

---

# ISAAC Newton: Input-based Approximate Curvature for Newton’s Method

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We present ISAAC (Input-baSed ApproximAte Curvature), a novel method that  
2 conditions the gradient using selected second-order information and has an asymp-  
3 totically vanishing computational overhead, assuming a batch size smaller than  
4 the number of neurons. We show that it is possible to compute a good conditioner  
5 based on only the input to a respective layer without a substantial computational  
6 overhead. The proposed method allows effective training even in small-batch  
7 stochastic regimes, which makes it competitive to first-order as well as quasi-  
8 Newton methods.

## 9 1 Introduction

10 While second-order optimization methods are traditionally much less explored than first-order  
11 methods in large-scale machine learning (ML) applications due to their memory requirements and  
12 prohibitive computational cost per iteration, they have recently become more popular in ML mainly  
13 due to their fast convergence properties when compared to first-order methods [1]. The expensive  
14 computation of an inverse Hessian (also known as pre-conditioning matrix) in the Newton step has  
15 also been tackled via estimating the curvature from the change in gradients. Loosely speaking, these  
16 algorithms are known as *quasi-Newton methods* and a comprehensive treatment can be found in  
17 the textbook [2]. In addition, various new approximations to the pre-conditioning matrix have been  
18 proposed in the recent literature [3]–[6]. From a theoretical perspective, second-order optimization  
19 methods are not nearly as well understood as first-order methods. It is an active research direction to  
20 fill this gap [7], [8].

21 Motivated by the task of training neural networks, and the observation that invoking local curvature  
22 information associated with neural network objective functions can achieve much faster progress  
23 per iteration than standard first-order methods [9]–[11], several methods have been proposed. One  
24 of these methods, that received significant attention, is known as *Kronecker-factored Approximate*  
25 *Curvature (K-FAC)* [12], whose main ingredient is a sophisticated approximation to the generalized  
26 Gauss-Newton matrix and the Fisher information matrix quantifying the curvature of the underlying  
27 neural network objective function, which then can be inverted efficiently.

28 Inspired by the K-FAC approximation and the Tikhonov regularization of the Newton method, we  
29 introduce a novel two parameter regularized Kronecker-factorized Newton update step. The proposed  
30 scheme disentangles the classical Tikhonov regularization and allows us to condition the gradient  
31 using selected second-order information and has an asymptotically vanishing computational overhead.  
32 While this property makes the presented method highly attractive from the computational complexity  
33 perspective, we show that its achieved empirical performance on complicated high-dimensional  
34 Machine Learning problems remains comparable to existing state-of-the-art methods.

35 The contributions of this paper can be summarized as follows: (i) we propose a novel two parameter  
36 regularized K-FAC approximated Gauss-Newton update step; (ii) we show that asymptotically—as

37 both regularization parameters vanish—our method recovers the classical K-FAC scheme and in  
 38 the opposite setting—as both regularization parameters grow—our method asymptotically reduces  
 39 to classical gradient descent; (iii) we prove that for an arbitrary pair of regularization parameters,  
 40 the proposed update direction is always a direction of decreasing loss; (iv) in the limit, as one  
 41 regularization parameter grows, we obtain an efficient and effective conditioning of the gradient with  
 42 an asymptotically vanishing overhead; (v) we empirically analyze the presented method and find that  
 43 our efficient conditioning method maintains the performance of its more expensive counterpart; (vi)  
 44 we demonstrate the effectiveness of the presented method in the setting of small-batch stochastic  
 45 regimes and observe that it is competitive to first-order as well as quasi-Newton methods.

## 46 2 Preliminaries

47 In this section, we review aspects of second-order optimization, with a focus on generalized Gauss-  
 48 Newton methods. In combination with Kronecker factorization, this leads us to a new regularized  
 49 update scheme. We consider the training of an  $L$ -layer neural network  $f(x; \theta)$  defined recursively as

$$z_i \leftarrow a_{i-1} W^{(i)} \quad (\text{pre-activations}), \quad a_i \leftarrow \phi(z_i) \quad (\text{activations}), \quad (1)$$

50 where  $a_0 = x$  is the vector of inputs and  $a_L = f(x; \theta)$  is the vector of outputs. Unless noted otherwise,  
 51 we assume these vectors to be row vectors (i.e., in  $\mathbb{R}^{1 \times n}$ ) as this allows for a direct extension to the  
 52 (batch) vectorized case (i.e., in  $\mathbb{R}^{b \times n}$ ) introduced later. For any layer  $i$ , let  $W^{(i)} \in \mathbb{R}^{d_{i-1} \times d_i}$  be a  
 53 weight matrix and let  $\phi$  be an element-wise nonlinear function. We consider a convex loss function  
 54  $\mathcal{L}(y, y')$  that measures the discrepancy between  $y$  and  $y'$ . The training optimization problem is then

$$\arg \min_{\theta} \mathbb{E}_{x,y} [\mathcal{L}(f(x; \theta), y)], \quad (2)$$

55 where  $\theta = [\theta^{(1)}, \dots, \theta^{(L)}]$  with  $\theta^{(i)} = \text{vec}(W^{(i)})$ .

56 The classical Newton method for solving (2) is expressed as the update rule

$$\theta' = \theta - \eta \mathbf{H}_{\theta}^{-1} \nabla_{\theta} \mathcal{L}(f(x; \theta), y), \quad (3)$$

57 where  $\eta > 0$  denotes the learning rate and  $\mathbf{H}_{\theta}$  is the Hessian corresponding to the objective function  
 58 in (2). The stability and efficiency of an estimation problem solved via the Newton method can be  
 59 improved by adding a Tikhonov regularization term [13] leading to a regularized Newton method

$$\theta' = \theta - \eta (\mathbf{H}_{\theta} + \lambda \mathbf{I})^{-1} \nabla_{\theta} \mathcal{L}(f(x; \theta), y), \quad (4)$$

60 where  $\lambda > 0$  is the so-called Tikhonov regularization parameter. It is well-known [14], [15], that  
 61 under the assumption of approximating the model  $f$  with its first-order Taylor expansion, the Hessian  
 62 corresponds with the so-called generalized Gauss-Newton (GGN) matrix  $\mathbf{G}_{\theta}$ , and hence (4) can be  
 63 expressed as

$$\theta' = \theta - \eta (\mathbf{G}_{\theta} + \lambda \mathbf{I})^{-1} \nabla_{\theta} \mathcal{L}(f(x; \theta), y). \quad (5)$$

64 A major practical limitation of (5) is the computation of the inverse term. A method that alleviates this  
 65 difficulty is known as Kronecker-Factored Approximate Curvature (K-FAC) [12] which approximates  
 66 the block-diagonal (i.e., layer-wise) empirical Hessian or GGN matrix. Inspired by K-FAC, there  
 67 have been other works discussing approximations of  $\mathbf{G}_{\theta}$  and its inverse [15]. In the following, we  
 68 discuss a popular approach that allows for (moderately) efficient computation.

69 The generalized Gauss-Newton matrix  $\mathbf{G}_{\theta}$  is defined as

$$\mathbf{G}_{\theta} = \mathbb{E} [(\mathbf{J}_{\theta} f(x; \theta))^{\top} \nabla_f^2 \mathcal{L}(f(x; \theta), y) \mathbf{J}_{\theta} f(x; \theta)], \quad (6)$$

70 where  $\mathbf{J}$  and  $\mathbf{H}$  denote the Jacobian and Hessian matrices, respectively. Correspondingly, the diagonal  
 71 block of  $\mathbf{G}_{\theta}$  corresponding to the weights of the  $i$ th layer  $W^{(i)}$  is

$$\mathbf{G}_{W^{(i)}} = \mathbb{E} [(\mathbf{J}_{W^{(i)}} f(x; \theta))^{\top} \nabla_f^2 \mathcal{L}(f(x; \theta), y) \mathbf{J}_{W^{(i)}} f(x; \theta)].$$

72 According to the backpropagation rule  $\mathbf{J}_{\theta^{(i)}} f(x; \theta) = \mathbf{J}_{z_i} f(x; \theta) a_{i-1}$ ,  $a^{\top} b = a \otimes b$ , and the  
 73 mixed-product property, we can rewrite  $\mathbf{G}_{W^{(i)}}$  as

$$\mathbf{G}_{W^{(i)}} = \mathbb{E} [((\mathbf{J}_{z_i} f(x; \theta) a_{i-1})^{\top} (\nabla_f^2 \mathcal{L}(f(x; \theta), y))^{1/2}) ((\nabla_f^2 \mathcal{L}(f(x; \theta), y))^{1/2} \mathbf{J}_{z_i} f(x; \theta) a_{i-1})] \quad (7)$$

$$= \mathbb{E} [(\bar{g}^{\top} a_{i-1})^{\top} (\bar{g}^{\top} a_{i-1})] = \mathbb{E} [(\bar{g} \otimes a_{i-1})^{\top} (\bar{g} \otimes a_{i-1})] = \mathbb{E} [(\bar{g}^{\top} \bar{g}) \otimes (a_{i-1}^{\top} \otimes a_{i-1})], \quad (8)$$

74 where

$$\bar{g} = (\mathbf{J}_{z_i} f(x; \theta))^{\top} (\nabla_f^2 \mathcal{L}(f(x; \theta), y))^{1/2}. \quad (9)$$

75 **Remark 1** (Monte-Carlo Low-Rank Approximation for  $\bar{g}^\top \bar{g}$ ). As  $\bar{g}$  is a matrix of shape  $m \times d_i$   
76 where  $m$  is the dimension of the output of  $f$ ,  $\bar{g}$  is generally expensive to compute. Therefore, [12] use  
77 a low-rank Monte-Carlo approximation to estimate  $\mathbf{H}_f \mathcal{L}(f(x; \theta), y)$  and thereby  $\bar{g}^\top \bar{g}$ . For this, we  
78 need to use the distribution underlying the probabilistic model of our loss  $\mathcal{L}$  (e.g., Gaussian for MSE  
79 loss, or a categorical distribution for cross entropy). Specifically, by sampling from this distribution  
80  $p_f(x)$  defined by the network output  $f(x; \theta)$ , we can get an estimator of  $\mathbf{H}_f \mathcal{L}(f(x; \theta), y)$  via the  
81 identity

$$\mathbf{H}_f \mathcal{L}(f(x; \theta), y) = \mathbb{E}_{\hat{y} \sim p_f(x)} [\nabla_f \mathcal{L}(f(x; \theta), \hat{y})^\top \nabla_f \mathcal{L}(f(x; \theta), \hat{y})]. \quad (10)$$

82 An extensive reference for this (as well as alternatives) can be found in Appendix A.2 of Dangel et  
83 al. [15]. The respective rank-1 approximation (denoted by  $\triangleq$ ) of  $\mathbf{H}_f \mathcal{L}(f(x; \theta))$  is

$$\mathbf{H}_f \mathcal{L}(f(x; \theta), y) \triangleq \nabla_f \mathcal{L}(f(x; \theta), \hat{y})^\top \nabla_f \mathcal{L}(f(x; \theta), \hat{y}),$$

84 where  $\hat{y} \sim p_f(x)$ . Respectively, we can estimate  $\bar{g}^\top \bar{g}$  using this rank-1 approximation with

$$\bar{g} \triangleq (\mathbf{J}_{z_i} f(x; \theta))^\top \nabla_f \mathcal{L}(f(x; \theta), \hat{y}) = \nabla_{z_i} \mathcal{L}(f(x; \theta), \hat{y}). \quad (11)$$

85 In analogy to  $\bar{g}$ , we introduce the gradient of training objective with respect to pre-activations  $z_i$  as

$$\mathbf{g}_i = (\mathbf{J}_{z_i} f(x; \theta))^\top \nabla_f \mathcal{L}(f(x; \theta), y) = \nabla_{z_i} \mathcal{L}(f(x; \theta), y). \quad (12)$$

86 In other words, for a given layer, let  $\mathbf{g} \in \mathbb{R}^{1 \times d_i}$  denote the gradient of the loss between an output and  
87 the ground truth and let  $\bar{g} \in \mathbb{R}^{m \times d_i}$  denote the derivative of the network  $f$  times the square root of  
88 the Hessian of the loss function (which may be approximated according to Remark 1), each of them  
89 with respect to the output  $z_i$  of the given layer  $i$ . Note that  $\bar{g}$  is not equal to  $\mathbf{g}$  and that they require one  
90 backpropagation pass each (or potentially many for the case of  $\bar{g}$ ). This makes computing  $\bar{g}$  costly.

91 Applying the K-FAC [12] approximation to (8) the expectation of Kronecker products can be  
92 approximated as the Kronecker product of expectations as

$$\mathbf{G} = \mathbb{E}((\bar{g}^\top \bar{g}) \otimes (\mathbf{a}^\top \mathbf{a})) \approx \mathbb{E}(\bar{g}^\top \bar{g}) \otimes \mathbb{E}(\mathbf{a}^\top \mathbf{a}), \quad (13)$$

93 where, for clarity, we drop the index of  $\mathbf{a}_{i-1}$  in (8) and denote it with  $\mathbf{a}$ ; similarly we denote  $\mathbf{G}_{W^{(i)}}$   
94 as  $\mathbf{G}$ . While the expectation of Kronecker products is generally not equal to the Kronecker product  
95 of expectations, this K-FAC approximation (13) has been shown to be fairly accurate in practice  
96 and to preserve the ‘‘coarse structure’’ of the GGN matrix [12]. The K-FAC decomposition in (13)  
97 is convenient as the Kronecker product has the favorable property that for two matrices  $A, B$  the  
98 identity  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$  which significantly simplifies the computation of an inverse.

99 In practice,  $\mathbb{E}(\bar{g}^\top \bar{g})$  and  $\mathbb{E}(\mathbf{a}^\top \mathbf{a})$  can be computed by averaging over a batch of size  $b$  as

$$\mathbb{E}(\bar{g}^\top \bar{g}) \simeq \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b, \quad \mathbb{E}(\mathbf{a}^\top \mathbf{a}) \simeq \mathbf{a}^\top \mathbf{a}/b, \quad (14)$$

100 where we denote batches of  $\mathbf{g}$ ,  $\bar{g}$  and  $\mathbf{a}$ , as  $\mathbf{g} \in \mathbb{R}^{b \times d_i}$ ,  $\bar{\mathbf{g}} \in \mathbb{R}^{r b \times d_i}$  and  $\mathbf{a} \in \mathbb{R}^{b \times d_{i-1}}$ , where our layer  
101 has  $d_{i-1}$  inputs,  $d_i$  outputs,  $b$  is the batch size, and  $r$  is either the number of outputs  $m$  or the rank of  
102 an approximation according to Remark 1. Correspondingly, the K-FAC approximation of the GGN  
103 matrix and its inverse are concisely expressed as

$$\mathbf{G} \approx (\bar{\mathbf{g}}^\top \bar{\mathbf{g}}) \otimes (\mathbf{a}^\top \mathbf{a})/b^2 \quad \mathbf{G}^{-1} \approx (\bar{\mathbf{g}}^\top \bar{\mathbf{g}})^{-1} \otimes (\mathbf{a}^\top \mathbf{a})^{-1} \cdot b^2. \quad (15)$$

104 Equipped with the standard terminology and setting, we now introduce the novel, regularized update  
105 step. First, inspired by the K-FAC approximation (13), the Tikhonov regularized Gauss-Newton  
106 method (5) can be approximated by

$$\theta^{(i)'} = \theta^{(i)} - \eta (\bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b + \lambda \mathbf{I})^{-1} \otimes (\mathbf{a}^\top \mathbf{a}/b + \lambda \mathbf{I})^{-1} \nabla_{\theta^{(i)}} \mathcal{L}(f(x; \theta)), \quad (16)$$

107 with regularization parameter  $\lambda > 0$ . A key observation, which is motivated by the structure of  
108 the above update, is to disentangle the two occurrences of  $\lambda$  into two independent regularization  
109 parameters  $\lambda_{\mathbf{g}}, \lambda_{\mathbf{a}} > 0$ . By defining the Kronecker-factorized Gauss-Newton update step as

$$\boldsymbol{\zeta} = \lambda_{\mathbf{g}} \lambda_{\mathbf{a}} (\bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b + \lambda_{\mathbf{g}} \mathbf{I})^{-1} \otimes (\mathbf{a}^\top \mathbf{a}/b + \lambda_{\mathbf{a}} \mathbf{I})^{-1} \nabla_{\theta^{(i)}} \mathcal{L}(f(x; \theta)), \quad (17)$$

110 we obtain the concise update equation  $\theta^{(i)'} = \theta^{(i)} - \eta^* \boldsymbol{\zeta}$ . (18)

111 This update (18) is equivalent to update (16) when in the case of  $\eta^* = \frac{\eta}{\lambda_{\mathbf{g}}\lambda_{\mathbf{a}}}$  and  $\lambda = \lambda_{\mathbf{g}} = \lambda_{\mathbf{a}}$ . This  
 112 equivalence does not restrict  $\eta^*$ ,  $\lambda_{\mathbf{g}}$ ,  $\lambda_{\mathbf{a}}$  in any way, and changing  $\lambda_{\mathbf{g}}$  or  $\lambda_{\mathbf{a}}$  does not mean that we  
 113 change our learning rate or step size  $\eta^*$ . Parameterizing  $\zeta$  in (17) with the multiplicative terms  $\lambda_{\mathbf{g}}\lambda_{\mathbf{a}}$   
 114 makes the formulation more convenient for analysis.

115 In this paper, we investigate the theoretical and empirical properties of the iterative update rule (18)  
 116 and in particular show how the regularization parameters  $\lambda_{\mathbf{g}}$ ,  $\lambda_{\mathbf{a}}$  affect the Kronecker-factorized  
 117 Gauss-Newton update step  $\zeta$ . When analyzing the Kronecker-factorized Gauss-Newton update step  
 118  $\zeta$ , a particularly useful tool is the vector product identity,

$$\left( (\bar{\mathbf{g}}^\top \bar{\mathbf{g}})^{-1} \otimes (\mathbf{a}^\top \mathbf{a})^{-1} \right) \text{vec}(\mathbf{g}^\top \mathbf{a}) = \text{vec} \left( (\bar{\mathbf{g}}^\top \bar{\mathbf{g}})^{-1} \mathbf{g}^\top \mathbf{a} (\mathbf{a}^\top \mathbf{a})^{-1} \right), \quad (19)$$

119 where the gradient with respect to the weight matrix is  $\mathbf{g}^\top \mathbf{a}$ .

### 120 3 Theoretical Guarantees

121 In this section, we investigate the theoretical properties of the Kronecker-factorized Gauss-Newton  
 122 update direction  $\zeta$  as defined in (17). We recall that  $\zeta$  introduces a Tikonov regularization, as it is  
 123 commonly done in implementations of second order-based methods. Not surprisingly, we show that  
 124 by decreasing the regularization parameters  $\lambda_{\mathbf{g}}$ ,  $\lambda_{\mathbf{a}}$  the update rule (18) collapses (in the limit) to the  
 125 classical Gauss-Newton method, and hence in the regime of small  $\lambda_{\mathbf{g}}$ ,  $\lambda_{\mathbf{a}}$  the variable  $\zeta$  describes the  
 126 Gauss-Newton direction. Moreover, by increasing the regularization strength, we converge (in the  
 127 limit) to the conventional gradient descent update step.

128 The key observation is that, as we disentangle the regularization of the two Kronecker factors  $\bar{\mathbf{g}}^\top \bar{\mathbf{g}}$   
 129 and  $\mathbf{a}^\top \mathbf{a}$ , and consider the setting where only one regularizer is large ( $\lambda_{\mathbf{g}} \rightarrow \infty$  to be precise),  
 130 we obtain an update direction that can be computed highly efficiently. We show that this setting  
 131 describes an approximated Gauss-Newton update scheme, whose superior numerical performance is  
 132 then empirically demonstrated in Section 4.

133 **Theorem 1** (Properties of  $\zeta$ ). *The K-FAC based update step  $\zeta$  as defined in (17) can be expressed as*

$$\zeta = \left( \mathbf{I}_m - \frac{1}{b\lambda_{\mathbf{g}}} \bar{\mathbf{g}}^\top \left( \mathbf{I}_b + \frac{1}{b\lambda_{\mathbf{g}}} \bar{\mathbf{g}} \bar{\mathbf{g}}^\top \right)^{-1} \bar{\mathbf{g}} \right) \cdot \mathbf{g}^\top \cdot \left( \mathbf{I}_b - \frac{1}{b\lambda_{\mathbf{a}}} \mathbf{a} \mathbf{a}^\top \left( \mathbf{I}_b + \frac{1}{b\lambda_{\mathbf{a}}} \mathbf{a} \mathbf{a}^\top \right)^{-1} \right) \cdot \mathbf{a}. \quad (20)$$

134 Moreover,  $\zeta$  admits the following asymptotic properties:

135 (i) *In the limit of  $\lambda_{\mathbf{g}}, \lambda_{\mathbf{a}} \rightarrow 0$ ,  $\frac{1}{\lambda_{\mathbf{g}}\lambda_{\mathbf{a}}}\zeta$  is the K-FAC approximation of the Gauss-Newton step, i.e.,*

$$136 \lim_{\lambda_{\mathbf{g}}, \lambda_{\mathbf{a}} \rightarrow 0} \frac{1}{\lambda_{\mathbf{g}}\lambda_{\mathbf{a}}}\zeta \approx \mathbf{G}^{-1} \nabla_{\theta^{(i)}} \mathcal{L}(f(x; \theta)), \text{ where } \approx \text{ denotes the K-FAC approximation (15).}$$

137 (ii) *In the limit of  $\lambda_{\mathbf{g}}, \lambda_{\mathbf{a}} \rightarrow \infty$ ,  $\zeta$  is the gradient, i.e.,  $\lim_{\lambda_{\mathbf{g}}, \lambda_{\mathbf{a}} \rightarrow \infty} \zeta = \nabla_{\theta^{(i)}} \mathcal{L}(f(x; \theta))$ .*

138 *The Proof is deferred to the Supplementary Material.*

139 We want to show that  $\zeta$  is well-defined and points in the correct direction, not only for  $\lambda_{\mathbf{g}}$  and  $\lambda_{\mathbf{a}}$   
 140 numerically close to zero because we want to explore the full spectrum of settings for  $\lambda_{\mathbf{g}}$  and  $\lambda_{\mathbf{a}}$ .  
 141 Thus, we prove that  $\zeta$  is a direction of increasing loss, independent of the choices of  $\lambda_{\mathbf{g}}$  and  $\lambda_{\mathbf{a}}$ .

142 **Theorem 2** (Correctness of  $\zeta$  is independent of  $\lambda_{\mathbf{g}}$  and  $\lambda_{\mathbf{a}}$ ).  *$\zeta$  is a direction of increasing loss,*  
 143 *independent of the choices of  $\lambda_{\mathbf{g}}$  and  $\lambda_{\mathbf{a}}$ .*

144 *Proof.* Recall that  $(\lambda_{\mathbf{g}}\mathbf{I}_m + \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b)$  and  $(\lambda_{\mathbf{a}}\mathbf{I}_n + \mathbf{a}^\top \mathbf{a}/b)$  are positive semi-definite (PSD) matrices by  
 145 definition. Their inverses  $(\lambda_{\mathbf{g}}\mathbf{I}_m + \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b)^{-1}$  and  $(\lambda_{\mathbf{a}}\mathbf{I}_n + \mathbf{a}^\top \mathbf{a}/b)^{-1}$  are therefore also PSD. As the  
 146 Kronecker product of PSD matrices is PSD, the conditioning matrix  $((\lambda_{\mathbf{g}}\mathbf{I}_m + \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b)^{-1} \otimes (\lambda_{\mathbf{a}}\mathbf{I}_n + \mathbf{a}^\top \mathbf{a}/b)^{-1} \approx \mathbf{G}^{-1})$  is PSD, and therefore the direction of the update step remains correct.  $\square$

148 From our formulation of  $\zeta$ , we can find that, in the limit for  $\lambda_{\mathbf{g}} \rightarrow \infty$ , Equation (21) does not depend  
 149 on  $\bar{\mathbf{g}}$ . This is computationally very beneficial as computing  $\bar{\mathbf{g}}$  is costly as it requires one or even  
 150 many additional backpropagation passes. In addition, it allows conditioning the gradient update by  
 151 multiplying a  $b \times b$  matrix between  $\mathbf{g}^\top$  and  $\mathbf{a}$ , which can be done very fast.

152 **Theorem 3** (Efficient Update Direction). *In the limit of  $\lambda_{\mathbf{g}} \rightarrow \infty$ , the update step  $\zeta$  converges to*  
 153  *$\lim_{\lambda_{\mathbf{g}} \rightarrow \infty} \zeta = \zeta^*$ , where*

$$\zeta^* = \mathbf{g}^\top \cdot \left( \mathbf{I}_b - \frac{1}{b\lambda_{\mathbf{a}}} \mathbf{a}\mathbf{a}^\top \left( \mathbf{I}_b + \frac{1}{b\lambda_{\mathbf{a}}} \mathbf{a}\mathbf{a}^\top \right)^{-1} \right) \cdot \mathbf{a}. \quad (21)$$

- 154 (i) *Here, the update direction  $\zeta^*$  is based only on the inputs and does not require computing  $\bar{\mathbf{g}}$*   
 155 *(which would require a second backpropagation pass), making it efficient.*  
 156 (ii) *The computational cost of computing the update  $\zeta^*$  lies in  $\mathcal{O}(bn^2 + b^2n + b^3)$ , where  $n$  is the*  
 157 *number of neurons in each layer. This comprises the conventional cost of computing the gradient*  
 158  *$\nabla = \mathbf{g}^\top \mathbf{x}$  lying in  $\mathcal{O}(bn^2)$ , and the overhead of computing  $\zeta^*$  instead of  $\nabla$  lying in  $\mathcal{O}(b^2n + b^3)$ .*  
 159 *The overhead is vanishing, assuming  $n \gg b$ . For  $b > n$  the complexity lies in  $\mathcal{O}(bn^2 + n^3)$ .*

160 *Proof.* We first show the property (21). Note that according to (22),  $\lambda_{\mathbf{g}} \cdot (\lambda_{\mathbf{g}} \mathbf{I}_m + \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b)^{-1}$  con-  
 161 verges in the limit of  $\lambda_{\mathbf{g}} \rightarrow \infty$  to  $\mathbf{I}_m$ , and therefore (21) holds.

162 (i) The statement follows from the fact that the term  $\bar{\mathbf{g}}$  does not appear in the equivalent characteriza-  
 163 tion (21) of  $\zeta^*$ .

164 (ii) We first note that the matrix  $\mathbf{a}\mathbf{a}^\top$  is of dimension  $b \times b$ , and can be computed in  $\mathcal{O}(b^2n)$  time.  
 165 Next, the matrix

$$\left( \mathbf{I}_b - \frac{1}{b\lambda_{\mathbf{a}}} \mathbf{a}\mathbf{a}^\top \left( \mathbf{I}_b + \frac{1}{b\lambda_{\mathbf{a}}} \mathbf{a}\mathbf{a}^\top \right)^{-1} \right)$$

166 is of shape  $b \times b$  and can be multiplied with  $\mathbf{a}$  in  $\mathcal{O}(b^2n)$  time. □

167 Notably, (21) can be computed with a vanishing computational overhead and with only minor  
 168 modifications to the implementation. Specifically, only the  $\mathbf{g}^\top \mathbf{a}$  expression has to be replaced by (21)  
 169 in the backpropagation step. As this can be done independently for each layer, this lends itself also to  
 170 applying it only to individual layers.

171 As we see in the experimental section, in many cases in the mini-batch regime (i.e.,  $b < n$ ), the  
 172 optimal (or a good) choice for  $\lambda_{\mathbf{g}}$  actually lies in the limit to  $\infty$ . This is a surprising result, leading to  
 173 the efficient and effective  $\zeta^* = \zeta_{\lambda_{\mathbf{g}} \rightarrow \infty}$  optimizer.

174 **Remark 2** (Relation between Update Direction  $\zeta$  and  $\zeta^*$ ). *When comparing the update direction*  
 175  *$\zeta$  in (20) without regularization (i.e.,  $\lambda_{\mathbf{g}} \rightarrow 0, \lambda_{\mathbf{a}} \rightarrow 0$ ) with  $\zeta^*$  (i.e.,  $\lambda_{\mathbf{g}} \rightarrow \infty$ ) as given in (21), it*  
 176 *can be directly seen that  $\zeta^*$  corresponds to a particular pre-conditioning of  $\zeta$ , since  $\zeta^* = M\zeta$  for*  
 177  *$M = \frac{1}{b\lambda_{\mathbf{g}}} \bar{\mathbf{g}}^\top \bar{\mathbf{g}}$ .*

178 As the last theoretical property of our proposed update direction  $\zeta^*$ , we show that in specific networks  
 179  $\zeta^*$  coincides with the Gauss-Newton update direction.

180 **Theorem 4** ( $\zeta^*$  is Exact for the Last Layer). *For the case of linear regression or, more generally, the*  
 181 *last layer of networks, with the mean squared error,  $\zeta^*$  is the Gauss-Newton update direction.*

182 *Proof.* The Hessian matrix of the mean squared error loss is the identity matrix. Correspondingly,  
 183 the expectation value of  $\bar{\mathbf{g}}^\top \bar{\mathbf{g}}$  is  $\mathbf{I}$ . Thus,  $\zeta^* = \zeta$ . □

184 **Remark 3.** *The direction  $\zeta^*$  corresponds to the Gauss-Newton update direction with an approxima-*  
 185 *tion of  $\mathbf{G}$  that can be expressed as  $\mathbf{G} \approx \mathbb{E}[\mathbf{I} \otimes (\mathbf{a}^\top \mathbf{a})]$ .*

186 **Remark 4** (Extension to the Natural Gradient). *In some cases, it might be more desirable to use the*  
 187 *Fisher-based natural gradient instead of the Gauss-Newton method. The difference to this setting is*  
 188 *that in (5) the GGN matrix  $\mathbf{G}$  is replaced by the empirical Fisher information matrix  $\mathbf{F}$ .*

189 *We note that our theory also applies to  $\mathbf{F}$ , and that  $\zeta^*$  also efficiently approximates the natural*  
 190 *gradient update step  $\mathbf{F}^{-1} \nabla$ . The  $i$ -th diagonal block of  $\mathbf{F}$  ( $\mathbf{F}_{\theta^{(i)}} = \mathbb{E}[(\mathbf{g}_i^\top \mathbf{g}_i) \otimes (a_{i-1}^\top \otimes a_{i-1})]$ ),  
 191 *has the same form as a block of the GGN matrix  $\mathbf{G}$  ( $\mathbf{G}_{\theta^{(i)}} = \mathbb{E}[(\bar{\mathbf{g}}_i^\top \bar{\mathbf{g}}_i) \otimes (a_{i-1}^\top \otimes a_{i-1})]$ ).*  
 192 *Thus, we can replace  $\bar{\mathbf{g}}$  with  $\mathbf{g}$  in our theoretical results to obtain their counterparts for  $\mathbf{F}$ .**

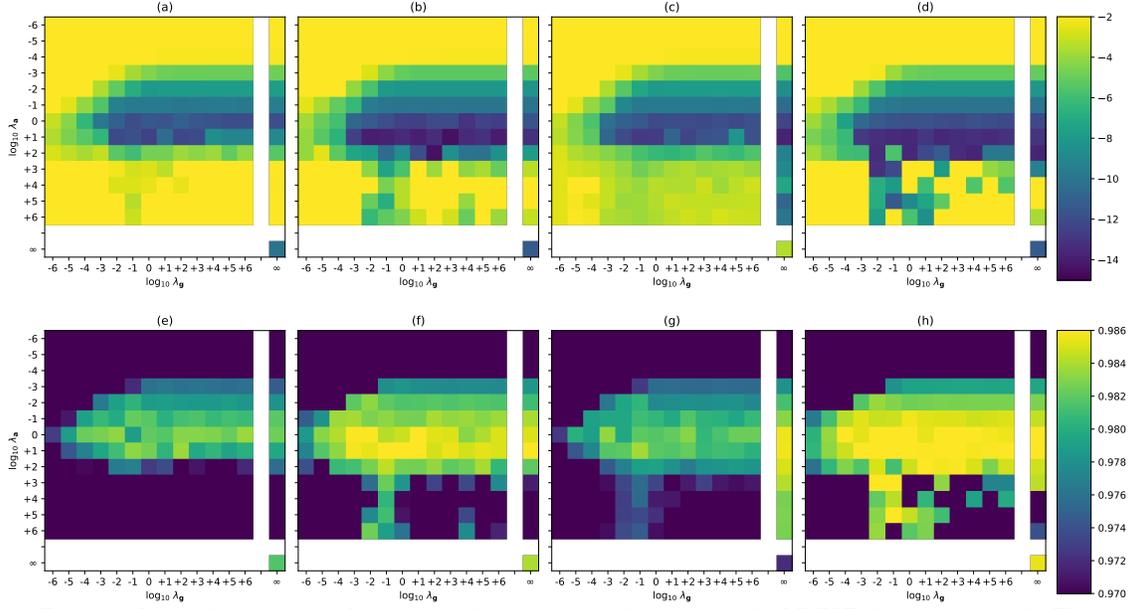


Figure 1: Logarithmic training loss (top) and test accuracy (bottom) on the MNIST classification task. The axes are the regularization parameters  $\lambda_g$  and  $\lambda_a$  in logarithmic scale with base 10. Training with a 5-layer ReLU activated network with 100 (left, a, e), 400 (center, b, c, f, g), and 1 600 (right, d, h) neurons per layer. The optimizer is SGD except for (c, g) where the optimizer is SGD with momentum. The top-left sector is  $\zeta$ , the top-right column is  $\zeta^*$ , and the bottom-right corner is  $\nabla$  (gradient descent). For each experiment and each of the three sectors, we use one learning rate, i.e.,  $\zeta$ ,  $\zeta^*$ ,  $\nabla$  have their own learning rate to make a fair comparison between the methods; within each sector the learning rate is constant. We can observe that in the limit of  $\lambda_g \rightarrow \infty$  (i.e., in the limit to the right) the performance remains good, showing the utility of  $\zeta^*$ .

## 193 4 Experiments

194 In the previous section, we discussed the theoretical properties of the proposed update directions  
 195  $\zeta$  and  $\zeta^*$  with the aspect that  $\zeta^*$  would actually be “free” to compute in the mini-batch regime. In  
 196 this section, we provide empirical evidence that  $\zeta^*$  is a good update direction, even in deep learning.  
 197 Specifically, we demonstrate that

- 198 (E1)  $\zeta^*$  achieves similar performance to K-FAC, while being substantially cheaper to compute.
- 199 (E2) The performance of our proposed method can be empirically maintained in the mini-batch  
 200 regime ( $n \gg b$ ).
- 201 (E3)  $\zeta^*$  may be used for individual layers, while for other layers only the gradient  $\nabla$  is used. This  
 202 still leads to improved performance.
- 203 (E4)  $\zeta^*$  also improves the performance for training larger models such as BERT and ResNet.
- 204 (E5) The runtime and memory requirements of  $\zeta^*$  are comparable to those of gradient descent.

### 205 E1: Impact of Regularization Parameters

206 For (E1), we study the dependence of the model’s performance on the regularization parameters  $\lambda_g$   
 207 and  $\lambda_a$ . Here, we train a 5-layer deep neural network on the MNIST classification task [16] with a  
 208 batch size of 60 for a total of 40 epochs or 40 000 steps.

209 The plots in Figure 1 demonstrate that the advantage of training by conditioning with curvature  
 210 information can be achieved by considering both layer inputs  $\mathbf{a}$  and gradients with respect to random  
 211 samples  $\tilde{\mathbf{g}}$ , but also using only layer inputs  $\mathbf{a}$ . In the plot, we show the performance of  $\zeta$  for different  
 212 choices of  $\lambda_g$  and  $\lambda_a$ , each in the range from  $10^{-6}$  to  $10^6$ . The right column shows  $\zeta^*$ , i.e.,  $\lambda_g = \infty$ ,  
 213 for different  $\lambda_a$ . The bottom-right corner is gradient descent, which corresponds to  $\lambda_g = \infty$  and  
 214  $\lambda_a = \infty$ .

215 Newton’s method or the general K-FAC approximation corresponds to the area with small  $\lambda_g$  and  $\lambda_a$ .  
 216 The interesting finding here is that the performance does not suffer by increasing  $\lambda_g$  toward  $\infty$ , i.e.,  
 217 from left to right in the plot.

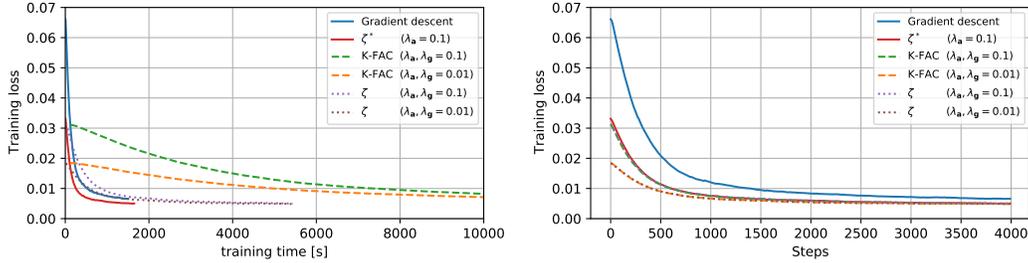


Figure 2: Training loss of the MNIST auto-encoder trained with gradient descent, K-FAC,  $\zeta$ , and  $\zeta^*$ . Comparing the performance per real-time (left) and per number of update steps (right). Runtimes are for a CPU core.

218 In addition, in Figure 3, we consider the case of regression with an auto-encoder trained with the  
 219 MSE loss on MNIST [16] and Fashion-MNIST [17]. Here, we follow the same principle as above  
 220 and also find that  $\zeta^*$  performs well.

221 In Figure 7, we compare the loss for different methods. Here, we distinguish  
 222 between loss per time (left) and loss per number of steps (right). We can observe that, for  $\lambda = 0.1$ , K-FAC,  $\zeta$ , and  
 223  $\zeta^*$  are almost identical per update step (right), while  $\zeta^*$  is by a large margin  
 224 the fastest, followed by  $\zeta$ , and the conventional K-FAC implementation is the  
 225 slowest (left). On the other hand, for  $\lambda = 0.01$  we can achieve a faster convergence  
 226 than with  $\lambda = 0.1$ , but here only the K-FAC and  $\zeta$  methods are numerically stable, while  $\zeta^*$  is unstable  
 227 in this case. This means in the regime of very small  $\lambda$ ,  $\zeta^*$  is not as robust as K-  
 228 FAC and  $\zeta$ , however, it achieves good performance with small but moderate  
 229  $\lambda$  like  $\lambda = 0.1$ . For  $\lambda < 0.01$ , also K-FAC and  $\zeta$  become numerically un-  
 230 stable in this setting and, in general, we observed that the smallest valid  $\lambda$  for  
 231 K-FAC is 0.01 or 0.001 depending on model and task. Under consideration  
 232 of the runtime,  $\zeta^*$  performs best as it is almost as fast as gradient descent while  
 233 performing equivalent to K-FAC and  $\zeta$ . Specifically, a gradient descent step is  
 234 only about 10% faster than  $\zeta^*$ .

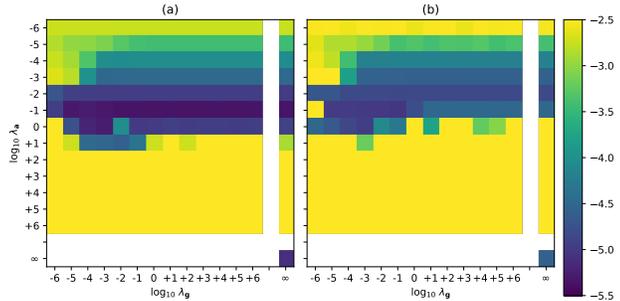


Figure 3: Training an auto-encoder on MNIST (left) and Fashion-MNIST (right). The model is the same as used by Botev *et al.* [18], i.e., it is a ReLU-activated 6-layer fully connected model with dimensions 784–1000–500–30–500–1000–784. Displayed is the logarithmic training loss.

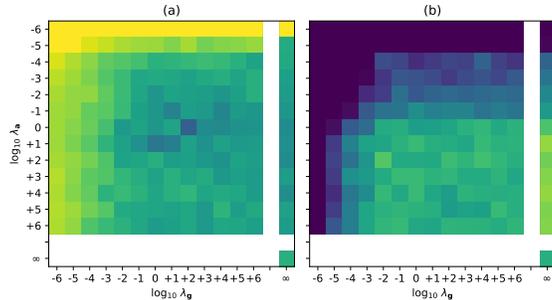


Figure 4: Training a 5-layer ReLU network with 400 neurons per layer on the MNIST classification task (as in Figure 1) but with the Adam optimizer [19].

## 250 E2: Minibatch Regime

251 For (E2), in Figure 1, we can see that training performs well for  $n \in \{100, 400, 1600\}$  neu-  
 252 rons per layer at a batch size of only 60. Also, in all other experiments, we use small batch sizes  
 253 of between 8 and 100.  
 254

## 256 E3: $\zeta^*$ in Individual Layers

257 In Figure 5, we train the 5-layer fully connected model with 400 neurons per layer. Here, we  
 258 consider the setting that we use  $\zeta^*$  in some of the layers while using the default gradient  $\nabla$   
 259 in other layers. Specifically, we consider the  
 260  
 261

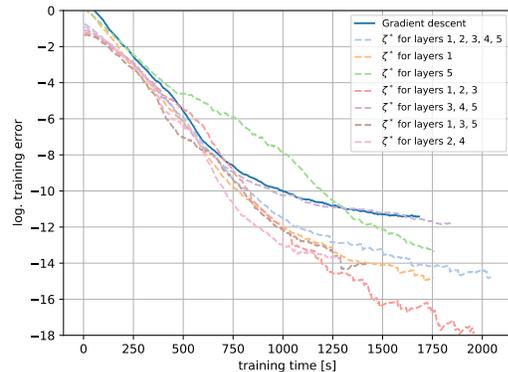


Figure 5: Training on the MNIST classification task using  $\zeta^*$  only in selected layers. Runtimes are for CPU.

Table 1: BERT results for fine-tuning pre-trained BERT-Base (B-B) and BERT-Mini (B-M) models on the CoLA, MRPC, and STSB text classification tasks. Larger values are better for all metrics. MCC is the Matthews correlation. Results averaged over 10 runs.

Method / Setting	CoLA (B-B)		MRPC (B-B)		STS-B (B-M)	
Metric	MCC	MCC	Acc.	F1	Pearson	Spearman
Gradient baseline	54.20 ± 7.56	21.08 ± 2.88	82.52 ± 1.22	87.88 ± 0.74	76.98 ± 1.10	76.88 ± 0.79
$\zeta^*$	57.62 ± 1.59	24.67 ± 2.62	83.28 ± 0.89	88.28 ± 0.70	81.09 ± 1.58	80.82 ± 1.57

262 settings, where all, the first, the final, the first three, the final three, the odd numbered, and the  
 263 even numbered layers are updated by  $\zeta^*$ . We observe that all settings with  $\zeta^*$  perform better than  
 264 plain gradient descent, except for “ $\zeta^*$  for layers 3,4,5” which performs approximately equivalent to  
 265 gradient descent.

266 **E4: Large-scale Models**

267 **BERT** To demonstrate the utility of  $\zeta^*$  also in large-scale models, we evaluate it for fine-tuning  
 268 BERT [20] on three natural language tasks. In Table 1, we summarize the results for the BERT  
 269 fine-tuning task. For the “Corpus of Linguistic Acceptability” (CoLA) [21] data set, we fine-tune  
 270 both the BERT-Base and the BERT-Mini models and find that we outperform the gradient descent  
 271 baseline in both cases. For the “Microsoft Research Paraphrase Corpus” (MRPC) [22] data set, we  
 272 fine-tune the BERT-Base model and find that we outperform the baseline both in terms of accuracy  
 273 and F1-score. Finally, on the “Semantic Textual Similarity Benchmark” (STS-B) [23] data set, we  
 274 fine-tune the BERT-Mini model and achieve higher Pearson and Spearman correlations than the  
 275 baseline. While for training with CoLA and MRPC, we were able to use the Adam optimizer [19]  
 276 (which is recommended for this task and model) in conjunction with  $\zeta^*$  in place of the gradient,  
 277 for STS-B Adam did not work well. Therefore, for STS-B, we evaluated it using the SGD with  
 278 momentum optimizer. For each method, we performed a grid search over the hyperparameters. We  
 279 note that we use a batch size of 8 in all BERT experiments.

280 **ResNet** In addition, we conduct an experiment  
 281 where we train the last layer of a ResNet with  
 282  $\zeta^*$ , while the remainder of the model is up-  
 283 dated using the gradient  $\nabla$ . Here, we train a  
 284 ResNet-18 [24] on CIFAR-10 [25] using SGD  
 285 with a batch size of 100 in a vanilla setting, i.e.,  
 286 without additional tricks employed in by He *et*  
 287 *al.* [24] and others. Specifically, we use (i) a  
 288 constant learning rate for each training (optimal  
 289 from (1, 0.3, 0.1, 0.03, 0.01)) and (ii) vanilla  
 290 SGD and not momentum-based SGD. The reason  
 291 behind this is that we want a vanilla exper-  
 292 iment and with aspects such as extensively tuning  
 293 multiple parameters of learning rate scheduler  
 294 would make the evaluation less transparent; how-  
 295 ever, therefore, all accuracies are naturally lower  
 296 than SOTA. In Figure 6, we plot the test accu-  
 297 racy against time. The results show that the pro-  
 298 posed method outperforms vanilla SGD when ap-  
 299 plied to the last layer of a ResNet-18. To vali-  
 date that the learning rate is not the cause for  
 the better performance, we also plot the neigh-  
 boring learning rates and find that even with a  
 too small or too large learning rate  $\zeta^*$  outper-  
 forms gradient descent with the optimal learning  
 rate.

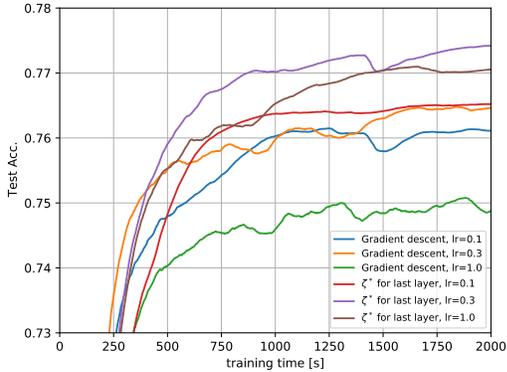


Figure 6: ResNet-18 trained on CIFAR-10. Runtimes are for a GPU. Results are averaged over 5 runs.

300 **E5: Runtime and Memory**

301 Finally, we also evaluate the runtime and memory requirements of each method. The runtime  
 302 evaluation is displayed in Table 2. We report both CPU and GPU runtime using PyTorch [26] and  
 303 (for K-FAC) the backpack library [15]. Note that the CPU runtime is more representative of the  
 304 pure computational cost, as for the first rows of the GPU runtime the overhead of calling the GPU  
 305 is dominant. When comparing runtimes between the gradient and  $\zeta^*$  on the GPU, we can observe  
 306 that we have an overhead of around 2.5 s independent of the model size. The overhead for CPU time  
 307 is also very small at less than 1% for the largest model, and only 1.3 s for the smallest model. In

308 contrast, the runtime of  $\zeta^*$  is around 4 times the runtime of the gradient, and K-FAC has an even  
 309 substantially larger runtime. Regarding memory,  $\zeta^*$  (contrasting the other approaches) also requires  
 310 only a small additional footprint.

311 **Remark 5** (Implementation). *The implementation of  $\zeta^*$  can be done by replacing the backpropagation*  
 312 *step of a respective layer by (21). As all “ingredients” are already available in popular deep learning*  
 313 *frameworks, it requires only little modification (contrasting K-FAC and  $\zeta$ , which require at least one*  
 314 *additional backpropagation.)*

Table 2: Runtimes and memory requirements for different models. Runtime is the training time per epoch on MNIST at a batch size of 60, i.e., for 1 000 training steps. The K-FAC implementation is from the backpack library [15]. The GPU is an Nvidia A6000.

Model	Gradient			K-FAC			$\zeta$			$\zeta^*$		
	CPU time	GPU time	Memory	CPU time	GPU t.	Memory	CPU time	GPU t.	Memory	CPU t.	GPU t.	Memory
5 layers w/ 100 n.	2.05 s	1.79 s	1.0 MB	62.78 s	17.63 s	11.5 MB	8.65 s	11.76 s	1.6 MB	3.34 s	4.07 s	1.0 MB
5 layers w/ 400 n.	23.74 s	1.84 s	4.8 MB	218.48 s	32.00 s	22.4 MB	38.67 s	12.62 s	7.7 MB	13.62 s	4.19 s	4.9 MB
5 layers w/ 1 600 n.	187.87 s	1.93 s	51.0 MB	6985.48 s	156.48 s	212.2 MB	665.80 s	12.53 s	85.8 MB	291.01 s	4.49 s	51.4 MB
5 layers w/ 6 400 n.	3439.59 s	8.22 s	691.0 MB	—	1320.81 s	3155.3 MB	9673 s	31.87 s	1197.8 MB	3451.61 s	10.24 s	692.5 MB
Auto-Encoder	78.61 s	2.20 s	16.2 MB	1207.58 s	74.09 s	70.7 MB	193.25 s	14.19 s	33.8 MB	87.39 s	4.93 s	16.5 MB

315 We will publish the source code of our implementation. In the appendix, we give a PyTorch [26]  
 316 implementation of the proposed method ( $\zeta^*$ ).

## 317 5 Related Work

318 Our methods are related to K-FAC by Martens and Grosse [12]. K-FAC uses the approximation  
 319 (13) to approximate the blocks of the Hessian of the empirical risk of neural networks. In most  
 320 implementations of K-FAC, the off-diagonal blocks of the Hessian are also set to zero. One of the  
 321 main claimed benefits of K-FAC is its speed (compared to stochastic gradient descent) for large-batch  
 322 size training. That said, recent empirical work has shown that this advantage of K-FAC disappears  
 323 once the additional computational costs of hyperparameter tuning for large batch training is accounted  
 324 for. There is a line of work that extends the basic idea of K-FAC to convolutional layers [27]. Botev *et*  
 325 *al.* [18] further extend these ideas to present KFLR, a Kronecker factored low-rank approximation,  
 326 and KFRA, a Kronecker factored recursive approximation of the Gauss-Newton step. Singh and  
 327 Alistarh [28] propose WoodFisher, a Woodbury matrix inverse-based estimate of the inverse Hessian,  
 328 and apply it to neural network compression. Yao *et al.* [29] propose AdaHessian, a second-order  
 329 optimizer that incorporates the curvature of the loss function via an adaptive estimation of the Hessian.  
 330 Frantar *et al.* [6] propose M-FAC, a matrix-free approximation of the natural gradient through a queue  
 331 of the (e.g., 1 000) recent gradients. These works fundamentally differ from our approach in that their  
 332 objective is to approximate the Fisher or Gauss-Newton matrix inverse vector products. In contrast,  
 333 this work proposes to approximate the Gauss-Newton matrix by only one of its Kronecker factors,  
 334 which we find to achieve good performance at a substantial computational speedup and reduction of  
 335 memory footprint. For an overview of this area, we refer to Kunstner *et al.* [30] and Martens [31].  
 336 For an overview of the technical aspects of backpropagation of second-order quantities, we refer to  
 337 Dangel *et al.* [15], [32]

338 Taking a step back, K-FAC is one of many Newton-type methods for training neural networks.  
 339 Other prominent examples of such methods include subsampled Newton methods [33], [34] (which  
 340 approximate the Hessian by subsampling the terms in the empirical risk function and evaluating the  
 341 Hessian of the subsampled terms) and sketched Newton methods [3]–[5] (which approximate the  
 342 Hessian by sketching, e.g., by projecting the Hessian to a lower-dimensional space by multiplying it  
 343 with a random matrix). The main features that distinguish K-FAC from this group of methods are  
 344 K-FAC’s superior empirical performance and K-FAC’s lack of theoretical justification.

## 345 6 Conclusion

346 In this work, we presented ISAAC Newton, a novel approximate curvature method based on layer-  
 347 inputs. We demonstrated it to be a special case of the regularization-generalized Gauss-Newton  
 348 method and empirically demonstrate its utility. Specifically, our method features an asymptotically  
 349 vanishing computational overhead in the mini-batch regime, while achieving competitive empirical  
 350 performance on various benchmark problems.

351 **References**

- 352 [1] N. Agarwal, B. Bullins, and E. Hazan, “Second-order stochastic optimization for machine  
353 learning in linear time,” *Journal on Machine Learning Research*, vol. 18, no. 1, pp. 4148–4187,  
354 2017.
- 355 [2] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2e. New York, NY, USA: Springer, 2006.
- 356 [3] A. Gonen and S. Shalev-Shwartz, “Faster SGD using sketched conditioning,” *arXiv preprint*,  
357 *arXiv:1506.02649*, 2015.
- 358 [4] M. Pilanci and M. J. Wainwright, “Newton sketch: A near linear-time optimization algorithm  
359 with linear-quadratic convergence,” *SIAM Journal on Optimization*, vol. 27, 2017.
- 360 [5] M. A. Erdogdu and A. Montanari, “Convergence rates of sub-sampled Newton methods,” in  
361 *Proc. Neural Information Processing Systems (NeurIPS)*, 2015.
- 362 [6] E. Frantar, E. Kurtic, and D. Alistarh, “M-FAC: Efficient matrix-free approximations of  
363 second-order information,” in *Proc. Neural Information Processing Systems (NeurIPS)*, 2021.
- 364 [7] N. Doikov and Y. Nesterov, “Convex Optimization based on Global Lower Second-order  
365 Models,” in *Proc. Neural Information Processing Systems (NeurIPS)*, Curran Associates, Inc.,  
366 2020.
- 367 [8] Y. Nesterov and B. T. Polyak, “Cubic regularization of Newton method and its global perfor-  
368 mance,” *Mathematical Programming*, vol. 108, 2006.
- 369 [9] S. Becker and Y. Lecun, “Improving the convergence of back-propagation learning with  
370 second-order methods,” 1989.
- 371 [10] T. Schaul, S. Zhang, and Y. LeCun, “No more pesky learning rates,” in *International Conference*  
372 *on Machine Learning (ICML)*, 2013.
- 373 [11] Y. Ollivier, “Riemannian metrics for neural networks i: Feedforward networks,” *Information*  
374 *and Inference*, vol. 4, pp. 108–153, Jun. 2015.
- 375 [12] J. Martens and R. Grosse, “Optimizing neural networks with Kronecker-factored approximate  
376 curvature,” in *International Conference on Machine Learning (ICML)*, 2015.
- 377 [13] A. N. Tikhonov and V. Y. Arsenin, *Solutions of Ill-posed problems*. W.H. Winston, 1977.
- 378 [14] P. Chen, “Hessian matrix vs. Gauss—Newton Hessian matrix,” *SIAM Journal on Numerical*  
379 *Analysis*, 2011.
- 380 [15] F. Dangel, F. Kunstner, and P. Hennig, “Backpack: Packing more into backprop,” in *Internat-*  
381 *ional Conference on Learning Representations*, 2020.
- 382 [16] Y. LeCun, C. Cortes, and C. Burges, “MNIST Handwritten Digit Database,” *ATT Labs*, 2010.
- 383 [17] H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-MNIST: A novel image dataset for benchmarking  
384 machine learning algorithms,” *arXiv*, 2017.
- 385 [18] A. Botev, H. Ritter, and D. Barber, “Practical Gauss-Newton optimisation for deep learning,”  
386 in *International Conference on Machine Learning (ICML)*, 2017.
- 387 [19] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Confer-*  
388 *ence on Learning Representations (ICLR)*, 2015.
- 389 [20] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional  
390 transformers for language understanding,” in *North American Chapter of the Association for*  
391 *Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2018.
- 392 [21] A. Warstadt, A. Singh, and S. R. Bowman, “Neural network acceptability judgments,” *Trans-*  
393 *actions of the Association for Computational Linguistics*, vol. 7, 2019.
- 394 [22] W. B. Dolan and C. Brockett, “Automatically constructing a corpus of sentential paraphrases,”  
395 in *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005.
- 396 [23] D. Cer, M. Diab, E. Agirre, I. Lopez-Gazpio, and L. Specia, “SemEval-2017 task 1: Semantic  
397 textual similarity multilingual and crosslingual focused evaluation,” in *Proceedings of the*  
398 *11th International Workshop on Semantic Evaluation (SemEval-2017)*, Vancouver, Canada:  
399 Association for Computational Linguistics, 2017.
- 400 [24] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in  
401 *Proc. International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 402 [25] A. Krizhevsky, V. Nair, and G. Hinton, “Cifar-10 (Canadian Institute for Advanced Research),”  
403 2009.
- 404 [26] A. Paszke, S. Gross, F. Massa, *et al.*, “Pytorch: An imperative style, high-performance deep  
405 learning library,” in *Proc. Neural Information Processing Systems (NeurIPS)*, 2019.

- 406 [27] R. Grosse and J. Martens, “A Kronecker-factored approximate Fisher matrix for convolution  
407 layers,” in *International Conference on Machine Learning (ICML)*, 2016.
- 408 [28] S. P. Singh and D. Alistarh, “Woodfisher: Efficient second-order approximation for neural  
409 network compression,” in *Proc. Neural Information Processing Systems (NeurIPS)*, 2020.
- 410 [29] Z. Yao, A. Gholami, S. Shen, M. Mustafa, K. Keutzer, and M. W. Mahoney, “Adahessian:  
411 An adaptive second order optimizer for machine learning,” in *AAAI Conference on Artificial  
412 Intelligence*, 2021.
- 413 [30] F. Kunstner, L. Balles, and P. Hennig, “Limitations of the empirical Fisher approximation for  
414 natural gradient descent,” in *Proc. Neural Information Processing Systems (NeurIPS)*, 2019.
- 415 [31] J. Martens, “New insights and perspectives on the natural gradient method,” *Journal of Machine  
416 Learning Research*, 2020.
- 417 [32] F. Dangel, S. Harmeling, and P. Hennig, “Modular block-diagonal curvature approximations  
418 for feedforward architectures,” in *International Conference on Artificial Intelligence and  
419 Statistics (AISTATS)*, 2020.
- 420 [33] F. Roosta-Khorasani and M. W. Mahoney, “Sub-Sampled Newton Methods I: Globally Con-  
421 vergent Algorithms,” *arXiv: 1601.04737*, 2016.
- 422 [34] P. Xu, J. Yang, F. Roosta, C. Ré, and M. W. Mahoney, “Sub-sampled Newton Methods with  
423 Non-uniform Sampling,” in *Proc. Neural Information Processing Systems (NeurIPS)*, 2016.

424 **Checklist**

- 425 1. For all authors...
- 426 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
427 contributions and scope? [Yes]
- 428 (b) Did you describe the limitations of your work? [Yes]
- 429 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 430 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them?  
431 [Yes]
- 432 2. If you are including theoretical results...
- 433 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 434 (b) Did you include complete proofs of all theoretical results? [Yes]
- 435 3. If you ran experiments...
- 436 (a) Did you include the code, data, and instructions needed to reproduce the main experimental  
437 results (either in the supplemental material or as a URL)? [Yes] / [No] We include a  
438 Python / PyTorch implementation of the method in the supplementary material. We will  
439 publicly release full source code for the experiments.
- 440 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were  
441 chosen)? [Yes]
- 442 (c) Did you report error bars (e.g., with respect to the random seed after running experiments  
443 multiple times)? [Yes]
- 444 (d) Did you include the total amount of compute and the type of resources used (e.g., type of  
445 GPUs, internal cluster, or cloud provider)? [Yes]
- 446 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 447 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 448 (b) Did you mention the license of the assets? [N/A]
- 449 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 450 (d) Did you discuss whether and how consent was obtained from people whose data you're  
451 using/curating? [N/A]
- 452 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
453 information or offensive content? [N/A]
- 454 5. If you used crowdsourcing or conducted research with human subjects...
- 455 (a) Did you include the full text of instructions given to participants and screenshots, if  
456 applicable? [N/A]
- 457 (b) Did you describe any potential participant risks, with links to Institutional Review Board  
458 (IRB) approvals, if applicable? [N/A]
- 459 (c) Did you include the estimated hourly wage paid to participants and the total amount spent  
460 on participant compensation? [N/A]

## 461 A PyTorch Implementation

462 We display a PyTorch [26] implementation of ISAAC for a fully-connected layer below. Here, we  
463 mark the important part (i.e., the part beyond the boilerplate) with a red rectangle.

```
import torch

class ISAACLinearFunction(torch.autograd.Function):
    @staticmethod
    def forward(ctx, input, weight, bias, la, inv_type):
        ctx.save_for_backward(input, weight, bias)
        ctx.la = la
        if inv_type == 'cholesky_inverse':
            ctx.inverse = torch.cholesky_inverse
        elif inv_type == 'inverse':
            ctx.inverse = torch.inverse
        else:
            raise NotImplementedError(inv_type)
        return input @ weight.T + (bias if bias is not None else 0)

    @staticmethod
    def backward(ctx, grad_output):
        input, weight, bias = ctx.saved_tensors
        if ctx.needs_input_grad[0]:
            grad_0 = grad_output @ weight
        else:
            grad_0 = None

        if ctx.needs_input_grad[1]:
            aaT = input @ input.T / grad_output.shape[0]
            I_b = torch.eye(aaT.shape[0], device=aaT.device, dtype=aaT.dtype)
            aaT_IaaT_inv = aaT @ ctx.inverse(aaT / ctx.la + I_b)
            grad_1 = grad_output.T @ (
                I_b - 1. / ctx.la * aaT_IaaT_inv
            ) @ input
        else:
            grad_1 = None

        return (
            grad_0,
            grad_1,
            grad_output.mean(0, keepdim=True) if bias is not None else None,
            None, None, None,
        )

class ISAACLinear(torch.nn.Linear):
    def __init__(self, in_features, out_features,
                 la, inv_type='inverse', **kwargs):
        super(ISAACLinear, self).__init__(
            in_features=in_features, out_features=out_features, **kwargs
        )
        self.la = la
        self.inv_type = inv_type

    def forward(self, input: torch.Tensor) -> torch.Tensor:
        return ISAACLinearFunction.apply(
            input, self.weight,
```

```

self.bias.unsqueeze(0) if self.bias is not None else None,
self.la,
self.inv_type
)

```

## 464 B Implementation Details

465 Unless noted differently, for all experiments, we tune the learning rate on a grid of  
466 (1, 0.3, 0.1, 0.03, 0.01, 0.003, 0.001). We verified this range to cover the full reasonable range of  
467 learning rates. Specifically, for every single experiment, we made sure that there is no learning rate  
468 outside this range which performs better.

469 For all language model experiments, we used the respective Huggingface PyTorch implementation.

470 All other hyperparameter details are given in the main paper.

471 The code will be made publicly available.

## 472 C Additional Proofs

473 *Proof of Theorem 1.* We first show, that  $\zeta$  as defined in (17) can be expressed as in (20). Indeed by  
474 using (19), the Woodbury matrix identity and by regularizing the inverses, we can see that

$$\begin{aligned}
\zeta &= \lambda_g \lambda_a (\bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b + \lambda_g \mathbf{I})^{-1} \otimes (\mathbf{a}^\top \mathbf{a}/b + \lambda_a \mathbf{I})^{-1} \mathbf{g}^\top \mathbf{a} \\
&= \lambda_g \lambda_a \cdot (\lambda_g \mathbf{I}_m + \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b)^{-1} \mathbf{g}^\top \mathbf{a} (\lambda_a \mathbf{I}_n + \mathbf{a}^\top \mathbf{a}/b)^{-1} \\
&= \lambda_g \lambda_a \cdot \left( \frac{1}{\lambda_g} \mathbf{I}_m - \frac{1}{b \lambda_g^2} \bar{\mathbf{g}}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_g} \bar{\mathbf{g}} \bar{\mathbf{g}}^\top \right)^{-1} \bar{\mathbf{g}} \right) \\
&\quad \mathbf{g}^\top \mathbf{a} \left( \frac{1}{\lambda_a} \mathbf{I}_n - \frac{1}{b \lambda_a^2} \mathbf{a}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_a} \mathbf{a} \mathbf{a}^\top \right)^{-1} \mathbf{a} \right) \\
&= \left( \mathbf{I}_m - \frac{1}{b \lambda_g} \bar{\mathbf{g}}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_g} \bar{\mathbf{g}} \bar{\mathbf{g}}^\top \right)^{-1} \bar{\mathbf{g}} \right) \cdot \mathbf{g}^\top \\
&\quad \cdot \mathbf{a} \cdot \left( \mathbf{I}_n - \frac{1}{b \lambda_a} \mathbf{a}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_a} \mathbf{a} \mathbf{a}^\top \right)^{-1} \mathbf{a} \right) \\
&= \left( \mathbf{I}_m - \frac{1}{b \lambda_g} \bar{\mathbf{g}}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_g} \bar{\mathbf{g}} \bar{\mathbf{g}}^\top \right)^{-1} \bar{\mathbf{g}} \right) \cdot \mathbf{g}^\top \\
&\quad \cdot \left( \mathbf{a} - \frac{1}{b \lambda_a} \mathbf{a} \mathbf{a}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_a} \mathbf{a} \mathbf{a}^\top \right)^{-1} \mathbf{a} \right) \\
&= \left( \mathbf{I}_m - \frac{1}{b \lambda_g} \bar{\mathbf{g}}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_g} \bar{\mathbf{g}} \bar{\mathbf{g}}^\top \right)^{-1} \bar{\mathbf{g}} \right) \cdot \mathbf{g}^\top \\
&\quad \cdot \left( \mathbf{I}_b - \frac{1}{b \lambda_a} \mathbf{a} \mathbf{a}^\top \left( \mathbf{I}_b + \frac{1}{b \lambda_a} \mathbf{a} \mathbf{a}^\top \right)^{-1} \right) \cdot \mathbf{a}
\end{aligned}$$

475 To show Assertion (i), we note that according to (17)

$$\begin{aligned}
&\lim_{\lambda_g, \lambda_a \rightarrow 0} \frac{1}{\lambda_g \lambda_a} \zeta \\
&= \lim_{\lambda_g, \lambda_a \rightarrow 0} (\bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b + \lambda_g \mathbf{I})^{-1} \otimes (\mathbf{a}^\top \mathbf{a}/b + \lambda_a \mathbf{I})^{-1} \mathbf{g}^\top \mathbf{a} \\
&= (\bar{\mathbf{g}}^\top \bar{\mathbf{g}})^{-1} \otimes (\mathbf{a}^\top \mathbf{a})^{-1} \mathbf{g}^\top \mathbf{a} \\
&\approx \mathbf{G}^{-1} \mathbf{g}^\top \mathbf{a},
\end{aligned}$$

476 where the first equality uses the definition of  $\zeta$  in (17). The second equality is due to the continuity of  
 477 the matrix inversion and the last approximate equality follows from the K-FAC approximation (15).

478 To show Assertion (ii), we consider  $\lim_{\lambda_{\mathbf{g}} \rightarrow \infty}$  and  $\lim_{\lambda_{\mathbf{a}} \rightarrow \infty}$  independently, that is

$$\begin{aligned} & \lim_{\lambda_{\mathbf{g}} \rightarrow \infty} \lambda_{\mathbf{g}} \cdot (\lambda_{\mathbf{g}} \mathbf{I}_m + \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b)^{-1} & (22) \\ &= \lim_{\lambda_{\mathbf{g}} \rightarrow \infty} \left( \mathbf{I}_m + \frac{1}{b\lambda_{\mathbf{g}}} \bar{\mathbf{g}}^\top \bar{\mathbf{g}} \right)^{-1} = \mathbf{I}_m, \end{aligned}$$

479 and

$$\begin{aligned} & \lim_{\lambda_{\mathbf{a}} \rightarrow \infty} \lambda_{\mathbf{a}} \cdot (\lambda_{\mathbf{a}} \mathbf{I}_n + \mathbf{a}^\top \mathbf{a}/b)^{-1} & (23) \\ &= \lim_{\lambda_{\mathbf{a}} \rightarrow \infty} \left( \mathbf{I}_n + \frac{1}{b\lambda_{\mathbf{a}}} \mathbf{a}^\top \mathbf{a} \right)^{-1} = \mathbf{I}_n. \end{aligned}$$

480 This then implies

$$\begin{aligned} & \lim_{\lambda_{\mathbf{g}}, \lambda_{\mathbf{a}} \rightarrow \infty} \lambda_{\mathbf{g}} (\lambda_{\mathbf{g}} \mathbf{I}_m + \bar{\mathbf{g}}^\top \bar{\mathbf{g}}/b)^{-1} \cdot \mathbf{g}^\top & (24) \\ & \quad \cdot \mathbf{a} \cdot \lambda_{\mathbf{a}} (\lambda_{\mathbf{a}} \mathbf{I}_n + \mathbf{a}^\top \mathbf{a}/b)^{-1} \\ &= \mathbf{I}_m \cdot \mathbf{g}^\top \mathbf{a} \cdot \mathbf{I}_n = \mathbf{g}^\top \mathbf{a}, \end{aligned}$$

481 which concludes the proof. □

482 **D Additional Experiments**

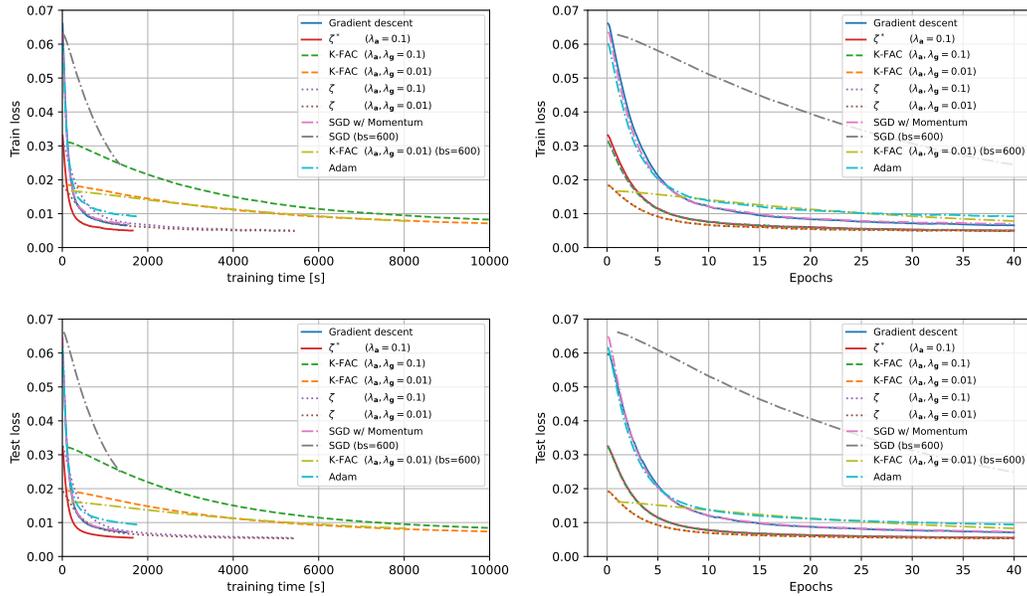


Figure 7: Training loss of the MNIST auto-encoder trained with gradient descent, K-FAC,  $\zeta$ ,  $\zeta^*$ , as well as SGD w/ momentum, SGD with a  $10\times$  larger batch size (600), K-FAC with a  $10\times$  larger batch size (600), and Adam. Comparing the performance per real-time (left) and per number of epochs (right). We display both the training loss (top) as well as the test loss (bottom) Runtimes are for a CPU core.

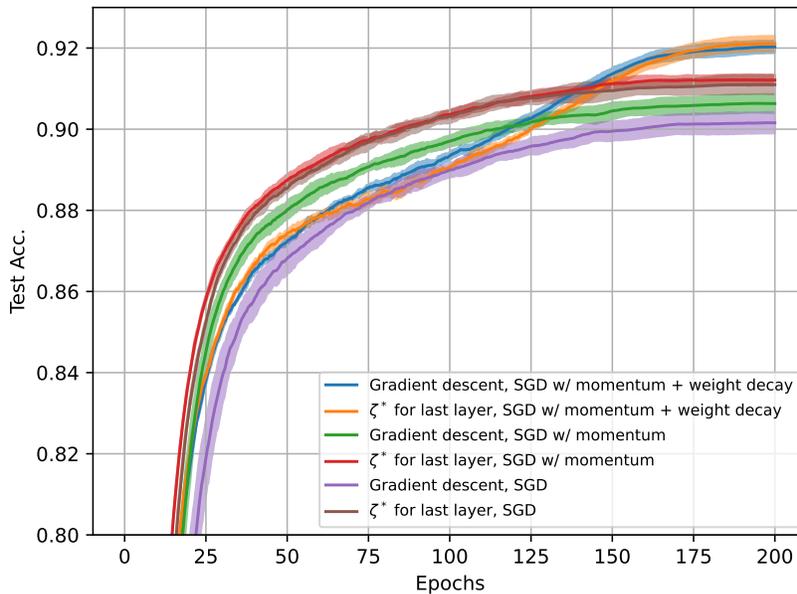


Figure 8: ResNet-18 trained on CIFAR-10 with image augmentation and a cosine learning rate schedule. The first line (blue) uses the hyperparameters of a public implementation. To ablate the optimizer, two additional settings are added, specifically, without weight decay and without momentum. Results are averaged over 5 runs and the standard deviation is indicated with the colored areas.

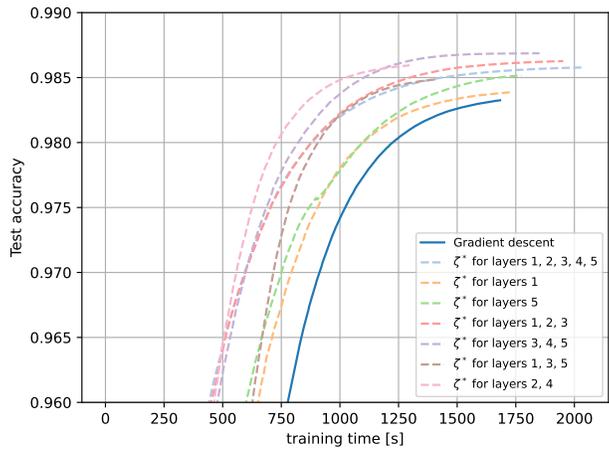


Figure 9: Test accuracy for training on the MNIST classification task using  $\zeta^*$  only in selected layers. Runtimes are for CPU.