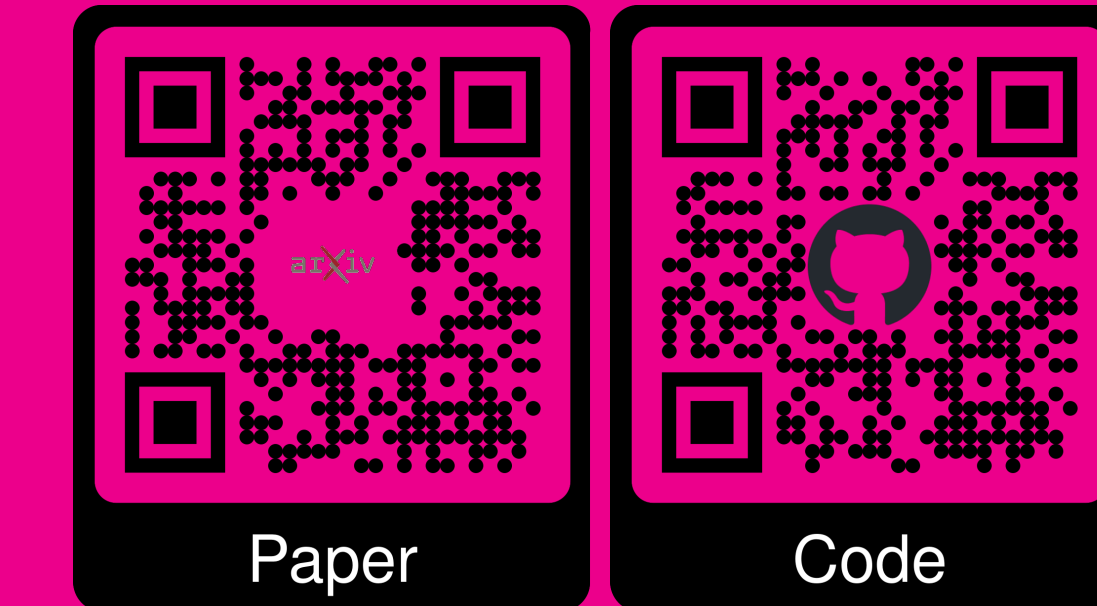


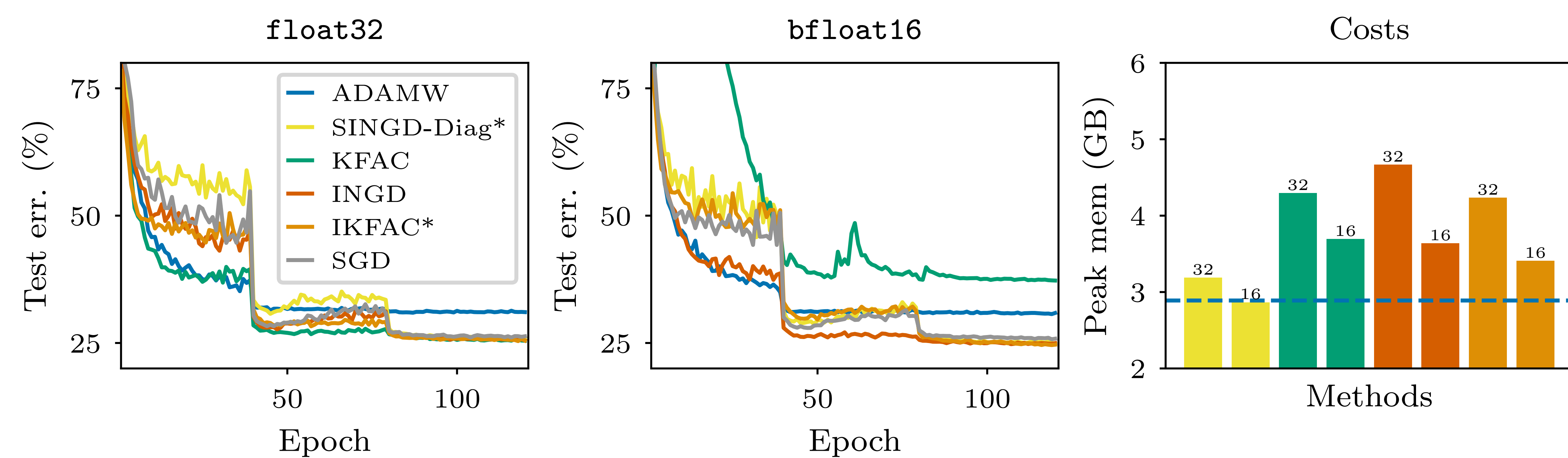
SINGD: Structured Inverse-Free Natural Gradient Descent

Wu Lin*, Felix Dangel*, Runa Eschenhagen, Kirill Neklyudov, Agustinus Kristiadi, Richard Turner, Alireza Makhzani

We present a **KFAC-style 2nd-order method** for DL that is **memory-efficient & numerically-stable in bfloat16** thanks to structured Kronecker factors and inverse-free update rule.



Inefficiencies of KFAC: Numerically Unstable In bfloat16 & Memory-Intensive due to Dense Kronecker Factors

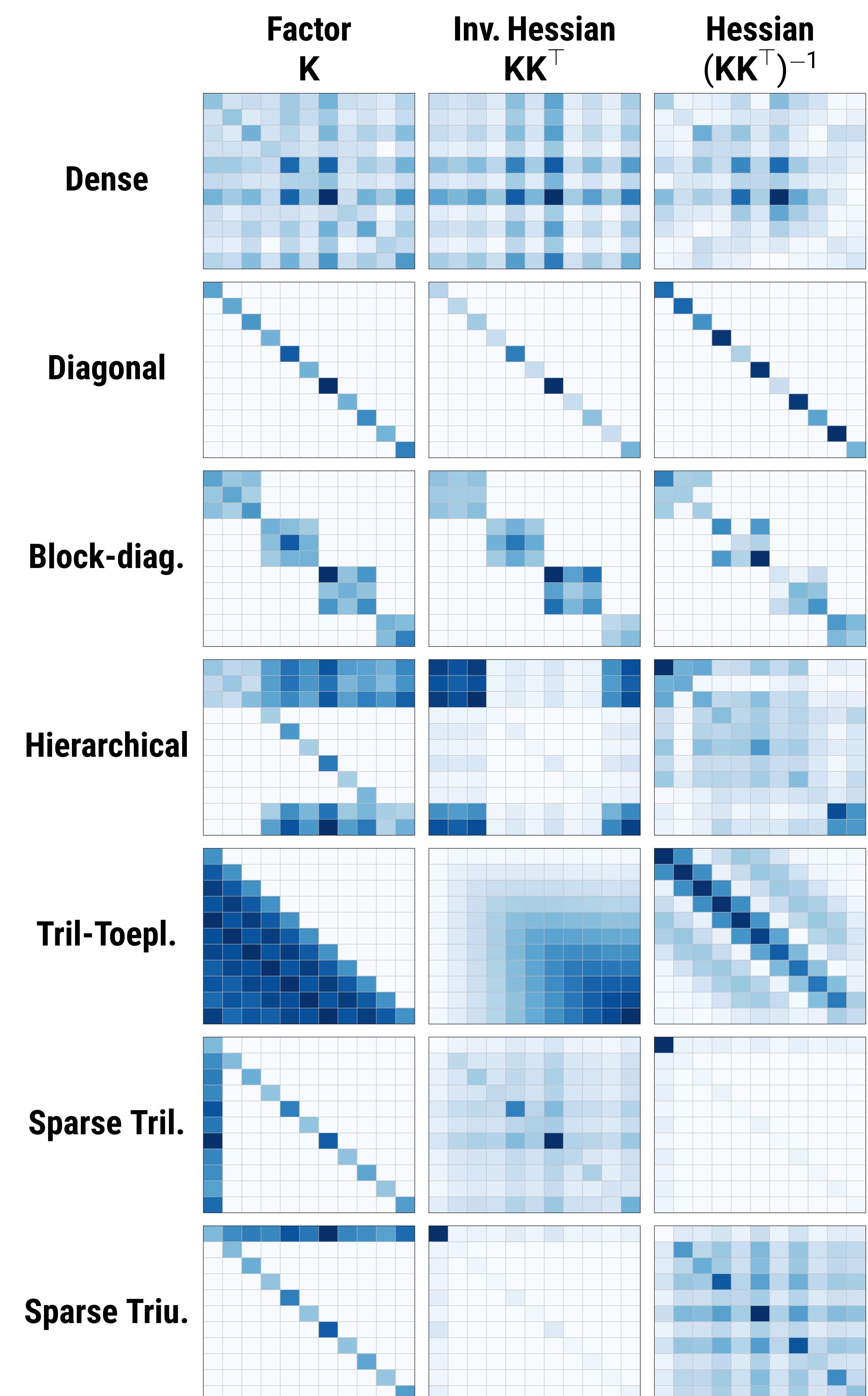


Computational efficiency on VGG + CIFAR-100

Method	Memory	Time
SGD, bfloat16	1.00 x	1.00 x
AdamW, bfloat16	1.02 x	1.07 x
SINGD-Diag*, bfloat16	1.02 x	1.29 x
IKFAC*, bfloat16	1.21 x	1.84 x
INGD, bfloat16	1.29 x	1.84 x
KFAC, float32	1.52 x	4.49 x

Available Structures

Different structures capture different information.

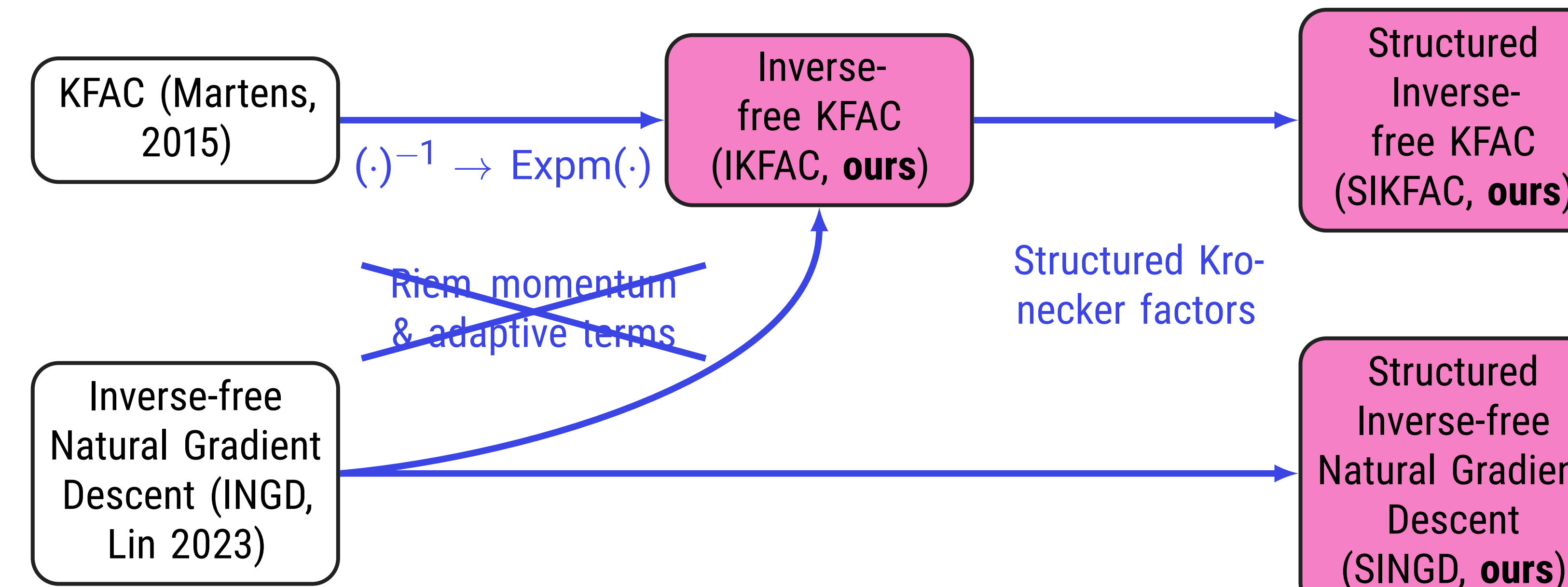


SINGD = Inverse-free KFAC + Structures

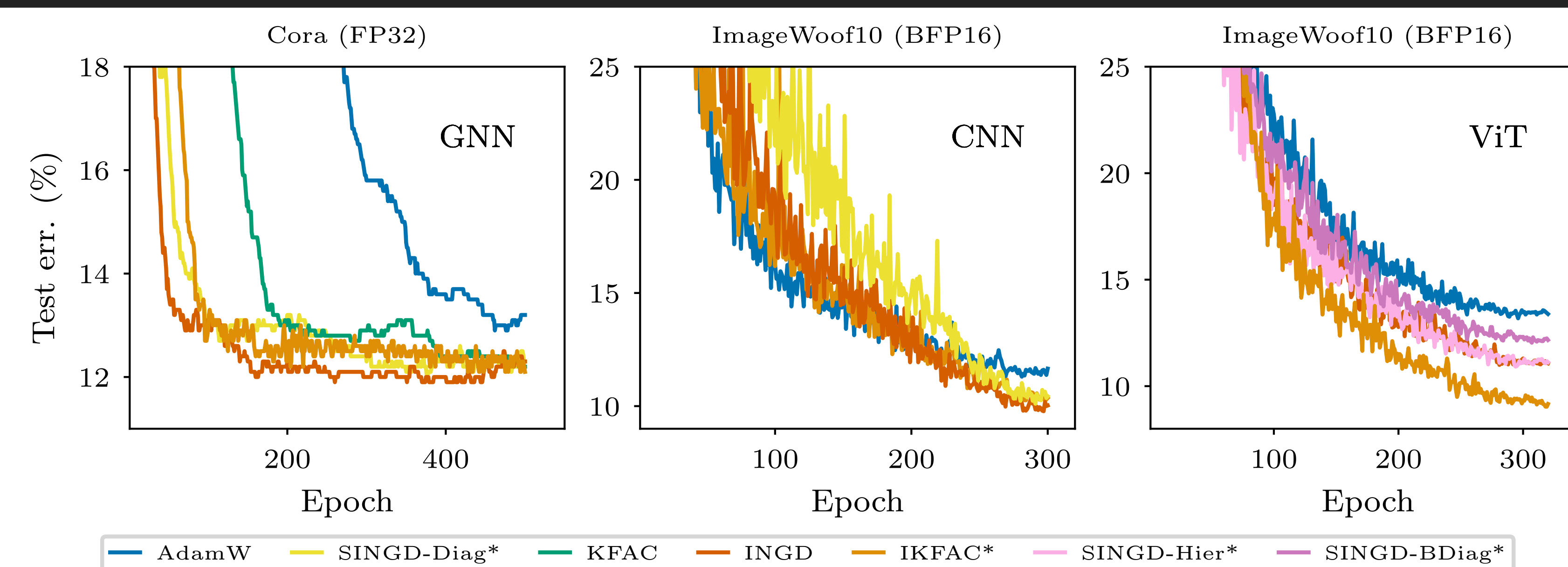
Goal: Approximate natural gradient $F^{-1}g$ with gradient g , inverse Fisher approximation F^{-1} , and damping λ . Per layer:

$$\begin{aligned} \text{KFAC: } F^{-1} &= (S_K + \lambda I)^{-1} \otimes (S_C + \lambda I)^{-1}, && \text{updates } S_K, S_C \\ \text{INGD: } F^{-1} &= KK^T \otimes CC^T, && \text{updates } K, C \\ \text{SINGD: } F^{-1} &= KK^T \otimes CC^T, && \text{with structured } K, C \end{aligned}$$

Theorem 1 (informal): At each iteration, **IKFAC's update** of the Kronecker factors K/C **matches KFAC's update** of $(S_{K/C} + \lambda I)^{-1}$ **up to first order** in the exponential moving average (default $\lesssim 10^{-2}$), assuming same initialization and sequence of layer statistics.



Performs Well on Many Problems & in bfloat16



Works like Any Optimizer (pip install singd)

```
1 from singd.optim.optimizer import SINGD
2 from torch_utils import CrossEntropyLoss, alexnet, new_batch
3
4 # 1. Setup
5 model, loss_func = alexnet(), CrossEntropyLoss()
6 X, y = new_batch()
7
8 optim = SINGD( # Optimizer setup
9     model,
10    kfac_like=False, # Use True for IKFAC
11    structures=("dense", "diagonal"), # Structure of (K, C)
12)
13
14 # 2. One training step
15 optim.zero_grad() # Clear previous gradients
16 loss_func(model(X), y).backward() # Forward + backward pass
17 optim.step() # Update pre-conditioner and parameters
```