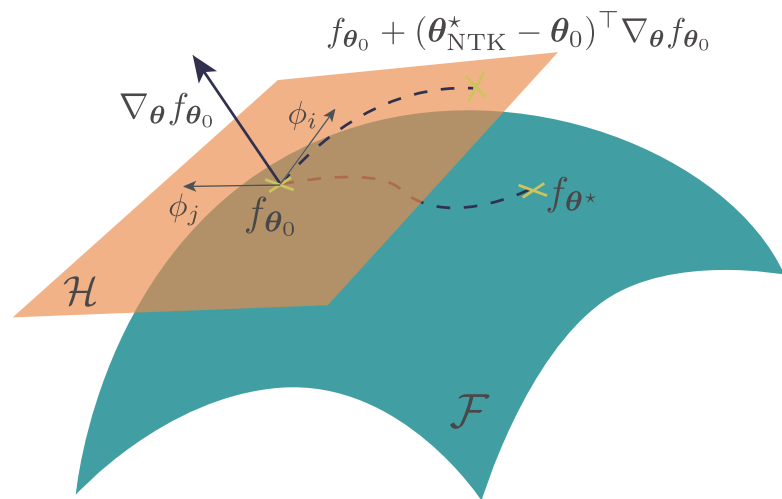


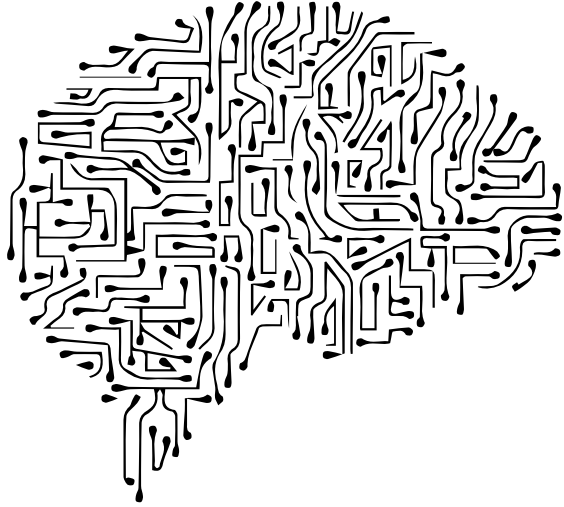
# What can linearized neural networks *actually* say about generalization ?

Guillermo Ortiz-Jiménez

Seyed-Mohsen Moosavi-Dezfooli

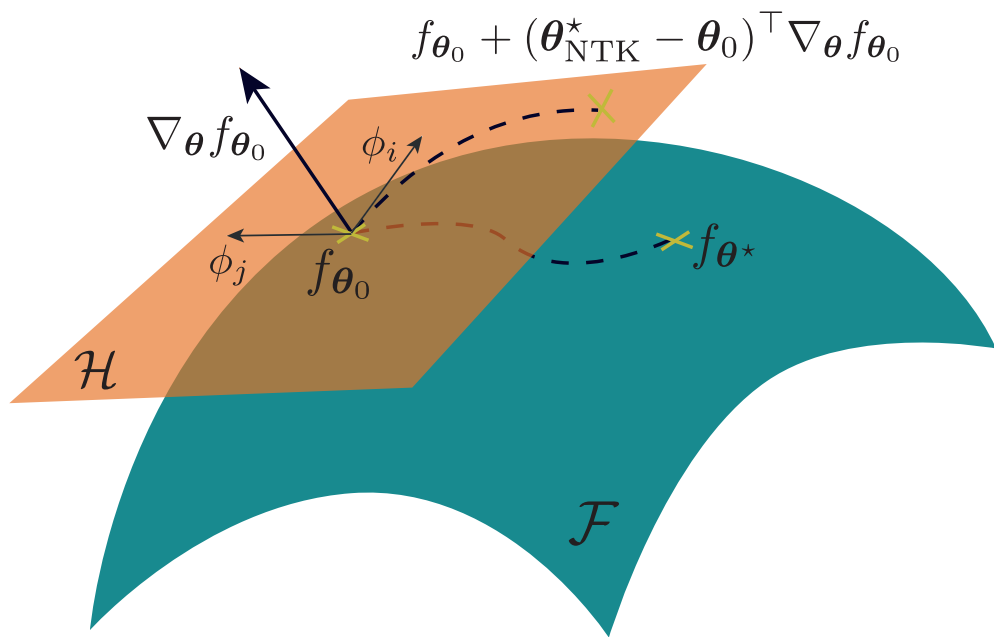
Pascal Frossard



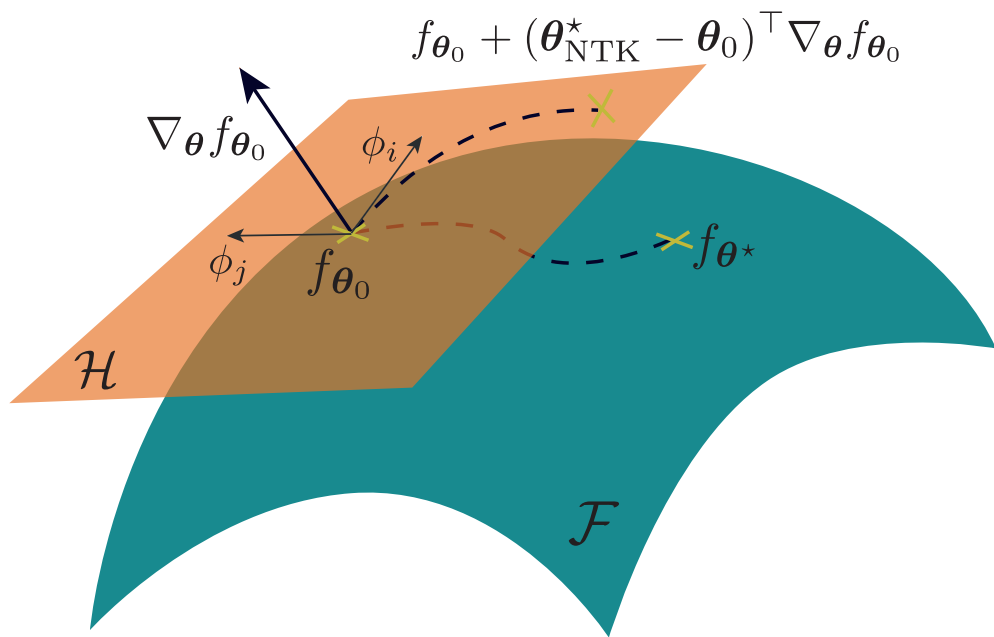


**Deep Learning?**





## Neural Tangent Kernel.



**NTK.**

$$f_{\theta}(\mathbf{x}) \approx f_{\theta_0}(\mathbf{x}) + (\theta - \theta_0)^{\top} \nabla_{\theta} f_{\theta_0}(\mathbf{x})$$

Non-linear

Linear

Neural tangent kernel (NTK):

$$\Theta(\mathbf{x}_1, \mathbf{x}_2) = \langle \nabla_{\theta} f_{\theta_0}(\mathbf{x}_1), \nabla_{\theta} f_{\theta_0}(\mathbf{x}_2) \rangle$$

It defines a metric over functions:

$$\|f\|_{\Theta}^2 = \sum_{j=1}^{\infty} \frac{1}{\lambda_j} (\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\phi_j(\mathbf{x}) f(\mathbf{x})])$$

Eigenfunction

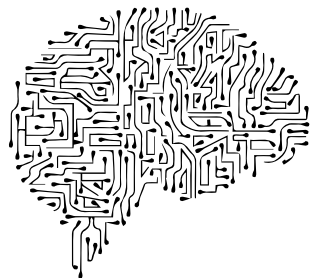
Eigenvalue

Proxy

$$\alpha(f) = \sum_{j=1}^{\infty} \lambda_j (\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\phi_j(\mathbf{x}) f(\mathbf{x})])$$

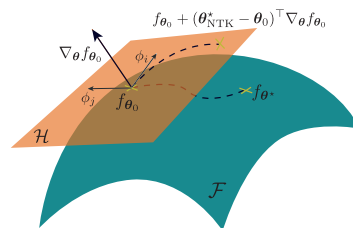


On many problems, provably...



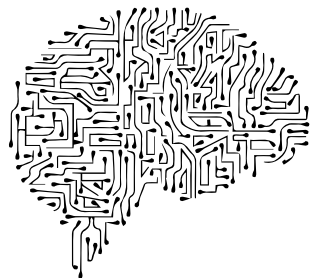
**Deep learning.**

$\neq$

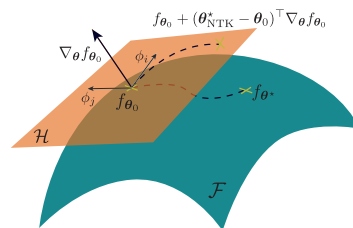


**NTK.**

On many problems, provably...



**Deep learning.**



**NTK.**



# An NTK perspective on [insert topic]



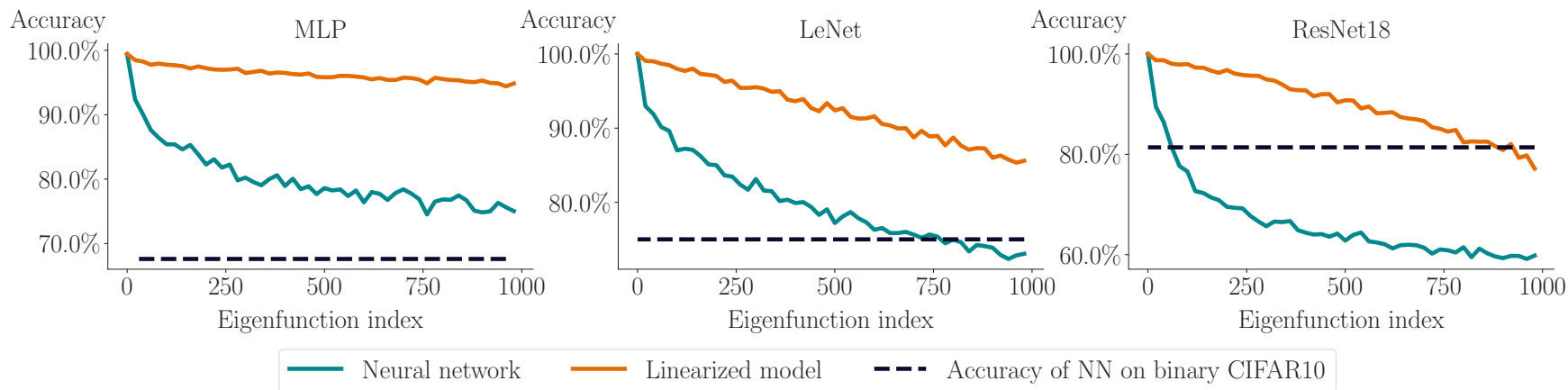
(Arora et al. 2020, Mobahi et al. 2020, Tancik et al 2020, Deshpande et al. 2021, Zancato et al. 2021, Gebhart et al. 2021, Bachman et al. 2021., Maddox et al. 2021, and more.)

**What can linearized neural networks  
*actually* say about generalization ?**

# A new benchmark

Generate different tasks defined **only** by labels:

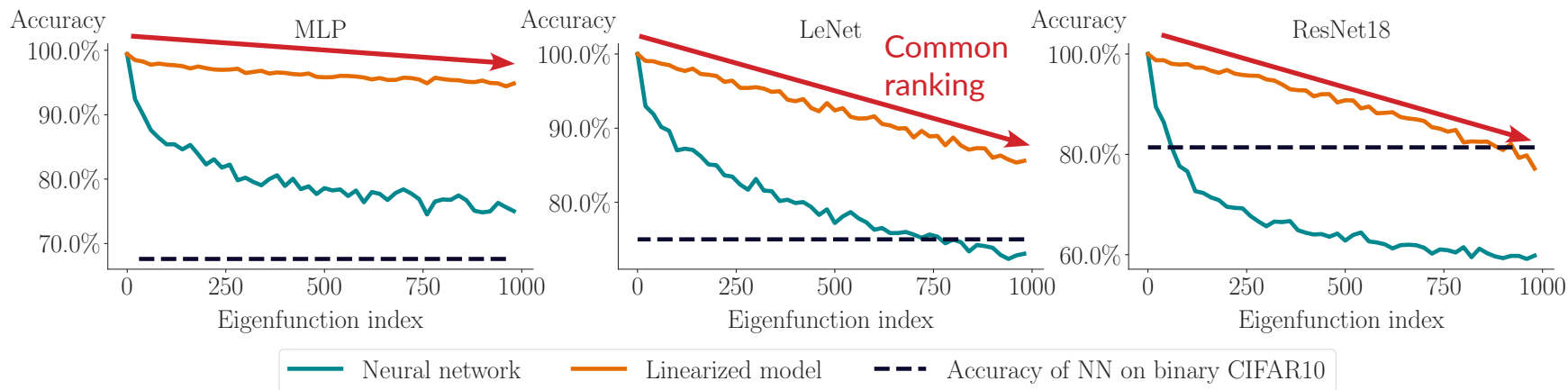
- Take CIFAR10 data
- Assign labels based on NTK:  $\mathbf{x} \mapsto \text{sign}(\phi_j(\mathbf{x}))$
- Train **linear** and **non-linear** models on each task.



# A new benchmark

Generate different tasks defined **only** by labels:

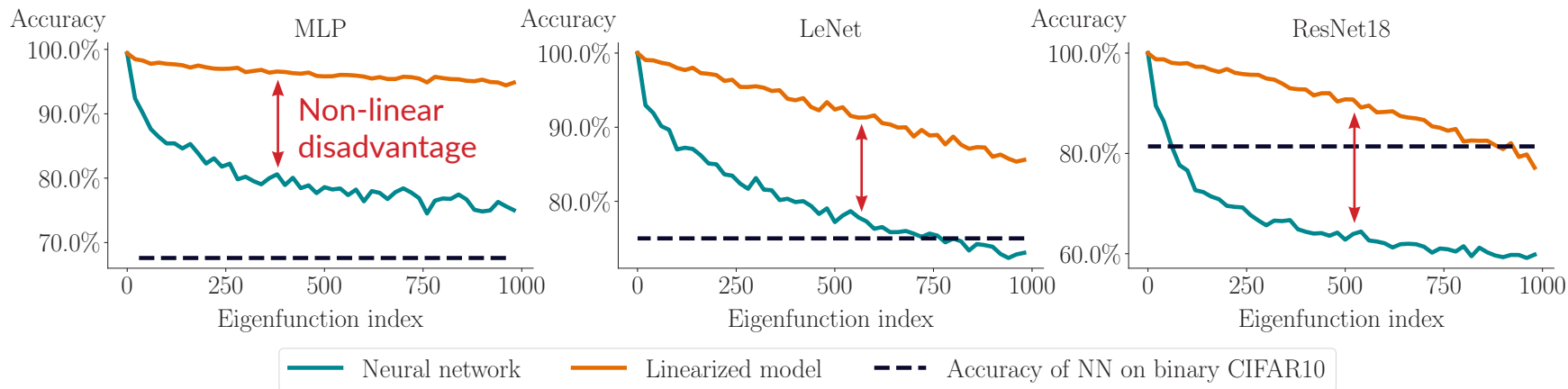
- Take CIFAR10 data
- Assign labels based on NTK:  $\mathbf{x} \mapsto \text{sign}(\phi_j(\mathbf{x}))$
- Train **linear** and **non-linear** models on each task.

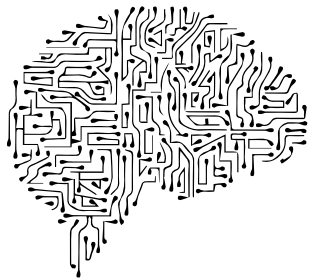


# A new benchmark

Generate different tasks defined **only** by labels:

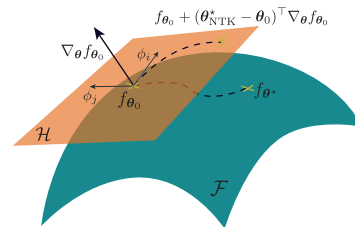
- Take CIFAR10 data
- Assign labels based on NTK:  $\mathbf{x} \mapsto \text{sign}(\phi_j(\mathbf{x}))$
- Train **linear** and **non-linear** models on each task.





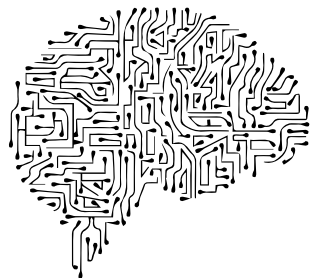
**Deep learning.**

Ordering

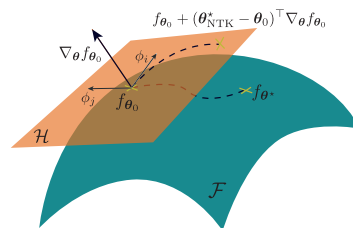


**NTK.**

On some problems...



**Deep learning.**



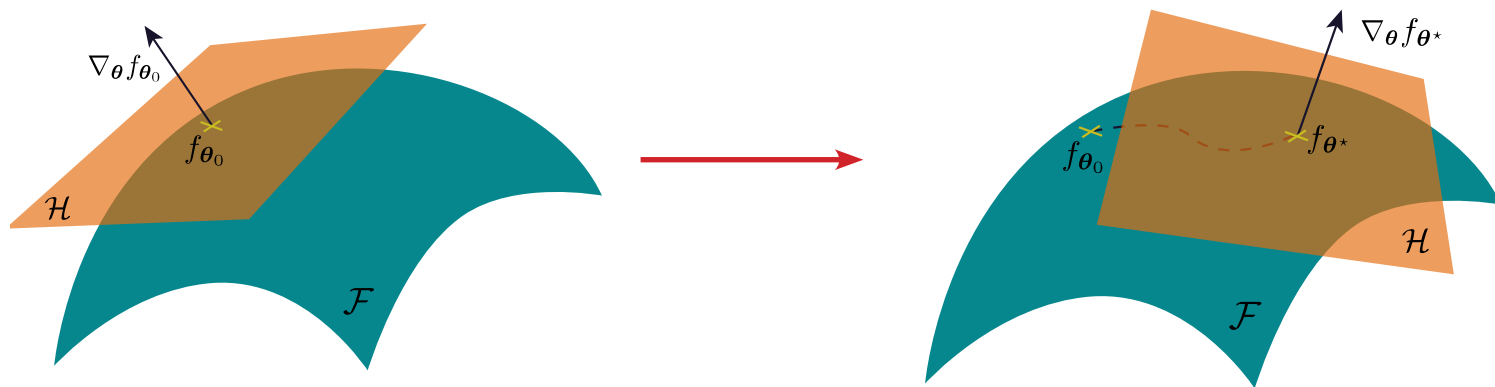
**NTK.**

**What is the source of the  
non-linear (dis)advantage ?**



# Neural network dynamics

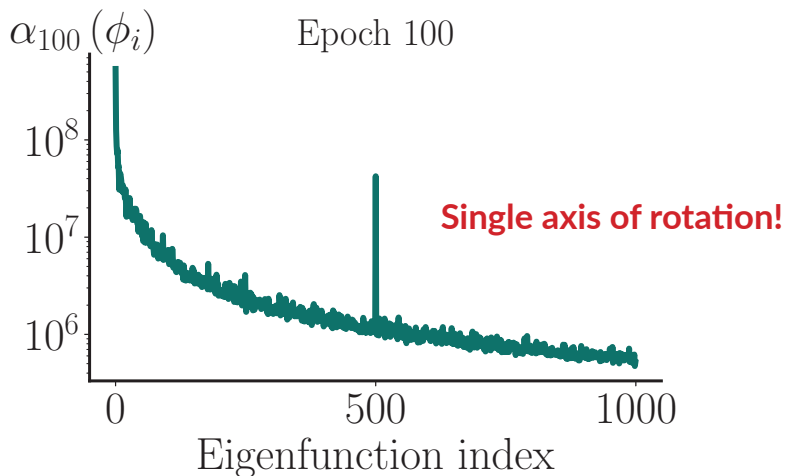
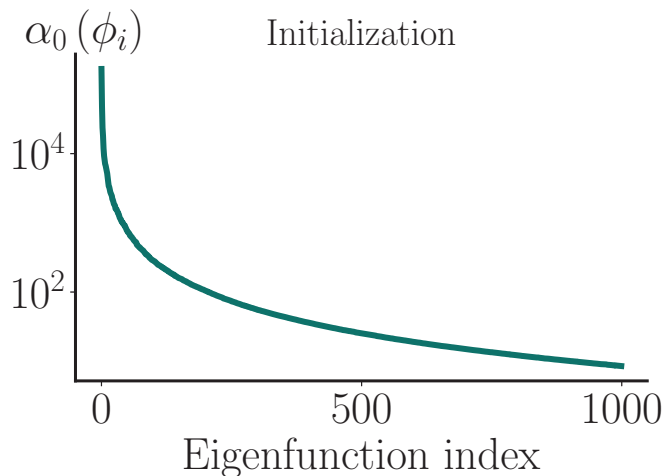
During training,  $\nabla_{\theta} f_{\theta_t}(\mathbf{x})$  evolves and so does  $\Theta_t(\mathbf{x}_1, \mathbf{x}_2) = \langle \nabla_{\theta} f_{\theta_t}(\mathbf{x}_1), \nabla_{\theta} f_{\theta_t}(\mathbf{x}_2) \rangle$



# Neural network dynamics

During training,  $\nabla_{\theta} f_{\theta_t}(\mathbf{x})$  evolves and so does  $\Theta_t(\mathbf{x}_1, \mathbf{x}_2) = \langle \nabla_{\theta} f_{\theta_t}(\mathbf{x}_1), \nabla_{\theta} f_{\theta_t}(\mathbf{x}_2) \rangle$

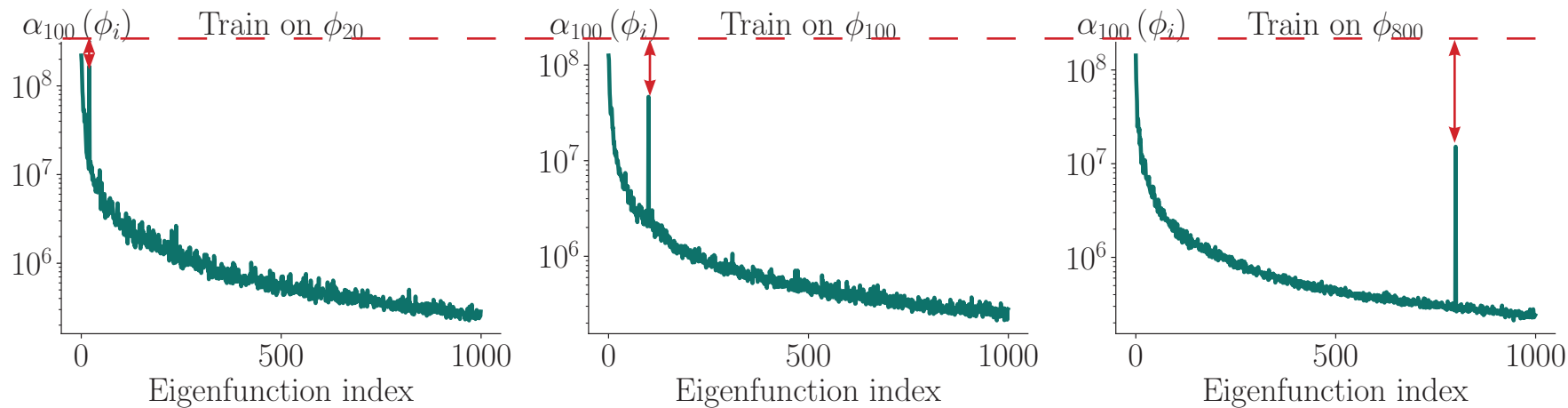
Alignment increases during training (Kopitkov & Indelman 2020, Paccolat et al. 2021, Baratin et al. 2021)

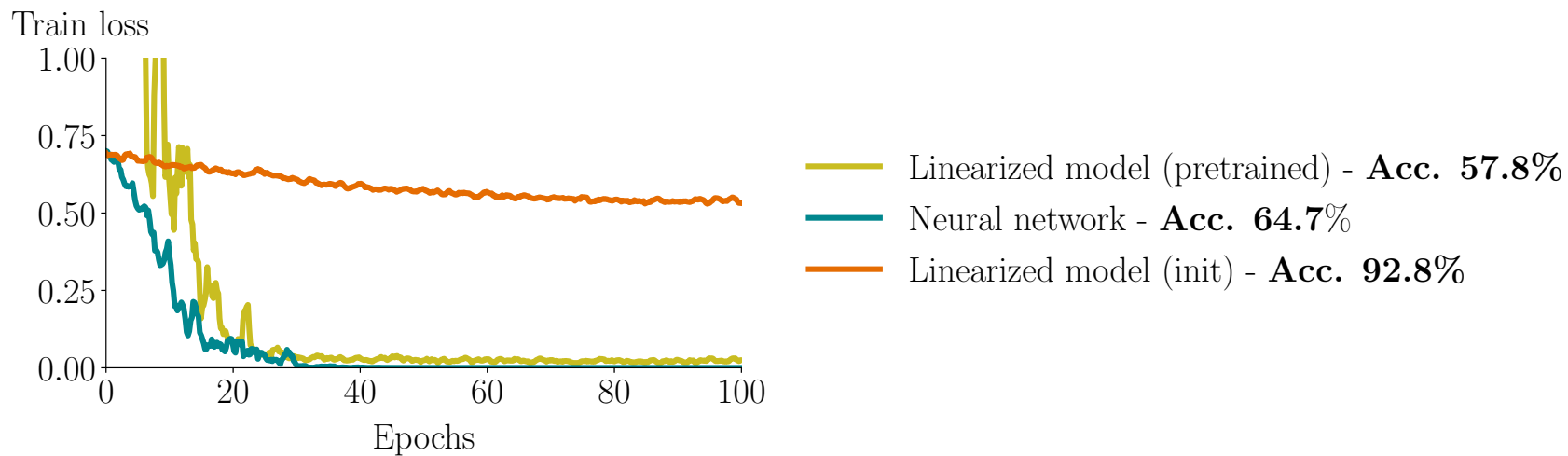


# Neural network dynamics

During training,  $\nabla_{\theta} f_{\theta_t}(\mathbf{x})$  evolves and so does  $\Theta_t(\mathbf{x}_1, \mathbf{x}_2) = \langle \nabla_{\theta} f_{\theta_t}(\mathbf{x}_1), \nabla_{\theta} f_{\theta_t}(\mathbf{x}_2) \rangle$

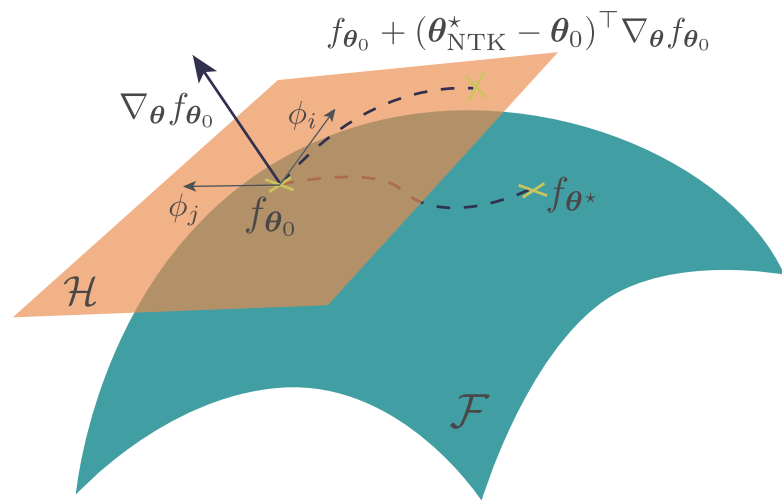
Alignment increases during training (Kopitkov & Indelman 2020, Paccolat et al. 2021, Baratin et al. 2021)





## Concluding remarks

- Linear and non-linear models agree on relative difficulty of different tasks
- Neural networks are not always better than kernel methods.
- Kernel adaptation can make neural networks overfit.

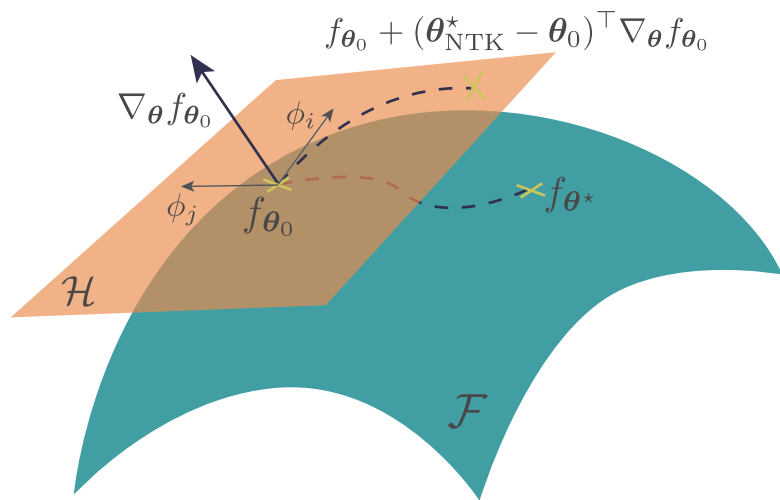


## Concluding remarks

- Linear and non-linear models agree on relative difficulty of different tasks
- Neural networks are not always better than kernel methods.
- Kernel adaptation can make neural networks overfit.

## Also in the paper...

- We delve deeper into the observations,
- use alignment to predict inductive bias,
- study role of training set size in non-linear (dis)advantage.



# What can linearized neural networks *actually* say about generalization ?

Guillermo Ortiz-Jiménez

Seyed-Mohsen Moosavi-Dezfooli

Pascal Frossard

