

DEBIASING MINI-BATCH QUADRATICS FOR APPLICATIONS IN DEEP LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Quadratic approximations form a fundamental building block of machine learning methods. E.g., second-order optimizers try to find the Newton step into the minimum of a local quadratic proxy to the objective function; and the second-order approximation of a network’s loss function can be used to quantify the uncertainty of its outputs via the Laplace approximation. When computations on the entire training set are intractable—typical for deep learning—the relevant quantities are computed on mini-batches. This, however, distorts and biases the shape of the associated *stochastic* quadratic approximations in an intricate way with detrimental effects on applications. In this paper, we (i) show that this bias introduces a systematic error, (ii) provide a theoretical explanation for it, (iii) explain its relevance for second-order optimization and uncertainty quantification via the Laplace approximation in deep learning, and (iv) develop and evaluate debiasing strategies.

1 INTRODUCTION

Quadratic approximations of the loss landscape are increasingly used by algorithms in deep learning, from pruning methods (Dong et al., 2017; Zeng & Urtasun, 2019) and influence functions (Koh & Liang, 2017) to second-order optimizers (Amari, 1998; Martens, 2010; Martens & Grosse, 2015; Grosse & Martens, 2016; Botev et al., 2017; Zhang et al., 2017; George et al., 2018; Martens et al., 2018; Osawa et al., 2019) and uncertainty quantification via the Laplace approximation (Ritter et al., 2018b;a; Kristiadi et al., 2020; Daxberger et al., 2021; Immer et al., 2021). When the computations are intractable on the entire training set—typical for deep learning—the quantities of interest are computed on mini-batches subsampled from the training data. The goal of this work is to highlight that mini-batching systematically biases the shape of a quadratic approximation.

A systematic bias? Figure 1 illustrates the phenomenon. It shows five *mini-batch* quadratics in their top-curvature 2D subspace for the fully trained ALL-CNN-C model on CIFAR-100 data. For comparison, the full-batch quadratic, where all quantities are evaluated on the *entire* training set, is projected into the *same* 2D subspace. Within that subspace, the two quadratics are quite different: The mini-batch quadratic is much “narrower” (exhibits larger curvature) than the full-batch version. Given that the full-batch quadratic is the “right” object to serve as the basis for, e.g., a Newton step or a Laplace approximation,¹ the mini-batch version is *not* a meaningful surrogate: Its Newton step is overly small and a Laplace approximation yields an overconfident uncertainty estimate.

Contributions. To enable stable and efficient stochastic second-order optimizers, as well as reliable techniques for uncertainty quantification, we analyze this phenomenon and develop strategies to mitigate it. More specifically, our contributions are as follows: (i) We study mini-batch quadratics empirically and show that their geometry is systematically biased; (ii) we provide an explanation for this phenomenon, explaining the bias as an instance of the classic regression to the mean (directions of extreme steepness/curvature for one particular mini-batch are less extreme for other mini-batches), (iii) we explain the relevance of this bias for second-order optimization and uncertainty quantification via the Laplace approximation, and (iv) develop and evaluate debiasing strategies.

¹We are *not* concerned with the approximation error arising from the quadratic approximation of the non-quadratic function $\mathcal{L}_{\text{reg}}(\cdot; \mathcal{D}) \approx q(\cdot; \mathcal{D})$ (see Equation (2)), but only with the consequences of replacing the full-batch quantities by their mini-batch counterparts $q(\cdot; \mathcal{D}) \approx q(\cdot; \mathcal{B})$.

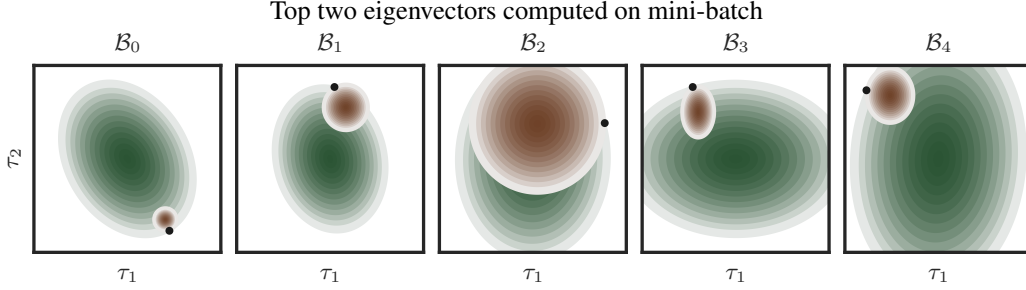


Figure 1: A systematic bias? We compute five mini-batch quadratics $q(\cdot; \mathcal{B}_m)$ with batch size $|\mathcal{B}_m| = 512$ for the loss landscape of the fully trained ALL-CNN-C model on CIFAR-100 data around $\theta_0 \leftarrow \theta_*$ (shown as \bullet). Each mini-batch quadratic defines a 2D subspace spanned by the top two eigenvectors $\mathbf{u}_1, \mathbf{u}_2$ of $\mathbf{H}_{\mathcal{B}_m}$, in which we evaluate (i) the quadratic $q(\theta_* + \tau_1 \mathbf{u}_1 + \tau_2 \mathbf{u}_2; \mathcal{B}_m)$ itself (shown in \square) and (ii) the full-batch quadratic $q(\theta_* + \tau_1 \mathbf{u}_1 + \tau_2 \mathbf{u}_2; \mathcal{D})$ (shown in \square). In that subspace, the mini-batch quadratic is much “narrower” than the full-batch version which leads to overly small Newton steps and overconfident uncertainty estimates via the Laplace approximation.

2 NOTATION & BACKGROUND

The regularized loss. We consider a general supervised learning problem, where we try to find the optimal parameters $\theta_* = \arg \min_{\theta \in \mathbb{R}^P} \mathcal{L}_{\text{reg}}(\theta; \mathcal{D})$ for the parameterized function $f_\theta: \mathbb{R}^D \rightarrow \mathbb{R}^C$ by minimizing the regularized loss \mathcal{L}_{reg} on N training examples $\mathbb{D} := \{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{R}^D \times \mathbb{R}^C\}_{n \in \mathcal{D}}$,

$$\mathcal{L}_{\text{reg}}(\theta; \mathcal{D}) := \mathcal{L}(\theta; \mathcal{D}) + r(\theta) \quad \text{with} \quad \mathcal{L}(\theta; \mathcal{D}) := \frac{1}{|\mathcal{D}|} \sum_{n \in \mathcal{D}} \ell(f_\theta(\mathbf{x}_n), \mathbf{y}_n), \quad \mathcal{D} = \{1, \dots, N\}. \quad (1)$$

The empirical risk \mathcal{L} measures the dissimilarity between the model’s predictions $f_\theta(\mathbf{x}_n)$ and the true outputs \mathbf{y}_n via a loss function $\ell: \mathbb{R}^C \times \mathbb{R}^C \rightarrow \mathbb{R}$. The regularizer $r: \mathbb{R}^P \rightarrow \mathbb{R}$, $r(\theta) := \beta/2 \|\theta\|_2^2$ with $\beta \in \mathbb{R}_{\geq 0}$, penalizes the “complexity” of the model.

2.1 LOCAL FULL-BATCH & MINI-BATCH QUADRATIC APPROXIMATIONS

Full-batch quadratic. A local quadratic approximation of the regularized loss around $\theta_0 \in \mathbb{R}^P$ is given by the second-order Taylor expansion

$$\mathcal{L}_{\text{reg}}(\theta; \mathcal{D}) \approx q(\theta; \mathcal{D}) := \frac{1}{2} (\theta - \theta_0)^\top \mathbf{H}_{\mathcal{D}} (\theta - \theta_0) + (\theta - \theta_0)^\top \mathbf{g}_{\mathcal{D}} + c_{\mathcal{D}}, \quad (2)$$

where $c_{\mathcal{D}} := \mathcal{L}_{\text{reg}}(\theta_0; \mathcal{D})$, $\mathbf{g}_{\mathcal{D}} := \nabla \mathcal{L}_{\text{reg}}(\theta_0; \mathcal{D})$ and $\mathbf{H}_{\mathcal{D}} := \nabla^2 \mathcal{L}_{\text{reg}}(\theta_0; \mathcal{D}) = \nabla^2 \mathcal{L}(\theta_0; \mathcal{D}) + \beta \mathbf{I}$ (some approximation of) the Hessian at θ_0 (all derivatives are with respect to the parameters θ unless stated otherwise). As all quantities are evaluated on the *entire* training set \mathcal{D} , we refer to this as the *full-batch* quadratic. It holds that $\nabla q(\theta; \mathcal{D}) = \mathbf{H}_{\mathcal{D}} (\theta - \theta_0) + \mathbf{g}_{\mathcal{D}}$ and $\nabla^2 q(\theta; \mathcal{D}) \equiv \mathbf{H}_{\mathcal{D}}$.

Mini-batch quadratic. When the computations are intractable on the entire training set, the quantities in Equation (2) are typically computed on a *mini-batch*—a small randomly drawn subset—of the training data $\mathcal{B} \subset \mathcal{D}$, $|\mathcal{B}| \ll N$, resulting in a *stochastic* quadratic approximation $q(\cdot; \mathcal{B}) \approx q(\cdot; \mathcal{D})$. As $c_{\mathcal{B}}$, $\mathbf{g}_{\mathcal{B}}$ and $\mathbf{H}_{\mathcal{B}}$ are unbiased estimates of $c_{\mathcal{D}}$, $\mathbf{g}_{\mathcal{D}}$ and $\mathbf{H}_{\mathcal{D}}$, this substitution may seem innocent, but, as we will see in Section 3, it affects the geometry of the quadratic approximation *substantially*.

Directional slope and curvature. Consider a cut r through the quadratic $q(\cdot; \mathcal{B})$ from $\theta_* \in \mathbb{R}^P$ along the normalized direction \mathbf{d} , $\|\mathbf{d}\| = 1$. It holds that (derivation in Appendix A.1) $r(\tau) := q(\theta_* + \tau \mathbf{d}; \mathcal{B}) = 1/2 \tau^2 \mathbf{d}^\top \nabla^2 q(\theta_*; \mathcal{B}) \mathbf{d} + \tau \mathbf{d}^\top \nabla q(\theta_*; \mathcal{B}) + \text{const}$. So, as a function of τ , $r: \mathbb{R} \rightarrow \mathbb{R}$ is a 1D parabola with derivatives $r'(\tau) = \tau \mathbf{d}^\top \nabla^2 q(\theta_*; \mathcal{B}) \mathbf{d} + \mathbf{d}^\top \nabla q(\theta_*; \mathcal{B})$ and $r''(\tau) \equiv \mathbf{d}^\top \nabla^2 q(\theta_*; \mathcal{B}) \mathbf{d}$. We denote the *directional slope* and *curvature* of the quadratic $q(\cdot, \mathcal{B})$ at θ_* in direction \mathbf{d} by

$$\partial_{\mathbf{d}} q(\theta_*; \mathcal{B}) := r'(0) = \mathbf{d}^\top \nabla q(\theta_*; \mathcal{B}) \quad \text{and} \quad \partial_{\mathbf{d}}^2 q(\theta_*; \mathcal{B}) := r''(0) = \mathbf{d}^\top \nabla^2 q(\theta_*; \mathcal{B}) \mathbf{d}.$$

The directional slope and curvature at θ_* are thus simply the projections of the quadratic’s gradient and Hessian at that location onto the direction.

Eigenvalues as directional curvatures. The directional curvature of the quadratic $q(\cdot; \mathcal{B})$ along one of $\mathbf{H}_{\mathcal{B}}$'s normalized eigenvectors \mathbf{u} coincides with the corresponding eigenvalue λ since $\partial_{\mathbf{u}}^2 q(\boldsymbol{\theta}_*; \mathcal{B}) = \mathbf{u}^\top \nabla^2 q(\boldsymbol{\theta}_*; \mathcal{B}) \mathbf{u} = \mathbf{u}^\top \mathbf{H}_{\mathcal{B}} \mathbf{u} = \lambda \|\mathbf{u}\|^2 = \lambda$. Thus, in the context of a quadratic, an eigenvalue of the Hessian $\mathbf{H}_{\mathcal{B}}$ has a geometric interpretation as the directional curvature along the respective eigenvector.

GGN & FIM. Next, we discuss second-order optimization methods and the Laplace approximation (LA) for neural networks. Both techniques rely on local quadratic approximations of the regularized loss function and require a *positive definite* curvature matrix $\mathbf{H}_{\mathcal{B}}$. The empirical risk's Hessian $\nabla^2 \mathcal{L}(\boldsymbol{\theta}_0; \mathcal{B})$ can be *indefinite* and is therefore typically replaced by (an approximation of) the positive semi-definite Generalized Gauss-Newton matrix (GGN) $\mathbf{G}_{\mathcal{B}}$ or Fisher information matrix (FIM) $\mathbf{F}_{\mathcal{B}}$. In fact, GGN and FIM are often identical (Martens, 2020, Sec. 9.2). The resulting $\mathbf{H}_{\mathcal{B}}$ is positive definite if $\beta > 0$, or when a damping term $\delta \mathbf{I}$, $\delta \in \mathbb{R}_{>0}$ is added (e.g., in trust-region methods).

2.2 SECOND-ORDER METHODS & CONJUGATE GRADIENTS

The Newton step. Due to the simple polynomial form of the quadratic $q(\cdot; \mathcal{B})$, its minimum can be derived in closed form and is given by the Newton step $\arg \min_{\boldsymbol{\theta}} q(\boldsymbol{\theta}; \mathcal{B}) = \boldsymbol{\theta}_0 - \mathbf{H}_{\mathcal{B}}^{-1} \mathbf{g}_{\mathcal{B}}$. This serves as the basis for second-order optimizers like L-BFGS (Liu & Nocedal, 1989; Nocedal, 1980), the Hessian-free approach (Martens, 2010) or K-FAC (Kronecker-factored approximate curvature) (Martens & Grosse, 2015; Grosse & Martens, 2016; Martens et al., 2018).

Conjugate gradients. We focus on the method of conjugate gradients (CG) (Hestenes & Stiefel, 1952), as it is a powerful method specifically designed to minimize quadratics with positive definite Hessians effectively (details in Appendix A.2). It is particularly useful in the context of large-scale optimization because it only requires access to matrix-vector products with the curvature matrix $\mathbf{v} \mapsto \mathbf{H}_{\mathcal{B}} \mathbf{v}$ that can be computed in a matrix-free manner (Pearlmutter, 1994; Schraudolph, 2002), i.e. without ever materializing the full matrix in memory—and it has been used successfully for training neural networks (Martens, 2010). Starting at $\boldsymbol{\theta}_0$, CG creates a sequence of iterates $(\boldsymbol{\theta}_0, \dots, \boldsymbol{\theta}_P)$. In each iteration p , two main steps are performed: (i) Given the current position $\boldsymbol{\theta}_p$ and a normalized search direction \mathbf{d}_p , the algorithm finds the minimum of the quadratic along $\boldsymbol{\theta}_p + \tau \mathbf{d}_p$, i.e.

$$\boldsymbol{\theta}_{p+1} = \boldsymbol{\theta}_p + \tau_p \mathbf{d}_p \quad \text{with} \quad \tau_p := \arg \min_{\tau \in \mathbb{R}} q(\boldsymbol{\theta}_p + \tau \mathbf{d}_p; \mathcal{B}) = -\frac{\partial_{\mathbf{d}_p} q(\boldsymbol{\theta}_p; \mathcal{B})}{\partial_{\mathbf{d}_p}^2 q(\boldsymbol{\theta}_p; \mathcal{B})}. \quad (3)$$

In the second step (ii), CG constructs the next search direction \mathbf{d}_{p+1} such that it is *conjugate* to all previous search directions, i.e. $\mathbf{d}_{p+1}^\top \mathbf{H}_{\mathcal{B}} \mathbf{d}_i = 0$ for all $i \in \{1, \dots, p\}$.

2.3 LAPLACE APPROXIMATION FOR NEURAL NETWORKS

Laplace approximation. The Laplace approximation (LA) turns a trained standard neural network into a Bayesian neural network in a post-hoc manner (MacKay, 1991; Ritter et al., 2018b; Kristiadi et al., 2020; Daxberger et al., 2021) (details in Appendix A.3). The idea is to reinterpret the regularized loss \mathcal{L}_{reg} as the negative unnormalized log posterior $\log p(\boldsymbol{\theta} \mid \mathbb{D})$ of a specific Bayesian model. This interpretation identifies the optimal parameters $\boldsymbol{\theta}_* = \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\text{reg}}(\boldsymbol{\theta}; \mathcal{D}) = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathbb{D})$ as the mode of the posterior, i.e. as the maximum a posteriori (MAP) estimate. A second-order approximation of the regularized loss around $\boldsymbol{\theta}_0 \leftarrow \boldsymbol{\theta}_*$ then translates into a Gaussian approximation of the posterior—the so-called Laplace approximation (MacKay, 1992), i.e.

$$\mathcal{L}_{\text{reg}}(\boldsymbol{\theta}; \mathcal{D}) \approx q(\boldsymbol{\theta}; \mathcal{D}) = \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_*)^\top \mathbf{H}_{\mathcal{D}}(\boldsymbol{\theta} - \boldsymbol{\theta}_*) + \text{const.} \quad \rightsquigarrow \quad p(\boldsymbol{\theta} \mid \mathbb{D}) \approx \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\theta}_*, \boldsymbol{\Sigma}_{\mathcal{D}}), \quad (4)$$

with $\boldsymbol{\Sigma}_{\mathcal{D}} := N^{-1} \mathbf{H}_{\mathcal{D}}^{-1}$ (or an approximation thereof). We obtain the predictive uncertainty $p(\mathbf{y}_{\diamond} \mid \mathbf{x}_{\diamond}, \mathbb{D})$ for some test input $\mathbf{x}_{\diamond} \in \mathbb{R}^D$ by propagating the parameter uncertainty $\mathcal{N}(\boldsymbol{\theta}_*, \boldsymbol{\Sigma}_{\mathcal{D}})$ to the model's outputs via the linearized network (Immer et al., 2021; Roy et al., 2024).

Full-batch vs. mini-batch LA. It is common practice to compute the LA on the *entire* training set. Depending on the curvature approximation, this comes at considerable computational costs. For example, a full-batch LA is prohibitive even for moderately sized models/datasets when a low-rank approximation of the Hessian is computed via repeated Hessian-vector products (each of which requires a full pass over the training data); and it is essentially infeasible for model classes like LLMs,

162 which are trained on massive datasets. However, what adds most to the costs is that it is standard
 163 to *tune* the prior precision (Daxberger et al., 2021). This adds another outer loop that requires the
 164 same procedure to be performed multiple times. Being able to emulate the behavior of the full-batch
 165 quadratic on a mini-batch would thus be useful. We therefore study the mini-batch setting, i.e. we
 166 replace $q(\cdot; \mathcal{D})$ by $q(\cdot; \mathcal{B})$ in Equation (4). We will see in Section 6.2 that, when mitigating the
 167 associated biases, a mini-batch LA can be a good proxy for the full-batch LA.

169 3 THE SHAPE OF A MINI-BATCH QUADRATIC

170
 171 This section studies the “shape” of a mini-batch quadratic $q(\cdot; \mathcal{B})$ and how it differs from $q(\cdot; \mathcal{D})$.
 172

173 3.1 EMPIRICAL STUDY OF THE DIRECTIONAL SLOPES AND CURVATURES

174
 175 **The high-curvature subspace is relevant for all common use cases.** In our empirical study, we
 176 focus on the *top*-curvature subspace. This subspace is particularly relevant for several reasons:
 177

- 178 1. By the Eckart-Young-Mirsky Theorem, a truncated SVD is Frobenius norm-optimal, i.e. a
 179 low-rank approximation using the top eigenvectors is ideal from a theoretical perspective.
 180 As the spectrum of the Hessian typically decays quickly (Ghorbani et al., 2019), the bulk of
 181 the curvature information is contained in the top-curvature subspace.
- 182 2. In the context of optimization, it has been observed that the gradient mainly lives in the
 183 high-curvature space (Gur-Ari et al., 2018; Dangel et al., 2022). Thus, it makes sense for
 184 a second-order method to operate mainly in that space, as steps outside of it can not be
 185 expected to reduce the objective function significantly.
- 186 3. In the context of the LA, directions of large curvature correspond to directions in the weight
 187 space with *low* variance, i.e. these directions capture what we *know* about the model’s
 188 parameters. Consequently, Daxberger et al. (2021, p. 19) describe a low-rank approximation
 189 of the Hessian based on its top eigenvectors.

190 **Experimental procedure.** We use the same setting as in Figure 1: The fully trained ALL-CNN-C
 191 model on the CIFAR-100 dataset. To isolate the effect of data subsampling, we eliminate *all other*
 192 sources of noise. Thus, we remove the dropout layers from the model. We use the cross-entropy loss,
 193 an ℓ^2 -regularizer, and train the model with SGD for 350 epochs, see Appendix B.1 for details.

194 We then pick a mini-batch \mathcal{B}_m of size 512 and compute the 100 eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_{100}$ to the 100
 195 *largest* eigenvalues of $\mathbf{H}_{\mathcal{B}_m} \leftarrow \mathbf{G}_{\mathcal{B}_m} + \beta \mathbf{I}$. That is the Hessian of the ℓ^2 -regularized mini-batch loss,
 196 where we replaced $\nabla^2 \mathcal{L}(\boldsymbol{\theta}_*; \mathcal{B}_m)$ with the GGN approximation. Next, we compute the directional
 197 slopes and curvatures for *all* mini-batch quadratics $q(\cdot; \mathcal{B}_{m'})$, $m' \in \{1, \dots, M\}$ along those 100
 198 eigenvectors. For a fixed eigenvector \mathbf{u}_p , the *average* of those directional slopes/curvatures over *all*
 199 mini-batches coincides with the directional slope/curvature of the full-batch quadratic, i.e.

$$200 \frac{1}{M} \sum_{m'=1}^M \underbrace{\partial_{\mathbf{u}_p} q(\boldsymbol{\theta}_*; \mathcal{B}_{m'})}_{\text{one } \bullet \text{ and many } \circ} = \underbrace{\partial_{\mathbf{u}_p} q(\boldsymbol{\theta}_*; \mathcal{D})}_{\text{+}} \quad \text{and} \quad \frac{1}{M} \sum_{m'=1}^M \underbrace{\partial_{\mathbf{u}_p}^2 q(\boldsymbol{\theta}_*; \mathcal{B}_{m'})}_{\text{one } \bullet \text{ and many } \circ} = \underbrace{\partial_{\mathbf{u}_p}^2 q(\boldsymbol{\theta}_*; \mathcal{D})}_{\text{+}} \quad (5)$$

203 derivation in in Appendix A.1. The colored markers in Equation (5) refer to Figure 2. We repeat this
 204 procedure for three mini-batches \mathcal{B}_m , $m \in \{0, 1, 2\}$.
 205

206 **Directional slopes and curvatures are biased.** Figure 2 reveals a systematic bias in the directional
 207 slopes and, more pronounced, in the directional curvatures: When eigenvectors and directional
 208 derivatives are computed on the *same* mini-batch, curvature is overestimated by roughly one order
 209 of magnitude compared to the curvature of the full-batch quadratic $q(\cdot; \mathcal{D})$. Within the space that
 210 actually carries curvature information—the top-eigenspace of the quadratic’s Hessian—the curvature
 211 magnitude is thus not at all representative of the true underlying loss landscape. For *other* mini-
 212 batches, the directional slopes and curvatures are similar to the full-batch quadratic. This is because
 213 the averages in Equation (5) are dominated by these unbiased samples.

214 We present additional results for batch size 2048 (instead of 512), the Hessian $\nabla^2 \mathcal{L}(\boldsymbol{\theta}_*; \mathcal{B}_m)$ (instead
 215 of its GGN approximation) and CG search directions (instead of eigenvectors) in Appendix B.4. The
 bias is present in *all* settings, but less pronounced for larger batch sizes. For the CG search directions,

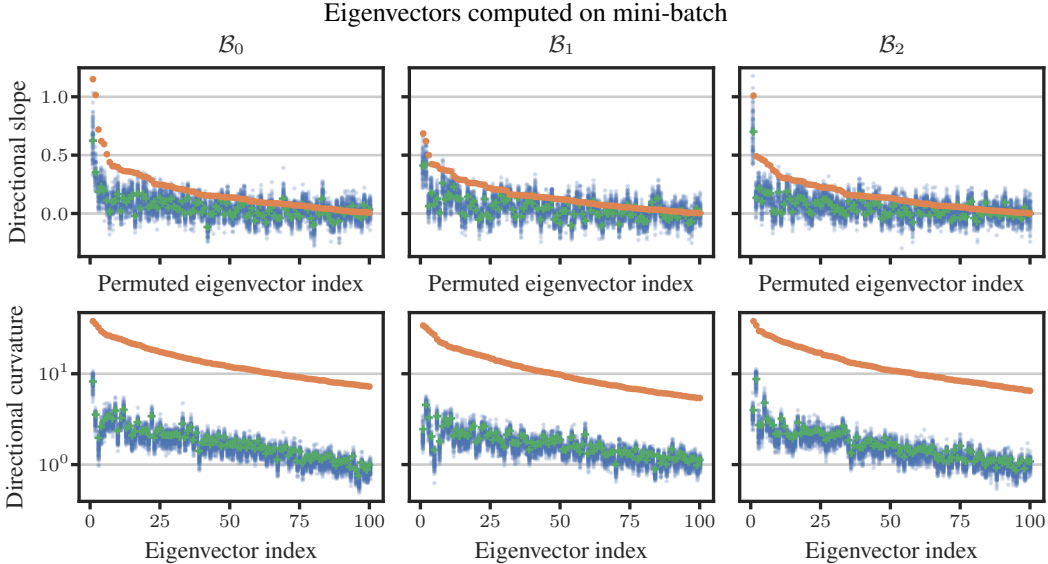


Figure 2: Directional slopes and curvatures are biased. We use the CIFAR-100 dataset with the fully trained ALL-CNN-C model and draw three mini-batches \mathcal{B}_m of size $|\mathcal{B}_m| = 512$ to compute the top 100 eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_{100}$. For each mini-batch/column, we show the directional slopes (*Top*) and curvatures (*Bottom*) evaluated on (i) $q(\boldsymbol{\theta}_*; \mathcal{B}_m)$ (i.e. on the *same* mini-batch of data) as \bullet , (ii) $q(\boldsymbol{\theta}_*; \mathcal{B}_{m'})$ for $m' \neq m$ (i.e. for all *other* mini-batches) as \bullet and (iii) the full-batch quadratic (the average of the orange and all blue dots, see Equation (5)) as \oplus . For the top panel, we switch the order and sign of the eigenvectors such that the orange dots are all above zero and in descending order. There is a strong, systematic bias, particularly in the curvature: Computing the eigenvectors and directional curvatures on the same data results in over-estimation by roughly one order of magnitude.

the bias in directional slope is much more distinct than for the eigenvectors: The biased slope is always negative, while the full-batch quadratic’s slope is in fact positive for most CG directions!

3.2 WHERE DO THE BIASES COME FROM?

In the following, we provide an explanation for the biases in the directional slopes and curvatures (i.e. the gap between the orange and a blue dot in Figure 2) from a theoretical perspective.

3.2.1 BIAS IN THE DIRECTIONAL SLOPE

Here, we focus on the CG search directions since the bias in the directional slope is even more pronounced for those directions than for the eigenvectors. This is not accidental! The CG directions are constructed from gradients—and gradients are directions that *maximize* the steepness \bullet for one particular mini-batch quadratic. For *other* mini-batch quadratics, the steepness along those directions (i.e. the directional slope \bullet) is therefore *less* extreme. We formalize this intuition in Appendix A.4.

3.2.2 BIAS IN THE DIRECTIONAL CURVATURE

Next, we consider the bias in the directional curvature along the eigenvectors of the curvature matrix.

Directional curvature along $\mathbf{H}_{\mathcal{B}}$ ’s eigenvectors. Let $\mathbf{u}_1, \dots, \mathbf{u}_P$ and $\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_P$ denote the eigenvectors of two mini-batch Hessians $\mathbf{H}_{\mathcal{B}}$ and $\mathbf{H}_{\tilde{\mathcal{B}}}$, respectively. Assume that the corresponding eigenvalues are in descending order, i.e. $\lambda_1 \geq \dots \geq \lambda_P$ and $\tilde{\lambda}_1 \geq \dots \geq \tilde{\lambda}_P$. We can write the directional curvatures on \mathcal{B} and $\tilde{\mathcal{B}}$ along an eigenvector \mathbf{u}_i as (details in Appendix A.5)

$$\underbrace{\partial_{\mathbf{u}_i}^2 q(\boldsymbol{\theta}_*; \tilde{\mathcal{B}})}_{\bullet} = \mathbf{u}_i^\top \mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{u}_i = \sum_{p=1}^P \tilde{\lambda}_p \Omega_{i,p}, \quad \text{and} \quad \underbrace{\partial_{\mathbf{u}_i}^2 q(\boldsymbol{\theta}_*; \mathcal{B})}_{\bullet} = \mathbf{u}_i^\top \mathbf{H}_{\mathcal{B}} \mathbf{u}_i = \lambda_i. \quad (6)$$

The weights $\{\Omega_{i,p} := (\mathbf{u}_i^\top \tilde{\mathbf{u}}_p)^2\}_{p=1}^P$ are non-negative and sum to one, i.e. $\sum_{p=1}^P \Omega_{i,p} = 1$. The bias in the directional curvature thus originates from misalignment of the eigenspaces of the two curvature matrices—captured by the weights $\Omega_{i,p}$ —and/or a systematic difference in their spectra.

Curvature bias is not due to differing spectra... Assume that the eigenspaces are perfectly aligned, i.e. $\Omega_{i,i} = 1$ and $\Omega_{i,p} = 0 \forall p \neq i$. In this case, we obtain $\partial_{\mathbf{u}_i}^2 q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}}) = \tilde{\lambda}_i$ from Equation (6), so the bias originates *exclusively* from the differences in the spectra. Figure 2 shows the eigenvalues (as directional curvatures) for three different CIFAR-100 mini-batches. As they are very similar, this cannot serve as the main explanation for the curvature bias.

... but due to misaligned eigenspaces. Now, assume that the spectra of $\mathbf{H}_{\mathcal{B}}$ and $\mathbf{H}_{\tilde{\mathcal{B}}}$ are identical, i.e. $\lambda_p = \tilde{\lambda}_p \forall p \in \{1, \dots, P\}$ which simplifies the unbiased estimate in Equation (6) to $\partial_{\mathbf{u}_i}^2 q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}}) = \sum_{p=1}^P \lambda_p \Omega_{i,p}$. Consider the curvature along \mathbf{u}_1 as an example. If $\mathbf{u}_1 = \tilde{\mathbf{u}}_1$, there is only one non-zero weight $\Omega_{1,1} = 1$ and the two directional curvatures are identical. However, if there is significant overlap with any other eigenvectors (i.e. some weight is distributed also on the *lower*-curvature directions), the resulting curvature $\partial_{\mathbf{u}_1}^2 q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})$ is *smaller* than $\partial_{\mathbf{u}_1}^2 q(\boldsymbol{\theta}_\bullet; \mathcal{B})$. Analogously, for \mathbf{u}_P , the directional curvature on $\tilde{\mathcal{B}}$ is *larger* than on \mathcal{B} if there is significant overlap between \mathbf{u}_P and some of $\mathbf{H}_{\tilde{\mathcal{B}}}$'s *higher*-curvature eigenvectors. This can be formalized by the following inequalities (derivation in Appendix A.5):

$$\underbrace{\partial_{\mathbf{u}_1}^2 q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})}_{\bullet} \leq \underbrace{\partial_{\mathbf{u}_1}^2 q(\boldsymbol{\theta}_\bullet; \mathcal{B})}_{\bullet} \quad \text{and} \quad \underbrace{\partial_{\mathbf{u}_P}^2 q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})}_{\bullet} \geq \underbrace{\partial_{\mathbf{u}_P}^2 q(\boldsymbol{\theta}_\bullet; \mathcal{B})}_{\bullet}. \quad (7)$$

In general, if the eigenspaces are not perfectly aligned such that weight is distributed on several eigenvectors, this leads to *less extreme* directional curvatures on mini-batch $\tilde{\mathcal{B}}$ on *both* ends of the spectrum. Figure 3 shows the weights $\Omega_{i,p}$ as pixels. The overlap between the eigenspaces is far from perfect, i.e. one eigenvector from \mathcal{B} overlaps with several eigenvectors from $\tilde{\mathcal{B}}$ to some extent. This explains the curvature overestimation in the top curvature subspace we observe in Figure 2 and identifies the misalignment of the eigenspaces as the dominating factor for the curvature bias.

3.2.3 SUMMARY OF THE THEORETICAL FINDINGS

The CG search directions and the eigenvectors of the curvature matrix are both *designed* to be “extreme” in some sense: The CG directions are based on gradients that maximize steepness of the quadratic $q(\cdot; \mathcal{B})$, while the top/bottom eigenvectors subsume directions of largest/smallest curvature. However, these directions are extreme only for one particular mini-batch \mathcal{B} . Another mini-batch $\tilde{\mathcal{B}}$ has its own extreme directions that typically differ from those of \mathcal{B} . Thus, **the extreme directions for \mathcal{B} are less extreme for $\tilde{\mathcal{B}}$** . Projecting both quadratics onto \mathcal{B} 's directions consequently leads to extreme steepness and curvatures for \mathcal{B} , but less extreme values for $\tilde{\mathcal{B}}$. This result is an instance of the classic **regression to the mean** (Galton, 1886): Using an algorithm to find the directions of most extreme steepness/curvature in one particular mini-batch, we must expect the steepness/curvature to be *less extreme* on most *other* batches (and thus also on the entire dataset). We can expect this phenomenon to occur for other datasets, models, and curvature proxies as well, since the underlying mechanism is the stochasticity of the geometric information.

4 DEBIASING MINI-BATCH QUADRATICS FOR APPLICATIONS

We here argue that the biases affect second-order applications, and develop debiasing strategies.

4.1 IMPLICATIONS FOR SECOND-ORDER OPTIMIZATION AND THE LAPLACE APPROXIMATION

Detrimental updates in second-order optimizers. In the context of second-order optimization, both the biases in the directional slopes and curvatures are relevant. From Equation (3), the CG update magnitude τ_p to reach the minimum of $q(\boldsymbol{\theta}_p + \tau \mathbf{d}_p; \mathcal{B}_m)$ is given by the negative ratio of the directional slope and curvature at the current iterate $\boldsymbol{\theta}_p$ in the direction \mathbf{d}_p —that is the 1D Newton step along that cut.² As both these quantities are biased, so is τ_p as shown in Figure 4: While the

²The exact same argument can also be made for the eigenvector directions as the Newton step can be decomposed into 1D Newton steps along the eigenvectors.

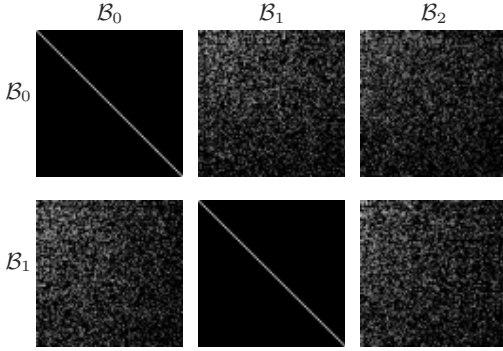


Figure 3: In practice, eigenspaces are misaligned. We reuse the setting of Figure 2 and compute the top 100 eigenvectors $U_m \in \mathbb{R}^{P \times 100}$ for $\{\mathcal{B}_m\}_{m \in \{0,1,2\}}$. The weights $\Omega_{i,j}$ are shown as a 100×100 greyscale image (color ranges from black for $\Omega_{i,j} \leq 10^{-8}$ to white for $\Omega_{i,j} = 1$) for $m \in \{0,1\}$, $m' \in \{0,1,2\}$. Clearly, the eigenspaces for different mini-batches are not perfectly aligned as eigenvectors from \mathcal{B}_m overlap with several eigenvectors from $\mathcal{B}_{m'}$.

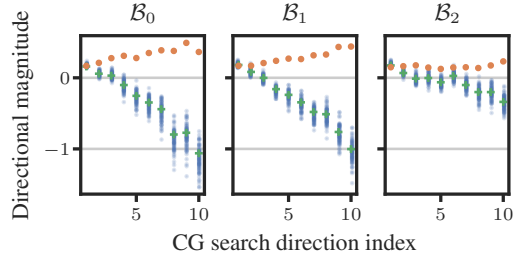


Figure 4: CG update magnitudes are biased. Same setting as Figure 9 (Bottom). We run CG on $\{\mathcal{B}_m\}_{m \in \{0,1,2\}}$ and show the directional update magnitudes τ_1, \dots, τ_{10} for the first 10 CG steps using (i) the same mini-batch \mathcal{B}_m (as \bullet), (ii) all other mini-batches (as \bullet) and (iii) the entire training set (as $\+$). The magnitudes are given by the negative ratio of the directional slope and curvature (see Equation (3)) and thus inherit the attached biases. Note that most of the update magnitudes that are based on a single mini-batch of data (\bullet) have the wrong sign resulting in detrimental updates in the wrong direction.

update magnitudes for \mathcal{B}_m are always positive (a property of CG), minimizing the other—*equally valid*—quadratics would require *negative* update magnitudes for most of the CG directions. In this sense, naive CG on a single mini-batch of data makes updates in the wrong direction. This can be attributed to the bias in the directional slope since the slope determines the sign of τ_p . Overestimation of the directional curvature is “accidentally beneficial” in this case as it leads to smaller steps.

Unreliable uncertainty quantification with the Laplace approximation. For the LA, only the bias in the directional curvature is relevant. After Section 2.3, the approximate posterior covariance over the network’s parameters is given by $\Sigma_{\mathcal{B}} = N^{-1} \mathbf{H}_{\mathcal{B}}^{-1}$. Via the eigendecomposition $\mathbf{H}_{\mathcal{B}} = \sum_{p=1}^P \lambda_p \mathbf{u}_p \mathbf{u}_p^T$, we obtain $\Sigma_{\mathcal{B}} = N^{-1} \sum_{p=1}^P \lambda_p^{-1} \mathbf{u}_p \mathbf{u}_p^T$. Due to the biases we describe in Section 3, the directional curvatures λ_p are not representative of the true underlying curvature, resulting in a *deformed* posterior covariance $\Sigma_{\mathcal{B}}$. Specifically, the overestimation of the curvature in the top-curvature subspace translates into an *underestimation* of the uncertainty in the posterior covariance (due to the inversion) with potentially severe consequences in safety-critical applications.

4.2 DEBIASING STRATEGIES

We now turn to strategies that mitigate the biases in second-order optimization and the LA. An empirical evaluation of these methods follows in Sections 6.1 and 6.2.

Decoupling directions from magnitudes with a two-batch strategy. The approaches we propose are simple yet effective. The idea is to commit to the imperfect directions from one mini-batch (as before) but use an *additional* mini-batch to estimate the directional derivatives. With this, we *decouple* the mechanism that determines the parameter subspace in which the method operates from the slope and curvature measurements within the space. Effectively, we use the blue dots from Figures 2 and 4 instead of the orange dots and thus obtain much more realistic estimates of the actual loss function’s geometry (within the subspace defined by the first mini-batch). Next, we describe in more detail how this strategy can be applied to debias CG and the LA.

4.2.1 DEBIASED CONJUGATE GRADIENTS

Debiased approach. For the debiased CG method, we need to implement two processes: (i) The first process applies $K \leq P$ CG iterations to the mini-batch quadratic $q(\cdot; \mathcal{B})$ and collects the search directions $(\mathbf{d}_1, \dots, \mathbf{d}_K)$. As before, this defines the subspace in which CG operates. (ii) The second process recomputes the trajectory $(\hat{\theta}_1, \dots, \hat{\theta}_K)$ within that subspace using debiased update

magnitudes that are computed on a different mini-batch $\tilde{\mathcal{B}}$, i.e. we use $\tilde{\theta}_0 := \theta_0$ and

$$\tilde{\theta}_{p+1} := \tilde{\theta}_p + \tilde{\tau}_p \mathbf{d}_p \quad \text{with} \quad \tilde{\tau}_p := - \frac{\partial_{\mathbf{d}_p} q(\tilde{\theta}_p; \tilde{\mathcal{B}})}{\partial_{\mathbf{d}_p}^2 q(\tilde{\theta}_p; \tilde{\mathcal{B}})} \quad \text{instead of} \quad \tau_p = - \frac{\partial_{\mathbf{d}_p} q(\theta_p; \mathcal{B})}{\partial_{\mathbf{d}_p}^2 q(\theta_p; \mathcal{B})}. \quad (8)$$

For $\tilde{\mathcal{B}} = \mathcal{B}$, $(\tilde{\theta}_1, \dots, \tilde{\theta}_K)$ is congruent with the original trajectory $(\theta_1, \dots, \theta_K)$ from the single-batch CG approach. If the two mini-batches are different, the debiased trajectory will use updates into the directional minima of $q(\cdot; \tilde{\mathcal{B}})$ instead of $q(\cdot; \mathcal{B})$. Processes (i) and (ii) can either be run side by side in an alternating fashion (offering more flexibility regarding e.g. the termination criterion at the cost of a small memory overhead, details in [Appendix A.2](#)) or one after the other.

4.2.2 DEBIASED LAPLACE APPROXIMATION

Laplace approximation with K-FAC. We briefly explore another popular curvature proxy: Kronecker-Factored Approximate Curvature (K-FAC), which is commonly used both in optimization ([Martens & Grosse, 2015](#); [Martens et al., 2018](#); [Eschenhagen et al., 2023](#)) and uncertainty quantification ([Ritter et al., 2018b](#)). It is a block-diagonal approximation to the FIM $\mathbf{F}_{\mathcal{B}} \approx \mathbf{K}_{\mathcal{B}} := \text{blockdiag}_{l=1, \dots, L}(\mathbf{K}_{\mathcal{B}}^{(l)})$, where each block $\mathbf{K}_{\mathcal{B}}^{(l)} := \mathbf{A}^{(l)} \otimes \mathbf{B}^{(l)}$ is approximated by the Kronecker product of two smaller matrices $\mathbf{A}^{(l)}$ and $\mathbf{B}^{(l)}$. Sampling from the respective LA with covariance $\Sigma_{\mathcal{B}} = N^{-1}(\mathbf{K}_{\mathcal{B}} + \beta \mathbf{I})^{-1}$ can be performed efficiently due to the specific structure of the K-FAC approximation. For details, see [Appendix A.3.2](#).

Debiased approach. For the debiased LA, we use two K-FAC approximations $\mathbf{K}_{\mathcal{B}}$ and $\mathbf{K}_{\tilde{\mathcal{B}}}$ computed on different mini-batches. Via the eigendecomposition $\mathbf{K}_{\mathcal{B}} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^{\top}$ with $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_P)$ and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_P)$, we can re-write the covariance matrix as $\Sigma_{\mathcal{B}} = N^{-1}(\mathbf{K}_{\mathcal{B}} + \beta \mathbf{I})^{-1} = N^{-1} \mathbf{U}(\mathbf{\Lambda} + \beta \mathbf{I})^{-1} \mathbf{U}^{\top}$. For the debiased approach, we keep the eigenspace defined by $\mathbf{K}_{\mathcal{B}}$ but recompute the directional curvatures based on $\mathbf{K}_{\tilde{\mathcal{B}}}$, i.e. we use

$$\tilde{\Sigma}_{\mathcal{B}} = \frac{1}{N} \mathbf{U} (\text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_P) + \beta \mathbf{I})^{-1} \mathbf{U}^{\top} \quad \text{with} \quad \tilde{\lambda}_p = \mathbf{u}_p^{\top} \mathbf{K}_{\tilde{\mathcal{B}}} \mathbf{u}_p \quad \text{instead of} \quad \lambda_p = \mathbf{u}_p^{\top} \mathbf{K}_{\mathcal{B}} \mathbf{u}_p. \quad (9)$$

4.2.3 COMPUTATIONAL COST OF DEBIASING

Both debiasing techniques roughly double runtime compared to their single-batch counterparts: Debiased CG performs one extra matrix-vector product with $\mathbf{H}_{\tilde{\mathcal{B}}}$ per CG iteration to compute the debiased update magnitude (details in [Appendix A.2](#)). Debiased LA requires an additional K-FAC approximation $\mathbf{K}_{\tilde{\mathcal{B}}}$; all subsequent operations to compute the debiased K-FAC can be carried out efficiently at Kronecker factor level (details in [Appendix A.3.3](#)). In [Section 6](#), we thus use the debiased versions at *half* the batch size, for a fair comparison. We will see that, at the resulting *similar* computational cost, the debiased versions clearly outperform the single-batch alternatives.

5 RELATED WORK

We here briefly list other, related forms of bias correction that have been suggested elsewhere.

A different two-batch approach. Other works have proposed to use different (not necessarily disjoint) mini-batches for the gradient and the Hessian ([Martens, 2010](#); [Byrd et al., 2011](#)). [Benzing \(2022, Sec. I.5\)](#) mentions the idea of using *independent* mini-batches for the gradient and the Hessian to obtain an, in some sense, unbiased estimate $-\mathbf{H}_{\mathcal{B}}^{-1} \mathbf{g}_{\tilde{\mathcal{B}}}$ of the exact Newton step. This does, however, not resolve the biases described in this paper. Via the eigendecomposition of the Hessian $\mathbf{H}_{\mathcal{B}} = \sum_{p=1}^P \lambda_p \mathbf{u}_p \mathbf{u}_p^{\top}$, we obtain $-\mathbf{H}_{\mathcal{B}}^{-1} \mathbf{g}_{\tilde{\mathcal{B}}} = -\sum_{p=1}^P \partial_{\mathbf{u}_p} q(\theta_{\star}; \tilde{\mathcal{B}}) (\partial_{\mathbf{u}_p}^2 q(\theta_{\star}; \mathcal{B}))^{-1} \mathbf{u}_p$. While the numerator yields an unbiased estimate of the directional slope (similar to a *blue* dot in the upper panel of [Figure 2](#)), the bias in the denominator remains since eigenvectors and directional curvatures are based on the same mini-batch (as for the *orange* dots in the bottom panel of [Figure 2](#)).

Another debiasing approach. EK-FAC (Eigenvalue-corrected Kronecker Factorization) ([George et al., 2018](#)) corrects the eigenvalues of the K-FAC approximation (similar to [Equation \(9\)](#)) which, provably, yields a more accurate approximation of the FIM than K-FAC (in Frobenius norm). This

correction, however, is designed to resolve a different kind of bias that is specific to the K-FAC approximation and does not address the biases described in this work.

Running averages. Other popular deep learning optimizers *aggregate* curvature information over multiple steps via exponential moving averages (see e.g. (Martens & Grosse, 2015)). This emulates larger mini-batch sizes and thus reduces the curvature biases. However, when curvature evolves rapidly, obsolete curvature estimates from past steps might slow down training. Aggregating more robust *debiased* curvature estimates instead might allow for shorter moving average windows and accelerate the optimization progress. We leave it to future work to explore these interactions.

6 EXPERIMENTS

In this section, we evaluate the effectiveness of the debiasing strategies from Section 4.2. In Appendix B, we provide the experimental details as well as additional empirical analyses. For instance, we show how the curvature biases with K-FAC depend on the mini-batch size (see Appendix B.8), how the biases evolve over the course of training (see Appendix B.9), and that the curvature biases become more pronounced for deeper/wider models (see Appendix B.10).

6.1 DEBIASED CONJUGATE GRADIENTS

We compare the standard single-batch CG method to the debiased version (see Section 4.2.1) on the fully trained ALL-CNN-C model without curvature damping. For a fair comparison, the single-batch approach uses one mini-batch of size 1024 while the debiased approach uses two mini-batches of size 512, such that a similar amount of data and runtime budget is used. Details in Appendix B.5.

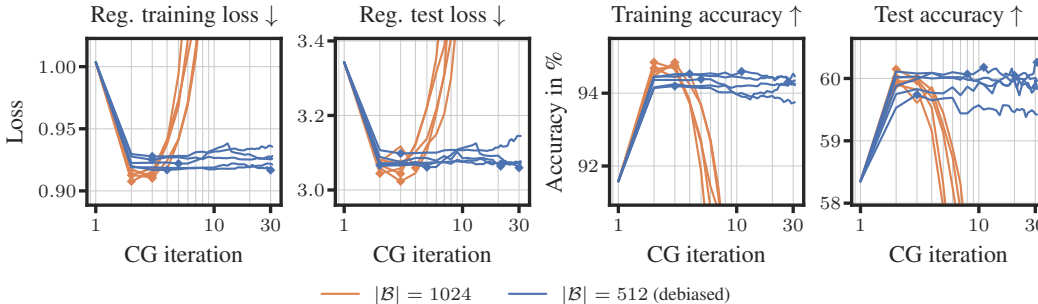


Figure 5: Debiased CG is much more stable than the single-batch approach. We compare CG runs without curvature damping ($\delta = 0$) with $K = 30$ iterations for the fully trained ALL-CNN-C model on the CIFAR-100 dataset in terms of training/test loss/accuracy at similar computational cost: The single-batch approach (shown as ---) uses one mini-batch of size 1024 while the debiased approach (shown as ---) uses two mini-batches of size 512 each. Both approaches use the GGN curvature proxy and are run 5 times on different mini-batches. The markers \blacklozenge and \blacklozenge are placed at peak performance. While the single-batch runs diverge quickly, the debiased CG runs are stable.

Results & discussion. All CG runs in Figure 5 achieve a significant improvement in all four performance metrics. The most striking difference between the two approaches is their stability: The single-batch runs quickly reach peak performance and then diverge. In contrast, although the steps tend to be *larger* (see Figure 4), the debiased CG runs are much more stable (without using any curvature damping). This suggests that the reason for the divergence in the single-batch approach is *not* the missing damping but the misinformed update magnitudes. The peak performance is slightly better for the single-batch runs which can be attributed to its more informative search directions (they were computed using *double* the data). The peak performance of the debiased runs could likely be improved (at the same computational cost) by using a larger batch size for the directions and a smaller one for the update magnitudes (as the former seem harder to estimate, see Section 3.2).

6.2 DEBIASED LAPLACE APPROXIMATION

Here, we use a fully trained ALL-CNN-C model on CIFAR-10 data and compare (i) the vanilla model without LA, (ii) the single-batch K-FAC LA approach, (iii) the debiased version (see Section 4.2.2), and (iv) the full-batch approach (where we compute K-FAC on the entire training set) in terms of accuracy, NLL and ECE. Again, we apply the debiased approach at half the batch size of the single-batch approach for a fair comparison. We use prior precisions between 10^{-4} and 10. Appendix B.6 contains the experimental details and additional results on the training and OOD data.

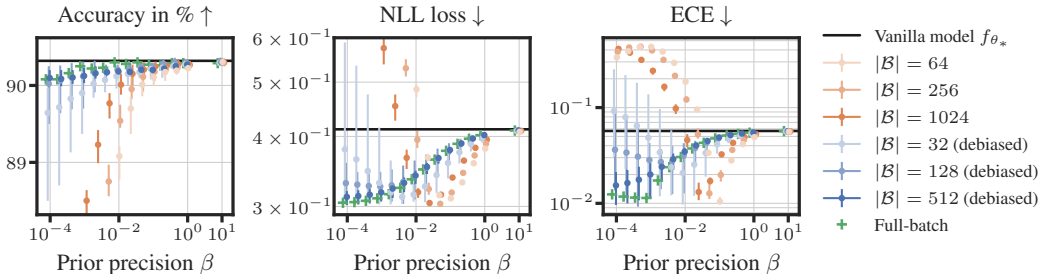


Figure 6: Debiased LA mimics the full-batch LA very well. We compare LAs for the fully trained ALL-CNN-C model on CIFAR-10 in terms of accuracy, negative log likelihood loss (NLL) and expected calibration error (ECE) on the CIFAR-10 test set. For each mini-batch size (lighter color indicating smaller batch size), we draw 5 mini-batches and report the mean performance as dot and min/max as vertical line. We also report the performance of the vanilla model without LA (shown as —) and the full-batch approach based on $K_{\mathcal{D}}$ (shown as +). In contrast to the single-batch approach, the debiased version mimics the behavior of the full-batch approach very well over the entire range of prior precisions.

Results & discussion. Figure 6 shows the results. If the prior precision is low (i.e. the LA relies mainly on the curvature information without the regularizing diagonal term), the single-batch version acts erratically due to the deformed curvature model—its performance drops drastically. The single-batch approach is also more sensitive to the choice of prior precision (in fact, it suggests a much larger prior precision than the full-batch approach). In contrast, the debiased approach achieves good performance over a wider range of prior precisions and mimics the full-batch approach ($N = 40,000$) very well despite using only a tiny fraction of the data.

In order to showcase our debiasing strategy’s scalability, we provide additional results on RESNET-50 and ViT LITTLE on the IMAGENET dataset in Appendix B.6.

Summary of the experimental results. The use of mini-batch quadratics is a simple way to keep the costs of second-order optimization and uncertainty quantification manageable. The resulting biases, however, severely degrade their value, requiring large mini-batches or algorithmic add-ons. Our experiments suggest that even simple debiasing strategies can largely mitigate this issue.

7 CONCLUSION

Our main takeaway is a general principle: **Quadratic approximations to the training loss computed on mini-batches of the training data provide a distorted representation of the true underlying loss landscape.** In particular, along the directions of large curvature, the mini-batch quadratics tend to strongly overestimate the curvature of the true loss. Our theoretical analysis shows that these biases can be traced back to the misalignment of the curvature matrices’ eigenspaces. These insights are highly relevant for applications: As we demonstrated empirically, the biases in the directional slope and curvature lead to severely misinformed updates in stochastic second-order optimizers, and cause unreliable uncertainty estimates with Laplace approximations. We also proposed simple two-batch strategies to mitigate these biases. Our experiments demonstrate their superiority over the single-batch approaches in terms of stability and quality at similar computational costs. Our findings reveal a design prerequisite for building better stochastic curvature-based methods, which should be generally considered, and further developed, for all methods using such curvature evaluations.

REFERENCES

- 540
541
542 Shun-ichi Amari. Natural gradient works efficiently in learning. *Neural Computation*, 1998.
- 543
544 Frederik Benzing. Gradient descent on neurons and its link to approximate second-order optimization.
545 In *International Conference on Machine Learning (ICML)*, 2022.
- 546 Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical Gauss-Newton optimisation for deep
547 learning. In *Proceedings of Machine Learning Research (PMLR)*, 2017.
- 548
549 Richard H. Byrd, Gillian M. Chin, Will Neveitt, and Jorge Nocedal. On the use of stochastic Hessian
550 information in optimization methods for machine learning. *SIAM Journal on Optimization*, 2011.
- 551
552 Felix Dangel, Frederik Kunstner, and Philipp Hennig. BackPACK: Packing more into backprop. In
553 *International Conference on Learning Representations (ICLR)*, 2020.
- 554
555 Felix Dangel, Lukas Tatzel, and Philipp Hennig. ViViT: Curvature access through the generalized
556 gauss-newton’s low-rank structure. *Transactions on Machine Learning Research (TMLR)*, 2022.
- 557
558 Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and
559 Philipp Hennig. Laplace redux—effortless Bayesian deep learning. In *Advances in Neural Informa-
tion Processing Systems (NeurIPS)*, 2021.
- 560
561 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
562 hierarchical image database. In *Conference on Computer Vision and Pattern Recognition (CVPR)*,
2009.
- 563
564 Nicki Detlefsen, Jiří Borovec, Justus Schock, Ananya Jha, Teddy Koker, Luca Di Liello, Daniel
565 Štancl, Quan Changsheng, Maxim Grechkin, and William Falcon. TorchMetrics - measuring
566 reproducibility in PyTorch. *The Journal of Open Source Software*, 2022.
- 567
568 Xin Dong, Shangyu Chen, and Sinno Pan. Learning to prune deep neural networks via layer-wise
569 optimal brain surgeon. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- 570
571 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
572 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
573 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
574 In *International Conference on Learning Representations*, 2021. URL [https://openreview.
net/forum?id=YicbFdNTTy](https://openreview.net/forum?id=YicbFdNTTy).
- 575
576 Runa Eschenhagen, Alexander Immer, Richard E. Turner, Frank Schneider, and Philipp Hennig.
577 Kronecker-Factored Approximate Curvature for modern neural network architectures. In *Advances
in Neural Information Processing Systems (NeurIPS)*, 2023.
- 578
579 Francis Galton. Regression towards mediocrity in hereditary stature. *The Journal of the Anthropolog-
ical Institute of Great Britain and Ireland*, 1886.
- 580
581 Thomas George, César Laurent, Xavier Bouthillier, Nicolas Ballas, and Pascal Vincent. Fast
582 approximate natural gradient descent in a kronecker factored eigenbasis. In *Advances in Neural
583 Information Processing Systems (NeurIPS)*, 2018.
- 584
585 Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization
586 via hessian eigenvalue density. In *Proceedings of Machine Learning Research (PMLR)*, 2019.
- 587
588 Roger Grosse and James Martens. A Kronecker-factored approximate Fisher matrix for convolution
589 layers. In *International Conference on Machine Learning (ICML)*, 2016.
- 590
591 Guy Gur-Ari, Daniel A. Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace. *arXiv*,
2018.
- 592
593 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

- 594 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
595 corruptions and perturbations. In *International Conference on Learning Representations (ICLR)*,
596 2019.
- 597 Magnus Rudolph Hestenes and Eduard Stiefel. Methods of conjugate gradients for solving linear
598 systems. *Journal of Research of the National Bureau of Standards*, 1952.
- 600 Alexander Immer, Maciej Korzepa, and Matthias Bauer. Improving predictions of Bayesian neural
601 nets via local linearization. In *International Conference on Artificial Intelligence and Statistics*
602 (*AISTATS*), 2021.
- 603 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
604 *Proceedings of Machine Learning Research (PMLR)*, 2017.
- 606 Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being Bayesian, even just a bit, fixes
607 overconfidence in ReLU networks. In *International Conference on Machine Learning (ICML)*,
608 2020.
- 609 Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009.
- 611 Dong C. Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization.
612 *Mathematical Programming*, 1989.
- 613 David J. C. MacKay. Bayesian model comparison and backprop nets. In *Advances in Neural*
614 *Information Processing Systems (NeurIPS)*, 1991.
- 616 David J. C. MacKay. The evidence framework applied to classification networks. *Neural Computation*,
617 1992.
- 618 James Martens. Deep learning via Hessian-free optimization. In *International Conference on Machine*
619 *Learning (ICML)*, 2010.
- 621 James Martens. New insights and perspectives on the natural gradient method. *arXiv*, 2020.
- 622 James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored approximate
623 curvature. In *International Conference on Machine Learning (ICML)*, 2015.
- 625 James Martens, Jimmy Ba, and Matt Johnson. Kronecker-factored curvature approximations for
626 recurrent neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- 627 Jorge Nocedal. Updating quasi-newton matrices with limited storage. *Mathematics of Computation*,
628 1980.
- 630 Jorge Nocedal and Stephen J. Wright. *Numerical optimization*. Springer, 2. ed. edition, 2006.
- 631 Kazuki Osawa, Yohei Tsuji, Yuichiro Ueno, Akira Naruse, Rio Yokota, and Satoshi Matsuoka. Large-
632 scale distributed second-order optimization using Kronecker-factored approximate curvature for
633 deep convolutional neural networks. In *Conference on Computer Vision and Pattern Recognition*
634 (*CVPR*), 2019.
- 636 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
637 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward
638 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,
639 Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep
640 learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- 641 Barak A. Pearlmutter. Fast exact multiplication by the Hessian. *Neural Computation*, 1994.
- 642 Hippolyt Ritter, Aleksandar Botev, and David Barber. Online structured laplace approximations
643 for overcoming catastrophic forgetting. In *Advances in Neural Information Processing Systems*
644 (*NeurIPS*), 2018a.
- 645 Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural
646 networks. In *International Conference on Learning Representations (ICLR)*, 2018b.

648 Hrittik Roy, Marco Miani, Carl Henrik Ek, Philipp Hennig, Marvin förtner, Lukas Tatzel, and Søren
649 Hauberg. Reparameterization invariance in approximate Bayesian inference. *arXiv*, 2024.
650

651 Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley
652 - benchmarking deep learning optimizers. In *International Conference on Machine Learning*
653 (*ICML*), 2021.

654 Frank Schneider, Lukas Balles, and Philipp Hennig. DeepOBS: A deep learning optimizer benchmark
655 suite. In *International Conference on Learning Representations (ICLR)*, 2019.
656

657 Nicol N. Schraudolph. Fast curvature matrix-vector products for second-order gradient descent.
658 *Neural Computation*, 2002.

659 Jost T. Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for
660 simplicity: The all convolutional net. *arXiv*, 2015.
661

662 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv*, 2017.

663 Wenyuan Zeng and Raquel Urtasun. MLPrune: Multi-layer pruning for automated neural network
664 compression, 2019.
665

666 Huishuai Zhang, Caiming Xiong, James Bradbury, and Richard Socher. Block-diagonal Hessian-free
667 optimization for training neural networks. *arXiv*, 2017.
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

SUPPLEMENTAL MATERIAL

Below, we provide additional details on the mathematical derivations, describe the experimental setup and present additional results.

A	Mathematical details	15
A.1	Directional derivatives of a quadratic	15
A.2	The method of conjugate gradients (CG)	16
A.3	Laplace approximation for neural networks	17
A.3.1	Derivation of the Laplace approximation for neural networks	17
A.3.2	Sampling from the K-FAC Laplace approximation	19
A.3.3	Debiased K-FAC Laplace approximation	20
A.4	Bias in the directional slope	21
A.5	Bias in the directional curvature	22
B	Experimental details	23
B.1	Test problems and training procedures	24
B.2	Matrix-vector products with the curvature matrix	25
B.3	Section 1 (Introduction): Figure 1	25
B.4	Section 3 (The shape of a mini-batch quadratic): Figures 2 to 4	25
B.5	Section 6.1 (Debiased conjugate gradients): Figure 5	27
B.6	Section 6.2 (Debiased Laplace approximation): Figure 6	32
B.7	Additional experiment: Biases for a WIDERESNET 40-4 model architecture	38
B.8	Additional experiment: K-FAC and the dependence of the biases on mini-batch size	38
B.9	Additional experiment: Development of the biases over the course of training	38
B.10	Additional experiment: The curvature bias increases with P	40

756 A MATHEMATICAL DETAILS

757 A.1 DIRECTIONAL DERIVATIVES OF A QUADRATIC

758 We claim in Section 2.1 that a cut r through the quadratic $q(\cdot; \mathcal{B})$ from $\theta_\bullet \in \mathbb{R}^P$ along the normalized
759 direction \mathbf{d} can be written as

$$760 \quad r(\tau) := q(\theta_\bullet + \tau \mathbf{d}; \mathcal{B}) = \frac{1}{2} \tau^2 \mathbf{d}^\top \nabla^2 q(\theta_\bullet; \mathcal{B}) \mathbf{d} + \tau \mathbf{d}^\top \nabla q(\theta_\bullet; \mathcal{B}) + \text{const.} \quad (10)$$

761 **Proof for Equation (10).** Here, we provide the derivation for Equation (10). Let $\theta_\bullet \in \mathbb{R}^P$ be a point
762 in parameter space and $\mathbf{d} \in \mathbb{R}^P$ a normalized direction, i.e. $\|\mathbf{d}\| = 1$. We consider the quadratic
763 approximation $q(\theta; \mathcal{B})$ around θ_0 and evaluate it along the cut $\theta_\bullet + \tau \mathbf{d}$ for $\tau \in \mathbb{R}$. We assume that
764 the Hessian (or its approximation) $\mathbf{H}_\mathcal{B}$ is symmetric and obtain

$$765 \quad \begin{aligned} r(\tau) &:= q(\theta_\bullet + \tau \mathbf{d}; \mathcal{B}) \\ &= \frac{1}{2} (\theta_\bullet + \tau \mathbf{d} - \theta_0)^\top \mathbf{H}_\mathcal{B} (\theta_\bullet + \tau \mathbf{d} - \theta_0) + (\theta_\bullet + \tau \mathbf{d} - \theta_0)^\top \mathbf{g}_\mathcal{B} + c_\mathcal{B} \\ &= \frac{1}{2} \tau^2 \mathbf{d}^\top \mathbf{H}_\mathcal{B} \mathbf{d} + \tau \mathbf{d}^\top \mathbf{H}_\mathcal{B} (\theta_\bullet - \theta_0) + \frac{1}{2} (\theta_\bullet - \theta_0)^\top \mathbf{H}_\mathcal{B} (\theta_\bullet - \theta_0) \\ &\quad + \tau \mathbf{d}^\top \mathbf{g}_\mathcal{B} + (\theta_\bullet - \theta_0)^\top \mathbf{g}_\mathcal{B} + c_\mathcal{B} \\ &= \frac{1}{2} \tau^2 \mathbf{d}^\top \mathbf{H}_\mathcal{B} \mathbf{d} + \tau \mathbf{d}^\top (\mathbf{H}_\mathcal{B} (\theta_\bullet - \theta_0) + \mathbf{g}_\mathcal{B}) + \text{const.} \end{aligned}$$

766 Since $\nabla q(\theta_\bullet; \mathcal{B}) = \mathbf{H}_\mathcal{B} (\theta_\bullet - \theta_0) + \mathbf{g}_\mathcal{B}$ and $\nabla^2 q(\theta_\bullet; \mathcal{B}) \equiv \mathbf{H}_\mathcal{B}$, we arrive at Equation (10). \square

767 **Proof for Equation (5).** Next, we show that the average directional slope/curvature over all mini-
768 batches in the training set coincides with the directional slope/curvature of the full-batch quadratic.
769 Let's assume that all M mini-batches $\{\mathcal{B}_{m'}\}_{m'=1}^M$ are disjoint, have the same size $|\mathcal{B}_1|$ and that their
770 union is the training set, i.e.

$$771 \quad |\mathcal{B}_{m'}| = |\mathcal{B}_1| \quad \forall m' \in \{1, \dots, M\} \quad \text{and} \quad \bigcup_{m'=1}^M \mathcal{B}_{m'} = \mathcal{D}. \quad (11)$$

772 This implies $M|\mathcal{B}_1| = |\mathcal{D}|$. It holds

$$773 \quad \begin{aligned} \frac{1}{M} \sum_{m'=1}^M \underbrace{\partial_{\mathbf{u}_p} q(\theta_\star; \mathcal{B}_{m'})}_{\text{one } \bullet \text{ and many } \bullet} &= \frac{1}{M} \sum_{m'=1}^M \mathbf{u}_p^\top \nabla q(\theta_\star; \mathcal{B}_{m'}) \quad \text{with } \nabla q(\theta_\star; \mathcal{B}_{m'}) = \mathbf{g}_{\mathcal{B}_{m'}} \text{ since } \theta_0 \leftarrow \theta_\star \\ &= \mathbf{u}_p^\top \frac{1}{M} \sum_{m'=1}^M \nabla \mathcal{L}_{\text{reg}}(\theta_\star; \mathcal{B}_{m'}) \\ &= \mathbf{u}_p^\top \left(\frac{1}{M} M \nabla r(\theta_\star) + \frac{1}{M} \sum_{m'=1}^M \nabla \mathcal{L}(\theta_\star; \mathcal{B}_{m'}) \right) \\ &= \mathbf{u}_p^\top \left(\nabla r(\theta_\star) + \frac{1}{M} \sum_{m'=1}^M \frac{1}{|\mathcal{B}_{m'}|} \sum_{n \in \mathcal{B}_{m'}} \nabla \ell(f_{\theta_\star}(\mathbf{x}_n), \mathbf{y}_n) \right) \\ &\stackrel{(11)}{=} \mathbf{u}_p^\top \left(\nabla r(\theta_\star) + \frac{1}{M|\mathcal{B}_1|} \sum_{m'=1}^M \sum_{n \in \mathcal{B}_{m'}} \nabla \ell(f_{\theta_\star}(\mathbf{x}_n), \mathbf{y}_n) \right) \\ &= \mathbf{u}_p^\top \nabla \mathcal{L}_{\text{reg}}(\theta_\star; \mathcal{D}) \\ &= \underbrace{\partial_{\mathbf{u}_p} q(\theta_\star; \mathcal{D})}_{+} \end{aligned}$$

800 As the mini-batch curvature $\nabla^2 \mathcal{L}(\theta_\star; \mathcal{B}_{m'})$ is also an *average* over the samples in the mini-batch
801 (this applies to both the actual Hessian and its GGN approximation), the same argument holds for the
802 associated directional curvatures. For K-FAC, the derivation above does *not* hold as $(1/2)(\mathbf{K}_\mathcal{B} +$
803 $\mathbf{K}_{\bar{\mathcal{B}}}) \neq \mathbf{K}_{\mathcal{B} \cup \bar{\mathcal{B}}}$. \square

A.2 THE METHOD OF CONJUGATE GRADIENTS (CG)

Using CG for minimizing a quadratic. The argmin of a quadratic $q(\cdot; \mathcal{B})$ is given by $\theta_0 - \mathbf{H}_{\mathcal{B}}^{-1} \mathbf{g}_{\mathcal{B}}$ (assuming that $\mathbf{H}_{\mathcal{B}}$ is symmetric and positive definite), see Section 2.2. We can use the method of conjugate gradients (see Algorithm 1) for computing the Newton step (or an approximation thereof) by setting $\mathbf{A} \leftarrow \mathbf{H}_{\mathcal{B}}$ and $\mathbf{b} \leftarrow -\mathbf{g}_{\mathcal{B}}$.

Algorithm 1: Method of conjugate gradients (CG), based on (Nocedal & Wright, 2006, Alg. 5.2)

Input. Access to matrix-vector products $\mathbf{v} \mapsto \mathbf{A}\mathbf{v}$, where $\mathbf{A} \in \mathbb{R}^{P \times P}$ is symmetric and positive definite, right-hand side $\mathbf{b} \in \mathbb{R}^P$, convergence tolerance $\epsilon \in \mathbb{R}_{>0}$, maximum number of iterations $P_{\max} \in \mathbb{N}$, $P_{\max} \leq P$

Output. Approximate solution \mathbf{x}_p to the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$

```

1 Initialize  $\mathbf{x}_0 = \mathbf{0}$ ,  $\mathbf{r}_0 \leftarrow -\mathbf{b}$  and  $\mathbf{s}_0 \leftarrow -\mathbf{r}_0$ 
2 for  $p = 0, 1, \dots, P$  do
3   if  $p = P_{\max}$  or  $\|\mathbf{r}_p\|_2 \leq \epsilon$  then Termination criteria
4     return (approximate) solution  $\mathbf{x}_p$ 
5    $\mathbf{t}_p \leftarrow \mathbf{A}\mathbf{s}_p$  Compute the matrix-vector product only once
6    $\alpha_p \leftarrow \frac{\mathbf{r}_p^\top \mathbf{r}_p}{\mathbf{s}_p^\top \mathbf{t}_p}$  Update magnitude  $\alpha_p = \frac{\mathbf{r}_p^\top \mathbf{r}_p}{\mathbf{s}_p^\top \mathbf{A}\mathbf{s}_p}$ 
7    $\mathbf{x}_{p+1} \leftarrow \mathbf{x}_p + \alpha_p \mathbf{s}_p$  Update along search direction  $\mathbf{s}_p$ 
8    $\mathbf{r}_{p+1} \leftarrow \mathbf{r}_p + \alpha_p \mathbf{t}_p$  Residual  $\mathbf{r}_p = \mathbf{A}\mathbf{x}_p - \mathbf{b}$  computed via recursion
9    $\beta_{p+1} \leftarrow \frac{\mathbf{r}_{p+1}^\top \mathbf{r}_{p+1}}{\mathbf{r}_p^\top \mathbf{r}_p}$ 
10   $\mathbf{s}_{p+1} \leftarrow -\mathbf{r}_{p+1} + \beta_{p+1} \mathbf{s}_p$  Construction of new (conjugate) search direction

```

Properties of CG & geometric interpretations. In the following, we provide some important properties and a geometric interpretation of the quantities involved in the CG method.

- **Residual \mathbf{r}_p .** Let $\mathbf{x}_p := \theta_p - \theta_0$. The residual \mathbf{r}_p can be written as $\mathbf{r}_p = \mathbf{A}\mathbf{x}_p - \mathbf{b} = \mathbf{H}_{\mathcal{B}}(\theta_p - \theta_0) + \mathbf{g}_{\mathcal{B}} = \nabla q(\theta_p; \mathcal{B})$, i.e. it coincides with the gradient of the quadratic at θ_p .
- **Update magnitude α_p .** Consider a cut through the quadratic from θ_p into the direction \mathbf{s}_p , i.e. $r(\tau) = q(\theta_p + \tau \mathbf{s}_p; \mathcal{B})$. Minimizing this 1D quadratic requires its first derivative to vanish. It holds (see Equation (10))

$$\begin{aligned}
r'(\tau) = 0 &\Leftrightarrow \tau \mathbf{s}_p^\top \nabla^2 q(\theta_p; \mathcal{B}) \mathbf{s}_p + \mathbf{s}_p^\top \nabla q(\theta_p; \mathcal{B}) = 0 \\
&\Leftrightarrow \tau = -\frac{\mathbf{s}_p^\top \nabla q(\theta_p; \mathcal{B})}{\mathbf{s}_p^\top \nabla^2 q(\theta_p; \mathcal{B}) \mathbf{s}_p} = -\frac{\mathbf{s}_p^\top \mathbf{r}_p}{\mathbf{s}_p^\top \mathbf{A} \mathbf{s}_p} = \frac{\mathbf{r}_p^\top \mathbf{r}_p}{\mathbf{s}_p^\top \mathbf{A} \mathbf{s}_p} = \alpha_p,
\end{aligned}$$

where the equality $-\mathbf{s}_p^\top \mathbf{r}_p = \mathbf{r}_p^\top \mathbf{r}_p$ is due to Equation (5.14a) and Theorem 5.2 in (Nocedal & Wright, 2006). So, the update magnitude α_p is chosen such that it minimizes the quadratic along the direction \mathbf{s}_p . Note that the update magnitude can also be written as a ratio of the negative directional slope and directional curvature at θ_p , i.e. a 1D directional Newton step. Let $\mathbf{d}_p := \mathbf{s}_p / \|\mathbf{s}_p\|$ denote the *normalized* search direction. It holds

$$\mathbf{x}_{p+1} = \mathbf{x}_p + \alpha_p \mathbf{s}_p = \mathbf{x}_p - \frac{\mathbf{s}_p^\top \mathbf{r}_p}{\mathbf{s}_p^\top \mathbf{A} \mathbf{s}_p} \mathbf{s}_p = \mathbf{x}_p - \frac{\mathbf{d}_p^\top \mathbf{r}_p}{\mathbf{d}_p^\top \mathbf{A} \mathbf{d}_p} \mathbf{d}_p = \mathbf{x}_p - \underbrace{\frac{\partial_{\mathbf{d}_p} q(\theta_d; \mathcal{B})}{\partial_{\mathbf{d}_p}^2 q(\theta_d; \mathcal{B})}}_{=: \tau_p} \mathbf{d}_p$$

i.e. $\tau_p = \alpha_p \|\mathbf{s}_p\|$. Via the shift $\theta_p = \theta_0 + \mathbf{x}_p$, we arrive at Equation (3).

- **Search direction \mathbf{s}_p .** CG constructs the search directions to be conjugate, i.e. $\mathbf{s}_{p+1}^\top \mathbf{A} \mathbf{s}_i = 0$ for all $i \leq p$. Note that this property also applies to the normalized search directions.

Efficient implementation of the debiased CG approach. For the debiased CG version, we need to re-evaluate the update magnitudes $\tilde{\tau}_p$ for a given set of search directions $\{\mathbf{d}_1, \dots, \mathbf{d}_P\}$ on a second mini-batch $\tilde{\mathcal{B}}$ (see Equation (8)), i.e.

$$\tilde{\tau}_p = -\frac{\partial_{\mathbf{d}_p} q(\tilde{\boldsymbol{\theta}}_p; \tilde{\mathcal{B}})}{\partial_{\mathbf{d}_p}^2 q(\tilde{\boldsymbol{\theta}}_p; \tilde{\mathcal{B}})} = -\frac{\mathbf{d}_p^\top \nabla q(\tilde{\boldsymbol{\theta}}_p; \tilde{\mathcal{B}})}{\mathbf{d}_p^\top \nabla^2 q(\tilde{\boldsymbol{\theta}}_p; \tilde{\mathcal{B}}) \mathbf{d}_p} = -\frac{\mathbf{d}_p^\top (\mathbf{H}_{\tilde{\mathcal{B}}}(\tilde{\boldsymbol{\theta}}_p - \boldsymbol{\theta}_0) + \mathbf{g}_{\tilde{\mathcal{B}}})}{\mathbf{d}_p^\top \mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{d}_p}$$

Implemented naively, this requires *two* matrix-vector products—one for the numerator and one for the denominator. However, we can use a recursive formula for the numerator:

$$\begin{aligned} \nabla q(\tilde{\boldsymbol{\theta}}_p; \tilde{\mathcal{B}}) &= \mathbf{H}_{\tilde{\mathcal{B}}}(\tilde{\boldsymbol{\theta}}_p - \boldsymbol{\theta}_0) + \mathbf{g}_{\tilde{\mathcal{B}}} \\ &= \mathbf{H}_{\tilde{\mathcal{B}}}(\tilde{\boldsymbol{\theta}}_{p-1} + \tilde{\tau}_{p-1} \mathbf{d}_{p-1} - \boldsymbol{\theta}_0) + \mathbf{g}_{\tilde{\mathcal{B}}} \\ &= \mathbf{H}_{\tilde{\mathcal{B}}}(\tilde{\boldsymbol{\theta}}_{p-1} - \boldsymbol{\theta}_0) + \mathbf{g}_{\tilde{\mathcal{B}}} + \tilde{\tau}_{p-1} \mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{d}_{p-1} \\ &= \nabla q(\tilde{\boldsymbol{\theta}}_{p-1}; \tilde{\mathcal{B}}) + \tilde{\tau}_{p-1} \mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{d}_{p-1}. \end{aligned} \quad (12)$$

So, if we store the gradient from the previous iteration $\nabla q(\tilde{\boldsymbol{\theta}}_{p-1}; \tilde{\mathcal{B}})$ and the Hessian vector product with the previous direction $\mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{d}_{p-1}$, the current gradient can be computed *without* an additional matrix-vector product. At iteration p , we thus (i) compute the gradient $\nabla q(\tilde{\boldsymbol{\theta}}_p; \tilde{\mathcal{B}})$ recursively from the cached vectors $\nabla q(\tilde{\boldsymbol{\theta}}_{p-1}; \tilde{\mathcal{B}})$ and $\mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{d}_{p-1}$ via Equation (12) (for $p = 0$, we have $\nabla q(\tilde{\boldsymbol{\theta}}_0; \tilde{\mathcal{B}}) = \mathbf{g}_{\tilde{\mathcal{B}}}$), (ii) compute the Hessian-vector product with the current direction $\mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{d}_p$, (iii) store both these vectors and (iv) compute the update magnitude $\tilde{\tau}_p$ (both the numerator and the denominator only require a simple dot product of two pre-computed vectors).

A.3 LAPLACE APPROXIMATION FOR NEURAL NETWORKS

A.3.1 DERIVATION OF THE LAPLACE APPROXIMATION FOR NEURAL NETWORKS

Preliminaries. The softmax function $\text{softmax} : \mathbb{R}^C \rightarrow \mathbb{R}^C$ is defined as

$$\text{softmax}(\mathbf{z}) := \left(\frac{\exp(z_1)}{\sum_{c'=1}^C \exp(z_{c'})}, \dots, \frac{\exp(z_C)}{\sum_{c'=1}^C \exp(z_{c'})} \right).$$

It maps an arbitrary vector $\mathbf{z} \in \mathbb{R}^C$ to a vector whose entries are non-negative and sum up to one. So, the output of the softmax function can be interpreted as a probability distribution over C classes. With this, the cross-entropy loss for a single datum (\mathbf{x}, \mathbf{y}) is given by

$$\ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) := -\sum_{c=1}^C \mathbf{y}_c \cdot \log(\text{softmax}(f_{\boldsymbol{\theta}}(\mathbf{x}))_c). \quad (13)$$

Here, we assume that $\mathbf{y} \in \{0, 1\}^C$ is a one-hot encoded vector representing the true class label. So, if $c_\star \in \{1, \dots, C\}$ is the correct class (i.e. $\mathbf{y}_{c_\star} = 1$), the cross-entropy loss is given by the negative logarithm of the probability the network assigns to this class $\ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) = -\log(\text{softmax}(f_{\boldsymbol{\theta}}(\mathbf{x}))_{c_\star})$. Finally, let

$$\text{Cat}(\mathbf{y}; \mathbf{p}) := \prod_{c=1}^C p_c^{\mathbf{y}_c} \quad (14)$$

denote the probability mass function of the categorical distribution, where $\mathbf{p} \in \mathbb{R}^C$ is a vector of probabilities (i.e. $p_c \geq 0$ and $\sum_c p_c = 1$), and $\mathbf{y} \in \{0, 1\}^C$ is a one-hot encoded vector representing a class label.

Probabilistic interpretation of the regularized loss. Recall from Equation (1) that the regularized loss function is given by

$$\mathcal{L}_{\text{reg}}(\boldsymbol{\theta}; \mathcal{D}) = \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) + r(\boldsymbol{\theta}) \quad \text{with} \quad \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) = \frac{1}{N} \sum_{n \in \mathcal{D}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_n), \mathbf{y}_n).$$

We assume a classification problem that uses the cross-entropy loss (see Equation (13)) and an ℓ^2 regularizer $r(\boldsymbol{\theta}) = \beta/2 \|\boldsymbol{\theta}\|_2^2$ with parameter $\beta \in \mathbb{R}_{>0}$. We use one-hot encoded labels $\mathbf{y}_n \in \{0, 1\}^C$ (with \mathbf{y}_{nc} , we denote the c -th entry of \mathbf{y}_n) and obtain

$$N \cdot \mathcal{L}_{\text{reg}}(\boldsymbol{\theta}; \mathcal{D}) = \sum_{n \in \mathcal{D}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_n), \mathbf{y}_n) + \frac{N\beta}{2} \|\boldsymbol{\theta}\|_2^2$$

$$\begin{aligned}
&\stackrel{(13)}{=} - \sum_{n \in \mathcal{D}} \sum_{c=1}^C \mathbf{y}_{nc} \cdot \log(\text{softmax}(f_{\boldsymbol{\theta}}(\mathbf{x}_n))_c) + \frac{N\beta}{2} \|\boldsymbol{\theta}\|_2^2 \\
&\stackrel{(14)}{=} - \sum_{n \in \mathcal{D}} \log(\text{Cat}(\mathbf{y}_n; \text{softmax}(f_{\boldsymbol{\theta}}(\mathbf{x}_n)))) - \left(-\frac{1}{2} \boldsymbol{\theta}^\top (N\beta \cdot \mathbf{I}) \boldsymbol{\theta}\right) \\
&= - \log \left(\underbrace{\prod_{n \in \mathcal{D}} \text{Cat}(\mathbf{y}_n; \text{softmax}(f_{\boldsymbol{\theta}}(\mathbf{x}_n)))}_{\text{likelihood } p(\mathbb{D} | \boldsymbol{\theta})} \right) - \log \left(\underbrace{\mathcal{N}\left(\boldsymbol{\theta}; \mathbf{0}, \frac{1}{N\beta} \mathbf{I}\right)}_{\text{prior } p(\boldsymbol{\theta})} \right) - Z,
\end{aligned}$$

where $Z := P/2 \log(2\pi/N\beta)$ absorbs the normalization constant of the Gaussian prior. The cross entropy loss \mathcal{L} for the training set is thus connected to the negative log categorical likelihood and the regularizer can be seen as a negative log Gaussian prior over the parameters. Note that a similar derivation is possible for other loss functions as well, e.g. the MSE loss (which is equivalent to a negative log Gaussian likelihood).

The derivation above shows that the (rescaled) regularized loss function can be interpreted as the negative unnormalized log posterior of a Bayesian model

$$N \cdot \mathcal{L}_{\text{reg}}(\boldsymbol{\theta}; \mathcal{D}) \stackrel{\text{c}}{=} -\log p(\mathbb{D} | \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) \stackrel{\text{c}}{=} -\log p(\boldsymbol{\theta} | \mathbb{D}) \quad (15)$$

with Gaussian prior and categorical likelihood. With $\stackrel{\text{c}}{=}$, we denote equality up to an additive constant.

Training as MAP estimation. Equation (15) allows re-interpreting the training procedure of a neural network as a maximum a posteriori (MAP) estimation problem since minimizing the regularized loss

$$\underbrace{\arg \min_{\boldsymbol{\theta} \in \mathbb{R}^P} \mathcal{L}_{\text{reg}}(\boldsymbol{\theta}, \mathcal{D})}_{= \boldsymbol{\theta}_*} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^P} N \cdot \mathcal{L}_{\text{reg}}(\boldsymbol{\theta}, \mathcal{D}) \stackrel{(15)}{=} \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^P} \log p(\boldsymbol{\theta} | \mathbb{D}) = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^P} p(\boldsymbol{\theta} | \mathbb{D}) \stackrel{=}{=} \boldsymbol{\theta}_*$$

is equivalent to maximizing the posterior $p(\boldsymbol{\theta} | \mathbb{D})$.

Laplace approximation. The idea of the Laplace approximation is to approximate the posterior distribution $p(\boldsymbol{\theta} | \mathbb{D}) \approx \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\theta}_*, \boldsymbol{\Sigma}_{\mathcal{D}})$ with a Gaussian at the mode $\boldsymbol{\theta}_*$ of the posterior, i.e. after training the network. For this, we approximate the log posterior by a second-order Taylor expansion around the mode $\boldsymbol{\theta}_0 \leftarrow \boldsymbol{\theta}_*$, i.e.

$$\log p(\boldsymbol{\theta} | \mathbb{D}) \stackrel{\text{c}}{=} -N \cdot \mathcal{L}_{\text{reg}}(\boldsymbol{\theta}; \mathcal{D}) \stackrel{(2)}{\approx} -N \cdot q(\boldsymbol{\theta}; \mathcal{D}) \stackrel{\text{c}}{=} -\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_*)^\top (N \cdot \mathbf{H}_{\mathcal{D}}) (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \quad (16)$$

where we assumed that $\boldsymbol{\theta}_*$ is a (local) minimum of the regularized loss function (i.e. $\mathbf{g}_{\mathcal{D}} = \nabla \mathcal{L}_{\text{reg}}(\boldsymbol{\theta}_*; \mathcal{D}) = 0$) such that the linear term $(\boldsymbol{\theta} - \boldsymbol{\theta}_*)^\top \mathbf{g}_{\mathcal{D}}$ in the quadratic approximation vanishes. The additive constants we did not state explicitly in Equation (16) turn into some multiplicative factor (denoted by Z^{-1}) when taking the exponential

$$p(\boldsymbol{\theta} | \mathbb{D}) \approx \frac{1}{Z} \exp \left(-\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_*)^\top (N \cdot \mathbf{H}_{\mathcal{D}}) (\boldsymbol{\theta} - \boldsymbol{\theta}_*) \right).$$

This immediately identifies Z as the normalization constant of the Gaussian and we obtain the Laplace approximation

$$p(\boldsymbol{\theta} | \mathbb{D}) \approx \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\theta}_*, \boldsymbol{\Sigma}_{\mathcal{D}}) \quad \text{with} \quad \boldsymbol{\Sigma}_{\mathcal{D}} := (N \cdot \mathbf{H}_{\mathcal{D}})^{-1} = \frac{1}{N} \mathbf{H}_{\mathcal{D}}^{-1}. \quad (17)$$

Mini-batch version of the Laplace approximation. In order to reduce the computational cost of the Laplace approximation, we can replace the full-batch quadratic $q(\boldsymbol{\theta}; \mathcal{D})$ in Equation (16) by a mini-batch quadratic $q(\boldsymbol{\theta}; \mathcal{B})$, i.e. we obtain $\boldsymbol{\Sigma}_{\mathcal{B}} = N^{-1} \mathbf{H}_{\mathcal{B}}^{-1}$.

Predictive uncertainty via the Laplace approximation. Ultimately, we want to use the Laplace approximation to equip the prediction for an unknown test input \mathbf{x}_{\circ} with uncertainty. Ideally, we would compute the expectation of the model likelihood under the approximate posterior, i.e.

$$p(\mathbf{y}_{\circ} | \mathbf{x}_{\circ}, \mathbb{D}) = \int p(\mathbf{y}_{\circ} | \mathbf{x}_{\circ}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbb{D}) d\boldsymbol{\theta}$$

$$\stackrel{(17)}{\approx} \int \text{Cat}(\mathbf{y}_\diamond; \text{softmax}(f_\theta(\mathbf{x}_\diamond))) \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\theta}_*, \boldsymbol{\Sigma}_B) d\boldsymbol{\theta}.$$

One way to approximate this integral is via Monte Carlo sampling, i.e.

$$p(\mathbf{y}_\diamond | \mathbf{x}_\diamond, \mathbb{D}) \approx \frac{1}{S} \sum_{s=1}^S \text{Cat}(\mathbf{y}_\diamond; \text{softmax}(f_{\theta^{(s)}}^{\text{lin}}(\mathbf{x}_\diamond))) \quad \text{where } \boldsymbol{\theta}^{(s)} \sim \mathcal{N}(\boldsymbol{\theta}_*, \boldsymbol{\Sigma}_B). \quad (18)$$

As suggested by Immer et al. (2021); Roy et al. (2024), we use the linearized network $f_{\theta}^{\text{lin}}(\mathbf{x}) := f_{\theta_*}(\mathbf{x}) + \nabla f_{\theta_*}(\mathbf{x})(\boldsymbol{\theta} - \boldsymbol{\theta}_*)$, where $\nabla f_{\theta_*}(\mathbf{x}) \in \mathbb{R}^{C \times P}$ is the network’s Jacobian at $\boldsymbol{\theta}_*$.

A.3.2 SAMPLING FROM THE K-FAC LAPLACE APPROXIMATION

The Monte Carlo (MC) approach from Equation (18) requires samples from the Gaussian $\mathcal{N}(\boldsymbol{\theta}_*, \boldsymbol{\Sigma}_B)$ (the following derivations work exactly the same when the full-batch LA based on the entire training set is used). However, $\boldsymbol{\Sigma}_B \in \mathbb{R}^{P \times P}$ is often too large to be built explicitly in memory and it requires inverting the $P \times P$ Hessian \mathbf{H}_B of the regularized loss function. One approach for that issue is to use the K-FAC curvature approximation $\mathbf{H}_B = \nabla^2 \mathcal{L}(\boldsymbol{\theta}; \mathcal{B}) + \beta \mathbf{I} \approx \mathbf{K}_B + \beta \mathbf{I}$, i.e. $\boldsymbol{\Sigma}_B \approx N^{-1}(\mathbf{K}_B + \beta \mathbf{I})^{-1}$. As we will see in the following, the K-FAC approximation enables us to sample efficiently from the corresponding LA.

Leveraging K-FAC’s block-diagonal structure. K-FAC is a block-diagonal curvature approximation, i.e. $\mathbf{K}_B = \text{blockdiag}_{l=1, \dots, L}(\mathbf{K}_B^{(l)})$ and $\mathbf{K}_B + \beta \mathbf{I}$ inherits this structure. The inverse of a block-diagonal matrix is also block-diagonal with the block inverses on the diagonal, i.e.

$$\begin{aligned} \boldsymbol{\Sigma}_B &\approx \frac{1}{N}(\mathbf{K}_B + \beta \mathbf{I})^{-1} \\ &= \frac{1}{N}(\text{blockdiag}_{l=1, \dots, L}(\mathbf{K}_B^{(l)} + \beta \mathbf{I}))^{-1} \\ &= \text{blockdiag}_{l=1, \dots, L} \left(\underbrace{\frac{1}{N}(\mathbf{K}_B^{(l)} + \beta \mathbf{I})^{-1}}_{=: \boldsymbol{\Sigma}^{(l)}} \right) \\ &= \text{blockdiag}_{l=1, \dots, L}(\boldsymbol{\Sigma}^{(l)}). \end{aligned}$$

To draw a sample $\mathbf{v} \in \mathbb{R}^P$ from $\mathcal{N}(\boldsymbol{\theta}_*, \boldsymbol{\Sigma}_B)$, we can thus simply sample from the block covariances $\mathbf{v}^{(l)} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{(l)})$, stack these samples and add the mean, i.e. $\mathbf{v} = \boldsymbol{\theta}_* + (\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(L)})$.

Leveraging the blocks’ Kronecker structure. To sample from $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}^{(l)})$, we can exploit the Kronecker structure of the blocks $\mathbf{K}_B^{(l)} = \mathbf{A} \otimes \mathbf{B}$ (we omit the layer index l for \mathbf{A} and \mathbf{B} for brevity). First, we compute the eigendecompositions of the Kronecker factors, i.e. $\mathbf{A} = \mathbf{U}_A \mathbf{S}_A \mathbf{U}_A^\top$ and $\mathbf{B} = \mathbf{U}_B \mathbf{S}_B \mathbf{U}_B^\top$. It follows

$$\begin{aligned} \mathbf{K}_B^{(l)} + \beta \mathbf{I} &= \mathbf{A} \otimes \mathbf{B} + \beta \mathbf{I} \\ &= \mathbf{U}_A \mathbf{S}_A \mathbf{U}_A^\top \otimes \mathbf{U}_B \mathbf{S}_B \mathbf{U}_B^\top + \beta \mathbf{I} \\ &= \underbrace{(\mathbf{U}_A \otimes \mathbf{U}_B)}_{=: \mathbf{U}} \underbrace{(\mathbf{S}_A \otimes \mathbf{S}_B)}_{=: \mathbf{S}} (\mathbf{U}_A \otimes \mathbf{U}_B)^\top + \beta \mathbf{I} \\ &= \mathbf{U} \mathbf{S} \mathbf{U}^\top + \beta \mathbf{I} \\ &= \mathbf{U}(\mathbf{S} + \beta \mathbf{I}) \mathbf{U}^\top. \end{aligned}$$

$\mathbf{U} = \mathbf{U}_A \otimes \mathbf{U}_B$ forms an orthogonal eigenbasis of the block and the diagonal matrix $\mathbf{S} + \beta \mathbf{I} = \mathbf{S}_A \otimes \mathbf{S}_B + \beta \mathbf{I}$ contains the corresponding eigenvalues (see e.g. (George et al., 2018)). It follows

$$\begin{aligned} \boldsymbol{\Sigma}^{(l)} &= \frac{1}{N}(\mathbf{K}_B^{(l)} + \beta \mathbf{I})^{-1} \\ &= \frac{1}{N} \mathbf{U}(\mathbf{S} + \beta \mathbf{I})^{-1} \mathbf{U}^\top \\ &= \underbrace{\frac{1}{\sqrt{N}} \mathbf{U}(\mathbf{S} + \beta \mathbf{I})^{-1/2}}_{=: \mathbf{V}} \frac{1}{\sqrt{N}} (\mathbf{S} + \beta \mathbf{I})^{-1/2} \mathbf{U}^\top \end{aligned}$$

$$= \mathbf{V}\mathbf{V}^\top.$$

So, in order to draw a sample $\mathbf{v}^{(l)} \sim \mathcal{N}(\mathbf{0}, \Sigma^{(l)})$, we first draw from a standard Gaussian $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and then transform the sample via $\mathbf{v}^{(l)} = \mathbf{V}\mathbf{w}$. The resulting vector $\mathbf{v}^{(l)}$ has mean $\mathbf{0}$ and covariance $\mathbf{V}\mathbf{I}\mathbf{V}^\top = \Sigma^{(l)}$. As Gaussians are closed under affine linear transformations, $\mathbf{v}^{(l)}$ is indeed distributed according to $\mathcal{N}(\mathbf{0}, \Sigma^{(l)})$.

The matrix-vector product $\mathbf{w} \mapsto \mathbf{V}\mathbf{w}$ can be computed efficiently without actually forming \mathbf{V} in memory. The first step is to multiply with $(\mathbf{S} + \beta\mathbf{I})^{-1/2}$. Since this is a diagonal matrix, we can simply multiply the vector \mathbf{w} element-wise with the inverse square root of the diagonal entries of $\mathbf{S} + \beta\mathbf{I}$. The second step is to multiply with $\mathbf{U} = \mathbf{U}_A \otimes \mathbf{U}_B$, which can be implemented efficiently by using the property $(\mathbf{U}_A \otimes \mathbf{U}_B) \text{vec}(\mathbf{W}) = \text{vec}(\mathbf{U}_B \mathbf{W} \mathbf{U}_A^\top)$. Finally, we scale by $N^{-1/2}$.

Summary. Drawing a sample from the K-FAC LA $\mathcal{N}(\boldsymbol{\theta}_*, N^{-1}(\mathbf{K}_B + \beta\mathbf{I})^{-1})$ can be done without ever forming the blocks of the covariance matrix explicitly. Using properties of the Kronecker product, we can efficiently transform a sample from the standard Gaussian to a sample $\mathbf{v}^{(l)} \sim \mathcal{N}(\mathbf{0}, \Sigma^{(l)})$. Stacking the samples from all blocks and adding the mean $\boldsymbol{\theta}_*$ yields a sample from the LA due to the block-diagonal structure of the covariance matrix.

A.3.3 DEBIASED K-FAC LAPLACE APPROXIMATION

The idea of the debiased K-FAC Laplace approximation is to construct one Kronecker-factored curvature matrix from two mini-batches \mathcal{B} and $\tilde{\mathcal{B}}$.

Eigenbasis of a K-FAC block. First, recall that the eigendecomposition of a block $\mathbf{K}_B^{(l)}$ from a K-FAC matrix can be constructed from the eigendecompositions of the Kronecker factors (see e.g. (George et al., 2018)). Consider the l -th block from both K-FAC approximations $\mathbf{K}_B^{(l)} = \mathbf{A} \otimes \mathbf{B}$ and $\mathbf{K}_{\tilde{B}}^{(l)} = \mathbf{C} \otimes \mathbf{D}$ (we omit the layer index l for \mathbf{A} , \mathbf{B} , \mathbf{C} , and \mathbf{D} for brevity). Again, let $\mathbf{A} = \mathbf{U}_A \mathbf{S}_A \mathbf{U}_A^\top$ and $\mathbf{B} = \mathbf{U}_B \mathbf{S}_B \mathbf{U}_B^\top$ denote the eigendecompositions of the Kronecker factors. It holds

$$\begin{aligned} \mathbf{K}_B^{(l)} &= \mathbf{A} \otimes \mathbf{B} \\ &= \mathbf{U}_A \mathbf{S}_A \mathbf{U}_A^\top \otimes \mathbf{U}_B \mathbf{S}_B \mathbf{U}_B^\top \\ &= \underbrace{(\mathbf{U}_A \otimes \mathbf{U}_B)}_{=: \mathbf{U}} \underbrace{(\mathbf{S}_A \otimes \mathbf{S}_B)}_{=: \mathbf{S}} (\mathbf{U}_A \otimes \mathbf{U}_B)^\top \\ &= \mathbf{U} \mathbf{S} \mathbf{U}^\top \end{aligned}$$

$\mathbf{U} = \mathbf{U}_A \otimes \mathbf{U}_B$ forms an orthogonal eigenbasis of the block and the diagonal matrix $\mathbf{S} = \mathbf{S}_A \otimes \mathbf{S}_B$ contains the eigenvalues.

Re-evaluation of the directional curvatures. For the debiased approach, we keep the block's eigenbasis \mathbf{U} , but instead of using the directional curvatures \mathbf{S} , we re-evaluate these measurements on the second mini-batch $\tilde{\mathcal{B}}$. First, consider the projection of the block $\mathbf{K}_{\tilde{B}}^{(l)}$ onto the eigenvectors \mathbf{U} , i.e.

$$\mathbf{U}^\top \mathbf{K}_{\tilde{B}}^{(l)} \mathbf{U} = (\mathbf{U}_A \otimes \mathbf{U}_B)^\top (\mathbf{C} \otimes \mathbf{D}) (\mathbf{U}_A \otimes \mathbf{U}_B) = \mathbf{U}_A^\top \mathbf{C} \mathbf{U}_A \otimes \mathbf{U}_B^\top \mathbf{D} \mathbf{U}_B.$$

The debiased directional curvatures are on the diagonal of $\mathbf{U}^\top \mathbf{K}_{\tilde{B}}^{(l)} \mathbf{U}$. For a square matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$, let $\text{Diag}(\mathbf{X})$ denote the operator that maps the matrix onto its diagonal, i.e. $\text{Diag} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n$ with $\text{Diag}(\mathbf{X})_i := X_{ii}$ for $i \in \{1, \dots, n\}$. It holds:

$$\begin{aligned} \text{Diag}(\mathbf{U}^\top \mathbf{K}_{\tilde{B}}^{(l)} \mathbf{U}) &= \text{Diag}(\mathbf{U}_A^\top \mathbf{C} \mathbf{U}_A \otimes \mathbf{U}_B^\top \mathbf{D} \mathbf{U}_B) \\ &= \underbrace{\text{Diag}(\mathbf{U}_A^\top \mathbf{C} \mathbf{U}_A)}_{=: \tilde{\mathbf{s}}_A} \otimes \underbrace{\text{Diag}(\mathbf{U}_B^\top \mathbf{D} \mathbf{U}_B)}_{=: \tilde{\mathbf{s}}_B} \\ &= \tilde{\mathbf{s}}_A \otimes \tilde{\mathbf{s}}_B \\ &=: \tilde{\mathbf{s}}. \end{aligned}$$

Construction of a debiased block. Now that we have the eigenbasis U and the debiased directional curvatures \tilde{s} , we construct the debiased block. Let $\tilde{S}_A := \text{diag}(\tilde{s}_A)$, $\tilde{S}_B := \text{diag}(\tilde{s}_B)$ and

$$\tilde{S} := \text{diag}(\tilde{s}) = \text{diag}(\tilde{s}_A \otimes \tilde{s}_B) = \text{diag}(\tilde{s}_A) \otimes \text{diag}(\tilde{s}_B) = \tilde{S}_A \otimes \tilde{S}_B.$$

The debiased block is given by $U\tilde{S}U^\top$. It can be written as the Kronecker product of two matrices \tilde{A} and \tilde{B} since

$$U\tilde{S}U^\top = (U_A \otimes U_B)(\tilde{S}_A \otimes \tilde{S}_B)(U_A \otimes U_B)^\top = \underbrace{(U_A \tilde{S}_A U_A^\top)}_{=: \tilde{A}} \otimes \underbrace{(U_B \tilde{S}_B U_B^\top)}_{=: \tilde{B}}.$$

This is important as it allows for efficient sampling as described in [Appendix A.3.2](#).

Computational cost. Since we need the eigendecompositions of the Kronecker factors $A = U_A \tilde{S}_A U_A^\top$ and $B = U_B \tilde{S}_B U_B^\top$ for sampling in any case (see [Appendix A.3.2](#)), the computational overhead for the debiased K-FAC approximation $\tilde{A} \otimes \tilde{B}$ consists of computing a K-FAC approximation $C \otimes D$ on another mini-batch, re-evaluating the directional curvatures $\tilde{s}_A = \text{Diag}(U_A^\top C U_A)$ and $\tilde{s}_B = \text{Diag}(U_B^\top D U_B)$ and finally computing $\tilde{A} = U_A \text{diag}(\tilde{s}_A) U_A^\top$ and $\tilde{B} = U_B \text{diag}(\tilde{s}_B) U_B^\top$.

From block-level to full matrix. So far, we have only considered the debiasing of a single block. However, correcting the blocks’s eigenvalues is sufficient because they coincide with the eigenvalues of the full matrix (due to the block-diagonal structure). Let $u^{(l)}$ denote an eigenvector of the l -th block $X^{(l)}$ of some block-diagonal matrix $X = \text{blockdiag}_{l=1,\dots,L}(X^{(l)})$ corresponding to the eigenvalue λ . Then, $u^\top := (\mathbf{0}^\top, \dots, u^{(l)\top}, \dots, \mathbf{0}^\top)$ is an eigenvector of X corresponding to the same eigenvalue λ , because

$$X \cdot u = \text{blockdiag}_{l=1,\dots,L}(X^{(l)}) \cdot u = \begin{pmatrix} X^{(1)} \cdot \mathbf{0} \\ \vdots \\ X^{(l)} \cdot u^{(l)} \\ \vdots \\ X^{(L)} \cdot \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \vdots \\ \lambda u^{(l)} \\ \vdots \\ \mathbf{0} \end{pmatrix} = \lambda u.$$

The eigenvalues of X thus coincide with the eigenvalues of its blocks; and X ’s eigenvectors can be constructed from the the eigenvectors of the blocks by filling them up with zeros.

Connection to Equation (9): The equivalence between [Equation \(9\)](#) and the approach we describe above may not be obvious. We thus show here that the directional curvature of the debiased matrix \hat{K} along an eigenvector u of K_B indeed coincides with the directional curvature $u^\top K_{\tilde{B}} u$ on \tilde{B} .

More concretely, let $\hat{K} = \text{blockdiag}_{l=1,\dots,L}(\hat{K}^{(l)})$ denote the debiased K-FAC approximation constructed from K_B and $K_{\tilde{B}}$ as described above. Also, let $u^\top := (\mathbf{0}^\top, \dots, u^{(l)\top}, \dots, \mathbf{0}^\top)$ denote an eigenvector of K_B , where $u^{(l)}$ is the i -th eigenvector of $K_B^{(l)}$ (i.e. the i -th column of U). It holds

$$u^\top \hat{K} u = u^{(l)\top} \hat{K}^{(l)} u^{(l)} = u^{(l)\top} U \tilde{S} U^\top u^{(l)} = e_i^\top \tilde{S} e_i = \tilde{s}_i = u^{(l)\top} K_{\tilde{B}}^{(l)} u^{(l)} = u^\top K_{\tilde{B}} u,$$

where e_i denotes the i -th eigenvector.

A.4 BIAS IN THE DIRECTIONAL SLOPE

Biased directional slopes along negative gradient directions. Assume that we are given a quadratic $q(\cdot; \mathcal{B})$ around the current parameters θ_0 . We consider the negative normalized gradient direction $\tilde{d} = -\nabla q(\theta_\bullet; \mathcal{B}) \cdot \|\nabla q(\theta_\bullet; \mathcal{B})\|^{-1}$ at some location θ_\bullet evaluated on mini-batch \mathcal{B} . We have

$$\underbrace{\partial_d q(\theta_\bullet; \tilde{\mathcal{B}})}_{\bullet} = \underbrace{\partial_d q(\theta_\bullet; \mathcal{B})}_{\bullet} + \|\nabla q(\theta_\bullet; \mathcal{B})\| (1 - \cos(\alpha)) \geq \underbrace{\partial_d q(\theta_\bullet; \mathcal{B})}_{\bullet}, \quad (19)$$

where $\alpha := \angle(\nabla q(\theta_\bullet; \mathcal{B}), \nabla q(\theta_\bullet; \tilde{\mathcal{B}}))$ and we assumed $\|\nabla q(\theta_\bullet; \tilde{\mathcal{B}})\| = \|\nabla q(\theta_\bullet; \mathcal{B})\|$ which is true at least in expectation. Projecting a gradient onto its negative direction—the right-hand side of the

inequality—will *always* result in a directional slope that is ≤ 0 since

$$\partial_{\mathbf{d}} q(\boldsymbol{\theta}_\bullet; \mathcal{B}) = \mathbf{d}^\top \nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B}) = -\frac{\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})^\top}{\|\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})\|} \nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B}) = -\|\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})\| \leq 0$$

and = 0 only if $\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B}) = \mathbf{0}$. Projecting a *different* gradient (of equal length) onto that direction—that is the left-hand side of the inequality—will result in a larger (possibly even positive) directional slope.

Proof of Equation (19). Equation (19) quantifies the bias in the directional slope along \mathbf{d} as a function of the alignment of the two mini-batch gradients. It holds

$$\begin{aligned} \underbrace{\partial_{\mathbf{d}} q(\boldsymbol{\theta}_\bullet; \mathcal{B})}_{\bullet} - \underbrace{\partial_{\mathbf{d}} q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})}_{\bullet} &= \mathbf{d}^\top \nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B}) - \mathbf{d}^\top \nabla q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}}) \\ &= \underbrace{\|\mathbf{d}\|}_{=1} \|\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})\| \underbrace{\cos(\pi)}_{=-1} - \underbrace{\|\mathbf{d}\|}_{=1} \|\nabla q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})\| \cos(\gamma) \\ &= \|\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})\|(-1 - \cos(\gamma)), \end{aligned}$$

where $\gamma = \angle(\mathbf{d}, \nabla q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}}))$ and we assumed that $\|\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})\| = \|\nabla q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})\|$ in the last step (which is true, at least, in expectation). Next, we re-write γ as

$$\gamma = \angle(\mathbf{d}, \nabla q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})) = \angle(-\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B}), \nabla q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})) = \pi - \underbrace{\angle(\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B}), \nabla q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}}))}_{=: \alpha}.$$

It follows $-1 - \cos(\gamma) = -1 - \cos(\pi - \alpha) = \cos(\alpha) - 1$. Substituting this into the expression for the bias, we arrive at

$$\underbrace{\partial_{\mathbf{d}} q(\boldsymbol{\theta}_\bullet; \mathcal{B})}_{\bullet} - \underbrace{\partial_{\mathbf{d}} q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})}_{\bullet} = \|\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})\|(-1 - \cos(\gamma)) = \|\nabla q(\boldsymbol{\theta}_\bullet; \mathcal{B})\|(\cos(\alpha) - 1),$$

from which Equation (19) follows. \square

The first CG search direction. Assume that the current parameters are $\boldsymbol{\theta}_0$ and we apply CG to the local quadratic approximation $q(\cdot; \mathcal{B})$. The very first search direction is the quadratic’s normalized negative gradient $\mathbf{d}_0 = -\nabla q(\boldsymbol{\theta}_0; \mathcal{B}) \cdot \|\nabla q(\boldsymbol{\theta}_0; \mathcal{B})\|^{-1}$ at $\boldsymbol{\theta}_0$. This is exactly the situation we describe above, where $\boldsymbol{\theta}_\bullet$ is set to $\boldsymbol{\theta}_0$. Equation (19) thus explains the bias for the first CG direction.

The subsequent CG search directions. For the subsequent CG search directions (\mathbf{d}_p with $p \geq 1$), the situation is more complex: Each direction \mathbf{d}_p is a linear combination of the current normalized negative residual $-\mathbf{r}_p$ (note that \mathbf{r}_p coincides with the gradient $\nabla q(\boldsymbol{\theta}_p; \mathcal{B})$, see Appendix A.2) and the previous search direction \mathbf{d}_{p-1} (see Appendix A.2) and the previous search direction \mathbf{d}_{p-1} (see Algorithm 1), i.e.

$$\mathbf{d}_p = \frac{\mathbf{s}_p}{\|\mathbf{s}_p\|} = \frac{1}{\|\mathbf{s}_p\|}(-\mathbf{r}_p) + \frac{\beta_p}{\|\mathbf{s}_p\|} \mathbf{s}_{p-1} = \underbrace{\frac{\|\mathbf{r}_p\|}{\|\mathbf{s}_p\|}}_{=: \eta_1 \geq 0} \frac{-\mathbf{r}_p}{\|\mathbf{r}_p\|} + \underbrace{\frac{\beta_p \|\mathbf{s}_{p-1}\|}{\|\mathbf{s}_p\|}}_{=: \eta_2 \geq 0} \mathbf{d}_{p-1} = \eta_1 \frac{-\mathbf{r}_p}{\|\mathbf{r}_p\|} + \eta_2 \mathbf{d}_{p-1}.$$

The additional correction term ensures conjugacy. The directional slope along \mathbf{d}_p thus also splits into two corresponding terms: The slope along the negative gradient direction and the slope along the previous search direction. It holds

$$\begin{aligned} \partial_{\mathbf{d}_p} q(\boldsymbol{\theta}_p; \mathcal{B}) &= \mathbf{d}_p^\top \nabla q(\boldsymbol{\theta}_p; \mathcal{B}) \\ &= \eta_1 \frac{-\mathbf{r}_p^\top}{\|\mathbf{r}_p\|} \nabla q(\boldsymbol{\theta}_p; \mathcal{B}) + \eta_2 \mathbf{d}_{p-1}^\top \nabla q(\boldsymbol{\theta}_p; \mathcal{B}) \\ &= \eta_1 \partial_{-\mathbf{r}_p/\|\mathbf{r}_p\|} q(\boldsymbol{\theta}_p; \mathcal{B}) + \eta_2 \partial_{\mathbf{d}_{p-1}} q(\boldsymbol{\theta}_p; \mathcal{B}). \end{aligned}$$

The bias in the first term—as explained by Equation (19)—also introduces a bias along the search direction \mathbf{d}_p .

A.5 BIAS IN THE DIRECTIONAL CURVATURE

Derivation for Equation (6). Let $\mathbf{H}_{\mathcal{B}} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$ and $\mathbf{H}_{\tilde{\mathcal{B}}} = \tilde{\mathbf{U}}\tilde{\boldsymbol{\Lambda}}\tilde{\mathbf{U}}^\top$ denote the eigendecompositions of the mini-batch Hessians, where $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_P)$, $\tilde{\mathbf{U}} = (\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_P) \in \mathbb{R}^{P \times P}$ contain

the orthonormal eigenvectors and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_P)$, $\tilde{\mathbf{\Lambda}} = \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_P)$ the respective eigenvalues in descending order, i.e. $\lambda_1 \geq \dots \geq \lambda_P$ and $\tilde{\lambda}_1 \geq \dots \geq \tilde{\lambda}_P$. The directional curvature along one of $\mathbf{H}_{\mathcal{B}}$'s eigenvectors \mathbf{u}_i on mini-batch \mathcal{B} is given by the corresponding eigenvalue since

$$\underbrace{\partial_{\mathbf{u}_i}^2 q(\boldsymbol{\theta}_\bullet; \mathcal{B})}_{\bullet} = \mathbf{u}_i^\top \mathbf{H}_{\mathcal{B}} \mathbf{u}_i = \mathbf{u}_i^\top \lambda_i \mathbf{u}_i = \lambda_i \|\mathbf{u}_i\|_2^2 = \lambda_i.$$

For $\tilde{\mathcal{B}}$, we obtain

$$\underbrace{\partial_{\mathbf{u}_i}^2 q(\boldsymbol{\theta}_\bullet; \tilde{\mathcal{B}})}_{\bullet} = \mathbf{u}_i^\top \mathbf{H}_{\tilde{\mathcal{B}}} \mathbf{u}_i = \mathbf{u}_i^\top \left(\sum_{p=1}^P \tilde{\lambda}_p \tilde{\mathbf{u}}_p \tilde{\mathbf{u}}_p^\top \right) \mathbf{u}_i = \sum_{p=1}^P \tilde{\lambda}_p \underbrace{(\tilde{\mathbf{u}}_p^\top \mathbf{u}_i)^2}_{=:\Omega_{i,p}} = \sum_{p=1}^P \tilde{\lambda}_p \Omega_{i,p}.$$

Weights sum up to one. The weights $\Omega_{i,p}$ are non-negative and sum up to one, i.e. $\sum_{p=1}^P \Omega_{i,p} = 1$. This is because

$$\sum_{p=1}^P \Omega_{i,p} = \sum_{p=1}^P (\tilde{\mathbf{u}}_p^\top \mathbf{u}_i)^2 = \|\tilde{\mathbf{U}}^\top \mathbf{u}_i\|_2^2 = \mathbf{u}_i^\top \underbrace{\tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top}_{=\mathbf{I}} \mathbf{u}_i = \|\mathbf{u}_i\|_2^2 = 1 \quad (20)$$

where we used that $\tilde{\mathbf{U}}$ is an orthogonal matrix and that \mathbf{u}_i is normalized.

Proof of Equation (7). Here, we assume the simplified case where the spectra of $\mathbf{H}_{\mathcal{B}}$ and $\mathbf{H}_{\tilde{\mathcal{B}}}$ are identical, i.e. $\lambda_p = \tilde{\lambda}_p \forall p \in \{1, \dots, P\}$.

First, consider \mathbf{u}_1 and assume that \mathbf{u}_1 and $\tilde{\mathbf{u}}_1$ are *not* perfectly aligned, i.e. $\Omega_{1,1} = (\mathbf{u}_1^\top \tilde{\mathbf{u}}_1)^2 < 1$. It holds

$$\underbrace{\partial_{\mathbf{u}_1}^2 q(\boldsymbol{\theta}_0, \tilde{\mathcal{B}})}_{\bullet} \stackrel{(6)}{=} \sum_{p=1}^P \lambda_p \Omega_{1,p} = \lambda_1 \Omega_{1,1} + \sum_{p=2}^P \lambda_p \Omega_{1,p}.$$

The sum can be bounded from above by putting all remaining weight $1 - \Omega_{1,1}$ (see Equation (20)) on the second-largest eigenvalue λ_2 , i.e.

$$\underbrace{\partial_{\mathbf{u}_1}^2 q(\boldsymbol{\theta}_0, \tilde{\mathcal{B}})}_{\bullet} \leq \lambda_1 \Omega_{1,1} + \lambda_2 (1 - \Omega_{1,1}) \stackrel{(*)}{\leq} \lambda_1 \Omega_{1,1} + \lambda_1 (1 - \Omega_{1,1}) = \lambda_1 = \underbrace{\partial_{\mathbf{u}_1}^2 q(\boldsymbol{\theta}_0, \mathcal{B})}_{\bullet}.$$

The inequality $(*)$ turns into $<$ if the top two eigenvalues are separated, i.e. $\lambda_1 > \lambda_2$. The proof for \mathbf{u}_P is similar. We assume that \mathbf{u}_P and $\tilde{\mathbf{u}}_P$ are *not* perfectly aligned, i.e. $\Omega_{P,P} = (\mathbf{u}_P^\top \tilde{\mathbf{u}}_P)^2 < 1$. It follows

$$\underbrace{\partial_{\mathbf{u}_P}^2 q(\boldsymbol{\theta}_0, \tilde{\mathcal{B}})}_{\bullet} \stackrel{(6)}{=} \sum_{p=1}^P \lambda_p \Omega_{P,p} = \sum_{p=1}^{P-1} \lambda_p \Omega_{P,p} + \lambda_P \Omega_{P,P}.$$

This time, we bound the sum from below by putting all remaining weight $1 - \Omega_{P,P}$ (see Equation (20)) on the second-smallest eigenvalue λ_{P-1} , i.e.

$$\underbrace{\partial_{\mathbf{u}_P}^2 q(\boldsymbol{\theta}_0, \tilde{\mathcal{B}})}_{\bullet} \geq \lambda_{P-1} (1 - \Omega_{P,P}) + \lambda_P \Omega_{P,P} \stackrel{(*)}{\geq} \lambda_P (1 - \Omega_{P,P}) + \lambda_P \Omega_{P,P} = \lambda_P = \underbrace{\partial_{\mathbf{u}_P}^2 q(\boldsymbol{\theta}_0, \mathcal{B})}_{\bullet}.$$

Again, the inequality $(*)$ turns into $>$ if the bottom two eigenvalues are separated, i.e. $\lambda_{P-1} > \lambda_P$. This concludes the proof of Equation (7). \square

B EXPERIMENTAL DETAILS

In the following, we provide information on the experimental setup and give detailed instructions on how to replicate all empirical results.

1242 B.1 TEST PROBLEMS AND TRAINING PROCEDURES
1243

1244 Throughout the paper, we use a series of test problems with different models, datasets, and training
1245 procedures which we describe in more detail in the following. We use DEEPOBS (Schneider et al.,
1246 2019) on top of PYTORCH (Paszke et al., 2019) as our general benchmarking framework as it provides
1247 easy access to a variety of datasets and model architectures.

1248 **Data.** We use the datasets CIFAR-10 (with $C = 10$ classes) and CIFAR-100 (with $C = 100$
1249 classes) (Krizhevsky, 2009). Each dataset contains 60,000 data points that are split into 40,000
1250 training samples, 10,000 validation samples and 10,000 test samples. For the experiments on out-of-
1251 distribution (OOD) data, we create the datasets CIFAR-10-C and CIFAR-100-C, each containing
1252 10,000 images, as described in (Hendrycks & Dietterich, 2019). In these datasets, each image is
1253 corrupted using one out of 15 different corruptions (chosen uniformly at random) at a specific severity
1254 level (a number between 1 and 5). We also use the IMAGENET dataset (Deng et al., 2009) which
1255 contains images from $C = 1000$ different classes.

1256 **Test problems.** All of the following test problems use the cross-entropy loss function.

- 1257
- 1258 (A) **ALL-CNN-C on CIFAR-100.** This test problem uses the ALL-CNN-C model architecture
1259 (Springenberg et al., 2015) (where we removed the dropout layers as explained in Section 3.1)
1260 and CIFAR-100 data. The training hyperparameters are taken from an existing benchmark
1261 (Schmidt et al., 2021): We train the model with SGD (learning rate 0.171234) with batch
1262 size 256 for 350 epochs. Weight decay $\beta = 0.0005$ is used on the weights but not the biases
1263 of the model.
 - 1264 (B) **ALL-CNN-C on CIFAR-10.** This test problem is similar to (A) but uses the CIFAR-10
1265 dataset. The model is trained with SGD (learning rate 0.025, momentum 0.9) with batch
1266 size 256 for 350 epochs. The learning rate is reduced by a factor of 10 at epochs 200, 250
1267 and 300, as suggested in (Springenberg et al., 2015). Weight decay $\beta = 0.001$ is used on
1268 the weights but not the biases of the model.
 - 1269 (C) **WIDERESNET 40-4 on CIFAR-100.** This is a test problem from DEEPOBS. For details
1270 on the architecture, see (Zagoruyko & Komodakis, 2017). We train this model with SGD
1271 (learning rate 0.1, momentum 0.9) with batch size 128 for 160 epochs. The learning rate is
1272 reduced by a factor of 5 after 60 and 120 epochs. Weight decay $\beta = 0.0005$ is used on the
1273 non-bias weights of the model.
 - 1274 (D) **CONVNET on CIFAR-10.** This test problem uses a simple convolutional neural network
1275 with variable depth d and width w (denoted by “model d - w ” in Figures 7 and 22). The
1276 *first* block of the model consists of a convolutional layer (kernel size 5, padding 2) with 3
1277 input and w output channels, a ReLU activation function, and a max-pooling layer (kernel
1278 size 2). The *last* block consists of a max-pooling layer (kernel size 2), a flatten layer, and a
1279 dense linear layer. In between those blocks, there are d *hidden* blocks each consisting of a
1280 convolutional layer (kernel size 5, padding 2) with w input and w output channels followed
1281 by a ReLU activation function. So, the depth d determines the number of hidden blocks,
1282 whereas w controls the number of parameters in the layers and is thus referred to as the width.
1283 We use $d \in \{1, 4, 7\}$ and $w \in \{32, 64, 128\}$ (resulting in 9 models) and train each model for
1284 100 epochs on the CIFAR-10 dataset using ADAM with standard hyperparameters (learning
1285 rate 0.001, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$) with batch size 256. No weight decay is used
1286 for this test problem.
 - 1287 (E) **RESNET-50 on IMAGENET.** This test problem uses the RESNET-50 model architecture
1288 (He et al., 2016) and the IMAGENET dataset. We use the IMAGENET1K_V1 weights pre-
1289 trained on IMAGENET-1K using SGD (initial learning rate 0.1, momentum 0.9) with batch
1290 size 256 for 90 epochs. For the pre-training, the learning rate was set by a multistep scheduler
1291 which multiplied the learning rate by a factor of 0.1 after every 30 epochs. Additionally,
1292 a weight decay of $\beta = 0.0001$ was used on all weights apart from biases and the learned
1293 batch normalization weights.
 - 1294 (F) **VIT LITTLE on IMAGENET.** This test problem uses the VIT LITTLE architecture of the
1295 VISION TRANSFORMER model family (Dosovitskiy et al., 2021) on the IMAGENET dataset.
We use the VIT_LITTLE_PATCH16_REG4_GAP_256.SBB_IN1K weights pre-trained on
IMAGENET-1K using NADAMW. For the exact hyperparameters used for the weights, refer
to the HuggingFace Model Card.

During training, we store the model’s parameters at 10 checkpoints spaced log-equidistantly between the first and last epoch. The training metrics for all test problems are shown in Figure 7.

B.2 MATRIX-VECTOR PRODUCTS WITH THE CURVATURE MATRIX

Access to curvature matrices via BACKPACK. In order to compute an eigendecomposition of the curvature matrix or apply the CG method to a quadratic, we need access to matrix-vector products with the curvature matrix $v \mapsto \mathbf{H}_B \cdot v$. When products with the full-batch curvature matrix are required, we accumulate the mini-batch quantities over manageable chunks of the training set. Our implementation uses BACKPACK (Dangel et al., 2020) that provides access to products with the Hessian $v \mapsto \nabla^2 \mathcal{L}(\theta; \mathcal{B}) \cdot v$ and GGN $v \mapsto \mathbf{G}_B \cdot v$ of the empirical risk \mathcal{L} as well as the K-FAC curvature approximation \mathbf{K}_B .

Eigenvectors of the curvature matrix. Given access to matrix-vector products with the curvature matrix, we can construct an instance of `scipy.sparse.linalg.LinearOperator` that can be used in `scipy.sparse.linalg.eigsh`.

B.3 SECTION 1 (INTRODUCTION): FIGURE 1

For the visual abstract in Figure 1, we use the fully trained ALL-CNN-C model from test problem (A) (see Appendix B.1). The experimental procedure below is repeated for 5 different mini-batches \mathcal{B}_m , $m \in \{0, 1, 2, 3, 4\}$ of size $|\mathcal{B}_m| = 512$ (we omit the index m in the following).

Experimental procedure. We use the GGN curvature proxy $\mathbf{H}_B \leftarrow \mathbf{G}_B + \beta \mathbf{I}$ and compute its top two eigenvectors \mathbf{u}_1 and \mathbf{u}_2 (see Appendix B.2). In order to evaluate the quadratic $q(\theta_\star; \mathcal{B})$ in the 2D space spanned by those eigenvectors efficiently, we use the following equation:

$$q(\theta_\star + \tau_1 \mathbf{u}_1 + \tau_2 \mathbf{u}_2; \mathcal{B}) = \frac{1}{2} \tau_1^2 [\mathbf{u}_1^\top \mathbf{H}_B \mathbf{u}_1] + \tau_1 \tau_2 [\mathbf{u}_1^\top \mathbf{H}_B \mathbf{u}_2] + \frac{1}{2} \tau_2^2 [\mathbf{u}_2^\top \mathbf{H}_B \mathbf{u}_2] \\ + \tau_1 [\mathbf{u}_1^\top \mathbf{g}_B] + \tau_2 [\mathbf{u}_2^\top \mathbf{g}_B] + [c_B].$$

We can compute all the terms on the right-hand side in brackets *once*, store them, and then evaluate the quadratic for arbitrary values of τ_1 and τ_2 at basically no cost. The derivation for the full-batch quadratic $q(\theta_\star + \tau_1 \mathbf{u}_1 + \tau_2 \mathbf{u}_2; \mathcal{D})$ is analogous.

B.4 SECTION 3 (THE SHAPE OF A MINI-BATCH QUADRATIC): FIGURES 2 TO 4

For the evaluation of the directional derivatives, we use the fully trained ALL-CNN-C model from test problem (A) (see Appendix B.1).

Three settings. Throughout this section, we consider three different settings: (i) Quadratics that use the Hessian of the empirical risk, i.e. $\mathbf{H}_B \leftarrow \nabla^2 \mathcal{L}(\theta_\star; \mathcal{B}) + \beta \mathbf{I}$ at batch size 512, (ii) quadratics that use the Hessian’s GGN approximation, i.e. $\mathbf{H}_B \leftarrow \mathbf{G}_B + \beta \mathbf{I}$ at batch size 512 and (iii) at batch size 2048. For each setting, the experimental procedure below is repeated for 3 mini-batches \mathcal{B}_m , $m \in \{0, 1, 2\}$. The prior precision β is set to 0.0005 (the same β was used for training, see Appendix B.1) and only acts on the weights of the model but not its biases.

Experimental procedure (biases). The experimental procedure consists of two steps: The computations of the directions on mini-batch \mathcal{B}_m and the evaluation of the directional derivatives on *all* mini-batches (of the same size) in the training dataset.

- Directions based on \mathcal{B}_m .** First, we use the quadratic $q(\theta_\star; \mathcal{B}_m)$ with the given curvature proxy and batch size to compute a set of 100 directions. This is either the top 100 eigenvectors of \mathbf{H}_B computed via `scipy.sparse.linalg.eigsh` (see Appendix B.2) or the first 100 CG search directions (see Appendix B.5). In the case of CG, we also store the trajectory of the iterates $\theta_\star = \theta_0, \theta_1, \dots, \theta_{100}$.
- Directional derivatives on all mini-batches.** Finally, we evaluate the directional slope and curvature for all directions from step 1 on *all* mini-batches $\mathcal{B}_{m'}$, $m' \in \{0, \dots, M - 1\}$ in the training dataset (where M denotes the number of mini-batches in the training data).
 - Directional derivatives along eigenvectors.** For an eigenvector \mathbf{u} the directional slope and curvature are given by $\partial_{\mathbf{u}} q(\theta_\star; \mathcal{B}_{m'})$ and $\partial_{\mathbf{u}}^2 q(\theta_\star; \mathcal{B}_{m'})$, respectively.

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

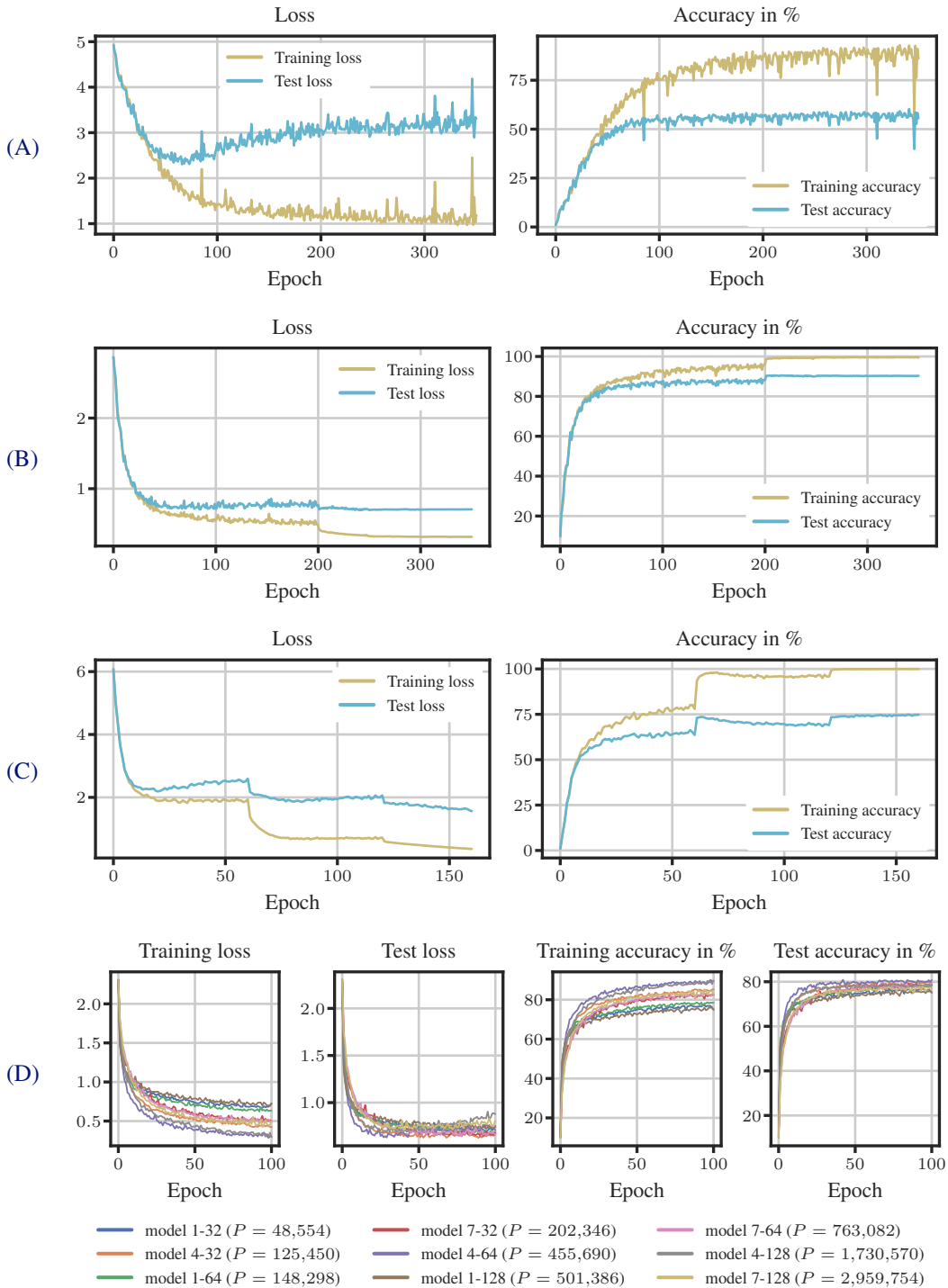


Figure 7: Training metrics for all test problems. The panels show the training/test loss/accuracy during training for all test problems (A) (ALL-CNN-C on CIFAR-100), (B) (ALL-CNN-C on CIFAR-10), (C) (WIDERESNET 40-4 on CIFAR-100) and (D) (CONVNET on CIFAR-10). Test problem (E) (RESNET-50 on IMAGENET) is not shown since a pre-trained model is used here.

- **Directional derivatives along CG search directions.** For CG, we evaluate the directional slope and curvature along each search direction \mathbf{d}_p at the respective iterate $\boldsymbol{\theta}_p$ of the trajectory, i.e. $\partial_{\mathbf{d}_p} q(\boldsymbol{\theta}_p; \mathcal{B}_{m'})$ for the slope and $\partial_{\mathbf{d}_p}^2 q(\boldsymbol{\theta}_p; \mathcal{B}_{m'})$ for the curvature. As the curvature $\nabla^2 q(\boldsymbol{\theta}; \mathcal{B}_{m'}) \equiv \mathbf{H}_{\mathcal{B}_{m'}}$ is independent of $\boldsymbol{\theta}$, this small distinction is irrelevant for the directional curvature. For the directional slope, however, it does make a difference.

Results (biases). The results for setting (i) are shown in Figure 8, for setting (ii) in Figures 2 and 9 and for (iii) in Figure 10. The biases in the directional slopes and curvatures can be observed across all scenarios. The biases in the slopes are even more pronounced along the CG search directions than along the eigenvectors of the curvature matrix. The opposite holds for the curvature biases that tend to be larger along the eigenvectors. This is consistent with our explanations from Section 3.2. Both the biases in the slope and curvature decrease with increasing mini-batch size. For setting (i), CG encounters a search direction with negative curvature (indicating that, as expected, the curvature matrix $\mathbf{H}_{\mathcal{B}}$ is indefinite) and thus terminates already after 4 iterations.

Experimental procedure (overlaps). Next, we compute the overlaps between the eigenspaces of the curvature matrices on different mini-batches \mathcal{B} and $\tilde{\mathcal{B}}$. More specifically, we compute $\Omega_{i,p} = (\mathbf{u}_i^\top \tilde{\mathbf{u}}_p)^2$, where $i, p \in \{1, \dots, 100\}$, \mathbf{u}_i is an eigenvector of $\mathbf{H}_{\mathcal{B}}$ and $\tilde{\mathbf{u}}_p$ is an eigenvector of $\mathbf{H}_{\tilde{\mathcal{B}}}$. The values $\Omega_{i,p}$ are bounded between 0 (\mathbf{u}_i and $\tilde{\mathbf{u}}_p$ are orthogonal) and 1 (\mathbf{u}_i and $\tilde{\mathbf{u}}_p$ are identical, up to their sign). We apply a log-transform to $\Omega_{i,p}$. In Figure 3 and Appendix B.4, values below -8 (i.e. $\Omega_{i,p} \leq 10^{-8}$) are shown in black, values equal to 0 (i.e. $\Omega_{i,p} = 10^0 = 1$) are shown in white.

Results (overlaps). The results are shown in Appendix B.4. They show that the eigenspaces of the curvature matrices are not perfectly aligned: Eigenvectors from \mathcal{B}_m typically overlap with several eigenvectors from $\tilde{\mathcal{B}}$. The eigenspaces are more aligned at the larger batch size 2048 than at batch size 512. This seems reasonable since random vectors in high-dimensional spaces tend to be orthogonal to each other—less stochasticity (due to a larger batch size) thus leads to better alignment.

B.5 SECTION 6.1 (DEBIASED CONJUGATE GRADIENTS): FIGURE 5

Here, we describe the experimental details for Figure 5 from Section 6.1. For the derivation and the mathematical details of the standard single-batch CG method and the debiased approach, see Appendix A.2. The experiment uses the fully trained ALL-CNN-C model from test problem (A) (see Appendix B.1) and the GGN curvature proxy $\mathbf{H}_{\mathcal{B}} \leftarrow \mathbf{G}_{\mathcal{B}} + \beta \mathbf{I}$ with $\beta = 0.0005$.

Experimental procedure. The experimental procedure consists of two steps: The computation of the CG trajectories and the evaluation of the four performance metrics.

1. **Computation of the CG trajectories.** For the single-batch CG approach, we use one mini-batch of size 1024, apply $K = 30$ CG iterations and store the trajectory $\boldsymbol{\theta}_\star = \boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{30}$. For the debiased approach (details in Appendix A.2), we use two mini-batches of size 512 to construct the trajectory. As the debiased approach uses a total of 60 GGN-vector products (30 for the search directions and 30 for the update magnitudes) at half the computational cost (since the cost for the GGN-vector product scales linearly with the mini-batch size), the total cost of both approaches are comparable. For both approaches, the above procedure is repeated for 5 different mini-batches/mini-batch pairs, such that 10 trajectories are computed in total. We compute two additional trajectories: The standard CG trajectory for the full-batch quadratic $q(\cdot; \mathcal{D})$ and the trajectory for a debiased full-batch approach, where we use half the training data for the directions and the other half for the update magnitudes.
2. **Evaluation of performance metrics.** For each of the 12 trajectories, we evaluate four performance metrics: (i) the regularized loss \mathcal{L}_{reg} (see Equation (1)) on the training set \mathcal{D} and (ii) on the test set $\mathcal{D}_{\text{test}}$ as well as the accuracy (i.e. the relative number of correctly classified samples) on the (iii) training and (iv) test dataset.

Results. The mini-batch results are shown in Figure 5 and discussed in Section 6.1. The results for the two full-batch CG variants are shown in Figure 12. Surprisingly, the full-batch CG approach diverges after roughly 10 iterations. If we consider the training data as a (large) sample from the data

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

Biases: Additional results for $H_B \leftarrow \nabla^2 \mathcal{L}(\theta_*; \mathcal{B}) + \beta I$ (batch size 512)

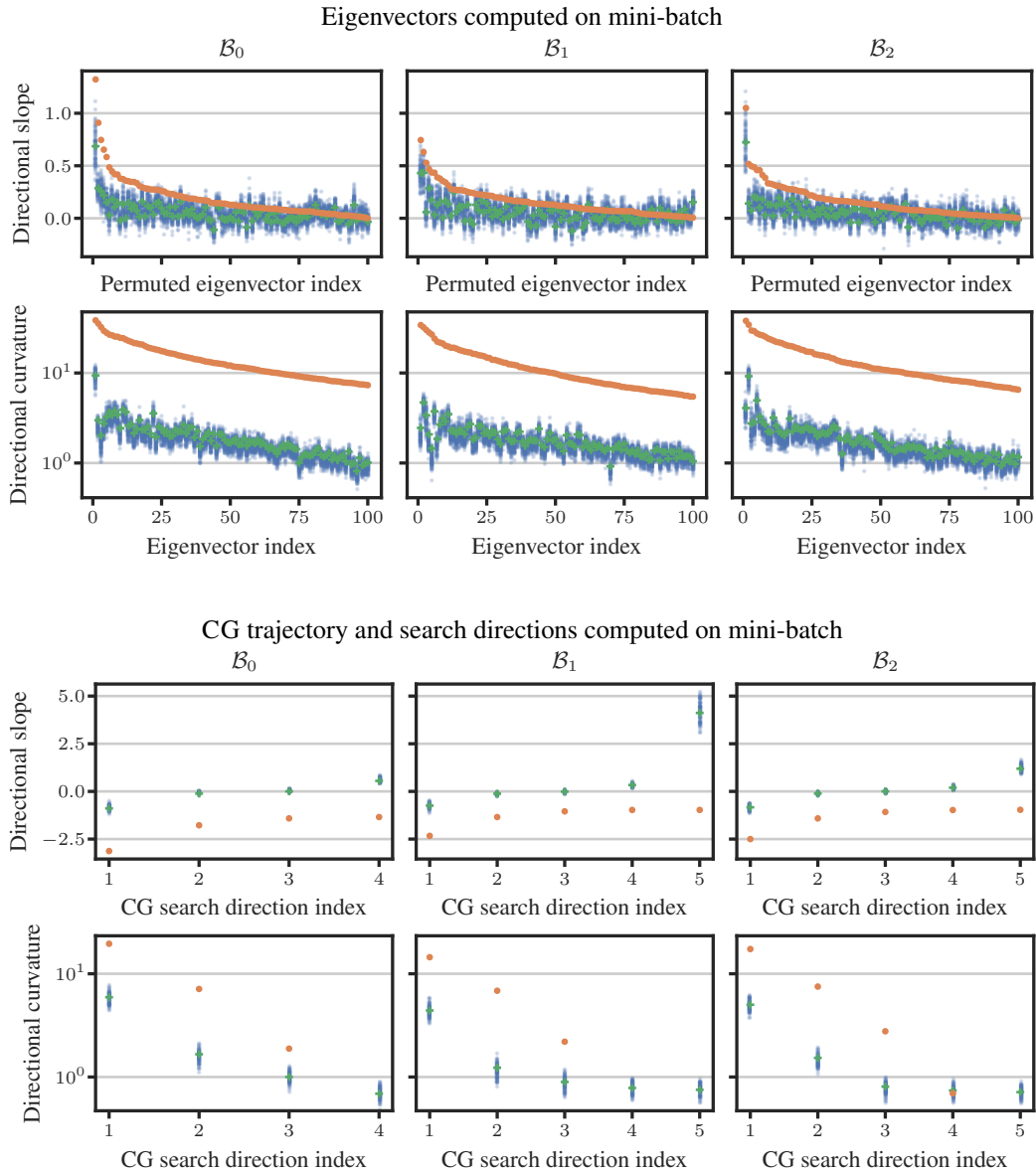


Figure 8: Directional slopes and curvatures are biased. The experimental setting is similar to Figure 2 but uses the Hessian of the empirical risk, i.e. $H_B \leftarrow \nabla^2 \mathcal{L}(\theta_*; \mathcal{B}) + \beta I$ at batch size 512. The upper plot (*Top*) shows the directional derivatives along the top 100 eigenvectors of H_B , the lower plot (*Bottom*) shows the directional derivatives along the first 100 CG search directions. For the top panel of the upper plot, we switch the order and sign of the eigenvectors such that the orange dots are all above zero and in descending order. There are strong systematic biases in the directional slopes and curvatures. The curvature biases are more pronounced along the eigenvectors, whereas the biases in the slope are larger along the CG directions.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

Biases: Additional results for $H_B \leftarrow G_B + \beta I$ (batch size 512)

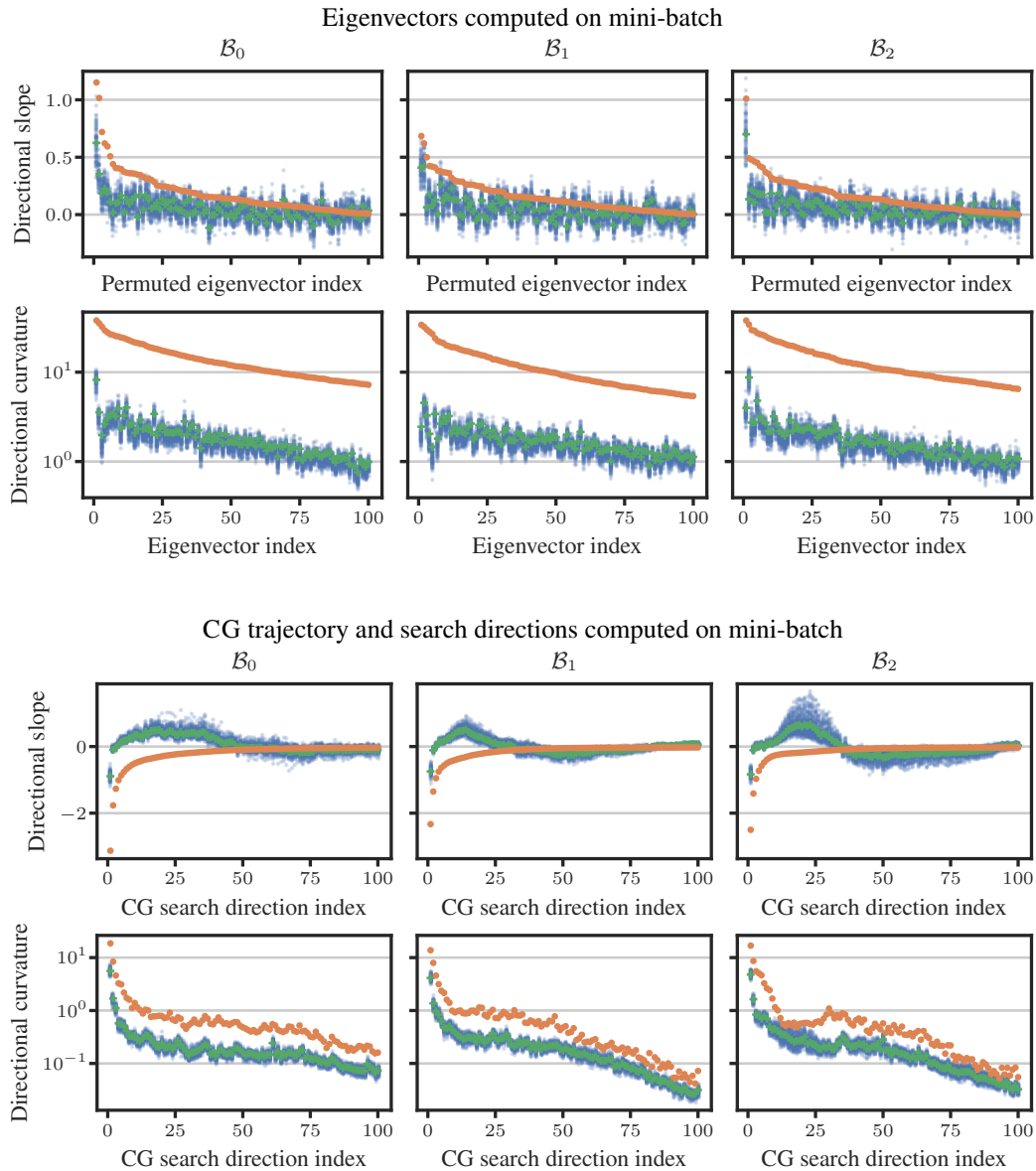


Figure 9: Directional slopes and curvatures are biased. The experimental setting is the same as in Figure 2: We use the GGN curvature proxy $H_B \leftarrow G_B + \beta I$ at batch size 512. The upper plot (*Top*) shows the directional derivatives along the top 100 eigenvectors of H_B , the lower plot (*Bottom*) shows the directional derivatives along the first 100 CG search directions. For the top panel of the upper plot, we switch the order and sign of the eigenvectors such that the orange dots are all above zero and in descending order. There are strong systematic biases in the directional slopes and curvatures. The curvature biases are more pronounced along the eigenvectors, whereas the biases in the slope are larger along the CG directions.

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

Biases: Additional results for $H_B \leftarrow G_B + \beta I$ (batch size 2048)

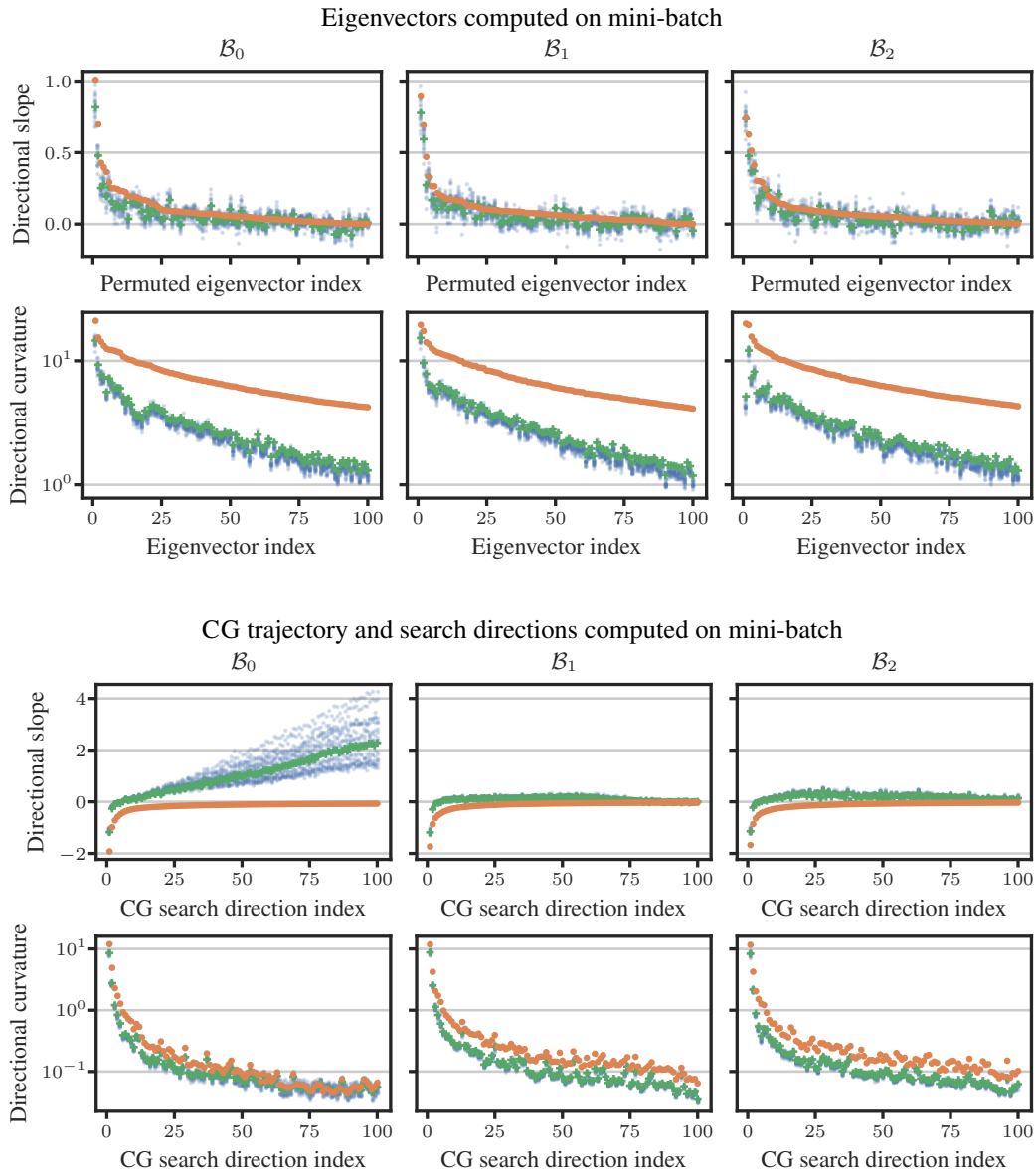
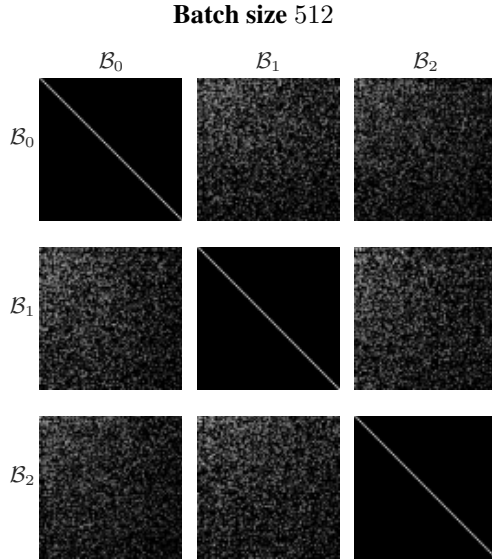


Figure 10: Directional slopes and curvatures are biased. The experimental setting is similar to Figure 2: We use the GGN curvature proxy $H_B \leftarrow G_B + \beta I$ but batch size 2048. The upper plot (Top) shows the directional derivatives along the top 100 eigenvectors of H_B , the lower plot (Bottom) shows the directional derivatives along the first 100 CG search directions. For the top panel of the upper plot, we switch the order and sign of the eigenvectors such that the orange dots are all above zero and in descending order. There are strong systematic biases in the directional slopes and curvatures. The curvature biases are more pronounced along the eigenvectors, whereas the biases in the slope are larger along the CG directions.

1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

Eigenspace overlaps: Additional result for $H_B \leftarrow \nabla^2 \mathcal{L}(\theta_*, \mathcal{B}) + \beta I$



Eigenspace overlaps: Additional results for $H_B \leftarrow G_B + \beta I$

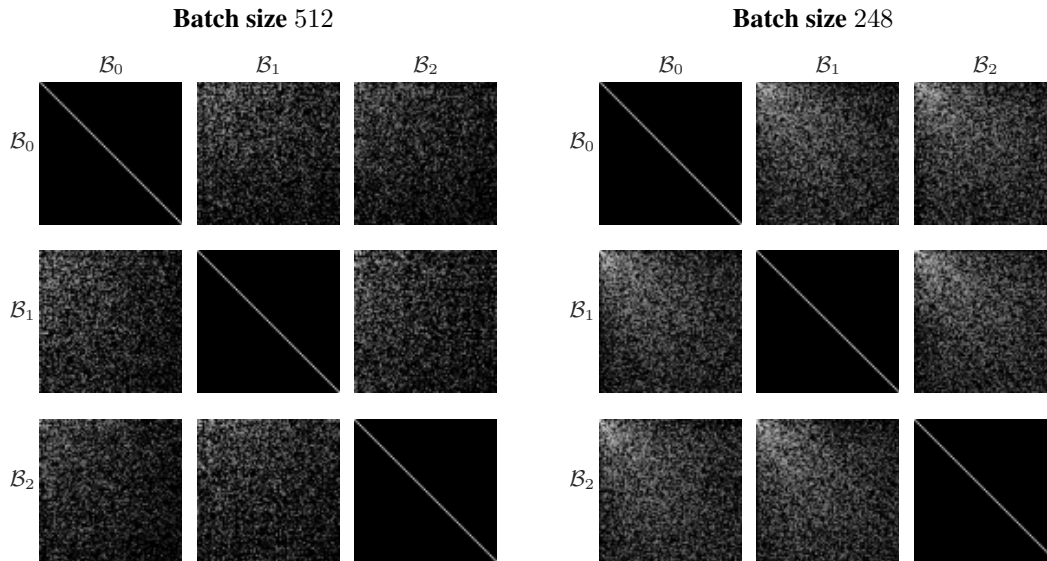


Figure 11: In practice, eigenspaces are misaligned. The experimental setting is similar to Figure 3. The top plot shows overlaps between three curvature matrices that use the Hessian of the empirical risk, i.e. $H_B \leftarrow \nabla^2 \mathcal{L}(\theta_*, \mathcal{B}) + \beta I$ at batch size 512. The bottom plots use the GGN curvature proxy $H_B \leftarrow G_B + \beta I$ at batch size 512 (Left) and 2048 (Right). The weights $\Omega_{i,j}$ are shown as a 100×100 greyscale image (color ranges from black for $\Omega_{i,j} \leq 10^{-8}$ to white for $\Omega_{i,j} = 1$) for $m, m' \in \{0, 1, 2\}$. Clearly, the eigenspaces for different mini-batches are not perfectly aligned as eigenvectors from \mathcal{B}_m overlap with several eigenvectors from $\mathcal{B}_{m'}$. At the larger batch size, the eigenspaces are more aligned.

distribution, the full-batch CG approach can be seen as an instance of the mini-batch approach with a very large mini-batch size. Thus, it might also exhibit the associated biases we describe in Section 3.2. Another explanation would be the approximation error in Equation (2): $\mathcal{L}_{\text{reg}}(\cdot; \mathcal{D}) \approx q(\cdot; \mathcal{D})$. If the trajectory moves far away from θ_* , it leaves the region where the approximation is valid (whereas the debiased full-batch version might take smaller steps and thus remains stable). We leave it to future work to investigate this further.

Additional results. We also provide the results for two other mini-batch sizes $|\mathcal{B}| \in \{512, 2048\}$ (the debiased approach uses two mini-batches of size 256 and 1024, respectively). The results are shown in Figure 13.

Runtime and memory consumption. Our current implementation uses the naive approach that requires *two* matrix-vector products with the GGN for computing one debiased update magnitude (see Appendix A.2 for details). As this would skew the runtime comparison, we refrain from providing concrete numbers here. However, as we describe in detail in Appendix A.2, it is possible to implement the debiased approach in a more efficient way that only requires one additional matrix-vector product for each CG direction to compute the debiased update magnitude. Using the debiased approach at half the mini-batch size, this brings the computational costs down to roughly the same level as the standard approach (assuming that the runtime of a GGN-vector product scales linearly with the mini-batch size).

As explained in detail in Appendix A.2, the memory overhead for the efficiently implemented debiased CG approach is two additional vectors of size P . For the experimental setting used for Figure 5, the number of parameters is $P = 1,387,108$, which corresponds to 5.55 MB in single precision. The computational overhead incurred by the debiased approach is thus about 11.1 MB.

B.6 SECTION 6.2 (DEBIASED LAPLACE APPROXIMATION): FIGURE 6

Figure 6 shows a comparison of the vanilla model, the mini-batch K-FAC LA, the debiased version, and the full-batch LA. Here, we describe the experiment in more detail and present additional results, in particular on OOD data. For the derivation and the mathematical details of these approaches, see Appendix A.3. The experiment uses the fully trained ALL-CNN-C model on CIFAR-10 data from test problem (B) (see Appendix B.1).

Experimental procedure. The experimental procedure consists of three steps: The computation of the K-FAC curvature approximations, the evaluation of the corresponding predictive class probabilities, and the computation of the performance metrics.

1. **K-FAC curvature approximations.** We consider a log-equidistant grid of 13 prior precisions β between 10^{-4} and 10^0 and add $\beta = 10$.

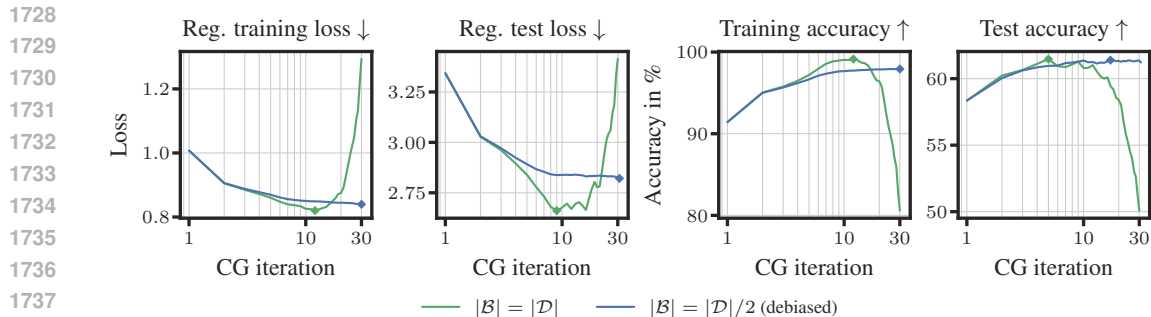
For each of those 14 values, we compute the K-FAC curvature approximation via BACKPACK (Dangel et al., 2020) (see Appendix B.2) using (i) a single batch of size 64, 256, and 1024, (ii) the two-batch debiased LA version (details in Appendix A.3.3) at batch sizes 32, 128, and 512, and (iii) the full-batch LA. For the latter, we accumulate mini-batch K-FAC approximations over the entire training set. As an additional baseline, we consider (iv) the vanilla model (without LA), which is independent of the prior precision β . For approaches (i) and (ii), we repeat the experiment for 5 different mini-batches/mini-batch pairs.

2. **Evaluation of predictive class probabilities.** The first step in the evaluation of the performance metrics is the computation of the predictive uncertainty. For a given test dataset of size N_\diamond , these can be represented as a matrix $\mathbf{P} \in \mathbb{R}^{N_\diamond \times C}$, where $\mathbf{P}_{n,c}$ is the probability that the n -th sample belongs to class c , i.e. rows of \mathbf{P} sum to 1. In all experiments, we draw $S = 40$ MC samples $\{\theta^{(s)}\}_{s=1}^S$ from the weight posterior following the procedure in Appendix A.3.2. The predictive class probabilities are then obtained via Equation (18). In the case of the vanilla model without LA, the sum in Equation (18) collapses to a single term corresponding to the MAP model.

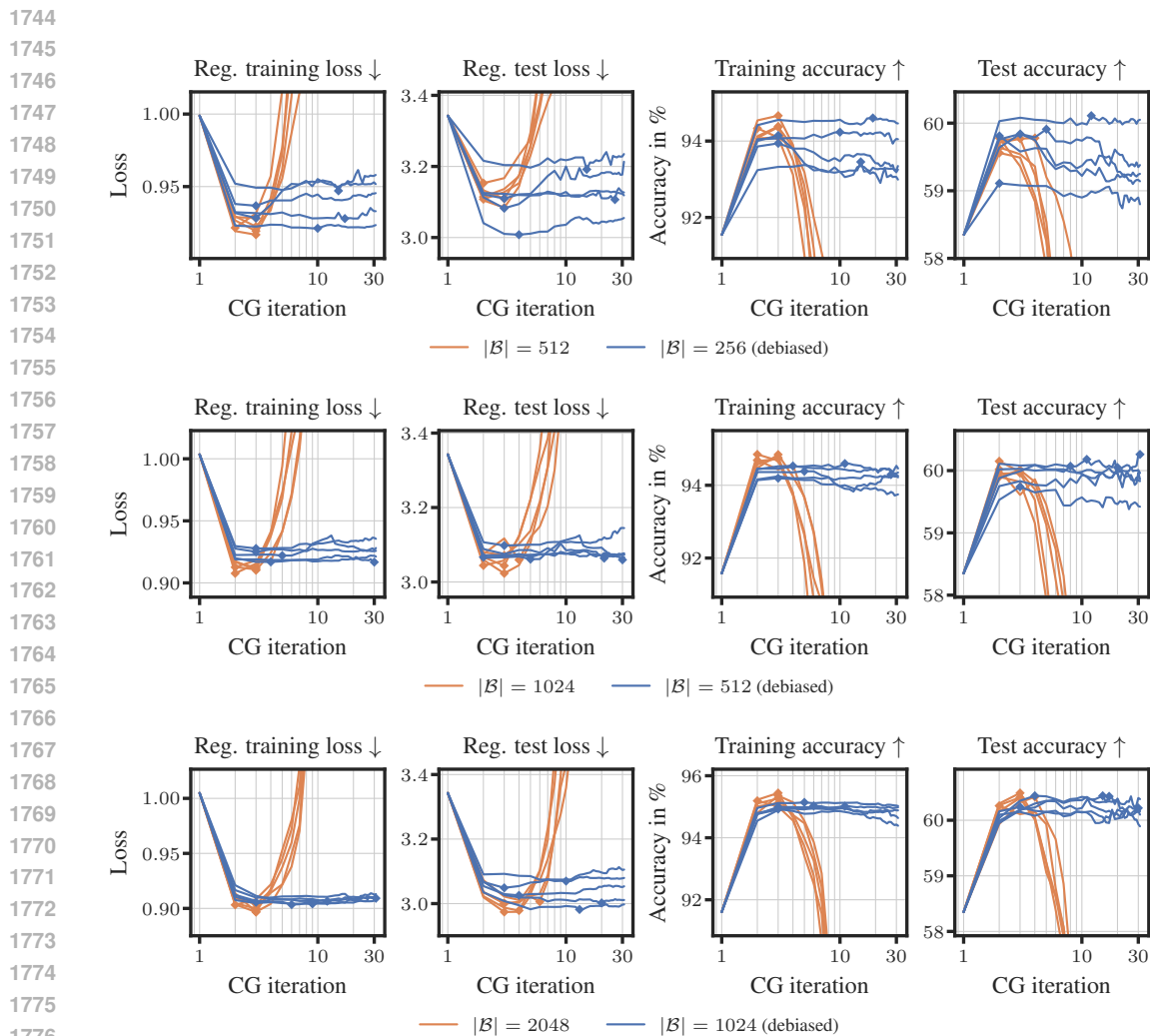
The evaluation procedure above can be applied to arbitrary test datasets. We consider the training data, test data and CIFAR-10-C (i.e. OOD datasets at severity 1 to 5).

3. **Performance metrics.** We consider the following performance metrics:

- **Accuracy.** The predictive classes are obtained from \mathbf{P} by extracting the class with the highest probability. The accuracy is the relative number of correctly classified samples.



1739 **Figure 12: Comparison of full-batch CG approaches.** We use the setting from Figure 5. The
1740 full-batch approach (shown as —) applies standard CG to $q(\cdot; \mathcal{D})$ while the debiased approach
1741 (shown as —) uses one half of the training data for the directions and the other half for the update
1742 magnitudes. The markers \blacklozenge and \blacklozenge are placed at peak performance. Surprisingly, the full-batch run
1743 diverges. The debiased CG run is stable.



1777 **Figure 13: Comparison of CG approaches at different batch sizes.** We use the same experimental
1778 setting as in Figure 5 but consider different mini-batch sizes $|\mathcal{B}| \in \{512, 1024, 2048\}$ (from (Top) to
1779 (Bottom)). The debiased approach uses two mini-batches of size 256, 512, and 1024, respectively.
1780 Across all mini-batch sizes, the single-batch CG runs diverge quickly while our debiased approach
1781 maintains stability.

- **Negative log-likelihood (NLL).** For each datum in the test set, we compute the negative log-probability of the true class. We then average over all samples in the test dataset. This coincides with the empirical risk from Equation (1), evaluated on the test dataset.
- **Expected calibration error (ECE).** The ECE is a measure of the calibration of the model’s predictive probabilities. It groups the classification confidences (i.e. the maximum entry in each of P ’s rows) into bins, and within these bins, compares the average confidence with the actual accuracy. We use `MulticlassCalibrationError` from `torchmetrics` (Detlefsen et al., 2022).
- **AUROC.** For the OOD datasets, we provide the Area Under the Receiver Operating Characteristic curve (AUROC). Our goal is to distinguish in-distribution (ID) and out-of-distribution (OOD) samples using the model’s uncertainty. We use the entropy of the predictive distribution $p(\mathbf{y}_\diamond | \mathbf{x}_\diamond, \mathbb{D})$ as our uncertainty estimate $u(\mathbf{x}_\diamond) \in \mathbb{R}$ for input \mathbf{x}_\diamond . The ground-truth binary labels are the ID/OOD indicators of the test inputs. Several thresholds ξ could be established to turn the scalar uncertainty estimates into binary predictions by using $\mathbf{1}_{\{u(\mathbf{x}_\diamond) > \xi\}}$. Instead of choosing a particular threshold and evaluating its accuracy, the AUROC metric directly evaluates $u(\cdot)$ by measuring the area under the plot of the true positive rate (TPR) against the false positive rate (FPR) for all possible thresholds.

Additional experimental details.

- **Uncertainty over weights but not biases.** As the prior acts only on the weights of the network but not its biases (see Appendix B.1), we only consider the uncertainty over the weights in the LA. This slightly reduces the size of the covariance matrix as it excludes the bias parameters.
- **Single vs. double precision.** Although the K-FAC factors are positive semi-definite by construction, their numerical eigenvalues can be negative due to numerical inaccuracies. To balance precision and computational cost, we use double precision for all computations until drawing the MC samples and single precision afterwards.

Results. The results are presented in Figures 6, 14 and 15. For the single-batch and debiased approach (where we use 5 mini-batches/mini-batch pairs each), we report the average performance as a dot and the min/max as a vertical line. The performance of the vanilla model is shown as a horizontal line as its performance does not depend on β , as explained above.

Across all datasets and performance metrics, the debiased LA mimics the behavior of the full-batch approach much better than the single-batch LA (although both approaches use a comparable amount of computational resources). In particular, for small prior precisions, where the covariance matrix relies almost exclusively on the K-FAC curvature information, the debiased LA maintains a good performance in contrast to the single-batch approach.

Additional results.

- **RESNET-50 on IMAGENET.** To showcase the scalability of our debiased LA approach, we repeat the experiment for a RESNET-50 model on the IMAGENET dataset (test problem (E) in Appendix B.1). The results are shown in Figure 16. Again, the debiased approach behaves similarly to the full-batch LA and maintains stability over the entire spectrum of prior precisions.
- **VIT LITTLE on IMAGENET.** Similar to the previous experiment, we also apply our LA debiasing approach to a VIT LITTLE model on the IMAGENET dataset (test problem (F) in Appendix B.1) to ensure that it is beneficial for other architectures as well. The results are shown in Figure 17. The debiased approach mimics the calibration behavior of the full-batch LA remarkably well and remains stable for lower prior precision values as well where the single-batch approach behaves erratically.

Runtime analysis (computational overhead of debiasing). We claim in Section 6.2 that, by using the debiased approach at half the mini-batch size that is used by the single-batch approach, we obtain a fair comparison (more details in Section 4.2.3). Here, we substantiate this claim by providing a detailed runtime comparison of the biased, debiased, and full-batch LA approaches on CIFAR-10.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

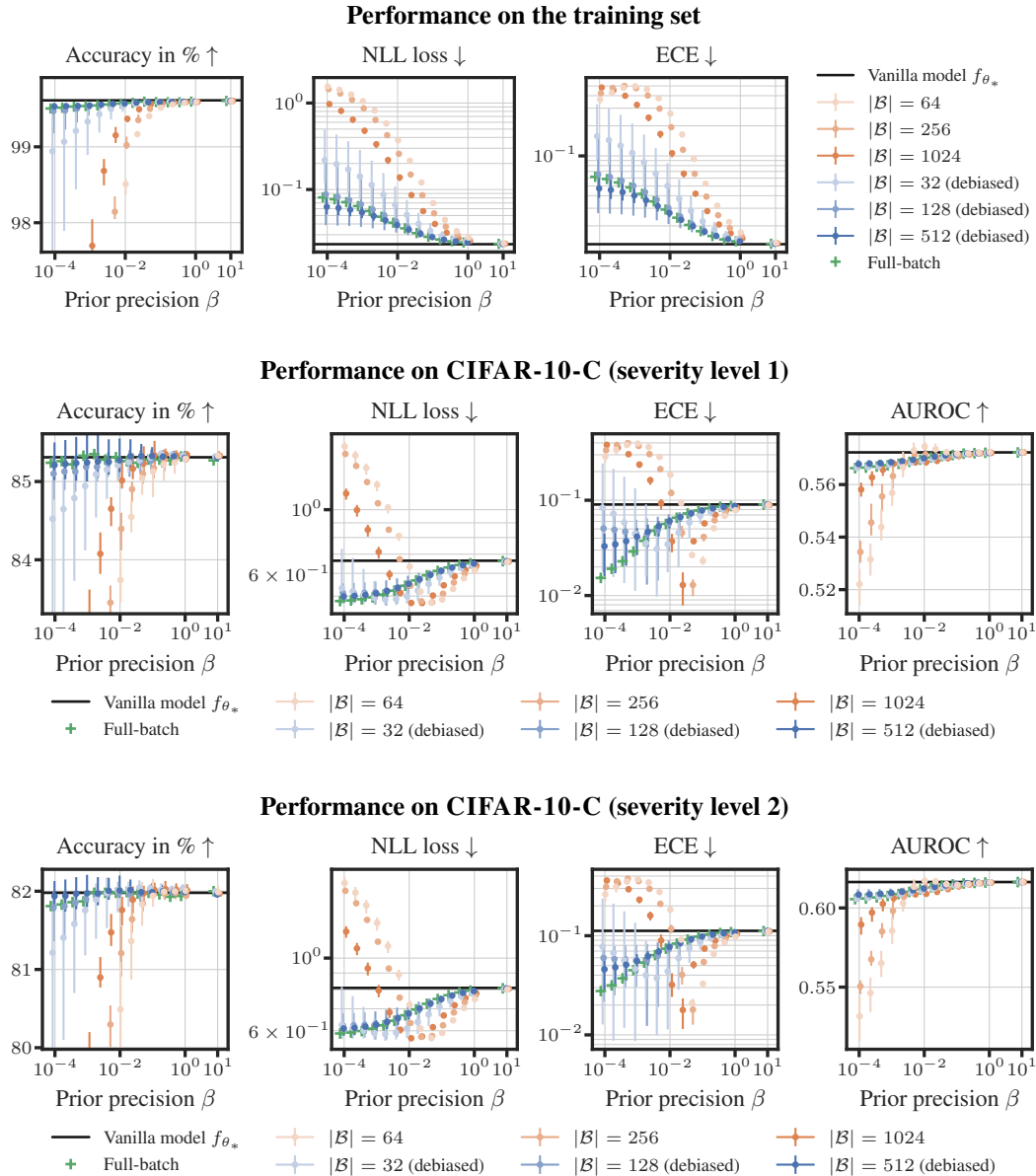


Figure 14: Debiased LA mimics the full-batch LA very well. The experimental setting is the same as in Figure 6, but we report results on additional datasets (training set and OOD datasets at severity level 1 and 2). For the OOD datasets, we report the AUROC metric in addition to the accuracy, NLL and ECE. In contrast to the single-batch approach, the debiased version mimics the behavior of the full-batch approach very well over the entire range of prior precisions and across datasets.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

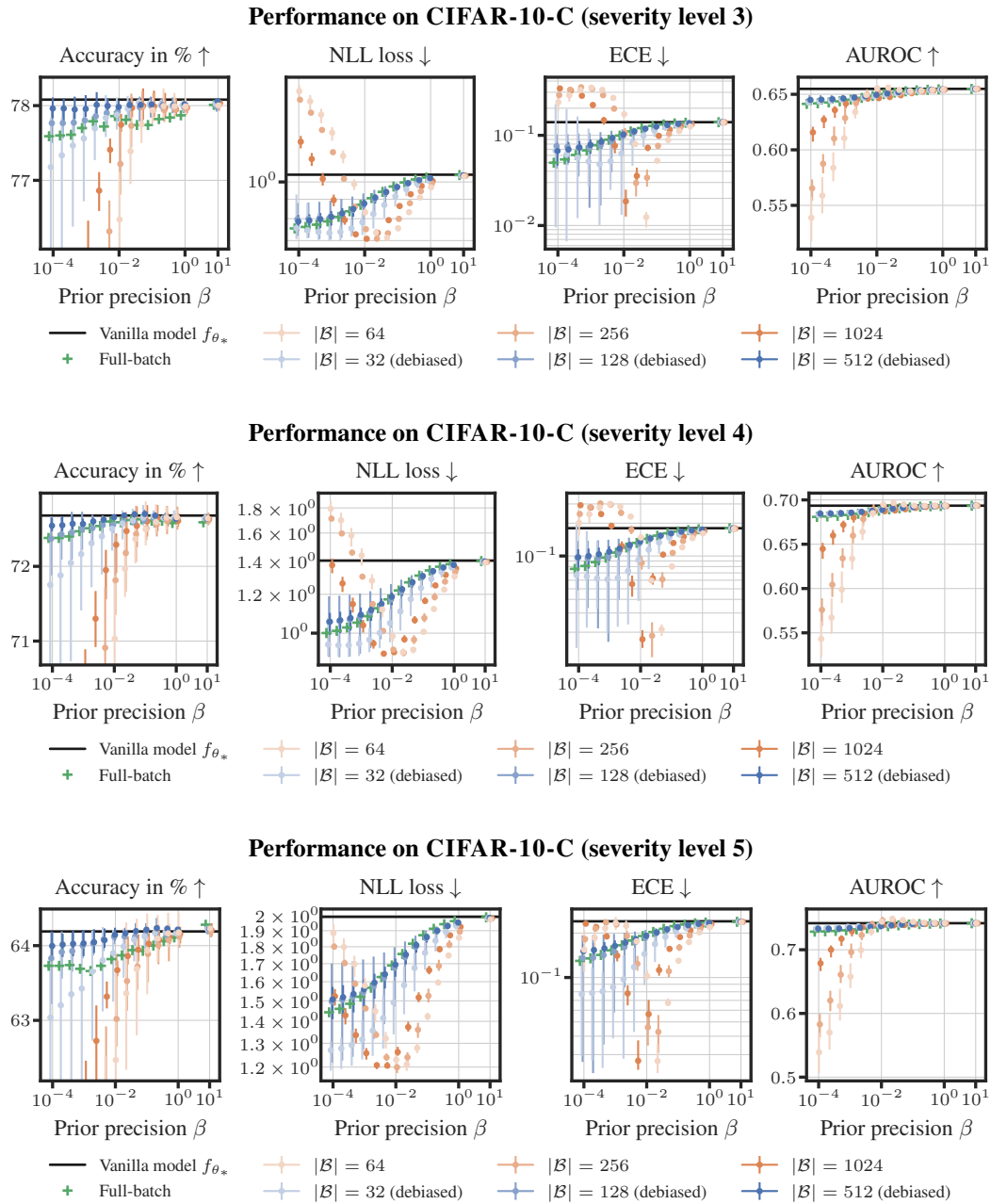


Figure 15: Debiased LA mimics the full-batch LA very well. The experimental setting is the same as in Figure 6, but we report results on additional datasets (OOD datasets at severity level 3, 4, and 5). We report the AUROC metric in addition to the accuracy, NLL, and ECE. In contrast to the single-batch approach, the debiased version mimics the behavior of the full-batch approach very well over the entire range of prior precisions and across datasets.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

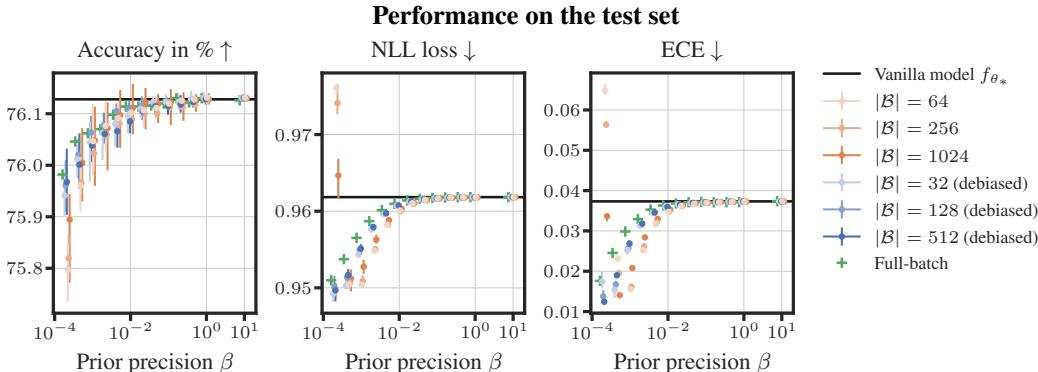


Figure 16: Our approach scales to RESNET-50 on IMAGENET. The experimental setting is the same as in Figure 6, but we use test problem (E) (see Appendix B.1) (RESNET-50 on IMAGENET). The results are consistent with the findings on CIFAR-100: The unbiased approach behaves similarly to the full-batch LA and maintains stability even for small prior precisions.

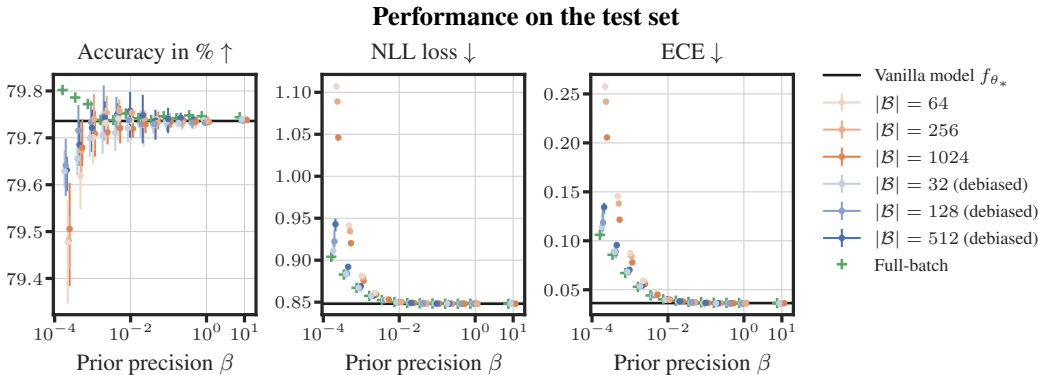


Figure 17: Curvature biases are also present for the ViT LITTLE architecture on IMAGENET. The experimental setting is the same as in Figure 6 and Figure 16, but we use test problem (F) (see Appendix B.1) (ViT LITTLE on IMAGENET). Again, the unbiased approach behaves similarly to the full-batch LA and maintains stability even for small prior precisions (where the covariance is almost exclusively based on the K-FAC curvature information).

Results (runtimes). We report the runtimes of the biased (orange square), unbiased (blue square), and full-batch (green square) K-FAC Laplace approximation schemes in Figure 18. We find that (i) both biased and unbiased approaches are orders of magnitude cheaper than the full-batch LA, and (ii) the unbiased approach introduces only negligible overhead and can, sometimes, even improve upon the runtime of the biased approach.

Memory consumption. The ALL-CNN-C model we use for the LA experiment in Section 6.2 has $P = 1,368,480$ trainable parameters (this amounts to 10.95 MB of memory in double precision). The K-FAC approximation requires storing 11,483,965 (91.87 MB) numbers to represent all Kronecker factors. The memory of a K-FAC approximation is thus roughly equivalent to that of 8.4 models. As the unbiased approach is based on *two* such approximations in the case of a naive implementation, the overhead of unbiased is another 91.87 MB of memory. However, this can be significantly improved by building up the K-FAC approximations block by block: Having two corresponding blocks available (based on two different mini-batches) already suffices to compute the respective unbiased block (see Appendix A.3.3). This way, we only have two blocks (represented by their Kronecker factors) in memory at the same time (instead of two entire K-FAC approximations) which greatly reduces the memory overhead of unbiased, down to the largest block. For the ALL-CNN-C model, this is two Kronecker factors of sizes 192×192 and 1728×1728 amounting to 3,022,848 numbers in total, i.e. the equivalent of 2.20 models or 24.18 MB in double precision. The statement that the unbiased approach roughly doubles the memory consumption is thus a worst-case scenario.

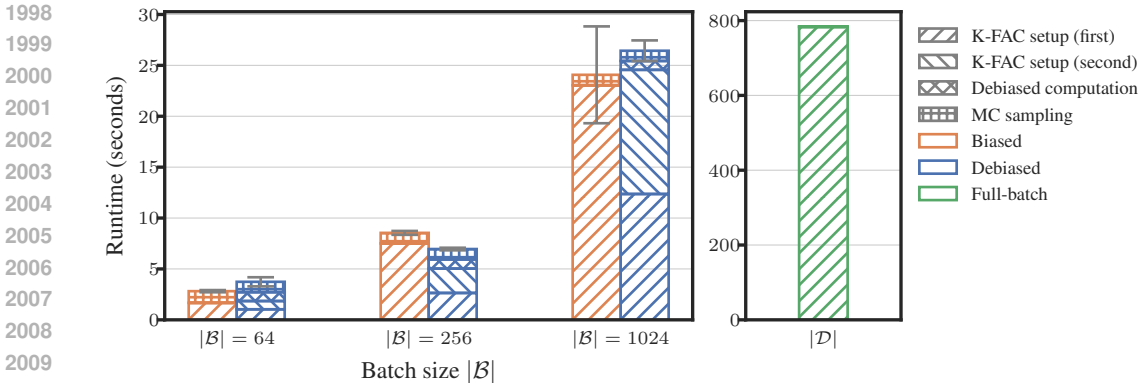


Figure 18: The computational overhead of debiasing is negligible for Laplace approximations. We follow the experimental setup of Section 6.2 and provide a detailed runtime comparison of the biased (orange), debiased (blue), and full-batch (green) K-FAC Laplace approximation schemes on CIFAR-10. The total runtime includes the construction of K-FAC’s Kronecker factors (on one or two mini-batches for the biased and debiased approach, respectively), the debiasing computations (see Appendix A.3.3), and the Monte Carlo sampling of $S = 40$ weights. All reported runtimes are averages over 5 runs. The error bars for the total runtime cover one standard deviation. The results show that the debiased approach introduces a negligible overhead. Compared to the full-batch variant, both mini-batch approaches are orders of magnitude faster. Note that the evaluation (i.e., the computation of the predictive uncertainty) is not included in the comparison, as the time requirements are the same for all approaches; however, this step requires the most computational resources.

B.7 ADDITIONAL EXPERIMENT: BIASES FOR A WIDERESNET 40-4 MODEL ARCHITECTURE

Experimental setting and results. Here, we extend our analysis from Section 3.1 (details in Appendix B.4) to test problem (C) (see Appendix B.1): A WIDE RESNET model on CIFAR-100. We use the GGN curvature proxy $\mathbf{H}_{\mathcal{B}} \leftarrow \mathbf{G}_{\mathcal{B}} + \beta \mathbf{I}$ at batch size 256 and compute the directional slopes and curvatures along the top 100 eigenvectors of $\mathbf{H}_{\mathcal{B}}$. The results are shown in Figure 19. Along the first few top-curvature directions, the curvature biases are even larger for this test problem than for the ALL-CNN-C model on CIFAR-100, but then they decay more quickly.

B.8 ADDITIONAL EXPERIMENT: K-FAC AND THE DEPENDENCE OF THE BIASES ON MINI-BATCH SIZE

Experimental setting. Here, we extend our analysis of test problem (A) (see Appendix B.1) from Section 3.1 to K-FAC and investigate the impact of the mini-batch size on the curvature biases.

Results & discussion. Figure 20 shows the directional curvatures of the K-FAC approximation for four different mini-batch sizes. When we compute the eigenvectors and directional curvatures on the same mini-batch, we observe a systematic curvature bias that decreases with increasing mini-batch size. There are two phenomena at play: With increasing mini-batch size, (i) the green crosses move upwards and (ii) the orange dots move downwards. Our intuition for this is as follows: With increasing mini-batch size, the eigenvectors become more meaningful such that, on other data, they also exhibit large curvature—this explains (i). Similarly, with increasing mini-batch size, it gets harder to find directions of *extreme* curvature as these directions have to exhibit large curvature on *all* data points within the mini-batch—this explains (ii).

B.9 ADDITIONAL EXPERIMENT: DEVELOPMENT OF THE BIASES OVER THE COURSE OF TRAINING

Experimental setting and results. Here, we investigate how the biases in the slope and curvature develop over the course of training for test problem (A) (see Appendix B.1). We use the same procedure as in Appendix B.4 but evaluate the biases at 10 different checkpoints during training (spread log-equidistantly between the first and last epoch). At each checkpoint, we draw 5 mini-

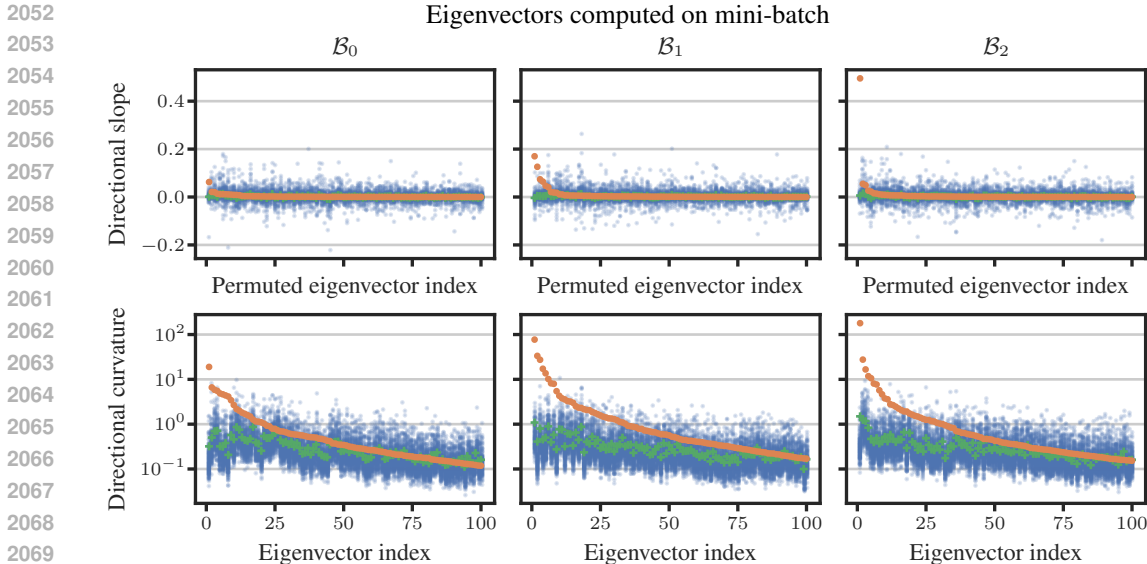


Figure 19: Directional slopes and curvatures are biased. The experimental setting is similar to Figure 2, but we use test problem (C) (WIDE RESNET on CIFAR-100) with the GGN curvature proxy $H_B \leftarrow G_B + \beta I$ at batch size 256. For the top panel, we switch the order and sign of the eigenvectors such that the orange dots are all above zero and in descending order. There is a strong, systematic bias, particularly in the curvature: Computing the eigenvectors and directional curvatures on the same data results in over-estimation that is very pronounced along the first few eigenvectors but then decays more quickly than in test problem (A).

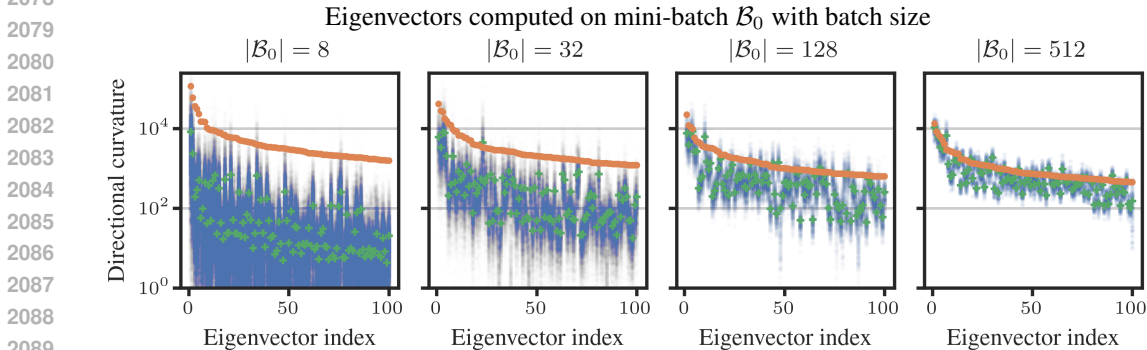


Figure 20: Directional curvatures with K-FAC. We use the CIFAR-100 dataset with the fully trained ALL-CNN-C model. For each mini-batch size $\in \{8, 32, 128, 512\}$, we draw one mini-batch and compute the top 100 eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_{100}$ of the K-FAC matrix \mathbf{K}_{B_0} . We show the directional curvatures $\mathbf{u}_p^\top \mathbf{K}_{B_m} \mathbf{u}_p, p \in \{1, \dots, 100\}$ on the same mini-batch ($m = 0$) as \bullet , all other mini-batches in the training set ($m \in \{1, 2, \dots\}$) as \bullet and their average as $\+$. Similar to Figure 2, we observe a systematic bias in the curvature. The bias decreases with increasing mini-batch size.

batches, compute the top 100 eigenvectors of the corresponding GGN-based quadratics, and finally evaluate the *relative* errors in the directional slopes and curvatures (where the ground truth is the full-batch quadratic’s slope/curvature). For each mini-batch, this distribution of 100 relative curvature and slope biases is represented as a dot (at the mean relative bias) and a vertical line ranging from the 25% to the 75% percentile in Figure 21. While the biases in the slope remain relatively stable over the course of training, the biases in the curvature increase over more than 3 orders of magnitude up to a relative error of order 10 in epoch 349 (which is consistent with the results in Figure 2). This suggests that the eigenspaces of different curvature matrices become more and more misaligned as training progresses steadily increasing the need for effective debiasing strategies.

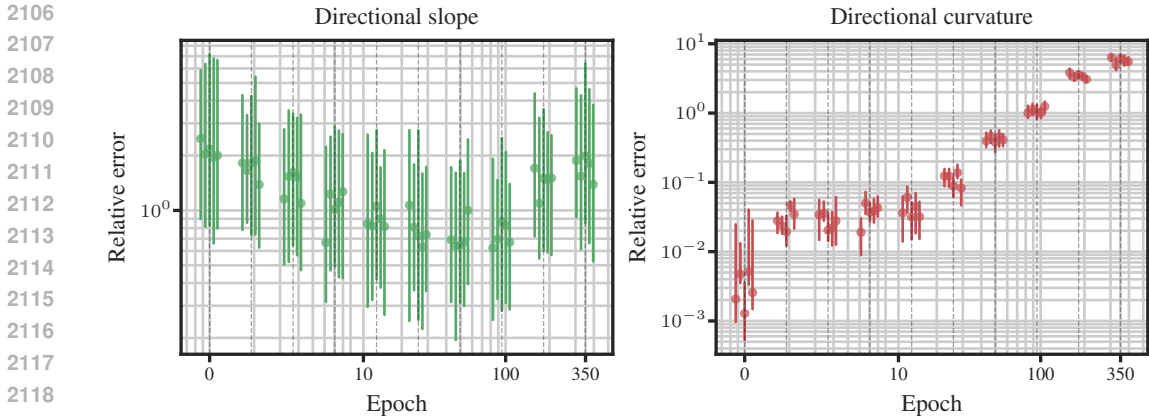
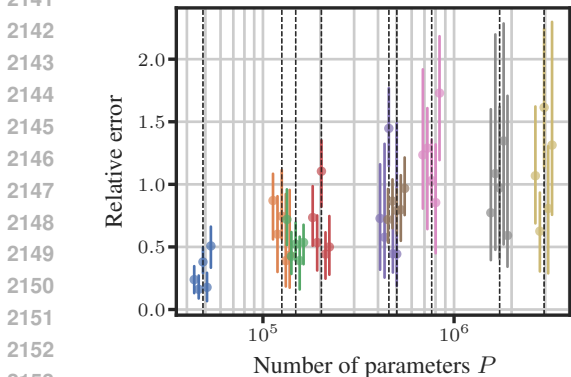


Figure 21: Curvature biases increase over the course of training. We evaluate the relative slope and curvature biases at different checkpoints during training for the ALL-CNN-C model on CIFAR-100. 5 mini-batches are drawn per checkpoint. For each mini-batch, the relative biases are represented as a dot ● (at the mean relative bias) and a vertical line — ranging from the 25% to the 75% percentile. While the biases in the slope remain relatively stable over the course of training, the biases in the curvature increase (over more than 3 orders of magnitude).

B.10 ADDITIONAL EXPERIMENT: THE CURVATURE BIAS INCREASES WITH P

In high-dimensional spaces, it becomes increasingly unlikely that random vectors are aligned. While the eigenspaces of the curvature matrices are not completely random, they are subject to noise. It is thus conceivable that their overlap decreases as the number of parameters P increases. If this hypothesis is true, the curvature biases should become more pronounced in large models.

Experimental procedure & results. To test this hypothesis, we use test problem (D) that implements a simple convolutional neural network with variable width and depth (for details, see Appendix B.1). For each fully trained model, we evaluate the relative error between $\partial_{\mathbf{u}_p}^2 q(\boldsymbol{\theta}_*; \mathcal{B})$ and $\partial_{\mathbf{u}_p}^2 q(\boldsymbol{\theta}_*; \mathcal{D})$ for the $\mathbf{G}_{\mathcal{B}}$'s top 100 eigenvectors at batch size 128. This procedure is repeated for 5 different mini-batches for each of the 9 models, resulting in a total of 45 error distributions (see Figure 22). The results confirm our hypothesis: The relative errors tend to increase with the number of parameters P . In the massively overparameterized regime, the biases might thus become even more relevant and effective debiasing strategies are needed.



● model 1-32 ● model 7-32 ● model 7-64
 ● model 4-32 ● model 4-64 ● model 4-128
 ● model 1-64 ● model 1-128 ● model 7-128

Figure 22: Increasing curvature biases with P . We train 9 convolutional neural networks with different widths and depths for 100 epochs on the CIFAR-10 dataset using ADAM with standard hyperparameters. For the fully trained model, we evaluate the relative error $|\partial_{\mathbf{u}_p}^2 q(\boldsymbol{\theta}_*; \mathcal{B}) - \partial_{\mathbf{u}_p}^2 q(\boldsymbol{\theta}_*; \mathcal{D})| \cdot |\partial_{\mathbf{u}_p}^2 q(\boldsymbol{\theta}_*; \mathcal{D})|^{-1}$ for $\mathbf{G}_{\mathcal{B}}$'s top 100 eigenvectors at batch size $|\mathcal{B}| = 128$. 5 different mini-batches are used per model resulting in a total of 45 error distributions (each consisting of 100 numbers). These distributions are represented by their median (as a dot ●) and the 25% and 75% percentiles (as a line segment —) P . The 5 distributions for each model are slightly spread along the x -axis for better visibility. The experiment confirms our hypothesis: The relative errors tend to increase with the number of parameters P .