



AMLD22 @EPFL

AI&Physics Track

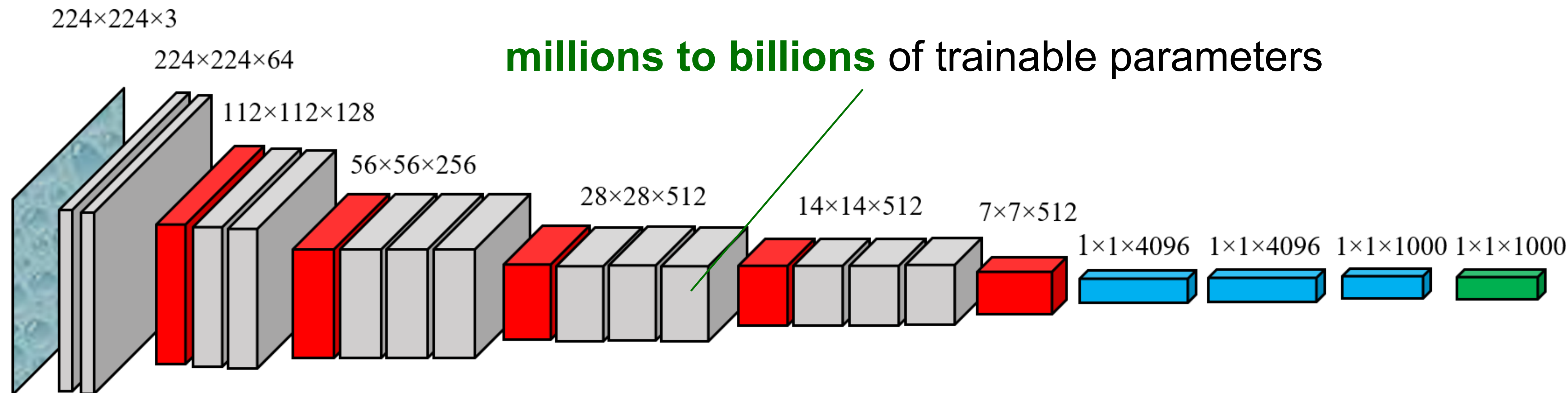
# UNIVERSAL MEAN FIELD UPPER BOUND FOR THE GENERALISATION GAP OF DEEP NEURAL NETWORKS

Pietro Rotondo — INFN, University of Milan

[S. Ariosto, R. Pacelli, F. Ginelli, M. Gherardi, PR; arXiv:2201.11022 (2022)]



# OVERPARAMETRISATION IN DEEP NETS: A BLESS FOR PRACTITIONERS, A CURSE FOR THEORISTS



**overparametrised regime**

number of trainable parameters  $\gg$  size of the training set

# STATISTICAL LEARNING THEORY IN A NUTSHELL: MAIN INGREDIENTS

◆  $P_{\mathcal{X}, \mathcal{Y}}(X, Y)$  **input/output joint probability distribution**

$\mathcal{X}$  input       $\mathcal{Y} = \{+1, -1\}$  output

$\mathcal{T}_P = \{X^\mu, Y^\mu\}_{\mu=1, \dots, P}$  **training set**

$f \in \mathcal{F}$  **hypothesis space**

◆  $\epsilon_g(f) = \langle \mathbf{1}_{f(X) \neq Y} \rangle_{P_{\mathcal{X}, \mathcal{Y}}}$

**generalisation error  
(true risk)**

$\epsilon_t(f) = \frac{1}{P} \sum_{\mu=1}^P \mathbf{1}_{f(X^\mu) \neq Y^\mu}$

**training error  
(empirical risk)**

# STATISTICAL LEARNING THEORY IN A NUTSHELL: (ONE) MAIN THEOREM

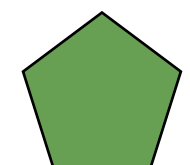
$$\Delta\epsilon(f) = \epsilon_g(f) - \epsilon_t(f) \quad \text{generalisation gap}$$

**VC dimension:** a compact measure of the expressivity of a model  $\mathcal{F}$

## **Theorem** [Vapnik-Chervonenkis]

For any  $\delta > 0$ , with probability at least  $1 - \delta$

$$\forall f \in \mathcal{F}, \quad \Delta\epsilon(f) \leq 2 \sqrt{2 \frac{d_{\text{VC}} \log \frac{eP}{d_{\text{VC}}} + \log \frac{2}{\delta}}{P}}$$



**uniform** in the functions of the model and **data-independent**

# STATISTICAL LEARNING THEORY: MAIN LIMITATIONS

$$\Delta\epsilon \lesssim \sqrt{\frac{d_{VC}}{P}}$$

the VC dimension of a DNN is (very) roughly proportional to the number of trainable parameters  $\sim 10^6 - 10^9$

the typical size of a training dataset in a supervised learning problem is of order  $\sim 10^4 - 10^6$



*“[...] Their derivation reveals many possible causes for their poor quantitative performance:*

- (i) Practical data distributions may lead to smaller deviations (between the expected and empirical classification error) than the worst possible data distribution.*
- (ii) Uniform bounds hold for all possible classification functions. Better bounds may hold when one restricts the analysis to functions that perform well on plausible training sets.”*

*(from L. Bottou, “Making Vapnik-Chervonenkis bounds accurate”)*

**MAIN GOAL: Improve this bound with Statistical Physics**

# THE OTHER MAJOR FRAMEWORK TO INVESTIGATE GENERALISATION: THE TEACHER-STUDENT SCENARIO

$$P(\mathbf{x}, y) = \rho(\mathbf{x}) \delta(y - f_T(\mathbf{x})) \quad f_T(\mathbf{x}) = \frac{1}{\sqrt{N_T}} \sum_{\alpha=1}^{N_T} t_{\alpha} \phi_{\alpha}^{(T)}(\mathbf{x})$$

input density distribution

a **teacher** provides the ground truth (the label)

$$f_S(\mathbf{x}) = \frac{1}{\sqrt{N_S}} \sum_{\alpha=1}^{N_S} v_{\alpha} \phi_{\alpha}^{(S)}(\mathbf{x})$$

a **student** optimises its weights to match the ground truth

**Goal:** compute the optimal generalisation and training errors for large  $N_S$  and  $P$

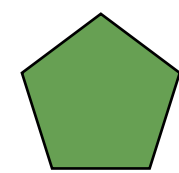
◆ **linear** teacher, **linear** student, **factorised** input density is a textbook exercise

◆ **polynomial** teacher, **polynomial** student, **factorised** input density

# RECENT RESULTS: GENERALISATION AND TRAINING ERRORS FOR GENERIC KERNELS/GENERIC (QUENCHED) FEATURES

[A. Canatar, B. Bordelon, C. Pehlevan; Nat. Comm. (2021)]

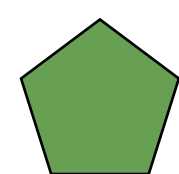
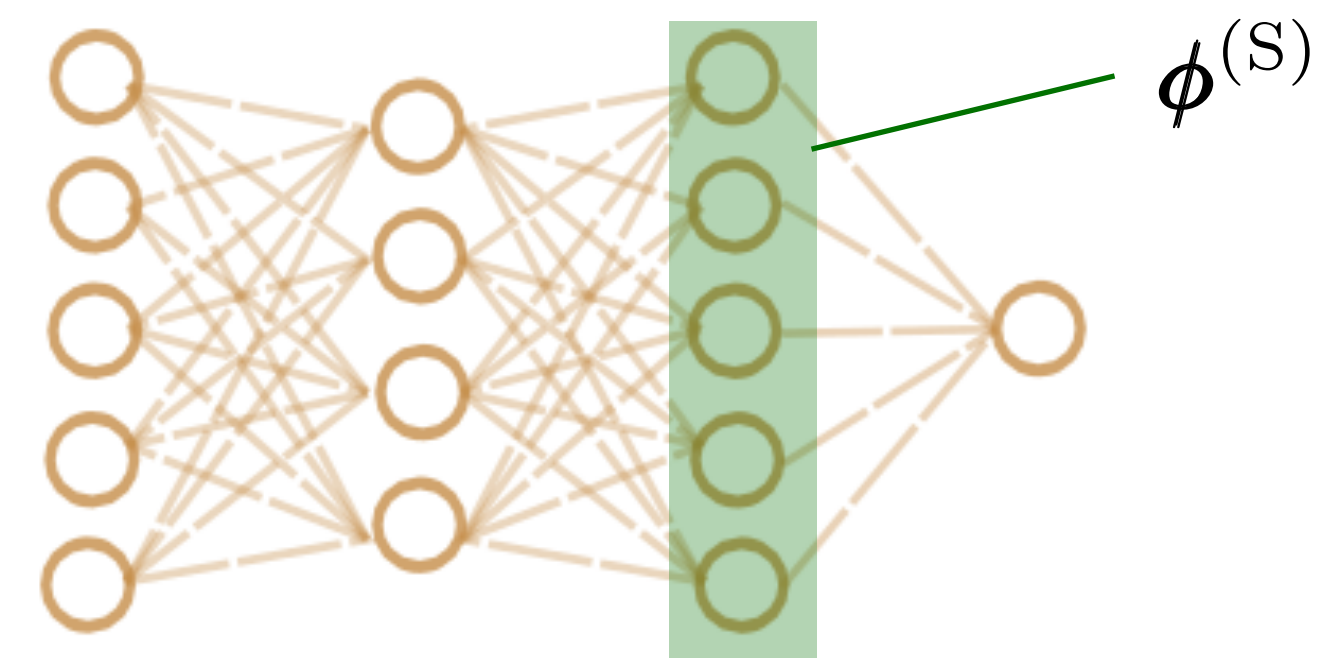
[B. Loureiro, C. Gerbelot, H. Cui, S. Goldt, F. Krzakala, M. Mézard, L. Zdeborová; NeurIPS (2021)]



Exact **formulas for the generalisation and training errors.**

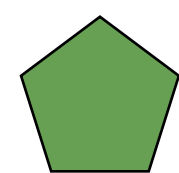
Particularly simple for regression problems and quadratic loss function

$$\epsilon_{g/t} \left( N_S, P, \phi^{(T)}, \phi^{(S)} \right)$$



These formulas **capture the learning curves of multilayer neural networks**

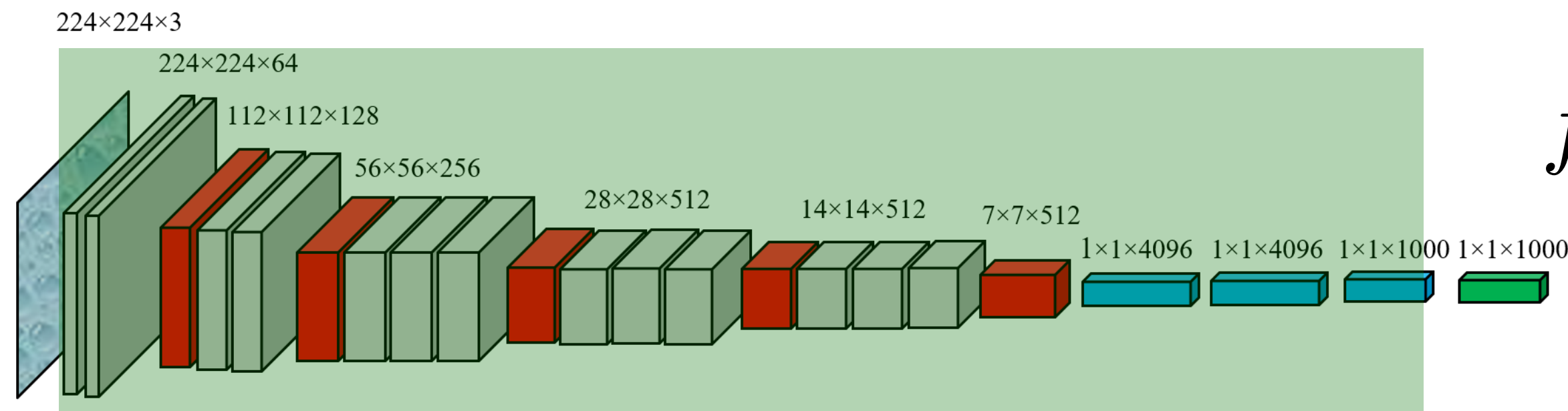
if we consider as features those obtained by pre-trained networks on realistic datasets!



Based on a conjecture: the **Gaussian Equivalence Principle**



# RESULTS: FROM QUENCHED FEATURES TO A UNIVERSAL MEAN FIELD UPPER BOUND FOR THE GENERALISATION GAP OF DNNs



$$\phi_{\alpha}^{\text{DNN}}(\mathbf{x}, \mathcal{W})$$

$$f_{\text{DNN}}(\mathbf{x}, \theta) = \frac{1}{N_{\text{out}}} \sum_{\alpha=1}^{N_{\text{out}}} v_{\alpha} \phi_{\alpha}^{\text{DNN}}(\mathbf{x}, \mathcal{W})$$

$$\theta = \{\mathcal{W}, \mathbf{v}\}$$

number of weights in the last layer

$$N_{\text{out}} \ll P$$

$10^2 - 10^3$        $10^4 - 10^5$

the equation holds for each realisation of the weights  $\mathcal{W}$  and it assumes perfect training over the last layer

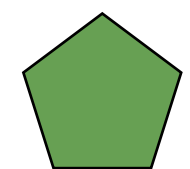
$$0 \leq \epsilon_{\text{g}}^{\text{R}}(\mathcal{W}) \leq T$$

$$\Delta \epsilon(\mathcal{W}) \simeq 2 \epsilon_{\text{g}}^{\text{R}}(\mathcal{W}) \frac{N_{\text{out}}}{P}$$

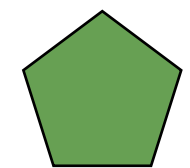
# RESULTS: FROM QUENCHED FEATURES TO A UNIVERSAL MEAN FIELD UPPER BOUND FOR THE GENERALISATION GAP OF DNNs

$$\Delta \tilde{\epsilon}(\mathcal{W}) = \frac{\Delta \epsilon(\mathcal{W})}{T} \leq \frac{2N_{\text{out}}}{P}$$

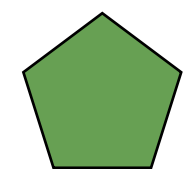
valid for any local minimum of the loss function of the DNN



the gap of fully-trained DNNs should **decrease at least as 1/P asymptotically**

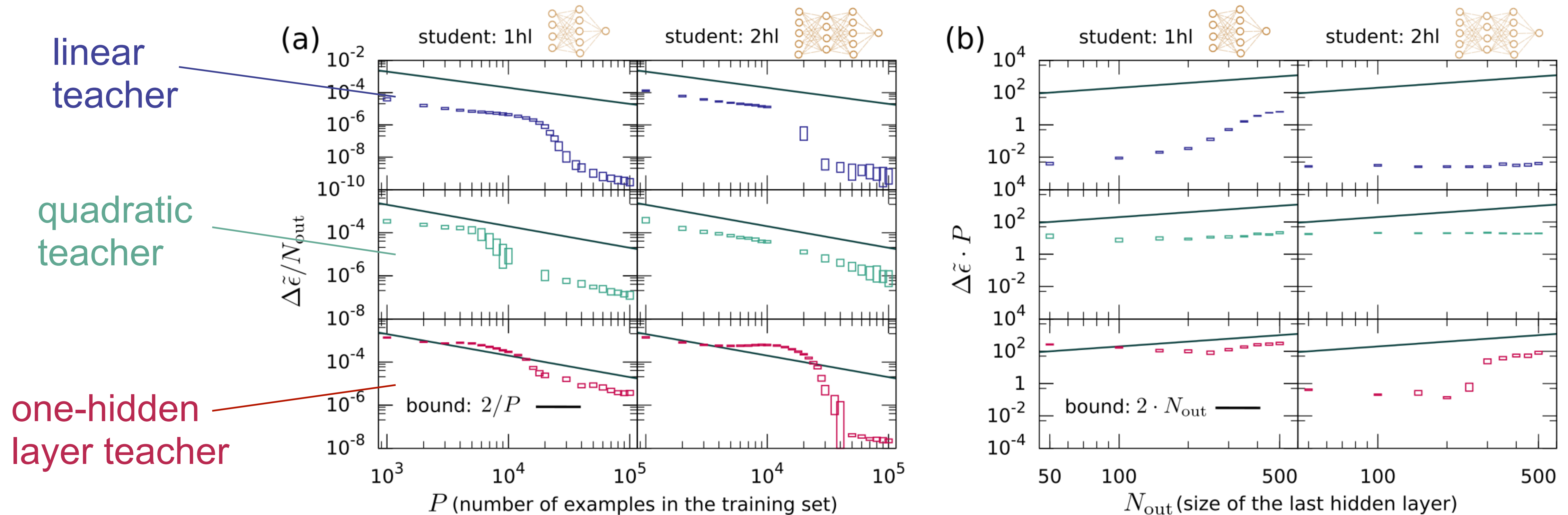


the **degradation** of the generalisation performance should be **at most linear as the size of the last layer is increased**



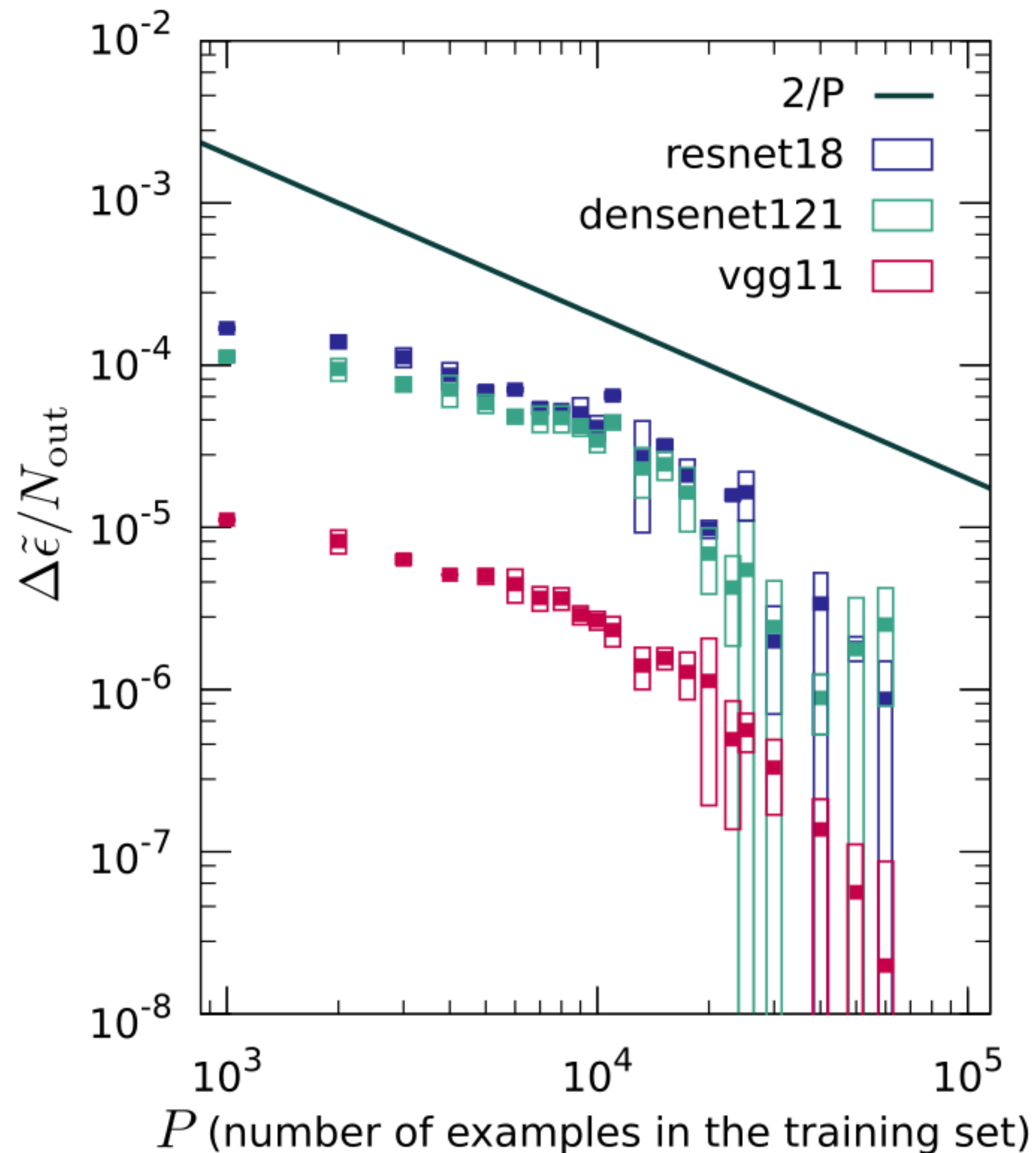
the bound **rules out** any asymptotic **linear or sub-linear dependence on the size of the hidden layers**

# RESULTS: GENERALISATION GAP FOR TOY FULLY CONNECTED STUDENTS TRAINED ON SYNTHETIC DATASETS



In all these experiments the input density is factorised over its coordinates

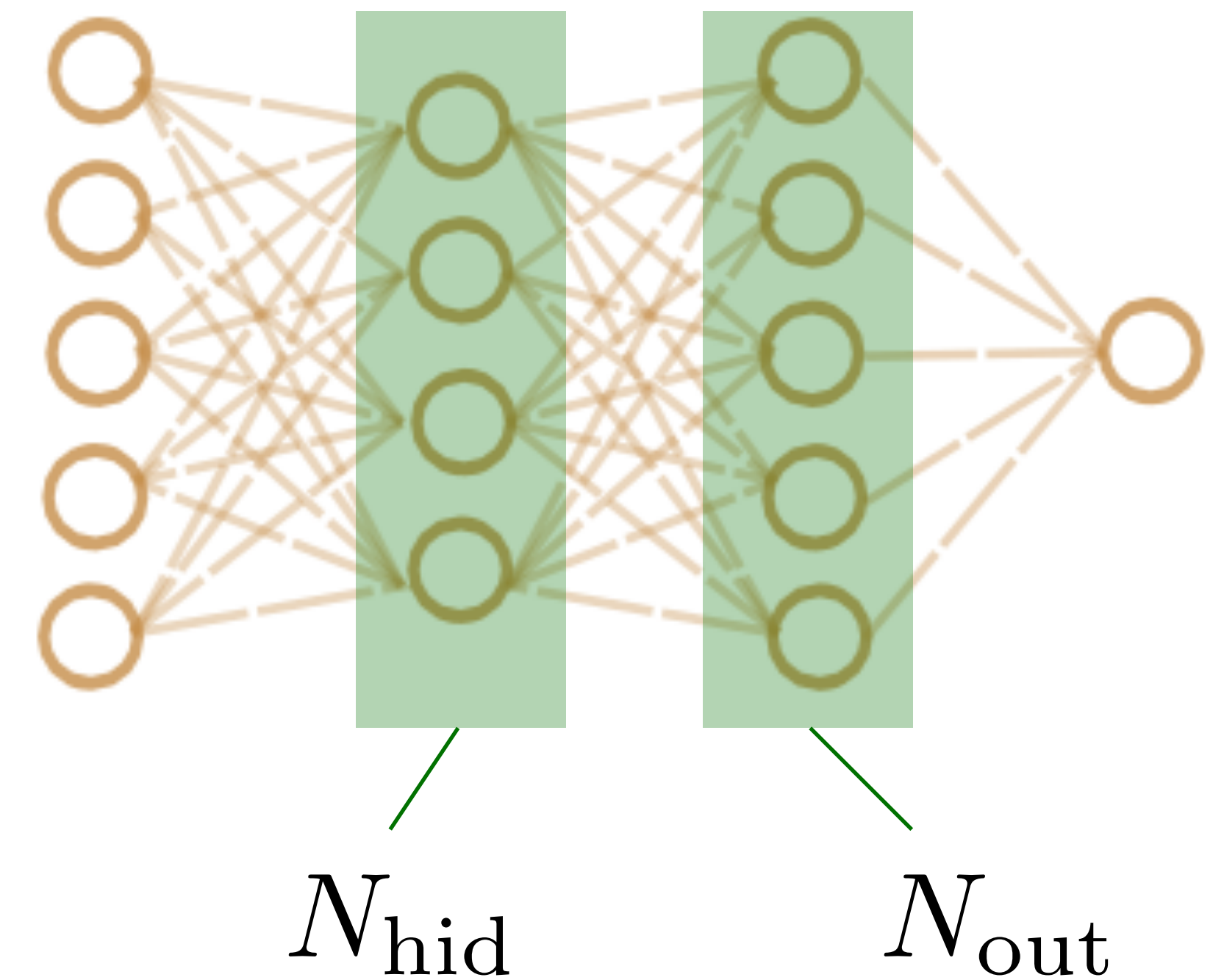
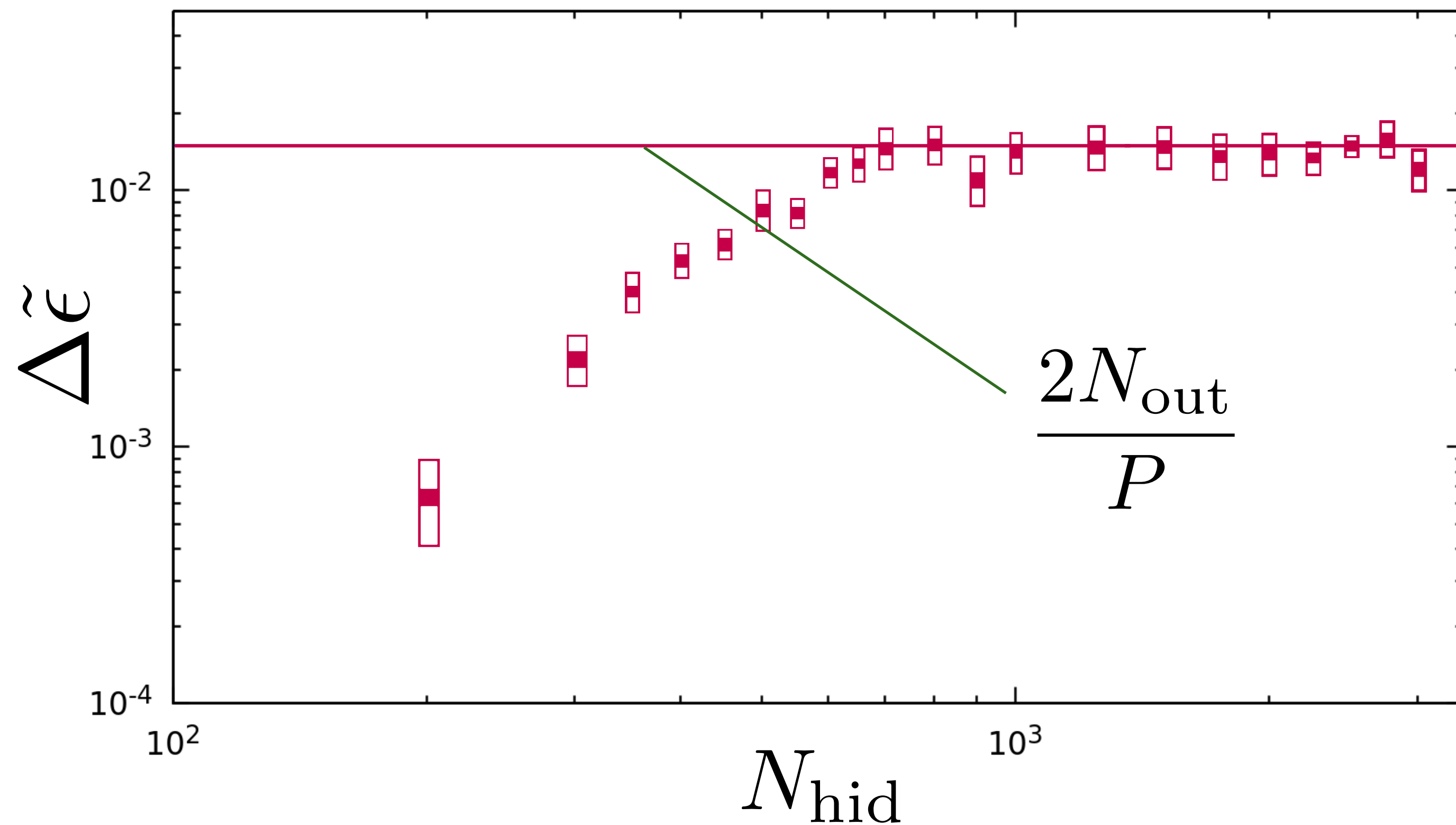
# RESULTS: GENERALISATION GAP FOR STATE-OF-THE-ART ARCHITECTURES TRAINED ON MNIST



**Remark:** the generalisation gap is defined for regression, **not** for classification

# RESULTS: THE FORM OF THE BOUND RULES OUT ANY LINEAR OR SUBLINEAR DEPENDENCE OF THE GAP ON THE SIZE OF HIDDEN LAYERS

two hidden layer student learns a one hidden layer teacher



Not only a universal upper bound, but a universal state equation?

# CONCLUSIONS AND FUTURE PERSPECTIVES

- A **more stringent** and **universal** asymptotic upper bound for the generalisation gap of DNNs
- One (out of many) possible next step: finding the most general set of hypotheses such that the deep gaussian equivalence principle holds

[S. Goldt, B. Loureiro, G. Reeves, F. Krzakala, M. Mézard, L. Zdeborová; PMLR (2021)]  
[A. Montanari, B. Saeed; arXiv:2202.08832 (2022)]

- The approach (at the moment) is **FAITH** in spirit

**FAITH** = **F**airly **A**ccurate **I**ntuition **T**hrough **H**andwaving

# THANKS!



**Sebastiano Ariosto**



**Rosalba Pacelli**



**Francesco Ginelli**



**Marco Gherardi**

[S. Ariosto, R. Pacelli, F. Ginelli, M. Gherardi, PR; arXiv:2201.11022 (2022)]