

SVRG AND BEYOND VIA POSTERIOR CORRECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Stochastic Variance Reduced Gradient (SVRG) and its variants aim to speed-up training by using gradient corrections, but have seen limited success in deep learning. Here, we show surprising new foundational connections of SVRG to a recently proposed Bayesian method called posterior correction. Specifically, we show that SVRG is recovered as a special case of posterior correction over the isotropic-Gaussian family, while novel extensions are automatically obtained by using more flexible exponential families. We derive two new SVRG variants by using Gaussian families: First, a Newton-like variant that employs novel Hessian corrections, and second, an Adam-like extension that improves (continual) pre-training and finetuning of Transformer language models. This is the first work to connect SVRG to Bayes and use it to boost variational training for deep networks.

1 INTRODUCTION

Variance Reduction is a powerful technique to speed-up stochastic optimization. For example, stochastic variance reduced gradient (SVRG) uses full-batch gradients to stabilize future mini-batch updates (Johnson & Zhang, 2013). The method originates in the works of Roux et al. (2012); Schmidt et al. (2017); Shalev-Shwartz & Zhang (2013) which require storing individual gradients over the full dataset. Since then, a large number of variants have been proposed exploring various aspects of this method (Nguyen et al., 2017; Fang et al., 2018; Cutkosky & Orabona, 2019). Variance reduction has become a useful tool to accelerate both convex and nonconvex optimization.

Our goal here is to explore new connections of SVRG to Bayes. Currently, no work exists in this space and no foundational connections are known. We are particularly interested in investigating whether variance reduction can be effective in non-traditional settings, for example, in variational continual pre-training and fine-tuning of language models (Shen et al., 2024). Variance reduction has not yet seen a lot of success in deep learning (Defazio & Bottou, 2019). Our exploration here is meant to assess its potential in non-traditional variational training.

In this paper, we present surprising new foundational connections of SVRG to a recently proposed Bayesian method called posterior correction (Khan, 2025). Specifically, we show that SVRG is recovered as a special case of posterior correction over the isotropic-Gaussian family, while novel extensions are automatically obtained by using more flexible exponential families (Fig. 1(left)). This result is surprising because posterior correction is not a variance reduction method, rather a knowledge transfer method. The result offers a new perspective of variance reduction as knowledge transfer, for instance, through frequent mega-batch gradient computations.

Using this result, we derive new SVRG variants that go beyond existing proposals. For example, by using full-covariance Gaussians, we obtain a new variance-reduction method that employs Hessian corrections within a variational Online Newton algorithm. This differs from most works on Newton steps that only use corrections for the gradient and never for the Hessian (Derezinski, 2023; Sadiev et al., 2024; Garg et al., 2024; Sun et al., 2025). Another Adam-like extension is obtained by using diagonal covariances, implementing posterior correction over the IVON optimizer (Shen et al., 2024). Empirically, the new variant boosts IVON’s speed in non-traditional settings and shows promising results at scales up to pretraining an LLM from scratch on 50B tokens; see the middle and right panels Fig. 1. In this setting, our method outperforms Adam too, which is unlike many other traditional applications where no effective gains in performance are observed. This result also far exceeds the scale of any prior work on SVRG-based methods that use mega-batches. We validate these findings with various small to large experiments on several architectures, including ResNet,

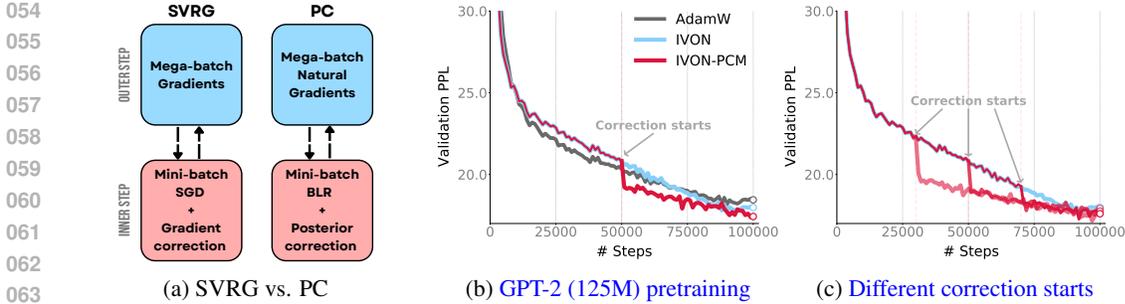


Figure 1: (a) We present a generalization of SVRG by using Posterior Correction (PC) where gradients used in SGD are replaced by natural gradients of VB objectives via the Bayesian Learning Rule (BLR). (b) Our new IVON-PCM (red) improves performance over IVON and AdamW when pretraining GPT2-125M from scratch on ca. 50B tokens from OpenWebText. Until 50K steps IVON-PC takes the same steps as IVON, and a huge boost is obtained when correction is started. Validation Perplexities at the end are 17.4, 18.0, 18.4. (c) We show three different IVON-PC runs where correction is started at a different iteration (pink to red). We see consistent improvements irrespective of the starting iteration.

Algorithm 1 SVRG

```

1: Initialize  $\theta_{in}$ 
2: while not converged do
3:    $\mathbf{g}_{out} \leftarrow \sum_{i=1}^N \nabla \ell_i(\theta_{in})$ 
4:    $\theta_{out} \leftarrow \theta_{in}$ 
5:   for  $t = 1, 2, \dots, m$  do
6:     Randomly pick  $i \in \{1, 2, \dots, N\}$ 
7:      $\mathbf{g}_{in} \leftarrow \nabla \ell_i(\theta_{in}) - \nabla \ell_i(\theta_{out}) + \frac{1}{N} \mathbf{g}_{out}$ 
8:      $\theta_{in} \leftarrow \theta_{in} - \eta \mathbf{g}_{in}$ 
9:   end for
10: end while

```

Algorithm 2 Posterior Correction (PC)

```

1: Initialize  $\lambda_{in}$ 
2: while not converged do
3:    $\tilde{\mathbf{g}}_{out} \leftarrow \sum_{i=1}^N \tilde{\nabla} \mathcal{L}_i(\lambda_{in})$ 
4:    $\lambda_{out} \leftarrow \lambda_{in}$ 
5:   for  $t = 1, 2, \dots, m$  do
6:     Randomly pick  $i \in \{1, 2, \dots, N\}$ 
7:      $\tilde{\mathbf{g}}_{in} \leftarrow \tilde{\nabla} \mathcal{L}_i(\lambda_{in}) - \tilde{\nabla} \mathcal{L}_i(\lambda_{out}) + \frac{1}{N} \tilde{\mathbf{g}}_{out}$ 
8:      $\lambda_{in} \leftarrow (1-\eta)\lambda_{in} - \eta N \tilde{\mathbf{g}}_{in}$ 
9:   end for
10: end while

```

Figure 2: Pseudo-code for SVRG (left) and our Bayesian generalization using PC (right). The latter replaces all instances of θ and $\nabla \ell_i$ in SVRG by the natural parameter λ (of the posterior q) and natural gradient $\tilde{\nabla} \mathcal{L}_i$. It also uses the BLR for the inner loop in line 8. We show that these differences disappear with an isotropic Gaussian q and SVRG is recovered as a special case of PC.

GPT and ViT. Overall, our work encourages further investigation of the new connection to obtain real improvements with SVRG in deep learning.

2 BACKGROUND ON SVRG AND VARIATIONAL BAYES

In SVRG, the goal is to minimize an empirical risk $\sum_{i=1}^N [l_i(\theta)/N]$ averaged over losses l_i for examples $i = 1, 2, \dots, N$. A regularizer, denoted by l_0 , is often added and handled by redefining the losses as $l_i + l_0/N$. Standard stochastic gradient descent (SGD) steps use stochastic gradients, for instance, the following update where the i 'th example is randomly sampled at each iteration,

$$\theta \leftarrow \theta - \eta \nabla l_i(\theta). \tag{1}$$

SVRG reduces the variance of such steps by using an outer loop where full-batch gradients are computed at parameters θ_{out} , and then parameters θ_{in} are updated using stochastic increments at a randomly sampled i at every iteration,

$$\theta_{in} \leftarrow \theta_{in} - \eta \left[\nabla l_i(\theta_{in}) - \nabla l_i(\theta_{out}) + \frac{1}{N} \sum_{j=1}^N \nabla l_j(\theta_{out}) \right]. \tag{2}$$

108 These steps can be seen as a full-batch gradient descent where old gradients $\nabla \ell_i(\theta_{\text{out}})$ are ‘corrected’
 109 by adding the fresh new gradients $\nabla \ell_i(\theta_{\text{in}})$. The use of the full-batch gradient can reduce variance
 110 and speed up future mini-batch steps. Pseudo-code is shown in Alg. 1 where m steps are taken in
 111 the inner loop and a constant learning rate η is used.

112 SVRG is built upon earlier ideas in Stochastic Average Gradient (SAG) (Roux et al., 2012; Schmidt
 113 et al., 2017) and Stochastic Dual Coordinate Descent (SDCA) (Shalev-Shwartz & Zhang, 2013)
 114 methods where the full batch gradient is stored for all examples separately and entries are re-
 115 freshed whenever corresponding examples are chosen during updating. Subsequently, many new
 116 practical variants have been proposed, for instance, SARAH (Nguyen et al., 2017) and SPI-
 117 DER (Fang et al., 2018), among many other proposals (Dubois-Taine et al., 2022; Lei et al., 2017;
 118 Babanezhad Harikandeh et al., 2015). Instead of full-batch, it is also possible to use large *mega-*
 119 *batches* which can be, for example, 10-50 times bigger than the mini-batches. It is also helpful to
 120 use Adam and down-weight the corrections (Yin et al., 2025). There is a large number of papers that
 121 show variance reduction is useful for accelerating both convex and nonconvex optimization.

122 Despite this, variance reduction has not seen success in deep learning yet. The ineffectiveness is
 123 studied extensively by Defazio & Bottou (2019) who find little to no gain in run-time for traditional
 124 deep learning. They show that variance can even increase if mega-batch gradients are not refreshed
 125 often, in the end giving no effective gain in speed. The matter is complicated by the presence of
 126 additional deep learning tricks, such as, mini-batching, momentum, learning-rate schedules, etc.
 127 Numerous works have suggested new improvements but SVRG’s ineffectiveness is not fixed yet
 128 (Yin et al., 2025; Tondji et al., 2021; Cutkosky & Orabona, 2019; Arnold et al., 2019; Ma & Yarats,
 129 2019).

130 Clearly, using SVRG to get a real speed-up in deep learning is not easy. Here, we ask if the same
 131 is true for non-traditional settings. We consider variational training of deep networks where the
 132 IVON optimizer has lead to promising results. Such Bayesian approaches hold promise to facilitate
 133 continual, federated, and active learning of deep networks. Use of mega-batches can enable us to
 134 build priors to speed up future learning. Currently, no such work exists in this space. With this in
 135 mind, we explore new connections between SVRG and Bayes and assess the potential of SVRG-
 136 style ideas for non-traditional deep learning.

137 2.1 VARIATIONAL BAYES (VB) AND BAYESIAN LEARNING RULE (BLR)

138 Throughout, we will use the variational-Bayesian (VB) generalization of empirical-risk minimiza-
 139 tion (ERM) which optimizes over a tractable set of distributions $q(\theta) \in \mathcal{Q}$ instead of a point estimate
 140 $\theta \in \Theta$. This formulation is crucial for us to connect and extend SVRG to Bayes. Below we show
 141 an ERM problem on the left and its VB formulation on the right,

$$142 \theta_* = \arg \min_{\theta \in \Theta} \sum_{i=0}^N \ell_i(\theta) \quad \text{vs} \quad q_* = \arg \min_{q \in \mathcal{Q}} \sum_{i=1}^N \mathbb{E}_q[\ell_i] + \mathbb{D}_{\text{KL}}[q \| p_0]. \quad (3)$$

143 The ERM above contains a regularizer, denoted by ℓ_0 , which is also used to define a valid prior
 144 distribution $p_0 \propto \exp(-\ell_0)$ that is added as a Kullback-Leibler (KL) divergence penalty. If ℓ_i is a
 145 proper likelihood, q_* produces a tractable approximation to the posterior $p_* \propto p_0 \prod_i \exp(-\ell_i)$.

146 Often we restrict \mathcal{Q} to be exponential-family (EF) distributions, for instance, Gaussians. An EF
 147 takes the following log-linear form with respect to a sufficient statistic, denoted by $\mathbf{T}(\theta)$, as shown
 148 below with an example of an isotropic Gaussian where $\mathbf{T}(\theta) = \theta$,

$$149 q(\theta) \propto h(\theta) \exp(\boldsymbol{\lambda}^\top \mathbf{T}(\theta)), \quad \mathcal{N}(\theta | \mathbf{m}, \mathbf{I}) \propto e^{-\frac{1}{2} \theta^\top \theta} \exp(\mathbf{m}^\top \theta). \quad (4)$$

150 The distribution is parameterized by the natural parameter $\boldsymbol{\lambda}$ and uses a base measure $h(\theta)$. For
 151 an isotropic Gaussian, $\boldsymbol{\lambda} = \mathbf{m}$ and $h(\theta) = \exp(-\frac{1}{2} \theta^\top \theta)$. Throughout, we will use the natural
 152 parameterization $\boldsymbol{\lambda}$ because updates over it naturally generalize those used in ERM for θ .

153 Natural gradients are convenient for optimizing Eq. 3, even though for ERM they can be computa-
 154 tionally expensive. The convenience is that by using the expectation parameter $\boldsymbol{\mu} = \mathbb{E}_q[\mathbf{T}(\theta)]$ we
 155 can easily compute the natural gradient of $\mathcal{L}_i(\boldsymbol{\lambda}) = \mathbb{E}_q[\ell_i]$ while avoid computing the Fisher $\mathbf{F}(\boldsymbol{\lambda})$,

$$156 \widetilde{\nabla} \mathcal{L}_i = \mathbf{F}(\boldsymbol{\lambda})^{-1} \nabla \mathcal{L}_i = \nabla_{\boldsymbol{\mu}} \mathcal{L}_i. \quad (5)$$

The natural gradient of the KL term can also be simplified. Thus, the natural-gradient descent (NGD) update takes a much simpler form which resembles Bayes’ rule. This is called the Bayesian Learning Rule (BLR) (Khan & Rue, 2023) and is shown below in two equivalent forms: using the NGD update (shown on the left) and in its Bayes’ form (shown on the right),

$$\boldsymbol{\lambda} \leftarrow (1 - \eta)\boldsymbol{\lambda} - \eta \sum_{i=0}^N \tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}) \quad \Leftrightarrow \quad q \leftarrow q^{1-\eta} \prod_{i=0}^N \exp\left(-\eta \hat{\ell}_i\right). \quad (6)$$

The NGD update is simply a moving average of the natural gradients. The Bayes’ form is obtained by substituting $\boldsymbol{\lambda}$ obtained with an NGD step into the EF form of Eq. 4 and by defining the *site* functions $\hat{\ell}_i(\boldsymbol{\theta}) = \tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda})^\top \mathbf{T}(\boldsymbol{\theta})$ of the losses ℓ_i . The form is equivalent to Bayes’ rule with prior $q^{1-\eta}$ and likelihood $\exp(-\eta \hat{\ell}_i)$.

The BLR is attractive not only because it closely mirrors Bayes’ rule, but also because it subsumes many ERM algorithms (Khan & Rue, 2023). For example, SGD (Eq. 1) can be derived with isotropic Gaussians q (Eq. 4). Since $\boldsymbol{\lambda} = \mathbf{m}$, $\boldsymbol{\mu} = \mathbb{E}_q[\boldsymbol{\theta}] = \mathbf{m}$, and $\nabla \log h(\boldsymbol{\theta}) = -\boldsymbol{\theta}$, we can write Eq. 6 as

$$\mathbf{m} \leftarrow \mathbf{m} - \eta \sum_{i=0}^N \nabla \mathcal{L}_i(\mathbf{m}) \quad \Leftrightarrow \quad q \leftarrow q^{1-\eta} \prod_{i=0}^N \exp\left(-\eta \nabla \mathcal{L}_i(\mathbf{m})^\top \boldsymbol{\theta}\right). \quad (7)$$

To derive SGD, we use the delta method to approximate $\nabla \mathcal{L}_i(\mathbf{m}) \approx \nabla \ell_i(\mathbf{m})$ and write the stochastic version by redefining $\ell_i + \ell_0/N$. We will use this same strategy to connect gradient correction in SVRG to natural-gradient correction in the BLR. When using $q = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \boldsymbol{\Sigma})$, we obtain the following Newton-like update, called Variational Online Newton (VON), where the pre-conditioner is the precision matrix $\mathbf{S} = \boldsymbol{\Sigma}^{-1}$,

$$\mathbf{m} \leftarrow \mathbf{m} - \eta \mathbf{S}^{-1} \sum_{i=0}^N \nabla \mathcal{L}_i(\mathbf{m}), \quad \text{where} \quad \mathbf{S} \leftarrow (1 - \eta)\mathbf{S} + \eta \sum_{i=0}^N \nabla^2 \mathcal{L}_i(\mathbf{m}) \quad (8)$$

An improved version of VON is proposed by Lin et al. (2020) and used in Shen et al. (2024) to train large deep networks (such as GPT-2) and obtain competitive results to the Adam optimizer. We will derive a variant of this algorithm later as an extension of SVRG.

3 SVRG AND BEYOND VIA POSTERIOR CORRECTION

We will now present a Bayesian approach that generalizes SVRG and enables us to derive new extensions to improve variational deep learning. The approach uses Posterior Correction (PC) (Khan, 2025) which unifies various knowledge transfer tasks, such as, continual learning, federated learning, and model merging. We will now show that it can also be used to recover SVRG.

3.1 POSTERIOR CORRECTION

Posterior correction aims to speed-up training by reusing previously computed posteriors, for example, in the form of old checkpoints Khan (2025, Sec. 4.2). We denote the checkpoint by natural parameters $\boldsymbol{\lambda}_{\text{out}}$ and compute the posterior via natural gradients at $\boldsymbol{\lambda}_{\text{out}}$ over the whole dataset:

$$\hat{q}_{\text{out}} \leftarrow \prod_{i=0}^N \exp\left(-\hat{\ell}_{i|\text{out}}\right) \quad \text{where} \quad \hat{\ell}_{i|\text{out}}(\boldsymbol{\theta}) = \tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{\text{out}})^\top \mathbf{T}(\boldsymbol{\theta}). \quad (9)$$

We can use this to ‘correct’ the future BLR updates by simply multiplying and dividing by $\hat{q}_{\text{out}}^\eta$ in the RHS of the Bayes’ form shown in Eq. 6, with added parts highlighted in red,

$$q \leftarrow q^{1-\eta} \hat{q}_{\text{out}}^\eta \prod_{i=0}^N \exp\left(-\eta \left[\hat{\ell}_i - \hat{\ell}_{i|\text{out}}\right]\right). \quad (10)$$

The update remains unchanged because we have simply multiplied it by $\mathbf{1}$, but the new form can be seen as Bayes’ rule on a modified model where the prior includes contributions from \hat{q}_{out} and the likelihood is ‘corrected’. Khan (2025) argue that quick adaptation is possible by reducing the correction and show that many existing schemes perform such corrections. Here, we will use posterior correction to generalize SVRG.

Algorithm 3 VSGD-PC: Variational SGD with posterior correction

Initialize: Number of inner steps m , learning rates η

- 1: Initialize \mathbf{m}_{in} .
- 2: **while** not converged **do**
- 3: $\mathbf{g}_{out} \leftarrow \sum_{i=1}^N \nabla \ell_i(\boldsymbol{\theta}_{in})$ **where** $\boldsymbol{\theta}_{in} = \mathbf{m}_{in} + \boldsymbol{\epsilon}$ **with** $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$
- 4: $\boldsymbol{\theta}_{out} \leftarrow \boldsymbol{\theta}_{in}$
- 5: **for** $t = 1, 2, \dots, m$ **do**
- 6: Randomly pick $i \in \{1, 2, \dots, N\}$
- 7: $\mathbf{g}_{in} \leftarrow \nabla \ell_i(\boldsymbol{\theta}_{in}) - \nabla \ell_i(\boldsymbol{\theta}_{out}) + \frac{1}{N} \mathbf{g}_{out}$ **where** $\boldsymbol{\theta}_{in} = \mathbf{m}_{in} + \boldsymbol{\epsilon}$ **with** $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$
- 8: $\mathbf{m}_{in} \leftarrow \mathbf{m}_{in} - \eta \mathbf{g}_{in}$
- 9: **end for**
- 10: **end while**

3.2 POSTERIOR CORRECTION GENERALIZES SVRG

To mirror the inner loop Eq. 2 of SVRG, we write a stochastic version of Eq. 10 where one example is sampled. Analogously to Eq. 2, we denote the inner loop iterate as q_{in} . We absorb the regularizer in $\ell_i \leftarrow (\ell_i + \ell_0/N)$, sample a random example i and weight the correction term by N to get

$$q_{in} \leftarrow q_{in}^{1-\eta} \hat{q}_{in}^\eta \exp\left(-\eta N \left[\hat{\ell}_{i|in} - \hat{\ell}_{i|out}\right]\right). \quad (11)$$

When written in terms of $\boldsymbol{\lambda}$, the update takes an identical form to Eq. 2, as shown below:

$$\boldsymbol{\lambda}_{in} \leftarrow (1 - \eta)\boldsymbol{\lambda}_{in} + \eta \hat{\boldsymbol{\lambda}}_{out} - \eta N \left[\tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{in}) - \tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{out})\right] \quad (12)$$

$$\implies \boldsymbol{\lambda}_{in} \leftarrow (1 - \eta)\boldsymbol{\lambda}_{in} - \eta N \left[\tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{in}) - \tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{out}) + \frac{1}{N} \sum_{j=1}^N \tilde{\nabla} \mathcal{L}_j(\boldsymbol{\lambda}_{out})\right]. \quad (13)$$

The second update is obtained by using Eq. 9 which shows that $\hat{\boldsymbol{\lambda}}_{out} = \sum_i \tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{out})$. The update is strikingly similar to Eq. 2 with all instances of gradients $\nabla \ell_i$ replaced by natural gradients $\tilde{\nabla} \mathcal{L}_i$. Alg. 2 shows this algorithm, where we denote $\hat{\boldsymbol{\lambda}}_{out}$ by $\tilde{\mathbf{g}}_{out}$ and use the BLR in the inner-loop update.

We will now show our main result that SVRG is a special case of Alg. 2, when we use PC over the isotropic-Gaussian family, that is, we set \mathcal{Q} to be a set of $q_{in} = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}_{in}, \mathbf{I})$.

Theorem 1 For isotropic-Gaussian family, Eq. 13 reduces to the following update of the mean \mathbf{m}_{in} ,

$$\mathbf{m}_{in} \leftarrow \mathbf{m}_{in} - \eta N \left[\mathbb{E}_{q_{in}}[\nabla \ell_i] - \mathbb{E}_{q_{out}}[\nabla \ell_i] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{out}}[\nabla \ell_j]\right]. \quad (14)$$

The proof follows similarly to Eq. 7 by plugging-in the definitions of $\boldsymbol{\lambda}$, $\boldsymbol{\mu}$, and $h(\boldsymbol{\theta})$, and using Bonnet’s theorem (Bonnet, 1964) which states that $\nabla_{\mathbf{m}} \mathcal{L}_i = \mathbb{E}_q[\nabla \ell_i]$. That is, the natural gradients can be computed by computing the expectation of $\ell_i(\boldsymbol{\theta})$ at sample $\boldsymbol{\theta} \sim q$, which is convenient for implementation. An algorithm to implement the update is given in Alg. 3, which we call VSGD-PC and where for simplicity we use one Monte Carlo (MC) sample. **VSGD-PC and using SVRG with SGD have almost the same cost, except for the additional parameter sampling which is fast to compute.** The algorithm is almost identical to SVRG (Alg. 1) with differences highlighted in red but can exactly match SVRG as shown below.

Theorem 2 SVRG (Alg. 1) is equivalent to VSGD-PC (Alg. 3) where we set $\boldsymbol{\epsilon} = 0$ in line 3 and 7.

Setting $\boldsymbol{\epsilon} = 0$ is equivalent to applying the delta method: $\mathbb{E}_{q_{in}}[\nabla \ell_i] \approx \nabla \ell_i(\mathbf{m})$ which we also use in Eq. 7 to recover SGD from the BLR. We will show in the next section that, by using other EF forms, we can derive new extensions that go beyond SVRG. We note that this is the first result of its kind connecting SVRG and Bayes. Previous works have applied SVRG to improve ELBO optimization (Zhang et al., 2018), but are merely using SVRG and unable to recover or extend it.

3.3 BEYOND SVRG: NEW EXTENSIONS DERIVED USING POSTERIOR CORRECTION

In this section, we derive new extensions to boost variational training of deep networks.

A Newton-like Variant: The original SVRG and many of its variants focus on stabilizing the gradient and little work has been done on methods that do the same for the Hessian. We derive such variants by using more flexible Gaussian forms than isotropic Gaussians, for example, the multi-variate Gaussian $q = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{S}^{-1})$, where the precision matrix \mathbf{S} is also estimated. Then, Alg. 2 reduces to a posterior correction version of the Variational-Online-Newton (VON) algorithm shown in Eq. 8. The update shown below employs a stochastic variance reduction for the Hessian.

Theorem 3 For Gaussian $q_{in} = \mathcal{N}(\mathbf{m}_{in}, \mathbf{S}_{in}^{-1})$, Eq. 13 reduces to a Newton-like update,

$$\mathbf{m}_{in} = \mathbf{m}_{in} - \eta N \mathbf{S}_{in}^{-1} \left[\mathbb{E}_{q_{in}}[\nabla \ell_i] - \mathbb{E}_{q_{out}}[\nabla \ell_i] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{out}}[\nabla \ell_j] + \mathbf{H}_{out \setminus i}(\mathbf{m}_{in} - \mathbf{m}_{out}) \right] \quad (15)$$

where we use a Stochastic Variance-Reduced Hessian (SVRH) estimate as the pre-conditioner

$$\mathbf{S}_{in} \leftarrow (1 - \eta) \mathbf{S}_{in} + \eta N \left[\mathbb{E}_{q_{in}}[\nabla^2 \ell_i] + \bar{\mathbf{H}}_{out \setminus i} \right]. \quad (16)$$

Here, $\bar{\mathbf{H}}_{out \setminus i} = \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{out}}[\nabla^2 \ell_j] - \mathbb{E}_{q_{out}}[\nabla^2 \ell_i]$ is the full-batch expected Hessian without ℓ_i .

The derivation uses the definition of natural parameter and gradient in Eq. 13 and is given in App. A. An algorithm is in Alg. 5 which we name VON-PC, as a posterior correction version of VON.

We are not aware of any other Newton variant that implements similar Hessian corrections as shown above. Most works on Newton steps only use corrections for the gradient and never for the Hessian (Derezinski, 2023; Sadiev et al., 2024; Garg et al., 2024; Sun et al., 2025). One surprising connection is that the term $\mathbf{H}_{out \setminus i}(\mathbf{m}_{in} - \mathbf{m}_{out})$ in Eq. 15 is also used in Chayti et al. (2024, Eqs. 11–12) but is derived differently via cubic-Newton. The method has rigorous guarantees and its relation to our Bayesian approach remains an interesting case to study. The term can also be viewed as forcing the inner iteration to stay close to the most recent outer iteration.

The PC method can be applied to any EF distribution and therefore yields novel extensions that go way beyond SVRG, for example, for binary neural networks (Meng et al., 2020) to yield SVRG-style updates for the Straight-Through Estimator (Bengio et al., 2013) via Bernoulli distributions. We omit the derivation because the procedure is similar to the ones we presented here. The PC method can yield novel extensions of SVRG by exploiting flexible EF distributions.

An Adam-like variant: A cheaper Adam-like variant is obtained by using diagonal covariance, for example, $q = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \text{diag}(\mathbf{s})^{-1})$. This implements posterior correction over IVON (Shen et al., 2024). A detailed derivation is in App. B and a pseudo-code of IVON-PCM is in Alg. 4, where we highlight additional computation on top of IVON. Similarly to IVON, IVON-PCM uses several practical tricks, for instance, weight decay, mini-batching, temperature, momentum, and debiasing. If momentum is switched off, i.e. $\rho_1 = \rho_2 = 1$, we simply refer to the method as IVON-PC.

Handling Mega-Batches and Weighted Corrections: The version shown in Alg. 4 also avoids expensive full-batch computations. Instead, it uses *mega*-batches (Defazio & Bottou, 2019) and builds online estimates of gradients and Hessians (line 3-4). Mega-batches can be tens of times the size of the inner loop mini-batches and can be used to slowly build an estimate of full-batch gradients and Hessians via an estimate of the site function, for example, for isotropic Gaussians:

$$\hat{q}_{out} \leftarrow \exp\left(-\boldsymbol{\theta}^\top \mathbf{g}_{out}\right), \quad \text{where } \mathbf{g}_{out} \leftarrow \rho_1 \mathbf{g}_{out} + (1 - \rho_1) \sum_{i \in \mathcal{M}} \mathbb{E}_{q_{out}}[\nabla \ell_i]. \quad (17)$$

A similar approach for full-Gaussians is in App. B.1. Posteriors constructed using mega-batches can be imperfect and downweighted by modifying the PC update of Eq. 11 using $\alpha < 1$,

$$q_{in} \leftarrow q_{in}^{1-\eta} \hat{q}_{out}^{\eta\alpha} \exp\left(-\eta N \left[\hat{\ell}_{i|in} - \alpha \hat{\ell}_{i|out} \right]\right), \quad (18)$$

For $\alpha = 0$, the update reduces to standard BLR, while for $\alpha = 1$ we use perfect corrections which are good for full-batch \hat{q}_{out} . When using mega-batches, a smaller value could be used and tuned.

Algorithm 4 IVON-PCM: IVON with Posterior Correction and Momentum. IVON-PC is obtained by removing momentum (differences to IVON highlighted in red)

Require: Learning rates $\{\eta_t\}$, $\beta_1 \in [0, 1)$, $\beta_2 \in [0, 1)$, $\delta > 0$, $\kappa > 0$, $h_0 > 0$, clip radius $\xi > 0$, mini-batch size B , **mega-batch size M** , **outer loop learning rate ρ_1, ρ_2** , **refresh rate α** .

```

1: Initialize:  $\mathbf{m}_{\text{in}} \leftarrow$  (NN weight init),  $\mathbf{h}_{\text{in}} \leftarrow h_0$ ,  $\sigma_{\text{in}} \leftarrow 1/\sqrt{\kappa(\mathbf{h}_{\text{in}} + \delta)}$ ,  $\mathbf{g} \leftarrow 0$ 
2: while not converged do
3:    $\hat{\mathbf{g}}_{\text{out}} \leftarrow \frac{1}{M} \sum_{i \in \mathcal{M}} \nabla \ell_i(\boldsymbol{\theta}_{\text{in}})$    where we sample a mega-batch  $\mathcal{M}$  and  $\boldsymbol{\theta}_{\text{in}} \sim \mathcal{N}(\mathbf{m}_{\text{in}}, \sigma_{\text{in}}^2)$ 
4:    $\mathbf{g}_{\text{out}} \leftarrow \rho_1 \mathbf{g}_{\text{out}} + (1 - \rho_1) \hat{\mathbf{g}}_{\text{out}}$ ,   and    $\mathbf{h}_{\text{out}} \leftarrow \rho_2 \mathbf{h}_{\text{out}} + (1 - \rho_2) \hat{\mathbf{g}}_{\text{out}}(\boldsymbol{\theta}_{\text{in}} - \mathbf{m}_{\text{in}})/\sigma_{\text{in}}^2$ 
5:    $\mathbf{m}_{\text{out}} \leftarrow \mathbf{m}_{\text{in}}$ ,  $\sigma_{\text{out}} \leftarrow \sigma_{\text{in}}$ 
6:   for  $t = 1, 2, \dots, m$  do
7:     Sample a mini-batch  $\mathcal{B}$ ,  $\boldsymbol{\theta}_{\text{in}} \sim \mathcal{N}(\mathbf{m}_{\text{in}}, \sigma_{\text{in}}^2)$ ,  $\boldsymbol{\theta}_{\text{out}} \sim \mathcal{N}(\mathbf{m}_{\text{out}}, \sigma_{\text{out}}^2)$ 
8:      $\hat{\mathbf{g}}_{\text{in}} \leftarrow \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla \ell_i(\boldsymbol{\theta}_{\text{in}})$    and    $\hat{\mathbf{h}}_{\text{in}} \leftarrow \hat{\mathbf{g}}_{\text{in}}(\boldsymbol{\theta}_{\text{in}} - \mathbf{m}_{\text{in}})/\sigma_{\text{in}}^2$ 
9:      $\hat{\mathbf{g}}_{\text{out}} \leftarrow \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla \ell_i(\boldsymbol{\theta}_{\text{out}})$    and    $\hat{\mathbf{h}}_{\text{out}} \leftarrow \hat{\mathbf{g}}_{\text{out}}(\boldsymbol{\theta}_{\text{out}} - \mathbf{m}_{\text{out}})/\sigma_{\text{out}}^2$ 
10:     $\hat{\mathbf{g}} \leftarrow \hat{\mathbf{g}}_{\text{in}} - \alpha(\hat{\mathbf{g}}_{\text{out}} - \mathbf{g}_{\text{out}})$    and    $\hat{\mathbf{h}} \leftarrow \hat{\mathbf{h}}_{\text{in}} - \alpha(\hat{\mathbf{h}}_{\text{out}} - \mathbf{h}_{\text{out}})$ 
11:     $\mathbf{g} \leftarrow \beta_1 \mathbf{g} + (1 - \beta_1) \hat{\mathbf{g}}$ 
12:     $\mathbf{h} \leftarrow \beta_2 \mathbf{h} + (1 - \beta_2) \hat{\mathbf{h}} + \frac{1}{2}(1 - \beta_2)^2(\mathbf{h} - \hat{\mathbf{h}})^2/(\mathbf{h} + \delta)$ 
13:     $\bar{\mathbf{g}} \leftarrow \mathbf{g}/(1 - \beta_1^t)$ 
14:     $\bar{\mathbf{g}} \leftarrow (\bar{\mathbf{g}} + \delta \mathbf{m}_{\text{in}} + \alpha(\mathbf{h}_{\text{out}} - \hat{\mathbf{h}}_{\text{out}})(\mathbf{m}_{\text{in}} - \mathbf{m}_{\text{out}}))/(\mathbf{h} + \delta)$ 
15:     $\mathbf{m}_{\text{in}} \leftarrow \mathbf{m}_{\text{in}} - \eta_t \text{clip}(\bar{\mathbf{g}}, \xi)$ 
16:     $\sigma_{\text{in}} \leftarrow 1/\sqrt{\kappa(\mathbf{h} + \delta)}$ 
17:   end for
18: end while
19: return  $\mathbf{m}, \sigma$ 

```

When applied to isotropic Gaussians, this reduces to α -SVRG (Yin et al., 2025), though their motivation to use α does not stem from the use of mega-batches (they appear to use full batches). Their motivation is to reduce variance early on via scheduling α , starting with a high value. From a Bayesian perspective, the \hat{q}_{out} estimates are expected to be less useful early on. Thus, it makes more sense to use a ‘burn-in’ period with $\alpha = 0$ and then turn on α . In our experiments, we use a constant α , which seems to work better for deep learning.

3.4 COMPUTATIONAL & MEMORY REQUIREMENTS

The computation and memory overheads of IVON-PCM are similar to the use of Adam to implement alpha-SVRG. Adam uses an accumulation of squared gradients, while here reparameterization trick is used to compute a Hessian estimate (line 4, 8, 9). Unlike alpha-SVRG, IVON-PCM use an additional Hessian correction as well, but this does not add a significant cost because the Hessian is already computed. We just need to store an additional \mathbf{h}_{out} in the outer loop (in line 9), as well as σ_{out} (line 5). Similarly to VSGD-PC, sampling is added in line 3 and 7. All of these costs are not significant. Just like in SVRG and alpha-SVRG, the major overhead is the mega-batch computation and the use of two gradient (line 10). An extra Hessian correction is used in line 10 as well.

4 EXPERIMENTS & RESULTS

4.1 LOGISTIC REGRESSION

We first illustrate our new method IVON-PC on simple convex logistic regression problems of varying dimensionality and number of data examples in Fig. 3. We show additional results on CIFAR-10 logistic regression in App. C. We compare IVON-PC against IVON, SVRG, and SGD. For all experiments, we used a large constant learning rate and a batch size of 5. For the IVON methods, we set $\rho_1 = \rho_2 = 1$ (i.e. no momentum) and downweight the extra terms in line 17 of Alg. 4 by 0.01. The effects of varying ρ and the additional term are explored in our ablation experiments in App. C.

Consistent with the original work by Johnson & Zhang (2013), our experiments show that SVRG drastically boosts the performance of SGD once the first outer loop—illustrated by the gray bars,

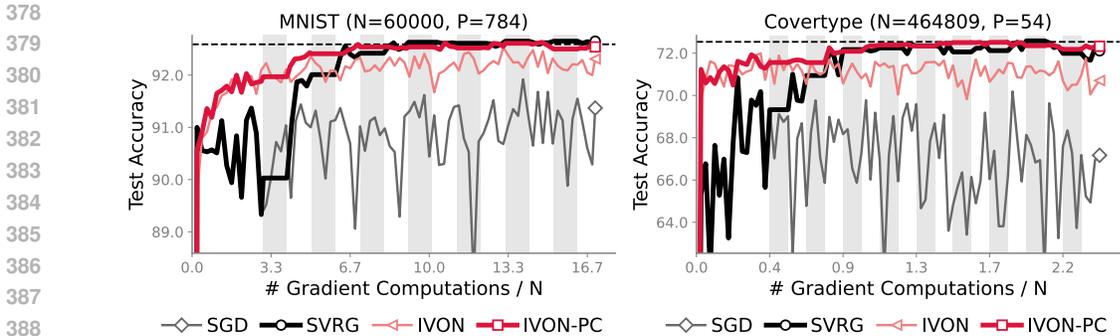


Figure 3: IVON-PC significantly boosts the convergence speed of IVON and performs much better than SVRG, here on three convex logistic regression problems of varying dimension and size. The horizontal dashed line indicates the performance at the minimum, the gray bars indicate outer gradient computations used in SVRG and IVON-PC.

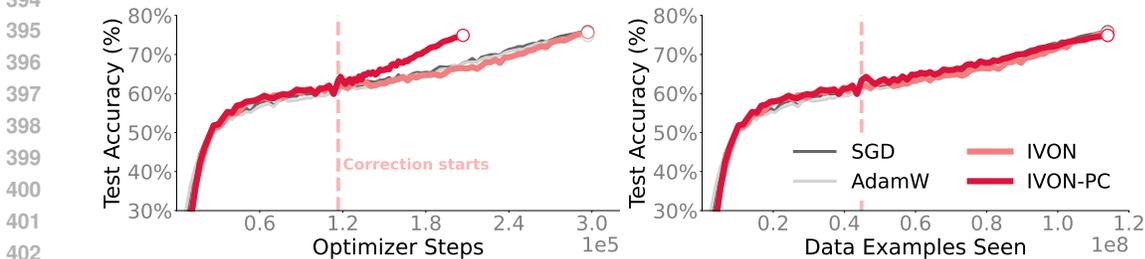


Figure 4: Performance on ImageNet for ResNet-50. When comparing by the number of optimization steps (left) IVON-PC gives clear improvements but not in terms of data examples seen (right). On the left, we zoom in on the final stage of training when correction is added.

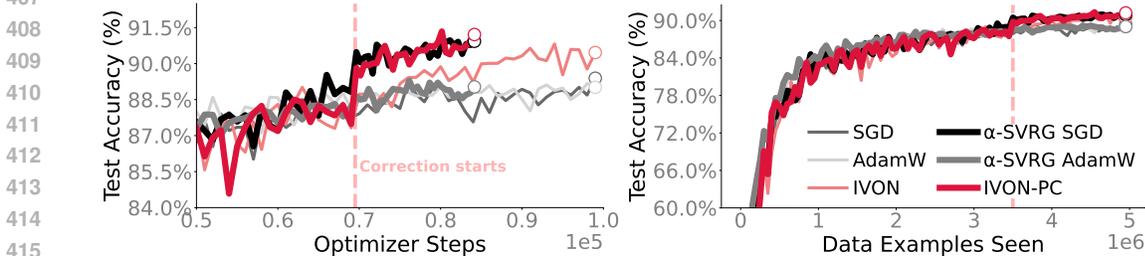


Figure 5: Comparison to α -SVRG on CIFAR-10 for ResNet-20. IVON-PC and α -SVRG with SGD give improvements, also when counting the number of data examples seen (right). At the end of training, IVON-PC performs best out of all methods.

which indicate the large outer batch size—is hit. Similarly, IVON-PC boosts the performance of IVON. Since IVON already outperforms SGD and is comparable to SVRG due to its inherent preconditioning and momentum, the further boost provided by IVON-PC makes it perform best.

4.2 IMAGE CLASSIFICATION WITH RESNETS

We first evaluate our proposed IVON-PC method on training a ResNet-50 on ImageNet. Fig. 4 shows the results. We can see that IVON-PC reaches a comparable accuracy to SGD, IVON and AdamW in much fewer optimization steps. Factoring in the outer loop gradient calculation, all methods perform similarly well. Variance reduction methods such as SVRG have classically struggled on training ResNets on ImageNet (Defazio & Bottou, 2019). We observe a similar effect but performance is greatly improved when just counting optimization steps as the recent work by Yin et al. (2025).

Table 1: IVON-PC can improve finetuning. All methods use the same number of steps. We indicate with “+” improvements over the IVON baseline.

	ViT-B/32				Qwen2.5-0.5B-it		Llama-3.1-8B
	Cars	DTD	GTSRB	RESISC45	XSUM (R-1, R-L)		GSM8k
IVON	79.5	72.9	99.9	95.2	48.7	23.6	30.4
+ PC	80.0_{+0.5}	73.4_{+0.5}	99.9_{+0.0}	96.1_{+0.9}	49.6_{+0.9}	23.8_{+0.2}	30.6_{+0.2}

In a second experiment on a smaller ResNet-20 in Fig. 5, we do not anneal the learning rate to zero but to a quarter of the starting learning rate. Then, the variance is not completely removed through learning rate annealing and both α -SVRG with SGD and IVON-PC bring improvements in accuracy over their respective baseline methods, with IVON-PC having the advantage of estimating a variational posterior distribution. All details for the two experiments are in App. D.1.

4.3 LANGUAGE MODEL PRETRAINING

Here, we present results when pretraining a 125M parameter GPT-2 model from scratch on 50B tokens from the OpenWebText dataset¹. We follow the set-up from Shen et al. (2024) and train each model for 100,000 steps using AdamW, IVON, and IVON-PCM. For IVON-PCM we use two different configurations: One refreshes a megabatch with size 10 times the minibatch size after 10 inner step starting correction after 50,000 steps and another where we start correction at varying points, namely, after 30,000, 50,000, and 70,000 steps with a megabatch factor and refresh rate of 20 as opposed to 10. Results are shown in Fig. 1b and Fig. 1c, respectively. We find that IVON-PCM provides improvements in terms of validation perplexity in both settings. Interestingly, correction can be started at varying intervals and benefits, especially a direct jump downwards of the validation perplexity, are maintained. Details are found in App. D.2.

4.4 CONTINUAL PRETRAINING

Next, we present results on continually pretraining the GPT-125M model from Shen et al. (2024) on 1B tokens from Fineweb-edu (Penedo et al., 2024). We compare AdamW and alpha-SVRG with AdamW against IVON, IVON-PC, and IVON-PCM. Both alpha-SVRG and the IVON-PC runs use $\alpha = 0.7$, 40 steps for the inner loop, and 1,000 warmup steps without correction. Results are in Fig. 6. Adding correction both with α -SVRG and IVON-PC helps but the gap between IVON, which struggles with larger learning rates, and IVON-PC is larger. Adding Momentum to IVON-PC improves results beyond those of α -SVRG and IVON-PCM converges to a better validation perplexity after the same number of steps. When comparing time, α -SVRG does not improve over AdamW but IVON-PCM improves over IVON. Details can be found in App. D.3.

4.5 FINETUNING

Here, we use IVON-PC for finetuning different Transformers, namely ViT-B/32 (Dosovitskiy et al., 2021), Qwen2.5-0.5B-Instruct (Yang et al., 2025), and Llama-3.1-8B (Dubey et al., 2024) on image classification, XSUM, and GSM8k, respectively. For the ViT models we finetune only the image coder, for Qwen we finetune the entire model, and for Llama-3.1 we use LoRA finetuning (Hu et al., 2022). Results are shown in Table 1 and show that IVON-PC can improve final performance over IVON when the models are trained with the same number of optimization steps. Details and hyperparameters used for these experiments can be found in App. D.4.

4.6 ABLATIONS

Influence of α : We study the influence of α for α -SVRG and IVON-PC on wikitext103 (Merity et al., 2017) by training a GPT-2-based model with 33M parameters from scratch for one epoch with batch size of 64 and 5000 warmup steps without correction. Results are shown in Fig. 6 (left): for both methods using correction improves performance over various values of α . Both times the

¹<https://huggingface.co/datasets/SkyLion007/openwebtext>

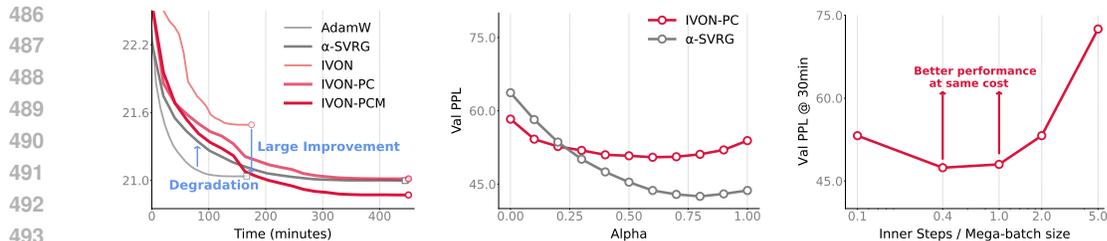


Figure 6: (left) IVON-PC and IVON-PCM also boost IVON’s performance when continually pre-training a GPT2-125M model on 1B tokens from Fineweb-Edu, here measured in terms of wallclock time. This is unlike AdamW that does not benefit from SVRG when using α -SVRG. (center) Tuning α has a large effect on performance on wikitext103 with a 33M Transformer trained from scratch. (right) In the same setting, increasing the number of refreshes can help but at some point becomes slow, with better performance at the same time budget obtained by fewer refreshes.

curve has a u-shape indicating that biasing too much towards large batches can harm performance. We provide details for these and the following analysis of the inner loop size in App. D.2.

Influence of Inner Loop Size Here, we fix the megabatch size to $50 \cdot 64$ in the same setting as above and vary the number of refreshes. Fig. 6 (center) shows that refreshing more often in general helps performance but when refreshing too often performance can get worse at the same compute budget.

5 CONCLUSION

In this paper, we present surprising new connections between two seemingly unrelated ideas: SVRG and posterior correction (PC). The connections allow us to derive new extensions of SVRG which can boost variational training of deep networks. We hope that the non-traditional settings considered in this paper encourage researchers who wish to see SVRG successfully applied to deep learning. Our work attempts to offer a new perspective on variance reduction as knowledge transfer, giving SVRG-style ideas a fresh restart. In the future, we hope to apply these ideas to other non-traditional problems where reusing past knowledge is crucial for reducing cost.

REFERENCES

- Sébastien Arnold, Pierre-Antoine Manzagol, Reza Babanezhad Harikandeh, Ioannis Mitliagkas, and Nicolas Le Roux. Reducing the variance in online optimization by transporting past gradients. *Advances in Neural Information Processing Systems*, 32, 2019. URL https://papers.neurips.cc/paper_files/paper/2019/hash/1dba5eed8838571e1c80af145184e515-Abstract.html.
- Reza Babanezhad Harikandeh, Mohamed Osama Ahmed, Alim Virani, Mark Schmidt, Jakub Konečný, and Scott Sallinen. Stopwasting my gradients: Practical svrg. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/a50abba8132a77191791390c3eb19fe7-Paper.pdf.
- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. URL <https://arxiv.org/abs/1308.3432>.
- G. Bonnet. Transformations des signaux aléatoires a travers les systemes non linéaires sans mémoire. In *Annales des Télécommunications*, volume 19, pp. 203–220. Springer, 1964. URL <https://doi.org/10.1007/BF03014720>.
- El Mahdi Chayti, Martin Jaggi, and Nikita Doikov. Unified convergence theory of stochastic and variance-reduced cubic newton methods. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=FCs5czlDTr>.

- 540 Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote Sensing Image Scene Classification: Bench-
541 mark and State of the Art. *Proceedings of the IEEE*, 105(10):1865–1883, Oct 2017. ISSN 1558-
542 2256. doi: 10.1109/jproc.2017.2675998. URL [http://dx.doi.org/10.1109/JPROC.](http://dx.doi.org/10.1109/JPROC.2017.2675998)
543 [2017.2675998](http://dx.doi.org/10.1109/JPROC.2017.2675998).
- 544 M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing Textures in the Wild. In
545 *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014. URL
546 <https://doi.org/10.1109/CVPR.2014.461>.
- 547 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
548 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
549 Schulman. Training verifiers to solve math word problems, 2021. URL [https://arxiv.](https://arxiv.org/abs/2110.14168)
550 [org/abs/2110.14168](https://arxiv.org/abs/2110.14168).
- 551 Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex sgd.
552 In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,
553 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/](https://proceedings.neurips.cc/paper_files/paper/2019/file/b8002139cdde66b87638f7f91d169d96-Paper.pdf)
554 [file/b8002139cdde66b87638f7f91d169d96-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/b8002139cdde66b87638f7f91d169d96-Paper.pdf).
- 555 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention:
556 Fast and memory-efficient exact attention with io-awareness. In S. Koyejo, S. Mo-
557 hamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural In-*
558 *formation Processing Systems*, volume 35, pp. 16344–16359. Curran Associates, Inc.,
559 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/](https://proceedings.neurips.cc/paper_files/paper/2022/file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf)
560 [file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf).
- 561 Aaron Defazio and Leon Bottou. On the ineffectiveness of variance reduced optimization for
562 deep learning. In *Advances in Neural Information Processing Systems*, volume 32. Curran
563 Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2019/file/84d2004bf28a2095230e8e14993d398d-Paper.pdf)
564 [paper/2019/file/84d2004bf28a2095230e8e14993d398d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/84d2004bf28a2095230e8e14993d398d-Paper.pdf).
- 565 Michal Dereziński. Stochastic variance-reduced newton: Accelerating finite-sum minimization
566 with large batches. In *OPT 2023: Optimization for Machine Learning, 2023*. URL [https:](https://openreview.net/forum?id=EUshjvvMcj)
567 [//openreview.net/forum?id=EUshjvvMcj](https://openreview.net/forum?id=EUshjvvMcj).
- 568 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
569 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-
570 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recogni-
571 tion at scale. In *International Conference on Learning Representations*, 2021. URL [https:](https://openreview.net/forum?id=YicbFdNTTy)
572 [//openreview.net/forum?id=YicbFdNTTy](https://openreview.net/forum?id=YicbFdNTTy).
- 573 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
574 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
575 *arXiv preprint arXiv:2407.21783*, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 576 Benjamin Dubois-Taine, Sharan Vaswani, Reza Babanezhad, Mark Schmidt, and Simon Lacoste-
577 Julien. Svrg meets adagrad: painless variance reduction. *Machine Learning*, 111(12):4359–4409,
578 2022.
- 579 Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-
580 convex optimization via stochastic path-integrated differential estimator. In S. Bengio,
581 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Ad-*
582 *vances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.,
583 2018. URL [https://proceedings.neurips.cc/paper_files/paper/2018/](https://proceedings.neurips.cc/paper_files/paper/2018/file/1543843a4723ed2ab08e18053ae6dc5b-Paper.pdf)
584 [file/1543843a4723ed2ab08e18053ae6dc5b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/1543843a4723ed2ab08e18053ae6dc5b-Paper.pdf).
- 585 Sachin Garg, Albert S Berahas, and Michał Dereziński. Second-order information promotes
586 mini-batch robustness in variance-reduced gradients. *arXiv:2404.14758*, 2024. URL [https:](https://arxiv.org/abs/2404.14758)
587 [//arxiv.org/abs/2404.14758](https://arxiv.org/abs/2404.14758).
- 588 Sebastian Houben, Johannes Stallkamp, Jan Salmen, Marc Schlipf, and Christian Igel. Detection
589 of traffic signs in real-world images: The German Traffic Sign Detection Benchmark. In *Interna-*
590 *tional Joint Conference on Neural Networks (IJCNN)*, 2013. URL [https://doi.org/10.](https://doi.org/10.1109/IJCNN.2013.6706807)
591 [1109/IJCNN.2013.6706807](https://doi.org/10.1109/IJCNN.2013.6706807).

- 594 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
595 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International*
596 *Conference on Learning Representations (ICLR)*, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=nZeVKeeFYf9)
597 [forum?id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 598 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,
599 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali
600 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/](https://doi.org/10.5281/zenodo.5143773)
601 [zenodo.5143773](https://doi.org/10.5281/zenodo.5143773).
- 602 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using pre-
603 dictive variance reduction. In *Advances in Neural Information Processing Sys-*
604 *tems*, 2013. URL [https://papers.nips.cc/paper_files/paper/2013/hash/](https://papers.nips.cc/paper_files/paper/2013/hash/acldd209cbcc5e5d1c6e28598e8cbbe8-Abstract.html)
605 [acldd209cbcc5e5d1c6e28598e8cbbe8-Abstract.html](https://papers.nips.cc/paper_files/paper/2013/hash/acldd209cbcc5e5d1c6e28598e8cbbe8-Abstract.html).
- 606 Mohammad Emtiyaz Khan. Knowledge adaptation as posterior correction. *arXiv:2506.14262*, 2025.
607 URL <https://www.arxiv.org/abs/2506.14262>.
- 608 Mohammad Emtiyaz Khan and Haavard Rue. The Bayesian learning rule. *Journal of Ma-*
609 *chine Learning Research*, 24(281):1–46, 2023. URL [http://jmlr.org/papers/v24/](http://jmlr.org/papers/v24/22-0291.html)
610 [22-0291.html](http://jmlr.org/papers/v24/22-0291.html).
- 611 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3D object representations for fine-grained
612 categorization. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2013. URL
613 <https://doi.org/10.1109/ICCVW.2013.77>. (Workshops).
- 614 Lihua Lei, Cheng Ju, Jianbo Chen, and Michael I Jordan. Non-convex finite-sum optimization via
615 scsg methods. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,
616 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Cur-
617 ran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2017/file/81ca0262c82e712e50c580c032d99b60-Paper.pdf)
618 [paper/2017/file/81ca0262c82e712e50c580c032d99b60-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/81ca0262c82e712e50c580c032d99b60-Paper.pdf).
- 619 Wu Lin, Mark Schmidt, and Mohammad Emtiyaz Khan. Handling the positive-definite constraint
620 in the Bayesian learning rule. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*
621 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*
622 *Research*, pp. 6116–6126. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v119/lin20d.html)
623 [press/v119/lin20d.html](https://proceedings.mlr.press/v119/lin20d.html).
- 624 Jerry Ma and Denis Yarats. Quasi-hyperbolic momentum and adam for deep learning. In *Internat-*
625 *ional Conference on Learning Representations*, 2019. URL [https://openreview.net/](https://openreview.net/forum?id=S1fUpOR5FQ)
626 [forum?id=S1fUpOR5FQ](https://openreview.net/forum?id=S1fUpOR5FQ).
- 627 Xiangming Meng, Roman Bachmann, and Mohammad Emtiyaz Khan. Training binary neural
628 networks using the Bayesian learning rule. In Hal Daumé III and Aarti Singh (eds.), *Pro-*
629 *ceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proced-*
630 *ings of Machine Learning Research*, pp. 6852–6861. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v119/meng20a.html)
631 [press/v119/meng20a.html](https://proceedings.mlr.press/v119/meng20a.html).
- 632 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mix-
633 ture models. In *International Conference on Learning Representations*, 2017. URL [https://openreview.net/](https://openreview.net/forum?id=Byj72udxe)
634 [forum?id=Byj72udxe](https://openreview.net/forum?id=Byj72udxe).
- 635 Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don’t give me the details, just the sum-
636 mary! topic-aware convolutional neural networks for extreme summarization. In Ellen Riloff,
637 David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of the 2018 Confer-*
638 *ence on Empirical Methods in Natural Language Processing*, pp. 1797–1807, Brussels, Bel-
639 gium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/
640 [D18-1206](https://aclanthology.org/D18-1206). URL <https://aclanthology.org/D18-1206>.
- 641 Lam M. Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for
642 machine learning problems using stochastic recursive gradient. In Doina Precup and Yee Whye
643 Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70
644 of *Proceedings of Machine Learning Research*, pp. 2613–2621. PMLR, 06–11 Aug 2017. URL
645 <https://proceedings.mlr.press/v70/nguyen17b.html>.

- 648 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin
649 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for
650 the finest text data at scale. In *The Thirty-eight Conference on Neural Information Processing
651 Systems Datasets and Benchmarks Track*, 2024. URL [https://openreview.net/forum?
652 id=n6SCkn2QaG](https://openreview.net/forum?id=n6SCkn2QaG).
- 653 Nicolas Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an expo-
654 nential convergence rate for finite training sets. In F. Pereira, C.J. Burges, L. Bottou, and
655 K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 25. Cur-
656 ran Associates, Inc., 2012. URL [https://proceedings.neurips.cc/paper_files/
657 paper/2012/file/905056c1ac1dad141560467e0a99e1cf-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/905056c1ac1dad141560467e0a99e1cf-Paper.pdf).
- 658 Abdurakhmon Sadiev, Aleksandr Beznosikov, Abdulla Jasem Almansoori, Dmitry Kamzolov,
659 Rachael Tappenden, and Martin Takáč. Stochastic gradient methods with preconditioned up-
660 dates. *Journal of Optimization Theory and Applications*, 201(2):471–489, 2024. URL [https:
661 //doi.org/10.1007/s10957-023-02365-3](https://doi.org/10.1007/s10957-023-02365-3).
- 662 Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic
663 average gradient. *Mathematical Programming*, 162(1):83–112, 2017. URL [https://doi.
664 org/10.1007/s10107-016-1030-6](https://doi.org/10.1007/s10107-016-1030-6).
- 665 Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized
666 loss. *The Journal of Machine Learning Research*, 14(1):567–599, 2013. URL [https://www.
667 jmlr.org/papers/v14/shalev-shwartz13a.html](https://www.jmlr.org/papers/v14/shalev-shwartz13a.html).
- 668 Yuesong Shen, Nico Daheim, Bai Cong, Peter Nickl, Gian Maria Marconi, Bazan Clement
669 Emile Marcel Raoul, Rio Yokota, Iryna Gurevych, Daniel Cremers, Mohammad Emtiyaz Khan,
670 and Thomas Möllenhoff. Variational learning is effective for large deep networks. In *Forty-first
671 International Conference on Machine Learning*, 2024. URL [https://openreview.net/
672 forum?id=cXBv07GKvk](https://openreview.net/forum?id=cXBv07GKvk).
- 673 Jingruo Sun, Zachary Frangella, and Madeleine Udell. SAPPHERE: preconditioned stochastic
674 variance reduction for faster large-scale statistical learning. *arXiv:2501.15941*, 2025. URL
675 <https://arxiv.org/abs/2501.15941>.
- 676 Lionel Tondji, Sergii Kashubin, and Moustapha Cisse. Variance reduction in deep learning: More
677 momentum is all you need. *arXiv preprint arXiv:2111.11828*, 2021. URL [https://arxiv.
678 org/abs/2111.11828](https://arxiv.org/abs/2111.11828).
- 679 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
680 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In
681 I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,
682 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30,
683 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
684 file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 685 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
686 Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
687 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger,
688 Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural
689 language processing. In Qun Liu and David Schlangen (eds.), *Proceedings of the 2020 Confer-
690 ence on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–
691 45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.
692 emnlp-demos.6. URL <https://aclanthology.org/2020.emnlp-demos.6>.
- 693 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
694 Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
695 Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang,
696 Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang,
697 Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,
698 Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL
699 <https://arxiv.org/abs/2412.15115>.

702 Yida Yin, Zhiqiu Xu, Zhiyuan Li, Trevor Darrell, and Zhuang Liu. A coefficient makes SVRG
703 effective. In *The Thirteenth International Conference on Learning Representations, 2025*. URL
704 <https://openreview.net/forum?id=twtTLZnG0B>.
705
706 Cheng Zhang, Judith Bütepage, Hedvig Kjellström, and Stephan Mandt. Advances in variational
707 inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):2008–2026,
708 2018. URL <https://doi.org/10.1109/TPAMI.2018.2889774>.
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

A DERIVATION OF THE NEWTON-LIKE SVRG EXTENSION

We start by writing q in an exponential-family form which in this case is convenient to write in terms of the precision $\mathbf{S} = \Sigma^{-1}$. This is shown below with sufficient statistics highlighted in red,

$$\mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{S}^{-1}) \propto \exp \left[\mathbf{m}^\top \mathbf{S} \boldsymbol{\theta} + \text{Tr} \left(\left(-\frac{1}{2}\mathbf{S}\right) \boldsymbol{\theta} \boldsymbol{\theta}^\top \right) \right]. \quad (19)$$

This is a log-linear form for a sufficient-statistics vector of $\boldsymbol{\theta}$ and $\boldsymbol{\theta} \boldsymbol{\theta}^\top$. We denote the natural parameter $\boldsymbol{\lambda} = (\mathbf{S} \mathbf{m}, -\frac{1}{2} \mathbf{S})$ consisting of two elements: a vector and a square matrix. The natural gradient can be written in terms of the gradient and Hessian of ℓ_i (Khan & Rue, 2023, Eq. 10-11),

$$\tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}) = \mathbb{E}_q \left[(\nabla \ell_i - \nabla^2 \ell_i \mathbf{m}, \frac{1}{2} \nabla^2 \ell_i) \right]. \quad (20)$$

This also has two elements, and uses gradient and Hessian computed at samples from q .

We will denote the $\boldsymbol{\lambda}_{\text{in}} = (\mathbf{S}_{\text{in}} \mathbf{m}_{\text{in}}, -\frac{1}{2} \mathbf{S}_{\text{in}})$ and $\boldsymbol{\lambda}_{\text{out}} = (\mathbf{S}_{\text{out}} \mathbf{m}_{\text{out}}, -\frac{1}{2} \mathbf{S}_{\text{out}})$. We first write the update for the second entry of $\boldsymbol{\lambda}_{\text{in}}$ which is $-\frac{1}{2} \mathbf{S}_{\text{in}}$,

$$\mathbf{S}_{\text{new}} \leftarrow (1 - \eta) \mathbf{S}_{\text{in}} + \eta N \left[\mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla^2 \ell_i] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_i] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_j] \right]. \quad (21)$$

We denoted the new value as \mathbf{S}_{new} to differentiate it with the old value. This will be useful to simplify the update for the first entry $\mathbf{S}_{\text{in}} \mathbf{m}_{\text{in}}$ which we show below,

$$\begin{aligned} \mathbf{S}_{\text{new}} \mathbf{m}_{\text{new}} &= (1 - \eta) \mathbf{S}_{\text{in}} \mathbf{m}_{\text{in}} - \eta N \left[\mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla \ell_i - \nabla^2 \ell_i \mathbf{m}_{\text{in}}] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_i - \nabla^2 \ell_i \mathbf{m}_{\text{out}}] \right. \\ &\quad \left. + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_j - \nabla^2 \ell_j \mathbf{m}_{\text{out}}] \right] \\ &= \left\{ \mathbf{S}_{\text{new}} - \eta N \left[\mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla^2 \ell_i] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_i] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_j] \right] \right\} \mathbf{m}_{\text{in}} + \\ &\quad - \eta N \left[\mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla \ell_i - \nabla^2 \ell_i \mathbf{m}_{\text{in}}] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_i - \nabla^2 \ell_i \mathbf{m}_{\text{out}}] \right. \\ &\quad \left. + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_j - \nabla^2 \ell_j \mathbf{m}_{\text{out}}] \right] \\ &= \mathbf{S}_{\text{new}} \mathbf{m}_{\text{in}} - \eta N \left[\cancel{\mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla^2 \ell_i] \mathbf{m}_{\text{in}}} - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_i] \mathbf{m}_{\text{in}} + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_j] \mathbf{m}_{\text{in}} \right. \\ &\quad \left. + \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla \ell_i - \nabla^2 \ell_i \mathbf{m}_{\text{in}}] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_i - \nabla^2 \ell_i \mathbf{m}_{\text{out}}] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_j - \nabla^2 \ell_j \mathbf{m}_{\text{out}}] \right] \\ &= \mathbf{S}_{\text{new}} \mathbf{m}_{\text{in}} - \eta N \left[\mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla \ell_i] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_i] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_j] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_i] \mathbf{m}_{\text{in}} \right. \\ &\quad \left. + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_j] \mathbf{m}_{\text{in}} + \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_i \mathbf{m}_{\text{out}}] - \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_j \mathbf{m}_{\text{out}}] \right] \\ &= \mathbf{S}_{\text{new}} \mathbf{m}_{\text{in}} - \eta N \left[\mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{in}}}} [\nabla \ell_i] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_i] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla \ell_j] \right. \\ &\quad \left. \left(\frac{1}{N} \sum_{j=1}^N \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_j] - \mathbb{E}_{q_{\boldsymbol{\lambda}_{\text{out}}}} [\nabla^2 \ell_i] \right) (\mathbf{m}_{\text{in}} - \mathbf{m}_{\text{out}}) \right] \end{aligned} \quad (22)$$

Algorithm 5 VON-PC: Variational Online Newton with Posterior Correction

Initialize: Number of inner steps m , learning rates α and β

- 1: Initialize $\mathbf{m}_{\text{in}}, \mathbf{S}_{\text{in}}$
- 2: **while** not converged **do**
- 3: $\mathbf{g}_{\text{out}} \leftarrow \sum_{i=1}^N \nabla \ell_i(\boldsymbol{\theta}_{\text{in}})$ where $\boldsymbol{\theta}_{\text{in}} = \mathbf{m}_{\text{in}} + \mathbf{S}_{\text{in}}^{-\frac{1}{2}} \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$
- 4: $\mathbf{H}_{\text{out}} \leftarrow \sum_{i=1}^N \nabla^2 \ell_i(\boldsymbol{\theta}_{\text{in}})$
- 5: $\boldsymbol{\theta}_{\text{out}} \leftarrow \boldsymbol{\theta}_{\text{in}}, \mathbf{m}_{\text{out}} \leftarrow \mathbf{m}_{\text{in}}$
- 6: **for** $t = 1, 2, \dots, m$ **do**
- 7: Randomly pick $i \in \{1, 2, \dots, N\}$
- 8: $\mathbf{g}_{\text{in}} \leftarrow \nabla \ell_i(\boldsymbol{\theta}_{\text{in}}) - \nabla \ell_i(\boldsymbol{\theta}_{\text{out}}) + \frac{1}{N} \mathbf{g}_{\text{out}}$ where $\boldsymbol{\theta}_{\text{in}} = \mathbf{m}_{\text{in}} + \mathbf{S}_{\text{in}}^{-\frac{1}{2}} \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$
- 9: $\mathbf{H}_{\text{out} \setminus i} \leftarrow \frac{1}{N} \mathbf{H}_{\text{out}} - \nabla^2 \ell_i(\boldsymbol{\theta}_{\text{out}})$
- 10: $\mathbf{H}_{\text{in}} \leftarrow \nabla^2 \ell_i(\boldsymbol{\theta}_{\text{in}}) + \mathbf{H}_{\text{out} \setminus i}$
- 11: $\mathbf{S}_{\text{in}} \leftarrow (1 - \beta) \mathbf{S}_{\text{in}} + \beta N \mathbf{H}_{\text{in}}$
- 12: $\mathbf{m}_{\text{in}} \leftarrow \mathbf{m}_{\text{in}} - \alpha N \mathbf{S}_{\text{in}}^{-1} [\mathbf{g}_{\text{in}} + \mathbf{H}_{\text{out} \setminus i} (\mathbf{m}_{\text{in}} - \mathbf{m}_{\text{out}})]$
- 13: **end for**
- 14: **end while**

An algorithm can be conveniently written by defining the following outer-loop quantities using the output $\boldsymbol{\lambda}_{\text{in}}$ of the inner loop,

$$\mathbf{g}_{\text{out}} \leftarrow \sum_{j=1}^N \mathbb{E}_{q_{\text{in}}} [\nabla \ell_j], \quad \mathbf{H}_{\text{out}} \leftarrow \sum_{j=1}^N \mathbb{E}_{q_{\text{in}}} [\nabla^2 \ell_j]. \quad (23)$$

We then set $\mathbf{m}_{\text{out}} \leftarrow \mathbf{m}_{\text{in}}$ and $\mathbf{S}_{\text{out}} \leftarrow \mathbf{S}_{\text{in}}$, and the corresponding natural parameter to be $\boldsymbol{\lambda}_{\text{out}}$. Using these, we can write the updates in the inner loop as follows (strictly in this order),

$$\begin{aligned} \mathbf{g}_{\text{in}} &\leftarrow \mathbb{E}_{q_{\text{in}}} [\nabla \ell_i] - \mathbb{E}_{q_{\text{out}}} [\nabla \ell_i] + \mathbf{g}_{\text{out}}/N \\ \mathbf{H}_{\text{out} \setminus i} &\leftarrow \mathbf{H}_{\text{out}}/N - \mathbb{E}_{q_{\text{out}}} [\nabla^2 \ell_i] \\ \mathbf{H}_{\text{in}} &\leftarrow \mathbb{E}_{q_{\text{in}}} [\nabla^2 \ell_i] + \mathbf{H}_{\text{out} \setminus i} \\ \mathbf{S}_{\text{in}} &\leftarrow (1 - \eta) \mathbf{S}_{\text{in}} + \eta N \mathbf{H}_{\text{in}} \\ \mathbf{m}_{\text{in}} &\leftarrow \mathbf{m}_{\text{in}} - \eta N \mathbf{S}_{\text{in}}^{-1} [\mathbf{g}_{\text{in}} + \mathbf{H}_{\text{out} \setminus i} (\mathbf{m}_{\text{in}} - \mathbf{m}_{\text{out}})] \end{aligned} \quad (24)$$

These steps are implemented in Alg. 5 by using one Monte-Carlo sample to evaluate the expectations. We highlight in red the new parts added on top of SVRG. We note two useful points regarding the implementation: first, \mathbf{S}_{in} need to be always updated before \mathbf{m}_{in} , and second, variance is further reduced if different example use different seeds (which is not explicitly written in the algorithm).

B DERIVATION OF THE ADAM-LIKE SVRG EXTENSION

To derive the SVRG extension of IVON, we will first write the VB objective in the form used by IVON; see Shen et al. (2024, Eq. 1). Essentially, they use mini-batches \mathcal{B} of size B , and to accommodate this they scale the expected loss by a constant $\kappa \mathbb{E}_q[\ell_i]$. Setting $\kappa = N$ gives back the ERM loss but it can also be set to other values. In IVON, we also treat the regularizer ℓ_0 explicitly by setting it to the weight decay. It is not merged in the losses ℓ_i as in the previous sections. With these changes, Sec. 3.2 can be written as where an explicit natural gradient of \mathcal{L}_0 is added,

$$\boldsymbol{\lambda}_{\text{in}} \leftarrow (1 - \eta) \boldsymbol{\lambda}_{\text{in}} - \eta \kappa \left[\frac{1}{B} \sum_{i \in \mathcal{B}} \left(\tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{\text{in}}) - \tilde{\nabla} \mathcal{L}_i(\boldsymbol{\lambda}_{\text{out}}) \right) + \frac{1}{\kappa} \sum_{j=1}^N \tilde{\nabla} \mathcal{L}_j(\boldsymbol{\lambda}_{\text{out}}) + \frac{1}{\kappa} \tilde{\nabla} \mathcal{L}_0(\boldsymbol{\lambda}_{\text{in}}) \right].$$

We then plug in the natural parameter and natural gradients of $q = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}, \text{diag}(\mathbf{s})^{-1})$, which yields a Newton-like update very similar to Thm. 3.

To derive the IVON-PC update, we assume a quadratic regularizer $\ell_0 = s_0 \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta}$ with $s_0 > 0$. The VONcorr update can be written as follows where the new parts compared to Eq. 15 are in red,

$$\begin{aligned} \mathbf{s}_{\text{in}} &\leftarrow (1 - \eta) \mathbf{s}_{\text{in}} + \eta \kappa \left[\frac{1}{B} \sum_{i \in \mathcal{B}} \left(\mathbb{E}_{q_{\lambda_{\text{in}}}} [\nabla^2 \ell_i] - \mathbb{E}_{q_{\lambda_{\text{out}}}} [\nabla^2 \ell_i] \right) + \frac{1}{\kappa} \sum_{j=1}^N \mathbb{E}_{q_{\lambda_{\text{out}}}} [\nabla^2 \ell_j] + \frac{s_0}{\kappa} \right], \\ \mathbf{m}_{\text{in}} &\leftarrow \mathbf{m}_{\text{in}} - \eta \kappa \frac{1}{\mathbf{s}_{\text{in}}} \left[\frac{1}{B} \sum_{i \in \mathcal{B}} \left(\mathbb{E}_{q_{\lambda_{\text{in}}}} [\nabla \ell_i] - \mathbb{E}_{q_{\lambda_{\text{out}}}} [\nabla \ell_i] \right) + \frac{1}{\kappa} \sum_{j=1}^N \mathbb{E}_{q_{\lambda_{\text{out}}}} [\nabla \ell_j] + \frac{s_0}{\kappa} \mathbf{m}_{\text{in}} \right. \\ &\quad \left. + \left(\frac{1}{\kappa} \sum_{j=1}^N \mathbb{E}_{q_{\lambda_{\text{out}}}} [\nabla^2 \ell_j] - \frac{1}{B} \sum_{i \in \mathcal{B}} \mathbb{E}_{q_{\lambda_{\text{out}}}} [\nabla^2 \ell_i] \right) (\mathbf{m}_{\text{in}} - \mathbf{m}_{\text{out}}) \right]. \end{aligned} \quad (25)$$

To write the update in IVON form, we make a few modifications.

1. For weight decay, we tune $\delta = s_0/\kappa$ directly.
2. We remove δ from the \mathbf{s}_{in} update and divide the whole update by κ . The resulting update is written in terms of \mathbf{h}_{in} such that $\mathbf{s}_{\text{in}} = \kappa(\mathbf{h}_{\text{in}} + \delta)$ and $\sigma_{\text{in}}^2 = 1/(\kappa(\mathbf{h}_{\text{in}} + \delta))$.
3. We use different learning rate for \mathbf{m}_{in} and \mathbf{h}_{in} updates. For \mathbf{m}_{in} , we use a scheduled η_t for iteration t . For \mathbf{h}_{in} , we use $\beta_2 \in [0, 1)$.
4. We use momentum with learning rate $\beta_1 \in [0, 1)$ and debiasing for \mathbf{g} (but not for \mathbf{h}).
5. We add a term for the update of \mathbf{h} which ensures positivity of \mathbf{h} . We initialize \mathbf{h} by a scalar constant $h_0 > 0$.
6. The scaling factor κ is often set to N but it can be different for cases when the effective number of examples is not immediately clear (for example, for LLM training).
7. We add an additional factor α in front of the outer gradients and Hessian, similarly to α -SVRG (Yin et al., 2025).

B.1 HANDLING MEGA-BATCHES FOR FULL-GAUSSIAN

For full-Gaussians, we can build the site as follows,

$$\begin{aligned} \hat{q}_{\text{out}} &\leftarrow \exp \left(-\boldsymbol{\theta}^\top \mathbf{g}_{\text{out}} - \frac{1}{2} (\boldsymbol{\theta} - \mathbf{m}_{\text{out}})^\top \mathbf{H}_{\text{out}} (\boldsymbol{\theta} - \mathbf{m}_{\text{out}}) \right), \\ &\text{where } \mathbf{g}_{\text{out}} \leftarrow \rho_1 \mathbf{g}_{\text{out}} + (1 - \rho_1) \sum_{i \in \mathcal{M}} \mathbb{E}_{q_{\text{out}}} [\nabla \ell_i] \\ &\quad \mathbf{H}_{\text{out}} \leftarrow \rho_2 \mathbf{H}_{\text{out}} + (1 - \rho_2) \sum_{i \in \mathcal{M}} \mathbb{E}_{q_{\text{out}}} [\nabla^2 \ell_i]. \end{aligned} \quad (26)$$

C ABLATIONS AND ADDITIONAL RESULTS ON LOGISTIC REGRESSION

Here, we run several ablations on the IVON-PC algorithm proposed in the main paper (Alg. 4).

Extra term in m-update. The leftmost plot in Fig. 7 shows how the new term in the mean-update can cause instabilities in the training when naively implemented. This is likely due to the noisy reparametrization-trick based Hessian estimate. Using exact diagonal Hessian fixes the problem and leads to more stable training. Since the diagonal Hessian is expensive to compute, one can also simply downweigh the extra term by 0.01. This works well throughout our experiments. Therefore, we follow this approach in all logistic regression and image classification experiments. Aggressive gradient clipping ($\xi = 10^{-4}$) also somewhat stabilizes the training, but does not work as well as downweighting on this example. For more frequent refreshing as used in the transformer experiments, gradient clipping was enough to stabilize the training and no downweighting was required.

Outer megabatch momentum. We also perform two ablations over the outer momentum parameter ρ . The middle plot in Fig. 7 shows that when a few large outer batches are used, it is best to use $\rho = 0$ (no momentum) which always uses the freshest outer batch. In contrast, when using small megabatches as in the right plot in Fig. 7, larger values of $\rho = 0.9$ are advantageous.

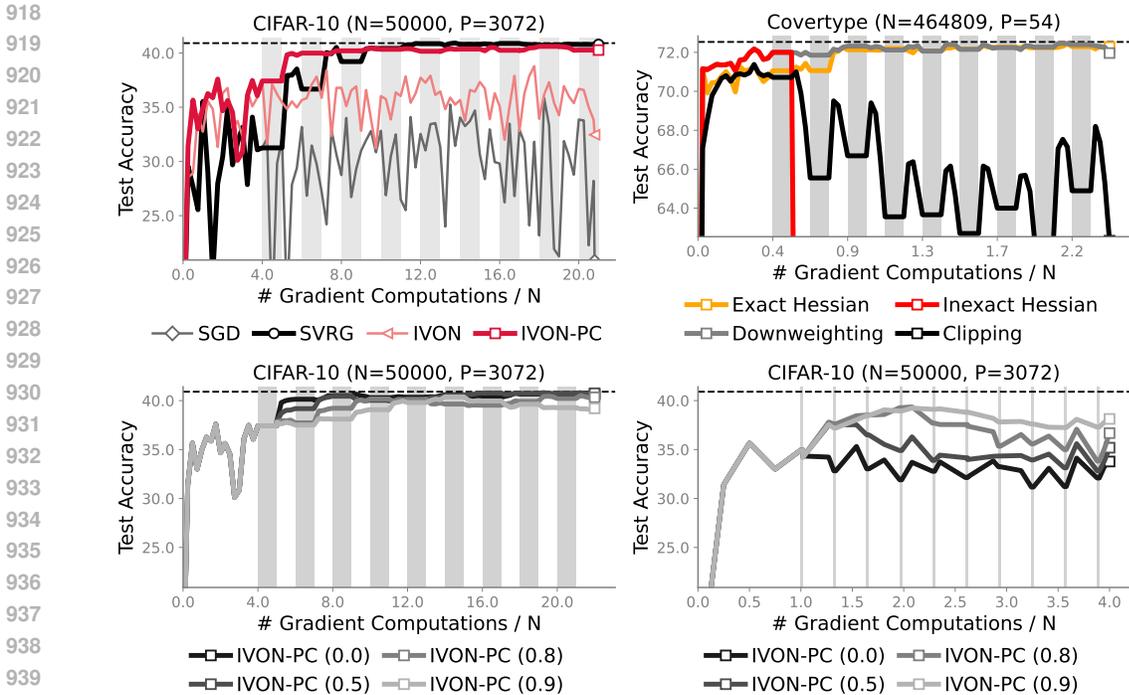


Figure 7: First plot shows additional result on CIFAR-10. For the second plot, we ablate the new red term in the m -update of Alg. 4 to show that it can be numerically unstable. Using exact Hessians or downweighting it improves the performance. Third plot shows that when large outer batches with infrequent refreshes are used, $\rho = 0$ performs best (no momentum). For small outer batches, small momentum works best ($\rho = 0.9$).

D HYPERPARAMETERS AND ADDITIONAL DETAILS

D.1 IMAGE CLASSIFICATION WITH RESNETS

The ImageNet experiments run for 90 epochs, using batch size 384 on 8 GPUs which took around 12-15 hours. The learning rate is annealed to zero using a cosine schedule. For the CIFAR-10 experiments, we train for 100 epochs and batch size 50 on a single GPU, and each run took around 1-2 hours. As described in the main text, the learning rate is annealed to a quarter of the starting learning rate for all methods.

The hyperparameters for SGD, AdamW and IVON are set identical to the ones used in (Shen et al., 2024), except for CIFAR-10 where we use an ess of $5 \cdot 10^5$. All SVRG methods and IVON-PC inherit the hyperparameters from their base algorithm. The variance reduced methods use a megabatch size 50 times larger than the minibatch size on ImageNet and 10 times larger than the minibatch size on CIFAR-10. We tuned α for α -SVRG and PC-IVON separately. The optimal α is 0.2 for AdamW, 0.4 for SGD and 0.1 for IVON-PC. On ImageNet, we used $\alpha = 0.6$ but it was not overly tuned due to computational constraints. IVON-PC uses outer momentum of $\rho_1 = \rho_2 = 0.3$ on CIFAR-10 and no momentum was tried on ImageNet.

D.2 PRETRAINING FROM SCRATCH

We first detail the experiments described in Sec. 4.3. We use the same set-up as in (Shen et al., 2024) which follows the nanoGPT repository found under <https://github.com/karpathy/nanoGPT>. We use an effective batch size of 480 achieved via 48 gradient accumulation steps to train a 125M parameter GPT-2 model from scratch on ca. 50B tokens from OpenWebtext in 100,000 steps. The AdamW run uses a learning rate of $6 \cdot 10^{-4}$, $(\beta_1, \beta_2) = (0.9, 0.95)$. The IVON run uses a learning rate of 0.3, $(\beta_1, \beta_2) = 0.9, 0.99999$, an effective sample size $\kappa = 1 \cdot 10^{10}$,

a Hessian initialization of 0.001 as well as element-wise clipping of 0.001. For IVON-PCM we use the same hyperparameters but starting at step 50,000 we add a correction with a megabatch that is 10 times the size of a minibatch and megabatch statistics that get updated after every 10 inner loop steps. We also use $\rho_1 = 0.6$ and $\rho_2 = 0.1$ for momentum. All these experiments are run on 8xA100 GPUs for up to one and a half days using bf16 and flash attention (Dao et al., 2022).

For the ablations we follow a similar recipe but rather use a smaller GPT-2 model which we downsize to ca. 33M parameters by using only 4 layers and 4 heads with an embedding dimension of 512. We run the experiments on wikitext103 and use the official train-validation splits for training and evaluation. Hyperparameters for IVON, IVON-PC, and AdamW are kept as above. We warmup IVON-PC and α -SVRG with IVON and AdamW for 5,000 steps for the ablation over α and for 1,000 steps for the ablation over the refresh rate of the outer estimates. We use a batch size of 64 and a short context length of 128.

D.3 CONTINUAL PRETRAINING

In this experiment we continually pretrain the GPT2-125M model from Shen et al. (2024) which is publicly available (<https://huggingface.co/team-approx-bayes/gpt2-small>) which was trained for 50B tokens on OpenWebText. All methods use a context length of 512 and a batch size of 80 and are trained on a single NVIDIA A100 80GB GPU with bf16 and flash attention to speed up training and reduce GPU memory utilization. We use the first 1B tokens from the Fineweb-edu-sample-10BT which is available openly under the following <https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu> and split of the final 2000 documents for a validation split. For all experiments the learning rates are annealed to zero.

For IVON and IVON-PC, we continue using the optimizer state from pretraining but change some of the hyperparameters. We change the ess to $3 \cdot 10^{10}$ and use $\beta_2 = 0.9999$. The learning rate is set to 0.04 instead of 0.3, which was used for pretraining. With larger learning rates we found IVON to become less stable. For AdamW we start the optimizer state freshly and use a learning rate of $1 \cdot 10^{-4}$, down from the $6 \cdot 10^{-4}$ which was used for AdamW-trained models in Shen et al. (2024). Above, we found that loss sharply increased early in training which could lead to more forgetting of the data it was trained on. We use $\beta_1 = 0.9$ and $\beta_2 = 0.999$ which are oft-used for finetuning and the default choice in huggingface transformers (Wolf et al., 2020), which we use to implement our experiments.

For α -SVRG and IVON-PC we warmstart training with 1,000 steps of AdamW and IVON, respectively, and use 40 inner steps before refreshing the outer gradient and Hessian estimates using 40 randomly sampled batches of size 80, i.e., 3200 examples and up to 1,638,400 tokens in total for estimating the gradients. We sample these batches randomly, so they need not overlap with the batches used in the following inner loop. We have found this to perform better in small experiments. Potentially, randomly sampling the batches reduces bias towards the same data used in the inner loop. For IVON-PC we set $\rho_1 = 0.3$ and $\rho_2 = 0.05$.

D.4 FINETUNING

We finetune various models following the Transformer architecture (Vaswani et al., 2017). First, we use Vision Transformers (Dosovitskiy et al., 2021) with ca. 88M parameters. Our experiment uses OpenCLIP (Ilharco et al., 2021) and we only train the vision encoder but not the text encoder which produces label embeddings to which the image embeddings are matched. We train for 5 epochs on Cars (Krause et al., 2013), DTD (Cimpoi et al., 2014), GTSRB (Houben et al., 2013) and RESISC45 (Cheng et al., 2017) and start correction for IVON-PC after just 50 steps. We use a batch size of 8, 32 warmup steps, $\beta_1 = 0.9$, $\beta_2 = 0.99999$, the Hessian is initialized to 0.1, ess equals $1 \cdot 10^{10}$, and the learning rate is initialized to 0.3 and annealed to zero.

Next, we finetune two LLMs. First, we finetune the full Qwen2.5-0.5B-Instruct model on the first 50% of the training set of XSUM (Narayan et al., 2018) and evaluate on the corresponding test split. We finetune for a single epoch with a learning rate of 0.01, $\beta_1 = 0.9$, $\beta_2 = 0.99999$, ess of $1 \cdot 10^{10}$ and a Hessian initialized to 0.001 as well as element-wise clipping to 0.001. We use $\alpha = 0.7$ and refresh the outer gradients and Hessians every 50 steps. We use the same hyperparameters for LoRA-finetuning of LLAMA-3.1-8B but increase the learning rate to 0.05. We train the model for 3

1026 epochs on GSM8k (Cobbe et al., 2021) and calculate the loss on both input and output tokens with
1027 a standard cross-entropy criterion. For all LLM methods we use greedy decoding and zero-shot
1028 prompting.
1029

1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079