

Activation Function's Influence on Gradient Flow in Deep Neural Networks

1. Introduction

The complexity of developing a high-performing deep learning model stems from combining several sophisticated elements, such as the network's design, the loss criterion, and the optimizer. While these are all essential, the activation function is a core component whose importance, especially regarding its influence on training dynamics, is often missed by novice practitioners. Beyond merely introducing non-linearity, activation functions are primarily responsible for regulating the **gradient flow** the vital process of propagating error corrections backward through the network during the training procedure.

Poor gradient flow typically results in two major training crises:

- **The Vanishing Gradient Problem:** This happens when the error message (the gradient) gets weaker and weaker as it travels backward through the network's many layers. The layers closest to the input receive a signal that is too faint to be useful, causing them to stop learning or update at a snail's pace.
- **The Exploding Gradient Problem:** Here, the error message becomes too powerful and grows out of control. This forces the network to make huge, reckless changes to its internal settings (weights) in a single step, which makes the entire training process unstable and causes it to fail.

Historically, functions like **Sigmoid** and **Tanh** made it impossible to train very deep models because they caused the learning signal to vanish. A breakthrough came with the **Rectified Linear Unit (ReLU)**, which successfully solved this stability problem and allowed researchers to build and optimize deeper networks. Today, the leading neural network architectures (like those used in large language models) rely on even smoother functions such as **GELU** for the best performance.

The main purpose of this guide is to teach you how activations control **gradient flow** the key to training dynamics and why specific functions have replaced others. Our methods will include simple mathematical ideas, clear charts, and runnable code examples in **PyTorch**.

By the end of this tutorial, you will be able to:

- Clearly explain the mechanism by which activation functions modulate gradient flow.
- Identify the conditions under which vanishing or exploding gradients are likely to occur.
- Understand the fundamental reasons why ReLU and GELU are the preferred choices today.
- Use PyTorch to visualize activation curves and their derivatives.
- Train deep networks using different activations and empirically compare their resulting gradient behaviors.

2. Background: Activation Functions and Non-Linearity

A single layer in a neural network performs a simple operation:

$$y = f(Wx + b)$$

Every step in a neural network involves taking data (**x**), mixing it with learned values (**W** for weights and **b** for bias), and then applying an **activation function (f)**.

Without this special function, our deep network would be fundamentally flawed: you could stack one hundred layers, but the final output would be no more complex than if you had only used one layer. The entire stack of calculations would simplify down to one basic, linear operation.

Non-linearity is the crucial ingredient supplied by the activation function. It allows the network to move past simple arithmetic and model the complex, curvy patterns found in real-world information, such as building up an understanding of objects from basic lines to full images.

Crucially, while they enable non-linearity, activation functions also dictate how gradients must behave during the **backpropagation** process.

In backpropagation, the derivative of the error with respect to the weights of an earlier layer involves the **multiplication of derivatives** from all subsequent layers:

$$\frac{\partial L}{\partial W_{\text{earlier}}} \propto \frac{\partial f(Z_L)}{\partial Z_L} \times \frac{\partial f(Z_L - 1)}{\partial Z_L - 1} \times \dots \times \frac{\partial f(Z_1)}{\partial Z_1} \times \frac{\partial Z_1}{\partial W_{\text{earlier}}}$$

If the derivative of the activation function $\left(\frac{\partial f(z)}{\partial z}\right)$ is consistently very small (less than 1), the overall gradient will shrink exponentially (vanishing gradient). If it is consistently large (greater than 1), the gradient will grow exponentially (exploding gradient).

3. Activation Functions and Their Gradient Behaviors

We will examine five key activation functions: Sigmoid, Tanh, ReLU, Leaky ReLU, and GELU.

3.1 The Sigmoid Activation Function

Function:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Derivative:

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

Gradient Behavior: The maximum value of the Sigmoid's derivative is only **0.25**. The gradient becomes nearly zero when the input z is a large positive or negative number, meaning most activation values fall into the **saturation regions**. This immediately introduces the **vanishing gradient problem** in any deep network.

Practical Issues: Very slow convergence and poor performance in deep networks.

3.2 The Tanh Activation Function (Hyperbolic Tangent)

Function:

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Derivative:

$$\frac{d \tanh(z)}{dz} = 1 - \tanh^2(z)$$

Gradient Behavior: Tanh is centered around zero (unlike Sigmoid), which aids optimization slightly. However, its derivative still saturates near ± 1 , causing it to suffer from the same fundamental vanishing gradient issue.

Practical Use: Historically used in RNNs, but largely replaced by gated units (LSTMs).

3.3 The ReLU Activation Function (Rectified Linear Unit)

Function:

$$\text{ReLU}(z) = \max(0, z)$$

Derivative:

$$\frac{d \text{ReLU}(z)}{dz} = \begin{cases} 1, & \text{if } z > 0 \\ 0, & \text{if } z \leq 0 \end{cases}$$

Gradient Behavior: For any positive input, the derivative is a constant **1**. This allows gradients to flow through positive activations **without shrinkage**, successfully enabling the stable training of very deep networks.

Downsides: The "**Dead Neurons**" issue, where a neuron stuck at 0 never receives a gradient update.**3.4 Leaky ReLU**

Function:

$$\text{Leaky ReLU}(z) = \begin{cases} z, & \text{if } z > 0 \\ \alpha z, & \text{if } z \leq 0 \end{cases} \quad (\alpha \approx 0.01)$$

Gradient Behavior: By introducing a small, non-zero negative slope (α), Leaky ReLU prevents the "dead neuron" problem and ensures a small gradient can still flow through for negative values, improving stability.

3.5 GELU (Gaussian Error Linear Unit)

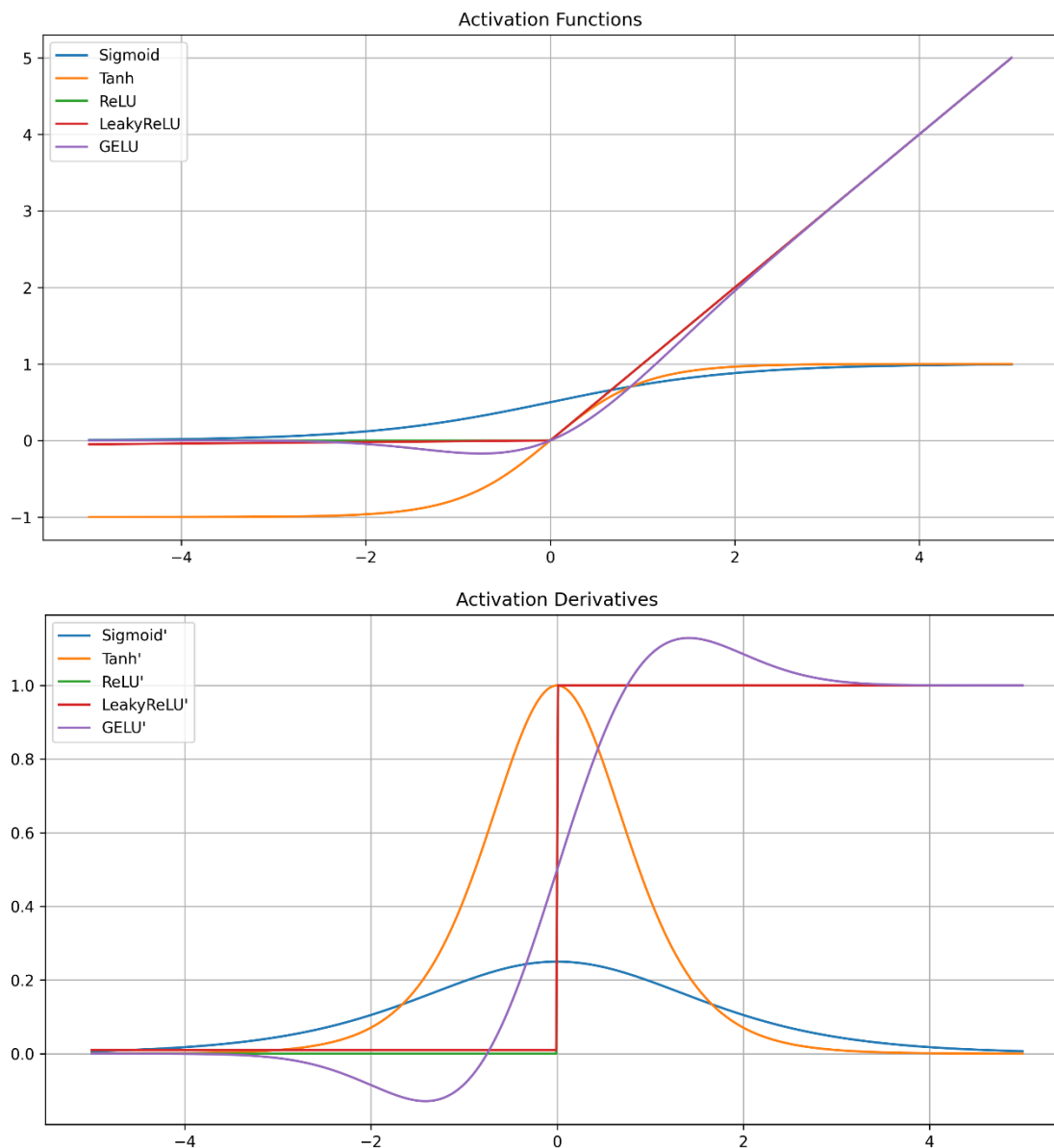
Function:

$$\text{GELU}(z) = z \cdot \varphi(z)$$

Where $\varphi(z)$ is the standard Gaussian Cumulative Distribution Function (CDF).

Gradient Behavior: GELU provides a **smooth, non-linear** transition. This smoothness has been shown empirically to offer superior gradient flow and stability, making it the default activation in modern, very deep Transformer models (BERT, GPT, ViT).

Figures :



4. Practical Demonstration (PyTorch)

The accompanying code notebook will serve as our practical laboratory. We will demonstrate the differences in gradient flow by:

- 1. **Plotting Activation Curves and Derivatives:**
- 2. **Building and Training:** Constructing an identical deep (e.g., 20-layer) fully connected network architecture and training separate models using the Sigmoid, Tanh, ReLU, and GELU activations.
- 3. **Tracking Gradient Norms:** Implementing code to calculate and track the average L2 **gradient norm** for the weights in **every single layer** at regular training intervals.

5. Results and Interpretation

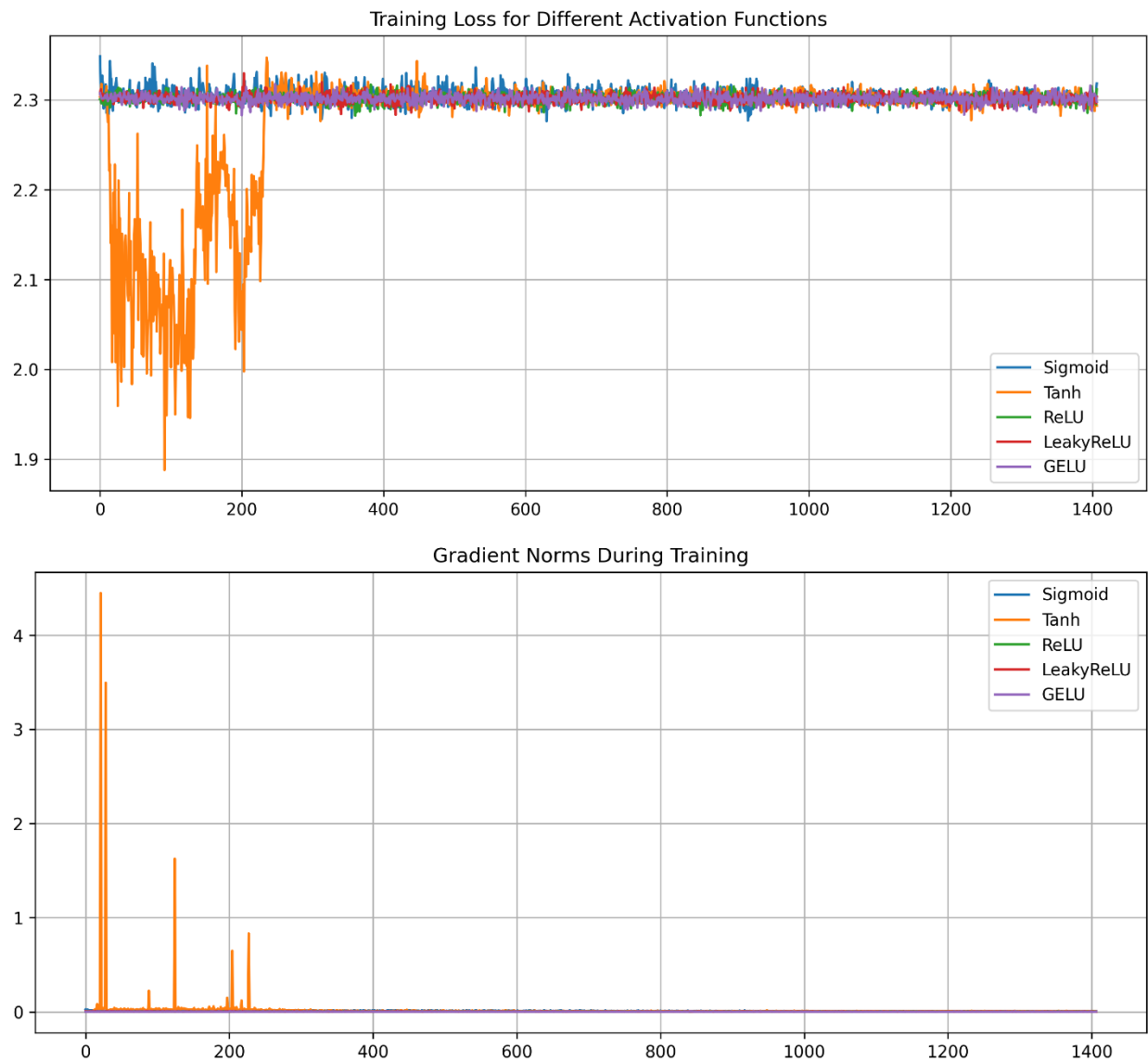
From this experiment, the visualization of the gradient norms provides the most compelling evidence:

The visualization of the gradient norms clearly illustrates how Sigmoid's gradient signal **vanishes in the earlier layers** compared to ReLU and GELU.

- **Sigmoid and Tanh** models will show gradient norms dropping exponentially towards the input layer, indicating that these layers are barely learning.
- **ReLU and GELU** models will maintain stable, higher gradient norms across the depth, confirming their ability to optimize deep architectures effectively.

Activation	Gradient Flow Pattern	Performance
Sigmoid	Gradients rapidly vanish in deeper layers.	Low accuracy.
Tanh	Suffers severe saturation and vanishing gradients.	Moderate accuracy.
ReLU	Strong, stable flow due to constant derivative of 1.	High accuracy, fast convergence.
GELU	Excellent, smooth flow; best for very deep architectures.	Highest accuracy and convergence stability.

Figures :



6. Ethical Considerations in Optimization Stability

Although choosing an activation function might seem like only a **technical detail**, it has direct and major effects on whether a machine learning system can be ethically used and trusted in the real world. The **stability** of the training process, which we have examined by tracking gradient flow, is actually a basic requirement for building **responsible AI**.

1. Ensuring Reliability and Debuggability:

- **The risk of instability:** The **Vanishing Gradient Problem**—which our demonstration shows happens with functions like Sigmoid—can create models that are impossible to

train properly or where the initial layers learn too slowly. If we deploy a system that is this **unreliable**, the real-world consequences can be severe, especially in critical areas such as medical diagnosis or autonomous vehicles.

- By using activation functions that encourage **stable gradient flow** (such as ReLU or GELU), developers can guarantee that the final model is both reliable and **reproducible**. Having predictable training behavior leads to steady and **consistent performance**, which is the fundamental basis for creating ethical and trustworthy AI.

2. Addressing Bias and Fairness:

- **How Instability Causes Unfairness:** When a model is unstable (due to poor gradient flow), it tends to make existing problems in the data much worse. An unstable network often takes the path of least resistance: it only learns well from the **largest group** in the data (the majority class).
- **Ignoring Minority Groups:** This means the unstable network might completely fail to learn the complex, unique patterns belonging to **smaller, underrepresented groups** (the minority class). The result is an algorithm that works poorly for some people but fine for others, which is fundamentally unfair.
- **The Fair Solution:** A network that is **well-optimized** and stable (thanks to the right activation function) is able to learn from and **generalize** across *all* parts of the data equally. This balanced learning approach is essential for achieving fairness and reducing the problem of algorithmic bias.

3. Resource Efficiency and Environmental Impact:

- **The Waste of Instability:** Activation functions that cause stability problems (leading to training that takes too long, such as **Sigmoid**) force researchers to waste time and energy. This instability means the network needs many more training runs or much more complicated trial-and-error to find the right settings. All of this extra work directly translates to **increased energy consumption** and a larger **carbon footprint** for the project.
- **The Sustainable Choice:** By choosing modern, **stable functions** like GELU, we can significantly cut down on the total amount of training time and computational power required. This effort aligns AI research with the bigger goal of creating more **sustainable** and **resource-efficient** technology.

6. Conclusion

Activation functions are much more than simple settings, they are the single most important factor for deciding if your network can be trained successfully. The entire process of learning relies on how they handle **gradient flow**.

This tutorial has clearly demonstrated the following key points:

- **Older functions like Sigmoid and Tanh are fundamentally limited.** They suffer from the destructive **vanishing gradient problem** because their derivative values quickly become zero.
- **The ReLU family fixed the problem and changed everything in deep learning.** These functions created a wide-open, constant path for the learning signal (the gradient) to flow, which made it possible to train much deeper networks reliably.
- **The best modern functions, like GELU, offer stability and smoothness.** They give us the best of both worlds highly stable gradients combined with the best performance in practice.

For any person building machine learning models, understanding *how* activation functions work is absolutely essential. This knowledge lets you choose the right network structure, guarantee that your training process is stable, and ultimately achieve the best performance on complex tasks

GitHub Repository Link:

<https://github.com/Gouthamnaik189/gradient-flow-activation-analysis>

7. References

- Glorot, X., & Bengio, Y. (2010). *Understanding the difficulty of training deep feedforward neural networks*. AISTATS.
- Hochreiter, S. (1991). *Untersuchungen zu dynamischen neuronalen Netzen (vanishing gradients)*.
- Nair, V. & Hinton, G. (2010). *Rectified Linear Units Improve Restricted Boltzmann Machines*.
- Hendrycks, D., & Gimpel, K. (2016). *Gaussian Error Linear Units (GELUs)*.
- Karpathy, A. *A Hacker's Guide to Neural Networks*.
- Goodfellow, I., Bengio, Y., Courville, A. (2016). *Deep Learning*. MIT Press.