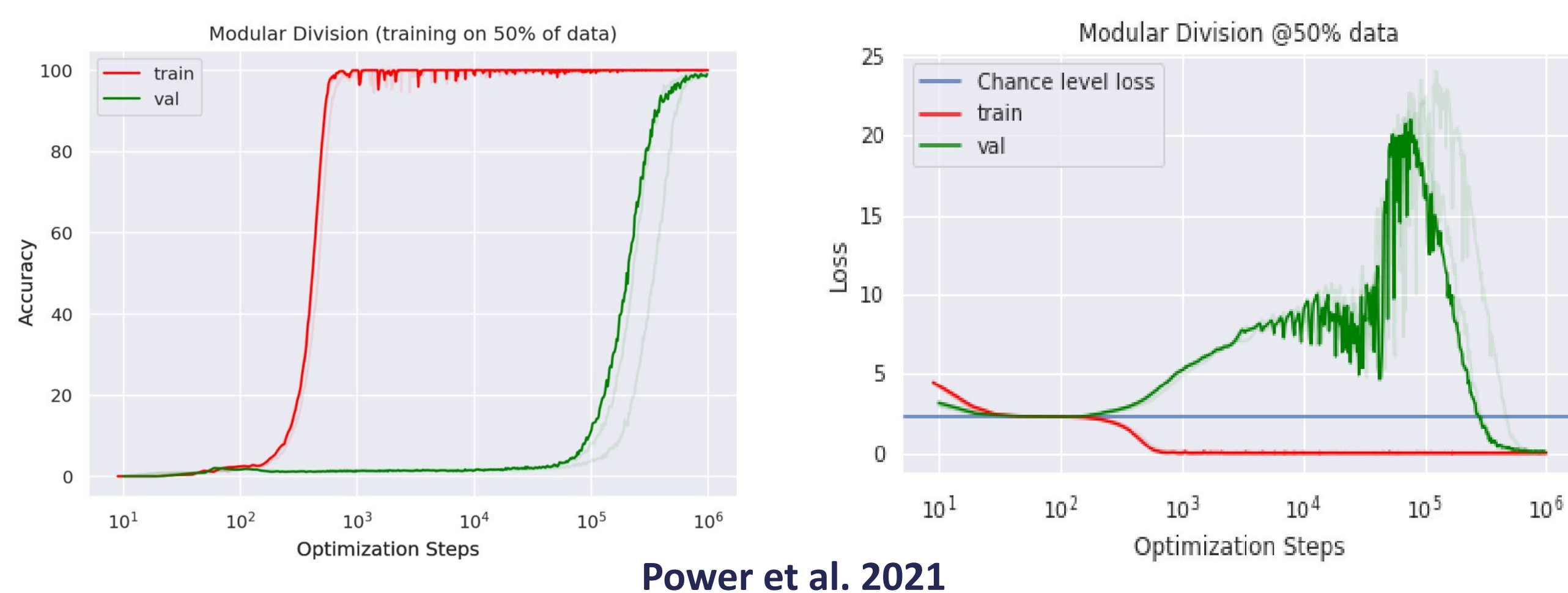


## TL;DR

Grokking may occur even in standard logistic regression. This happens if the data is almost linearly separable. In this case, the model may overfit for an arbitrarily long time, before converging to the ground truth.

## What is Grokking?

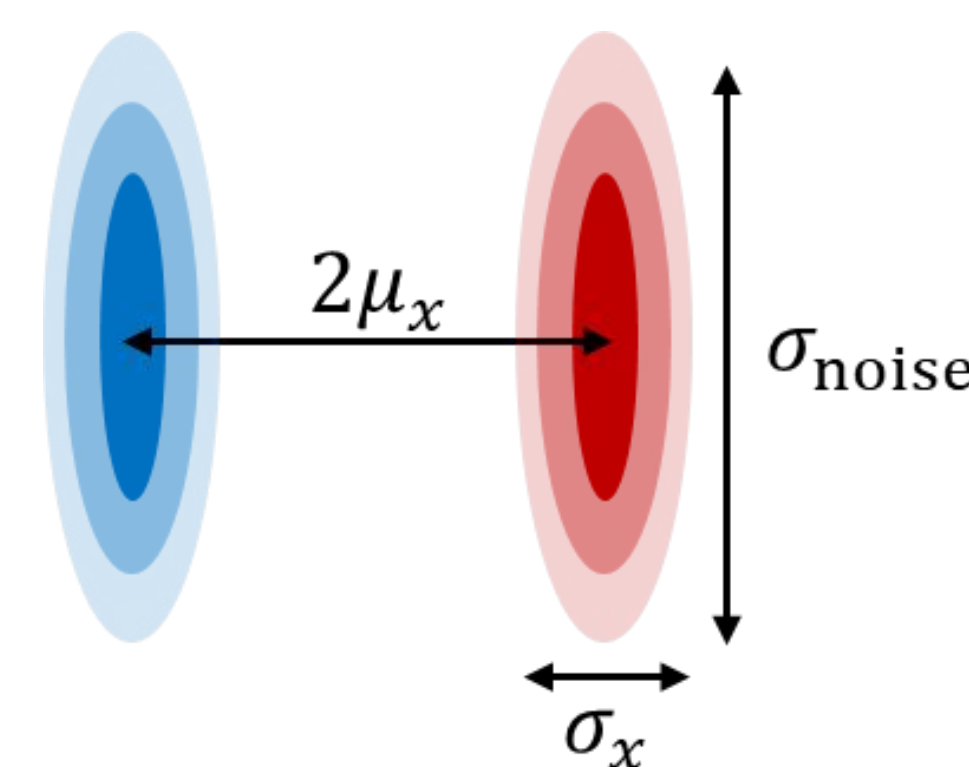


Delayed generalization + non-monotonic test loss. **Counterintuitive!** Overfitting usually implies never generalizing.

**The goal:** Understanding the mechanism in a **minimal** setup

## Binary Classification Setup

- d-dimensional Gaussian data.
- Separation: x-axis.
- Noise in all other directions.
- $\lambda = d/N, N \rightarrow \infty$
- $\sigma_x \ll \mu_x, \mu_x \lesssim \sigma_{\text{noise}}$



Equivalent to a (d-1) dimensional model + bias:

$$\mathbf{x}_i \sim \mathcal{N}(0, \sigma \mathbf{I}_d), \sigma = \sigma_x^2 / \mu_x^2 \gtrsim 1 \quad \tilde{\mathbf{x}}_i = (1, \mathbf{x}_i)$$

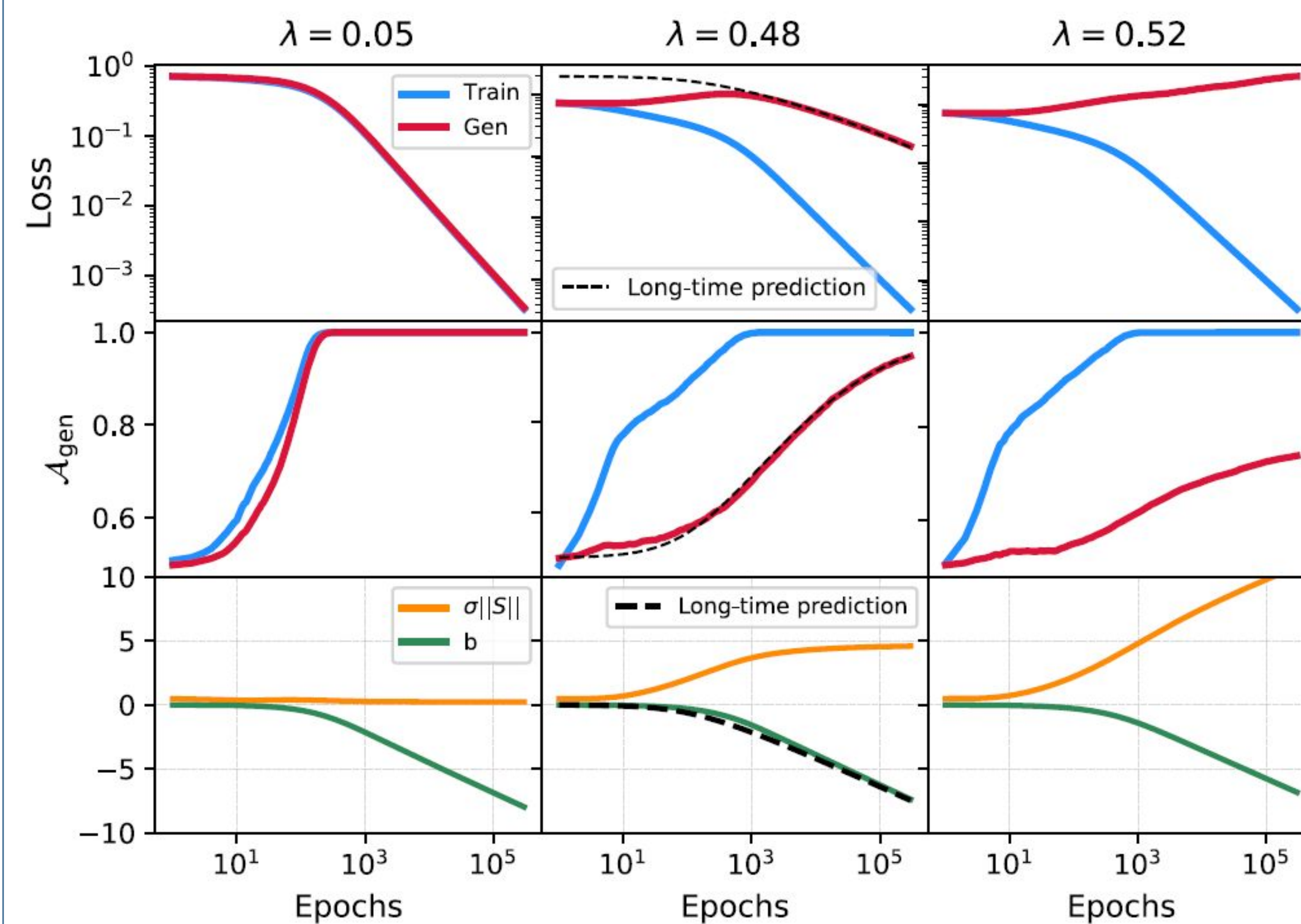
Same label for all points! (for example, -1)

Linear model:  $f(\mathbf{x}_i) = \mathbf{S} \cdot \mathbf{x}_i + b, \mathbf{S} \in \mathbb{R}^{d-1}, b \in \mathbb{R}$

$$\text{CE-loss: } \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{-y_i f(\mathbf{x}_i)}), y_i = -1$$

**Generalization:** only if  $b \rightarrow -\infty$

## Empirical results (GD)



- Cannot generalize for  $\lambda > 0.5$
- **Grokking** for  $\lambda \rightarrow 0.5^-$

## Edge of Linear Separability

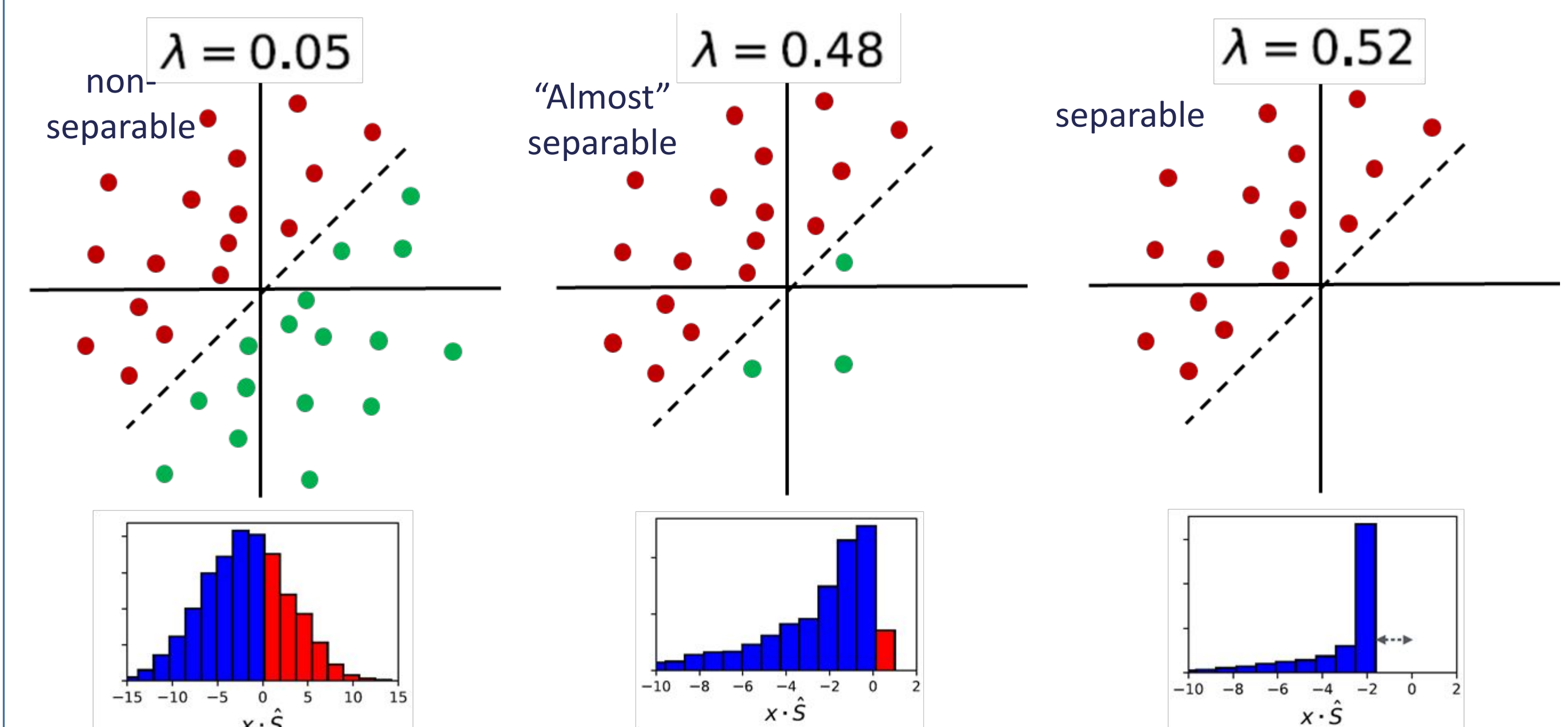
$$\mathcal{A}_{\text{gen}}(\mathbf{S}, b) = \frac{1}{2} \left[ 1 - \text{erf} \left( \frac{1}{\sqrt{2}} \frac{b}{\sigma \|\mathbf{S}\|} \right) \right]$$

Grokking: first  $\|\mathbf{S}\|$  increase (memorization), then saturates while  $b \rightarrow -\infty$  (generalization).

What is special about  $\lambda \approx 0.5^-$ ?

**Theorem 1:** Training data is linearly separable (from the origin) iff  $\lambda > 0.5$

**Proof:** Wendel's theorem



## Mechanism

**Theorem 2:** Generalization iff training data is linearly separable, in particular:

$$\lambda < 0.5 \quad \lambda > 0.5$$

$\|\mathbf{S}_\infty\| = C$ , diverges as  $\lambda \rightarrow 0.5^-$   
Proof: Exponential loss + "conformal time":

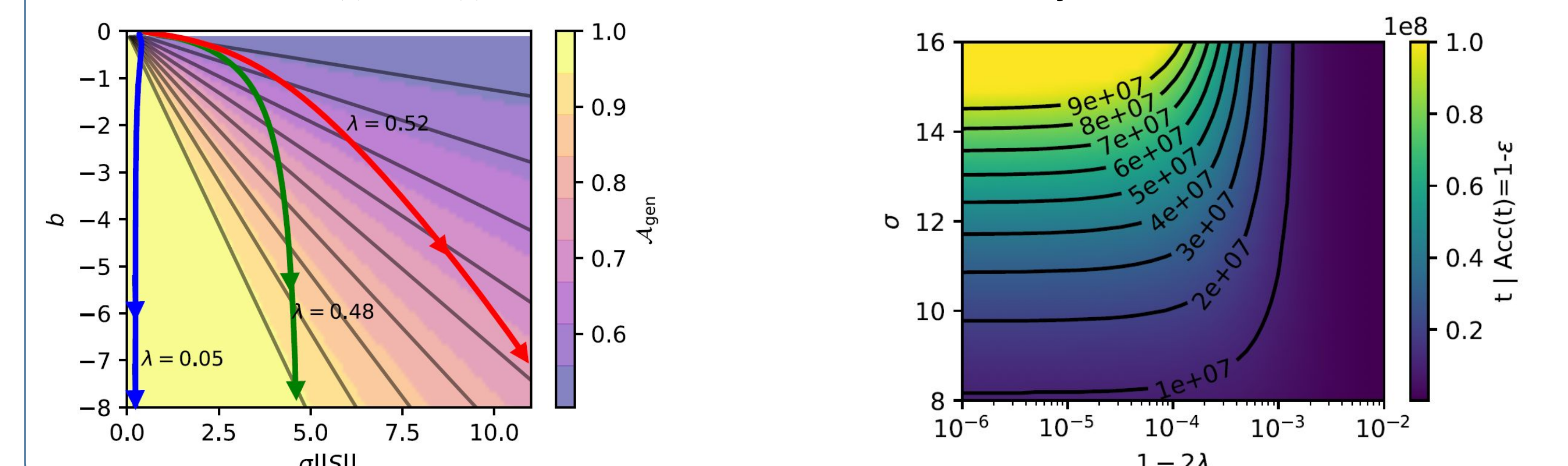
$\|\mathbf{S}_\infty\| \rightarrow \infty$ ,  $\mathcal{A}_{\text{gen}} = \frac{1}{2} [1 + \text{erf}(\frac{1}{\sigma M \sqrt{2}})]$   
Proof: Soudry et al.

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N e^{\mathbf{S}^T \mathbf{x}_i + b} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{S}} = -\frac{\eta}{N} \sum_{i=1}^N e^{\mathbf{S}^T \mathbf{x}_i} \mathbf{x}_i \quad \mathbf{w}(t) = \mathbf{w}_{\text{SVM}} \log(t) + \boldsymbol{\rho}(t)$$

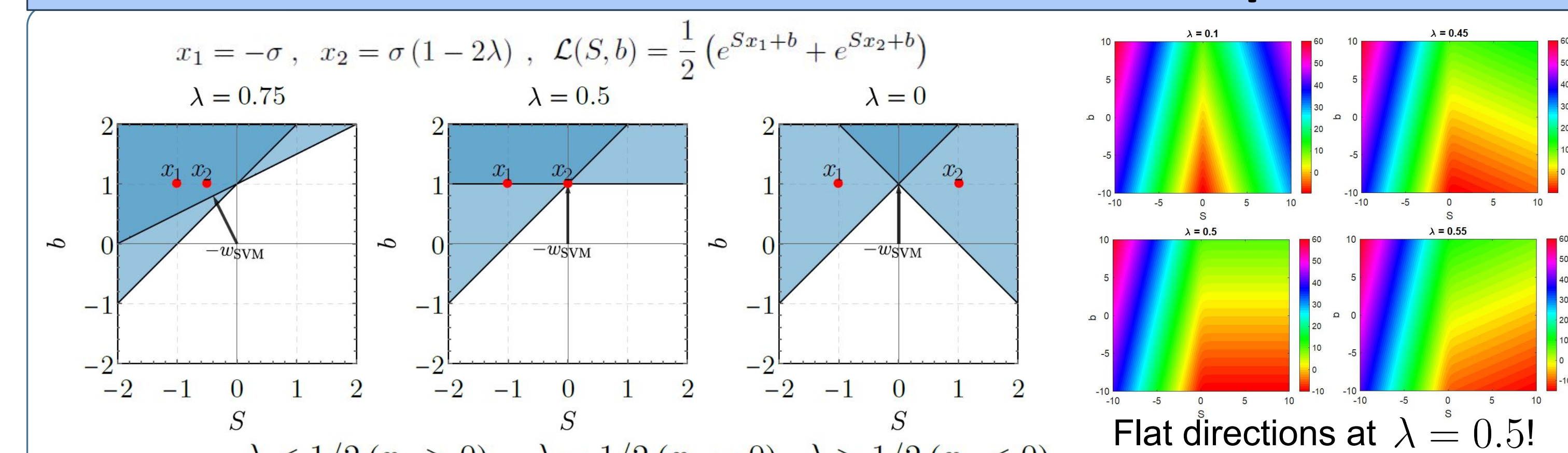
$$\tau(t) = \int_0^t e^{b(t')} dt' \quad \frac{\partial \mathcal{L}}{\partial \tau} = -\frac{\eta}{N} \sum_{i=1}^N e^{\mathbf{S}^T \mathbf{x}_i} \quad \mathbf{w}_{\text{SVM}} = \argmin_{(\mathbf{S}, b)} \{ \|\mathbf{S}\|^2 + b^2 \text{ s.t. } \mathbf{S}^T \mathbf{x}_i + b \leq -1 \}$$

$$\text{Exponential loss without bias: } \mathcal{L} = \frac{1}{N} \sum_{i=1}^N e^{\mathbf{S}^T \mathbf{x}_i} \quad \mathbf{w}_{\text{SVM}} = \left( \frac{M}{1 + M^2} \mathbf{S}_{\text{SVM}}, -\frac{1}{1 + M^2} \right)$$

**Collecting the pieces:** For large enough  $\sigma$ , dynamics first reach  $\|\mathbf{S}_\infty\|$  while  $b \rightarrow -\infty$  only follows later!

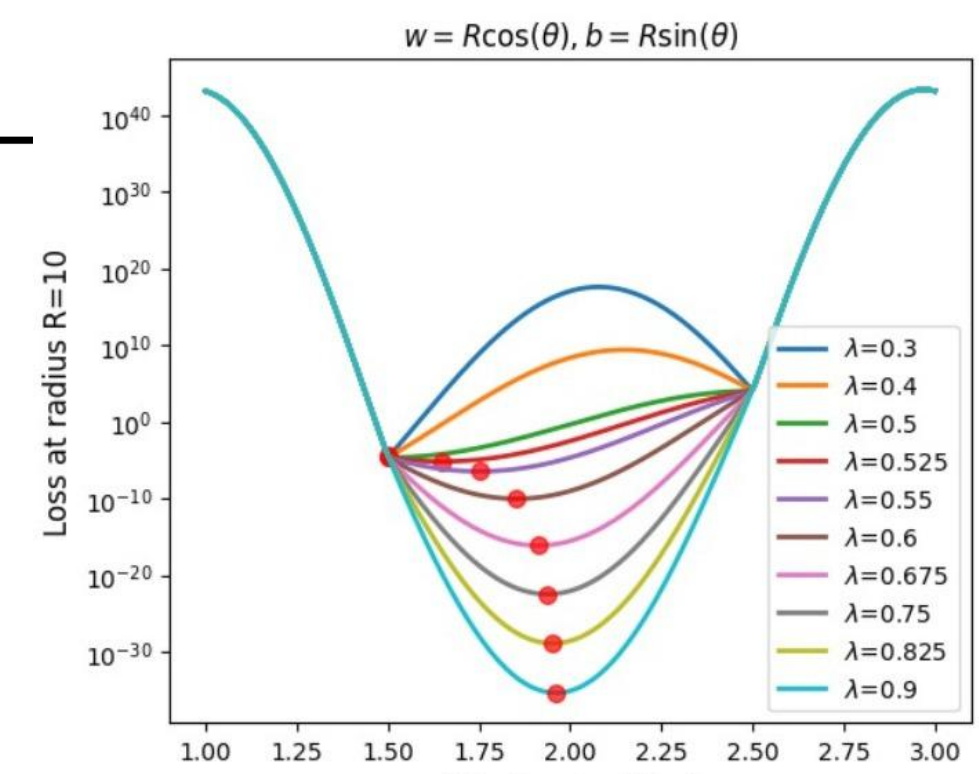


## Effective 1d Solvable Description



$$\begin{aligned} b(t \gg 1) &= -\log(t) & -\log(t) &= -\log(t) & -\frac{1}{1+M^2} \log(t) \\ \|\mathbf{S}\|(t \gg 1) &= \frac{1}{2(1-\lambda)} \log\left(\frac{1}{1-2\lambda}\right) & \log(\log(t)) &= \log(\log(t)) & \frac{M}{1+M^2} \log(t) \end{aligned}$$

Similar to **critical phenomena** in physical systems (phase transitions)  
Future work: Universality classes?



**Check out our paper:**

\*alombk2@gmail.com  
†noam.levi@epfl.ch  
‡ybarsinai@gmail.com

