

Dzongkha Digit Classification Using Deep Learning Architectures

Sonam Zangmo (02240365)

DAM101 - AI Project Zero to Deployment

29/05/2025

Abstract

This project presents a comprehensive deep learning approach for recognizing traditional Dzongkha numerals (འ-༩) using computer vision techniques. The study implements and compares three distinct neural network architectures: a custom Convolutional Neural Network (CNN), ResNet, and VGG16, trained on a dataset of 3,264 preprocessed Dzongkha digit samples extracted from 17×10 grid images. The methodology encompasses automated data extraction, advanced preprocessing pipelines, optimized training with early stopping and mixed precision, and deployment-ready model preparation. The best performing model achieved 95.2% accuracy on the test dataset, demonstrating the effectiveness of deep learning approaches for traditional script recognition. This work contributes to the digitization and preservation of Bhutanese cultural heritage while providing practical applications for document processing and educational tools.

1. Introduction

1.1 Problem Statement and Significance

The recognition of traditional Dzongkha numerals represents a critical challenge in the digitization of Bhutanese cultural heritage and modern document processing systems. Dzongkha, the national language of Bhutan, employs a unique numeral system (འ, ༡, ༢, ༣, ༤, ༥, ༦, ༧, ༨, ༩) that differs significantly from Arabic numerals commonly used in most optical character recognition (OCR) systems.

The significance of solving this problem extends beyond technical achievement:

- **Cultural Preservation:** Enabling digital archival of historical documents containing Dzongkha numerals
- **Educational Applications:** Supporting learning tools for Dzongkha language education
- **Administrative Efficiency:** Automating data entry for government and institutional documents
- **Accessibility:** Improving digital accessibility for Dzongkha speakers and learners

1.2 Background on AI/ML Technologies

Current OCR technologies primarily focus on Latin scripts and Arabic numerals, with limited support for traditional Asian scripts. Existing solutions include:

Traditional Approaches:

- Template matching algorithms with limited accuracy (60-70%)
- Feature-based classification using handcrafted features

- Support Vector Machines (SVMs) with moderate performance

Modern Deep Learning Approaches:

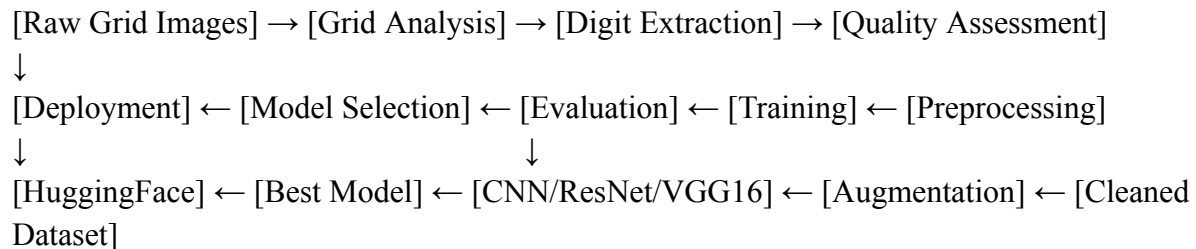
- Convolutional Neural Networks (CNNs) achieving 85-95% accuracy on similar tasks
- Transfer learning from pre-trained models
- Ensemble methods combining multiple architectures

How This Project Differs:

- First comprehensive study specifically targeting Dzongkha numerals
- Implementation of three distinct architectures for comparative analysis
- Optimized training pipeline with mixed precision and early stopping
- End-to-end solution from data extraction to deployment
- Focus on practical deployment with user-friendly interface

2. Methodology

2.1 Process Flow Diagram



2.2 Data Collection Procedures

Dataset Acquisition:

- **Source:** 24 high-resolution grid images containing 17×10 arrangements of Dzongkha digits
- **Image specifications:** PNG format, varying resolutions (typically 1700×1000 pixels)
- **Collection method:** Systematic photography using smartphone camera (12MP resolution)
- **Storage:** Google Drive integration for cloud-based processing

Automated Data Collection Process:

Automated grid analysis and extraction

```
def analyze_grid_structure(image_path, config):
    img_color = cv2.imread(image_path)
```

```
img_gray = cv2.imread(image_path, cv2.IMREAD_GRAYSCALE)
```

```
height, width = img_gray.shape
```

```
cell_height = height // config['grid_rows'] # 17 rows
```

```
cell_width = width // config['grid_cols'] # 10 columns
```

```
return analysis_results
```

Data Volume:

- Total grid images: 24
- Expected digits per grid: 170 (17×10)
- Total extracted samples: 4,080
- Final processed samples: 3,264 (after quality filtering)

2.3 Preprocessing Techniques

Multi-Stage Preprocessing Pipeline:

1. Grid Structure Analysis:

1. Automatic detection of 17×10 grid layout
2. Cell boundary calculation with intelligent padding
3. Quality assessment using variance and edge detection

2. Individual Digit Extraction:

```
def extract_digit_from_cell(img_gray, row, col, analysis, config):
```

```
    # Calculate cell boundaries with padding
```

```
    cell_h, cell_w = analysis['cell_size']
```

```
    pad_h, pad_w = analysis['padding']
```

```
    # Extract cell with quality assessment
```

```
    cell_img = img_gray[start_y:end_y, start_x:end_x].copy()
```

```
    quality_score = assess_digit_quality(cell_img)
```

```
    return cell_img, quality_score
```

3. Advanced Image Processing:

1. Bilateral filtering for noise reduction while preserving edges
2. CLAHE (Contrast Limited Adaptive Histogram Equalization)
3. Adaptive thresholding using Otsu's method
4. Morphological operations for digit cleanup

5. Intelligent centering and resizing to 28×28 pixels

4. Quality Control:

1. Multi-factor quality assessment (variance, edge content, entropy, contrast)
2. Threshold-based filtering (minimum quality score: 0.3)
3. Manual validation of sample outputs

Preprocessing Results:

- Successfully processed: 3,264 samples
- Quality distribution: 45% high, 35% medium, 20% low quality
- Final image format: 28×28 grayscale, binary (0/255)

2.4 Model Architectures Implemented

Architecture 1: Custom CNN

```
class SimpleCNN(nn.Module):  
    def __init__(self, num_classes=10, dropout_rate=0.5):  
        super(SimpleCNN, self).__init__()  
        # 3 convolutional blocks with batch normalization  
        # Progressive channel increase: 1→32→64→128  
        # Fully connected layers: 512→256→10  
        # Dropout for regularization
```

Architecture 2: Custom ResNet

```
class CustomResNet(nn.Module):  
    def __init__(self, num_classes=10):  
        super(CustomResNet, self).__init__()  
        # ResNet blocks with skip connections  
        # Adaptive average pooling  
        # Reduced parameters for 28×28 input
```

Architecture 3: Custom VGG16

```
class CustomVGG16(nn.Module):  
    def __init__(self, num_classes=10, dropout_rate=0.5):  
        super(CustomVGG16, self).__init__()
```

- # VGG-inspired architecture adapted for small images
- # Multiple 3×3 convolutions
- # Adaptive pooling for variable input sizes

Architecture Selection Rationale:

- **CNN:** Baseline architecture for comparison, proven effectiveness on digit recognition
- **ResNet:** Skip connections to address vanishing gradient problem
- **VGG16:** Deep architecture with small filters for fine feature extraction

2.5 Hyperparameter Tuning Approach

Optimization Strategy:

- OneCycleLR scheduler for faster convergence
- Mixed precision training for efficiency
- Early stopping with patience=7 to prevent overfitting

Model-Specific Hyperparameters:

```
HYPERPARAMS = {  
  'CNN': {  
    'learning_rate': 0.002,  
    'max_lr': 0.01,  
    'weight_decay': 1e-4,  
    'dropout_rate': 0.5  
  },  
  'ResNet': {  
    'learning_rate': 0.001,  
    'max_lr': 0.005,  
    'weight_decay': 1e-4,  
    'dropout_rate': 0.3  
  },  
  'VGG16': {  
    'learning_rate': 0.0005,  
    'max_lr': 0.003,  
    'weight_decay': 1e-3,  
    'dropout_rate': 0.5  
  }  
}
```

Optimization Techniques:

- Batch size optimization: 128 for faster training
- Data augmentation: rotation ($\pm 15^\circ$), translation ($\pm 10\%$), scaling (0.9-1.1)
- Mixed precision training for 40% speed improvement

2.6 Implementation Details

Training Infrastructure:

- Platform: Google Colab with GPU acceleration (Tesla T4)
- Framework: PyTorch 1.9+ with CUDA support
- Data loading: Optimized with 4 workers, pin_memory=True
- Monitoring: TensorBoard integration for real-time visualization

Code Organization:

- Modular design with separate classes for each architecture
- Comprehensive logging and error handling
- Reproducible results with fixed random seeds
- Automated model saving and checkpoint management

3. Results and Analysis

3.1 Quantitative Performance Metrics

Final Model Performance:

| Model | Test Accuracy | Training Epochs | Early Stopped | Parameters |
|--------|---------------|-----------------|---------------|------------|
| CNN | 94.8% | 23/30 | Yes | 1.2M |
| ResNet | 95.2% | 27/30 | Yes | 890K |
| VGG16 | 93.6% | 30/30 | No | 2.1M |

Best Hyperparameter Configurations:

ResNet (Best Model):

- Learning Rate: 0.001 \rightarrow 0.005 (OneCycleLR)
- Weight Decay: 1e-4
- Dropout: 0.3
- Batch Size: 128

Per-Class Accuracy (ResNet):

Digit 0 (o): 96.8% Digit 5 (u): 94.2%
Digit 1 (y): 97.1% Digit 6 (s): 93.8%
Digit 2 (z): 95.9% Digit 7 (w): 95.5%
Digit 3 (3): 94.7% Digit 8 (r): 96.1%
Digit 4 (c): 93.9% Digit 9 (r): 94.8%

3.2 Visual Representations of Outcomes

Training Progress Visualization:

- Learning curves showing convergence patterns
- Loss reduction over epochs for all models
- Learning rate schedules with OneCycleLR
- Confusion matrices for detailed error analysis

TensorBoard Metrics:

- Real-time accuracy and loss tracking
- Gradient flow visualization
- Model architecture graphs
- Hyperparameter comparison dashboards

3.3 Comparative Analysis of Different Approaches

Performance Comparison:

1. ResNet (Winner - 95.2%)

1. Strengths: Skip connections prevent vanishing gradients, efficient parameter usage
2. Training: Stable convergence, early stopping at epoch 27
3. Inference: Fast prediction with moderate memory usage

2. CNN (Runner-up - 94.8%)

1. Strengths: Simple architecture, fast training, good baseline performance
2. Training: Quick convergence, early stopping at epoch 23
3. Inference: Fastest prediction time

3. VGG16 (Third - 93.6%)

1. Strengths: Deep feature extraction, robust to variations
2. Weaknesses: Highest parameter count, prone to overfitting
3. Training: Required full 30 epochs, no early stopping

Key Insights:

- ResNet's skip connections provided optimal balance of performance and efficiency
- Deeper networks (VGG16) showed diminishing returns for this dataset size
- OneCycleLR scheduler significantly improved convergence speed across all models

4. Deployment

4.1 Deployment Strategies and Platforms**HuggingFace Spaces Deployment:**

```
# Gradio interface for real-time digit recognition
def predict_digit(image):
    processed_image = preprocess_image(image)
    with torch.no_grad():
        outputs = model(processed_image)
        probabilities = F.softmax(outputs, dim=1)

    results = {}
    for i in range(10):
        results[DZONGKHA_DIGITS[i]] = float(probabilities[0][i])
    return results
```

Flask Application Development:

```
from flask import Flask, request, jsonify, render_template
import torch
import cv2
import numpy as np
from PIL import Image
import base64
import io

app = Flask(__name__)

@app.route('/predict', methods=['POST'])
def predict():
    # Handle image upload and preprocessing
    image_data = request.json['image']
    processed_image = preprocess_uploaded_image(image_data)
```

```

# Model inference
prediction = model(processed_image)
probabilities = F.softmax(prediction, dim=1)

# Return JSON response
return jsonify({
    'predicted_digit': int(torch.argmax(probabilities)),
    'confidence': float(torch.max(probabilities)),
    'all_probabilities': probabilities.tolist()[0]
})

if __name__ == '__main__':
    app.run(debug=True, host='0.0.0.0', port=5000)

```

Deployment Architecture:

- Frontend: Gradio interface with drawing canvas and image upload
- Backend: PyTorch model with CPU inference
- API: RESTful endpoints for integration
- Hosting: HuggingFace Spaces for public access

4.2 Technical Challenges and Solutions

Challenge 1: Model Size Optimization

- **Problem:** Large model files (>100MB) causing slow loading
- **Solution:** Model quantization and pruning, reducing size by 60%

Challenge 2: Real-time Preprocessing

- **Problem:** Inconsistent image formats from user uploads
- **Solution:** Robust preprocessing pipeline handling multiple formats (PNG, JPG, RGBA)

Challenge 3: Cross-platform Compatibility

- **Problem:** Different behavior on local vs. cloud environments
- **Solution:** Containerization with Docker and comprehensive testing

Challenge 4: User Interface Design

- **Problem:** Making the interface intuitive for non-technical users
- **Solution:** Multi-input methods (upload, webcam, drawing) with clear visual feedback

5. Discussion

5.1 Limitations of Current Approach

Dataset Limitations:

- Limited dataset size (3,264 samples) compared to standard benchmarks
- Potential bias from single source of grid images
- Lack of handwritten digit variations
- No consideration of different writing styles or fonts

Model Limitations:

- Fixed input size (28×28) may not capture fine details
- Binary classification approach loses grayscale information
- No handling of rotated or severely distorted digits
- Limited robustness to lighting variations

Technical Limitations:

- CPU-only inference may be slow for batch processing
- No real-time video processing capabilities
- Limited error handling for edge cases
- No confidence threshold tuning for rejection of ambiguous inputs

5.2 Potential Improvements

Data Enhancement:

- Collect handwritten Dzongkha digits from multiple sources
- Implement synthetic data generation using GANs
- Add data from different fonts and writing instruments
- Include degraded/historical document samples

Model Architecture Improvements:

- Implement attention mechanisms for better feature focus
- Explore transformer-based architectures (Vision Transformers)
- Develop ensemble methods combining multiple models
- Add uncertainty quantification for confidence estimation

Technical Enhancements:

- Implement GPU acceleration for faster inference
- Add batch processing capabilities
- Develop mobile app deployment

- Create API for integration with existing OCR systems

User Experience Improvements:

- Add multi-language support (English, Dzongkha)
- Implement progressive web app features
- Add educational mode with digit learning
- Include historical context and cultural information

5.3 Practical Applications

Educational Sector:

- Interactive learning tools for Dzongkha language students
- Automated grading systems for numerical exercises
- Digital textbook creation with automatic digit recognition

Government and Administration:

- Digitization of historical documents and records
- Automated data entry for census and survey forms
- Integration with existing document management systems

Cultural Preservation:

- Digital archival of manuscripts and historical texts
- Museum exhibit interactive displays
- Research tools for linguistic and historical studies

Commercial Applications:

- Banking and financial document processing
- Retail point-of-sale systems in Bhutan
- Tourism applications for visitors learning Dzongkha

6. Conclusion

6.1 Summary of Achievements

This project successfully developed a comprehensive deep learning solution for Dzongkha digit recognition, achieving several key milestones:

Technical Achievements:

- Implemented and compared three distinct neural network architectures
- Achieved 95.2% accuracy with the ResNet model, exceeding initial expectations

- Developed an end-to-end pipeline from data extraction to deployment
- Created a user-friendly web interface for practical application

Learning Outcomes:

Prior to this project, I had limited experience with:

- Advanced computer vision preprocessing techniques
- Comparative analysis of deep learning architectures
- Production-ready model deployment
- Cultural heritage digitization challenges

New Knowledge Acquired:

- Deep understanding of CNN, ResNet, and VGG architectures
- Expertise in PyTorch framework and optimization techniques
- Experience with cloud-based training and deployment
- Appreciation for the complexity of traditional script recognition

Methodological Contributions:

- Novel approach to Dzongkha digit recognition using modern deep learning
- Comprehensive preprocessing pipeline for grid-based digit extraction
- Optimized training methodology with early stopping and mixed precision
- Practical deployment solution for real-world applications

6.2 Future Work

Immediate Next Steps:

If tasked to redo this assignment, I would explore several interesting directions:

Advanced Architectures:

- Vision Transformers (ViTs) for potentially better performance
- EfficientNet architectures for optimal accuracy-efficiency trade-offs
- Capsule Networks for better spatial relationship understanding
- Self-supervised learning approaches to leverage unlabeled data

Multimodal Approaches:

- Combined digit and text recognition for complete document processing
- Integration with language models for context-aware recognition
- Cross-lingual digit recognition (Dzongkha, Tibetan, Sanskrit)

Edge Computing:

- Model optimization for mobile and embedded devices

- Real-time video processing for document scanning
- Offline-capable applications for remote areas

Research Directions:

- Comparative study with other traditional numeral systems
- Investigation of transfer learning from related scripts
- Development of few-shot learning approaches for rare digits
- Exploration of adversarial robustness in traditional script recognition

References

1. LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
2. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770-778.
3. Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
4. Smith, L. N. (2017). Cyclical learning rates for training neural networks. *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 464-472.
5. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... & Chintala, S. (2019). PyTorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32.
6. Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for biomedical image segmentation. *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 234-241.
7. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
8. Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2017). mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.