

# **A Novel Training Method for Deep Neural Networks Leveraging Attribution Graphs**

Deep neural networks have achieved remarkable success across a multitude of domains, demonstrating unparalleled capabilities in tasks ranging from image recognition and natural language processing to complex decision-making processes<sup>1</sup>. The increasing complexity and scale of modern deep learning models have enabled them to tackle increasingly intricate problems, pushing the boundaries of artificial intelligence<sup>2</sup>. While the backpropagation algorithm has served as the cornerstone for training these deep networks, its inherent limitations in terms of efficiency, interpretability, and robustness are becoming more apparent as models grow in size and are applied to more critical applications<sup>3</sup>. The computational cost and training time associated with backpropagation can be substantial, hindering the development and deployment of very large models<sup>2</sup>. Furthermore, the opacity of the learned representations and the difficulty in understanding the model's decision-making process pose significant challenges for interpretability<sup>5</sup>. The vulnerability of deep networks trained with backpropagation to adversarial attacks and their sensitivity to noisy data raise concerns about their reliability and robustness in real-world scenarios<sup>6</sup>. To address these challenges, alternative training methodologies are being explored, with the potential to offer improvements in these critical aspects of deep learning.

One promising avenue of research involves the use of attribution graphs, which provide a means to understand the inner workings of neural networks by assigning importance scores to different components of the model or the input data<sup>5</sup>. By visualizing these importance scores in a graph structure, researchers can gain insights into which features or neurons are most influential in the network's predictions<sup>8</sup>. This approach offers the potential to move beyond the global weight updates of backpropagation towards more targeted interventions within the network. This report outlines a research proposal for a novel training method that leverages attribution graphs to perform targeted neuron updates. The objective of this research is to develop a training paradigm that can enhance the efficiency, interpretability, and robustness of deep neural networks compared to traditional backpropagation by selectively updating critical neurons identified through attribution analysis.

## **Understanding Attribution Graphs**

In the context of deep learning, "attribution" refers to the process of assigning importance scores to the inputs or internal components of a neural network with respect to a specific output<sup>7</sup>. Essentially, it is a method of credit assignment that highlights which dimensions of the input data or which elements within the network

had the most significant influence on the model's prediction for a given instance <sup>5</sup>. This can be viewed as determining the effect of an input feature on the prediction function's output, inherently addressing a causal question about the model's behavior <sup>9</sup>. For example, in image classification, attribution methods might highlight the pixels most responsible for the classification of an image as a cat <sup>5</sup>.

Building upon the concept of attribution, "attribution graphs" represent a way to visualize and analyze these importance scores within the structure of a graph <sup>8</sup>. In a deep attribution graph, nodes can correspond to inputs, features, or even individual neurons within the network <sup>7</sup>. The edges in these graphs typically denote relationships or connections, weighted by the attribution scores, indicating the strength and direction of influence between the nodes <sup>8</sup>. For instance, in the SUMMIT system, an attribution graph visualizes highly activated neurons as vertices, with edges representing the most influential connections between them, thereby summarizing crucial neuron associations and substructures that contribute to the model's outcomes <sup>8</sup>. Another example is the Deep Attribution Graph for Deep Knowledge Transferability (DEPARA), where nodes represent inputs based on their vectorized attribution maps, and edges denote the relatedness between these inputs as measured by their feature similarity within a pre-trained deep neural network <sup>7</sup>.

The significance of attribution graphs lies in their ability to provide interpretability to the often opaque decision-making processes of deep neural networks <sup>5</sup>. By highlighting the input dimensions that are influential to a network's prediction, attribution methods can expose the statistical regularities learned by the model <sup>5</sup>. If these patterns align with human intuition, it can bolster confidence in the model's predictions <sup>5</sup>. Conversely, if attributions reveal that the model is exploiting spurious correlations or violating common sense, they can serve as a valuable debugging tool, allowing for the identification and potential correction of issues in the dataset or model architecture <sup>5</sup>. In fields like drug discovery, the interpretability offered by attribution graphs can increase user trust in the model's output <sup>5</sup>. Furthermore, attribution graphs can reveal not only the importance of individual features but also the relationships and interactions between neurons within the network <sup>8</sup>. This can shed light on potential biases learned by the model and provide a deeper understanding of its internal representations <sup>10</sup>. The diverse applications of attribution graphs, ranging from understanding model behavior in image and text processing to analyzing graph neural networks and even facilitating knowledge transfer between models, underscore their potential as a powerful tool in deep learning research <sup>5</sup>.

## **Limitations of Traditional Backpropagation**

Despite its widespread success, traditional backpropagation exhibits several limitations that become increasingly pertinent in the context of modern deep learning<sup>2</sup>. One significant challenge is the issue of vanishing and exploding gradients, particularly in very deep networks<sup>2</sup>. During the backward pass, as gradients are propagated from the output layer back to the earlier layers, they can become progressively smaller (vanishing) or larger (exploding), hindering effective learning in the initial layers of the network<sup>2</sup>. This makes training deep networks difficult and can lead to suboptimal performance.

Another limitation of backpropagation is its inherent sequential nature<sup>2</sup>. The algorithm requires a complete forward pass to compute the error at the output, followed by a backward pass to calculate the gradients for all the weights in the network<sup>12</sup>. This sequential processing can be inefficient, especially in distributed training scenarios where parallel updates could potentially accelerate the learning process<sup>3</sup>. Furthermore, the "locking problem" in backpropagation means that a layer cannot be updated until the subsequent layers have completed their forward and backward computations, preventing independent and potentially more asynchronous learning across different parts of the network<sup>3</sup>.

Backpropagation also lacks inherent interpretability<sup>3</sup>. While it effectively adjusts the weights of the network to minimize the error between the predicted and actual outputs, it does not directly provide insights into which input features are important for the model's decisions or how the network arrives at its conclusions<sup>12</sup>. The weight updates are driven by error signals, making it challenging to understand the semantic meaning or the contribution of individual neurons to the overall function of the network.

Moreover, backpropagation can be susceptible to getting stuck in local minima, especially in non-convex loss landscapes that are common in deep learning<sup>3</sup>. The algorithm relies on gradient descent, which can converge to a suboptimal solution if the error surface contains many local minima.

From a biological perspective, backpropagation is also considered implausible<sup>6</sup>. The algorithm requires bidirectional synaptic weight transport, where information needs to flow backward through the same connections used for the forward pass<sup>13</sup>. This is not consistent with the unidirectional flow of information observed in biological neurons<sup>6</sup>. Additionally, the update locking mechanism in backpropagation, where weight updates are delayed until the entire forward and backward passes are complete, is also not biologically realistic<sup>13</sup>. Finally, the performance of backpropagation is highly sensitive to the choice of hyperparameters, such as the learning rate, requiring careful

tuning which can be a time-consuming and computationally expensive process <sup>14</sup>.

The multifaceted limitations of backpropagation, spanning optimization challenges, efficiency bottlenecks, lack of inherent interpretability, and biological implausibility, provide a strong motivation for exploring alternative training methods. The potential to address these limitations at a more granular, neuron-level, guided by insights from attribution analysis, could lead to significant advancements in deep learning.

## Methods for Generating Attribution Graphs

Several methods exist for generating attribution graphs, each with its own strengths, weaknesses, computational cost, and interpretability <sup>7</sup>. These methods can be broadly categorized into gradient-based, perturbation-based, and propagation-based approaches.

**(a) Gradient-based methods:** These methods leverage the gradients of the output with respect to the input features to determine their importance.

- **Saliency Maps:** This approach calculates the gradient of the class score of interest with respect to the input pixels <sup>16</sup>. The magnitude of the gradient for each pixel is then visualized as a saliency map, highlighting the regions in the input that have the most influence on the model's prediction <sup>17</sup>. For example, in image classification, a saliency map might show which parts of an image (pixels) were most important for the network to classify it correctly <sup>19</sup>. Gradient-based methods like saliency maps are generally computationally efficient <sup>16</sup>. However, they can suffer from noisy gradients and sensitivity to small perturbations in the input, potentially leading to unreliable explanations <sup>15</sup>. Furthermore, the sign of the gradient, indicating whether increasing the pixel value would increase or decrease the class probability, is not always directly interpretable <sup>16</sup>. A known issue with basic saliency maps is the saturation problem, especially when using ReLU activation functions, where the gradient becomes zero for saturated neurons, potentially masking their importance <sup>16</sup>.
- **Integrated Gradients (IG):** To address some of the limitations of basic gradient methods, Integrated Gradients computes the attribution of each input feature by integrating the gradients along a path from a defined baseline input to the current input <sup>21</sup>. The baseline is often chosen as a neutral or zero-value input <sup>23</sup>. By accumulating the gradients along this path, IG aims to provide a more comprehensive attribution that is less sensitive to local gradient saturation <sup>21</sup>. This method satisfies important theoretical axioms such as sensitivity and completeness <sup>22</sup>. Sensitivity ensures that if an input and baseline differ in one feature and result in different predictions, the differing feature gets a non-zero

attribution<sup>24</sup>. Completeness states that the sum of attributions should equal the difference in the model's output between the input and the baseline<sup>22</sup>. While IG offers a more robust attribution compared to simple gradients, its results can be sensitive to the choice of the baseline, and the computational cost is higher due to the integration process, which typically involves multiple gradient computations<sup>21</sup>.

**(b) Perturbation-based methods:** These methods estimate feature importance by observing how the model's output changes when parts of the input are perturbed or removed.

- **LIME (Local Interpretable Model-agnostic Explanations):** LIME works by approximating the behavior of the complex deep learning model with a simpler, interpretable model (like a linear regression or decision tree) in the vicinity of a specific input instance<sup>26</sup>. It does this by generating perturbed samples around the original input, feeding them into the black-box model, and then learning a weighted local model based on how the predictions change<sup>28</sup>. The weights are assigned based on the proximity of the perturbed samples to the original input<sup>29</sup>. The coefficients of the local interpretable model then provide an explanation of the feature importance for that particular prediction<sup>27</sup>. LIME is model-agnostic, meaning it can be applied to any machine learning model<sup>28</sup>. It also provides local explanations, focusing on individual predictions<sup>28</sup>. However, LIME can be computationally expensive, especially for high-dimensional inputs, as it requires generating and classifying many perturbed samples<sup>28</sup>. The stability of the explanations can also depend on the sampling strategy and the choice of the local interpretable model<sup>25</sup>.
- **SHAP (SHapley Additive exPlanations):** SHAP is a model-agnostic method based on game theory that aims to explain the output of any machine learning model by calculating the contribution of each feature to the prediction<sup>30</sup>. It considers each feature as a "player" in a game where the "payout" is the prediction, and it calculates the Shapley values, which represent the average marginal contribution of each feature across all possible feature combinations<sup>30</sup>. SHAP provides consistent and interpretable feature importance scores, offering both local explanations for individual predictions and global insights into the overall feature importance<sup>30</sup>. For deep learning models, SHAP often uses specialized explainers like DeepExplainer, which leverages the properties of neural networks to estimate Shapley values more efficiently<sup>32</sup>. Despite its strong theoretical foundation and desirable properties, calculating exact Shapley values can be computationally intensive, especially for complex models and large datasets, as it may involve evaluating the model on many different feature subsets

<sup>30</sup>. Approximations are often used to make it more tractable.

**(c) Layer-wise Relevance Propagation (LRP):** LRP is a propagation-based method that aims to explain the prediction of a neural network by propagating the relevance score backward through the layers of the network, starting from the output layer <sup>33</sup>. It uses specific local propagation rules at each layer to redistribute the relevance received by a neuron in a higher layer to the neurons in the lower layer that contributed to its activation <sup>15</sup>. This process continues until the input layer is reached, providing a relevance score for each input feature, indicating its contribution to the final prediction <sup>35</sup>. LRP can highlight the positive contributions of input features to the network's classification <sup>36</sup>. It also adheres to a layer-wise conservation property, where the total relevance is conserved as it is propagated backward <sup>33</sup>. LRP is efficient as it typically involves a single backward pass <sup>15</sup>. However, the choice of appropriate propagation rules for different types of layers and network architectures can be crucial and might require careful tuning <sup>34</sup>. The method is also inherently tied to the network's structure, unlike model-agnostic methods <sup>34</sup>.

**Table 1: Comparison of Attribution Graph Generation Methods**

Method	Strengths	Weaknesses	Computational Cost	Interpretability	Relevant Snippets
<b>Saliency Maps</b>	Computationally efficient <sup>16</sup>	Noisy gradients, sensitive to perturbations, sign of gradient not always interpretable, saturation issues <sup>15</sup>	Low	Visual explanation of important regions <sup>16</sup>	<sup>15</sup>
<b>Integrated Gradients</b>	Satisfies sensitivity and completeness axioms <sup>22</sup> , implementation invariance <sup>24</sup> , strong theoretical justification <sup>24</sup>	Sensitive to baseline choice, higher computational cost than basic gradients <sup>21</sup>	Moderate	Provides feature importance scores relative to a baseline <sup>24</sup>	<sup>15</sup>
<b>LIME</b>	Model-agnostic <sup>28</sup> , provides local explanations <sup>28</sup>	Computationally expensive, potential instability due to sampling, explanations can disagree with other methods <sup>25</sup>	High	Provides interpretable local models (e.g., linear regression coefficients) <sup>28</sup>	<sup>26</sup>

<b>SHAP</b>	Model-agnostic <sup>30</sup> , consistent and interpretable feature importance scores <sup>30</sup> , provides both local and global interpretability <sup>30</sup> , handles complex interactions <sup>30</sup>	Computationally expensive, complexity of interpretation can be high, approximation methods often used <sup>25</sup>	High	Quantifies the contribution of each feature to the prediction <sup>32</sup>	<sup>30</sup>
<b>LRP</b>	Highlights positive contributions <sup>36</sup> , layer-wise conservation of relevance <sup>33</sup> , computationally efficient <sup>34</sup>	Dependent on network architecture, parameter selection for propagation rules can be challenging, potential sensitivity to rules <sup>34</sup>	Low	Provides relevance scores for input features, indicating their contribution to the prediction <sup>35</sup>	<sup>15</sup>



## Works cited

1. Deep Learning With Spiking Neurons: Opportunities and Challenges - Frontiers, accessed March 27, 2025, <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2018.00774/full>
2. Why we need a better learning algorithm than Backpropagation in Deep Learning - Medium, accessed March 27, 2025, <https://medium.com/towards-data-science/why-we-need-a-better-learning-algorithm-than-backpropagation-in-deep-learning-2faa0e81f6b>
3. What is Backpropagation Algorithm? Benefits & Drawbacks, accessed March 27, 2025, <https://www.deepchecks.com/glossary/backpropagation-algorithm/>
4. The Problem with Back-Propagation | by Anthony Repetto | TDS Archive - Medium, accessed March 27, 2025, <https://medium.com/towards-data-science/the-problem-with-back-propagation-13aa84aabd71>
5. proceedings.neurips.cc, accessed March 27, 2025, <https://proceedings.neurips.cc/paper/2020/file/417fbbf2e9d5a28a855a11894b2e795a-Paper.pdf>
6. [Discussion] What are the problems of the backpropagation algorithm? - Reddit, accessed March 27, 2025, [https://www.reddit.com/r/MachineLearning/comments/70tz1n/discussion\\_what\\_are\\_the\\_problems\\_of\\_the/](https://www.reddit.com/r/MachineLearning/comments/70tz1n/discussion_what_are_the_problems_of_the/)
7. openaccess.thecvf.com, accessed March 27, 2025, [https://openaccess.thecvf.com/content\\_CVPR\\_2020/papers/Song\\_DEPARA\\_Deep\\_Attribution\\_Graph\\_for\\_Deep\\_Knowledge\\_Transferability\\_CVPR\\_2020\\_paper.pdf](https://openaccess.thecvf.com/content_CVPR_2020/papers/Song_DEPARA_Deep_Attribution_Graph_for_Deep_Knowledge_Transferability_CVPR_2020_paper.pdf)
8. SUMMIT: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations, accessed March 27, 2025, <https://par.nsf.gov/servlets/purl/10183516>
9. Neural Network Attributions: A Causal Perspective - Proceedings of Machine Learning Research, accessed March 27, 2025, <http://proceedings.mlr.press/v97/chattopadhyay19a/chattopadhyay19a.pdf>
10. Using attribution to decode binding mechanism in neural network models for chemistry | PNAS, accessed March 27, 2025, <https://www.pnas.org/doi/10.1073/pnas.1820657116>
11. GOAt: Explaining Graph Neural Networks via Graph Output Attribution - arXiv, accessed March 27, 2025, <https://arxiv.org/html/2401.14578v1>
12. Backpropagation - Engati, accessed March 27, 2025, <https://www.engati.com/glossary/back-propagation>
13. Learning Without Feedback: Fixed Random Learning Signals Allow for Feedforward Training of Deep Neural Networks - Frontiers, accessed March 27, 2025, <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2021.629892/full>
14. Understand the Impact of Learning Rate on Neural Network Performance -

- MachineLearningMastery.com, accessed March 27, 2025,  
<https://machinelearningmastery.com/understand-the-dynamics-of-learning-rate-on-deep-learning-neural-networks/>
15. A unified view of gradient-based attribution methods for Deep Neural Networks - Heatmapping.org, accessed March 27, 2025,  
<http://www.interpretable-ml.org/nips2017workshop/papers/02.pdf>
  16. 28 Saliency Maps – Interpretable Machine Learning, accessed March 27, 2025,  
<https://christophm.github.io/interpretable-ml-book/pixel-attribution.html>
  17. An Introduction to Saliency Maps in Deep Learning - MarkTechPost, accessed March 27, 2025,  
<https://www.marktechpost.com/2022/03/07/an-introduction-to-saliency-maps-in-deep-learning/>
  18. Saliency map - Wikipedia, accessed March 27, 2025,  
[https://en.wikipedia.org/wiki/Saliency\\_map](https://en.wikipedia.org/wiki/Saliency_map)
  19. What is Saliency Map? - GeeksforGeeks, accessed March 27, 2025,  
<https://www.geeksforgeeks.org/what-is-saliency-map/>
  20. Explainable machine learning #3: Saliency Maps - YouTube, accessed March 27, 2025,  
<https://www.youtube.com/watch?v=eudf1wQmXnc>
  21. IG2: Integrated Gradient on Iterative Gradient Path for Feature Attribution - arXiv, accessed March 27, 2025,  
<https://arxiv.org/html/2406.10852v1>
  22. Integrated Gradients — Alibi 0.9.7.dev0 documentation, accessed March 27, 2025,  
<https://docs.seldon.io/projects/alibi/en/latest/methods/IntegratedGradients.html>
  23. Integrated gradients | TensorFlow Core, accessed March 27, 2025,  
[https://www.tensorflow.org/tutorials/interpretability/integrated\\_gradients](https://www.tensorflow.org/tutorials/interpretability/integrated_gradients)
  24. Interpreting Deep Neural Networks using Integrated Gradients ..., accessed March 27, 2025,  
<https://towardsdatascience.com/interpreting-deep-neural-networks-using-integrated-gradients-f9b8ecdd3c57/>
  25. Machine learning interpretability with feature attribution - Christian Garbin's personal blog, accessed March 27, 2025,  
<https://cgarbin.github.io/machine-learning-interpretability-feature-attribution/>
  26. Understand Network Predictions Using LIME - MathWorks, accessed March 27, 2025,  
<https://www.mathworks.com/help/deeplearning/ug/understand-network-predictions-using-lime.html>
  27. LIME: Local Interpretable Model-Agnostic Explanations - C3 AI, accessed March 27, 2025,  
<https://c3.ai/glossary/data-science/lime-local-interpretable-model-agnostic-explanations/>
  28. From Opaque to Transparent: Understanding Machine Learning ..., accessed March 27, 2025,  
<https://ishwaryasriraman.medium.com/from-opaque-to-transparent-understanding-machine-learning-models-with-lime-3f2a2d147642>
  29. LIME: explain Machine Learning predictions | by Giorgio Visani | TDS Archive - Medium, accessed March 27, 2025,

<https://medium.com/towards-data-science/lime-explain-machine-learning-predictions-af8f18189bfe>

30. How to explain neural networks using SHAP | Your Data Teacher, accessed March 27, 2025,  
<https://www.yourdatateacher.com/2021/05/17/how-to-explain-neural-networks-using-shap/>
31. Deep Learning Model Explainability with SHAP - Paperspace Blog, accessed March 27, 2025,  
<https://blog.paperspace.com/deep-learning-model-interpretability-with-shap/>
32. Can we use SHAP values to explain the performance of a Neural Network ? | ResearchGate, accessed March 27, 2025,  
[https://www.researchgate.net/post/Can\\_we\\_use\\_SHAP\\_values\\_to\\_explain\\_the\\_performance\\_of\\_a\\_Neural\\_Network](https://www.researchgate.net/post/Can_we_use_SHAP_values_to_explain_the_performance_of_a_Neural_Network)
33. 10 Layer-Wise Relevance Propagation: An Overview - Fraunhofer Heinrich-Hertz-Institut, accessed March 27, 2025,  
<https://iphome.hhi.de/samek/pdf/MonXAI19.pdf>
34. (PDF) Layer-Wise Relevance Propagation: An Overview, accessed March 27, 2025,  
[https://www.researchgate.net/publication/335708351\\_Layer-Wise\\_Relevance\\_Propagation\\_An\\_Overview](https://www.researchgate.net/publication/335708351_Layer-Wise_Relevance_Propagation_An_Overview)
35. Layer-wise Relevance Propagation for Neural Networks with Local Renormalization Layers, accessed March 27, 2025,  
<https://arxiv.org/abs/1604.00825>
36. Layer-Wise Relevance Propagation for Explaining Deep Neural Network Decisions in MRI-Based Alzheimer's Disease Classification, accessed March 27, 2025,  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC6685087/>
37. Gradient-Based Attribution Methods - CGL @ ETHZ, accessed March 27, 2025,  
<https://cgl.ethz.ch/Downloads/Publications/Papers/2019/Anc19c/Anc19c.pdf>
38. Evaluating Attribution Methods in Machine Learning Interpretability, accessed March 27, 2025, <https://par.nsf.gov/servlets/purl/10385266>
39. Explaining Deep Learning Models for Tabular Data Using Layer-Wise Relevance Propagation - MDPI, accessed March 27, 2025,  
<https://www.mdpi.com/2076-3417/12/1/136>