# Scalable Marginal Likelihood Estimation for Model Selection in Deep learning

## ICML 2021

Alexander Immer

ETH Zurich

Matthias Bauer

Deepmind

Vincent Fortuin

ETH Zurich

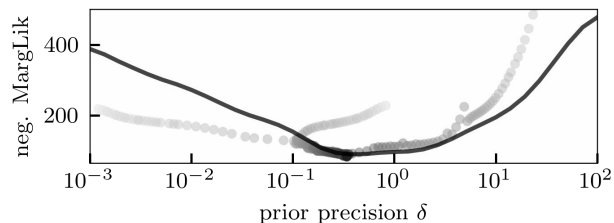Gunnar Rätsch

ETH Zurich

Emtiyaz Khan

RIKEN AIP

# Model Selection in Deep Learning

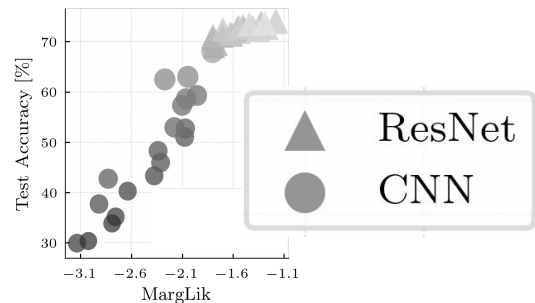**(1)** hyperparameters (regularization) and **(2)** model architecture (ResNet vs CNN).

But **validation data** might be **unavailable** (e.g. in continual learning).

**We show the training marginal likelihood is viable for model selection in DL!**

**(1)** Differentiable hyperparameters *during training*



**(2)** Architecture Selection *after training*

# Marginal Likelihood Estimation for Deep Learning

① **Laplace approximation** [1] to the log marginal likelihood

$$\mathbf{H}_\theta = \nabla_{\theta\theta}^2 \log p(\mathcal{D}, \theta | \mathcal{M})$$

$$\log p(\mathcal{D}|\mathcal{M}) \approx \underbrace{\log p(\mathcal{D}|\theta_*, \mathcal{M})}_{\text{Training data fit}} + \underbrace{\log p(\theta_*|\mathcal{M}) - \frac{1}{2}\log\left|\frac{1}{2\pi}\mathbf{H}_{\theta_*}\right|}_{\text{Complexity penalty}}$$
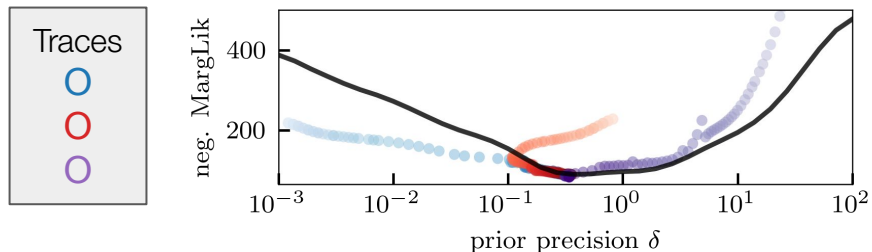
② **Scalable approximations** to the Hessian

| | Approx. types | | Correlation captured |
|---|---|---|---|
| $\mathbf{H}_{\theta_*} \approx$ | Gauss-Newton | | Full |
| | Fisher Information | $\times$ | KFAC (block-diagonal) [3, 4] |
| | Empirical Fisher | | Diagonal |

[1] MacKay. "A practical Bayesian framework for backpropagation networks." Neural computation (1992).
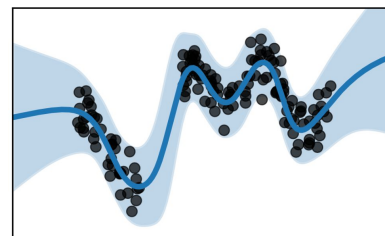
# Maximizing the Marginal Likelihood during Training

**Optimize Hyperparameters during Training
(e.g. regularization)**

**Compare Architectures
(e.g. #layers)**



Traces
○ (blue)
○ (red)
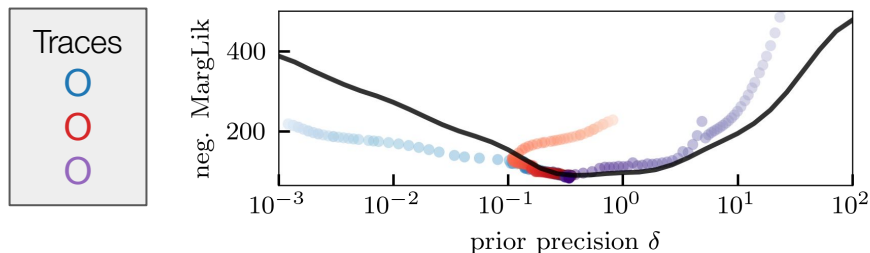○ (purple)

3 layers, 5221 params
MargLik = $-88$

1 layer, 151 params
MargLik = $-110$

Every epoch:
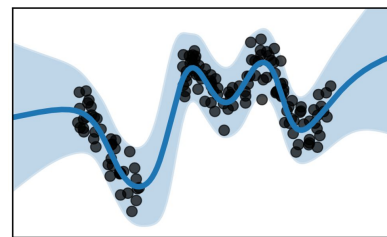    Update network parameters (e.g. Adam)
    Differentiate MargLik wrt. hyperparameters
    Update differentiable hyperparameters
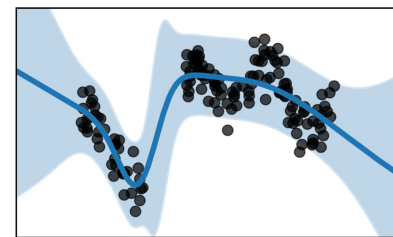
# Maximizing the Marginal Likelihood during Training

**Optimize Hyperparameters during Training
(e.g. regularization)**

**Compare Architectures
(e.g. #layers)**



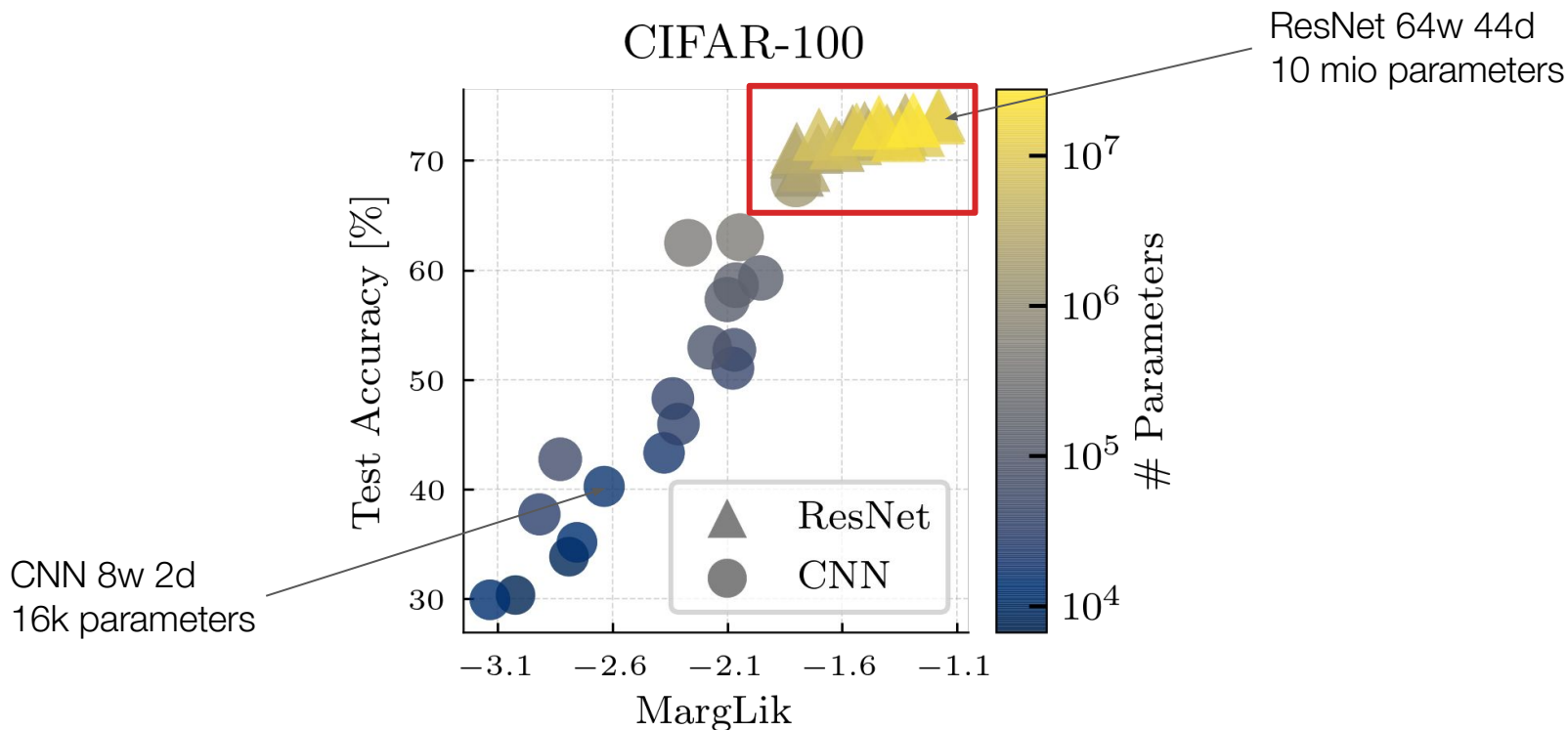3 layers, 5221 params
MargLik $= -88$

1 layer, 151 params
MargLik $= -110$

- On par or better than cross-validation
  - UCI regression/classification, image classification
- Several hundred hyperparameters at once
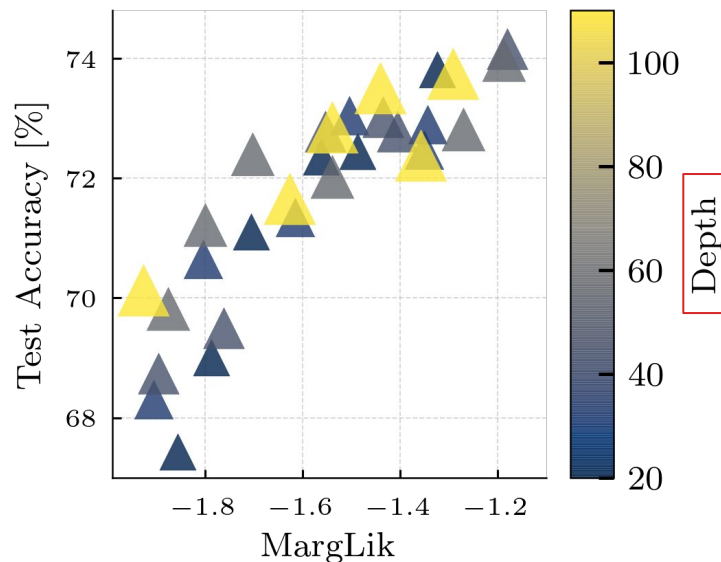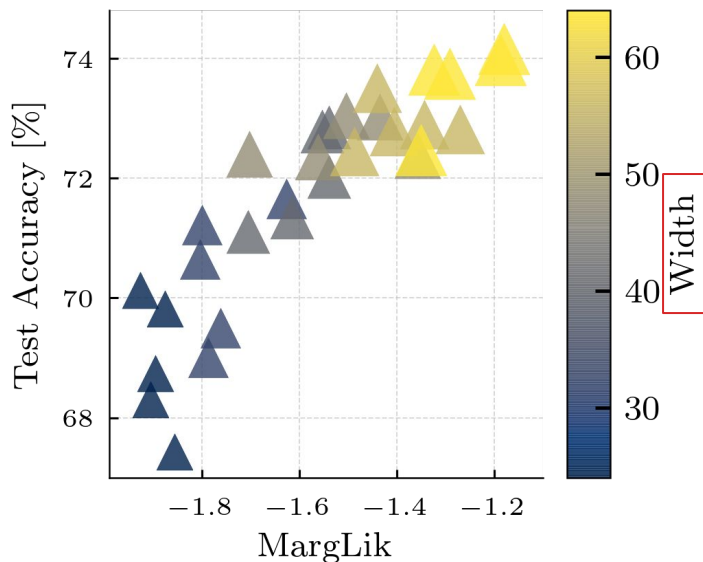  - No overhead for some approximations

# Marginal Likelihood for Architecture Comparison

Two architectures (**CNN, ResNet**) + varying **width** (≤ 64) and **depth** (≤ 110)



CIFAR-100

ResNet 64w 44d
10 mio parameters

CNN 8w 2d
16k parameters

# Marginal Likelihood for Architecture Comparison

ResNets of varying width (≤ 64) and depth (≤ 110)



→ In line with proposed Wide ResNet architecture [5]

# Summary

- Marginal likelihood viable for model selection in DL **without validation data**

- Optimize margLik: **hundreds of hyperparameters during training**

- **Model comparison across architectures** seems possible

# References (abbreviated)

[1] MacKay: *"A practical Bayesian framework for backpropagation networks"*

[2] Rasmussen, et al.: *"Occam's razor"*

[3] Martens, et al.: *"Optimizing neural networks with KFAC"*

[4] Ritter, et al.: *"A Scalable Laplace Approximation for Neural Networks"*

[5] Zagoruyko, et al.: *"Wide Residual Networks"*