
Fishy: Layerwise Fisher Approximation for Higher-order Neural Network Optimization

Abel L. Peirson*

Google Brain & Stanford University
alpv95@stanford.edu

Ehsan Amid*

Google Brain
eamid@google.com

Yatong Chen

Google Brain & UC Santa Cruz
ychen592@ucsc.edu

Vlad Feinberg

Google Brain
vladf@google.com

Manfred Warmuth

Google Research
manfred@google.com

Rohan Anil

Google Brain
rohananil@google.com

Abstract

We introduce *Fishy*, a local approximation of the Fisher information matrix at each layer for natural gradient descent training of deep neural networks. The true Fisher approximation for deep networks involves sampling labels from the model’s predictive distribution at the output layer and performing a full backward pass – *Fishy* defines a Bregman exponential family distribution at each layer, performing the sampling locally. Local sampling allows for model parallelism when forming the preconditioner, removing the need for the extra backward pass. We demonstrate our approach through the Shampoo optimizer, replacing its preconditioner gradients with our locally sampled gradients. Our training results on deep autoencoder and VGG16 image classification models indicate the efficacy of our construction.

1 Introduction

Natural gradient descent (NGD) [1] is a second-order update rule that preconditions the gradient direction with the inverse of the model’s Fisher information matrix (FIM). NGD corresponds to the steepest descent direction in the Riemannian space associated with the FIM and is equivariant to any differentiable reparameterization of the model weights. However, this property only holds in the small step limit and whether the benefits extend to finite steps remains an open question [16]. Nevertheless, NGD has proven to be remarkably effective in training deep neural networks, and many standard optimizers in deep learning can be seen as approximations of NGD [17, 11].

NGD is challenging to implement for training deep neural networks. Firstly, the FIM requires samples from the model’s predictive distribution, with an additional backward pass to calculate the gradients. Secondly, NGD requires calculating the inverse FIM of the whole network, which immediately becomes infeasible even for medium-sized models. K-FAC [17] uses a block diagonal Kronecker factor approximation to address the latter obstacle, considering the FIM for each layer as a Kronecker product of two matrices which can more easily be stored and inverted. However, the Kronecker approximation used in K-FAC becomes challenging to generalize for layers other than fully-connected layers. Moreover, the sampling and additional backward passes keep K-FAC expensive. The Shampoo optimizer [8, 5] approximates the full-matrix AdaGrad preconditioner using the same minibatch gradients of the training examples. Shampoo preconditioning avoids additional sampling or extra backward passes but is no longer related to true Fisher [14, 16] and, as we show experimentally, is outperformed when replaced with a true Fisher approximation. Using a Shampoo approximation of the true Fisher is related to [19]; primary differences are the choice of exponents and

*Equal Contribution

dampening term used for per-dimension preconditioners. Our work follows Shampoo’s Kronecker product approximation of the preconditioner, which has the advantage of being agnostic to the type of layers. Sun and Nielsen [22] recently proposed a layerwise relative FIM, but do not consider the full form of layers’ predictive distributions nor demonstrate scaling to practical architectures.

This paper presents a new layerwise FIM approximation approach, extending the popular Shampoo optimizer to NGD. First, we update Shampoo to an approximate NGD method, sampling from the output layer distribution to construct the preconditioners. Second, inspired by the matching loss [9, 4] construction for LocoProp and layerwise representation learning [3, 2], we define Bregman exponential family distributions for the outputs of each layer based on their activation functions. Sampling locally from these distributions allows preconditioner calculation to be performed in parallel across layers. Third, we compare our approach to standard Shampoo and its Fisher counterpart on a benchmark MNIST deep autoencoder problem and CIFAR100 classification with a VGG16 [20].

2 Methodology

Fisher Information Matrix and Natural Gradient Descent Given a probabilistic model $P(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})$ with input random variable \mathbf{x} and target random variable \mathbf{y} , the Fisher Information Metric (FIM) is defined in terms of a local Riemannian metric as an approximation to the KL divergence [18],

$$D_{\text{KL}}(P(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}), P(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta} + d\boldsymbol{\theta})) = \int_{\mathbf{x}, \mathbf{y}} p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})}{p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta} + d\boldsymbol{\theta})} d\mathbf{x} d\mathbf{y} \\ \approx 1/2 d\boldsymbol{\theta}^\top \underbrace{\mathbb{E}_{\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}}[\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})^\top]}_{:= \mathbf{F}(\boldsymbol{\theta}) \text{ Fisher Information matrix}} d\boldsymbol{\theta},$$

where the expectation is with respect to the model distribution $P(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})$. The result is identical when flipping the order of the arguments in the KL divergence. The FIM corresponds to the Riemannian metric that locally explains the model in the parameter space [15]. This probabilistic model approximates an underlying joint data distribution $P_d(\mathbf{x}, \mathbf{y})$ from which the data is sampled. Thus, training the model corresponds to finding optimal parameters $\boldsymbol{\theta}^*$ that minimize the KL divergence between the data and model distributions. For a discriminative model that generates $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$, the marginal distribution $P(\mathbf{x}|\boldsymbol{\theta})$ is assumed to be the same as the underlying input data distribution $P_d(\mathbf{x})$. Thus, the loss can be written as

$$L(\boldsymbol{\theta}) = \int_{\mathbf{x}, \mathbf{y}} p_d(\mathbf{x}, \mathbf{y}) \log \frac{p_d(\mathbf{x}, \mathbf{y})}{p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})} d\mathbf{x} d\mathbf{y} = \underbrace{\mathbb{E}_{\mathbf{x}, \mathbf{y}}[\log p_d(\mathbf{y}|\mathbf{x})]}_{\text{constant}} - \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})], \quad (1)$$

where the expectations are applied with respect to $P_d(\mathbf{x}, \mathbf{y})$. In practice, the empirical loss is formed by a Monte Carlo approximation of the expectation,

$$L(\boldsymbol{\theta}|\mathcal{X}) = -\frac{1}{n} \sum_i \log P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}),$$

using training examples $\mathcal{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$ that are sampled from $P_d(\mathbf{x}, \mathbf{y})$.

Given the loss $L(\boldsymbol{\theta}|\mathcal{X})$ on the batch of examples \mathcal{X} , the *natural gradient* step [1] is defined as the minimizer of the following linearized objective,

$$d\boldsymbol{\theta}_{\text{NGD}} = \arg \min_{d\boldsymbol{\theta}} \{L(\boldsymbol{\theta}|\mathcal{X}) + d\boldsymbol{\theta}^\top \nabla L(\boldsymbol{\theta}|\mathcal{X}) + 1/2\gamma d\boldsymbol{\theta}^\top \mathbf{F}(\boldsymbol{\theta}) d\boldsymbol{\theta}\} \\ = -\gamma \mathbf{F}(\boldsymbol{\theta})^{-1} \nabla L(\boldsymbol{\theta}|\mathcal{X}), \quad (2)$$

in which the loss at $L(\boldsymbol{\theta} + d\boldsymbol{\theta}|\mathcal{X})$ is approximated by its first-order Taylor expansion. The non-negative multiplier $\gamma > 0$ in the regularizer term controls the step size. Specifically, Equation (2) corresponds to minimizing a linear approximation of the loss at the current parameter $\boldsymbol{\theta}$ while the amount of deviation is measured in terms of the local Mahalanobis distance induced by the FIM.

Fishy: Layerwise Fisher Construction Let $\hat{\mathbf{a}}_\ell$ denote the pre-activation at layer $\ell \in [L]$ of an L -layer network. For a fully-connected layer with weights \mathbf{W}_ℓ , we can write $\hat{\mathbf{a}}_\ell = \mathbf{W}_\ell \hat{\mathbf{y}}_{\ell-1}$ and

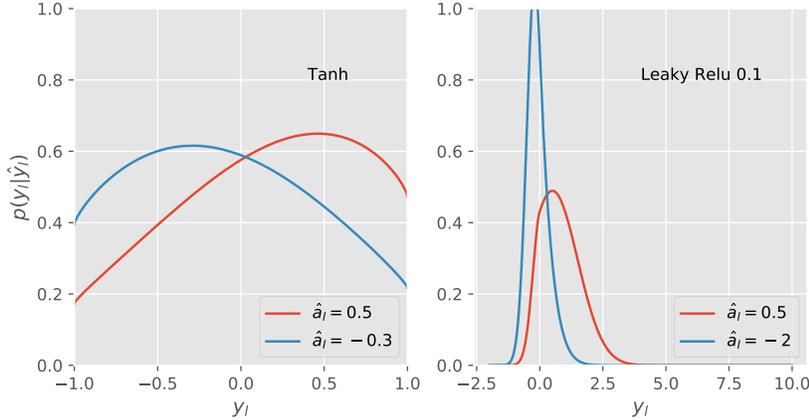


Figure 1: Example of layerwise predictive distributions for tanh (left) and leaky ReLU (right) activation functions.

the post-activation $\hat{\mathbf{y}}_\ell = f_\ell(\hat{\mathbf{a}}_\ell)$. Here, f_ℓ denotes the (elementwise) *strictly-increasing* activation function in layer $\ell \in [L]$. (Note that using this notation, $\hat{\mathbf{y}}_0 = \mathbf{x}$.) Fishy aims to approximate the Fisher \mathbf{F} for each layer *independently* by defining a predictive distribution for \mathbf{y}_ℓ , the layer’s output random variable, as

$$p_\ell(\mathbf{y}_\ell|\hat{\mathbf{y}}_\ell) = \frac{\exp(-D_{F_\ell^*}(\mathbf{y}_\ell, \hat{\mathbf{y}}_\ell))}{Z(\hat{\mathbf{y}}_\ell)}, \quad (3)$$

where $\hat{\mathbf{y}}_\ell$ is the post-activation of the layer given the input $\hat{\mathbf{y}}_{\ell-1}$. The term $Z(\hat{\mathbf{y}}_\ell)$ is the normalizing partition function that ensures the distribution sums to one. The distribution in Equation (3) belongs to a Bregman exponential family [6] in which the Bregman divergence is induced by the convex function F_ℓ^* . We apply the layerwise loss construction in [3] to obtain $D_{F_\ell^*}$, which corresponds to the matching loss [9, 13, 4] of the activation function f_ℓ induced by its integral function. Intuitively, $D_{F_\ell^*}$ encodes the layer’s local geometry; thus, correcting for the gradient direction using a local Fisher approximation can stabilize training and improve convergence [22]. For the linear, sigmoid, and softmax activations, Equation (3) yields the standard Gaussian, Bernoulli, and categorical distributions, respectively. Figure 1 gives some example distributions for tanh and leaky ReLU activations.

To calculate the layer’s FIM, we need to calculate the gradient of the *score* function $\log p_\ell(\mathbf{y}_\ell|\hat{\mathbf{y}}_\ell)$. Taking the gradient with respect to the pre-activation $\hat{\mathbf{a}}_\ell$, we have

$$\nabla_{\hat{\mathbf{a}}_\ell} \log p_\ell(\mathbf{y}_\ell|\hat{\mathbf{y}}_\ell) = \mathbf{y}_\ell - \hat{\mathbf{y}}_\ell - \nabla_{\mathbf{a}_\ell} \log Z(\boldsymbol{\theta}) = \mathbf{y}_\ell - \hat{\mathbf{y}}_\ell - \mathbb{E}_{P_\ell(\mathbf{y}_\ell|\hat{\mathbf{y}}_\ell)} [\mathbf{y}_\ell - \hat{\mathbf{y}}_\ell]. \quad (4)$$

The last term in Equation (4) is non-zero in general for asymmetric distributions induced by activation functions such as tanh and leaky ReLU (see Figure 1). Nevertheless, in our approximation, we discard the last term in Equation (4) to form an upper bound of the local FIM approximation. For a fully-connected layer with $\boldsymbol{\theta}_\ell = \text{vec}(\mathbf{W}_\ell)$, we can write

$$\mathbf{F}(\boldsymbol{\theta}_\ell) \preceq \mathbb{E}_{P(\mathbf{x}, \mathbf{y}_\ell)} [((\mathbf{y}_\ell - \hat{\mathbf{y}}_\ell) \otimes \hat{\mathbf{y}}_{\ell-1})((\mathbf{y}_\ell - \hat{\mathbf{y}}_\ell) \otimes \hat{\mathbf{y}}_{\ell-1})^T], \quad (5)$$

where $P(\mathbf{x}, \mathbf{y}_\ell) = P_d(\mathbf{x}) P_\ell(\mathbf{y}_\ell|\hat{\mathbf{y}}_\ell(\mathbf{x}|\boldsymbol{\theta}_\ell))$. The FIM approximation in Equation (5) corresponds to a block-diagonal approximation of the full FIM as the layers are treated independently. This is a common procedure for calculating FIM for large neural networks [17, 22]. With Fishy, we have the flexibility of sampling at every layer or dividing the network into multiple disjoint sections and performing the sampling at the output layer of each section. The sampled gradients are then backpropagated through the corresponding section. In the extreme case where the entire network is treated as a single section, Equation (5) reduces to a block diagonal approximation of the true Fisher by sampling from the model’s output predictive distribution [17]. Next, we discuss using the Shampoo approximation to apply the inverse of Equation (5) as the preconditioner at each layer.

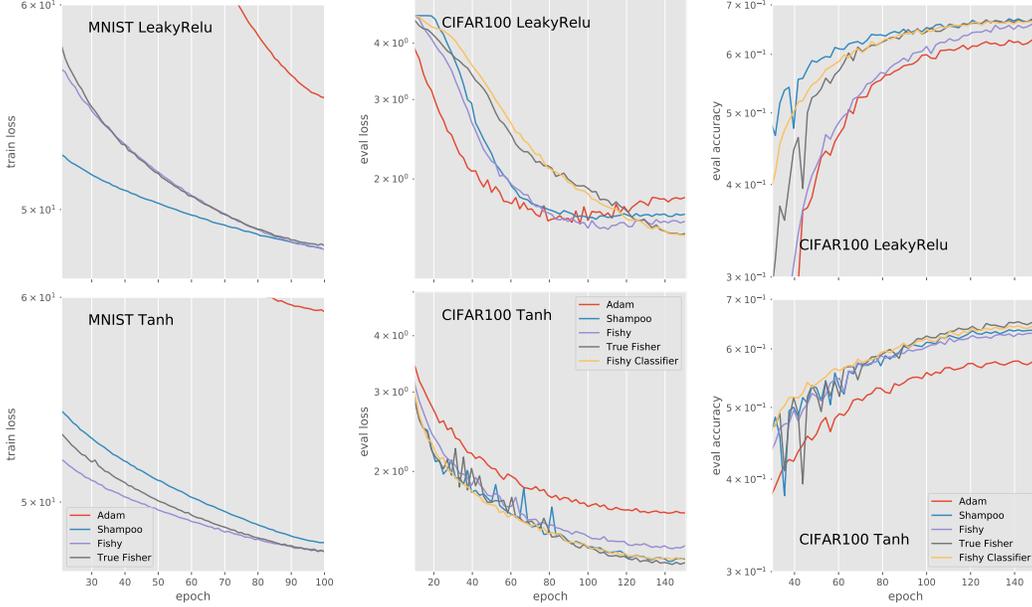


Figure 2: Experimental results on deep autoencoder for MNIST (left column) and VGG16 convolutional network on CIFAR100 for image classification (middle and right columns).

Fishy via Shampoo Approximation For a fully-connected layer with gradient $\mathbf{G} \in \mathbb{R}^{d_i \times d_o}$, Shampoo approximates the outer-product of the vectorized gradients $\mathbf{g} = \text{vec}(\mathbf{G}) \in \mathbb{R}^{d_i d_o}$ as a Kronecker product of two factors

$$\varepsilon \mathbf{I}_{d_i d_o} + \frac{1}{r} \mathbf{g} \mathbf{g}^\top \preccurlyeq \mathbf{L}^{\frac{1}{2}} \otimes \mathbf{R}^{\frac{1}{2}} = (\varepsilon \mathbf{I}_{d_i} + \mathbf{G} \mathbf{G}^\top)^{\frac{1}{2}} \otimes (\varepsilon \mathbf{I}_{d_o} + \mathbf{G}^\top \mathbf{G})^{\frac{1}{2}}, \quad (6)$$

where r is an upper-bound on the rank of the gradients (see [8] for further details), and \mathbf{L} and \mathbf{R} correspond to the left and right preconditioners. The construction in Equation (6) is more general and can be applied to many types of layers. However, notice that \mathbf{g} in Equation (6) corresponds to the average gradient of a batch (and not the per example gradient).

To approximate Equation (5) via Shampoo, we first sample $\mathbf{y}_\ell^{(i)}$ from the predictive distribution in Equation (3) given the input sample $\mathbf{x}^{(i)}$ in the batch $\{\mathbf{x}^{(i)}\}_{i=1}^n$ and the corresponding post-activations $\hat{\mathbf{y}}_{\ell-1}$ and $\hat{\mathbf{y}}_\ell$ at layer $\ell-1$ and ℓ , respectively. Next, instead of the per example sampled gradient $\mathbf{g}_s^{(i)} = (\mathbf{y}_\ell^{(i)} - \hat{\mathbf{y}}_\ell^{(i)}) \otimes \hat{\mathbf{y}}_{\ell-1}^{(i)}$, we pass the average sampled gradient $\mathbf{g}_s = 1/n \sum_i \mathbf{g}_s^{(i)}$ into Shampoo to form the preconditioners. A similar averaging is performed in K-FAC [17] on the sampled gradients before forming the preconditioners. Shampoo then applies the (EMA of the) left \mathbf{L}_s and right \mathbf{R}_s Fishy preconditioners on the gradient of the training examples \mathbf{G} ,

$$\mathbf{G}_{\text{Fishy-NGD}} = \mathbf{L}_s^{-\frac{1}{2}} \mathbf{G} \mathbf{R}_s^{-\frac{1}{2}}. \quad (7)$$

Notice that the exponents in Equation (7) are different than the original Shampoo formulation ($-1/4$) as we are interested in the inverse of the local FIM (and not its inverse square root). In practice, the exponent in Shampoo is treated as a tunable hyperparameter.

3 Experiments

We perform experiments on two architectures, and datasets, a deep MNIST autoencoder [10] with fully connected hidden layers [1000, 500×8 , 250, 30, 250, 500×8 , 1000], a standard benchmark for second-order methods, and a VGG16 CNN architecture [20] applied to CIFAR100. We compare Adam [12], standard Shampoo [5], FIM Shampoo (true Fisher), and our method Fishy. For VGG16, we include a variant of Fishy that only approximates the local FIM for the final fully-connected

layers (called Fishy classifier), letting earlier convolutional layers inherit gradients from the first fully-connected layer’s predictive distribution. Hyperparameters for each method are tuned over 400 trials with the Vizier Bayesian optimization toolbox [21, 7], and the trial with the best final objective value is displayed. All Shampoo and Fishy variants share the same tunable hyperparameter ranges. We find significant improvements over standard Shampoo with Fishy or its variants on most of the tasks. Fishy even outperforms FIM Shampoo on some problems. FIM Shampoo can be considered Fishy applied to the output layer only. We plan to open-source the code to reproduce these experiments.

4 Conclusion

We introduce Fishy, a layerwise Fisher construction method that allows model parallelism for approximating the FIM. We apply our construction to Shampoo, which only requires simple adjustments but can improve performance on several problems. In the future, we would like to perform more extensive experiments comparing K-FAC, derive a full expectation version of Fishy without requiring sampling, and develop Fishy for the K-FAC approximation. Further treatment of layerwise predictive distributions on their own can also yield interesting research directions.

References

- [1] Shun-ichi Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):251–276, February 1998.
- [2] Ehsan Amid, Rohan Anil, Christopher Fifty, and Manfred K. Warmuth. Layerwise Bregman Representation Learning with Applications to Knowledge Distillation, September 2022. arXiv:2209.07080 [cs].
- [3] Ehsan Amid, Rohan Anil, and Manfred Warmuth. LocoProp: Enhancing BackProp via Local Loss Optimization. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, pages 9626–9642. PMLR, May 2022. ISSN: 2640-3498.
- [4] Ehsan Amid, Manfred K. K Warmuth, Rohan Anil, and Tomer Koren. Robust Bi-Tempered Logistic Loss Based on Bregman Divergences. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [5] Rohan Anil, Vineet Gupta, Tomer Koren, Kevin Regan, and Yoram Singer. Scalable second order optimization for deep learning. *arXiv preprint arXiv:2002.09018*, 2020.
- [6] Arindam Banerjee, Srujana Merugu, Inderjit S Dhillon, and Joydeep Ghosh. Clustering with Bregman divergences. *Journal of machine learning research*, 6(Oct):1705–1749, 2005.
- [7] Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Karro, and D. Sculley. Google vizier: A service for black-box optimization. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*, pages 1487–1495. ACM, 2017.
- [8] Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned Stochastic Tensor Optimization. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1842–1850. PMLR, July 2018.
- [9] D.P. Helmbold, J. Kivinen, and M.K. Warmuth. Relative loss bounds for single neurons. *IEEE Transactions on Neural Networks*, 10(6):1291–1304, November 1999.
- [10] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.
- [11] Mohammad Emtiyaz Khan and Håvard Rue. The Bayesian learning rule. *arXiv preprint arXiv:2107.04562*, 2021.
- [12] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

- [13] J. Kivinen and M. K. Warmuth. Relative Loss Bounds for Multidimensional Regression Problems. *Machine Learning*, 45(3):301–329, December 2001.
- [14] Frederik Kunstner, Philipp Hennig, and Lukas Balles. Limitations of the empirical Fisher approximation for natural gradient descent. *Advances in neural information processing systems*, 32, 2019.
- [15] John M Lee. *Riemannian manifolds: an introduction to curvature*, volume 176. Springer Science & Business Media, 2006.
- [16] James Martens. New Insights and Perspectives on the Natural Gradient Method. *Journal of Machine Learning Research*, 21(146):1–76, 2020.
- [17] James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15*, pages 2408–2417, Lille, France, July 2015. JMLR.org.
- [18] Frank Nielsen. An elementary introduction to information geometry. *Entropy*, 22(10):1100, 2020.
- [19] Yi Ren and Donald Goldfarb. Tensor normal training for deep learning models. *Advances in Neural Information Processing Systems*, 34:26040–26052, 2021.
- [20] Karen Simonyan and Andrew Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition, April 2015. arXiv:1409.1556 [cs].
- [21] Xingyou Song, Sagi Perel, Chansoo Lee, Greg Kochanski, and Daniel Golovin. Open source vizier: Distributed infrastructure and API for reliable and flexible blackbox optimization. In *Automated Machine Learning Conference, Systems Track (AutoML-Conf Systems)*, 2022.
- [22] Ke Sun and Frank Nielsen. Relative Fisher Information and Natural Gradient for Learning Large Modular Models. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3289–3298. PMLR, July 2017.