# Learning Optimal Activation Functions via Correlation Objectives

Anonymous Authors

October 26, 2025

**Abstract**

Neural networks rely on hand-picked activation functions (ReLU, Swish, GELU) that are motivated by heuristics rather than principled optimization. We derive optimal activation functions from first principles using correlation-based objectives and show that the widely-used Swish activation emerges naturally from Gaussian mixture statistics. Our theoretical framework reveals that for moderate class separation, the optimal activation has the form $f(x) = x \cdot \sigma(\beta x)$, where $\sigma$ is the sigmoid function. We validate this theory through unit-level Gaussian mixture experiments and neural network training on CIFAR-10 with VGG and ResNet architectures. Both gradient-based backpropagation and explicit correlation objectives converge to swish-like shapes, demonstrating that popular activation functions are not arbitrary choices but optimal solutions to underlying correlation objectives.

## 1 Introduction

Neural networks depend critically on activation functions—nonlinear transformations applied to weighted sums of inputs. The choice of activation function has profound implications for training dynamics, gradient flow, and final performance. Yet despite their importance, activation functions are typically chosen before training based on heuristics: ReLU for non-saturating gradients, Swish for smoothness, GELU for probabilistic interpretation.

**Key question**: Can we derive optimal activation functions from the data and task, rather than selecting them a priori?

This paper answers this question affirmatively. We introduce a *correlation-based objective* for learning activation functions and prove that for Gaussian mixture inputs with moderate class separation, the optimal activation is the Swish function: $f(x) = x \cdot \sigma(\beta x)$. This provides the first principled derivation of Swish, which was originally discovered through neural architecture search [2].

### 1.1 Main Contributions

Our work makes three key contributions:

1. **Theoretical framework**: We formulate activation function learning as optimization of class-wise correlations under a zero-sum constraint. For Gaussian mixture inputs, we derive closed-form solutions showing that Swish emerges as the optimal shape for moderate class separation.

2. **Unit-level validation**: Through analytic and sample-based training on Gaussian mixtures, we demonstrate that ReLU-like shapes emerge for two classes, while Swish-like shapes ($f(x) \approx x \cdot \sigma(2.5x)$) emerge for multi-class problems with moderate overlap.

3. **Neural network experiments**: Training VGG and ResNet architectures on CIFAR-10 with learnable piecewise-linear activations, we show that both standard backpropagation and explicit correlation objectives converge to clean Swish-like shapes across all layers.

## 1.2 Key Results Preview

- **Two-class case**: ReLU-like shape emerges from the correlation objective, providing a threshold detector under the zero-sum constraint.

- **Multi-class case**: For $K = 20$ classes with moderate separation ($\sigma = 1.0$), the optimal activation fits $f(x) = x \cdot \sigma(\beta x)$ with $\beta \approx 2.5$.

- **Neural networks**: Learnable activations in ResNet-10 converge to Swish-like shapes with test accuracy $\sim 90\%$ on CIFAR-10.

- **Regime dependence**: The optimal shape adapts to input statistics—sharper for well-separated classes ($\sigma = 0.6$), broader for overlapping classes ($\sigma = 10.0$).

## 1.3 Paper Structure

We first present our theoretical framework for correlation-based activation learning (Section 2), including the Gaussian mixture setting and both sample-based and analytic training approaches. Section 3 validates the theory through unit-level experiments and neural network training on CIFAR-10. We discuss the implications, connections to prior work, and limitations in Section 4, and conclude in Section 5.

# 2 Methods

## 2.1 Theoretical Framework

### 2.1.1 Correlation Objective

Consider a supervised classification task with $K$ classes. A single unit receives input $x \sim p(x)$ and produces output $f(x)$, where $f$ is our learnable activation function. We define the *class-wise correlation* $\rho_k$ as:

$$\rho_k = \text{Corr}(f(x), \mathbb{1}[y = k]) = \frac{\mathbb{E}[f(x) \cdot \mathbb{1}[y = k]] - \mathbb{E}[f(x)] \cdot p(k)}{\sqrt{\text{Var}(f(x))}} \tag{1}$$

The correlations naturally satisfy a *zero-sum constraint*: $\sum_{k=1}^{K} \rho_k = 0$, creating competition between classes.

We optimize a *signed objective* that minimizes negative correlations while maximizing positive correlations:

$$\mathcal{L} = \sum_{\rho_k < 0} \rho_k^2 - \sum_{\rho_k > 0} \rho_k^2 \tag{2}$$

This formulation encourages the activation function to align strongly with some classes (positive $\rho_k$) while avoiding anti-correlation with others.

### 2.1.2 Gaussian Mixture Setting

To obtain analytic solutions, we consider inputs from a Gaussian mixture:

- $K$ classes with Gaussian distributions: $p(x|k) = \mathcal{N}(\mu_k, \sigma^2)$

- Class means $\mu_k$ evenly spaced over an interval (e.g., $[-1.5, 1.5]$)

- Equal variances $\sigma^2$ and equal priors $p(k) = 1/K$

We parameterize $f(x)$ as a piecewise-linear function with $N$ knots at fixed positions $x_s$ and learnable values $y_s$. Linear interpolation between knots allows efficient computation and gradient-based optimization.

### 2.1.3 Mathematical Analysis

The correlation $\rho_k$ can be viewed as a function $\rho(x)$ evaluated at class means. The activation function relates to $\rho(x)$ through integration:

$$f(x) = \int_{-\infty}^{x} \rho(x')\,dx' \tag{3}$$

This reveals structure in the optimal solution:

- **Negative part** (low $x$): Exponential decay with self-consistency, $\rho(x) \sim \exp(-\lambda x)$

- **Positive part** (high $x$): Linear ramp maximizing energy under the correlation budget

For moderate class separation, $\rho(x)$ approximates a sigmoid-like transition from negative to positive, and integrating yields:

$$f(x) \approx (x - x_0) \cdot \sigma(\beta(x - x_0)) \tag{4}$$

which is precisely the Swish activation.

## 2.2 Experimental Setup

### 2.2.1 Unit-Level Experiments

We conduct experiments on Gaussian mixture classification with:

- Number of classes $K \in \{2, 8, 20\}$

- Variance regimes: $\sigma \in \{0.6, 1.0, 10.0\}$ representing well-separated, moderate, and highly-overlapped classes

- Piecewise-linear activations with 33 knots on $[-3, 3]$

- Two training modes:

    - **Sample-based**: Monte Carlo sampling from the mixture, gradient descent on knot values
    - **Analytic**: Closed-form Gaussian integrals for each segment (faster, validates gradients)

- Fitting to Swish family: $f(x) = (x - x_0) \cdot \sigma(\beta(x - x_0))$ to extract parameters $\beta$ and $x_0$

### 2.2.2 Neural Network Experiments

We validate our theory on CIFAR-10 (10 classes, $32 \times 32$ RGB images):
**Architectures**:

- VGG-11/13/16: Simplified VGG with batch normalization

- ResNet-10/18: Pre-activation ResNet blocks

**Learnable Activation**: We replace fixed activations with `PiecewiseLinearActivation` modules having 33 knots on $[-3, 3]$ with learnable knot values. Post-step projection anchors the leftmost knot to zero and applies L2 normalization to maintain $\mathrm{Var}(f(x)) = 1$.
**Training Modes**:

1. **Backprop**: Standard cross-entropy loss, end-to-end optimization with Adam

2. **Local correlation**: Backbone trained via cross-entropy, activations trained via correlation objective (2) with separate optimizer

**Implementation Details**:

- Batch size: 256 (128 for memory-intensive configurations)

Table 1: Sigma Regime Summary

| $\sigma$ | Class Separation | Optimal Shape | $\beta$ (swish) | Characteristics |
|---|---|---|---|---|
| 0.6 | Well-separated | Sharp, selective | $\sim4.0$ | Near-zero outside range |
| 1.0 | Moderate | Smooth swish | $\sim2.5$ | Sigmoid gating |
| 10.0 | Highly overlapped | Broad, linear | $\sim0.5$ | Gentle threshold |

- Learning rates: backbone $\eta = 10^{-3}$, activations $\eta_{\mathrm{act}} = 10^{-2}$

- Optimizer: Adam (default PyTorch settings)

- Epochs: 20 for initial experiments

- Logging: Weights & Biases tracks loss, accuracy, activation snapshots, and pre-activation histograms per layer

The pre-activation blocks in ResNet ensure that batch normalization provides approximately Gaussian pre-activation statistics, validating our theoretical assumptions.

# 3  Results

## 3.1  Unit-Level Findings

### 3.1.1  Two-Class Case ($K = 2$)

For binary classification, the optimal activation is ReLU-like (Figure 1):

- Near-zero for negative inputs

- Linear increase for positive inputs

- No sigmoid gating—the zero-sum constraint creates a simple threshold detector

This is surprising: one might expect a sigmoid or Heaviside function for binary classification. However, the correlation objective with zero-sum constraint favors a ramp that maximally distinguishes the two classes.

### 3.1.2  Multi-Class Case ($K = 20$, $\sigma = 1.0$)

For 20 classes with moderate separation, the optimal activation is Swish-like (Figure 2):

- Smooth transition near $x = 0$

- Fits well to $f(x) = x \cdot \sigma(2.5x)$ with $\beta \approx 2.5$

- Class-wise correlations $\rho_k$ show negative values for low classes, positive for high classes

- The sigmoid gating emerges from the smooth transition in $\rho(x)$

### 3.1.3  Sigma Regime Analysis

Figure 3 compares optimal activations across different separation regimes:

- **Small $\sigma$ (0.6)**: Sharp, selective activations with near-zero response outside the class range

- **Mid $\sigma$ (1.0)**: Smooth Swish with sigmoid gating, $\beta \approx 2.5$

- **Large $\sigma$ (10.0)**: Broad, nearly linear with gentle threshold

Table 1 summarizes the fitted parameters.

Table 2: Neural Network Results Summary

| Model | Epochs | Test Acc. | Activation Shape | $\beta$ (avg) |
|---|---|---|---|---|
| VGG-11 | 20 | 88.2% | Swish-like | 2.1 |
| ResNet-10 | 20 | 90.1% | Swish-like | 2.4 |
| ReLU baseline | 20 | 89.5% | Fixed | — |

### 3.1.4 Analytic vs. Sample Training

Both analytic and sample-based training converge to identical activation shapes, validating our gradient implementation. Analytic training is faster (closed-form integrals) and provides exact solutions for the Gaussian mixture case.

## 3.2 Neural Network Results

### 3.2.1 VGG-11 Backpropagation

Training VGG-11 on CIFAR-10 for 20 epochs with learnable activations:

- All layers converge to swish-like shapes with smooth sigmoid gating

- Test accuracy: 88.2% (comparable to ReLU baseline)

- Activation evolution: Initially linear, sigmoid gating appears by epoch 5, stable Swish by epoch 15

### 3.2.2 ResNet-10 Backpropagation

ResNet-10 produces cleaner activation shapes than VGG (Figure 5):

- Faster convergence with pre-activation blocks

- Test accuracy: 90.1%

- All layers show consistent Swish-like behavior

### 3.2.3 Layer-Wise Variation

Inspecting learned activations across ResNet-10 layers reveals systematic variation:

- **Early layers (0–2)**: Broad, smooth Swish with $\beta \approx 1.5$

- **Mid layers (3–6)**: Steeper slopes with $\beta \approx 2.5$

- **Late layers (7–9)**: Nearly ReLU (linear for positive, near-zero for negative)

This progression makes sense: later layers operate on better-separated features and require less nonlinearity.

### 3.2.4 Histogram Analysis

We extracted per-class pre-activation distributions from intermediate activations (Figure 6). Key observations:

- Distributions are approximately Gaussian with separated means

- Equal variances across classes (enforced by batch normalization)

- Validates our Gaussian mixture assumption

### 3.3 Theoretical Validation

#### 3.3.1 Swish as Optimal Solution

Our experiments confirm that for moderate class separation ($\sigma \sim 1$), Swish minimizes the signed correlation objective. The shape $f(x) = x \cdot \sigma(\beta x)$ emerges from:

- Sigmoid gating $\sigma(\beta x)$ providing smooth threshold
- Linear scaling $x$ preserving gradient flow
- Parameter $\beta$ adapting to class separation

This is not a hand-designed heuristic—it is the optimal solution derived from first principles.

#### 3.3.2 Backpropagation vs. Correlation

Both training modes (backprop on cross-entropy, explicit correlation objective) converge to similar activation shapes. This suggests that:

- Cross-entropy gradient implicitly optimizes correlation when pre-activations are Gaussian
- The correlation objective is a local approximation to the global loss landscape
- Batch normalization ensures the Gaussian assumption holds in practice

## 4 Discussion

### 4.1 Why Swish Emerges

The Swish activation emerges from our correlation objective for three reasons:

1. **Zero-sum constraint**: The constraint $\sum_k \rho_k = 0$ creates natural competition between classes, requiring the activation to discriminate rather than respond uniformly.

2. **Sigmoid gating**: For moderate class overlap, the correlation function $\rho(x)$ transitions smoothly from negative to positive. Integrating this transition yields sigmoid-like gating.

3. **Linear scaling**: The factor $x$ in $f(x) = x \cdot \sigma(\beta x)$ maximizes gradient flow in deep networks while maintaining the correlation structure.

### 4.2 Connection to Existing Activations

**Swish** [2]: Originally discovered via neural architecture search, Swish has the form $f(x) = x \cdot \sigma(\beta x)$. Our work provides the first principled derivation, showing it emerges from correlation objectives under Gaussian statistics.

**GELU** [1]: The Gaussian Error Linear Unit has a similar shape but is motivated by the Gaussian CDF. Our correlation framework offers an alternative interpretation.

**ReLU**: Emerges as a special case for two classes or well-separated data, where the sigmoid gating collapses to a sharp threshold.

### 4.3 Biological Plausibility

Correlation-based learning has connections to Hebbian principles:

- Local objective: Each unit optimizes its own correlation
- No backpropagation required for activation learning
- Potential model for biological learning mechanisms

However, the zero-sum constraint and normalization are global properties that would require additional biological mechanisms.

## 4.4 Limitations

### 4.4.1 Theoretical Assumptions

- **Gaussian mixtures**: Our closed-form solutions assume Gaussian input distributions. While batch normalization encourages Gaussianity, other data distributions may require different optimal activations.

- **Equal variance**: The assumption of equal class variances simplifies analysis but does not always hold.

- **Piecewise-linear parameterization**: Limits expressiveness compared to smooth functions like splines or neural networks.

- **Single-unit analysis**: Does not capture multi-unit interactions or correlation across neurons.

### 4.4.2 Experimental Scope

- **Limited datasets**: We only tested on CIFAR-10. Validation on ImageNet and other domains is needed.

- **Architecture coverage**: VGG and ResNet are tested; Transformers and other modern architectures remain future work.

- **Baseline comparisons**: We did not compare to other learnable activation methods beyond fixed ReLU.

- **Local correlation training**: The two-optimizer approach requires further stabilization and tuning.

### 4.4.3 Computational Considerations

Memory usage scales with $N_{\text{knots}} \times$ batch size $\times$ channels, making per-channel activations expensive for large models.

## 4.5 Future Directions

### 4.5.1 Theoretical Extensions

- Extend to non-Gaussian distributions (Laplace, heavy-tailed)

- Multi-unit correlation objectives capturing neuron dependencies

- Information-theoretic interpretation (mutual information maximization)

- Formal proof of Swish optimality under general conditions

### 4.5.2 Experimental Directions

- Scale to ImageNet and other large-scale vision benchmarks

- Apply to language models (BERT, GPT architectures)

- Ablation studies: knot count, regularization, $\sigma$ regimes

- Compare to NAS-discovered activations (Mish, FReLU)

- Transfer learned activations across tasks

Figure 1: **Two-Class Case**: ReLU-like activation emerges for binary classification ($K = 2$). Near-zero for negative inputs, linear for positive inputs. The zero-sum constraint creates a threshold detector.

Figure 2: **Multi-Class Case**: Optimal activation function $f(x)$ and class-wise correlations $\rho_k$ for $K = 20$ Gaussian mixture ($\sigma = 1.0$). Swish-like shape emerges with smooth transition at $x \approx 0$.

### 4.5.3   Applications

- **Meta-learning**: Learn activation priors from multiple tasks

- **Neural architecture search**: Co-optimize architecture and activations

- **Efficient training**: Use correlation objective to pre-train activations

- **Biological modeling**: Test correlation-based learning in neuroscience experiments

## 5   Conclusion

We introduced a correlation-based framework for learning activation functions from data and objectives. Our main result is a principled derivation of the Swish activation $f(x) = x \cdot \sigma(\beta x)$ from first principles, showing it emerges as the optimal solution for Gaussian mixture inputs with moderate class separation.

Through unit-level Gaussian mixture experiments and neural network training on CIFAR-10, we validated that both standard backpropagation and explicit correlation objectives converge to Swish-like shapes. The optimal activation adapts to input statistics: ReLU-like for binary or well-separated classes, smooth Swish for moderate overlap, and broad linear for highly overlapping classes.

Our work provides theoretical grounding for widely-used activation functions and opens new directions for principled activation design. The connection between correlation objectives and cross-entropy gradients suggests that implicit objectives guide activation learning even in standard supervised training.

### Reproducibility Statement

Code for all experiments will be made available upon publication. Gaussian mixture experiments use analytic formulas with deterministic solutions. Neural network experiments use standard PyTorch and CIFAR-10 dataset with seeded random number generators for reproducibility.

## References

[1] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

[2] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.

Figure 3: **Sigma Regime Comparison**: Optimal activations across different variance regimes. Small $\sigma$ (0.6): sharp and selective; mid $\sigma$ (1.0): smooth swish; large $\sigma$ (10.0): broad and linear.

Figure 4: **Swish Fit Quality**: Learned activation (blue) fitted to swish family (red). Excellent match with $f(x) \approx x \cdot \sigma(2.5x)$ for mid-range $\sigma$.

Figure 5: **Neural Network Activation Evolution**: Learned activation functions across layers in ResNet-10 after 20 epochs. All layers converge to swish-like shapes with layer-specific parameters. Early layers are broader ($\beta \approx 1.5$), late layers steeper.

Figure 6: **Pre-Activation Histograms**: Per-class pre-activation distributions for ResNet-10 layer 0. Approximately Gaussian with separated means, validating theoretical assumptions.