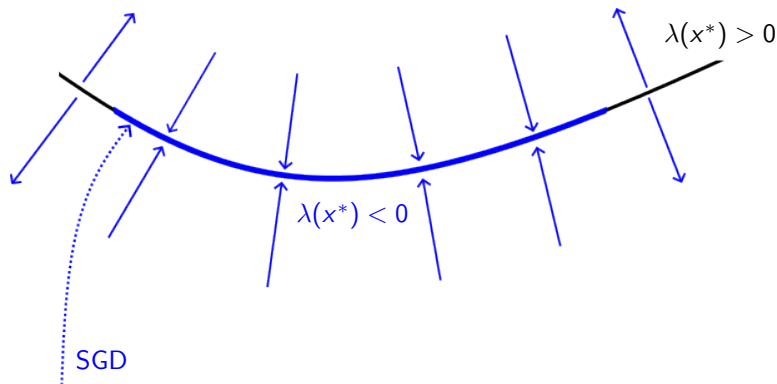
Let  $x^* \in \mathcal{M}$  be a “regular minimum”. Then

$$\lambda(x^*) < 0 \Rightarrow x^* \in \text{supp}(X_{\text{lim}}^{\text{SGD}}),$$

$$\lambda(x^*) > 0 \Rightarrow x^* \notin \text{supp}(X_{\text{lim}}^{\text{SGD}}).$$

## Proof strategies

- ▶ For  $\lambda(x^*) < 0$ , **Quenched**: “Pathwise dynamics” for typical  $\omega$
- ▶ For  $\lambda(x^*) > 0$ , **Annealed**: construct Lyapunov function for Markov process away from  $x^*$



## B) Random Quadratic Form on a Sphere: Synchronization by Common Noise

[E./SHALOVA, ARXIV:2603.06187, 2026+]

## Transformer models [VASWANI ET AL. 2017]

- ▶ Sentence of length  $d$  represented via *tokens*  $(x_i)_{1 \leq i \leq d}$ ,  $x_i \in \mathbb{R}^n$ .

## Transformer models [VASWANI ET AL. 2017]

- ▶ Sentence of length  $d$  represented via *tokens*  $(x_i)_{1 \leq i \leq d}$ ,  $x_i \in \mathbb{R}^n$ .
- ▶ Neural map at  $k$ -th transformer layer

$$x_i^{k+1} = \hat{\Phi}_k((x_i)_{1 \leq i \leq d}),$$

with **self-attention**, **feed-forward** and projection to  $\mathbb{S}^{n-1}$ .

## Transformer models [VASWANI ET AL. 2017]

- ▶ Sentence of length  $d$  represented via *tokens*  $(x_i)_{1 \leq i \leq d}$ ,  $x_i \in \mathbb{R}^n$ .
- ▶ Neural map at  $k$ -th transformer layer

$$x_i^{k+1} = \hat{\Phi}_k((x_i)_{1 \leq i \leq d}),$$

with **self-attention**, **feed-forward** and projection to  $\mathbb{S}^{n-1}$ .

### Continuous-time model [GESHKOVSKI ET AL. 2024, 2025]

$$\begin{aligned} \dot{x}_i &= P_{x_i}(\text{FF}(x_i) + \text{Attn}(x_i; x_1, x_2 \dots x_d)), \\ \text{FF}(x_i) &= \sigma(Mx_i + B), \quad P_{x_i} : \mathbb{R}^n \rightarrow \mathbb{S}^{n-1}. \end{aligned}$$

## Transformer models [VASWANI ET AL. 2017]

- ▶ Sentence of length  $d$  represented via *tokens*  $(x_i)_{1 \leq i \leq d}$ ,  $x_i \in \mathbb{R}^n$ .
- ▶ Neural map at  $k$ -th transformer layer

$$x_i^{k+1} = \hat{\Phi}_k((x_i)_{1 \leq i \leq d}),$$

with **self-attention**, **feed-forward** and projection to  $\mathbb{S}^{n-1}$ .

### Continuous-time model [GESHKOVSKI ET AL. 2024, 2025]

$$\begin{aligned}\dot{x}_i &= P_{x_i}(\text{FF}(x_i) + \text{Attn}(x_i; x_1, x_2 \dots x_d)), \\ \text{FF}(x_i) &= \sigma(Mx_i + B), \quad P_{x_i} : \mathbb{R}^n \rightarrow \mathbb{S}^{n-1}.\end{aligned}$$

**Simplified model** ( $\text{Attn} = 0 = B$ ) for random parametrization at each  $k$ :

## Transformer models [VASWANI ET AL. 2017]

- ▶ Sentence of length  $d$  represented via *tokens*  $(x_i)_{1 \leq i \leq d}$ ,  $x_i \in \mathbb{R}^n$ .
- ▶ Neural map at  $k$ -th transformer layer

$$x_i^{k+1} = \hat{\Phi}_k((x_i)_{1 \leq i \leq d}),$$

with **self-attention**, **feed-forward** and projection to  $\mathbb{S}^{n-1}$ .

### Continuous-time model [GESHKOVSKI ET AL. 2024, 2025]

$$\begin{aligned}\dot{x}_i &= P_{x_i}(\text{FF}(x_i) + \text{Attn}(x_i; x_1, x_2 \dots x_d)), \\ \text{FF}(x_i) &= \sigma(Mx_i + B), \quad P_{x_i} : \mathbb{R}^n \rightarrow \mathbb{S}^{n-1}.\end{aligned}$$

**Simplified model** ( $\text{Attn} = 0 = B$ ) for random parametrization at each  $k$ :

$$\begin{aligned}dX_i(t) &= P_{X_i(t)} \partial Q(t) X_i(t), & \textbf{(Random Quadratic Form (RQF))} \\ Q_t &= \frac{1}{2}(B_t + B_t^T),\end{aligned}$$

where  $\{B_t^{ij} : i, j \in 1 \dots n\}$  are independent Brownian motions.

Synchronization of tokens purely by common noise?

## Synchronization of tokens purely by common noise?

Consider the RDS solutions  $\varphi_\omega^t(x)$  of

$$dX(t) = P_{X(t)} \partial Q(t) X(t), \quad X(0) = x \in \mathbb{S}^{n-1}. \quad (\mathbf{RQF})$$

## Synchronization of tokens purely by common noise?

Consider the RDS solutions  $\varphi_\omega^t(x)$  of

$$dX(t) = P_{X(t)} \partial Q(t) X(t), \quad X(0) = x \in \mathbb{S}^{n-1}. \quad (\mathbf{RQF})$$

### Theorem (E./Shalova 2026+)

For any  $x, y \in \mathbb{S}^{n-1}$

$$\text{dist}(\varphi_\omega^t(y), \varphi_\omega^t(x)) \rightarrow 0 \quad \text{or} \quad \text{dist}(\varphi_\omega^t(y), -\varphi_\omega^t(x)) \rightarrow 0,$$

$\Omega$ -almost surely.

## Synchronization of tokens purely by common noise?

Consider the RDS solutions  $\varphi_\omega^t(x)$  of

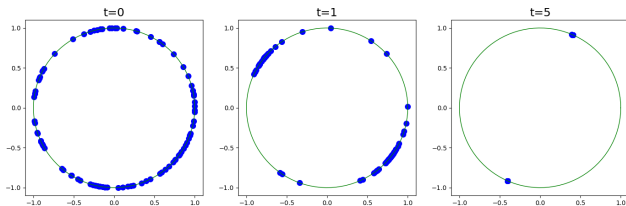
$$dX(t) = P_{X(t)} \partial Q(t) X(t), \quad X(0) = x \in \mathbb{S}^{n-1}. \quad (\text{RQF})$$

### Theorem (E./Shalova 2026+)

For any  $x, y \in \mathbb{S}^{n-1}$

$$\text{dist}(\varphi_\omega^t(y), \varphi_\omega^t(x)) \rightarrow 0 \quad \text{or} \quad \text{dist}(\varphi_\omega^t(y), -\varphi_\omega^t(x)) \rightarrow 0,$$

$\Omega$ -almost surely.



## Proof strategies

### First Proof:

- ▶ For solutions  $X_t, Y_t$  of (RQF), write  $Z_t = \langle X_t, Y_t \rangle \in [-1, 1]$ .

## Proof strategies

### First Proof:

- ▶ For solutions  $X_t, Y_t$  of (RQF), write  $Z_t = \langle X_t, Y_t \rangle \in [-1, 1]$ .
- ▶ Derive that  $Z_t$  has same law as process

$$d\hat{Z}_t = 2\hat{Z}_t(1 - \hat{Z}_t^2)dt + \sqrt{2}(1 - \hat{Z}_t^2)\partial B_t,$$

## Proof strategies

### First Proof:

- ▶ For solutions  $X_t, Y_t$  of (RQF), write  $Z_t = \langle X_t, Y_t \rangle \in [-1, 1]$ .
- ▶ Derive that  $Z_t$  has same law as process

$$d\hat{Z}_t = 2\hat{Z}_t(1 - \hat{Z}_t^2)dt + \sqrt{2}(1 - \hat{Z}_t^2)\partial B_t,$$

- ▶ Use **Feller criterion** to show that

$$\mathbb{P}(\hat{Z}_t \rightarrow 1) + \mathbb{P}(\hat{Z}_t \rightarrow -1) = 1.$$

## Proof strategies

### First Proof:

- ▶ For solutions  $X_t, Y_t$  of (RQF), write  $Z_t = \langle X_t, Y_t \rangle \in [-1, 1]$ .
- ▶ Derive that  $Z_t$  has same law as process

$$d\hat{Z}_t = 2\hat{Z}_t(1 - \hat{Z}_t^2)dt + \sqrt{2}(1 - \hat{Z}_t^2)\partial B_t,$$

- ▶ Use **Feller criterion** to show that

$$\mathbb{P}(\hat{Z}_t \rightarrow 1) + \mathbb{P}(\hat{Z}_t \rightarrow -1) = 1.$$

### Second Proof:

- ▶ Show that **Lyapunov exponent**  $\lambda < 0$  for RDS induced by (RQF)

## Proof strategies

### First Proof:

- ▶ For solutions  $X_t, Y_t$  of (RQF), write  $Z_t = \langle X_t, Y_t \rangle \in [-1, 1]$ .
- ▶ Derive that  $Z_t$  has same law as process

$$d\hat{Z}_t = 2\hat{Z}_t(1 - \hat{Z}_t^2)dt + \sqrt{2}(1 - \hat{Z}_t^2)\partial B_t,$$

- ▶ Use **Feller criterion** to show that

$$\mathbb{P}(\hat{Z}_t \rightarrow 1) + \mathbb{P}(\hat{Z}_t \rightarrow -1) = 1.$$

### Second Proof:

- ▶ Show that **Lyapunov exponent**  $\lambda < 0$  for RDS induced by (RQF)
- ▶ Write process on projective space  $\mathbb{RP}^n$

## Proof strategies

### First Proof:

- ▶ For solutions  $X_t, Y_t$  of (RQF), write  $Z_t = \langle X_t, Y_t \rangle \in [-1, 1]$ .
- ▶ Derive that  $Z_t$  has same law as process

$$d\hat{Z}_t = 2\hat{Z}_t(1 - \hat{Z}_t^2)dt + \sqrt{2}(1 - \hat{Z}_t^2)\partial B_t,$$

- ▶ Use **Feller criterion** to show that

$$\mathbb{P}(\hat{Z}_t \rightarrow 1) + \mathbb{P}(\hat{Z}_t \rightarrow -1) = 1.$$

### Second Proof:

- ▶ Show that **Lyapunov exponent**  $\lambda < 0$  for RDS induced by (RQF)
- ▶ Write process on projective space  $\mathbb{RP}^n$
- ▶ Use synchronization criteria [[BAXENDALE 1991](#)], [[FLANDOLI ET AL. 2017](#)]