

Deep Learning Techniques for Image Classification

Samuel R. Patel

GMIU, Bhavnagar

<http://doi.org/10.64643/JATIRVIII-140022-001>

Abstract- Image classification has rapidly advanced due to deep learning. This paper surveys and compares core deep learning techniques used for image classification — convolutional neural networks (CNNs), residual networks, efficient architectures, and Vision Transformers (ViTs) — and presents a reproducible experimental framework for fair comparison. We explain architecture design choices, training procedures, common datasets, and evaluation metrics. We propose a set of controlled experiments (data preprocessing, augmentation, training hyperparameters) and provide a baseline implementation template. Finally, we analyze strengths and weaknesses of each approach and give practical recommendations for researchers and practitioners. All text in this manuscript is original.

Index-Terms: image classification, convolutional neural networks, ResNet, EfficientNet, Vision Transformer, data augmentation, reproducible experiments

I. INTRODUCTION

Automated image classification—assigning a label to an input image—has been transformed by deep learning. Since the emergence of deep convolutional networks, performance on standard benchmarks has improved dramatically. Today’s powerful models include deep CNN variants with architectural innovations (e.g., residual connections, depthwise separable convolutions), efficient scaling strategies, and transformer-based models that borrow attention mechanisms from natural language processing.

Despite progress, practitioners face trade-offs between accuracy, computational cost, data efficiency, and robustness. This paper aims to (1) summarize key deep learning techniques for image classification, (2) present a reproducible experimental setup to compare representative models, and (3) provide practical guidance on choosing and training models for different constraints (compute, dataset size, latency).

II. BACKGROUND AND KEY CONCEPTS

2.1 Convolutional Neural Networks (CNNs)

CNNs use convolutional layers to learn spatially local features. Early, widely-used blocks include convolution + nonlinearity (ReLU), pooling, and batch normalization. Deeper networks learn hierarchical features from edges to object parts.

2.2 Residual Connections

Residual (skip) connections add inputs to outputs of blocks, enabling training of very deep networks by alleviating vanishing gradients. Residual networks (ResNets) let optimization focus on residual functions and have become a backbone in many vision tasks.

2.3 Efficient Architectures and Scaling

Designs such as depthwise separable convolutions (used in MobileNet) and compound scaling (used in EfficientNet) enable high accuracy with lower FLOPs and parameters, important for mobile and embedded deployment.

2.4 Vision Transformers (ViT)

ViTs split an image into patches, embed them, and apply transformer encoders with self-attention. They can achieve competitive performance, especially when pretrained on large datasets. ViTs tend to be data-hungry but scale well.

2.5 Training Techniques

Key training strategies that strongly affect outcomes include:

- Data augmentation (random crops, flips, color jitter, CutMix, MixUp)
- Regularization (weight decay, dropout, stochastic depth)
- Learning rate schedules (cosine decay, step decay, warmup)
- Optimization algorithms (SGD with momentum, AdamW)
- Transfer learning and fine-tuning from pretrained models

III. REPRESENTATIVE ARCHITECTURES (OVERVIEW)

3.1 Baseline CNN — SimpleNet (example)

A straightforward CNN with repeated Conv-BatchNorm-ReLU blocks and pooling. Useful as a baseline for small datasets.

3.2 AlexNet-style / VGG-style

Deeper stacks of convolutions followed by fully connected layers. VGG-style networks are easy to implement but computationally heavy.

3.3 ResNet family

Residual blocks (basic or bottleneck) enable training of 18–152 layer networks. ResNet-50 is a common strong baseline.

3.4 EfficientNet family

Compound scaling balances width, depth and resolution for better accuracy per FLOP. EfficientNet variants (B0–B7) trade off size and performance.

3.5 MobileNet / MobileNetV2 / MnasNet

Mobile-first designs use depthwise separable convolutions and inverted residuals for mobile deployment.

3.6 Vision Transformer (ViT) and hybrid CNN+ViT

ViT processes non-overlapping patches with transformer encoders. Hybrid models use convolutional stems to improve inductive bias on small datasets.

IV. EXPERIMENTAL DESIGN (REPRODUCIBLE FRAMEWORK)

4.1 Goals

- Compare accuracy-efficiency trade-offs across architectures.
- Use a single, consistent training pipeline to limit confounding variables.
- Provide reproducible hyperparameters and dataset splits.

4.2 Datasets

Suggested datasets for experimentation:

- CIFAR-10 / CIFAR-100: small-scale, 32×32 images, good for quick experimentation.
- Tiny ImageNet: 200 classes, 64×64 images, moderate scale.
- ImageNet-1K (ILSVRC2012): large-scale (1,000 classes), for state-of-the-art comparisons and transfer learning.
- Domain-specific datasets: when evaluating real-world use cases (medical, satellite), include domain datasets and note required ethical approvals.

4.3 Preprocessing & augmentation

- Resize/center crop images to network input resolution (e.g., 224×224).
- Normalize per-channel mean and std from training set or standard ImageNet statistics.
- Use augmentations: random horizontal flip, random resized crop, color jitter. Optionally use MixUp or CutMix for improved generalization.
- For small datasets, stronger augmentations and regularization prevent overfitting.

4.4 Training hyperparameters (recommended baseline)

- Optimizer: SGD with momentum = 0.9 (or AdamW for transformers)
- Learning rate: `base_lr` (e.g., 0.1 for batch size 256) with linear scaling rule, warmup for first 5 epochs, cosine annealing schedule
- Weight decay: 1e-4 (CNNs), 0.05 for AdamW (transformers)
- Batch size: choose by GPU memory (try 128 or 256 where possible)
- Epochs: 90 for ImageNet-like runs; 200–300 for smaller datasets when appropriate
- Label smoothing: 0.1 optional
- Mixed precision training (FP16) to reduce memory and speed training

4.5 Evaluation metrics

- Top-1 and Top-5 accuracy (for ImageNet)
- Precision, recall, F1-score for class-imbalanced tasks
- Inference latency (ms) and throughput (images/sec)
- Model size (parameters) and FLOPs (compute cost)
- Calibration metrics (expected calibration error) for models where probabilistic predictions matter

4.6 Cross-validation and statistical significance

- For small datasets, perform k-fold cross validation or repeated runs with different seeds to assess variance.
- Report mean \pm standard deviation of metrics across runs.

V. IMPLEMENTATION TEMPLATE (PSEUDOCODE / COMMANDS)

Below is an example training outline you can adapt (PyTorch-like pseudocode):

```
# Pseudocode outline (PyTorch)
model = create_model('resnet50', pretrained=False, num_classes=NUM_CLASSES)
optimizer = SGD(model.parameters(), lr=BASE_LR, momentum=0.9, weight_decay=1e-4)
scheduler = CosineAnnealingLR(optimizer, T_max=EPOCHS)
criterion = CrossEntropyLoss(label_smoothing=0.1)

for epoch in range(EPOCHS):
    model.train()
    for images, labels in train_loader:
        images = augment(images) # random crop, flip, MixUp/CutMix optional
        outputs = model(images)
        loss = criterion(outputs, labels)
        loss.backward()
        clip_grad_norm_(model.parameters(), max_norm=5.0)
        optimizer.step()
```

```

optimizer.zero_grad()
scheduler.step()
validate(model, val_loader)

```

Notes:

- Use `torch.cuda.amp` for mixed-precision.
- Save the best model by validation Top-1 accuracy.
- For ViT models, use AdamW and a longer warmup.

VI. EXAMPLE EXPERIMENT PLAN

Conduct the following controlled experiments:

1. Baseline comparison on CIFAR-100
 - Models: SimpleNet, ResNet-34, EfficientNet-B0, ViT-small (pretrained on ImageNet if available)
 - Same augmentation pipeline, epochs = 200, SGD for CNNs, AdamW for ViT
 - Report Top-1 accuracy, parameters, FLOPs.
2. Scaling experiment on Tiny ImageNet
 - Vary input resolution (64, 96, 128) and observe accuracy vs. compute.
 - Test EfficientNet compound scaling vs. ResNet scaling depth/width.
3. Transfer learning on limited data
 - Fine-tune ImageNet-pretrained ResNet-50 and ViT on a small domain dataset (e.g., 1,000 labeled images). Measure sample efficiency.
4. Robustness to distribution shift
 - Test on corruptions (e.g., gaussian noise, blur) or dataset shifts to assess robustness and calibration.

For each experiment, run 3 seeds and report mean \pm std.

VII. ANALYSIS AND EXPECTED OBSERVATIONS (GUIDANCE)

- ResNets: Strong, reliable baselines; perform well across dataset sizes; good generalization when trained properly.
- EfficientNets: Provide better accuracy-per-FLOP; often preferable when inference budget is limited.
- MobileNets: Best for extremely constrained devices; trade accuracy for small size.
- ViTs: Shine at large-scale pretraining; may underperform CNNs when trained from scratch on small datasets but can be competitive with appropriate data augmentation and pretraining.
- Augmentation & regularization: Often produce larger gains than minor architecture tweaks—invest heavily in augmentation pipelines.

- Hyperparameter sensitivity: Transformers are sensitive to learning rates, weight decay, and warmup; CNNs are sensitive to batch size and learning rate scaling.

VIII. PRACTICAL RECOMMENDATIONS

- For research prototyping: start with ResNet-50 or EfficientNet-B0 pretrained on ImageNet; fine-tune with strong augmentations.
- For mobile deployment: prefer MobileNetV2/EdgeTPU-optimized networks or quantized EfficientNet-lite variants.
- For high accuracy at scale: ViT or hybrid models pretrained on very large datasets often yield top performance.
- For limited labels: use transfer learning from large pretrained models and apply semi-supervised techniques (pseudo-labeling, consistency regularization).
- For explainability & calibration: add temperature scaling and evaluate expected calibration error; consider attention-based visualizations or Grad-CAM for interpretability.

IX. ETHICAL CONSIDERATIONS

- Dataset biases can lead to unfair or unsafe predictions; evaluate models across demographic slices.
- For sensitive domains (medical, legal), obtain domain expert review and appropriate approvals before deployment.
- Consider privacy-preserving techniques (federated learning, differential privacy) when training on private data.

X. CONCLUSION

Deep learning for image classification involves many architecture and training choices. Residual CNNs remain strong, EfficientNets offer favourable efficiency, and Vision Transformers bring a powerful attention-based paradigm when sufficient data or pretraining is available. Careful experimental design, augmentation, and hyperparameter tuning are often decisive. The reproducible framework and templates provided here can be used to benchmark models and guide real-world selection based on constraints.

REFERENCE

- [1] LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. — *Gradient-based learning applied to document recognition* (1998).
- [2] Krizhevsky, A., Sutskever, I., & Hinton, G. — *ImageNet classification with deep convolutional neural networks* (AlexNet).

- [3] He, K., Zhang, X., Ren, S., & Sun, J. — *Deep residual learning for image recognition* (ResNet).
- [4] Tan, M., & Le, Q. — *EfficientNet: Rethinking model scaling for convolutional neural networks*.
- [5] Dosovitskiy, A., et al. — *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale* (ViT).
- [6] Howard, A.G., et al. — *MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications*.