
Preconditioners for the Stochastic Training of Implicit Neural Representations

Shin-Fang Chng*
Australian Institute of Machine Learning
University of Adelaide

Hemanth Saratchandran*
Australian Institute of Machine Learning
University of Adelaide

Simon Lucey
Australian Institute of Machine Learning
University of Adelaide

Abstract

Implicit neural representations have emerged as a powerful technique for encoding complex continuous multidimensional signals as neural networks, enabling a wide range of applications in computer vision, robotics, and geometry. While Adam is commonly used for training due to its stochastic proficiency, it entails lengthy training durations. To address this, we explore alternative optimization techniques for accelerated training without sacrificing accuracy. Traditional second-order optimizers like L-BFGS are suboptimal in stochastic settings, making them unsuitable for large-scale data sets. Instead, we propose stochastic training using curvature-aware diagonal preconditioners, showcasing their effectiveness across various signal modalities such as images, shape reconstruction, and Neural Radiance Fields (NeRF).

1 Introduction

Implicit neural representations (INRs) have gained significant traction in the fields of machine learning and computer graphics [10, 12, 28, 35, 55, 56, 57, 60]. Unlike conventional grid-based representations that rely on explicit geometric primitives or parameterized functions [8, 19, 59, 68], INRs take a distinct approach. They define shapes and objects implicitly through trained neural networks. These representations have garnered attention for their ability to efficiently encode complex and highly detailed continuous signals, making them suitable for tasks such as image synthesis [12, 57], shape modeling [1, 11, 35, 38, 56] and robotics [13, 62]. INRs also offer desirable properties, including continuous and differentiable outputs, allowing seamless integration with gradient-based optimizers.

Classically INRs were trained with a ReLU activation, leading to low fidelity reconstructions due to such networks admitting spectral bias [43]. To overcome spectral bias, positional embedding layers (PE) [61] are often added but this increases the parameter count of the network and can lead to noisy gradients. A more recent approach is to use non-traditional activations such as sine [55], Gaussian [44], and wavelets [53]. The use of such activations enable high-frequency encoding without positional embedding layers. The major benefit of these activations over positional embedding layers is their well-behaved gradients [44, 55].

While INRs have opened up new possibilities for realistic and efficient signal representation and manipulation, the training of these architectures still presents a significant research challenge. Curiously, the prevailing approach to training such networks has become almost axiomatic - practitioners

* Equal contribution.

Correspondence to: Shin-Fang Chng <shinfang.chng@adelaide.edu.au>, Hemanth Saratchandran <hemanth.saratchandran@adelaide.edu.au>

typically rely on the Adam optimizer. In the realm of architectures that rely on extensive datasets, the adoption of a stochastic training scheme becomes essential. Surprisingly, the preference for using Adam as the optimizer over a more traditional vanilla optimizer like SGD remains unexplained, despite its widespread usage. In [54], the authors demonstrated remarkable acceleration in training INRs by utilizing second-order optimizers such as L-BFGS [37]. However, an obstacle encountered by [54] was the limited performance of second-order optimizers, such as L-BFGS, in stochastic settings, rendering their methods inefficient for stochastic training of INRs.

In this paper, we investigate the stochastic training of INRs from the perspective of preconditioning. Preconditioners are essential components in optimization and iterative solvers, aiming to enhance convergence and efficiency in numerical methods [5, 37, 50]. They specifically address the condition number, which measures the sensitivity of the solution to small input perturbations. By applying suitable transformations, preconditioners effectively reduce the condition number, facilitating rapid convergence in iterative solvers [16].

We delve into the relationship between the Adam optimizer and preconditioned gradient descent algorithms. Specifically, we demonstrate that Adam can be viewed as employing a diagonal preconditioner based on the Gauss-Newton matrix. This sheds light on why Adam performs remarkably well for INRs compared to vanilla SGD. Furthermore, we propose the use of curvature-aware preconditioners for training INRs, especially for architectures with non-traditional activations such as sine [55], Gaussian [44], or wavelets [53]. We show that employing curvature-aware preconditioners can significantly accelerate training times in stochastic settings.

Our analysis focuses on specific algorithms that leverage these preconditioners, including equilibrated stochastic gradient descent (ESGD) [16], AdaHessian [66], Shampoo [24], and K-FAC [33]. By exploring the potential of these preconditioners, we gain valuable insights into improving the training process for INRs in stochastic settings.

Our main contributions are:

1. An analysis of preconditioning for training INRs unveiling the theoretical underpinnings and illuminating the precise scenarios where preconditioning becomes useful.
2. An empirical validation of our theory by demonstrating the superiority of curvature-aware preconditioning algorithms, such as ESGD, over Adam across diverse signal modalities, including images, shape reconstruction, and neural radiance fields.

2 Related Work

Implicit neural representations (INRs) encode signals as weights within a neural network using low-dimensional coordinates, aiming to preserve smoothness in the outputs [68]. The seminal work of Mildenhall et al. [35] demonstrated the power of INRs by showing their remarkable generalization abilities for novel view synthesis. Since this work, there have been several works on INRs finding applications in image view synthesis, shape reconstruction and robotics [9, 14, 15, 17, 20, 36, 39, 40, 41, 45, 51, 52, 55, 63, 67, 69]. Conventionally these works employ a positional embedding layer to their networks, in order to overcome the spectral bias of a ReLU activation [43]. More recently, several new activations have been employed for INRs, such as sine [55], Gaussian [44], wavelets [53]. These new non-traditional activations are showing superior performance on several vision tasks [15, 32, 44, 53, 55, 65].

Preconditioners were first introduced in the context of numerical linear algebra [22, 50], so as to transform ill-conditioned linear systems into better-conditioned forms, enabling faster convergence of iterative methods. Since then they have found their way into machine learning, where they have been widely employed to enhance the efficiency and effectiveness of optimization algorithms. The use of preconditioning techniques in machine learning can be traced back to early works on optimization, such as the conjugate gradient method and its variants [37]. Over time, researchers have explored various preconditioning strategies tailored to the specific characteristics of machine learning problems, including sparse matrices, non-convex objectives, and large-scale datasets [4, 27, 42, 46]. Several well known optimizers such as ADMM [30], mirror descent [2], quasi-Newton methods [7, 21, 37], adaptive optimizers [16, 18, 24, 33, 66] all employ some form of preconditioning.

3 Analysis

3.1 Preconditioners for Stochastic Training

Preconditioning enhances optimization convergence and efficiency by regulating curvature, particularly in stochastic gradient descent (SGD) [16, 42]. High curvature directions in SGD require a small learning rate to avoid overshooting, leading to slow progress in low curvature directions.

For general objective functions f on parameters $\mathbf{x} \in \mathbb{R}^n$. Preconditioning is done by employing a linear change of parameters $\tilde{\mathbf{x}} = D^{\frac{1}{2}} \mathbf{x}$, where in practise $D^{\frac{1}{2}}$ is a non-singular matrix. With this change of parameters a new objective function \tilde{f} , defined on $\tilde{\mathbf{x}}$, can be defined as

$$\tilde{f}(\tilde{\mathbf{x}}) = f(\mathbf{x}). \quad (1)$$

Denoting the gradient of the original objective function by ∇f and its Hessian by $H(f) = \nabla^2 f(\mathbf{x})$. The chain rule gives

$$\nabla \tilde{f}(\tilde{\mathbf{x}}) = D^{-\frac{1}{2}} \nabla f \quad (2)$$

$$H(\tilde{f})(\tilde{\mathbf{x}}) = (D^{-\frac{1}{2}})^T (H(f)(\mathbf{x})) (D^{-\frac{1}{2}}). \quad (3)$$

Equation (2) exhibits the key idea of training with a preconditioner. If we consider gradient descent on the transformed parameter space, a gradient descent iteration takes the form

$$\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} - \eta \nabla_{\tilde{\mathbf{x}}_{t-1}} \tilde{f} \quad (4)$$

Applying equations (1) and (2) this corresponds to doing the update

$$\mathbf{x}_t = \mathbf{x}_{t-1} - \eta D^{-\frac{1}{2}} \nabla_{\mathbf{x}_{t-1}} f. \quad (5)$$

The main endeavour of preconditioning for optimization is to find a preconditioning matrix D such that the new Hessian, given by equation (3), has equal curvature in all directions. This will allow the update (5) to move faster along small curvature directions, leading to faster convergence. Algorithm 1 gives the basic pseudocode for preconditioned SGD.

Algorithm 1 Preconditioned SGD

Require: initial point x_0 , iterations N , learning rate η , preconditioner D .

```

for  $t = 0, \dots, N$  do
   $g_t \leftarrow \nabla_{x_t} f$ 
   $D_t \leftarrow D(x_t)$ 
   $x_{t+1} \leftarrow x_t - \eta D(x_t)^{-\frac{1}{2}} g_t$ 
end for

```

From equation (3) we see that the best preconditioner to use would be the absolute value of the Hessian itself. In the case of a positive definite Hessian, this is precisely what Newton's method does [37]. The main issue with this preconditioner is that the cost of storing and inverting the preconditioner will be $\mathcal{O}(n^3)$ for an objective function of n parameters. This will be extremely costly for the training of overparameterized neural networks.

We list some of the main preconditioners that are being used in the literature.

Diagonal Preconditioners: are matrices that represent some form of a diagonal approximation of the Hessian. Examples include the Jacobi preconditioner, which is given by $D^J = |\text{diag}(H)|$, where $|\cdot|$ is element-wise absolute value and H is the Hessian. In practical settings, the absolute value is taken so that the preconditioner is positive semidefinite. The equilibration preconditioner is a diagonal matrix which has the 2-norm of each row (or column) on the diagonal: $D^E = \|H_{i,\cdot}\|$. Diagonal preconditioners are often used because of their ability to be stored and inverted in linear time. Diagonal preconditioners are also known as adaptive learning rates in the literature.

The algorithm Equilibrated Stochastic Gradient Descent (ESGD) [16] employs D^E as a preconditioner to condition stochastic gradient descent. In [16], it is shown that applying ESGD leads to

Algorithm 2 ESGD

Require: initial point x_0 , iterations N , learning rate η .

```

 $D \leftarrow 0$ 
for  $t = 0, \dots, N$  do
   $g_t \leftarrow \nabla_{x_t} f$ 
   $D_t \leftarrow D_{t-1} + D_t^E$ 
   $x_{t+1} \leftarrow x_t - \eta D_{t-1}^{-\frac{1}{2}} g_t$ 
end for

```

faster training times for autoencoder tasks when compared to vanilla SGD, RMSProp [48] and a preconditioned SGD using D^J as preconditioner. The pseudocode for ESGD is given in algorithm 2.

To efficiently calculate D^E during each iteration, it is unnecessary to compute the entire Hessian and then determine the 2-norm of each row, which would consume $\mathcal{O}(n^2)$ memory for an objective function with n parameters. Alternatively, one leverages the following matrix free estimator result of Bradley and Murray [6]

$$\|H_{i,\cdot}\| = \mathbb{E}[(Hv)^2] \text{ for } v \sim \mathcal{N}(0, 1). \quad (6)$$

Eqn (6) allows one to compute D^J at each iteration with a cost equivalent to one backpropagation.

AdaHessian [66] is a recent algorithm that applies a diagonal preconditioner in the setting of moving averages and thus can be seen as a way of enhancing the Adam optimizer [26]. Given an objective function f , AdaHessian, like Adam and RMSProp, computes a first order moment m_t and a second order moment v_t :

$$m_t = \frac{(1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} g_i}{1 - \beta_1^t} \quad (7)$$

$$v_t = \left(\sqrt{\frac{(1 - \beta_2) \sum_{i=1}^t D_i D_i}{1 - \beta_2^t}} \right)^k \quad (8)$$

where g_i denotes the gradient iteration i , D_i denotes the diagonal of the Hessian at iteration i and $0 < \beta_1, \beta_2 < 1$ are hyperparameters that are also used in the Adam optimizer. When $k = 1$, equation (8) defines a moving average preconditioner that uses the diagonal of the Hessian, which implies in such cases AdaHessian is using the Jacobi preconditioner D^J . Algorithm 3 gives the pseudocode for AdaHessian.

Algorithm 3 AdaHessian

Require: initial point x_0 , iterations N , learning rate η , exponential decay rates β_1, β_2 , power k

```

 $v_0 \leftarrow 0$ 
 $m_0 \leftarrow 0$ 
for  $t = 1, \dots, N$  do
   $g_t \leftarrow \nabla_{x_t} f$ 
   $D_t \leftarrow \text{diag}(H_t) = \text{diagonal of current Hessian}$ 
  Update  $v_t$  and  $m_t$  using eqns (7), (8)
   $x_{t+1} \leftarrow x_t - \eta \frac{m_t}{v_t}$ 
end for

```

To efficiently calculate D^J during each iteration, one can use the following matrix free estimator result of Bekas et al. [3]

$$\text{Diag}(H) = \mathbb{E}[v \odot Hv] \text{ for } v \sim \mathcal{R} \quad (9)$$

where \mathcal{R} represents a Rademacher distribution.

Kronecker Factored Preconditioners: are preconditioners of the form $D = A^T A$, where $A = A_1 \otimes \dots \otimes A_k$ and \otimes denotes the Kronecker product. The algorithm K-FAC [33] employs a Kronecker factored preconditioner.

Others: There are several second order optimisation algorithms that use preconditioners that do not fall into the above three categories such as, The Gauss-Newton method, The Conjugate Gradient

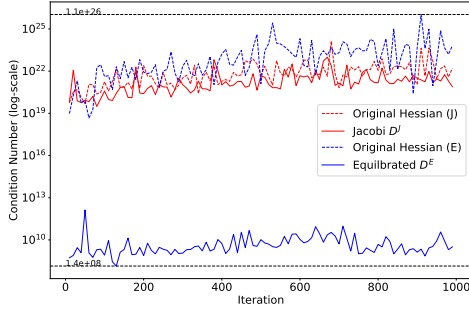


Figure 1: Comparison of condition number (10) before and after applying equilibrated and Jacobi preconditioners on a 1D signal regression. The condition number of the equilibrated preconditioner has significantly reduced (solid blue vs. dotted blue) than Jacobi (solid red vs. dotted red)

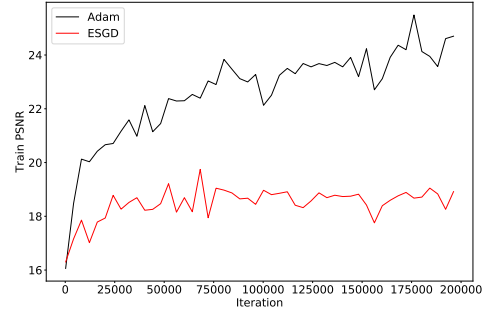


Figure 2: Comparison of training convergence (PSNR) between Adam and ESGD on a NeRF task with ReLU-PE. The ReLU-PE INR trained with Adam significantly outperformed the one trained with ESGD, suggesting curvature aware preconditioners offer no benefit for ReLU-PE.

method, Limited memory Broyden-Fisher-Golfarb- Shanno (L-BFGS) [37]. It was shown in [54] that these algorithms do not train well in stochastic settings. Another algorithm of importance for this work is Shampoo, which applies a left and right preconditioning matrix, see [24] for details.

3.2 Condition Number

In Sec. 3.1, we explained that SGD will move faster through the loss landscape when the curvatures in each direction are similar. Quantitatively the regularity of the loss landscape at a point is measured via the condition number of the Hessian.

Definition 1. Given a matrix A we define the condition number of A by

$$\kappa(A) := \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}. \quad (10)$$

Using eqn (3) we seek a preconditioner D such that $\kappa((D^{-\frac{1}{2}})^T(H(f)(\mathbf{x}))(D^{-\frac{1}{2}}))$ is closer to 1 than $\kappa(H)$. This implies that the preconditioner D is regularising the curvatures of the landscape to be equal in each direction.

In [23], Guggenheimer et al. prove the following upper bound on the conditioner number of an $N \times N$ Hessian matrix.

$$\kappa(H) < \frac{2}{|\det(H)|} \left(\frac{\|H\|_F}{\sqrt{N}} \right). \quad (11)$$

Furthermore, they show that a tight bound is expected when all the singular values, except for the smallest, are roughly equal. In [16], it is shown that by applying an equilibrated preconditioner the above bound can be tightened by a factor of $\det(D^E) \left(\frac{\|H\|_F}{\sqrt{N}} \right)^N$.

In our empirical validation (based on [16]), we examined a one-dimensional INR scenario. The INR consisted of two hidden layers, each with 16 neurons, aimed at fitting the signal $f(x) = \sin(2\pi x) + \sin(6\pi x)$. During optimization using gradient descent, we employed an equilibrated preconditioner D^E and a Jacobi preconditioner D^J . At each optimization step, we computed both the original Hessian of the gradient descent algorithm and the transformed Hessian using equation (3). Fig. 1 illustrates the experimental results, revealing that the equilibrated descent algorithm exhibits a trajectory with a significantly lower-conditioned Hessian compared to the Jacobi preconditioner and the original Hessian of a gradient descent.

3.3 Why INRs train so well with Adam compared to SGD

Given an objective function $L(\theta) = |f(\theta)|^2$ for a non-linear function f . The Jacobian J_L of L can be computed via the chain rule: $J_L = 2J_f^T f$, where J_f denotes the Jacobian of f . The Hessian of L ,

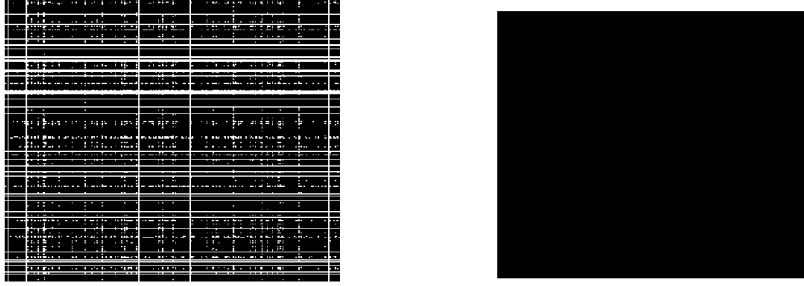


Figure 3: Visualization of the sparsity patterns in the Hessian-vector product matrices for the penultimate layer of ReLU-PE (*Left*) and Gaussian networks (*Right*): white pixels represent **zero values**, and black pixels denotes non-zero values.

H_L , can be computed by one more application of the chain rule

$$H_L = 2J_f^T J_f + 2H_f f. \quad (12)$$

We can then approximate the Hessian H_L by the term $2J_f^T J_f$. This leads to

$$H_L \approx 2J_f^T J_f = \frac{1}{2|f|^2} J_L^T J_L = \frac{1}{2L} J_L^T J_L. \quad (13)$$

Eqn. 13 yields the basic idea of the Gauss-Newton method. By taking the Gauss-Newton matrix $J_L^T J_L$ as an approximation to the Hessian H_L , we can perform an optimization update that uses the step: $\Delta_{GN} = -\frac{1}{J^T J} \nabla L$. The Gauss-Newton matrix is in general only positive semi-definite. Therefore, when implementing the algorithm one often adds a small damping factor λ so the implemented step is $\Delta_{GN} = -\frac{1}{J^T J + \lambda} \nabla L$. A similar derivation can be made for more general loss functions used in non-linear least squares type regression problems [37]. It follows that the Gauss-Newton method can be seen as a preconditioned SGD that uses a Hessian approximate preconditioner given by the Gauss-Newton matrix $J^T J$.

The Adam optimizer involves two key steps. The first is to compute a first moment estimate of the form $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$, where g_t is the gradient of a given objective function f at iteration t , and $m_0 = 0$. The second is to compute a second moment estimate $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$, where $g_t^2 = g_t \odot g_t$, \odot denotes the Hadamard product, and $v_0 = 0$. The numbers $0 < \beta_1 < \beta_2 < 1$ are hyperparameters. The algorithm then applies a bias correction to the above first and second moment estimates, and then takes a moving average to perform an update step based on $-\frac{m_t}{\sqrt{v_t}}$, see [26] for details.

For the following discussion, we will only focus on analyzing the second raw moment estimate, as this is related to the Gauss-Newton matrix. Given an objective function f , observe that we can express $g_t^2 = g_t \odot g_t = \text{Diag}(J^T J)$, where J denotes the Jacobian of f . Thus the second raw moment estimate of Adam can be recast as:

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) \text{Diag}(J_t^T J_t) \quad (14)$$

where J_t denotes the Jacobian at iteration t . Thus we see that Adam applies a preconditioner given by the diagonal of the Gauss-Newton matrix. Since the Gauss-Newton matrix can be seen as a first order approximation to the Hessian, Adam can be seen as using curvature information to adapt the learning rate for faster optimization. This suggests that Adam can outperform SGD which employs no preconditioner.

We empirically verified this insight in Fig. 4 on a binary occupancy and NeRF experiment. The details of the architectures and experiments are given in Sec. 4.2 and 4.3. It is clear from this figure that Adam outperformed SGD in both signal modalities. Similar results hold in the case of a ReLU/ReLU-PE INR, see supp. material.

3.4 Activation matters for curvature preconditioning

In Sec. 3.1, we highlighted that many implemented preconditioners in practice do not directly compute a Hessian matrix. Instead, they rely on approximations through Hessian vector products. Consequently, comprehending the Hessian vector product’s structure becomes crucial for INRs.

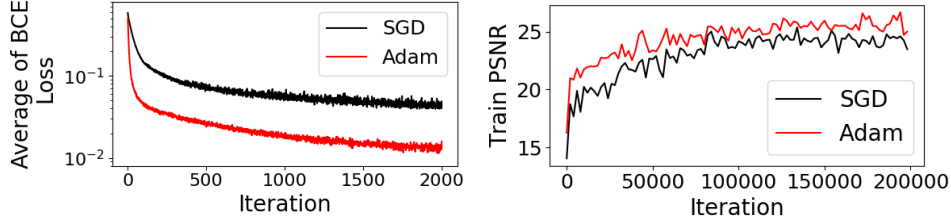


Figure 4: Comparison of training convergence between Adam and SGD on optimizing a 3D occupancy field (Left) and a NeRF (Right). Binary cross-entropy (BCE) loss measures the accuracy of the shape reconstruction, while PSNR loss assesses the scene fidelity.

Theorem 3.1. *Let F denote a sine- or Gaussian-activated INR. Let $\mathcal{L}(\theta)$ denote the MSE loss associated to F and a training set (X, Y) . Let H denote the Hessian of \mathcal{L} at a fixed parameter point θ . Then given a non-zero vector v the Hessian vector product Hv will be a dense vector.*

While it may seem rather costly for a sine- or Gaussian-activated INR to have dense Hessian vector products when applied to the MSE loss function, we point out that the cost is much less than storing a dense Hessian.

Theorem 3.2. *Let F denote a ReLU/ReLU-PE-activated INR. Let $\mathcal{L}(\theta)$ denote the MSE loss associated to F and a training set (X, Y) . Let H denote the Hessian of \mathcal{L} at a fixed parameter point θ . Then given a non-zero vector v the Hessian vector product Hv will be a sparse vector.*

Theorems 3.1 and 3.2 also hold for the binary cross-entropy loss.

Theorem 3.2 implies that when applying a preconditioner, built from a Hessian vector product, to a ReLU/ReLU-PE INR, such as ESGD [16] or AdaHessian [66], only a few components of the gradient will be scaled during the update step. This means that the preconditioner will not be adding much to the optimization process. On the other hand, theorem 3.1 implies that applying a curvature preconditioner to the gradient of the MSE loss of a sine/Gaussian-activated INR should help optimization as many of the gradient components will get scaled by a quantity depending on the curvature of the loss landscape.

Fig. 3 shows the results of a Hessian vector product of the MSE loss of a ReLU-PE and Gaussian INR with respect to the weights in the penultimate layer, see Sec. 4.1 for the network configuration. The random vector was chosen from a Rademacher distribution. Hessian vector product with respect to weights in the other layers of the neural network exhibited similar properties. As can be seen from the experiment the ReLU-PE network produces a sparse vector, as implied by theorem 3.2, and the Gaussian one produces a dense vector, as implied by theorem 3.1.

Theorem 3.1 suggests that training a ReLU/ReLU-PE INR with a curvature aware preconditioner that uses Hessian vector products would not lead to any significant advantage. In order to empirically verify this, we trained a ReLU-PE INR on a NeRF task, see Sec. 4.3 for details on the configuration, using ESGD and Adam. As shown in Fig. 2, the ReLU-PE INR trained with Adam significantly outperforms the one trained with ESGD. This suggests that for ReLU/ReLU-PE activated INRs, such curvature aware preconditioners offer no benefit.

4 Experiments

We compare different preconditioners in training of INR on various popular benchmarks [25, 31, 54, 55]: 2D image reconstruction, 3D shape reconstruction and novel view synthesis using NeRF. We focus on analyzing Gaussian networks, with variance denoted by σ , with preconditioners, which are ESGD [16]¹, AdaHessian [66]², Shampoo [24]³, and Kronecker-based preconditioner [33]⁴, see supp. material for analysis for other non-traditional activations and the reproducibility details.

¹<https://github.com/crowsonkb/esgd>

²<https://github.com/amirgholami/adahessian>

³https://github.com/google-research/google-research/tree/master/scalable_shampoo/pytorch

⁴<https://github.com/Thrandis/EKFAC-pytorch>

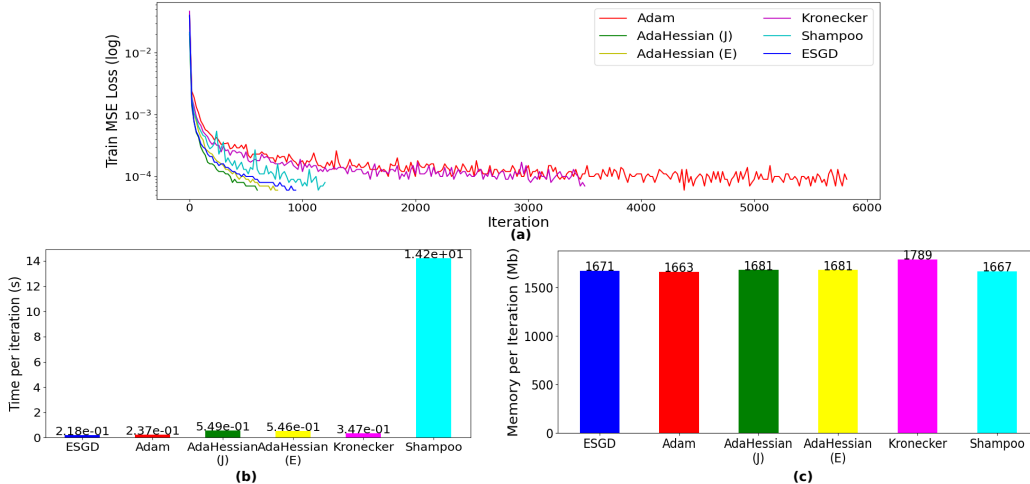


Figure 5: Comparison of training convergence in Mean Squared Error (MSE) loss (a) and computational complexity of various preconditioner (b and c). ESGD demonstrates superior performance compared to other curvature-aware preconditioners on an image reconstruction task, striking a balance between accuracy and computational efficiency.

Remark: We observed that in the implementation of AdaHessian, an equilibrated preconditioner D^E is utilized, even though the paper primarily focuses on using a Jacobi preconditioner [16]. Therefore, we introduced a Jacobi preconditioner D^J into AdaHessian, which we refer to as **AdaHessian(J)**. We use the term **AdaHessian(E)** for the original implementation. Since the local Hessian computed with minibatches is generally noisy [66], we opted to recompute the D^E for every N iterations for computational efficiency, we use $N = 50$ for all the experiments. Additionally, we also employed a moving average technique [26, 66] to mitigate the impact of noise.

4.1 Image Reconstruction

Given pixel coordinates $\mathbf{x} \in \mathbb{R}^2$, we optimized an INR to regress the associated RGB values $c \in \mathbb{R}^3$. We use a 5-layer Gaussian-activated network with 256 neurons and $\sigma = 0.05$. The network is trained stochastically using mini-batches of size 512. As seen in Fig. 5(a), while both AdaHessian(J) and AdaHessian(E) demonstrate faster convergence than ESGD in terms of iterations (Fig. 5b), it is important to note that their training time per iteration is longer than ESGD.

4.2 3D Shape Reconstruction

Given 3D point clouds $\mathbf{x} \in \mathbb{R}^3$, we optimize an INR to represent the occupancy field of 3D shapes.⁵ For training, we sampled 12 million 3D points – one-third of the points were sampled uniformly within the volume, and remaining two-thirds of the points were sampled near the mesh surface and perturbed with random Gaussian noise using sigma of 0.1 and 0.01, respectively [25]. We used a 5-layer Gaussian network with 256 neurons and $\sigma = 0.09$. The network was trained stochastically using mini-batches of size 50000. We optimized a binary cross-entropy loss between the predicted occupancy and the true occupancy. Fig. 6(a) shows that ESGD has superior convergence than other preconditioner-type optimizers, see more results in supp. material. Interestingly, when dealing with higher mode signals, AdaHessian(J) and AdaHessian(E) no longer showed superior convergence, as in Sec. 4.1. We speculated that the local Hessians might become noisy in higher mode signals, thereby impacted the gradient updates. However, this noise did not have a significant impact on ESGD’s convergence since the preconditioner updates were performed every N iterations.

⁵<http://graphics.stanford.edu/data/3Dscanrep/>



Figure 6: Comparison of training convergence (a) in average of Binary Cross-Entropy (BCE) loss of various preconditioners on optimizing a 3D occupancy field. We show the intermediate results (b) at epoch 500 for Adam [26](Left) and ESGD [16](Right). Compared to INR trained with Adam, ESGD has reconstructed the shapes with significantly improved fidelities (*zoomed in* for better visibility).

Scene	Train PSNR	Iteration ↓	Time (s) ↓	Adam Test			Iteration ↓	Time (s) ↓	ESGD Test		
				PSNR ↑	SSIM ↑	LPIPS ↓			PSNR ↑	SSIM ↑	LPIPS ↓
fern	25.38	200K	368.43	24.38	0.74	0.28	120K	270.81	24.41	0.74	0.30
flower	29.17	200K	469.31	25.67	0.78	0.12	120K	263.93	25.65	0.78	0.13
orchids	23.3	200K	501.90	19.56	0.61	0.24	120K	269.01	19.74	0.61	0.21
fortress	30.18	140K	333.70	28.93	0.81	0.15	140K	334.56	28.6	0.80	0.17
leaves	21.36	200K	482.34	19.56	0.61	0.24	120K	271.25	19.43	0.59	0.25

Table 1: Quantitative results of NeRF on instances from the LLFF dataset [34]. Note that the report time solely refers to the optimization update step.

4.3 NeRF

NeRF is a compelling strategy for utilizing an INR to model 3D objects and scenes using multi-view 2D images. This approach shows promise for generating high-fidelity reconstruction in novel view synthesis task [15, 29, 35, 47]. Given 3D points $\mathbf{x} \in \mathbb{R}^3$ and viewing direction, NeRF estimates the radiance field of a 3D scene which maps each input 3D coordinate to its corresponding volume density $\sigma \in \mathbb{R}$ and directional emitted color $\mathbf{c} \in \mathbb{R}^3$ [15, 29, 35]. We used a 8-layer Gaussian-activated network with 256 neurons and $\sigma = 0.1$. Table 1 shows that ESGD achieves on average faster convergence than Adam, see more results in the supp. material.

5 Conclusion

We investigated the use of preconditioners in stochastic training of INRs. We found that employing a curvature-aware preconditioner significantly speeds up training for INRs admitting non-traditional activations such as Gaussian functions. Experimental validation across signal modalities, including images, shape reconstruction, and neural radiance fields, using ESGD and AdaHessian algorithms supported our findings. However, for INRs with ReLU/ReLU-PE activations, curvature-aware preconditioners offered no significant advantage over Adam. This was confirmed in a comparison study on a NeRF task, where Adam outperformed ESGD for a ReLU-PE INR.

6 Limitations

Our investigation into condition numbers, as detailed in Sec. 3.2, revealed that an equilibrated preconditioner effectively reduces the condition number of the loss landscape compared to a Jacobi preconditioner. This theoretical insight was established in [16] and further validated through empirical experiments involving autoencoder tasks. Surprisingly, our image regression experiment (see Figure 5) demonstrated that a Jacobi preconditioner yielded faster training results than an equilibrated preconditioner. These findings suggest that when employing a preconditioned gradient descent algorithm for optimizing implicit neural representations, the influence of the condition number on convergence may be more complex than initially anticipated. Further research in this area is warranted to gain a deeper understanding.

References

- [1] Matan Atzmon, David Novotny, Andrea Vedaldi, and Yaron Lipman. Augmenting implicit neural shape representations with explicit deformation fields. *arXiv preprint arXiv:2108.08931*, 2021.
- [2] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [3] Costas Bekas, Effrosyni Kokiopoulou, and Yousef Saad. An estimator for the diagonal of a matrix. *Applied numerical mathematics*, 57(11-12):1214–1229, 2007.
- [4] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM review*, 60(2):223–311, 2018.
- [5] Andrew M Bradley. Algorithms for the equilibration of matrices and their application to limited-memory quasi-newton methods. Technical report, STANFORD UNIV CA, 2010.
- [6] Andrew M Bradley and Walter Murray. Matrix-free approximate equilibration. *arXiv preprint arXiv:1110.2805*, 2011.
- [7] Charles George Broyden. The convergence of a class of double-rank minimization algorithms 1. general considerations. *IMA Journal of Applied Mathematics*, 6(1):76–90, 1970.
- [8] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 333–350. Springer, 2022.
- [9] Anpei Chen, Zexiang Xu, Xinyue Wei, Siyu Tang, Hao Su, and Andreas Geiger. Factor fields: A unified framework for neural fields and beyond. *arXiv preprint arXiv:2302.01226*, 2023.
- [10] Boyuan Chen, Robert Kwiatkowski, Carl Vondrick, and Hod Lipson. Fully body visual self-modeling of robot morphologies. *Science Robotics*, 7(68):eabn1944, 2022.
- [11] Wei-Ting Chen, Wang Yifan, Sy-Yen Kuo, and Gordon Wetzstein. Dehazenerf: Multiple image haze removal and 3d shape reconstruction using neural radiance fields. *arXiv preprint arXiv:2303.11364*, 2023.
- [12] Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8628–8638, 2021.
- [13] Yun-Chun Chen, Adithyavairavan Murali, Balakumar Sundaralingam, Wei Yang, Animesh Garg, and Dieter Fox. Neural motion fields: Encoding grasp trajectories as implicit value functions. *arXiv preprint arXiv:2206.14854*, 2022.
- [14] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019.
- [15] Shin-Fang Chng, Sameera Ramasinghe, Jamie Sherrah, and Simon Lucey. Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, pages 264–280. Springer, 2022.
- [16] Yann Dauphin, Harm De Vries, and Yoshua Bengio. Equilibrated adaptive learning rates for non-convex optimization. *Advances in neural information processing systems*, 28, 2015.
- [17] Boyang Deng, J. P. Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Nasa neural articulated shape approximation. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 612–628, Cham, 2020. Springer International Publishing.
- [18] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [19] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.
- [20] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

- [21] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of computation*, 24(109):23–26, 1970.
- [22] Anne Greenbaum. *Iterative methods for solving linear systems*. SIAM, 1997.
- [23] Heinrich W Guggenheimer, Alan S Edelman, and Charles R Johnson. A simple estimate of the condition number of a linear system. *The College Mathematics Journal*, 26(1):2–5, 1995.
- [24] Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In *International Conference on Machine Learning*, pages 1842–1850. PMLR, 2018.
- [25] Amir Hertz, Or Perel, Raja Giryes, Olga Sorkine-Hornung, and Daniel Cohen-Or. Sape: Spatially-adaptive progressive encoding for neural optimization. *Advances in Neural Information Processing Systems*, 34: 8820–8832, 2021.
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] Xi-Lin Li. Preconditioned stochastic gradient descent. *IEEE transactions on neural networks and learning systems*, 29(5):1454–1466, 2017.
- [28] Yunzhu Li, Shuang Li, Vincent Sitzmann, Pulkit Agrawal, and Antonio Torralba. 3d neural scene representations for visuomotor control. In *Conference on Robot Learning*, pages 112–123. PMLR, 2022.
- [29] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021.
- [30] Tianyi Lin, Shiqian Ma, and Shuzhong Zhang. Iteration complexity analysis of multi-block admm for a family of convex minimization without strong convexity. *Journal of Scientific Computing*, 69:52–81, 2016.
- [31] David B Lindell, Dave Van Veen, Jeong Joon Park, and Gordon Wetzstein. Bacon: Band-limited coordinate networks for multiscale scene representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16252–16262, 2022.
- [32] Andy Lomas. Synthesis of abstract dynamic quasiperiodic 3d forms using sirens. *Proceedings of EVA London 2022*, pages 08–12, 2022.
- [33] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417. PMLR, 2015.
- [34] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019.
- [35] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [36] Rotem Mulayoff, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability: A view from function space. *Advances in Neural Information Processing Systems*, 34:17749–17761, 2021.
- [37] Jorge Nocedal and Stephen J Wright. *Numerical optimization*. Springer, 1999.
- [38] Marco Orsingher, Paolo Zani, Paolo Medici, and Massimo Bertozzi. Learning neural radiance fields from multi-view geometry. *arXiv preprint arXiv:2210.13041*, 2022.
- [39] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [40] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021.
- [41] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021.

- [42] Zhaonan Qu, Wenzhi Gao, Oliver Hinder, Yinyu Ye, and Zhengyuan Zhou. Optimal diagonal preconditioning. *arXiv preprint arXiv:2209.00809*, 2022.
- [43] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pages 5301–5310. PMLR, 2019.
- [44] S. Ramasinghe and S. Lucey. Beyond Periodicity: Towards a Unifying Framework for Activations in Coordinate-MLPs. In *ECCV*, 2022.
- [45] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14153–14161, June 2021.
- [46] Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. Hogwild!: A lock-free approach to parallelizing stochastic gradient descent. *Advances in neural information processing systems*, 24, 2011.
- [47] C. Reiser, S. Peng, Y. Liao, and A. Geiger. KiloNeRF: Speeding Up Neural Radiance Fields With Thousands of Tiny MLPs. In *ICCV*, 2021.
- [48] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [49] W Rudin. *Functional analysis tata mcgraw*, 1973.
- [50] Yousef Saad. *Iterative methods for sparse linear systems*. SIAM, 2003.
- [51] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [52] Vishwanath Saragadam, Jasper Tan, Guha Balakrishnan, Richard G Baraniuk, and Ashok Veeraraghavan. Miner: Multiscale implicit neural representation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIII*, pages 318–333. Springer, 2022.
- [53] Vishwanath Saragadam, Daniel LeJeune, Jasper Tan, Guha Balakrishnan, Ashok Veeraraghavan, and Richard G Baraniuk. Wire: Wavelet implicit neural representations. *arXiv preprint arXiv:2301.05187*, 2023.
- [54] Hemanth Saratchandran, Shin-Fang Chng, Sameera Ramasinghe, Lachlan MacDonald, and Simon Lucey. Curvature-aware training for coordinate networks. *arXiv preprint arXiv:2305.08552*, 2023.
- [55] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, G., and Wetzstein. Implicit Neural Representations with Periodic Activation Functions. In *NIPS*, 2020.
- [56] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems*, 32, 2019.
- [57] Ivan Skorokhodov, Savva Ignatyev, and Mohamed Elhoseiny. Adversarial generation of continuous images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10753–10764, 2021.
- [58] Elias M Stein and Rami Shakarchi. *Measure theory, integration, and hilbert spaces*, 2005.
- [59] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022.
- [60] Yu Sun, Jiaming Liu, Mingyang Xie, Brendt Wohlberg, and Ulugbek S Kamilov. Coil: Coordinate-based internal learning for imaging inverse problems. *arXiv preprint arXiv:2102.05181*, 2021.
- [61] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems*, 33: 7537–7547, 2020.

- [62] Julen Urain, An T Le, Alexander Lambert, Georgia Chalvatzaki, Byron Boots, and Jan Peters. Learning implicit priors for motion optimization. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7672–7679. IEEE, 2022.
- [63] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf–: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- [64] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. NeRF—: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- [65] Yitong Xia, Hao Tang, Radu Timofte, and Luc Van Gool. Sinerf: Sinusoidal neural radiance fields for joint pose estimation and scene reconstruction. *arXiv preprint arXiv:2210.04553*, 2022.
- [66] Zhewei Yao, Amir Gholami, Sheng Shen, Mustafa Mustafa, Kurt Keutzer, and Michael Mahoney. Adahessian: An adaptive second order optimizer for machine learning. In *proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 10665–10673, 2021.
- [67] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4578–4587, June 2021.
- [68] Jianqiao Zheng, Sameera Ramasinghe, Xueqian Li, and Simon Lucey. Trading positional complexity vs deepness in coordinate networks. In *Computer Vision – ECCV 2022*, pages 144–160, Cham, 2022. Springer Nature Switzerland.
- [69] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12786–12796, 2022.

7 Supplementary Material

7.1 Code

We provide a demo code in Google Colab to showcase the superior convergence of ESGD versus Adam in the optimization of a Gaussian-INR for a 2D image reconstruction task. Please find the code in the following link: <https://colab.research.google.com/drive/1PYsS4UCOVHn2A5qze61WzgGTY9uuwMst?usp=sharing>

7.2 Algorithms

In **Sec. 3.1 of the main paper**, we spoke about Kronecker factored preconditioners and mentioned that there were other algorithms that used other types of preconditioning techniques such as Gauss-Newton, L-BFGS and Shampoo. In this section we discuss the K-FAC algorithm that makes use of a Kronecker factored preconditioner and the Shampoo algorithm that makes use of left and right preconditioning matrices.

7.2.1 A summary of K-FAC

Kronecker Factored Approximate Curvature (K-FAC) is a second order preconditioning algorithm introduced in [33]. The algorithm seeks to approximate the Hessian of an objective function via the Fisher information matrix [33]. In general, the Fisher information matrix is dense matrix, thus for implementation purposes this would require a huge memory cost for large parameter objective functions. Therefore, the authors impose that the approximation should factor into Kronecker products, where matrix forming the Kronecker product comes from the layer of the neural network.

For this section we fix an objective function f that we want to minimize. Let $G = J^T J$ denote the Gauss-Newton matrix of f , where $J = \nabla f^T$ is the Jacobian of f . When implementing K-FAC one must resort to the empirical Fisher information matrix, which is defined as as the expected value of G

$$F = \mathbb{E}(G) \approx \mathbb{E}(H). \tag{15}$$

Since G is a first order approximation of the hessian H of f , we see that the empirical Fisher information matrix approximates the expected value of the Hessian.

Given a neural network Φ with l layers. F will be a block diagonal matrix with $l \times l$ blocks. Each block \tilde{F}_{ij} of F is given by

$$\mathbb{E}(\nabla_{\theta^i} f \otimes \nabla_{\theta^j} f) \quad (16)$$

where i and j are layers in the network and θ_i the parameters in layer i .

The idea of the algorithm is to now impose a Kronecker product structure on each such block and form the approximation

$$\mathbb{E}(\nabla_{\theta^i} f \otimes \nabla_{\theta^j} f) \approx \mathbb{E}(\tilde{\phi}^{i-1}(\tilde{\phi}^{j-1})^T) \otimes \mathbb{E}(\delta^i(\delta^j)^T) \quad (17)$$

where $\tilde{\phi}^i$ is the activation values in layer i with an appended 1 in the last position of the vector i.e. $(\tilde{\phi}^i)^T = [(\phi^i)^T 1]$. This is done as we are treating the weights and biases of the network together. The terms δ^i arise from standard backpropagation formulas.

We thus see that the empirical Fisher information matrix can be approximated by Kronecker products $\mathbb{E}(\tilde{\phi}^{i-1}(\tilde{\phi}^{j-1})^T) \otimes \mathbb{E}(\delta^i(\delta^j)^T)$. Abusing notation and denoting this approximation by F as well, the K-FAC algorithm update is

$$\theta_{k+1} = \theta_k - \eta F^{-1} \nabla_{\theta^i} f. \quad (18)$$

We thus see that K-FAC is a preconditioned gradient descent algorithm with preconditioner given by a Kronercker factored approximation to the empirical Fisher information matrix. For more details on the K-FAC algorithm we refer the reader to [33].

7.2.2 A summary of Shampoo

Shampoo is a preconditioning algorithm that optimizes objective functions over tensor spaces, introduced in [24]. In general, the algorithm works to optimize parameters that can be represented as general tensors. In this regard, the optimizer is extremely effective when working on general neural architectures such as convnets. Implicit neural representations (INRs) are represented by feed forward neural networks, and hence their parameters are represented by a 2-tensor i.e. a matrix. Therefore, we will be analyzing Shampoo in the setting that our parameters are represented as matrices. However, we point out to the reader that the algorithm does work on more general tensors, see [24] for details.

As explained in section 3.1, a preconditioner is in general a matrix that acts on the gradient before an update is performed. This action is done by matrix multiplication on the left of the gradient. Let us denote the preconditioner by L , then the action is given by

$$L \cdot g \quad (19)$$

We can also act on the gradient on the right. Given a matrix R , we could also do

$$g \cdot R. \quad (20)$$

Combining (19) and (20) we obtain the left and right preconditioning transformation

$$L \cdot g \cdot A. \quad (21)$$

If we fix an objective function $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, e.g. the MSE loss of a neural network. The shampoo algorithm for f (in the case of 2-tensors) takes as L_t the Gauss-Newton matrix $J^T J$, an $m \times m$ -matrix, and as R_t the matrix $J J^T$, an $n \times n$ -matrix. The pseudocode for Shampoo (in the case of 2-tensors) is given in algorithm 4.

For the experiments carried out in this paper (including supplementary material) we will make use of the Shampoo algorithm for 2-tensors, as INRs are feedforward networks. For details on Shampoo for general tensors we refer the reader to the paper [24].

7.3 Activation matters for curvature preconditioning

In this section, we present the proof of theorems 3.1 and 3.2 from **Sec. 3.4 of the main paper**.

Algorithm 4 Shampoo

Require: initial point x_0 , iterations N , learning rate η , $\epsilon > 0$.

```
 $L_0 \leftarrow \epsilon I_m$   
 $R_0 \leftarrow \epsilon I_n$   
for  $t = 0, \dots, N$  do  
   $g_t \leftarrow \nabla_{x_t} f$   
   $L_t \leftarrow L_{t-1} + J_t^T J_t$   
   $R_t \leftarrow R_{t-1} + J_t J_t^T$   
   $x_{t+1} \leftarrow x_t - \eta L_t^{-1/4} g_t R_t^{-1/4}$   
end for
```

7.3.1 Preliminaries

We fix a feedforward network of L layers and widths $\{n_0, n_1, \dots, n_L\}$ $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L}$ defined by

$$f : x \rightarrow T_L \circ \psi \circ T_{L-1} \circ \dots \circ \psi \circ T_1(x) \quad (22)$$

where $T_i : x_i \rightarrow \theta_i x_i + b_i$ is an affine transformation with trainable weights $\theta_i \in \mathbb{R}^{n_i \times n_{i-1}}$, $b_i \in \mathbb{R}^{n_i}$, and ψ is a non-linear activation acting component wise. The number n_0 is the input dimension and the number n_L is the output dimension. The composition $\psi \circ T_k \circ \dots \circ \psi \circ T_1$ gives the first k -layers of the neural network function.

So as to make things easier we will assume our network only has weight parameters and no bias.

We fix a data set, which we denote by $(X, Y) \in \mathbb{R}^{n_0 \times N} \times \mathbb{R}^{n_L \times N}$, consisting of input data $X \in \mathbb{R}^{n_0 \times N}$ and targets $Y \in \mathbb{R}^{n_L \times N}$. X will have feature dimension n_0 and the number of training samples will be given by N . Thus we can think of (X, Y) as a collection of training samples $\{(x_i, y_i)\}_{i=1}^N$, with each pair (x_i, y_i) a training point.

An implicit neural representation can thus be viewed as a map

$$f : \mathbb{R}^{n_0 \times N} \times \mathbb{R}^p \rightarrow \mathbb{R}^{n_L \times N} \quad (23)$$

where p denotes the parameter dimension. In general, parameter vectors $\theta \in \mathbb{R}^p$ will be denoted in coordinates by as $\theta = (\theta(1), \dots, \theta(L))$, where each $\theta(i) \in \mathbb{R}^{n_i \times n_{i-1}}$ corresponds to the parameters of the i -th-layer. Each such parameter $\theta(i)$ represents a matrix of trainable weights given by:

$$\theta(i) = \begin{bmatrix} \theta(i)_{11} & \dots & \theta(i)_{1n_{i-1}} \\ \vdots & \vdots & \vdots \\ \theta(i)_{n_i 1} & \dots & \theta(i)_{n_i n_{i-1}} \end{bmatrix}$$

Thus with the notation above parameters are represented as matrices. However, when we take derivatives with respect to parameter variables, we will to flatten the parameter matrices. We will do so by flattening each row to a column, thereby obtaining a column vector. Thus the above $\theta(i)$ will be flattened to the vector

$$\text{Vec}(\theta(i)) = (\theta(i)_{11}, \dots, \theta(i)_{1n_{i-1}}, \dots, \theta(i)_{n_i 1}, \dots, \theta(i)_{n_i n_{i-1}})^T. \quad (24)$$

We will always assume such flattening of matrices has been done and therefore we will not explicitly use the *vec* notation. At times we will have to go back and forth between the actual parameter matrix and its associated flattened vector. The context should make it clear as to which representation we are dealing with.

Feedforward neural networks are composed of their layers. We express this by writing

$$f = f_L \circ \dots \circ f_1 \quad (25)$$

where for each $1 \leq k \leq L$

$$f_k : \mathbb{R}^{n_{k-1} \times N} \times \mathbb{R}^{n_k \times n_{k-1}} \times \dots \times \mathbb{R}^{n_L \times n_{L-1}} \rightarrow \mathbb{R}^{n_k \times N} \times \mathbb{R}^{n_{k+1} \times n_k} \times \dots \times \mathbb{R}^{n_L \times n_{L-1}} \quad (26)$$

is the map defined by

$$f_k(Z, \theta(k), \dots, \theta(L)) = (\psi(\theta(k) \cdot Z), \theta(k+1), \dots, \theta(L))^T. \quad (27)$$

The quantity $\theta(k) \cdot Z$ is the matrix in $\mathbb{R}^{n_k \times N}$ given by applying $\theta(k)$ viewed as a $n_k \times n_{k-1}$ matrix acting on a $n_{k-1} \times N$ matrix Z . Furthermore, the latter entries $(\theta(k+1), \dots, \theta(L))$ are all viewed as flattened vectors. This is an example of how we have had to use parameters both in their matrix form and their flattened vector form.

Using (27), the map f_1 is a map

$$f_1 : \mathbb{R}^{n_0 \times N} \times \mathbb{R}^{n_1 \times n_0} \times \dots \times \mathbb{R}^{n_L \times n_{L-1}} \rightarrow \mathbb{R}^{n_1 \times N} \times \mathbb{R}^{n_2 \times n_1} \times \dots \times \mathbb{R}^{n_L \times n_{L-1}} \quad (28)$$

where \mathbb{R}^{n_0} is the input space, the space in which our input data resides. As we will not be taking any derivatives with respect to data, and our data set has already been fixed, we will make life easier by viewing

$$f_1 : \mathbb{R}^{n_1 \times n_0} \times \dots \times \mathbb{R}^{n_L \times n_{L-1}} \rightarrow \mathbb{R}^{n_1 \times N} \times \mathbb{R}^{n_2 \times n_1} \times \dots \times \mathbb{R}^{n_L \times n_{L-1}} \quad (29)$$

defined by

$$f_1(\theta(1), \dots, \theta(L)) = (\psi(\theta(1) \cdot X), \theta(2), \dots, \theta(L))^T \quad (30)$$

where X is our fixed input data in $\mathbb{R}^{n_0 \times N}$.

One final notation we introduce is the following. Given a matrix $A \in \mathbb{R}^{m \times n}$, viewed as a $m \times n$ matrix, we let A^j denote the j th-row of A and A_j denote the j th column of A . With this notation, we observe the following: Given a parameter vector $\theta(k) \in \mathbb{R}^{n_k \times n_{k-1}}$, viewed as a $n_k \times n_{k-1}$ matrix, and a $Z \in \mathbb{R}^{n_{k-1} \times N}$, the product $\theta(k) \cdot Z$ is flattened to the vector

$$(\theta^1(k) \cdot Z_1, \dots, \theta^1(k) \cdot Z_N, \dots, \theta^{n_k}(k) \cdot Z_1, \dots, \theta^{n_k}(k) \cdot Z_N)^T$$

where each $\theta^j(k)$ is a row vector with n_{k-1} entries and each Z_j is a column vector with n_{k-1} entries.

For this section, we will fix our loss function as the MSE loss function. We write this loss function as

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} |f(x_i, \theta) - y_i|^2. \quad (31)$$

Observe that the MSE loss can be written as the composition $c \circ f$, where c is a convex cost function given by the average squared error:

$$c : \mathbb{R}^{n_L \times N} \rightarrow \mathbb{R} \quad (32)$$

defined by

$$c(z_1, \dots, z_N) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} |z_i - y_i|^2 \quad (33)$$

where we remind the reader that our data set consists of point $\{(x_i, y_i)\}_{i=1}^N$, with each $x_i \in \mathbb{R}^{n_0}$ and $y_i \in \mathbb{R}^{n_L}$. With this notation we can easily see that the MSE loss is given by the composition $c \circ f$. Theorems 3.1 and 3.2 from section 3.4 of the paper also hold for more general loss functions such as Binary Cross Entropy (BCE) loss.

We now compute the differential of the MSE loss function, from which a simple transpose gives the gradient. Observe that the loss function is a map

$$\mathcal{L} : \mathbb{R}^p \rightarrow \mathbb{R} \quad (34)$$

and therefore its differential is a map

$$D\mathcal{L} : \mathbb{R}^p \rightarrow \text{Lin}(\mathbb{R}^p, \mathbb{R}) \cong \mathbb{R}^{p \times 1} \quad (35)$$

where $\text{Lin}(\mathbb{R}^p, \mathbb{R})$ denotes the linear maps from \mathbb{R}^p to \mathbb{R} , which is linearly isomorphic to the space of $1 \times p$ matrices which we can identify as $\mathbb{R}^{p \times 1}$.

The chain rule gives $D\mathcal{L}(\theta) = Dc(f(\theta)) \cdot Df(\theta)$. Therefore, in order to compute the differential of the loss, it suffices to compute the differential of the cost function c and the differential of the neural network function f . The following proposition is an easy consequence of equation (33).

Proposition 7.1. $Dc(Z) = \frac{1}{N}(Z - Y)$, where Y denotes the matrix of targets from our data set.

The next step is to compute the differential of the neural network function f . The differential will be a map

$$Df : \mathbb{R}^p \rightarrow \text{Lin}(\mathbb{R}^p, \mathbb{R}^{n_L N}) \cong \mathbb{R}^{p \times n_L N}. \quad (36)$$

By equation (25), and the chain rule, we have that for a vector $\theta \in \mathbb{R}^p$

$$Df(\theta) = Df_L(f_{L-1} \circ \dots \circ f_1(\theta)) \cdot Df_{L-1}(f_{L-2} \circ \dots \circ f_1(\theta)) \cdots Df_2(f_1(\theta)) \cdot Df_1(\theta). \quad (37)$$

We therefore need to compute the differential of each f_k .

We use the following notation to denote partial derivatives with respect to the space variable Z and the parameter variable $\theta(j)$, for $k \leq j \leq L$. The partial derivative $\frac{\partial}{\partial Z}$ will denote derivatives with respect to the variable Z , and $\frac{\partial}{\partial \theta^i(j)}$ will denote derivatives with respect to the i th row of the parameter variable $\theta(j)$ for $k \leq j \leq L$. Thus for example, we have that

$$\frac{\partial}{\partial Z} \psi(\theta^j(k) \cdot Z) = \left(\frac{\partial}{\partial Z} \psi(\theta^j(k) \cdot Z_1), \dots, \frac{\partial}{\partial Z} \psi(\theta^j(k) \cdot Z_N) \right)^T. \quad (38)$$

The following lemma shows how variable derivatives of variable indices interact with each other.

Lemma 7.2. $\frac{\partial}{\partial Z_j} \psi(\theta^i(k) Z_i) = 0$ for all $j \neq i$ and $\frac{\partial}{\partial \theta^j(k)} \psi(\theta^i(k) Z_i) = 0$ for all $j \neq i$.

Proof. The result follows immediately from noting that the term $\psi(\theta^i(k) Z_i)$ does not depend on the variable Z_j and $\theta^j(k)$ for $j \neq i$. \square

Proposition 7.3. The differential Df_k is a $(Nn_k + n_{k+1}n_k \cdots + n_L n_{L-1}) \times (n_{k-1}N + n_k n_{k-1} + \cdots + n_L n_{L-1})$ matrix given by the following representation

$$\begin{bmatrix} \frac{\partial}{\partial Z} \psi(\theta^1(k) \cdot Z) & \frac{\partial}{\partial \theta^1(k)} \psi(\theta^1(k) \cdot Z) & 0 & \cdots & 0 & 0 & \cdots & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \ddots & \vdots \\ \frac{\partial}{\partial Z} \psi(\theta^{n_k}(k) \cdot Z) & 0 & 0 & \cdots & \frac{\partial}{\partial \theta^{n_k}(k)} \psi(\theta^{n_k}(k) \cdot Z) & 0 & \cdots & \cdots & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & I_{n_{k+1}n_k} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots & I_{n_L n_{L-1}} \end{bmatrix}$$

where I_i denotes an $i \times i$ identity matrix.

Proof. The matrix representation follows from a straight forward calculation of partial derivatives which we explain. First of all we remind the reader that the image of f_k is to be thought of as a flattened column vector, where remember that we flatten each row to a column. By definition

$$f_k(Z, \theta(k), \dots, \theta(L)) = (\psi(\theta(k) \cdot Z), \theta(k+1), \dots, \theta(L))^T.$$

We now observe that the terms $\theta(k+1), \dots, \theta(L)$ will give zero when we apply $\frac{\partial}{\partial Z}$ and $\frac{\partial}{\partial \theta^i(k)}$ for $1 \leq i \leq n_k$. This leads to the zeros in the bottom left of the big matrix. The term $\psi(\theta(k) \cdot Z)$ will give zero when we apply the derivatives $\frac{\partial}{\partial \theta^{k+1}}, \dots, \frac{\partial}{\partial \theta^L}$. This leads to the zeros in the top left of the big matrix.

The matrix representation of Df_k now follows from the derivatives $\frac{\partial}{\partial Z} \psi(\theta^i(k) \cdot Z)$ and $\frac{\partial}{\partial \theta^i(k)} \psi(\theta^i(k) \cdot Z)$, noting that by lemma 7.2 we have that $\frac{\partial}{\partial \theta^j(k)} \psi(\theta^i(k) \cdot Z) = 0$ for $i \neq j$. \square

The Hessian of the loss function \mathcal{L} can be computed by applying the fact that $\mathcal{L} = c \circ f$ and the chain rule. We will use the notation $D^2 \mathcal{L}$ to denote the Hessian of the loss, which is to be thought of as the second differential of \mathcal{L} .

Proposition 7.4. Given a point $\theta \in \mathbb{R}^p$ we have that

$$D^2 \mathcal{L}(\theta) = Df(\theta)^T \cdot \left(\frac{1}{N} I \right) \cdot Df(\theta) + \frac{1}{N} (f(\theta) - Y) \cdot D^2 f \quad (39)$$

where $\frac{1}{N} I$ denotes the identity matrix with $1/N$ on its diagonal and recall that Y is a matrix consisting of the targets from the fixed data set.

Proof. This follows by applying the chain and product rule to $D\mathcal{L}$. Given a point θ , we have that

$$D\mathcal{L}(\theta) = Dc(f(\theta)) \cdot Df(\theta) \quad (40)$$

Differentiating once more, and applying the chain and product rules, we get

$$D^2\mathcal{L}(\theta) = Df(\theta)^T \cdot D^2c(f(\theta)) \cdot Df(\theta) + Dc(f(\theta)) \cdot D^2f(\theta). \quad (41)$$

By proposition 7.1, we have that $Dc(f(\theta)) = \frac{1}{N}(f(\theta) - Y)$ and then differentiating once more we get $D^2c(f(\theta)) = \frac{1}{N}I$ and the result follows. \square

Proposition 7.4 implies that in order to compute the Hessian of the loss, we need to compute the Hessian of the neural network function. This can be done by the chain rule and induction.

Lemma 7.5. *We have the following decomposition for the Hessian of f*

$$\begin{aligned} D^2f &= (Df_1)^T \cdots (Df_{L-1})^T D^2f_L (Df_{L-1}) \cdots (Df_1) \\ &\quad + (Df_1)^T \cdots (Df_{L-2})^T Df_L D^2f_{L-1} (Df_{L-2}) \cdots (Df_1) \\ &\quad + (Df_1)^T \cdots (Df_{L-3})^T Df_L Df_{L-1} D^2f_{L-2} (Df_{L-3}) \cdots (Df_1) \\ &\quad + \cdots + (Df_1)^T Df_L Df_{L-1} \cdots Df_3 D^2f_2 D^2f_2 (Df_1) \\ &\quad + Df_L Df_{L-1} \cdots Df_2 D^2f_1 \end{aligned}$$

Proof. The proof of this follows by induction on the layers. \square

Lemma 7.5 implies that in order to compute the Hessian of the neural network, we need to compute the Hessian and the differential of each layer. The differential of each of the layers was already computed in proposition 7.3. We will now give a matrix formula for the Hessian of each layer f_k .

In order to compute the Hessian D^2f_k , we will flatten Df_k so that it is a map

$$Df_k : \mathbb{R}^{n_{k-1} \times N} \times \mathbb{R}^{n_k \times n_{k-1}} \times \cdots \times \mathbb{R}^{n_L \times n_{L-1}} \rightarrow \mathbb{R}^{(n_k \times N + n_{k+1} \times n_k + \cdots + n_L \times n_{L-1})(n_{k-1} \times N + n_k \times n_{k-1} + \cdots + n_L \times n_{L-1})}. \quad (42)$$

Then for a point $(Z, \theta(k), \dots, \theta(L))$, we have that $D^2f_k(Z, \theta(k), \dots, \theta(L))$ will be a $((n_k \times N + n_{k+1} \times n_k + \cdots + n_L \times n_{L-1})(n_{k-1} \times N + n_k \times n_{k-1} + \cdots + n_L \times n_{L-1})) \times (n_{k-1} \times N + n_k \times n_{k-1} + \cdots + n_L \times n_{L-1})$ matrix. In fact, it can be thought of as a collection of $(n_k \times N + n_{k+1} \times n_k + \cdots + n_L \times n_{L-1})$ square $(n_{k-1} \times N + n_k \times n_{k-1} + \cdots + n_L \times n_{L-1}) \times (n_{k-1} \times N + n_k \times n_{k-1} + \cdots + n_L \times n_{L-1})$ matrices stacked on top of each other.

Each such square matrix arises from the rows of the matrix representation of Df_k , see proposition 7.3. We start by computing these square matrices.

Lemma 7.6. *Given the matrix representation of Df_k in proposition 7.3 and $1 \leq i \leq n_k N$, we have that the derivative of the i th-row of Df_k is given by*

$$\begin{bmatrix} \frac{\partial^2}{\partial Z \partial Z} \psi(\theta^i(k) \cdot Z) & 0 & \cdots & 0 & \frac{\partial^2}{\partial \theta^i(k) \partial Z} \psi(\theta^i(k) \cdot Z) & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 \\ \frac{\partial^2}{\partial Z \partial \theta^i(k)} \psi(\theta^i(k) \cdot Z) & 0 & \cdots & 0 & \frac{\partial^2}{\partial \theta^i(k) \partial \theta^i(k)} \psi(\theta^i(k) \cdot Z) & 0 & \cdots & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \end{bmatrix}$$

Furthermore if $n_k N < i \leq n_{k+1} \times n_k + \cdots + n_L \times n_{L-1}$ then the derivative of the i th-row of Df_k will be a matrix of zeros.

Proof. The proof follows by inspecting each row of the matrix representation of Df_k given in proposition 7.3. We observe that if $1 \leq i \leq n_k N$, then the i th row of Df_k is given by

$$\left[\frac{\partial}{\partial Z} \psi(\theta^i(k) \cdot Z) \quad 0 \quad 0 \quad \cdots \quad 0 \quad \frac{\partial}{\partial \theta^i(k)} \psi(\theta^i(k) \cdot Z) \quad 0 \quad \cdots \quad \cdots \quad 0 \right]$$

We then observe that the derivatives $\frac{\partial}{\partial \theta^j(k)}$ of any element in the above row will be zero for $j \neq i$ as none of the elements in the row depend on the variable $\theta^j(k)$ when $j \neq i$. This means the only derivatives that could possibly be non-zero for such a row will come from $\frac{\partial}{\partial Z}$ and $\frac{\partial}{\partial \theta^i(k)}$. This proves the first part of the proposition. To prove the second part, we simply observe that the i th-rows of Df_k for $n_k N < i \leq n_{k+1} \times n_k + \dots + n_L \times n_{L-1}$ have only one non-zero entry which will be a 1. When differentiated with respect to any of the variables this will give zero, and thus we simply get the zero matrix for such a row. This proves the second part of the proposition. \square

Using lemma 7.6, we can compute a full matrix representation of the Hessian of f_k .

Proposition 7.7. *A matrix representation of the Hessian of f_k is given by*

$$\begin{bmatrix} \frac{\partial^2}{\partial Z \partial Z} \psi(\theta^1(k) \cdot Z) & \frac{\partial^2}{\partial \theta^1(k) \partial Z} \psi(\theta^1(k) \cdot Z) & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \frac{\partial^2}{\partial Z \partial \theta^1(k)} \psi(\theta^1(k) \cdot Z) & \frac{\partial^2}{\partial \theta^1(k) \partial \theta^1(k)} \psi(\theta^1(k) \cdot Z) & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ \frac{\partial^2}{\partial Z \partial \theta^{n_k}(k)} \psi(\theta^{n_k}(k) \cdot Z) & 0 & 0 & \dots & 0 & \frac{\partial^2}{\partial \theta^{n_k}(k) \partial \theta^{n_k}(k)} \psi(\theta^{n_k}(k) \cdot Z) & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \frac{\partial^2}{\partial Z \partial \theta^{n_k}(k)} \psi(\theta^{n_k}(k) \cdot Z) & 0 & 0 & \dots & 0 & \frac{\partial^2}{\partial \theta^{n_k}(k) \partial \theta^{n_k}(k)} \psi(\theta^{n_k}(k) \cdot Z) & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix}$$

Proof. The proof follows by taking each square matrix from lemma 7.6 and stacking them on top of each other. \square

Proposition 7.3 shows that the gradient of the k th-layer neural function has some sparsity in the matrix representation. Furthermore, proposition 7.7 shows that the Hessian of the k th-layer neural function f_k has some sparsity, regardless of the activation. We will not consider such structures as sparse as they are shared by all implicit neural representations. Rather, we will show in the next section, that much more sparsity can be obtained in the Hessian of the loss when the neural network uses a ReLU activation, by showing that the derivative terms in 7.7 vanish with high probability.

7.3.2 Proofs of theorem 3.1 and 3.2 from the main paper

From propositions 7.4, 7.3, 7.7 and lemma 7.5 we see that in order to compute the Hessian of the MSE loss \mathcal{L} , we need to understand the first and second derivatives of our activation function.

Proposition 7.8.

1. Let $\phi(x) = \sin(\omega x)$. Then $\phi'(x) = \omega \cos(\omega x)$ and $\phi''(x) = -\omega^2 \sin(\omega x)$.
2. Let $\phi(x) = e^{-\frac{\sigma^2 x^2}{2}}$. Then $\phi'(x) = -\sigma^2 x e^{-\frac{\sigma^2 x^2}{2}}$ and $\phi''(x) = -\sigma^2 e^{-\frac{\sigma^2 x^2}{2}} + \sigma^4 x^2 e^{-\frac{\sigma^2 x^2}{2}}$.

The proof of the above proposition is straightforward from standard differentiation rules along with the chain rule.

In the paper [53], a Gabor wavelet is used as an activation function for implicit neural representations. The Gabor wavelet $\psi(x; \omega_0; s_0)$ is a complex valued function, and its real part is given by $\cos(\omega_0 x) e^{-s_0^2 x^2}$ and its imaginary part by $\sin(\omega_0 x) e^{-s_0^2 x^2}$. Applying the chain rule and proposition 7.8 we get similar derivative formulae for the real and imaginary part of the Gabor wavelet. Therefore, for the following we will focus on sine, Gaussian and ReLU activations with the knowledge that the

proofs for the sine and Gaussian case go through for the Gabor wavelet case by restricting to the real and imaginary parts of the wavelet.

We will use some basic concepts from measure theory such as sets of measure zero. The reader who is not familiar with such material may consult the book [58].

Proposition 7.9. *Let $\phi(x)$ denote a sine, Gaussian or Gabor wavelet function. On any finite interval I in \mathbb{R} , the set of zeros of ϕ , ϕ' and ϕ'' have Lebesgue measure zero. In particular, normalizing the Lebesgue measure to have measure 1 on I , so that we obtain a probability measure, we have that the zeros of ϕ , ϕ' and ϕ'' have probability zero.*

Proof. The proof of this follows by first observing that on any finite interval ϕ only has finitely many zeros. Using proposition 7.8, we see that ϕ' and ϕ'' also have only finite zeros on a finite interval. The result follows. \square

Proposition 7.10. 1. $\frac{d}{dx} ReLU(x) = H(x)$, where H is the step function centered at 0.

2. Viewing $\frac{d}{dx} ReLU(x)$ as a distribution, we have that $\frac{d^2}{dx^2} ReLU = \delta$, where δ denotes a Dirac delta distribution centered at 0.

Proof. By definition $ReLU(x) = \max(0, x)$, therefore for $x < 0$ it is clear that $\frac{d}{dx} ReLU(x) = 0$. For $x \geq 0$, we have that $ReLU(x) = x$ and therefore $\frac{d}{dx} ReLU(x) = 1$ for such points. It follows that one can represent the derivative $\frac{d}{dx} ReLU(x) = H(x)$ distributionally, with a discontinuity at the origin.

We move on to proving the second identity. Given a function $f \in C_c^\infty(\mathbb{R})$ the distribution H is defined by

$$\langle H, f \rangle = \int_{-\infty}^{\infty} H(x)f(x)dx = \int_0^{\infty} f(x)dx.$$

The derivative of H is then given by (see [49] for preliminaries on derivatives of distributions)

$$\langle H', f \rangle = -\langle H, f' \rangle = -\int_0^{\infty} f'(x)dx = f(0) = \langle \delta, f \rangle$$

where the second equality comes from the fundamental theorem of calculus and the fact that f is compactly supported. The final equality follows by definition of the Dirac delta distribution. It thus follows that $\frac{d^2}{dx^2} ReLU = \delta$ as distributions. \square

Proposition 7.11. *Let ϕ denote a ReLU activation. If I is an interval of finite non-zero Lebesgue measure that contains a sub-interval of negative numbers of non-negative measure. Then the zero set of ϕ and ϕ' has non-zero measure. If I is any interval of finite non-zero Lebesgue measure then the zero set of ϕ'' has non-zero Lebesgue measure. By normalizing the Lebesgue measure of I to be 1, we get that the probability of the zero set of ϕ and ϕ' is non-zero when I contains a sub-interval of negative numbers of non-zero probability. Furthermore, the probability of the zero set of ϕ'' is non-zero.*

Proof. The proof of this follows from proposition 7.10. \square

Proof of theorem 3.1 from main paper. By propositions 7.8, 7.9 we get that the derivative terms in the formula 7.7, will be non-zero with probability 1. Then from propositions 7.4, 7.3, we have that the Hessian will in general be a dense matrix. Therefore, given any non-zero vector v , we have that the Hessian vector product Hv will be a dense vector with high probability. \square

Proof of theorem 3.2 from main paper. By propositions 7.10 and 7.11 we get that the the derivative terms in the formula 7.7, will be zero with high probability. Then from propositions 7.4, 7.3 we have that the Hessian will in general be a very sparse matrix. Therefore, given any non-zero vector v , we have that the Hessian vector product Hv will be a sparse vector with high probability. \square

7.4 Condition Number

In this section, we provide the reproducibility and implementation details for the condition number experiment in **Sec. 3.2 of the main paper**.

Data. We consider a 1D signal $f(x) = \sin(2\pi x) + \sin(6\pi x)$, where x denotes the input coordinates which scale from -1 to 1. For this experiment, we opted to use a 1D problem because this enables us to compute the **exact Hessian matrix**, which typically scales exponentially with the network’s parameters.

Architecture. We employed a 3-layer Gaussian network with a hidden layer size of 16.

Initialization. We used the default PyTorch initialization (Kaiming Uniform).

Hardware. We ran the experiment on a NVIDIA 3080 GPU with 12Gb of memory.

Hyperparameters. Unless specified, we maintained the default parameter values for each optimizer. We used a learning rate of 2. We incorporated a warmup strategy. Specifically, we implemented the equilibrated/jacobi gradient preconditioning for the first 20 iterations as a warmup phase, followed by subsequent updates every 10 iterations.

7.5 Sparsity of Hessian

In this section, we provide the reproducibility and implementation details for the experiment in **Sec. 3.4 of the main paper**. Additionally, we provide an analysis of the sparsity of the Hessian-vector products of all layers during the training process of a 2D image reconstruction task using ESGD in Sec. 7.5.2. We cover various INRs, including ReLU, ReLU-PE, sine and Gaussian.

7.5.1 Reproducibility & Implementation Details

Data. We use image from DIV2K dataset ⁶. For our experiment, we downsampled the original images, initially captured at a resolution of 512×512 pixels, to a resolution of 256×256 pixels. We used minibatch size of 512 during training. We use the *lion* image for this analysis.

Architectures. We employed a 5-layer network with a hidden layer size of 256. We used $\mathbf{x} \in \mathbb{R}^d$ to denote the input coordinates. For ReLU-PE INR, we embedded low-dimensional data inputs \mathbf{x} into a higher-dimensional space using an embedding layer $\gamma : \mathbb{R}^d \rightarrow \mathbb{R}^{d+L}$ [35]. We analyze the sparsity of the Hessian for different L , we used $L = 4$ and $L = 8$. For Sine-INR, we used a sine activation function $\mathbf{x} \rightarrow \sin(2\pi\omega\mathbf{x})$ [55]. We used $\omega = 30$. For Gaussian-INR, we employed a Gaussian activation function $\mathbf{x} \rightarrow \exp\left(\frac{|\mathbf{x}-\mu|^2}{2\sigma^2}\right)$ [44]. We employed $\mu = 0$ and $\sigma = 0.05$.

Initialization. We employed the principled initialization scheme proposed by Sitzmann et al. for sine-INR [55]. We used the default PyTorch initialization (Kaiming Uniform) for ReLU/ReLU-PE as well as Gaussian-INRs in all the experiments.

Hardware. We ran all experiments on a NVIDIA 3080 GPU with 12Gb of memory.

Hyperparameters. Unless specified, we maintained the default parameter values for each optimizer. For ReLU and ReLU-PE INRs, we used a learning rate of 0.1. As for sine and Gaussian INRs, we used a learning rate of 0.15. We incorporated a warmup strategy for all the INRs. Specifically, we implemented the equilibrated gradient preconditioning for the first 50 iterations as a warmup phase, followed by subsequent updates every N iterations, with N set to 100.

7.5.2 Sparsity of Hessian-vector products for different activations

In **Sec. 3.4 of the main paper** we theoretically analyzed Hessian vector products for different activation functions. In this section we conduct an empirical verification of Thms. 3.1 and 3.2 while

⁶<https://paperswithcode.com/dataset/div2k>

training an INR on an image regression task using Gaussian, sine, ReLU and ReLU-PE activations. To visualize the results, we generated Figure 7, which illustrates the percentage of zero entries in the Hessian matrix of each INR, averaged over all iterations of the optimization process. Our observations clearly demonstrate that the ReLU/ReLU-PE INRs consistently exhibit sparse Hessians throughout training, while the Gaussian/sine INRs consistently display dense Hessians throughout training. These findings align precisely with the predictions derived from Thms 3.1 and 3.2, as discussed in Section 3.4 of the paper.

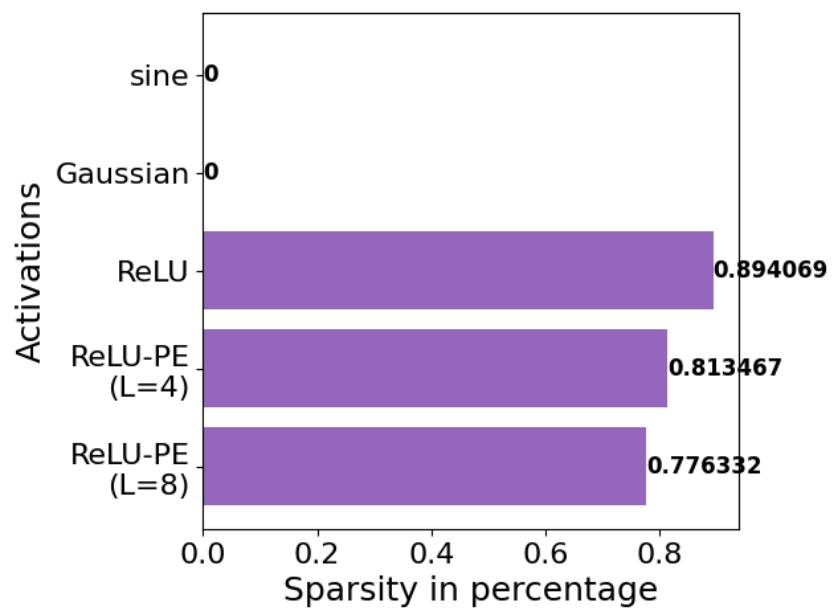


Figure 7: Sparsity of Hessian-vector product matrices in percentage for **all** layers **throughout the training** process of a 2D image reconstruction task using ESGD. The x-axis ranges from 0 to 1, with values closer to 1 indicating higher sparsity.

7.6 2D Image Reconstruction

In this section, we provide the reproducibility and implementation details for the experiments in **Sec. 4.1 of the main paper**. We also present empirical analysis for other non-traditional activations like Sine [55] and wavelet [53], see Sec. 7.6.2 below. We present more qualitative results for additional instances from the DIV2K dataset in Sec. 7.6.3.

7.6.1 Reproducibility & Implementation Details

Data. We use image from DIV2K dataset. For our experiment, we downsampled the original images, initially captured at a resolution of 512×512 pixels, to a resolution of 256×256 pixels. We used minibatch size of 512 during training. We use the *lion* image for this analysis.

Architecture. We employed a 5-layer network with a hidden layer size of 256. For Gaussian-INR, we used $\sigma = 0.05$. For sine-INR presented in (Sec.7.6.2), we used $\omega = 30$. As for wavelet-INR (Sec. 7.6.2), we employed a wavelet activation $\mathbf{x} \rightarrow \cos(\omega_0 \mathbf{x}) \exp^{-s_0^2 \mathbf{x}^2}$. We used $\omega_0 = 10$ and $s_0 = 10$.

Initialization. We used default PyTorch initialisation (Kaiming Uniform) for Gaussian-INR and wavelet-INR in all the experiments. As for sine-INR presented in (Sec.7.6.2), we used the Sitzmann et al. initialization scheme [55].

Hyperparameters. Unless specified, we maintained the default parameter values for each optimizer. For the Adam optimizer, we used a learning rate of 1×10^{-4} for both Gaussian and wavelet-INRs, and a lower learning rate of 5×10^{-5} for sine-INR. For AdaHessian(J) and AdaHessian(E), we set the learning rate to 0.15 and used a hessian power of 1. For Kronecker-based preconditioner, we used a learning rate of 0.1 and performed the inverse update every 100 iterations. For Shampoo, we used a learning rate of 0.1 and 0.01 for Gaussian/wavelet-INRs and sine-INR, respectively. As for ESGD, we chose a learning rate of 0.15. We incorporated a warmup strategy. Specifically, we implemented the equilibrated gradient preconditioning for the first 50 iterations as a warmup phase, followed by subsequent updates every N iterations, with N set to 100.

Hardware. We ran all experiments on a NVIDIA 3080 GPU with 12Gb of memory.

7.6.2 Other activations: sine & wavelet

Fig. 8 and 9 demonstrates that curvature-aware preconditioners offer significant advantage for the training convergence of non-traditional activations like sine and wavelet, compared to Adam.

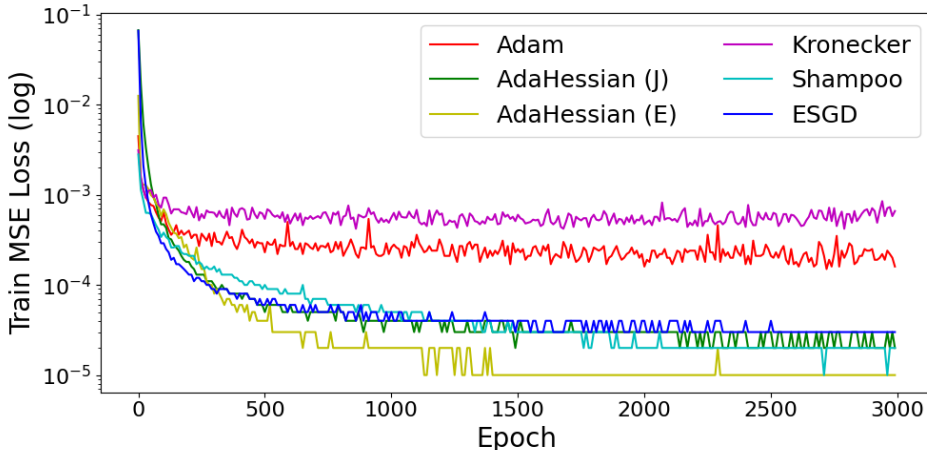


Figure 8: Comparison of training convergence using sine-INR in a 2D image reconstruction task, measured in Mean Squared Error (MSE) loss, of various preconditioners.

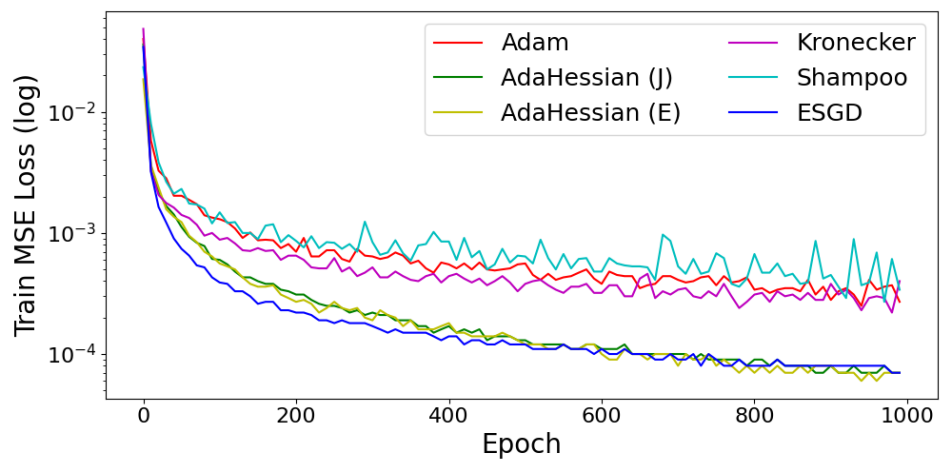
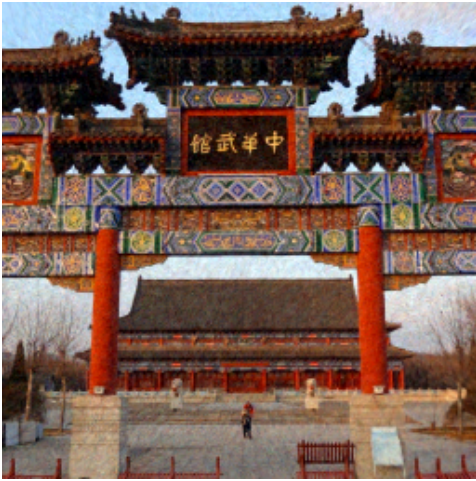


Figure 9: Comparison of training convergence using wavelet-INR in a 2D image reconstruction task, measured in Mean Squared Error (MSE) loss, of various preconditioners.

7.6.3 Additional Results

We showcase additional qualitative comparisons of Gaussian INR trained with Adam and ESGD optimizers with other instances from the DIV2K dataset.



(a)



(b)



(a)



(b)



(a)



(b)



(a)



(b)



(a)



(b)

Figure 14: Compared to INR trained with Adam (**left panel**), ESGD (**right panel**) has achieved superior fidelities at epoch 500.

7.7 3D Shape Reconstruction

In this section, we provide the reproducibility and implementation details for the experiments in **Sec. 4.2 of the main paper**. We also present the training convergence analysis for Gaussian-INR in Sec. 7.7.2. Besides, we also present comprehensive comparative analysis of different activations (ReLU/ReLU-PE and sine) when trained using Adam and SGD in Sec. 7.7.3. Finally, we showcase the qualitative results of the mesh reconstruction using the trained binary occupancy field during intermediate training in Sec. 7.7.4

7.7.1 Reproducibility & Implementation Details

Data. We use the *Armadillo*, *Bimba*, *Gargoyle* and *Dragon* from the Stanford 3D Scanning Repository ⁷. For training, we sampled 12 million 3D points – one-third of the points were sampled uniformly within the volume, and remaining two-thirds of the points were sampled near the mesh surface and perturbed with random Gaussian noise using sigma of 0.1 and 0.01, respectively.

Architecture. We employed a 5-layer network with a hidden layer size of 256. For Gaussian-INR, we used $\sigma = 0.09$. For sine-INR, we used $\omega = 30$. As ReLU-PE INR, we consider two different embedding sizes, $L = 4$ and 8.

Initialization. We use default PyTorch initialisation (Kaiming Uniform) for Gaussian-INR and ReLU/ReLU-PE INRs in all the experiments. As for sine-INR, we used the Sitzmann et al. initialization scheme [55].

Hyperparameters. Unless specified, we maintained the default parameter values for each optimizer. For the Adam optimizer, we used a learning rate of 1×10^{-4} . For AdaHessian(J) and AdaHessian(E), we set the learning rate to 0.15 and used a hessian power of 1. For Kronecker-based preconditioner, we used a learning rate of 0.1, Tikhonov regularization parameter of 0.01, and performed the inverse update every 100 iterations. For Shampoo, we used a learning rate of 0.1. As for ESGD, we chose a learning rate of 0.15. We incorporated a warmup strategy. Specifically, we implemented the equilibrated gradient preconditioning for the first 50 iterations as a warmup phase, followed by subsequent updates every N iterations, with N set to 100.

Hardware. We ran all experiments on a NVIDIA 3080 GPU with 12Gb of memory.

7.7.2 Training convergence analysis for Gaussian-INR

Fig. 15 compares the training convergence in terms of both time and epoch of various preconditioners.

⁷<http://graphics.stanford.edu/data/3Dscanrep>

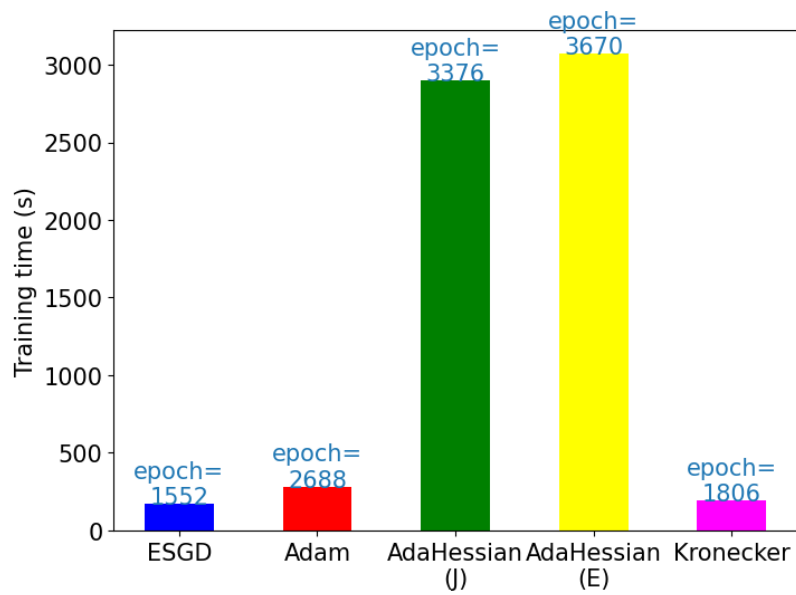


Figure 15: Comparison of training convergence of Gaussian-INR in terms of both time (**y-axis**) and number of epoch (represented in **blue** text) of various preconditioner in a 3D shape reconstruction task on an *armadillo* instance. We omit *Shampoo* from the comparison here as it did not achieve the desired convergence.

7.7.3 Comparative Analysis of other activations when trained with Adam and SGD

Fig. 16, Fig. 17, Fig. 18 demonstrate that Adam outperforms SGD for all activations such as sine, ReLU and ReLU-PE INRs. This suggest the significance of preconditioner in accelerating training convergence.

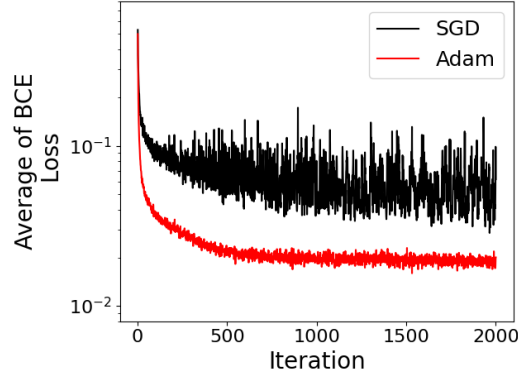


Figure 16: Comparison of training convergence in terms of Binary Cross-Entropy (BCE) loss between Adam and SGD optimizers for optimizing a 3D occupancy field of an *armadillo* using **sine-INR**.

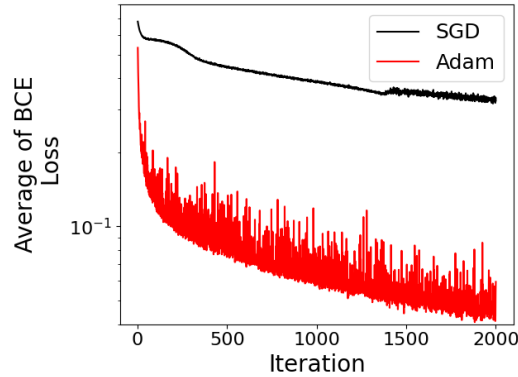


Figure 17: Comparison of training convergence in terms of Binary Cross-Entropy (BCE) loss between Adam and SGD optimizers for optimizing a 3D occupancy field of an *armadillo* using **ReLU-INR**.

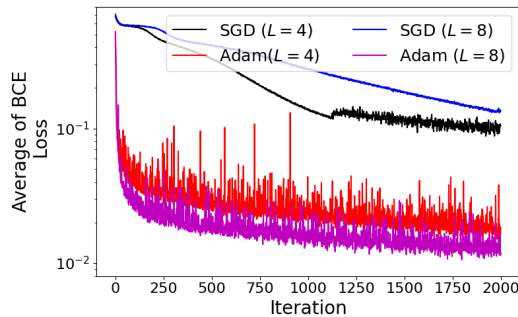


Figure 18: Comparison of training convergence in terms of Binary Cross-Entropy (BCE) loss between Adam and SGD optimizers for optimizing a 3D occupancy field of an *armadillo* using **ReLU-PE-INR** with different embedding sizes L .

7.7.4 Qualitative Results

We showcase additional qualitative comparisons of Gaussian INR trained with Adam and ESGD optimizers.

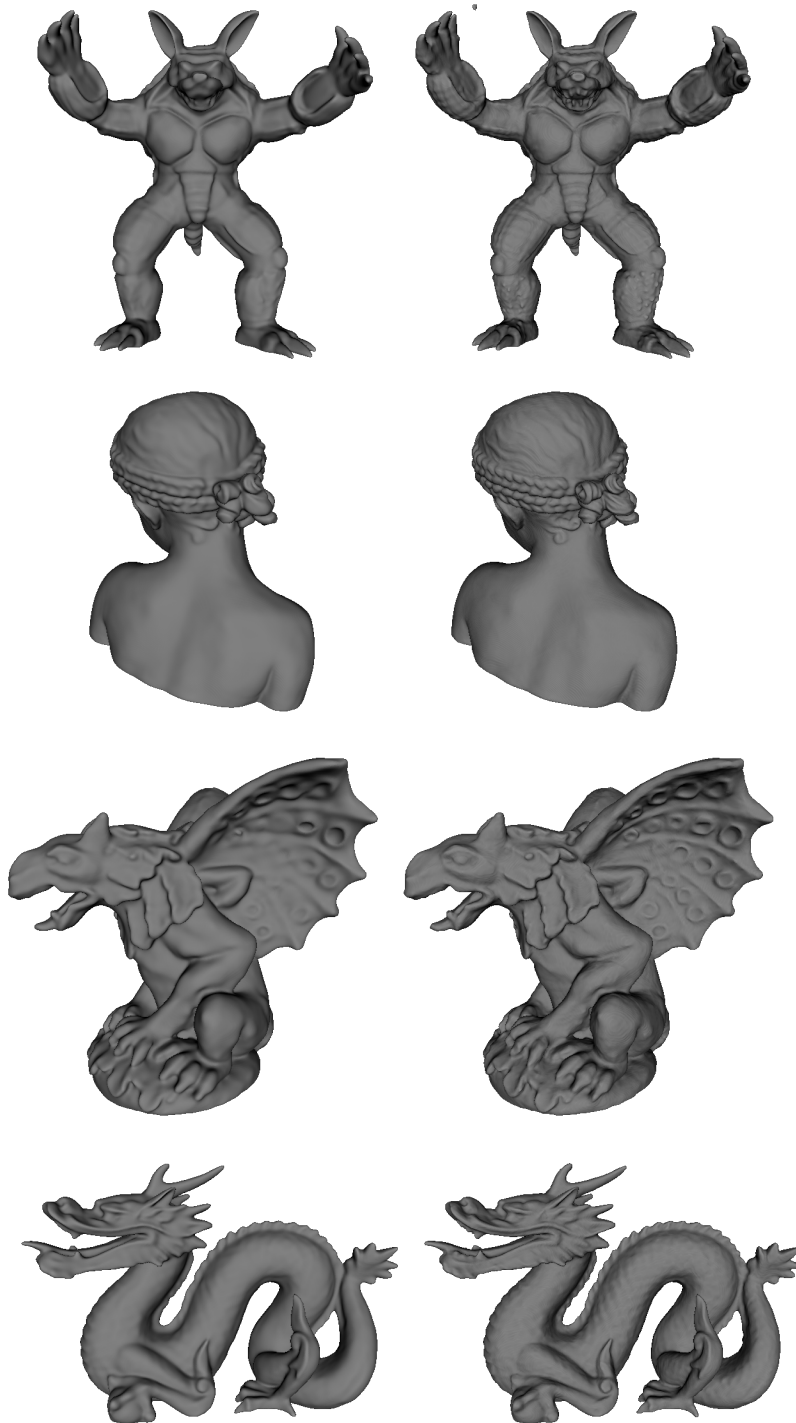


Figure 19: Compared to INR trained with Adam (**left panel**), ESGD (**right panel**) has reconstructed the shapes with significantly improved fidelities at epoch 500 (*zoomed in* for better visibility).

7.8 Neural Radiance Fields (NeRF)

In this section, we provide the reproducibility and implementation details for the experiments in **Sec. 4.3 of the main paper**. Additionally, we present a comprehensive comparative analysis of different activations (ReLU/ReLU-PE) when trained using Adam and SGD optimizers (Sec. 7.8.2), including reproducibility details. Although ReLU is not commonly used in NeRF due to spectral bias, we include the analysis for the sake of completeness. We also offer the comparison of ReLU-PE when trained with Adam and SGD (Sec. 7.8.3, along with reproducibility details. Additionally, we showcase additional results for NeRF using Gaussian-INR on other datasets, such as Blender 360 [35] and Blender Forward-Facing Dataset (BLEFF) [64] (Sec. 7.8.4. We also present the novel view synthesis results for the experiments in **Sec. 4.3 of the main paper** in Sec. 7.8.5.

7.8.1 Reproducibility & Implementation details in Sec. 4.3 of the main paper.

Training details. We resize the images to 480×640 pixels. We train all models for 200K iterations. As in [15, 29], we trained a single model without additional hierarchical sampling. During each optimization step, we randomly sample 2048 pixel rays, and integrated 128 points along each ray.

Hyperparameters. Unless specified, we maintained the default parameter values for each optimizer. For the Adam optimizer, we used a learning rate of 1×10^{-4} and customised step learning rate schedule. For the ESGD optimizer, we started with a learning rate of 1 and exponentially decayed it to 0.01.

Hardware. The experiments are ran on a Tesla v100 GPU.

7.8.2 Comparative Analysis of other activations when trained with Adam and SGD

Training details. We resize the images to 480×640 pixels. We train all models for 200k iterations. Following the approach in [15, 29], we trained a single model without additional hierarchical sampling. For each optimization step, we randomly sample 2048 pixel rays, and integrated 128 points along each ray.

Hyperparameters. For the Adam optimizer, we started with a learning rate of 1×10^{-3} and exponentially decayed it to 1×10^{-4} . As for the SGD optimizer, we started from a learning rate of 2 and exponentially decayed it to 1×10^{-4} .

Hardware. The experiments for ReLU-NeRF and ReLU-PE-NeRF are ran on a Tesla v100 GPU and an NVIDIA a100 GPU, respectively.

Results. Table 2 and 3 are correlated with the analysis in Sec. 7.7.3. This again suggest the significance of preconditioner in accelerating training convergence.

Scene	Train PSNR \uparrow	Adam Test			Train PSNR \uparrow	SGD Test		
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
fern	24.35	23.51	0.70	0.37	20.93	20.57	0.55	0.72
flower	26.29	24.59	0.74	0.21	21.37	20.73	0.52	0.59
orchids	20.84	19.28	0.56	0.32	17.38	16.64	0.34	0.68
fortress	28.65	28.09	0.75	0.24	23.80	23.94	0.53	0.64
leaves	19.98	19.04	0.56	0.38	17.11	16.61	0.33	0.77

Table 2: Quantitative results of for **ReLU NeRF** on instances from the LLFF dataset [34]. As SGD was not able to achieve the same training convergence as Adam, we presented the result based on a *fixed number of iterations* (200k) in this experiment. Note that the report time solely refers to the optimization update step. Training convergence is determined based on the training curve displayed on TensorBoard, with a smoothing factor of 0.8 – 0.9 applied for better visualization.

Scene	Adam				SGD			
	Train PSNR \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Train PSNR \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
fern	26.02	23.27	0.71	0.29	20.27	18.59	0.45	0.81
flower	27.54	22.89	0.65	0.26	20.53	17.84	0.42	0.70
orchids	22.26	16.97	0.50	0.32	17.45	14.89	0.28	0.66
fortress	30.67	26.02	0.78	0.19	9.00	9.47	0.46	0.78
leaves	18.81	14.13	0.24	0.57	15.44	14.36	0.20	0.93

Table 3: Quantitative results of for **ReLU-PE NeRF** on instances from the LLFF dataset [34]. As SGD was not able to achieve the same training convergence as Adam, we presented the result based on a *fixed number of iterations* (200k) in this experiment. Note that the report time solely refers to the optimization update step. Training convergence is determined based on the training curve displayed on TensorBoard, with a smoothing factor of 0.8 – 0.9 applied for better visualization.

7.8.3 Activation matters: Comparison of ReLU-PE trained on Adam and ESGD

Training details. We resize the images to 480×640 pixels. We train all models for 200k iterations. Following the approach in [15, 29], we trained a single model without additional hierarchical sampling. For each optimization step, we randomly sample 2048 pixel rays, and integrated 128 points along each ray.

Hyperparameters. For the Adam optimizer, we started with a learning rate of 1×10^{-3} and exponentially decayed it to 1×10^{-4} . As for the ESGD optimizer, we started from a learning rate of 0.1 and exponentially decayed it to 5×10^{-3} .

Hardware. The experiments are ran on a NVIDIA a100 GPU with 48Gb of memory.

Results. Table 4 shows that ReLU-PE trained with Adam significantly outperform the one trained with ESGD, which indicates that curvature-aware preconditioner does not offer any significant advantage to ReLU-PE INR on NeRF.

Scene	Train PSNR \uparrow	Adam Test			Train PSNR \uparrow	ESGD Test		
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
fern	26.02	23.27	0.71	0.29	22.63	21.84	0.58	0.55
flower	27.54	22.89	0.65	0.26	22.61	20.61	0.49	0.55
orchids	22.26	16.97	0.50	0.32	18.87	17.43	0.43	0.49
fortress	30.67	26.02	0.78	0.19	26.63	23.05	0.62	0.43
leaves	18.81	14.13	0.24	0.57	16.73	14.43	0.24	0.69

Table 4: Quantitative results of for **ReLU-PE NeRF** on instances from the LLFF dataset [34]. As ESGD was not able to achieve the same training convergence as Adam, we presented the result based on a *fixed number of iterations* (200k) in this experiment. Note that the report time solely refers to the optimization update step. Training convergence is determined based on the training curve displayed on TensorBoard, with a smoothing factor of 0.8 – 0.9 applied for better visualization.

7.8.4 Additional Results for Gaussian-INR on NeRF

Training details. We resize the images to 400×400 pixels. We train all models for 200K iterations. Following the approach in [15, 29], we trained a single model without additional hierarchical sampling. For each optimization step, we randomly sample 2048 pixel rays, and integrated 128 points along each ray.

Hyperparameters. Unless specified, we maintained the default parameter values for each optimizer. For the Adam optimizer, we started with a learning rate of 1×10^{-4} and exponentially decayed it to 1×5^{-6} . For the ESGD optimizer, we started with a learning rate of 1 and exponentially decayed it to 0.01.

Hardware. The experiments are ran on a Tesla v100 GPU.

Results. Table 5 shows that Gaussian-INR trained with ESGD outperforms the one trained with Adam.

Scene	Train PSNR	Iteration ↓	Time (s) ↓	Adam Test			Iteration ↓	Time (s) ↓	ESGD Test		
				PSNR ↑	SSIM ↑	LPIPS ↓			PSNR ↑	SSIM ↑	LPIPS ↓
lego	26.44	200K	435.09	25.37	0.87	0.12	80K	179.59	25.94	0.89	0.09
chair	31.33	200K	432.64	33.82	0.97	0.04	100K	227.33	34.10	0.98	0.03
drums	25.16	180K	411.10	24.62	0.88	0.13	100K	239.46	24.85	0.89	0.12
hotdog	25.2	200K	417.46	31.82	0.95	0.07	150K	324.04	32.86	0.96	0.06
figus	26.53	200K	481.04	25.25	0.91	0.08	80K	183.90	26.95	0.93	0.06

Table 5: Quantitative results of **Gaussian-NeRF** on instances from the **360 BLENDER** dataset [35]. Note that the report time solely refers to the optimization update step. Training convergence is determined based on the training curve displayed on TensorBoard, with a smoothing factor of 0.8 – 0.9 applied for better visualization.

Training details. We resize the images to 480×640 pixels. We train all models for 200K iterations. Following the approach in [15, 29], we trained a single model without additional hierarchical sampling. For each optimization step, we randomly sample 2048 pixel rays, and integrated 128 points along each ray.

Hyperparameters. Unless specified, we maintained the default parameter values for each optimizer. For the Adam optimizer, we started with a learning rate of 1×10^{-4} and a customised step learning rate scheduler. For the ESGD optimizer, we started with a learning rate of 1 and a customised step learning rate scheduler.

Hardware. The experiments are ran on a Tesla v100 GPU.

Results. Table 6 shows that Gaussian-INR trained with ESGD is on average faster than to the one trained with Adam.

Scene	Train PSNR	Adam					ESGD				
		Iteration ↓	Time (s) ↓	Test			Iteration ↓	Time (s) ↓	Test		
				PSNR ↑	SSIM ↑	LPIPS ↓			PSNR ↑	SSIM ↑	LPIPS ↓
balls	37.78	200K	522.56	37.11	0.93	0.08	150K	392.65	37.15	0.93	0.09
deer	47.50	150K	359.30	46.73	0.99	0.02	150K	360.76	46.53	0.99	0.02
chair	36.38	180K	492.08	36.27	0.90	0.17	180K	485.32	36.02	0.89	0.18
root	38.58	200K	501.13	37.4	0.98	0.03	150K	394.45	37.38	0.98	0.03
roundtable	48.25	160K	367.38	47.36	1.00	0.01	150K	341.11	46.90	0.99	0.01

Table 6: Quantitative results of **Gaussian-NeRF** on instances from the **BLEFF** dataset [64]. Note that the report time solely refers to the optimization update step. Training convergence is determined based on the training curve displayed on TensorBoard, with a smoothing factor of 0.8 – 0.9 applied for better visualization.

7.8.5 Novel View Synthesis Results for Sec. 4.3 in the main paper



(a) 1st View



(b) 50th View



(c) 1st View



(d) 50th View



(e) 1st View



(f) 50th View

Figure 20: Novel view synthesis results using Gaussian-INR trained with ESGD.