

# Optimizing neural networks for calibration during training

Recent research reveals that modern neural networks, despite achieving high accuracy, systematically produce overconfident predictions that don't reflect true probabilities. (ACM Digital Library +5) Traditional post-hoc calibration methods like temperature scaling, (GitHub) while effective, only address symptoms rather than root causes. (Stack Exchange) This report presents cutting-edge techniques for training inherently well-calibrated neural networks through modified objectives, differentiable calibration metrics, and specialized training strategies.

## The theoretical foundation: Why cross-entropy fails

Cross-entropy loss, the standard for neural network training, inherently encourages overconfidence through its connection to maximum likelihood estimation. (ACM Digital Library +2) The loss function minimizes  $-\log(p_{\text{correct}})$ , creating an incentive to push predicted probabilities toward 1 for the correct class. This logarithmic nature makes the loss more sensitive to changes in low-confidence predictions, encouraging models to quickly resolve uncertainty by increasing confidence. (Medium) Even after achieving optimal classification accuracy, models continue reducing negative log-likelihood by becoming more confident, creating a disconnect between minimizing classification error and producing calibrated probabilities.

Research shows this overconfidence stems from the gradient structure of cross-entropy:  $\partial L / \partial z_k = p_k - y_k$ , which creates stronger updates for higher-confidence incorrect predictions. (Medium) Combined with high model capacity and batch normalization, this leads to systematically miscalibrated models. (arXiv) (arXiv Vanity) The problem is exacerbated in modern architectures where increased depth and width allow models to more easily overfit to training distributions. (arXiv)

## Modified training objectives beyond cross-entropy

**Focal loss emerges as the most effective single modification** for improving calibration. Originally designed for object detection, focal loss modifies cross-entropy by down-weighting easy examples:

$$FL(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$

For calibration, optimal  $\gamma$  values range from 1.0–3.0 (lower than object detection's typical 5.0). The focusing parameter acts as entropy regularization, reducing overconfidence while maintaining accuracy. Empirical results show **30-50% ECE reduction** compared to standard cross-entropy.

(Jishnu Mukhoti +3)

The recently proposed **Dual Focal Loss** addresses both over and under-confidence: (arXiv)

$$DFL = -\alpha_1(1 - p_{gt})^{\gamma_1} \log(p_{gt}) - \alpha_2(1 - p_{sec})^{\gamma_2} \log(1 - p_{sec})$$

This formulation maximizes the gap between ground truth and highest non-ground truth logits, achieving superior calibration by balancing confidence across all predictions. (arXiv +3)

**Margin-based losses** directly optimize calibration metrics. The Maximum Mean Calibration Error (MMCE) loss uses kernel mean embeddings to create a naturally differentiable calibration measure:

(GitHub +5)

$$\text{MMCE}(f) = \|\mu_P - \mu_Q\|_H^2$$

When combined with cross-entropy ( $L_{\text{total}} = L_{\text{CE}} + \beta \cdot \text{MMCE}$ ), typical  $\beta$  values of 8.0-9.0 maintain high-confidence predictions while reducing miscalibration. (GitHub +2) MMCE provides theoretical guarantees—it equals zero if and only if the model is perfectly calibrated. (Efs-opensource +2)

## Multi-objective optimization balancing accuracy and calibration

Modern approaches treat accuracy and calibration as competing objectives requiring careful balance.

**Weighted multi-objective formulations** combine losses:

$$L_{\text{multi}} = \lambda_1 \cdot L_{\text{accuracy}} + \lambda_2 \cdot L_{\text{calibration}} + \lambda_3 \cdot L_{\text{uncertainty}}$$

Where  $L_{\text{accuracy}}$  uses cross-entropy or focal loss,  $L_{\text{calibration}}$  employs ECE or MMCE, and  $L_{\text{uncertainty}}$  adds entropy penalties. Weight selection strategies include fixed validation-based weights, adaptive scheduling during training, or meta-learning approaches. (Number Analytics)

**Pareto optimization** reveals explicit trade-offs between objectives. Rather than arbitrary weight selection, this approach generates a frontier of solutions, allowing practitioners to select based on application requirements. For critical applications like medical diagnosis, this enables informed decisions about accuracy-calibration trade-offs. (NCBI)

## Making ECE differentiable for direct optimization

Expected Calibration Error traditionally uses hard binning operations that block gradient flow.

(Lightning AI +3) Recent innovations enable direct ECE optimization through **soft binning approaches**:

(arXiv)

The Differentiable ECE (DECE) replaces indicator functions with continuous approximations using temperature-scaled sigmoid functions or Gaussian kernels for smooth probability assignments.

(GitHub) (OpenReview) This maintains differentiability while approximating standard ECE behavior. The Meta-Calibration framework demonstrates how DECE enables hyperparameter optimization for calibration. (GitHub +2)

Implementation requires careful handling of gradient flow through soft bins, with typical temperature values of 0.1-1.0 for sigmoid-based approximations. Despite computational overhead, DECE achieves

competitive results with post-processing methods while training end-to-end.

## Alternative calibration metrics for direct optimization

Beyond ECE, several metrics offer natural differentiability. **MMCE** stands out with its kernel-based approach requiring no binning: [arXiv +2](#)

$$\text{MMCE} = \sqrt{(\sum_{i,j \in D} 1/N^2 (c_i - r_j)(c_j - r_i)k(r_i, r_j))} \quad \text{Efs-opensource}$$

Using Laplacian kernels  $k(r_i, r_j) = \exp(-2.5|r_i - r_j|)$ , MMCE provides consistent convergence rates and fast estimation. [Efs-opensource](#) [mlr](#) Batch sizes around 100 suffice for reliable estimation with only 10% training overhead. [Proceedings of Machine Lear...](#)

**KS-ECE** eliminates binning entirely through Kolmogorov-Smirnov statistical tests, comparing cumulative distributions via differentiable spline approximations. This approach consistently outperforms traditional ECE methods while avoiding binning bias. [ResearchGate](#) [OpenReview](#)

The **Brier Score decomposition** offers another path: Brier = Reliability - Resolution + Uncertainty. [Wikipedia](#) The reliability component directly measures calibration error and can be optimized independently as a differentiable loss. [Wikipedia](#)

## Label smoothing: Simple yet powerful

Label smoothing remains one of the most effective calibration techniques, modifying targets as:

[GeeksforGeeks](#)

$$y'_i = (1 - \epsilon)y_i + \epsilon/K$$

**Optimal parameters vary by dataset:**  $\epsilon=0.1$  works universally well, with 0.05-0.1 for smaller datasets (CIFAR) and 0.1-0.2 for larger ones (ImageNet). [github](#) The technique reduces overconfidence by preventing extreme predictions and encouraging smaller logit gaps between classes. [arXiv](#)

Advanced variants include **Confidence-Calibrated Label Smoothing** that adapts smoothing based on prediction confidence, and **class-dependent smoothing** that customizes parameters for different class difficulties. [ScienceDirect](#) [Ophthalmologyscience](#) When combined with focal loss, label smoothing provides robust calibration improvements with negligible computational overhead.

## Beta calibration and histogram binning during training

**Beta calibration** uses the flexibility of Beta distributions to handle skewed confidence distributions common in neural networks. [Proceedings of Machine Lear...](#) By fitting Beta parameters jointly with model training, this approach addresses calibration without post-processing. The method particularly excels with naturally skewed outputs from certain architectures. [Proceedings of Machine Lear...](#)

**Differentiable histogram methods** approximate hard-binning operations through clever engineering. One approach uses "relu at 1" threshold functions with approximation error of just 0.000158 compared to NumPy histograms. Alternative methods employ CNN operations for continuous histogram layers with learnable bin centers and widths. (arXiv +2)

## Training strategies alternating objectives

Successful calibration-aware training often employs **strategic alternation** between objectives:

**Epoch-based switching** trains for accuracy initially, then introduces calibration objectives:

(Wiley Online Library)

```
python
```

```
for epoch in range(total_epochs):
    if epoch < accuracy_epochs:
        loss = ce_loss(logits, targets)
    else:
        loss = ce_loss(logits, targets) + beta * calibration_loss(logits, targets)
```

**Interleaved training** alternates objectives more frequently, using separate validation sets for calibration evaluation. This prevents overfitting to either objective while maintaining both accuracy and calibration throughout training. (arXiv)

**Gradient accumulation approaches** compute gradients for both objectives separately, then combine with appropriate weighting before parameter updates. This ensures neither objective dominates early training dynamics.

## PyTorch implementations and practical code

Several high-quality implementations enable immediate adoption:

**MDCA-Calibration** (<https://github.com/mdca-loss/MDCA-Calibration>) provides complete PyTorch implementation of Multi-Class Difference in Confidence and Accuracy loss:

```
bash
```

```
python train.py --dataset cifar10 --model resnet56 --loss FL+MDCA --gamma 1.0 --beta 1.0
```

**Focal Calibration** ([https://github.com/torrvision/focal\\_calibration](https://github.com/torrvision/focal_calibration)) implements focal loss variants with adaptive gamma selection and supports multiple calibration losses including MMCE and Brier score.

(github)

**TorchMetrics** offers differentiable calibration metrics: [Lightning AI](#)

```
python
```

```
from torchmetrics.classification import CalibrationError
metric = CalibrationError(task="multiclass", num_classes=10, norm='l1')
ece = metric(preds, targets)
```

Complete training loops typically follow this pattern:

```
python
```

```
def train_calibration_aware(model, train_loader, val_loader):
    for epoch in range(epochs):
        for batch in train_loader:
            logits = model(batch.x)
            ce_loss = F.cross_entropy(logits, batch.y)
            cal_loss = compute_mmce(F.softmax(logits, dim=1), batch.y)
            total_loss = ce_loss + beta * cal_loss
            total_loss.backward()
            optimizer.step()
```

## Key implementation insights

Through extensive research and experimentation, several best practices emerge. **Start with focal loss ( $\gamma=2.0$ )** as it provides consistent improvements with minimal overhead. [Jishnu Mukhoti](#) Add **label smoothing ( $\epsilon=0.1$ )** for additional regularization. For applications requiring maintained high-confidence predictions, use **MMCE with  $\beta=8.0-9.0$** . [GitHub](#)

Monitor calibration throughout training using reliability diagrams and multiple metrics (ECE, MCE, Brier score). [Wttech +2](#) The computational overhead varies: focal loss adds less than 5%, MMCE increases training time by 10-20%, while multi-objective approaches can double training time.

Critical hyperparameters require careful tuning. Focal loss  $\gamma$  values of 1.0-3.0 work well for calibration (lower than object detection). [Jishnu Mukhoti](#) MMCE coefficients between 8.0-9.0 balance calibration with maintained confidence. [GitHub](#) [GitHub](#) Label smoothing around 0.1 proves robust across architectures and datasets. [github](#)

## Conclusion

Training neural networks for calibration requires fundamental changes to standard practices. By understanding why cross-entropy produces overconfident models and implementing calibration-aware objectives, we can train models that provide reliable uncertainty estimates from the start. The

combination of focal loss, label smoothing, and differentiable calibration metrics represents the current state-of-the-art, with implementations readily available for immediate adoption. As neural networks deploy in increasingly critical applications, these techniques become essential for trustworthy AI systems.