





Visualization and Analysis of the Loss Landscape in Graph Neural Networks

Samir Moustafa^{1,2} , Lorenz Kummer^{1,2} , Simon Fetzl¹, Nils M. Kriege^{1,3} , and Wilfried N. Gansterer¹ 

¹ Faculty of Computer Science, University of Vienna, Vienna, Austria

² UniVie Doctoral School Computer Science, University of Vienna, Vienna, Austria

³ Research Network Data Science, University of Vienna, Vienna, Austria
{samir.moustafa, lorenz.kummer, simonf93, nils.kriege,
wilfried.gansterer}@univie.ac.at

Abstract. Graph Neural Networks (GNNs) are powerful models for graph-structured data, with broad applications. However, the interplay between GNN parameter optimization, expressivity, and generalization remains poorly understood. We address this by introducing an efficient learnable dimensionality reduction method for visualizing GNN loss landscapes, and by analyzing the effects of over-smoothing, jumping knowledge, quantization, sparsification, and preconditioner on GNN optimization. Our learnable projection method surpasses the state-of-the-art PCA-based approach, enabling accurate reconstruction of high-dimensional parameters with lower memory usage. We further show that architecture, sparsification, and optimizer’s preconditioning significantly impact the GNN optimization landscape and their training process and final prediction performance. These insights contribute to developing more efficient designs of GNN architectures and training strategies.

Keywords: Graph Neural Networks · Loss Landscape .

1 Introduction

Graph Neural Networks (GNNs) are tailored for graph-structured data and excel in tasks like network analysis, molecular property prediction, and recommendation systems [26]. Compared to Deep Neural Networks (DNNs), they have fewer parameters but higher computational costs due to large input graphs [18, 31]. The relationship between GNNs optimization dynamics and graph structure, architectural design, or numerical precision is poorly understood [26].

GNN research has focused chiefly on theoretical expressivity—the ability to distinguish non-isomorphic graphs or structurally distinct nodes [17]. While deeper networks are more expressive, over-smoothing [26] limits practical depth by making node embeddings indistinguishable. Although theoretical expressivity results guarantee the existence of suitable parameters, standard optimization techniques may fail to recover

This is a preprint of the work accepted to the International Conference on Artificial Neural Networks Workshops, LNCS 16072, pp. 1–13. doi.org/10.1007/978-3-032-04552-2_9

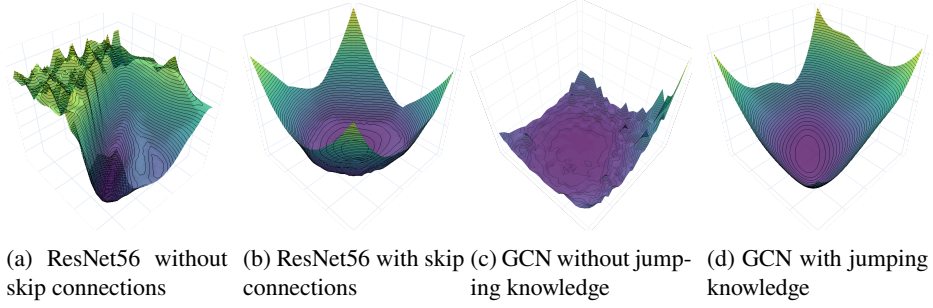


Fig. 1: 3D loss landscapes of ResNet56 on CIFAR-10 [13] and GCN on Cora [5], with and without architectural modifications.

them, revealing a gap between theory and practice. The interplay among expressivity, generalization, and optimization remains insufficiently understood [17].

Analyzing loss landscapes and training trajectories can illuminate the gap between theoretical and practical GNN behavior, yet it remains underexplored. The loss landscape maps parameters to loss values, offering a geometric view of model behavior, and is typically highly non-convex and complex, particularly under varying data distributions [14]. Understanding the loss landscape yields insights into optimization, generalization, and the stability of learned representations [1]. Loss landscape analysis for GNNs remains unexplored, hindered by challenges inherent to graph-structured data, message passing, and architectural diversity [26]. Recent advancements in GNNs training techniques, aimed at improving generalization [28], reducing inference time [29, 31], or optimizing parameters [11], have further complicated the understanding of the GNN loss landscape by obscuring the relationship between model parameters and loss.

This paper addresses this gap and sheds light on the importance and methods of analyzing and visualizing the loss landscape in GNNs.

Contribution of the Paper: Our contributions are threefold: **(1)** We propose a novel learnable projection technique for loss landscape visualization, mapping high-dimensional parameters into 2D and restoring them, while supporting batching for controlling memory usage. **(2)** We provide a comprehensive analysis of GNN training techniques and their effects on the loss landscape, considering over-smoothing [26], jumping knowledge [28], quantization [18], sparsification [29], and optimizer preconditioning [11].

Preliminaries: A graph G is a pair (V, E) , where V is a finite set of nodes and E is a finite set of edges. Let A denote the adjacency matrix of G with order $|V|$. Each node $v_i \in V$ has an associated feature vector $x_i \in \mathbb{R}^f$, stacked into a matrix $X \in \mathbb{R}^{|V| \times f}$. We consider a classification setting, assuming each node has an expected class, and in total there are C classes, represented by $Y \in \mathbb{R}^{|V| \times C}$.

A GNN is parameterized by θ , defining a function $\mathcal{F}_\theta : (A, X) \rightarrow Y$. GNNs employ message passing, where each node aggregates neighbor embeddings, implemented via multiplication of the adjacency matrix A and feature matrix X . At the k^{th} layer, $X^{(k)}$ is computed as $AX^{(k-1)}\theta^{(k)}$, where $X^{(k-1)}$ is the embedding from layer $k - 1$ [26].

Table 1: Standard aggregation and update functions in the k^{th} layer of three GNN architectures: GCN, GAT, and GIN. These functions gather (aggregation) and transform (update) operations by collecting information from each node’s neighborhood.

Architecture	Aggregation Function	Update Function
GCN [12]	$h_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{d_i d_j}} x_j^{(k-1)}$	$x_i^{(k)} = \text{ReLU} \left(\theta^{(k)} h_i^{(k)} \right)$
GAT [25]	$h_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \alpha_{i,j}^{(k)} x_j^{(k-1)}$	$x_i^{(k)} = \text{ReLU} \left(\theta^{(k)} h_i^{(k)} \right)$
GIN [27]	$h_i^{(k)} = \left(1 + \epsilon^{(k)} \right) x_i^{(k-1)} + \sum_{j \in \mathcal{N}(i)} x_j^{(k-1)}$	$x_i^{(k)} = \text{MLP}^{(k)} \left(h_i^{(k)} \right)$

Training of a GNN aims to adjust θ so that $\mathcal{F}_\theta(X)$, given G , approximates the target Y . During training, the loss function ε measures the discrepancy between the expected and actual outputs; we use multi-class cross-entropy, though any differentiable loss is applicable. During training, gradients of $\varepsilon(\mathcal{F}_\theta(A, X), Y)$ with respect to θ are computed to update parameters and minimize the error.

A loss landscape visualization depicts how the loss varies with the model parameters. Ideally, a well-fitted model yields a smooth, convex landscape, as shown in Figures 1b and 1d.

GNN Architectures: This work focuses on three fundamental GNN architectures selected for their pivotal role in the evolution and comprehension of GNNs. Graph Convolution Network (GCN) [12] applies convolution in the spectral domain to capture local graph structures invariant to graph isomorphisms. Graph Attention Network (GAT) [25] utilizes an attention mechanism that allows nodes to weigh their neighbors’ influence dynamically, enhancing model flexibility and performance on graph-structured data. Graph Isomorphism Network (GIN) [27] is more powerful in distinguishing different graph structures than GCN and GAT by using a multilayer perceptron (MLP) and a learnable parameter for feature updating, achieving theoretical equivalence with the Weisfeiler-Lehman test. Table 1 presents the mathematical formulations for GCN, GIN, and GAT, detailing aggregation and update functions based on hidden states h_i , output embeddings x_i , neighborhood structures $\mathcal{N}(i)$, and the $\alpha_{ij} = \frac{e_{i,j}}{\sum_{k \in \mathcal{N}(i)} e_{i,k}}$, such that $e_{i,j} = e^{\text{LeakyReLU}(a^\top [\theta h_i \| \theta h_j])}$, and a is the attention learnable parameter.

Problem Definition: The loss landscape $L(\theta)$ depends on the architecture, input data, target output, parameters, and error function. Visualizing $L(\theta)$ helps reveal how the error adapts to parameter changes, where $L(\theta)$ reduces to $\varepsilon(\mathcal{F}_\theta(A, X), Y)$ with other factors fixed. This landscape captures the model’s error across parameter configurations. Optimization seeks parameters minimizing ε , but is hindered by non-convexity and local minima. Due to the high dimensionality of θ , a full mathematical analysis of $L(\theta)$ is typically intractable.

2 Related Work

In this section, we delve into the definition of the loss landscape, the limitations of GNNs, explore various methods that have been adapted specifically for GNNs, and discuss the prevalent optimization challenges associated with them.

Loss Landscape Visualization: Despite the non-convexity of $L(\theta)$, DNNs often train efficiently by converging to local minima. This raises questions [6] about whether training avoids local minima or if saddle points dominate optimization dynamics. [6] study $L(\lambda, \theta_i + (1-\lambda)\theta^*)$ for $\lambda \in \mathbb{R}$, where θ_i and θ^* are parameters from epoch i and the final training epoch, respectively. They observe elevated loss regions between solutions, consistent with [10]. Visualization techniques include barycentric and bilinear interpolation, and projections along random or PCA directions [14, 16]. Moreover, [10] analyzed how optimization trajectories change when the algorithm switches mid-training, typically after loss stagnation. Post-switch trajectories consistently diverge, suggesting optimizers follow distinct paths at saddle points. Loss landscapes remain similar across runs with the same optimizer but different initializations. Despite differing final parameters, loss and accuracy are often comparable—a result also observed in [23], where varying the learning rate led to multiple equally performant optima [7].

Over-smoothing and Jumping Knowledge: A key challenge in GNN training is over-smoothing, where node embeddings become increasingly similar with depth, diminishing node identity and discriminative power [26]. Over-smoothing occurs as the row of $X^{(l)}$ converges to the same vector as $l \rightarrow \infty$ [26]. This can be quantified via similarity or distance measures between rows of $X^{(l)}$.

Jumping Knowledge (JK) mitigates over-smoothing by allowing node representations to aggregate multi-distance neighborhood information across layers. The final representation is given by $X^{(l)} = \phi(X^{(0)}, X^{(1)}, \dots, X^{(l-1)})$, where ϕ combines outputs from all previous layers. JK parallels residual connections in CNNs [22], enhancing trainability and addressing vanishing gradients.

Quantization and Sparsification: Quantization lowers memory, computation, and power demands by converting weights or activations to lower-precision formats, enabling efficient deployment at the cost of approximation errors [4]. [31] proposes a GNN-specific method that learns per-node bit-widths based on aggregation values, capturing topological variance. This adds complexity, requiring group-specific learning rates.

Sparsification reduces GNN memory and computation by constructing sparse graphs that retain essential structure. GraphSAINT [29] samples subgraphs to preserve batch connectivity, using normalization and diverse strategies to control bias and variance. Its effectiveness depends critically on sampling quality [26].

Optimization Challenges and Strategies: DNN optimization has shifted focus from local minima, which may support generalization [7], to saddle points as the main challenge [10, 20]. Low-error solutions exhibit parameter space symmetry, with local minima often matching global optima in fully connected networks [1, 19]. Under conditions like identity mappings or over-parameterization, stochastic gradient descent (SGD) can reach

global optima [15, 24]. In high dimensions, saddle points prevail, and methods like natural gradient descent (NGD) are effective for escaping them [2]. [11] introduces an information-geometric optimization method for GNNs using NGD, with the Fisher Information Matrix approximated via KFAC to avoid second-order derivatives. This improves efficiency and outperforms ADAM and SGD.

Overall, the unique structural characteristics of GNNs, alongside their phenomena and methodologies, present unexplored dimensions in understanding how they shape the loss landscape and govern optimization trajectories.

3 Methodology

Consider a k -layer graph neural network with parameters $\theta = (\theta^0, \theta^1, \dots, \theta^k)$, where each θ^i is a vector or matrix. Let $\Theta \in \mathbb{R}^d$ denote the flattened concatenation of all elements in θ , with d as the parameter dimension. The optimized flattened concatenated parameters are denoted $\Theta^* \in \mathbb{R}^d$.

Dimensionality Reduction for Visualization: To visualize the high-dimensional parameter space \mathbb{R}^d , dimensionality reduction to 2D is applied using two directions $b^1, b^2 \in \mathbb{R}^d$. The mapping from 2D to \mathbb{R}^d is defined by the function

$$M : \mathbb{R}^2 \rightarrow \mathbb{R}^d \text{ with } (x, y) \mapsto \Theta^* + x \cdot b^1 + y \cdot b^2, \quad (1)$$

where x and y are the visualization's axes. The visualization is then created via

$$\Psi : (x, y) \mapsto L(M(x, y)) \quad (2)$$

The directions b^1, b^2 can be selected via various methods; one common approach is random initialization [14]. Each direction vector (b_0, b_1, \dots, b_d) is partitioned into segments $\hat{b}_l = (b_m, \dots, b_k)$, corresponding to specific network parameters (e.g., weight matrices or bias vectors). These segments are normalized using their magnitude and that of the corresponding optimal parameters $\hat{\Theta}_l^* = (\Theta_m^*, \dots, \Theta_k^*)$, ensuring each direction is scaled relative to its associated optimized values. This scaling ensures that each direction vector segment matches the magnitude of its corresponding optimal parameter segment. The normalization is formalized as:

$$\hat{b}'_j := \frac{\hat{b}_l}{\|\hat{b}_l\|} \|\hat{\Theta}_l^*\| \quad (3)$$

Normalized random directions allow a more balanced and representative visualization of the high-dimensional parameters.

Visualizing the Optimizer Trajectory: To visualize optimizer trajectories, directions b^1 and b^2 are selected to align with the subspace spanned by parameters Θ_i , $0 \leq i \leq n$, across all n epochs. Using Θ^* as the origin, differences $D_i = \Theta_i - \Theta^*$ form a matrix $D \in \mathbb{R}^{n \times d}$. PCA on the covariance of D yields b^1 and b^2 [16]. However, computing the covariance matrix and its eigenvectors is memory-intensive for large networks [14], despite $n \ll d$ reducing the impact of D itself.

Projection into Visualization Space: To visualize trajectories, a point $p \in \mathbb{R}^d$ must be projected into the 2D visualization space (x, y) by solving:

$$p = \Theta^* + x \cdot b^1 + y \cdot b^2 + r \quad (4)$$

where r is the reconstruction error vector, which should be minimized. Basis pairs can be compared via their reconstruction error r , whereby a smaller r indicates a preferable basis. If the basis vectors b^1, b^2 are orthonormal [14], the dot product can be used to calculate the coordinates as $x = \langle p - \Theta^*, b^1 \rangle$ and $y = \langle p - \Theta^*, b^2 \rangle$. In general, b^1, b^2 are not orthonormal and not orthogonal. Therefore, an underdetermined linear equation system needs to be solved.

$$(b^1 \ b^2) \begin{pmatrix} x \\ y \end{pmatrix} = p - \Theta^* \quad (5)$$

The solution only approximates p ; the true loss may differ from $\Psi(x, y)$ in equation (2). Hence, 2D trajectory visualizations such as contour plots are preferred.

Learnable Projection: We propose a novel dimensionality reduction method, the *learnable projection model*, defined by a matrix $P \in \mathbb{R}^{2 \times d}$, where $P = (b^1 \ b^2)^\top$. P defines the projection unambiguously, encodes high-dimensional parameters into 2D, and decodes them back. Given input matrix D from Section 3, each point is encoded as $z_i = D_i P^\top \in \mathbb{R}^2$, assuming orthonormal b^1, b^2 (see Section 3). Decoding is performed via $z_i P$, equivalent to the mapping function M in equation (1).

The learnable projection model is trained to minimize the Euclidean distance between D and zP via Mean Squared Error (MSE), thereby minimizing the reconstruction error r from equation (4) and yielding a 2D basis with low projection error. For optimization, the problem is formulated as:

$$\min_{b^1, b^2 \in \mathbb{R}^d} \text{MSE}(D, DP^\top P), \quad D \in \mathbb{R}^{n \times d}, \quad P = (b^1 \ b^2)^\top \in \mathbb{R}^{2 \times d}. \quad (6)$$

Simplifying the objective yields an upper bound: minimizing $\|I_d - P^\top P\|_2^2$, where I_d is the $d \times d$ identity matrix. This reformulation shifts the focus to minimizing the squared spectral norm of $I_d - P^\top P$. This formulation relates to the orthogonal Procrustes problem, though it omits explicit dependence on Θ^* and Θ_i encoded in D . The ℓ_1 reconstruction error $\|r\|_1$ enables comparison with PCA-based methods, as both aim to minimize the same error term from equation (4), which highlights the effectiveness of our learnable projection relative to PCA, despite both sharing the same objective.

Theoretically, both our approach and PCA aim to minimize the *reconstruction error*, as defined in equation 7 from [3]:

$$\min_{V_q} \sum_{i=1}^N \left\| (x_i - \bar{x} - V_q V_q^\top (x_i - \bar{x})) \right\|^2 \quad (7)$$

where x_i are the data samples, \bar{x} is the mean and V_q is a matrix consisting of q orthonormal columns. Unlike PCA, which computes V_q via eigenvectors, our approach uses gradient-based optimization, avoiding covariance matrix computation and associated memory overhead. Grouping rows of D into batches enables the learnable projection to train on D_i in a batched manner, reducing memory usage even further.

Table 2: Characteristics of node classification datasets across different domains.

	Cora [5]	Citeseer [5]	Pubmed [21]	OGB-Arxiv [8]	OGB-MAG [8]
Number of Nodes	2,708	3,327	19,717	169,343	1,939,743
Number of Edges	10,556	9,104	88,648	1,166,243	21,111,007
Features per Node	1,433	3,703	500	128	128
Number of Classes	7	6	3	23	349

4 Experimental Setup

This section summarizes the hardware, datasets, and trajectory visualization process, including PCA-based projection [14] and training of our learnable projection method. Experiments⁴ were conducted on an Intel Xeon Gold 6130 (64 cores, 256GB RAM, x86_64 architecture). For evaluation and visualization, we used standard benchmark datasets, summarized in Table 2.

Projection Methods Pipelines and Memory Requirements: During training, GNN parameters from each epoch are stored to construct the matrix D and train the learnable projection (Section 3). Models are trained for up to 1000 epochs with early stopping, using ADAM (learning rate 0.01, weight decay 0.002). To ensure fair runtime comparison with PCA, our method is evaluated on CPU only, despite being GPU capable. PCA is computed via LOBPCG on CPU, and its memory cost grows with model size due to the covariance matrix.

Both methods require the matrix $D \in \mathbb{R}^{n \times d}$, but PCA additionally computes a $d \times d$ covariance matrix, while the learnable projection supports batching with batch size B , operating on $B \times d$ subsets of D . Thus, PCA incurs $\frac{d}{B}$ times the amount of memory compared to the learnable projection.

5 Evaluation and Analysis

Reconstruction error is computed between ground-truth GNN parameters and their reconstructions from the 2D landscape (x, y) . It serves as a proxy for how well the learned directions capture training dynamics. Given the difficulty of quantitatively evaluating visualization quality, reconstruction error offers an intuitive assessment metric.

Figure 2 shows mean reconstruction error (y-axis) and runtime (x-axis) for PCA-based projection [14] and our learnable projection, averaged over ten runs. For the GCN model, averaging over the number of layers, the learnable projection consistently lowers reconstruction error but alters runtime: on Cora the mean error drops from 297.67 to 215.46, on CiteSeer the error drops from 640.89 to 489.6, and on PubMed the error drops from 109.68 to 83.12, with GIN showing analogous trends. In contrast, for GAT on Cora, PCA yields slightly better reconstruction error when averaging over the number of layers (1887.09 vs. 2141.08), while the learnable projection reduces runtime from

⁴ Code is available at <https://github.com/SamirMoustafa/torch-loss-landscape>

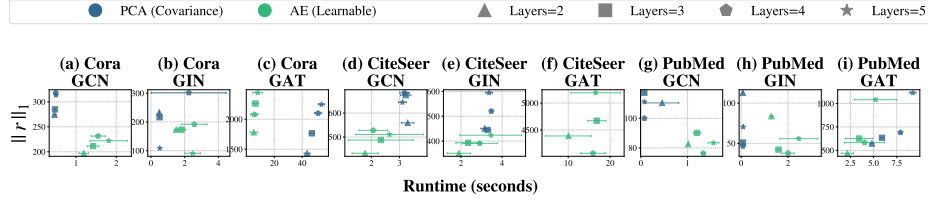


Fig. 2: Reconstruction errors $\|r\|_1$ versus the time taken to compute the projection directions for GNN architectures, respectively, with 2, 3, 4, and 5 layers on different datasets. Points represent the mean of 10 runs, and the line is the standard deviation.

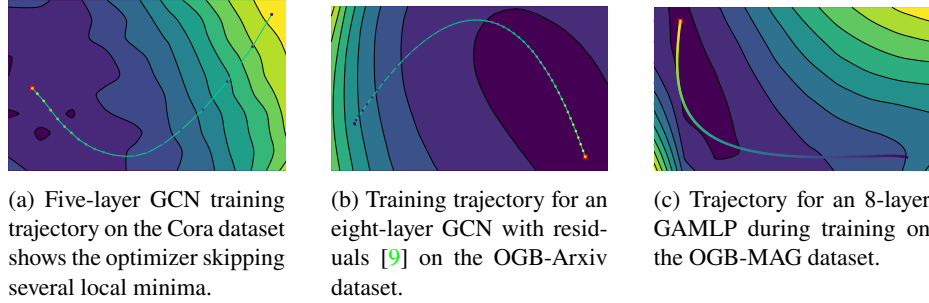


Fig. 3: Learning trajectories of GNNs shown via the learnable projection method.

48.94 seconds to 6.45 seconds (7.6 \times faster). PCA results for GAT on CiteSeer were omitted (exceeded 256 GB memory) due to the dense graph and GAT’s large parameter count.

Impact of Skip Connections and JK on Loss Landscapes: Figure 1 presents 3D loss landscape visualizations to illustrate the impact of architectural changes. For ResNet56 on CIFAR-10 [22], the landscape without skip connections (Figure 1a) exhibits multiple local minima and saddle points. In contrast, adding skip connections (Figure 1b) smooths the landscape, mitigating vanishing gradients and facilitating optimization. Similarly, Jumping Knowledge (JK), analogous to skip connections in GNNs, alters the loss landscape. As shown in Figure 1d, JK appears to reduce the non-convexity of the GCN loss surface, potentially leading to a smoother optimization landscape. Without JK (Figure 1c), the surface flattens and is noisy due to over-smoothing [26], obscuring the minimum.

Figure 3a illustrates the optimizer trajectory toward Θ^* (marked \times), the final optimizer point intended to represent the optimum. However, near this identified final point, at least three local minima exist. This suboptimality arises from over-smoothing: many nodes aggregate similar information during propagation, causing their features to become indistinguishable.

Consequently, the final GNN layer produces identical predictions for these nodes. As a result, different nodes incur the same loss value, creating a flat loss surface near the

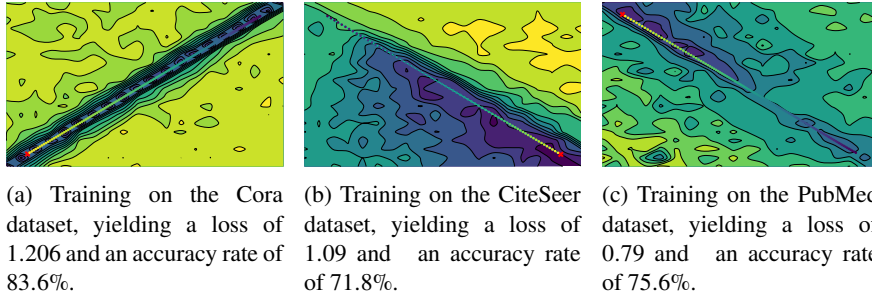


Fig. 4: Visualizing the training trajectory of a quantized GAT architecture across three different datasets, using our learnable projection method. The trajectory path appears as a straight line due to the bending in the loss landscape and its progression from a significantly far non-optimal point to the optimal point.

non-optimal point. This flat region disrupts optimization by having identical gradients, preventing further progress toward the true optimum.

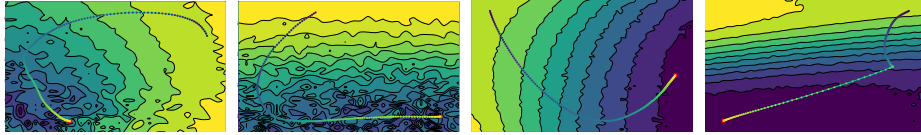
Table 3 shows that increasing GCN depth degrades accuracy and increases loss, with a sharp decline beyond four layers. This supports the observation that deeper architectures induce flatter, non-convex loss landscapes. In the 6-layer case, the optimizer fails to converge to a local minimum, hindering regularization and reducing accuracy. Jumping Knowledge variants can mitigate this degradation. [9] introduces residual connections with an error correlation mechanism that propagates residuals from training to test data. As shown in Figure 3b, their 8-layer GCN on OGB-Arxiv achieves 73.5% test accuracy with 0.727 loss, converging to a minimum without flat surrounding regions.

Figure 3c visualizes the loss trajectory of an 8-layer Graph Attention Multilayer Perceptron (GAMLP) [30] trained on OGB-MAG. Despite the model and dataset size, the learnable projection effectively visualizes the training process.

Quantization and Sparsification: While prior work has analyzed the loss landscapes of quantized DNNs across various architectures, no study has, to our knowledge, examined this for quantized GNNs. Addressing this gap, we adopt the state-of-the-art mixed-precision GNN quantization method from [31], which maintains performance across diverse datasets and tasks at both node and graph levels. Figure 4 visualizes the loss landscape of a quantized GAT, chosen for its sensitivity to perturbations in attention parameters. Quantization introduces approximation noise, increasing landscape ruggedness and local minima. Despite this, the optimizer achieves high accuracy, attributed to the method in [31], which uses group-wise learning rates within ADAM. However, this approach requires extensive hyperparameter tuning. Figures 4a and 4c show that, despite bypassing local minima, the model converges to accuracy-maximizing solutions.

Table 3: Mean and standard deviation of loss and accuracy over 10 runs for GCN on Cora.

Layers	Loss ↓	Acc. (%) ↑
2	0.86 ± 0.2	81.02 ± 0.9
3	0.85 ± 0.2	81.5 ± 1.0
4	0.86 ± 0.2	80.1 ± 0.9
5	0.88 ± 0.2	79.1 ± 1.4
6	1.33 ± 0.5	77.9 ± 2.0



(a) Native and complete graph structure is utilized. (b) Sparsification randomly omits parts of the graph. (c) Without preconditioner, and the final loss value is 0.31. (d) With preconditioner after the 50-th epoch, loss is 0.03.

Fig. 5: Training trajectories for GIN (a, b) and GCN (c, d) architecture during the training on Cora dataset, illustrating the effects of sparsification and preconditioning techniques.

Concerning sparsification, as noted in Section 2, local minima do not inherently hinder DNN training. In GNNs, graph sparsification can act as regularization, improving generalization and training efficiency by simplifying the graph and encouraging convergence to diverse local minima. We use the GIN architecture, known for matching the expressivity of the Weisfeiler-Lehman test [27], as its performance is highly sensitive to structural changes like sparsification, which can notably increase loss landscape ruggedness. Figure 5b shows the impact of sparsification on the GIN loss landscape for the Cora dataset. Using GraphSAINT [29], we compare sparsified and non-sparsified training. Sparsification increases the number of local minima, resulting in a more rugged loss landscape.

KFAC Preconditioning: As noted in Section 2, the KFAC preconditioner enhances optimization by adapting step sizes to local curvature, increasing the likelihood of reaching better minima in fewer iterations. This can lower final loss and improve post-training performance. Figure 5d illustrates the impact of applying the KFAC preconditioner during the final 50 training epochs, compared to early application. This highlights its influence on the optimizer’s trajectory. Notably, the pre-KFAC trajectory resembles that in Figure 5c, indicating consistent behavior prior to preconditioning.

6 Conclusion

We introduced a learnable projection method for visualizing GNN optimization trajectories, which outperformed the PCA-based approach in reconstruction error across most experiments. The negligible runtime increase is offset by lower memory usage, which is proportional to the batch size, enabling our method to scale to larger architectures compared to the PCA-based approach.

We analyzed the effects of over-smoothing, jumping knowledge, quantization, sparsification, and preconditioning on GNN optimization. Each impacts the loss landscape, trajectories, or both. While aligning with DNN findings, our results provide GNN-specific insights, highlighting the role of architectural choices in efficient training. Specifically, we identified that increasing GNN depth exacerbates over-smoothing, yielding flatter, non-convex loss landscapes and reduced performance. Quantization and sparsification introduce ruggedness and additional local minima. Preconditioning effectively alters optimizer trajectories, significantly aiding loss reduction.

In summary, GNN architecture, quantization, sparsification, and preconditioning substantially affect optimization trajectories, underscoring the importance of their careful consideration in GNN design and training.

Acknowledgment: Nils Kriege was supported by the Vienna Science and Technology Fund (WWTF) [10.47379/VRG19009].

References

1. Choromanska, A., Henaff, M., Mathieu, M., Arous, G.B., LeCun, Y.: The Loss Surfaces of Multilayer Networks. In: International Conference on Artificial Intelligence and Statistics. pp. 192–204 (2015), <https://proceedings.mlr.press/v38/choromanska15.pdf>
2. Dauphin, Y., et al.: Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In: Conference on Neural Information Processing Systems (2014), https://proceedings.neurips.cc/paper_files/paper/2014/file/04192426585542c54b96ba14445be996-Paper.pdf
3. Franklin, J.: The elements of statistical learning: data mining, inference and prediction. Springer Science and Business Media LLC, 2 edn. (2005). <https://doi.org/10.1007/bf02985802>
4. Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M.W., Keutzer, K.: A Survey of Quantization Methods for Efficient Neural Network Inference. CoRR pp. 291–326 (2022). <https://doi.org/10.1201/9781003162810-13>
5. Giles, L., Bollacker, K., Lawrence, S.: CiteSeer: An Automatic Citation Indexing System. In: ACM International Conference on Digital Libraries (1998). <https://doi.org/10.1145/276675.276685>
6. Goodfellow, I., Vinyals, O.: Qualitatively characterizing neural network optimization problems. In: International Conference on Learning Representations (2015). <https://doi.org/10.48550/arXiv.1412.6544>
7. Hochreiter, S., Schmidhuber, J.: Flat Minima. Neural Computation (1997). <https://doi.org/10.1162/neco.1997.9.1.1>
8. Hu, W., et al.: Open Graph Benchmark: Datasets for Machine Learning on Graphs. In: Conference on Neural Information Processing Systems (2020), <https://papers.neurips.cc/paper/2020/file/fb60d411a5c5b72b2e7d3527cfc84fd0-Paper.pdf>
9. Huang, Q., He, H., Singh, A., Lim, S., Benson, A.: Combining Label Propagation and Simple Models out-performs Graph Neural Networks. In: International Conference on Learning Representations (2021), <https://openreview.net/pdf?id=8E1-f3VhX1o>
10. Im, D.J., Tao, M., Branson, K.: An empirical analysis of the optimization of deep network loss surfaces. arXiv preprint arXiv:1612.04010 (2016). <https://doi.org/10.48550/arXiv.1612.04010>
11. Izadi, M.R., Fang, Y., Stevenson, R., Lin, L.: Optimization of Graph Neural Networks with Natural Gradient Descent. In: IEEE International Conference on Big Data (Big Data). pp. 171–179. IEEE (2020). <https://doi.org/10.1109/bigdata50022.2020.9378063>
12. Kipf, T., Welling, M.: Semi-supervised Classification with Graph Convolutional Networks. In: International Conference on Learning Representations (2017), <https://openreview.net/pdf?id=SJU4ayYgl>
13. Krizhevsky, A., Nair, V., Hinton, G.: Canadian Institute for Advanced Research (2009), <http://www.cs.toronto.edu/~kriz/cifar.html>

14. Li, H., Xu, Z., Taylor, G., Studer, C., Goldstein, T.: Visualizing the Loss Landscape of Neural Nets. In: Conference on Neural Information Processing Systems (2018), https://proceedings.neurips.cc/paper_files/paper/2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf
15. Li, Y., Yuan, Y.: Convergence Analysis of Two-layer Neural Networks with ReLU Activation. In: Conference on Neural Information Processing Systems (2017), https://proceedings.neurips.cc/paper_files/paper/2017/file/a96b65a721e561e1e3de768ac819ffbb-Paper.pdf
16. Lorch, E.: Visualizing Deep Network Training Trajectories with PCA. In: International Conference on Machine Learning Visualization Workshop (2016)
17. Morris, C., et al.: Weisfeiler and Leman go Machine Learning: The Story so far. Journal of Machine Learning Research (2023), <https://www.jmlr.org/papers/volume24/22-0240/22-0240.pdf>
18. Moustafa, S., Kriege, N., Gansterer, W.: Efficient Mixed Precision Quantization in Graph Neural Networks. In: IEEE 41st International Conference on Data Engineering (2025), <https://www.computer.org/csdl/proceedings-article/icde/2025/360300e038/26FZCgBUG4U>
19. Nguyen, Q., Hein, M.: The Loss Surface of Deep and Wide Neural Networks. In: International Conference on Machine Learning (2017), <https://dl.acm.org/doi/10.5555/3305890.3305950>
20. Saxe, A., McClelland, J., Ganguli, S.: Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In: International Conference on Learning Representations (2014), https://openreview.net/forum?id=_wzZwKpTDF_9C
21. Sen, P., Namata, G., Bilgic, M., Getoor, L., Gallagher, B., Eliassi-Rad, T.: Collective Classification in Network Data (2008). <https://doi.org/10.1609/aimag.v29i3.2157>
22. Shafiq, M., Gu, Z.: Deep Residual Learning for Image Recognition: A Survey. In: Applied Sciences. p. 8972. MDPI AG (2022). <https://doi.org/10.3390/app12188972>
23. Smith, L., Topin, N.: Exploring loss function topology with cyclical learning rate. arXiv preprint arXiv:1702.04283 (2017). <https://doi.org/10.48550/arXiv.1702.04283>
24. Soltanolkotabi, M., Javanmard, A., Lee, J.D.: Theoretical Insights Into the Optimization Landscape of Over-Parameterized Shallow Neural Networks. IEEE Trans. Inf. Theory pp. 742–769 (2019). <https://doi.org/10.1109/tit.2018.2854560>
25. Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph Attention Networks. In: International Conference on Learning Representations (2018), <https://openreview.net/pdf?id=rJXmpikCZ>
26. Wu, L., Cui, P., Pei, J., Zhao, L., Guo, X.: Graph Neural Networks: Foundation, Frontiers and Applications. In: Knowledge Discovery and Data Mining (2022)
27. Xu, K., Hu, W., Leskovec, J., Jegelka, S.: How Powerful are Graph Neural Networks? In: International Conference on Learning Representations (2019), <https://openreview.net/pdf?id=ryGs6iA5Km>
28. Xu, K., Li, C., Tian, Y., Sonobe, T., Kawarabayashi, K., Jegelka, S.: Representation Learning on Graphs with Jumping Knowledge Networks. In: International Conference on Machine Learning (2018), <https://proceedings.mlr.press/v80/xu18c/xu18c.pdf>
29. Zeng, H., Zhou, H., Srivastava, A., Kannan, R., Prasanna, V.: GraphSAINT: Graph Sampling Based Inductive Learning Method. In: International Conference on Learning Representations (2020), <https://openreview.net/pdf?id=BJe8pkHfWS>
30. Zhang, W., et al.: Graph Attention Multi-layer Perceptron. In: Knowledge Discovery and Data Mining (2022), <https://dl.acm.org/doi/10.1145/3534678.3539121>
31. Zhu, Z.A., et al.: A²Q: Aggregation-aware Quantization for Graph Neural Networks. In: International Conference on Learning Representations (2023), <https://openreview.net/pdf?id=7L2mgi0TNEP>