

# VIVIT: CURVATURE ACCESS THROUGH THE GENERALIZED GAUSS-NEWTON’S LOW-RANK STRUCTURE

## SUPPLEMENTARY MATERIAL

<b>A Mathematical details</b>	<b>12</b>
A.1 Reducing the GGN eigenvalue problem to the Gram matrix . . . . .	12
A.2 Relation between GGN and Gram matrix eigenvectors . . . . .	12
<b>B Experimental details</b>	<b>13</b>
B.1 Performance evaluation . . . . .	13
B.2 Noise analysis during training . . . . .	15
B.3 Bootstrap damping for second-order methods . . . . .	19
B.4 Derivation of noisy quadratic model . . . . .	19
B.5 Experiment on noisy quadratic . . . . .	23
<b>C Implementation details</b>	<b>23</b>
C.1 Optimized Gram matrix computation for linear layers . . . . .	24
C.2 Implicit multiplication with the inverse (block-diagonal) GGN . . . . .	24

## A MATHEMATICAL DETAILS

### A.1 REDUCING THE GGN EIGENVALUE PROBLEM TO THE GRAM MATRIX

For Equation (4), consider the left hand side of the GGN’s characteristic polynomial  $\det(\mathbf{G} - \lambda \mathbf{I}_D) = 0$ . Inserting the ViViT factorization (Equation (3)) and using the matrix determinant lemma yields

$$\begin{aligned}
 \det(-\lambda \mathbf{I}_D + \mathbf{G}) &= \det(-\lambda \mathbf{I}_D + \mathbf{V} \mathbf{V}^\top) && \text{(Low-rank structure (3))} \\
 &= \det(\mathbf{I}_{NC} + \mathbf{V}^\top (-\lambda \mathbf{I}_D)^{-1} \mathbf{V}) \det(-\lambda \mathbf{I}_D) && \text{(Matrix determinant lemma)} \\
 &= \det\left(\mathbf{I}_{NC} - \frac{1}{\lambda} \mathbf{V}^\top \mathbf{V}\right) (-\lambda)^D \\
 &= \left(-\frac{1}{\lambda}\right)^{NC} \det(\mathbf{V}^\top \mathbf{V} - \lambda \mathbf{I}_{NC}) (-\lambda)^D \\
 &= (-\lambda)^{D-NC} \det(\tilde{\mathbf{G}} - \lambda \mathbf{I}_{NC}) . && \text{(Gram matrix)}
 \end{aligned}$$

Setting the above expression to zero reveals that the GGN’s spectrum decomposes into  $D - NC$  zero eigenvalues and the Gram matrix spectrum obtained from  $\det(\tilde{\mathbf{G}} - \lambda \mathbf{I}_{NC}) = 0$ .

### A.2 RELATION BETWEEN GGN AND GRAM MATRIX EIGENVECTORS

Assume the nontrivial Gram matrix spectrum  $\tilde{\mathbb{S}}_+ = \{(\lambda_k, \tilde{\mathbf{e}}_k) \mid \lambda_k \neq 0, \tilde{\mathbf{G}} \tilde{\mathbf{e}}_k = \lambda_k \tilde{\mathbf{e}}_k\}_{k=1}^K$  with orthonormal eigenvectors  $\tilde{\mathbf{e}}_j^\top \tilde{\mathbf{e}}_k = \delta_{jk}$  ( $\delta$  represents the Kronecker delta) and  $K = \text{rank}(\mathbf{G})$ . We now show that  $\mathbf{e}_k = 1/\sqrt{\lambda_k} \mathbf{V} \tilde{\mathbf{e}}_k$  are normalized eigenvectors of  $\mathbf{G}$  and inherit orthogonality from  $\tilde{\mathbf{e}}_k$ .

To see the first, consider right-multiplication of the GGN with  $e_k$ , then expand the low-rank structure,

$$\begin{aligned}
Ge_k &= \frac{1}{\sqrt{\lambda_k}} V V^\top V \tilde{e}_k && \text{(Equation (3) and definition of } e_k) \\
&= \frac{1}{\sqrt{\lambda_k}} V \tilde{G} \tilde{e}_k && \text{(Gram matrix)} \\
&= \lambda_k \frac{1}{\sqrt{\lambda_k}} V \tilde{e}_k && \text{(Eigenvector property of } \tilde{e}_k) \\
&= \lambda_k e_k.
\end{aligned}$$

Orthonormality of the  $e_k$  results from the Gram matrix eigenvector orthonormality,

$$\begin{aligned}
e_j^\top e_k &= \left( \frac{1}{\sqrt{\lambda_j}} \tilde{e}_j^\top V^\top \right) \left( \frac{1}{\sqrt{\lambda_k}} V \tilde{e}_k \right) && \text{(Definition of } e_j, e_k) \\
&= \frac{1}{\sqrt{\lambda_j \lambda_k}} \tilde{e}_j^\top \tilde{G} \tilde{e}_k && \text{(Gram matrix)} \\
&= \frac{\lambda_k}{\sqrt{\lambda_j \lambda_k}} \tilde{e}_j^\top \tilde{e}_k && \text{(Eigenvector property of } \tilde{e}_k) \\
&= \delta_{jk}. && \text{(Orthonormality)}
\end{aligned}$$

## B EXPERIMENTAL DETAILS

### B.1 PERFORMANCE EVALUATION

**Hardware information:** Results presented in this section, as well as Section 2, were generated on a workstation with the following hardware:

- **CPU:** Intel Core i7-8700K CPU @ 3.70 GHz  $\times$  12 (32 GB)
- **GPU:** NVIDIA GeForce RTX 2080 Ti (11 GB)

We will use their shorthands to indicate the device that executed the computation.

**Settings:** Performance is evaluated with different GGN approximations, parameterized by the used mini-batch samples (full, frac), and the backpropagated loss Hessian representation (exact, MC):

- **exact, full:** Backpropagate the exact loss Hessian representation for all mini-batch samples ( $NC$  vectors).
- **MC, full:** Backpropagate an MC approximation of the loss Hessian (using a single MC sample) for all mini-batch samples ( $N$  vectors).
- **exact, frac:** Backpropagate the exact loss Hessian representation for a fraction ( $1/8$ , as in Zhang et al. (2017)) of mini-batch samples ( $N/8C$  vectors).
- **exact, MC:** Backpropagate an MC approximation of the loss Hessian (using a single MC sample) for a fraction ( $1/8$ , as in Zhang et al. (2017)) of mini-batch samples ( $N/8$  vectors).

In addition to computing the target quantity (GGN spectrum, damped Newton step) with BACKPACK, a standard gradient backpropagation on the full mini-batch is always performed in PYTORCH’s backward pass. Performance is evaluated on convolutional neural nets from DEEPOBS (Schneider et al., 2019): 3C3D on CIFAR-10, 2C2D on FASHION-MNIST, and ALL-CNN-C on CIFAR-100.

**GGN spectra:** To obtain the spectra of Figure 1, Figure S.4, and Figure S.5, we initialize the respective architecture, then draw a mini-batch and evaluate the GGN eigenvalues under the described approximations, clipping the Gram matrix eigenvalues at  $10^{-4}$ . Mini-batch sizes correspond to the default value for training where possible (CIFAR-10 3C3D:  $N = 128$ , FASHION-MNIST 2C2D:  $N = 128$ ). Only on CIFAR-100 ALL-CNN-C (trained with  $N = 256$ ), we reduce the batch size to  $N = 64$  to fit the exact computation on the full mini-batch, used as baseline, into memory.

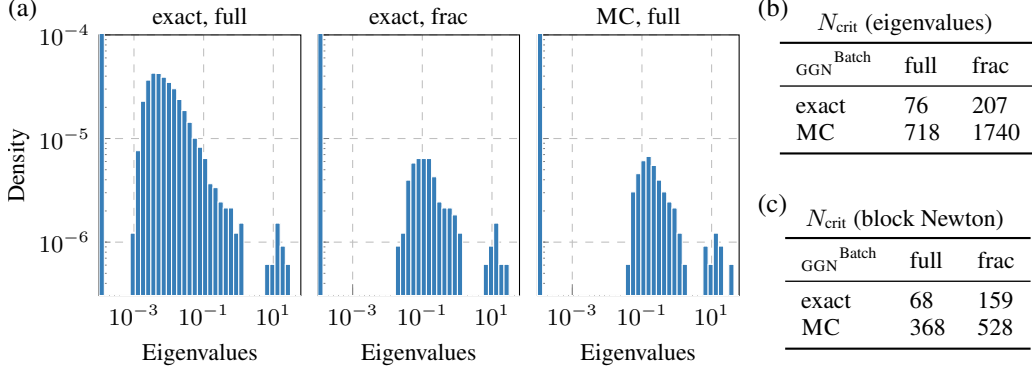


Figure S.4: **Scalability evaluation on FASHION-MNIST 2C2D** ( $D = 3,274,634$ ,  $C = 10$ ). (a) GGN eigenvalue distribution with different costs on a mini-batch of size  $N = 128$ . From left to right: Exact GGN on the full batch, exact GGN on a batch fraction ( $1/8$ , as in Zhang et al. (2017)), MC-approximation of the GGN on the full batch. (b) Maximum batch size  $N_{\text{crit}}$  (GPU) for a standard gradient computation and the GGN spectrum and (c) computing exact Newton steps with layer-wise parameter groups.

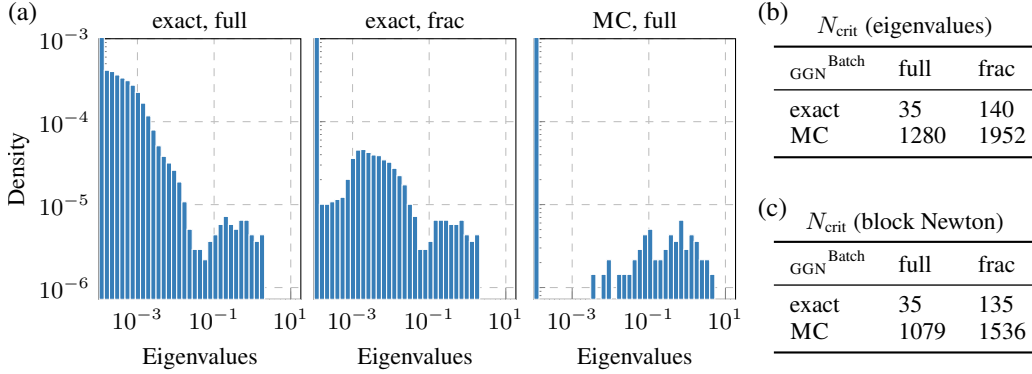


Figure S.5: **Scalability evaluation on CIFAR-100 ALL-CNN-C** ( $D = 1,387,108$ ,  $C = 100$ ). (a) GGN eigenvalue distribution with different costs on a mini-batch of size  $N = 64$ . From left to right: Exact GGN on the full batch, exact GGN on a batch fraction ( $1/8$ , as in Zhang et al. (2017)), MC-approximation of the GGN on the full batch. (b) Maximum batch size  $N_{\text{crit}}$  (GPU) for a standard gradient computation and the GGN spectrum and (c) computing exact Newton steps with layer-wise parameter groups.

**Critical batch sizes:** Similar to the GGN spectra, we repeat their computation and vary the mini-batch size until the device runs out of memory. The largest mini-batch size that can be handled is denoted as  $N_{\text{crit}}$ , the critical batch size. Figure 1b, Figure S.4b, Figure S.5b, and Table S.1a present additional results.

The critical batch sizes in Figure 1c, Figure S.4c, Figure S.5b, and Table S.1b, employ a block-diagonal GGN approximation with groups consisting of weights and bias terms in each layer (see Appendix C). For each block we compute a damped Newton step (first term in Equation (9), using Gram matrix eigenvalues larger than  $10^{-4}$ ) with constant damping  $\delta = 1$ .

As explained in Section 2, the GGN eigenvalues only require the Gram matrix. Newton steps from a block-diagonal approximation additionally require the directional derivatives (Equation (8)), which involve individual gradients. Second-order directional derivatives  $\lambda_{nk}$  (Equation (8b)) are evaluated on the same samples as the GGN eigenvectors, but we *always* use all mini-batch samples to compute the directional gradients  $\gamma_{nk}$  (Equation (8a)). As gradients are cheaper to compute, this suggests

Table S.1: **Critical batch sizes for eigenvalues and Newton steps with different approximations.** Additional results that complement Figure 1b,c, Figure S.4b,c, and Figure S.5 are shown column-wise for each architecture. From top to bottom, we report the critical batch sizes for computing (a) the GGN eigenvalue spectrum on CPU, (b) damped Newton steps with a block-diagonal GGN approximation corresponding to individual layers on CPU, and (c,d) damped Newton steps with the full GGN matrix on CPU and GPU. Interpretations and procedure details are provided in the text.

	CIFAR-10 3C3D			FASHION-MNIST 2C2D			CIFAR-100 ALL-CNN-C		
(a)	$N_{\text{crit}}$ (eigenvalues, CPU)			$N_{\text{crit}}$ (eigenvalues, CPU)			$N_{\text{crit}}$ (eigenvalues, CPU)		
	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac
	exact	1166	3004	exact	231	615	exact	94	519
	MC	8430	14908	MC	2208	5495	MC	3984	5859
(b)	$N_{\text{crit}}$ (block Newton, CPU)			$N_{\text{crit}}$ (block Newton, CPU)			$N_{\text{crit}}$ (block Newton, CPU)		
	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac
	exact	1046	2423	exact	210	487	exact	95	504
	MC	4997	6838	MC	1137	1643	MC	3360	3920
(c)	$N_{\text{crit}}$ (full Newton, CPU)			$N_{\text{crit}}$ (full Newton, CPU)			$N_{\text{crit}}$ (full Newton, CPU)		
	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac
	exact	667	2215	exact	202	487	exact	43	309
	MC	3473	5632	MC	1107	1639	MC	2015	2865
(d)	$N_{\text{crit}}$ (full Newton, GPU)			$N_{\text{crit}}$ (full Newton, GPU)			$N_{\text{crit}}$ (full Newton, GPU)		
	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac	GGN <sup>Batch</sup>	full	frac
	exact	208	727	exact	66	159	exact	13	87
	MC	1055	1816	MC	362	528	MC	640	959

evaluating them on more samples compared to curvature, see e.g. Zhang et al. (2017). The overhead thus leads to smaller critical batch sizes in comparison to computing the GGN spectrum.

For completeness, Table S.1c,d shows critical batch sizes when the block-diagonal GGN is replaced by its full representation (full Newton). In contrast to Newton steps with a block-diagonal matrix that can discard the stage-wise matrices  $\mathbf{V}^{(i)}$  during backpropagation,  $\mathbf{V}$  must now be stored until all parameters have been traversed. This leads to higher memory consumption, and hence smaller critical batch sizes, but avoids multiple Gram matrix inversions (one per parameter group).

In summary, we find that there always exists a combination of approximations which allows for critical batch sizes larger than the traditional size used for training (some architectures even permit exact computation). Different accuracy-cost trade-offs may be preferred, depending on the application and the computational budget. By the presented approximations, ViViT’s representation is capable to adapt over a wide range.

## B.2 NOISE ANALYSIS DURING TRAINING

**Procedure:** We train the following DEEPOBS (Schneider et al., 2019) architectures with SGD and ADAM: 3C3D on CIFAR-10 ( $N = 128$ ), 2C2D on FASHION-MNIST ( $N = 128$ ), and ALL-CNN-C on CIFAR-100 ( $N = 256$ ). To assert successful training, we use the hyperparameters from Dangel

Table S.2: **Hyperparameter settings for training runs to analyze noise.** For both SGD and ADAM, we report their learning rates  $\alpha$  (taken from the baselines in [Dangel et al. \(2020\)](#)) and link to their visualization. Momentum for SGD was fixed to 0.9, and ADAM uses the default parameters  $(\beta_1, \beta_2) = (0.99, 0.999)$ .

Problem	SGD	ADAM
CIFAR-10 3C3D	$\alpha \approx 3.79 \cdot 10^{-3}$ (Figure 2)	$\alpha \approx 2.98 \cdot 10^{-4}$ (Figure S.6)
FASHION-MNIST 2C2D	$\alpha \approx 2.07 \cdot 10^{-2}$ (Figure S.7)	$\alpha \approx 1.27 \cdot 10^{-4}$ (Figure S.8)
CIFAR-100 ALL-CNN-C	$\alpha \approx 4.83 \cdot 10^{-1}$ (Figure S.9)	$\alpha \approx 6.95 \cdot 10^{-4}$ (Figure S.10)

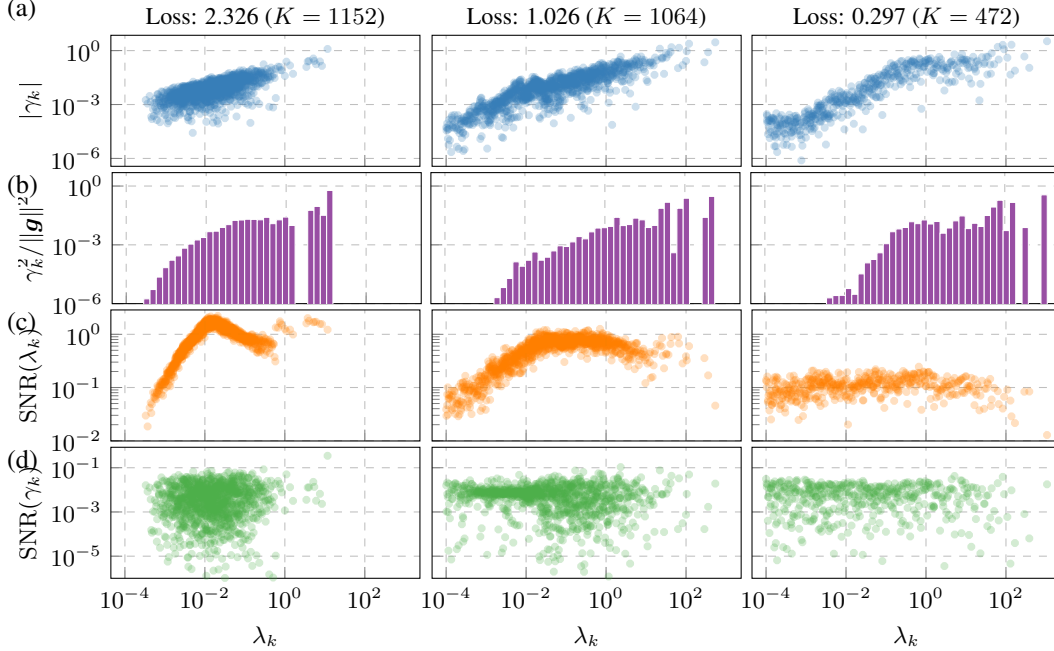


Figure S.6: **Gradient, curvature and noise during training of CIFAR-10 3C3D with ADAM.** Individual columns show the architecture’s state at initialization (*left*), an early (epoch 5, *center*), and advanced (epoch 68, *right*) stage of training. For each direction  $k$ , characterized by its curvature  $\lambda_k$ , we monitor (a) the directional gradient magnitude; (b) gradient-eigenvector alignment; (c,d) signal-to-noise ratios of curvatures and gradients.

[et al. \(2020\)](#) (see Table S.2), but turn off regularization as it would alter gradients and curvature, and their respective noise.

Metrics are computed on a single held-out mini-batch during training, using the same batch size if possible (3C3D on CIFAR-10:  $N = 128$ , 2C2D FASHION-MNIST:  $N = 128$ ), or a smaller value to fit the computation into memory (ALL-CNN-C on CIFAR-100:  $N = 64$ ). We focus on the exact GGN without further approximations and use the full mini-batch for the directional derivatives.

**Signal-to-noise ratios:** From the empirical mini-batch distributions  $\{\gamma_{nk}\}, \{\lambda_{nk}\}$  we compute

$$\text{SNR}(\chi_k) = \frac{\mathbb{E}[\chi_k]^2}{\text{Var}[\chi_k]} = \frac{\mathbb{E}[\chi_k]^2}{\mathbb{E}[\chi_k^2] - \mathbb{E}[\chi_k]^2} \quad \text{with} \quad \mathbb{E}[\chi_k] = \frac{1}{N} \sum_{n=1}^N \chi_{nk} \quad \text{for} \quad \chi \in \{\lambda, \gamma\}.$$

**Summary:** All analyzed runs exhibit similar behaviors as described in Section 4.1.

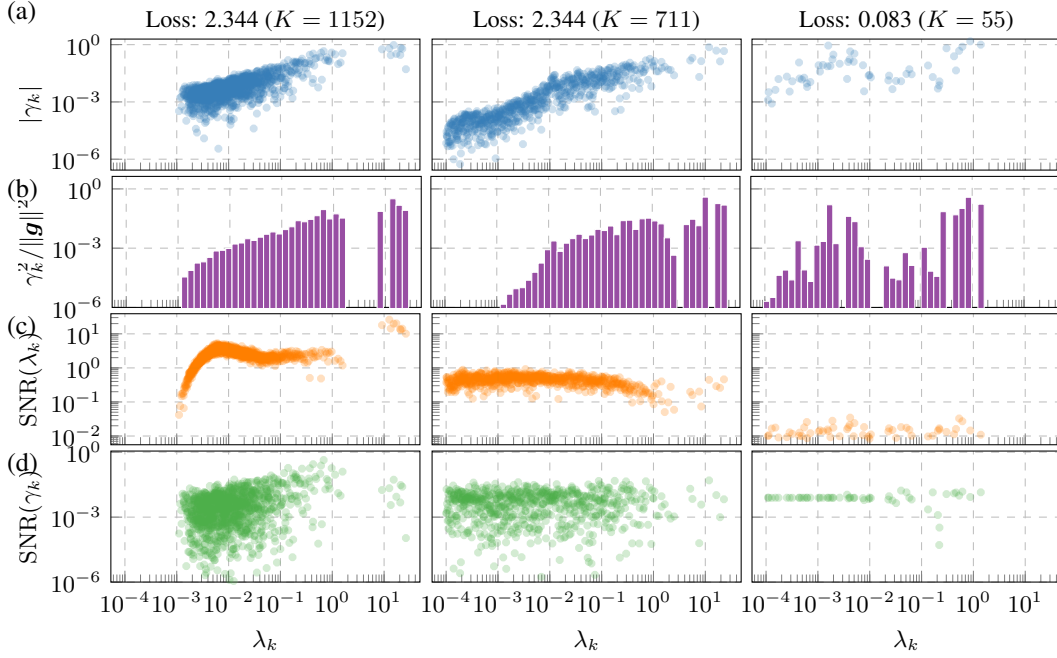


Figure S.7: **Gradient, curvature and noise during training of FASHION-MNIST 2C2D with SGD.** Individual columns show the architecture’s state at initialization (*left*), an early (after 100 steps, *center*), and advanced (epoch 57, *right*) stage of training. For each direction  $k$ , characterized by its curvature  $\lambda_k$ , we monitor (a) the directional gradient magnitude; (b) gradient-eigenvector alignment; (c,d) signal-to-noise ratios of curvatures and gradients.

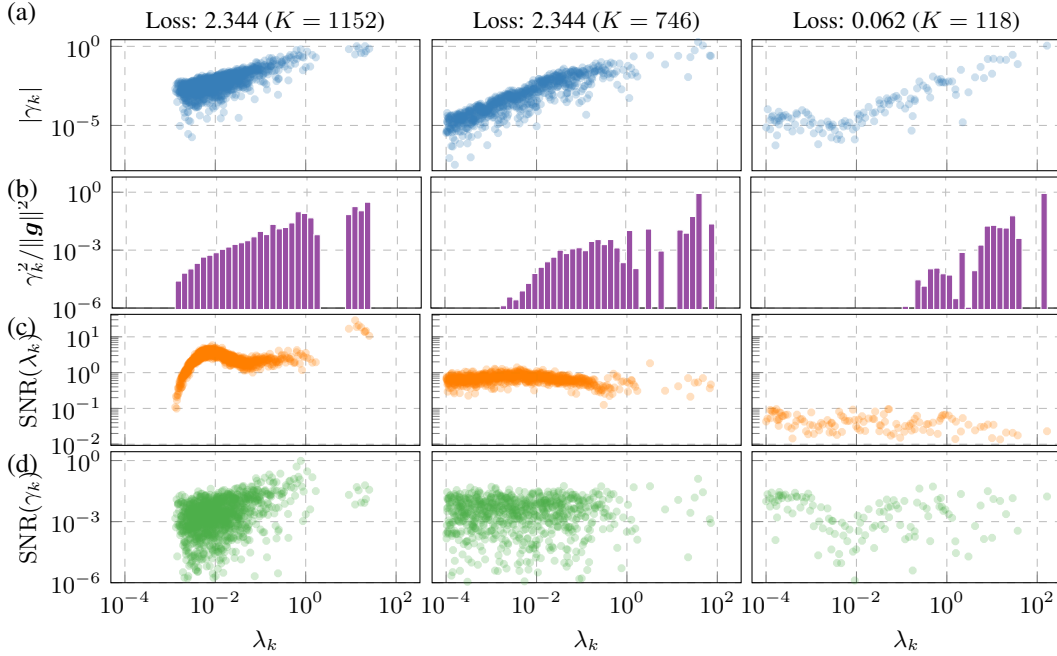


Figure S.8: **Gradient, curvature and noise during training of FASHION-MNIST 2C2D with ADAM.** Individual columns show the architecture’s state at initialization (*left*), an early (after 100 steps, *center*), and advanced (epoch 57, *right*) stage of training. For each direction  $k$ , characterized by its curvature  $\lambda_k$ , we monitor (a) the directional gradient magnitude; (b) gradient-eigenvector alignment; (c,d) signal-to-noise ratios of curvatures and gradients.

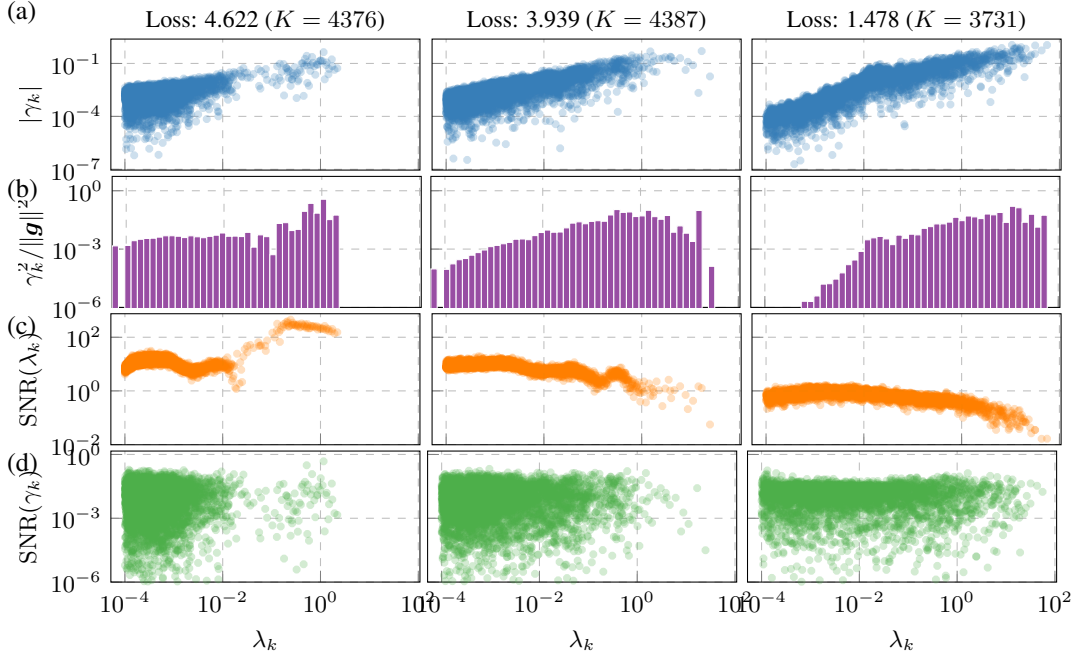


Figure S.9: **Gradient, curvature and noise during training of CIFAR-100 ALL-CNN-C with SGD.** Individual columns show the architecture’s state at initialization (*left*), an early (epoch 5, *center*), and advanced (epoch 311, *right*) stage of training. For each direction  $k$ , characterized by its curvature  $\lambda_k$ , we monitor (a) the directional gradient magnitude; (b) gradient-eigenvector alignment; (c,d) signal-to-noise ratios of curvatures and gradients.

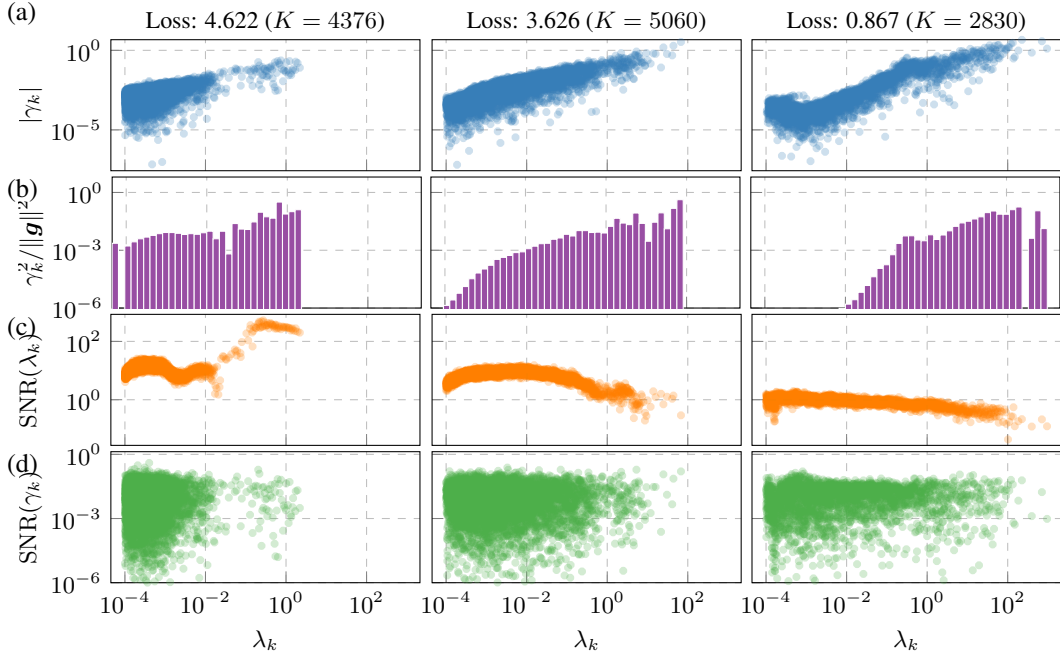


Figure S.10: **Gradient, curvature and noise during training of CIFAR-100 ALL-CNN-C with ADAM.** Individual columns show the architecture’s state at initialization (*left*), an early (epoch 5, *center*), and advanced (epoch 311, *right*) stage of training. For each direction  $k$ , characterized by its curvature  $\lambda_k$ , we monitor (a) the directional gradient magnitude; (b) gradient-eigenvector alignment; (c,d) signal-to-noise ratios of curvatures and gradients.



### B.3 BOOTSTRAP DAMPING FOR SECOND-ORDER METHODS

**Bootstrap damping:** The starting point for our directional bootstrap damping is Equation (10),

$$\mathcal{R}(\delta_k) = q(\theta) - q(\theta + s_k e_k) = -s_k \left( \frac{1}{N} \sum_{n=1}^N \gamma_{nk} \right) - \frac{1}{2} s_k^2 \left( \frac{1}{N} \sum_{n=1}^N \lambda_{nk} \right).$$

It describes the reduction of the quadratic mini-batch model  $q$  when taking a step  $s_k = -\frac{\gamma_k}{\lambda_k + \delta_k} \in \mathbb{R}$  in direction  $e_k$ . Since our ultimate goal is to minimize the training loss, we would like to choose a damping such that its corresponding update not only reduces  $q$ , but *consistently* decreases the loss over all other mini-batch models as well.

One way to assess the step  $s_k$  in this regard is the non-parametric bootstrap (Efron, 1979). The general idea is that we can simulate additional samples for the derivatives in direction  $e_k$  on other batches by resampling from the data  $\{\gamma_{nk}\}_{n=1}^N, \{\lambda_{nk}\}_{n=1}^N$ . First, we draw  $N$  indices  $i_1, \dots, i_N \in \{1, \dots, N\}$  with replacement. By taking the mean of  $\gamma_{i_1 k}, \dots, \gamma_{i_N k}$ , we can simulate the first directional derivative of an alternative, equally valid quadratic model (and similarly for the second directional derivative by taking the mean of  $\lambda_{i_1 k}, \dots, \lambda_{i_N k}$ ). Replacing the directional derivatives  $\frac{1}{N} \sum_{n=1}^N \gamma_{nk}$  and  $\frac{1}{N} \sum_{n=1}^N \lambda_{nk}$  in the equation above by these new averages  $\frac{1}{N} \sum_{j=1}^N \gamma_{i_j k}$  and  $\frac{1}{N} \sum_{j=1}^N \lambda_{i_j k}$  yields a new sample for  $\mathcal{R}_k(\delta_k)$ . This allows to create an arbitrary number of samples for any  $\delta_k$ .

For a given damping  $\delta_k$ , such samples indicate what reduction in training loss to expect with the respective update  $s_k$ . Taking the 5% percentile of the bootstrap-generated samples, we obtain a *confident* lower bound to this reduction. We then choose the  $\delta_k$  that maximizes this lower bound from candidates on a discrete grid. This approach is repeated for all non-trivial directions.

### B.4 DERIVATION OF NOISY QUADRATIC MODEL

Here, we reverse-engineer a noisy quadratic model used by the optimizer to minimize an inaccessible objective function such that we have full control over the directional noisy first- and second-order derivatives observed by the optimizer through automatic differentiation (directions themselves are not subject to noise). Not only do we formulate the optimization problem mathematically, but also derive an equivalent BACKPACK-compatible neural network training procedure, such that we can use ViViT to compute directional derivatives during a backward pass.

**One-dimensional case:** Consider a one-dimensional objective function  $f : \mathbb{R} \rightarrow \mathbb{R}, \vartheta \mapsto f(\vartheta)$  which we want to optimize,

$$\min_{\vartheta} f(\vartheta).$$

At the current iterate  $\vartheta_0 \in \mathbb{R}$ , an optimizer constructs a local model of  $f$  to update its solution.

We choose the family of local models to be convex quadratics. Let  $\varphi_{\text{true}} : \mathbb{R} \rightarrow \mathbb{R}, \vartheta \mapsto \varphi_{\text{true}}(\vartheta)$  denote one instance of a model, defined by its local curvature  $\lambda_{\text{true}}(\vartheta_0) \in \mathbb{R}_+$ , gradient  $\gamma_{\text{true}}(\vartheta_0) \in \mathbb{R}$ , and function value  $\varphi_{\text{true}}(\vartheta_0) \in \mathbb{R}$ ,

$$\varphi_{\text{true}}(\vartheta) = \frac{1}{2} \lambda_{\text{true}}(\vartheta_0) (\vartheta - \vartheta_0)^2 + \gamma_{\text{true}}(\vartheta_0) (\vartheta - \vartheta_0) + \varphi_{\text{true}}(\vartheta_0). \quad (\text{S.11})$$

By construction,  $\nabla_{\vartheta} \varphi_{\text{true}}(\vartheta_0) = \gamma_{\text{true}}(\vartheta_0)$  and  $\nabla_{\vartheta}^2 \varphi_{\text{true}}(\vartheta_0) = \lambda_{\text{true}}(\vartheta_0)$ . For a ‘good’ local description of  $f$ , those values should be representative of  $\nabla_{\vartheta} f(\vartheta_0), \nabla_{\vartheta}^2 f(\vartheta_0)$ , but may sometimes deviate to guarantee model properties such as convexity.

Both the objective  $f$  and its local model  $\varphi_{\text{true}}$ , assumed by the optimizer, are inaccessible in practice. A noisy version  $\varphi : \mathbb{R} \rightarrow \mathbb{R}, \vartheta \mapsto \varphi(\vartheta)$  of Equation (S.11) substitutes true curvature ( $\lambda_{\text{true}}(\vartheta_0) \rightarrow \lambda(\vartheta_0)$ ) and gradient ( $\gamma_{\text{true}}(\vartheta_0) \rightarrow \gamma(\vartheta_0)$ ) with random variables (the offset  $\varphi_{\text{true}}(\vartheta_0)$  is not crucial for optimization, and will thus not be perturbed with noise in this presentation),

$$\varphi(\vartheta) = \frac{1}{2} \lambda(\vartheta_0) (\vartheta - \vartheta_0)^2 + \gamma(\vartheta_0) (\vartheta - \vartheta_0) + \varphi_{\text{true}}(\vartheta_0). \quad (\text{S.12})$$

To observe  $\varphi(\vartheta)$  multiple times, Equation (S.12) is evaluated on samples  $\{(\lambda_n(\vartheta_0), \gamma_n(\vartheta_0)) \in \mathbb{R}_+ \times \mathbb{R}\}_{n=1}^N$  drawn from the joint distribution of  $(\lambda(\vartheta_0), \gamma(\vartheta_0))$ . Neglecting the offset, sample  $n$  gives rise to

$$\varphi_n(\vartheta) = \frac{1}{2} \lambda_n(\vartheta_0) (\vartheta - \vartheta_0)^2 + \gamma_n(\vartheta_0) (\vartheta - \vartheta_0). \quad (\text{S.13a})$$



It is common to batch-process multiple samples, compute the average

$$\bar{\varphi}(\vartheta) = \frac{1}{N} \sum_{n=1}^N \varphi_n(\vartheta), \quad (\text{S.13b})$$

and use automatic differentiation to compute the quantities employed by an optimizer. Recap that Equation (S.13b) produces the correct noisy first- and second-order partial derivatives, i.e.  $\nabla_{\vartheta} \varphi_n(\vartheta_0) = \gamma_n(\vartheta_0)$  and  $\nabla_{\vartheta}^2 \varphi_n(\vartheta_0) = \lambda_n(\vartheta_0)$ . Because our work not only relies on automatic differentiation, but computes directional derivatives through ViViT, we translate Equation (S.13b) into the training procedure of a sequential neural network with specifically engineered labeled data.

**BACKPACK-compatible one-dimensional formulation:** To reformulate Equation (S.13) as train loss of a sequential neural net, we complete the square in Equation (S.12),

$$\begin{aligned} \varphi(\vartheta) &= \frac{1}{2} \lambda(\vartheta_0) \left( \vartheta - \vartheta_0 + \frac{\gamma(\vartheta_0)}{\lambda(\vartheta_0)} \right)^2 - \frac{1}{2} \frac{\gamma(\vartheta_0)^2}{\lambda(\vartheta_0)} + \varphi_{\text{true}}(\vartheta_0) \\ &= \left( \sqrt{\frac{\lambda(\vartheta_0)}{2}} (\vartheta - \vartheta_0) + \frac{\gamma(\vartheta_0)}{\sqrt{2\lambda(\vartheta_0)}} \right)^2 - \frac{1}{2} \frac{\gamma(\vartheta_0)^2}{\lambda(\vartheta_0)} + \varphi_{\text{true}}(\vartheta_0) \\ &= \left( \sqrt{\frac{\lambda(\vartheta_0)}{2}} \vartheta - \left( \sqrt{\frac{\lambda(\vartheta_0)}{2}} \vartheta_0 - \frac{\gamma(\vartheta_0)}{\sqrt{2\lambda(\vartheta_0)}} \right) \right)^2 - \frac{1}{2} \frac{\gamma(\vartheta_0)^2}{\lambda(\vartheta_0)} + \varphi_{\text{true}}(\vartheta_0) \\ &= (x\vartheta - y)^2 + \text{const.} \end{aligned} \quad (\text{S.14a})$$

with

$$x = \sqrt{\frac{\lambda(\vartheta_0)}{2}} \in \mathbb{R}, \quad (\text{S.14b})$$

$$y = \sqrt{\frac{\lambda(\vartheta_0)}{2}} \vartheta_0 - \frac{\gamma(\vartheta_0)}{\sqrt{2\lambda(\vartheta_0)}} \in \mathbb{R}, \quad (\text{S.14c})$$

$$\text{const.} = -\frac{1}{2} \frac{\gamma(\vartheta_0)^2}{\lambda(\vartheta_0)} + \varphi_{\text{true}}(\vartheta_0) \in \mathbb{R}. \quad (\text{S.14d})$$

Note that the data, in which the noisy observations are embedded, must depend on the current location  $\vartheta_0$ . Such a dependence is not common in practical tasks. But it is a consequence of the design, because we want to achieve full control of gradient and curvature noise that the optimizer is exposed to at any time during optimization.

With that, we conclude that minimizing Equation (S.12) through noisy observations of the form Equation (S.13) is equivalent to the following neural network training: Assume the optimizer's current iterate to be  $\vartheta_0$ . Then, the following steps define a training iteration:

1. Generate data

- (a) Draw curvature and gradient samples  $\{(\lambda_n(\vartheta_0), \gamma_n(\vartheta_0))\}_{n=1}^N$
- (b) Compute inputs  $\{x_n = \sqrt{\lambda_n(\vartheta_0)/2}\}_{n=1}^N$  and labels  $\{y_n = \sqrt{\lambda_n(\vartheta_0)/2} \vartheta_0 - \gamma_n(\vartheta_0)/\sqrt{2\lambda_n(\vartheta_0)}\}_{n=1}^N\}$

2. Forward pass

- (a) Feed  $\{x_n\}_{n=1}^N$  through a linear layer ( $\mathbb{R} \rightarrow \mathbb{R}$ , no bias) with trainable weight  $\vartheta \in \mathbb{R}$ , set to  $\vartheta_0$
- (b) Feed the output  $\{x_n \vartheta\}_{n=1}^N$  through the mean-squared error (MSE) with labels  $\{y_n\}_{n=1}^N$ ,

$$\text{MSE}(\{x_n \vartheta\}, \{y_n\}) = \frac{1}{N} \sum_{n=1}^N (x_n \vartheta - y_n)^2 = \bar{\varphi}(\vartheta)$$

- 3. Backward pass: Compute the first- and second-order directional derivatives  $\{\gamma_n(\vartheta_0)\}_{n=1}^N$  and  $\{\lambda_n(\vartheta_0)\}_{n=1}^N$  with ViViT during a backward pass. By construction, this reproduces the sampled gradients and curvatures from Item 1a
- 4. Optimizer step: Update the value of  $\vartheta$  using  $\{\gamma_n(\vartheta_0)\}_{n=1}^N$  and  $\{\lambda_n(\vartheta_0)\}_{n=1}^N$ , set  $\vartheta_0$  to  $\vartheta$ 's new value

**Multi-dimensional case:** Next, we extend the one-dimensional case to multiple dimensions. Consider a multi-dimensional objective function  $F : \mathbb{R}^D \rightarrow \mathbb{R}, \theta \mapsto F(\theta)$  which we want to optimize,

$$\min_{\theta} F(\theta).$$

The local model of  $F$  at  $\theta_0 \in \mathbb{R}^D$ , employed by the optimizer, is a convex  $D$ -dimensional quadratic. Let  $\phi_{\text{true}} : \mathbb{R}^D \rightarrow \mathbb{R}, \theta \mapsto \phi_{\text{true}}(\theta)$  denote one instance of a model, defined by its local Hessian spectrum  $\{(\lambda_{d,\text{true}}(\theta_0), e_d(\theta_0))\}_{d=1}^D$  with local curvatures  $\{\lambda_{d,\text{true}}(\theta_0) \in \mathbb{R}_+\}_{d=1}^D$  such that the Hessian is  $\nabla_{\theta}^2 \phi(\theta_0) = \sum_{d=1}^D \lambda_{d,\text{true}}(\theta_0) e_d(\theta_0) e_d(\theta_0)^\top$ . It also requires local gradients  $\{\gamma_{d,\text{true}}(\theta_0) \in \mathbb{R}\}_{d=1}^D$  along the directions  $\{e_d(\theta_0) \in \mathbb{R}^D\}_{d=1}^D$ , such that  $\nabla_{\theta} \phi(\theta_0) = \sum_{d=1}^D \gamma_{d,\text{true}}(\theta_0) e_d(\theta_0)$ , and the local function value at  $\theta_0 \in \mathbb{R}^D$ ,

$$\begin{aligned} \phi_{\text{true}}(\theta) &= \frac{1}{2}(\theta - \theta_0)^\top \nabla_{\theta}^2 \phi(\theta_0)(\theta - \theta_0) + (\theta - \theta_0)^\top \nabla_{\theta} \phi(\theta_0) + \phi_{\text{true}}(\theta_0) \\ &= \left[ \sum_{d=1}^D \frac{1}{2} \lambda_{d,\text{true}}(\theta_0) (e_d(\theta_0)^\top (\theta - \theta_0))^2 + \gamma_{d,\text{true}}(\theta_0) e_d(\theta_0)^\top (\theta - \theta_0) \right] + \phi_{\text{true}}(\theta_0). \end{aligned} \quad (\text{S.15})$$

Note that expanding  $\theta = \sum_{d=1}^D \vartheta_d e_d(\theta_0)$  and  $\theta_0 = \sum_{d=1}^D \vartheta_{d,0} e_d(\theta_0)$  in terms of scalar coordinates  $\vartheta_d = e_d(\theta_0)^\top \theta$  and  $\vartheta_{d,0} = e_d(\theta_0)^\top \theta_0$  decouples Equation (S.15) into one-dimensional quadratics,

$$\begin{aligned} \phi_{\text{true}}(\theta) &= \left[ \sum_{d=1}^D \frac{1}{2} \lambda_{d,\text{true}}(\theta_0) (\vartheta_d - \vartheta_{d,0})^2 + \gamma_{d,\text{true}}(\theta_0) (\vartheta_d - \vartheta_{d,0}) \right] + \phi_{\text{true}}(\theta_0) \\ &= \sum_{d=1}^D \varphi_{d,\text{true}}(\vartheta_d) + \phi_{\text{true}}(\theta_0), \end{aligned} \quad (\text{S.16})$$

which is (up to constant offsets that are negligible for optimization) the sum of  $D$  one-dimensional quadratic functions of the form Equation (S.11), with  $\varphi_{d,\text{true}}$  defined by  $(\lambda_{d,\text{true}}(\theta_0), \gamma_{d,\text{true}}(\theta_0), \vartheta_{d,0})$  along each direction.

A noisy version  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}, \theta \mapsto \phi(\theta)$  of Equation (S.15) replaces the true curvatures  $(\{\lambda_{d,\text{true}}(\theta_0)\}_{d=1}^D \rightarrow \{\lambda_d(\theta_0)\}_{d=1}^D)$  and local gradients  $(\{\gamma_{d,\text{true}}(\theta_0)\}_{d=1}^D \rightarrow \{\gamma_d(\theta_0)\}_{d=1}^D)$  at  $\theta_0$  with random variables (the offset  $\phi_{\text{true}}(\theta_0)$  is not crucial for optimization, and will thus not be perturbed),

$$\phi(\theta) = \left[ \sum_{d=1}^D \frac{1}{2} \lambda_d(\theta_0) (e_d(\theta_0)^\top (\theta - \theta_0))^2 + \gamma_d(\theta_0) e_d(\theta_0)^\top (\theta - \theta_0) \right] + \phi_{\text{true}}(\theta_0). \quad (\text{S.17})$$

The optimizer observes  $\phi(\theta)$  through  $N$  samples  $\{(\lambda_{nd}(\theta_0), \gamma_{nd}(\theta_0)) \in \mathbb{R}_+ \times \mathbb{R}\}_{n=1, d=1}^{N,D}$  drawn from the joint distribution of  $(\lambda_d(\theta_0), \gamma_d(\theta_0))$ . Neglecting the offset, sample  $\{(\lambda_{nd}(\theta_0), \gamma_{nd}(\theta_0))\}_{d=1}^D$  gives rise to

$$\phi_n(\theta) = \sum_{d=1}^D \frac{1}{2} \lambda_{nd}(\theta_0) (e_d(\theta_0)^\top (\theta - \theta_0))^2 + \gamma_{nd}(\theta_0) e_d(\theta_0)^\top (\theta - \theta_0). \quad (\text{S.18a})$$

It is common to batch-process multiple samples and compute the average

$$\bar{\phi}(\theta) = \frac{1}{N} \sum_{n=1}^N \phi_n(\theta), \quad (\text{S.18b})$$

which will give the correct first- and second-order directional derivatives in an automatic differentiation engine, i.e.  $e_d(\theta_0)^\top \nabla_{\theta} \phi_n(\theta_0) = \gamma_{nd}(\theta_0)$  and  $e_d(\theta_0)^\top \nabla_{\theta}^2 \phi_n(\theta_0) e_d(\theta_0) = \lambda_{nd}(\theta_0)$ . We now phrase minimizing  $\phi(\theta)$  by observing  $\bar{\phi}(\theta)$  as neural network training.

**BACKPACK-compatible multi-dimensional formulation:** Completing the square in Equation (S.15), and in analogy to Equation (S.14), gives

$$\begin{aligned}
\phi(\boldsymbol{\theta}) &= \left[ \sum_{d=1}^D \frac{1}{2} \lambda_d(\boldsymbol{\theta}_0) \left( \mathbf{e}_d(\boldsymbol{\theta}_0)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \frac{\gamma_d(\boldsymbol{\theta}_0)}{\lambda_d(\boldsymbol{\theta}_0)} \right)^2 - \frac{1}{2} \frac{\gamma_d(\boldsymbol{\theta}_0)^2}{\lambda_d(\boldsymbol{\theta}_0)} \right] + \phi_{\text{true}}(\boldsymbol{\theta}_0) \\
&= \left[ \sum_{d=1}^D \left( \sqrt{\frac{\lambda_d(\boldsymbol{\theta}_0)}{2}} \mathbf{e}_d(\boldsymbol{\theta}_0)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \frac{\gamma_d(\boldsymbol{\theta}_0)}{\sqrt{2\lambda_d(\boldsymbol{\theta}_0)}} \right)^2 - \frac{1}{2} \frac{\gamma_d(\boldsymbol{\theta}_0)^2}{\lambda_d(\boldsymbol{\theta}_0)} \right] + \phi_{\text{true}}(\boldsymbol{\theta}_0) \\
&= \left[ \sum_{d=1}^D \left( \sqrt{\frac{\lambda_d(\boldsymbol{\theta}_0)}{2}} \mathbf{e}_d(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta} - \left( \sqrt{\frac{\lambda_d(\boldsymbol{\theta}_0)}{2}} \mathbf{e}_d(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta}_0 - \frac{\gamma_d(\boldsymbol{\theta}_0)}{\sqrt{2\lambda_d(\boldsymbol{\theta}_0)}} \right) \right)^2 \right. \\
&\quad \left. - \frac{1}{2} \frac{\gamma_d(\boldsymbol{\theta}_0)^2}{\lambda_d(\boldsymbol{\theta}_0)} \right] + \phi_{\text{true}}(\boldsymbol{\theta}_0) \\
&= \frac{1}{D} \left[ \sum_{d=1}^D \left( \sqrt{\frac{\lambda_d(\boldsymbol{\theta}_0)D}{2}} \mathbf{e}_d(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta} - \left( \sqrt{\frac{\lambda_d(\boldsymbol{\theta}_0)D}{2}} \mathbf{e}_d(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta}_0 - \sqrt{\frac{D}{2\lambda_d(\boldsymbol{\theta}_0)}} \gamma_d(\boldsymbol{\theta}_0) \right) \right)^2 \right. \\
&\quad \left. - \frac{1}{2} \frac{D\gamma_d(\boldsymbol{\theta}_0)^2}{\lambda_d(\boldsymbol{\theta}_0)} \right] + \phi_{\text{true}}(\boldsymbol{\theta}_0) \\
&= \frac{1}{D} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2 + \text{const.}
\end{aligned} \tag{S.19a}$$

with

$$\mathbf{X} = \sqrt{\frac{D}{2}} \begin{pmatrix} \sqrt{\lambda_1} \mathbf{e}_1(\boldsymbol{\theta}_0)^\top \\ \vdots \\ \sqrt{\lambda_D} \mathbf{e}_D(\boldsymbol{\theta}_0)^\top \end{pmatrix} \in \mathbb{R}^{D \times D}, \tag{S.19b}$$

$$\mathbf{y} = \sqrt{\frac{D}{2}} \begin{pmatrix} \sqrt{\lambda_1(\boldsymbol{\theta}_0)} \mathbf{e}_1(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta}_0 - \frac{\gamma_1(\boldsymbol{\theta}_0)}{\sqrt{\lambda_1(\boldsymbol{\theta}_0)}} \\ \vdots \\ \sqrt{\lambda_D(\boldsymbol{\theta}_0)} \mathbf{e}_D(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta}_0 - \frac{\gamma_D(\boldsymbol{\theta}_0)}{\sqrt{\lambda_D(\boldsymbol{\theta}_0)}} \end{pmatrix} \in \mathbb{R}^D, \tag{S.19c}$$

$$\text{const.} = \left( \sum_{d=1}^D -\frac{1}{2} \frac{\gamma_d(\boldsymbol{\theta}_0)^2}{\lambda_d(\boldsymbol{\theta}_0)} \right) + \phi_{\text{true}}(\boldsymbol{\theta}_0) \in \mathbb{R}. \tag{S.19d}$$

Note that we extracted a fraction  $1/D$ . This is due to our rephrasing goal as neural network training.

With that, we conclude that minimizing Equation (S.17) through noisy observations of the form Equation (S.18) is equivalent to the following neural network training: Assume the optimizer's current iterate to be  $\boldsymbol{\theta}_0$ . Then, the following sequence defines a training iteration:

1. Generate data

- (a) Generate orthonormal eigenvectors  $\{\mathbf{e}_d(\boldsymbol{\theta}_0)\}_{d=1}^D$  with  $\mathbf{e}_d(\boldsymbol{\theta}_0)^\top \mathbf{e}_{d'}(\boldsymbol{\theta}_0) = \delta_{dd'}$
- (b) Draw curvature and gradient samples  $\{(\lambda_{nd}(\boldsymbol{\theta}_0), \gamma_{nd}(\boldsymbol{\theta}_0))\}_{n=1, d=1}^{N, D}$
- (c) Compute inputs and labels

$$\begin{aligned}
\left\{ \mathbf{X}_n = \sqrt{\frac{D}{2}} \begin{pmatrix} \sqrt{\lambda_{n1}(\boldsymbol{\theta}_0)} \mathbf{e}_1(\boldsymbol{\theta}_0)^\top \\ \vdots \\ \sqrt{\lambda_{nD}(\boldsymbol{\theta}_0)} \mathbf{e}_D(\boldsymbol{\theta}_0)^\top \end{pmatrix} \right\}_{n=1}^N, \\
\left\{ \mathbf{y}_n = \frac{D}{2} \begin{pmatrix} \sqrt{\lambda_{n1}(\boldsymbol{\theta}_0)} \mathbf{e}_1(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta}_0 - \frac{\gamma_{n1}(\boldsymbol{\theta}_0)}{\sqrt{\lambda_{n1}(\boldsymbol{\theta}_0)}} \\ \vdots \\ \sqrt{\lambda_{nD}(\boldsymbol{\theta}_0)} \mathbf{e}_D(\boldsymbol{\theta}_0)^\top \boldsymbol{\theta}_0 - \frac{\gamma_{nD}(\boldsymbol{\theta}_0)}{\sqrt{\lambda_{nD}(\boldsymbol{\theta}_0)}} \end{pmatrix} \right\}_{n=1}^N
\end{aligned}$$

## 2. Forward pass

- (a) Feed  $\{\mathbf{X}_n\}$  through a linear layer ( $\mathbb{R}^{\dots \times D} \rightarrow \mathbb{R}^{\dots \times 1}$  where  $\dots$  denotes free axes preserved by the affine transformation<sup>6</sup>) with trainable weight  $\boldsymbol{\theta}^\top$  and no bias
- (b) Feed the output  $\{\mathbf{X}_n \boldsymbol{\theta}\}_{n=1}^N$  through the mean-squared error (MSE) with labels  $\{\mathbf{y}_n\}_{n=1}^N$ ,

$$\text{MSE}(\{\mathbf{X}_n \boldsymbol{\theta}\}, \{\mathbf{y}_n\}) = \frac{1}{N} \frac{1}{D} \sum_{n=1}^N \|\mathbf{X}_n \boldsymbol{\theta} - \mathbf{y}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \phi_n(\boldsymbol{\theta}) = \bar{\phi}(\boldsymbol{\theta}). \quad (\text{S.20})$$

Note that the factor  $1/D$  is required as the MSE implementation in common machine learning libraries averages the squared residuals over all components.<sup>7</sup>

3. Backward pass: Compute first- and second-order directional derivatives  $\{\gamma_{nd}(\boldsymbol{\theta}_0)\}_{n=1, d=1}^{N, D}$  and  $\{\lambda_{nd}(\boldsymbol{\theta}_0)\}_{n=1, d=1}^{N, D}$  with ViViT during a backward pass. They are of same value as the curvature and gradient samples defined in Item 1b
4. Optimizer step: Update the value of  $\boldsymbol{\theta}$  using  $\{\gamma_{nd}(\boldsymbol{\theta}_0)\}_{n=1, d=1}^{N, D}$  and  $\{\lambda_{nd}(\boldsymbol{\theta}_0)\}_{n=1, d=1}^{N, D}$ , set  $\boldsymbol{\theta}_0$  to  $\boldsymbol{\theta}$ 's new value

## B.5 EXPERIMENT ON NOISY QUADRATIC

We consider a quadratic loss function  $\mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{G} \boldsymbol{\theta}$  with  $G_{ii} = i^2$  for  $i \in \{1, \dots, D = 20\}$  and use the noisy quadratic network from Appendix B.4 to gain full control over the directions  $\mathbf{e}_k$  as well as directional derivative samples  $\gamma_{nk}, \lambda_{nk}$  observed by the optimizers at each step.

We set  $\mathbf{e}_k$  to the  $k$ -th unit vector for  $k \in \{1, \dots, D\}$ ,  $N = 128$  and sample  $\gamma_{nk}$  from a Normal distribution with mean  $\mathbf{e}_k^\top \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$  (the actual directional first derivative of  $\mathcal{L}$ ) and constant variance 5,000. Note that this variance implies a gradient SNR of at least 2 over all directions at  $\boldsymbol{\theta}_{\text{init}}$ . Since the gradient vanishes when moving towards the minimum, this SNR becomes arbitrarily low. When sampling  $\lambda_{nk}$ , we have to make sure that these samples are non-negative, since they correspond to projections of positive semi-definite matrices  $\mathbf{G}_n$  (compare Equation (8b)). This constraint can be incorporated by sampling  $\lambda_{nk}$  from a Gamma distribution with mean  $\mathbf{e}_k^\top \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{e}_k$  and constant variance 50. This globally limits the curvature SNR to values above 0.02. We compare SGD (learning rate  $10^{-3}$ ) to the second-order optimizers with global damping  $\delta \in \{10^{-4}, 10^{-3}, \dots, 10^2\}$  and directional damping  $\delta_k$  using a log-equidistant grid from  $10^{-4}$  to  $10^2$  with 200 grid points and 100 bootstrap samples. We run all optimizers for 20 steps. To obtain a reliable estimation of optimizer stability, this procedure is repeated 100 times. The results are shown in Figure 3.

## C IMPLEMENTATION DETAILS

**Layer view of backpropagation:** Consider a single layer  $T_{\boldsymbol{\theta}^{(i)}}^{(i)}$  that transforms inputs  $\mathbf{z}_n^{(i-1)} \in \mathbb{R}^{h^{(i-1)}}$  into outputs  $\mathbf{z}_n^{(i)} \in \mathbb{R}^{h^{(i)}}$  by means of a parameter  $\boldsymbol{\theta}^{(i)} \in \mathbb{R}^{d^{(i)}}$ . During backpropagation for  $\mathbf{V}$ , the layer receives vectors  $\mathbf{s}_{nc}^{(i)} = (\mathbf{J}_{\mathbf{z}_n^{(i)}} f_n)^\top \mathbf{s}_{nc}$  from the previous stage (recall  $\nabla_f^2 \ell_n = \sum_{c=1}^C \mathbf{s}_{nc} \mathbf{s}_{nc}^\top$ ). Parameter contributions  $\mathbf{v}_{nc}^{(i)}$  to  $\mathbf{V}$  are obtained by application of its Jacobian,

$$\begin{aligned} \mathbf{v}_{nc}^{(i)} &= (\mathbf{J}_{\boldsymbol{\theta}^{(i)}} f_n)^\top \mathbf{s}_{nc} \\ &= \left( \mathbf{J}_{\boldsymbol{\theta}^{(i)}} \mathbf{z}_n^{(i)} \right)^\top \left( \mathbf{J}_{\mathbf{z}_n^{(i)}} f_n \right)^\top \mathbf{s}_{nc} && \text{(Chain rule)} \\ &= \left( \mathbf{J}_{\boldsymbol{\theta}^{(i)}} \mathbf{z}_n^{(i)} \right)^\top \mathbf{s}_{nc}^{(i)}. && \text{(Definition of } \mathbf{s}_{nc}^{(i)} \text{)} \end{aligned} \quad (\text{S.21})$$

Consequently, the contribution of  $\boldsymbol{\theta}^{(i)}$  to  $\mathbf{V}$ , denoted by  $\mathbf{V}^{(i)} \in \mathbb{R}^{d^{(i)} \times NC}$ , is

$$\mathbf{V}^{(i)} = \frac{1}{\sqrt{N}} \begin{pmatrix} \mathbf{v}_{11}^{(i)} & \mathbf{v}_{12}^{(i)} & \dots & \mathbf{v}_{NC}^{(i)} \end{pmatrix} \quad \text{with} \quad \mathbf{v}_{nc}^{(i)} = (\mathbf{J}_{\boldsymbol{\theta}^{(i)}} f_n)^\top \mathbf{s}_{nc}. \quad (\text{S.22})$$

<sup>6</sup>This complies to the implementation in PyTorch, see `torch.nn.Linear`.

<sup>7</sup>See e.g. the PyTorch implementation of `torch.nn.MSELoss`.

### C.1 OPTIMIZED GRAM MATRIX COMPUTATION FOR LINEAR LAYERS

Our goal is to efficiently extract  $\theta^{(i)}$ 's contribution to the Gram matrix  $\tilde{\mathbf{G}}$ , given by

$$\tilde{\mathbf{G}}^{(i)} = \mathbf{V}^{(i)\top} \mathbf{V}^{(i)} \in \mathbb{R}^{NC \times NC}. \quad (\text{S.23})$$

**Gram matrix via expanding  $\mathbf{V}^{(i)}$ :** One way to construct  $\mathbf{G}^{(i)}$  is to first expand  $\mathbf{V}^{(i)}$  (Equation (S.22)) via the Jacobian  $\mathbf{J}_{\theta^{(i)}} \mathbf{z}_n^{(i)}$ , then contract it (Equation (S.23)). This can be a memory bottleneck for large linear layers which are common in many architectures close to the network output. However if only the Gram matrix rather than  $\mathbf{V}$  is required, structure in the Jacobian can be used to construct  $\tilde{\mathbf{G}}^{(i)}$  without expanding  $\mathbf{V}^{(i)}$  and thus reduce this overhead.

**Optimization for linear layers:** Now, let  $T_{\theta^{(i)}}^{(i)}$  be a linear layer with weights  $\mathbf{W}^{(i)} \in \mathbb{R}^{h^{(i)} \times h^{(i-1)}}$ , i.e.  $\theta^{(i)} = \text{vec}(\mathbf{W}^{(i)}) \in \mathbb{R}^{d^{(i)}=h^{(i)}h^{(i-1)}}$  with column stacking convention for vectorization,

$$T_{\theta^{(i)}}^{(i)} : \mathbf{z}_n^{(i)} = \mathbf{W}^{(i)} \mathbf{z}_n^{(i-1)}.$$

The Jacobian is

$$\mathbf{J}_{\theta^{(i)}} \mathbf{z}_n^{(i)} = \mathbf{z}_n^{(i-1)\top} \otimes \mathbf{I}_{h^{(i)}}. \quad (\text{S.24})$$

Its structure can be used to directly compute entries of the Gram matrix without expanding  $\mathbf{V}^{(i)}$ ,

$$\begin{aligned} \left[ \tilde{\mathbf{G}}^{(i)} \right]_{(nc)(n'c')} &= \mathbf{v}_{nc}^{(i)\top} \mathbf{v}_{n'c'}^{(i)} && (\text{Equation (S.23)}) \\ &= \mathbf{s}_{nc}^{(i)\top} \left( \mathbf{J}_{\theta^{(i)}} \mathbf{z}_n^{(i)} \right) \left( \mathbf{J}_{\theta^{(i)}} \mathbf{z}_{n'}^{(i)} \right)^\top \mathbf{s}_{n'c'}^{(i)} \\ &= \mathbf{s}_{nc}^{(i)\top} \left( \mathbf{z}_n^{(i-1)\top} \otimes \mathbf{I}_{h^{(i)}} \right) \left( \mathbf{z}_{n'}^{(i-1)\top} \otimes \mathbf{I}_{h^{(i)}} \right)^\top \mathbf{s}_{n'c'}^{(i)} && (\text{Equation (S.24)}) \\ &= \mathbf{s}_{nc}^{(i)\top} \left( \mathbf{z}_n^{(i-1)\top} \mathbf{z}_{n'}^{(i-1)} \otimes \mathbf{I}_{h^{(i)}} \right) \mathbf{s}_{n'c'}^{(i)} && (\text{Equation (S.21)}) \\ &= \mathbf{z}_n^{(i-1)\top} \mathbf{z}_{n'}^{(i-1)} \mathbf{s}_{nc}^{(i)\top} \mathbf{I}_{h^{(i)}} \mathbf{s}_{n'c'}^{(i)} && (\mathbf{z}_n^{(i-1)\top} \mathbf{z}_{n'}^{(i-1)} \in \mathbb{R}) \\ &= \left( \mathbf{z}_n^{(i-1)\top} \mathbf{z}_{n'}^{(i-1)} \right) \left( \mathbf{s}_{nc}^{(i)\top} \mathbf{s}_{n'c'}^{(i)} \right). \end{aligned}$$

We see that the Gram matrix is built from two Gram matrices based on  $\{\mathbf{z}_n^{(i-1)}\}_{n=1}^N$  and  $\{\mathbf{s}_{nc}^{(i)}\}_{n=1, c=1}^{N, C}$ , that require  $\mathcal{O}(N^2)$  and  $\mathcal{O}((NC)^2)$  memory, respectively. In comparison, the naïve approach via  $\mathbf{V}^{(i)} \in \mathbb{R}^{d^{(i)} \times NC}$  scales with the number of weights, which is often comparable to  $D$ . For instance, the 3C3D architecture on CIFAR-10 has  $D = 895,210$  and the largest weight matrix has  $d^{(i)} = 589,824$ , whereas  $NC = 1,280$  during training (Schneider et al., 2019).

### C.2 IMPLICIT MULTIPLICATION WITH THE INVERSE (BLOCK-DIAGONAL) GGN

**Inverse GGN-vector products:** Consider the damped Newton step of Equation (9) that requires multiplication by  $(\mathbf{G} + \delta \mathbf{I}_D)^{-1}$ .<sup>8</sup> By means of Equation (3) and the matrix inversion lemma,

$$\begin{aligned} (\delta \mathbf{I}_D + \mathbf{G})^{-1} &= (\delta \mathbf{I}_D + \mathbf{V} \mathbf{V}^\top)^{-1} && (\text{Equation (3)}) \\ &= \frac{1}{\delta} \left( \mathbf{I}_D + \frac{1}{\delta} \mathbf{V} \mathbf{V}^\top \right)^{-1} \\ &= \frac{1}{\delta} \left[ \mathbf{I}_D - \frac{1}{\delta} \mathbf{V} \left( \mathbf{I}_{NC} + \mathbf{V}^\top \frac{1}{\delta} \mathbf{V} \right)^{-1} \mathbf{V}^\top \right] && (\text{Matrix inversion lemma}) \\ &= \frac{1}{\delta} \left[ \mathbf{I}_D - \mathbf{V} (\delta \mathbf{I}_{NC} + \mathbf{V}^\top \mathbf{V})^{-1} \mathbf{V}^\top \right] && (\text{Gram matrix}) \\ &= \frac{1}{\delta} \left[ \mathbf{I}_D - \mathbf{V} (\delta \mathbf{I}_{NC} + \tilde{\mathbf{G}})^{-1} \mathbf{V}^\top \right]. && (\text{S.25}) \end{aligned}$$

<sup>8</sup> $\delta \mathbf{I}_D$  can be replaced by other easy-to-invert matrices.

Inverse GGN-vector products require inversion of the damped Gram matrix as well as applications of  $\mathbf{V}, \mathbf{V}^\top$  for the transformations between Gram and parameter space.

**Inverse block-diagonal GGN-vector products:** Next, we replace the full GGN by its block diagonal approximation  $\mathbf{G} \approx \mathbf{G}_{\text{BDA}} = \text{diag}(\mathbf{G}^{(1)}, \mathbf{G}^{(2)}, \dots)$  with

$$\mathbf{G}^{(i)} = \mathbf{V}^{(i)} \mathbf{V}^{(i)\top} \in \mathbb{R}^{d^{(i)} \times d^{(i)}}$$

and  $\mathbf{V}^{(i)}$  as in Equation (S.22). Then, inverse multiplication reduces to each block,

$$\mathbf{G}_{\text{BDA}}^{-1} = \text{diag}(\mathbf{G}^{(1)-1}, \mathbf{G}^{(2)-1}, \dots).$$

If again a damped Newton step is considered, we can reuse Equation (S.25) with the substitutions

$$(\mathbf{G}, D, \mathbf{V}, \mathbf{V}^\top, \tilde{\mathbf{G}}) \leftrightarrow (\mathbf{G}^{(i)}, d^{(i)}, \mathbf{V}^{(i)}, \mathbf{V}^{(i)\top}, \tilde{\mathbf{G}}^{(i)})$$

to apply the inverse and immediately discard the ViViT factors: At backpropagation of layer  $T_{\theta^{(i)}}^{(i)}$

1. Compute  $\mathbf{V}^{(i)}$  using Equation (S.22).
2. Compute  $\tilde{\mathbf{G}}^{(i)}$  using Equation (S.23).
3. Compute  $(\delta \mathbf{I}_{NC} + \tilde{\mathbf{G}}^{(i)})^{-1}$ .
4. Apply the inverse in Equation (S.25) with the above substitutions to the target vector.
5. Discard  $\mathbf{V}^{(i)}, \mathbf{V}^{(i)\top}, \tilde{\mathbf{G}}^{(i)}$ , and  $(\delta \mathbf{I}_{NC} + \tilde{\mathbf{G}}^{(i)})^{-1}$ . Proceed to layer  $i - 1$ .

Note that the above scheme should only be used for parameters that satisfy  $d^{(i)} > NC$ , i.e.  $\dim(\mathbf{G}^{(i)}) > \dim(\tilde{\mathbf{G}}^{(i)})$ . Low-dimensional parameters can be grouped with others to increase their joint dimension, and to control the block structure of  $\mathbf{G}_{\text{BDA}}$ .