



AdaFisher: Adaptive Second Order Optimization via Fisher Information



**Damien Martins Gomes^{1,2}, Yanlei Zhang³, Eugene Belilovsky^{1,3},
Guy Wolf^{3,4}, Mahdi S. Hosseini^{1,3}**



¹Concordia University, ²IPSA Toulouse, ³Quebec AI Institute – Mila,
⁴Université de Montréal

ICLR 2025 – Singapore



ICLR

Motivation

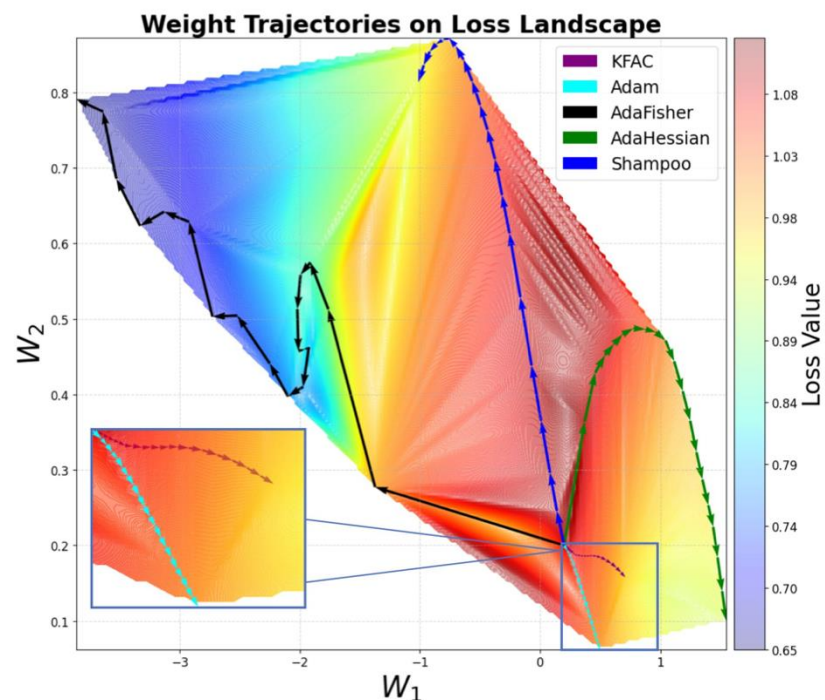
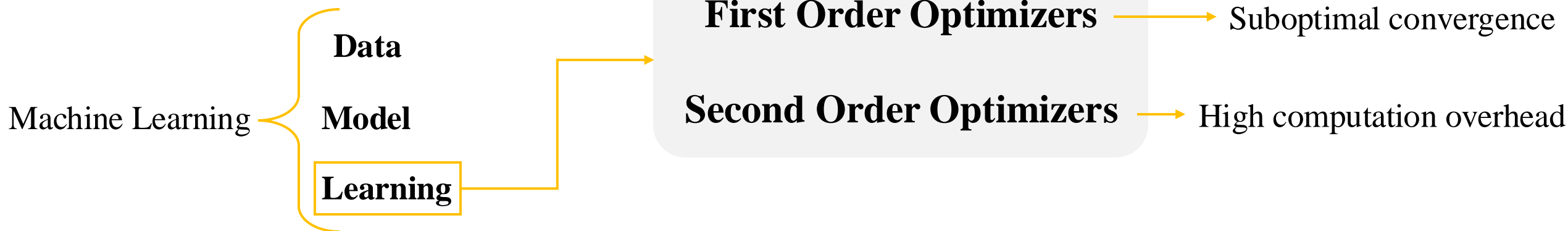


Figure 1: Visualizing optimization trajectories for various optimizers overlaid a loss landscape.

Can we design an optimizer that balances computational complexity (time and space) with strong generalization (fast convergence)?

AdaFisher

Methodology

Fisher Matrix Insights

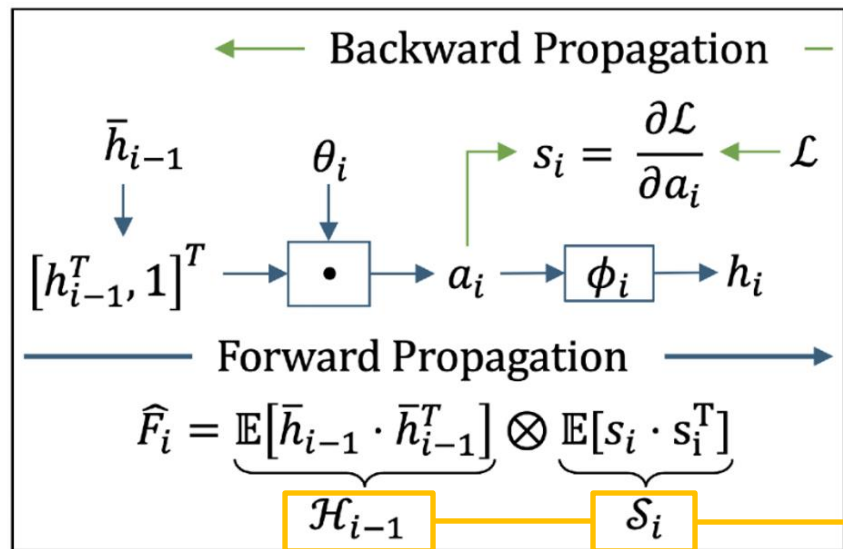


Figure 2: Illustration of EFIM computation using K-FAC [1] for a given layer i

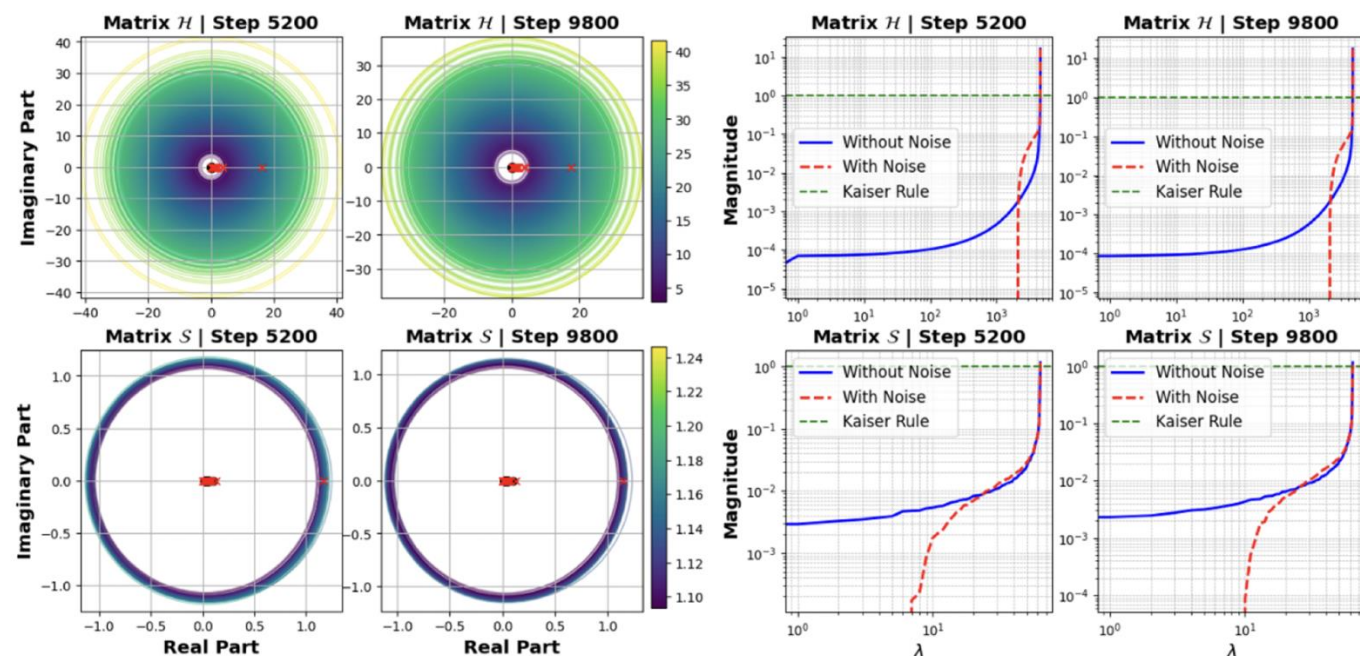
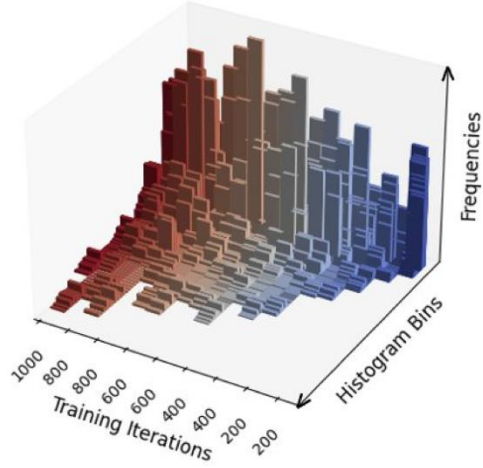


Figure 3: Gershgorin disks and eigenvalue perturbations in the 37th Convolutional layer of a ResNet-18 at steps 5200 (middle of training) And 9800 (end of training). Left: Gershgorin circles; Right: Eigenvalue spectrum with/without noise.

New efficient approximation to the FIM

Methodology

Histogram of FIM Diagonal for Adam



Histogram of FIM Diagonal for AdaFisher

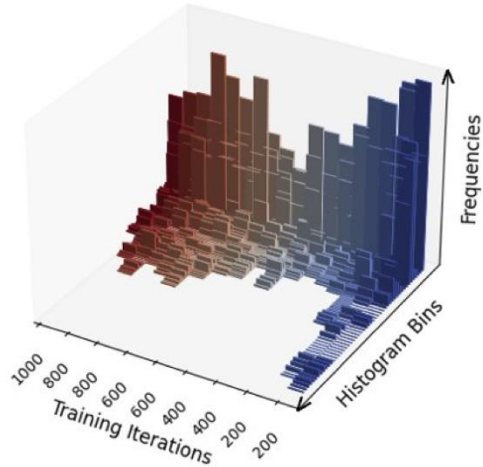


Figure 4: Comparison of FIM diagonal histograms during ResNet18 training on CIFAR10: The figure displays the FIM diagonal elements for the first convolutional layer with Adam and AdaFisher over 1,000 training iterations.

**Better
Generalization
than Adam**

Algorithm 1 AdaFisher optimization algorithm. Good default settings for the tested machine learning problems are $\alpha = 0.001$ (learning rate), $\lambda = 0.001$ (Tikhonov damping parameter), $\gamma = 0.8$ (Exponentially decaying factor). [Default parameters are: $\beta = 0.9$ (Exponentially decaying factor of Adam), κ (weight decay) (Kingma & Ba (2015), Loshchilov & Hutter (2019b))].

Require: Step size α ; Exponential decay rate for KFs $\gamma \in [0, 1]$; Tikhonov damping parameter λ ; Exponential decay rate for first moments β in $[0, 1]$; Initial parameters θ

Initialize 1st moment variable $m = 0$; FIM $\tilde{F}_{D_i} = \mathbf{I}$; time step $t = 0$

- 1: **while** stopping criterion not met **do**
- 2: Sample a minibatch of M examples from the training set $\{(x_n, y_n)\}_{n=1}^M$
- 3: Compute $\mathcal{H}_{D_{i-1}}, \mathcal{S}_{D_i}$ for $i \in \{1, \dots, L\}$ using Section A.3 (notice that: $\mathcal{H}_{D_0} = x$)
- 4: Compute EMAs of $\mathcal{H}_{D_{i-1}}$ and \mathcal{S}_{D_i} using Eq. (3)
- 5: Compute \tilde{F}_{D_i} for $i \in \{1, \dots, L\}$ using Eq. (4)
- 6: $g^{(t)} \leftarrow \frac{1}{M} \sum_n \nabla_{\theta^{(t)}} \mathcal{L}(f(x_n; \theta^{(t)}), y_n)$ (Compute gradient)
- 7: $m^{(t+1)} \leftarrow \frac{\beta m^{(t)} + (1-\beta)g^{(t)}}{1-\beta^t}$ (Update and correct biased first moment)
- 8: **Case AdaFisher:** $\Delta\theta^{(t)} = -\alpha(\tilde{F}_D^{(t)})^{-1}m^{(t)}$
 Case AdaFisherW: $\Delta\theta^{(t)} = -\alpha((\tilde{F}_D^{(t)})^{-1}m^{(t)} + \kappa\theta^{(t)})$
- 9: $\theta^{(t+1)} \leftarrow \theta^{(t)} + \Delta\theta^{(t)}$ (Apply update)
- 10: $t \leftarrow t + 1$
- 11: **end while**

Results

Image Classification

Table 2: Validation of ImageNet-1K / ResNet50 by different optimizers reported on Top-1 and Top-5 accuracy.

Optimizers	Batch size	Top-1	Top-5
Adam	256	67.78	88.37
K-FAC	256	70.96	89.44
Shampoo	256	72.82	91.42
AdaFisher	256	76.95	93.39
AdaFisher	512	77.01	93.45
AdaFisher	1024	77.09	93.56
SGD Goyal et al. (2017)	256	76.40	-
AdamW Chen et al. (2024)	1024	76.34	-
LAMB You et al. (2019)	16K	76.66	93.22
SGD You et al. (2019)	16K	75.20	-
LARS Huo et al. (2021)	16K	75.1	-

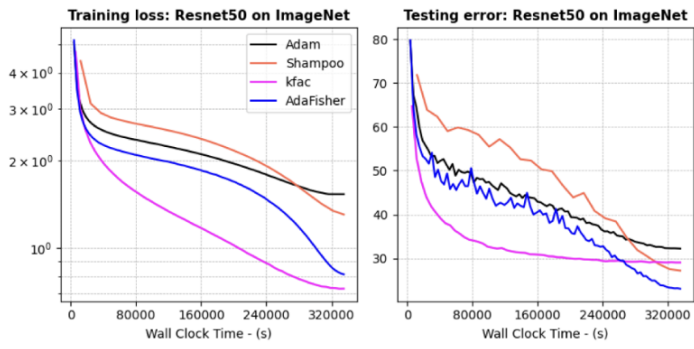


Figure 5: Training loss and validation error of ResNet-50 on ImageNet-1k. AdaFisher consistently achieves lower test error as compared to its counterparts.

Table 3: Performance metrics (mean, std) of different networks and optimizers on CIFAR10 and CIFAR100 using batch size 256 with a 200-epoch AdaFisher training cutoff

Network	CIFAR10						CIFAR100					
	SGD	Adam	AdaHessian	K-FAC	Shampoo	AdaFisher	SGD	Adam	AdaHessian	K-FAC	Shampoo	AdaFisher
ResNet18	95.64 _{0.1}	94.85 _{0.1}	95.44 _{0.1}	95.17 _{0.2}	94.08 _{0.2}	96.25_{0.2}	76.56 _{0.2}	75.74 _{0.1}	71.79 _{0.2}	76.03 _{0.3}	76.78 _{0.2}	77.28_{0.2}
ResNet50	95.71 _{0.1}	94.45 _{0.2}	95.54 _{0.1}	95.66 _{0.1}	94.59 _{0.1}	96.34_{0.2}	78.01 _{0.1}	74.65 _{0.5}	75.81 _{0.3}	77.40 _{0.4}	78.07 _{0.4}	79.77_{0.4}
ResNet101	95.98 _{0.2}	94.57 _{0.1}	95.29 _{0.6}	96.01 _{0.1}	94.63 _{0.1}	96.39_{0.1}	78.89 _{0.2}	75.56 _{0.3}	73.38 _{0.2}	77.01 _{0.4}	78.83 _{0.2}	80.65_{0.4}
DenseNet121	96.09 _{0.1}	94.86 _{0.1}	96.11 _{0.1}	96.12 _{0.1}	95.66 _{0.1}	96.72_{0.1}	80.13 _{0.4}	75.87 _{0.4}	74.80 _{0.9}	79.79 _{0.2}	80.24 _{0.3}	81.36_{0.3}
MobileNetV3	94.43 _{0.2}	93.32 _{0.1}	92.86 _{3.1}	94.34 _{0.1}	93.81 _{0.2}	95.28_{0.1}	73.89 _{0.3}	70.62 _{0.3}	56.58 _{4.5}	73.75 _{0.3}	70.85 _{0.3}	77.56_{0.1}
Tiny Swin	82.34 _{0.2}	87.37 _{0.6}	84.15 _{0.2}	64.79 _{0.5}	63.91 _{0.4}	88.74_{0.4}	54.89 _{0.4}	60.21 _{0.4}	56.86 _{0.5}	34.45 _{0.4}	30.39 _{1.2}	66.05_{0.5}
FocalNet	82.03 _{0.2}	86.23 _{0.1}	64.18 _{0.2}	38.94 _{0.8}	37.96 _{0.7}	87.90_{0.1}	47.76 _{0.3}	52.71 _{0.5}	32.33 _{0.3}	9.98 _{0.6}	9.18 _{0.1}	53.69_{0.3}
CCT-2/3×2	78.76 _{0.3}	83.89 _{0.4}	—	33.08 _{2.3}	35.16 _{0.4}	84.94_{0.3}	54.05 _{0.4}	59.78 _{0.5}	—	7.17 _{0.2}	8.60 _{0.1}	62.91_{0.5}

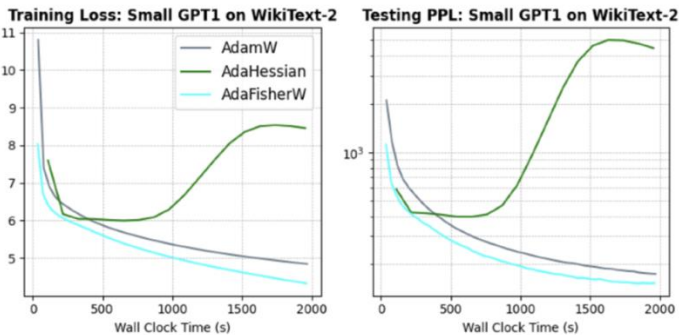
*Note that Adam and AdaFisher were used for all CNN architectures, while AdamW and AdaFisherW were applied for all ViT experiments.

Table 4: Performance comparison of different networks and optimizers on CIFAR10 and CIFAR100 using ImageNet-1K pretrained weights. Evaluation is based on wall clock time of 50 training epochs with AdaFisher.

Network	CIFAR10						CIFAR100					
	SGD	Adam	AdaHessian	K-FAC	Shampoo	AdaFisher	SGD	Adam	AdaHessian	K-FAC	Shampoo	AdaFisher
ResNet50	96.50 _{0.2}	96.45 _{0.2}	96.35 _{0.3}	96.45 _{0.1}	96.03 _{0.4}	97.13_{0.2}	82.12 _{0.1}	82.01 _{0.4}	80.64 _{0.9}	80.55 _{0.4}	81.70 _{0.2}	82.23_{0.2}
ResNet101	97.07 _{0.2}	96.70 _{0.1}	96.65 _{0.2}	96.84 _{0.1}	96.63 _{0.1}	97.22_{0.1}	84.01 _{0.1}	82.43 _{0.2}	81.36 _{0.8}	82.26 _{0.3}	82.65 _{0.2}	84.47_{0.2}
DenseNet121	94.80 _{0.1}	94.77 _{0.1}	93.08 _{0.1}	94.41 _{0.2}	94.76 _{0.1}	95.03_{0.1}	75.98 _{0.2}	75.65 _{0.3}	71.06 _{0.9}	76.10 _{0.3}	76.08 _{0.2}	76.92_{0.3}
MobileNetV3	91.76 _{0.3}	90.92 _{0.3}	86.45 _{2.5}	91.72 _{0.2}	91.39 _{0.3}	92.78_{0.2}	71.86 _{0.4}	66.11 _{0.8}	59.69 _{2.3}	69.85 _{0.4}	68.87 _{0.3}	72.38_{0.4}

Table 5: Language Modeling performance (PPL) on Wikitest-2 and PTB test dataset (lower is better)

Optimizer	Test PPL	
	WikiText-2	PTB
AdamW	175.06	44.70
AdaHessian	407.69	59.43
Shampoo	1727.75	—
AdaFisherW	152.72	41.15



Language Modeling

Stability Analysis

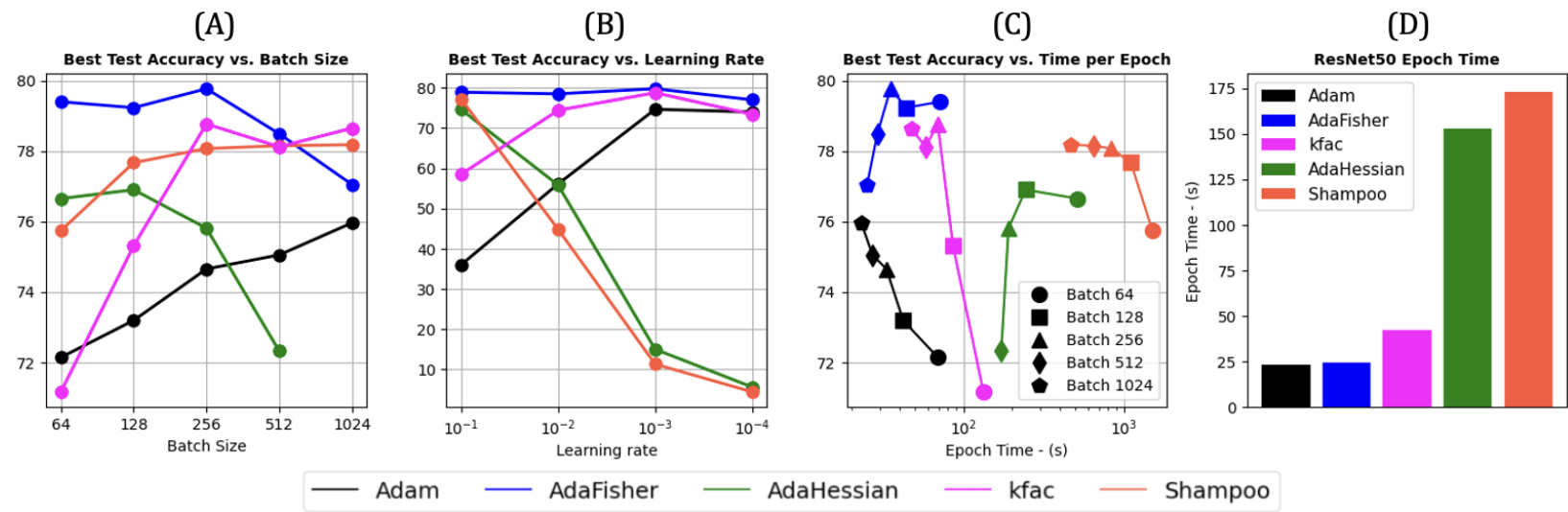


Figure 6: Performance comparison of AdaFisher and other optimizers using the ResNet50 network on the CIFAR100 dataset. (A) Test accuracy by batch size. (B) Accuracy vs. learning rates. (C) Accuracy related to epoch time across batch sizes. (D) Epoch time for different optimizers with a batch size of 256.

Code

Paper

AdaFisher: Adaptive Second Order Optimization vis Fisher

- ✓ Convergence, Generalization Capabilities
- ✓ Computational Efficiency

