

AI-DRIVEN OCULAR DISEASE DETECTION

Comprehensive Technical Data Report

Project Date: November 2025

Dataset: ODIR-5K with External Augmentation

Final Model: DenseNet-121 Transfer Learning Architecture

EXECUTIVE SUMMARY

This project successfully developed a multi-label ocular disease classification system using deep learning on fundus images. The final DenseNet-121 model achieved a **Macro F1-Score of 0.7871** and **Test AUC of 0.9074**, significantly outperforming the baseline CNN model (Macro F1: 0.5386, Test AUC: 0.8682).

Key Achievement: The model can simultaneously detect 8 distinct ocular pathologies from a single fundus image with clinically relevant accuracy.

1. BUSINESS UNDERSTANDING

1.1 Project Background

Ocular diseases such as Diabetic Retinopathy, Glaucoma, and Cataracts are leading causes of preventable blindness globally. Current diagnosis relies on manual examination of retinal fundus images by ophthalmologists, facing critical challenges:

- **Scalability Crisis:** Global shortage of ophthalmologists, especially in remote regions
- **Time-Intensive Process:** Manual screening consumes significant specialist time
- **Human Variability:** Subject to fatigue and inter-observer inconsistency

1.2 Problem Statement

The manual screening process is inefficient, unscalable, and inaccessible, leading to preventable vision loss. Healthcare providers require an automated screening tool that can analyze retinal fundus images and accurately identify multiple pathologies simultaneously.

1.3 Project Objectives

Primary Goal: Develop a Clinical Decision Support System (CDSS) that integrates retinal scan data to serve as an automated first-pass screening tool.

Specific Objectives:

- Build a multi-modal CNN model detecting 8 ocular pathologies
- Prioritize patient triage by flagging high-risk images
- Enhance efficiency by automating normal scan screening
- Deploy an accessible web application interface

Success Criteria:

- **Primary Metric:** Mean AUC-ROC ≥ 0.90 (Achieved: 0.9074)
 - **Secondary Metrics:** Per-class F1-Score, Precision, Recall
 - **Deployment:** Functional web-based application with human-readable outputs
-

2. DATA UNDERSTANDING

2.1 Dataset Overview

Metric	Value
Total Images	37,649
Original ODIR-5K Images	6,392
Augmented Images (Dataset 1)	4952
Augmented Images (Dataset 2)	10449
Training Daraset Final	15856
Image Format	RGB (512×512 pixels)
File Type	JPG
Disease Classes	8 (Multi-label)

2.2 Disease Label Distribution

Disease Category	Code	Sample Count	Percentage
Normal	N	4,291	24.8%
Diabetes (Diabetic Retinopathy)	D	2,123	12.3%
Glaucoma	G	397	2.3%
Cataract	C	402	2.3%
AMD (Age-related Macular Degeneration)	A	319	1.8%
Hypertension	H	203	1.2%
Myopia	M	306	1.8%
Other Abnormalities	O	1,588	9.2%

Key Observation: Significant class imbalance exists, with Normal and Diabetes dominating. Rare classes (Glaucoma, Cataract, AMD, Hypertension, Myopia) represent <2.5% each.

2.3 Patient Demographics (ODIR-5K Subset)

Age Distribution

- **Mean Age:** 57.8 years
- **Standard Deviation:** 11.7 years
- **Range:** 14-91 years
- **Median (Q2):** 58 years
- **Interquartile Range:** 51-66 years
- **Distribution:** Near-normal, slightly left-skewed toward older ages

Clinical Insight: Dataset focuses on middle-aged to elderly populations at higher risk for age-related ocular diseases.

Sex Distribution

- **Male:** 53.6% (3,424 patients)
- **Female:** 46.4% (2,968 patients)

Balance Assessment: Well-balanced gender distribution reduces potential model bias.

2.4 Image Quality Analysis

Sample analysis of 500 random images revealed:

Quality Metric	Finding
Most Common Size	512×512 pixels (standardized)
Color Mode	RGB (3 channels)
Format Consistency	100% JPG
Blurry Images (ODIR-5K)	927/6,392 (14.5%)
Corrupted Files	0
Size Variability	High in augmented datasets (requires resizing)

3. DATA PREPARATION

3.1 Data Integration Process

Three datasets were merged using a standardized label mapping system:

Label Map (8-element one-hot encoding):

[N, D, G, C, A, H, M, O]

Example: [1,0,0,0,0,0,0,0] = Normal

[0,1,0,0,0,0,0,0] = Diabetes

Integration Results:

- Original ODIR-5K: 6,392 images → Full paths constructed
- External Dataset 1: 7,986 images → Label-mapped from folder structure
- External Dataset 2: 7,415 images → Label-mapped from folder structure
- **Final Combined Dataset:** 21,793 images (100% with valid labels)

3.2 Data Quality Assurance

File Validation:

- All file paths validated for existence

- Non-existent files removed during generator initialization
- Missing files: 0 after validation

Label Integrity:

- Zero missing values in target labels
- All labels converted to NumPy float32 arrays
- Multi-label format maintained throughout pipeline

3.3 Data Splitting Strategy

Split	Image s	Percentag e	Purpose
Training	13,945	64%	Model learning
Validation	3,486	16%	Hyperparameter tuning
Test	4,362	20%	Final unbiased evaluation

- **Random Seed:** 42 (for reproducibility)
 - **Stratification:** Not applied (multi-label complexity)
 - **Shuffle:** Applied before splitting
-

4. EXPLORATORY DATA ANALYSIS

4.1 Univariate Analysis

Patient Age Distribution

- **Key Finding:** Bell-shaped distribution centered around 58 years
- **Clinical Relevance:** Aligns with age-related disease prevalence
- **Few Young Patients:** <20 years represent minimal proportion

Disease Prevalence

- **Most Common:** Normal (2,873 samples), Diabetes (2,123)
- **Least Common:** Hypertension (203), AMD (319), Myopia (306)
- **Challenge Identified:** Severe class imbalance requiring mitigation strategies

4.2 Bivariate Analysis

Age vs. Disease

Key Patterns:

- Older patients (>50 years) dominate: Diabetes, Cataract, Glaucoma, AMD
- Younger patients more common in: Normal, Myopia
- **Conclusion:** Age strongly linked to disease occurrence

Sex vs. Disease

Key Patterns:

- Males slightly higher: Normal, Diabetes
- Females marginally higher: Cataract, Myopia
- **Conclusion:** Sex is NOT a major determinant in this dataset

Age vs. Sex

- Female mean age: ~59 years
- Male mean age: ~57 years
- **Minimal difference:** Both groups span similar age ranges

4.3 Multivariate Analysis

Disease Co-occurrence Heatmap

Correlation Analysis:

- **Strongest Negative Correlations:**
 - Normal ↔ Diabetes: -0.49
 - Normal ↔ Other: -0.40
- **Interpretation:** Normal class acts as opposite of disease presence
- **Low Positive Correlations:** Most diseases occur independently
- **Multi-disease Cases:** Rare but present (e.g., Hypertension + Diabetes)

Clinical Insight: Ocular diseases rarely co-occur in same patient, except systemic conditions (Diabetes, Hypertension).

5. MODELING APPROACH

5.1 Baseline Model: Standard CNN

Architecture

- **Type:** Convolutional Neural Network from scratch
- **Blocks:** 3 convolutional blocks (32, 64, 128 filters)
- **Structure:** Conv2D → ReLU → MaxPooling2D
- **Classifier:** Flatten → Dense(128, ReLU) → Dropout(0.5) → Dense(8, Sigmoid)
- **Total Parameters:** ~2.5M trainable

Training Configuration

- **Optimizer:** Adam (lr=0.001)
- **Loss Function:** Binary Crossentropy
- **Metrics:** Binary Accuracy, AUC
- **Epochs:** 15 (with EarlyStopping patience=10)
- **Batch Size:** 32

Baseline Performance (Test Set)

Metric	Value
--------	-------

Test Loss	0.1424
-----------	--------

Test Accuracy	0.8911
---------------	--------

Test AUC	0.8682
----------	--------

Macro F1-Score	0.5386
----------------	--------

Micro F1-Score	0.5649
----------------	--------

Per-Class Performance:

Class	Precision	Recall	F1-Score	Support
Normal	0.47	0.92	0.62	1,878
Diabetes	0.49	0.74	0.59	1,762

Class	Precision	Recall	F1-Score	Support
Glaucoma	0.53	0.70	0.60	1,136
Cataract	0.86	0.67	0.75	944
AMD	0.45	0.26	0.33	588
Hypertension	0.55	0.23	0.33	586
Myopia	0.80	0.72	0.76	629
Other	0.45	0.36	0.40	1,076

Baseline Insights:

- Strong performance on Cataract (F1: 0.75), Myopia (F1: 0.76)
 - Weak detection of AMD (F1: 0.33), Hypertension (F1: 0.33)
 - High false positive rate for Normal class (Precision: 0.47)
-

5.2 Final Model: DenseNet-121 Transfer Learning

Architecture Selection

- **Base Model:** DenseNet-121 pre-trained on ImageNet
- **Rationale:** Dense connections improve gradient flow and feature reuse
- **Modification:** Custom classification head for 8-class multi-label output

Custom Classification Head

DenseNet-121 (frozen)

↓

GlobalAveragePooling2D

↓

Dropout(0.5)

↓

Dense(512, ReLU)

↓

Dropout(0.5)

↓

Dense(8, Sigmoid)

Two-Phase Training Strategy

Phase 1: Feature Extraction (5 epochs)

- DenseNet-121 base: FROZEN
- Train only classification head
- Learning Rate: 1e-4
- Purpose: Adapt new head to fundus image features

Phase 2: Fine-Tuning (25 epochs)

- DenseNet-121 base: UNFROZEN (except BatchNorm layers)
- Train entire network
- Learning Rate: 1e-5 (10× lower)
- Purpose: Gentle adaptation of deep layers to ocular dataset

Training Configuration

- **Optimizer:** Adam (adaptive learning rate)
- **Loss Function:** Binary Crossentropy
- **Metrics:** Binary Accuracy, Multi-label AUC
- **Callbacks:**
 - EarlyStopping (monitor='val_loss', patience=3)
 - ModelCheckpoint (save_best_only=True)
- **Total Epochs:** 30 (Phase 1: 5, Phase 2: 25)
- **Batch Size:** 32
- **Image Size:** 224×224×3

6. MODEL EVALUATION RESULTS

6.1 Overall Performance Metrics

Metric	Baseline CNN	DenseNet-121	Improvement
Test Loss	0.1424	0.2475	-73.7% ⚠
Test Accuracy	89.11%	89.11%	0%
Test AUC	0.8682	0.9074	+4.52% ✓
Macro F1-Score	0.5386	0.7871	+46.1% ✓
Weighted F1-Score	0.5810	0.7761	+33.6% ✓
Micro F1-Score	0.5649	0.7827	+38.6% ✓

Key Observation: Despite higher loss (likely due to more aggressive regularization), DenseNet-121 achieves superior discrimination capability (AUC) and balanced performance (F1-scores).

6.2 Per-Class Performance Comparison

Standard Threshold (0.5):

Class	Precision	Recall	F1-Score	Support	vs Baseline F1
Normal	0.75	0.84	0.79	1,813	+0.17 ✓
Diabetes	0.85	0.67	0.75	1,777	+0.16 ✓
Glaucoma	0.78	0.82	0.80	1,151	+0.20 ✓
Cataract	0.93	0.87	0.90	950	+0.15 ✓
AMD	0.84	0.79	0.81	614	+0.48 ✓✓
Hypertension	0.86	0.79	0.83	597	+0.50 ✓✓
Myopia	0.90	0.83	0.86	664	+0.10 ✓
Other	0.80	0.42	0.55	1,058	+0.15 ✓

Weighted Average: Precision=0.82, Recall=0.75, F1=0.78

6.3 Threshold Optimization Results

Optimal Thresholds Identified:

Class	Optimal Threshold	Rationale
Normal	0.514	Slightly stricter to reduce false positives
Diabetes	0.300	More sensitive detection (low threshold)
Glaucoma	0.462	Balanced precision-recall
Cataract	0.445	Already high performance
AMD	0.501	Near-default threshold
Hypertension	0.413	Moderate sensitivity boost
Myopia	0.494	Near-default threshold
Other	0.256	Highly sensitive (lowest threshold)

Performance with Custom Thresholds:

Metric	Standard (0.5)	Custom Improvement	
Macro F1-Score	0.7871	0.7871	0%
Normal F1	0.79	0.79	0%
Diabetes F1	0.75	0.75	0%
Other F1	0.55	0.55	0%
Other Recall	0.42	0.42	0%

Note: Threshold optimization showed minimal improvement in this case, suggesting the model is well-calibrated at 0.5 threshold.

6.4 Multi-Disease Prediction Analysis

Metric	True Labels	Predictions (0.5)	Predictions (Custom)
Multi-disease Samples	1,018	927	1,469

Interpretation:

- Standard threshold: Slightly conservative (underestimates multi-disease)
 - Custom thresholds: More aggressive (overestimates multi-disease)
 - **Clinical Trade-off:** Higher sensitivity valuable for screening applications
-

7. MODEL STRENGTHS & WEAKNESSES

7.1 Best Performing Classes

1. Cataract (F1: 0.90)

- Excellent precision (0.93) and recall (0.87)
- Likely due to distinctive lens opacity patterns
- Clinical utility: High confidence predictions

2. Myopia (F1: 0.86)

- Strong balanced performance
- Well-represented in training data despite small size
- Distinct retinal morphology aids detection

3. Hypertension (F1: 0.83)

- Remarkable improvement over baseline (+0.50)
- Transfer learning effectively learned vascular changes
- Small sample size (597) but high performance

7.2 Challenging Classes

1. Other Abnormalities (F1: 0.55)

- **Low Recall (0.42):** Model misses many cases
- **Challenge:** Heterogeneous class with diverse pathologies
- **Root Cause:** Lack of specific visual patterns
- **Recommendation:** Consider splitting into sub-categories

2. Diabetes (F1: 0.75)

- Good precision (0.85) but moderate recall (0.67)
- **Challenge:** Variable presentation of retinopathy stages

- **False Negatives:** Early-stage DR may lack visible lesions

7.3 Class Imbalance Impact

Observed Patterns:

- Well-represented classes (Normal, Diabetes): Strong performance
 - Rare classes (Hypertension, AMD): Surprisingly good performance
 - **Conclusion:** Transfer learning mitigated class imbalance effectively
-

8. COMPARISON WITH OTHER ARCHITECTURES

During experimentation, multiple architectures were evaluated:

8.1 Tested Models Summary

Model	Test AUC	Macro F1	Key Observation
Baseline CNN	0.8682	0.5386	Weak on rare classes
MobileNetV2	~0.78	~0.78*	Fast but lower accuracy
ResNet50	~0.78	~0.78*	Good balance
InceptionNetB 2	~0.57	~0.56*	Underperformed expectations
DenseNet-121	0.9074	0.7871	Best overall

*Approximate values from classification report images

8.2 Why DenseNet-121 Excelled

1. Dense Connections:

- Every layer receives feature maps from all previous layers
- Improves gradient flow during backpropagation
- Reduces vanishing gradient problem

2. Feature Reuse:

- Lower layers learn basic features (edges, textures)
- Higher layers access and reuse these features

- More efficient than ResNet's residual connections

3. Parameter Efficiency:

- Fewer parameters than ResNet50 (~8M vs ~25M)
- Lower overfitting risk
- Faster training and inference

4. Medical Image Performance:

- Proven track record on fundus image datasets
- Dense connections ideal for subtle pathology detection
- Better gradient propagation for small lesions

9. CLINICAL IMPLICATIONS

9.1 Deployment Readiness Assessment

Criterion	Status	Notes
Primary Metric (AUC ≥ 0.90)	✓ ACHIEVED	0.9074 exceeds target
Balanced Performance	✓ GOOD	Macro F1: 0.79 across 8 classes
Rare Disease Detection	✓ ACCEPTABLE	AMD, Hypertension F1 >0.80
Safety (High Recall)	⚠ CAUTION	Other class recall only 0.42
Clinical Validation	⚠ PENDING	Requires ophthalmologist review

9.2 Recommended Use Cases

✓ Appropriate Applications:

1. **Pre-screening Triage:** Flag abnormal cases for urgent review
2. **Normal Scan Filtering:** Automate healthy eye identification
3. **Second Opinion Tool:** Support junior ophthalmologists
4. **Remote Screening:** Enable diagnosis in underserved areas

X Not Recommended:

1. **Sole Diagnostic Tool:** Should not replace expert examination
2. **Treatment Planning:** Insufficient detail for surgical decisions
3. **Medico-legal Applications:** Model limitations require human oversight

9.3 Performance by Use Case

Screening Triage (High Sensitivity):

- Use custom low thresholds (e.g., 0.3 for Diabetes)
- Prioritize recall over precision
- **Trade-off:** More false positives, but no missed cases

Diagnostic Support (Balanced):

- Use standard 0.5 threshold
- Optimize F1-score for balanced errors
- **Trade-off:** Acceptable false positives and negatives

Specialist Workload Reduction (High Specificity):

- Use higher thresholds (e.g., 0.7) for normal cases
- Only flag high-confidence abnormalities
- **Trade-off:** Some abnormalities may be missed

10. LIMITATIONS & RECOMMENDATIONS

10.1 Current Limitations

Data-Related:

1. **Class Imbalance:** Rare classes <2% of dataset
2. **Geographic Bias:** ODIR-5K primarily Asian patients
3. **Image Quality:** 14.5% blurry images in original ODIR-5K
4. **Label Quality:** "Other" class is heterogeneous catch-all

Model-Related:

1. **Black Box Nature:** Limited interpretability despite Grad-CAM potential

2. **Multi-Disease Prediction:** Tends to under-predict co-morbidities
3. **Threshold Sensitivity:** Performance varies with probability cutoff
4. **Computational Cost:** Requires GPU for real-time inference

Clinical-Related:

1. **No Severity Grading:** Binary presence/absence only
2. **Limited Context:** No patient history integration
3. **Validation Gap:** Not tested on external clinical datasets
4. **Regulatory Status:** Not FDA/CE approved

10.2 Future Improvements

Short-Term (3-6 months):

1. **Implement Grad-CAM Visualization:**
 - Provide heatmaps showing which image regions influenced predictions
 - Build clinician trust through explainability
2. **Expand "Other" Class:**
 - Split into sub-categories (e.g., Retinal Detachment, Macular Edema)
 - Collect more diverse abnormality examples
3. **External Validation:**
 - Test on Messidor-2, EyePACS, APTOS datasets
 - Assess generalization to different populations

Medium-Term (6-12 months):

1. **Multi-Modal Integration:**
 - Incorporate patient metadata (age, hypertension status)
 - Fuse image + clinical features for improved accuracy
2. **Severity Grading:**
 - Train ordinal classification for DR stages (mild, moderate, severe, PDR)
 - Implement AMD severity levels (early, intermediate, advanced)
3. **Active Learning Pipeline:**

- Collect uncertain predictions for expert labeling
- Continuously improve model with feedback loop

Long-Term (1-2 years):

1. Prospective Clinical Trial:

- Deploy in real-world clinic settings
- Measure impact on patient outcomes and workflow efficiency

2. Model Ensemble:

- Combine DenseNet-121 with ResNet50 and EfficientNet
- Boost performance through prediction averaging

3. Edge Deployment:

- Optimize model for mobile/tablet devices
- Enable point-of-care screening in remote areas

4. Regulatory Approval:

- Pursue FDA 510(k) or De Novo pathway
 - Obtain CE marking for European deployment
-

11. TECHNICAL SPECIFICATIONS

11.1 Hardware Requirements

Training Environment:

- **GPU:** NVIDIA Tesla P100 (16GB VRAM) or equivalent
- **RAM:** 32GB minimum
- **Storage:** 50GB for datasets + models
- **Training Time:** ~4 hours (30 epochs, 21K images)

Inference Environment:

- **GPU:** NVIDIA GTX 1060 (6GB) or higher (recommended)
- **CPU Fallback:** Possible but 10-20x slower
- **RAM:** 8GB minimum

- **Storage:** 500MB for model weights

11.2 Software Dependencies

```
# Core Deep Learning
```

```
tensorflow==2.13.0
```

```
keras==2.13.0
```

```
numpy==1.24.3
```

```
# Image Processing
```

```
Pillow==10.0.0
```

```
opencv-python==4.8.0
```

```
# Data Manipulation
```

```
pandas==2.0.3
```

```
scikit-learn==1.3.0
```

```
# Visualization
```

```
matplotlib==3.7.2
```

```
seaborn==0.12.2
```

11.3 Model Artifacts

Saved Files:

- densenet121_best_model_phase2.keras.weights.h5 (230MB)
- my_ocular_model_densenet121.keras (235MB, full model)
- optimal_thresholds.json (custom thresholds per class)

Input Specifications:

- **Format:** RGB image (JPEG/PNG)
- **Size:** 224×224 pixels (resized automatically)
- **Normalization:** Pixel values ÷ 255 (range [0,1])

- **Color Space:** RGB (not BGR)

Output Format:

```
{  
    "predictions": [  
        {"class": "Normal", "probability": 0.823, "threshold": 0.514, "positive": true},  
        {"class": "Diabetes", "probability": 0.234, "threshold": 0.300, "positive": false},  
        # ... 6 more classes  
    ],  
    "multi_disease": false,  
    "confidence": "high" # if max_prob > 0.8  
}
```

12. REPRODUCIBILITY

12.1 Random Seeds

All random operations used seed=42:

- NumPy random state
- TensorFlow random seed
- Train/validation/test split
- Data generator shuffling

12.2 Training Reproducibility Checklist

- ✓ Fixed random seeds
- ✓ Documented hyperparameters
- ✓ Version-controlled code
- ✓ Saved best model weights
- ✓ Recorded training history
- ✓ Documented data preprocessing steps

12.3 Known Variability Sources

- GPU computation order (CUDA non-determinism)
- Data augmentation randomness (if enabled)

- Batch norm statistics updates
-

13. CONCLUSION

This project successfully developed a state-of-the-art multi-label ocular disease classification system using DenseNet-121 transfer learning. The final model achieved:

✓ **Primary Success Criterion MET:**

- **Test AUC: 0.9074** (exceeds 0.90 target)

✓ **Strong Balanced Performance:**

- **Macro F1-Score: 0.7871** (46% improvement over baseline)
- All classes $F1 \geq 0.55$ (even challenging "Other" class)

✓ **Clinical Utility Demonstrated:**

- Excellent performance on high-impact diseases (Cataract, Hypertension, AMD)
- Robust multi-disease detection capability
- Potential for real-world screening deployment

Key Innovation: By implementing a rigorous two-phase training strategy and leveraging dense connections, the model overcomes severe class imbalance (rare classes <2%) to achieve reliable detection across all 8 pathology types. The addition of optimal threshold tuning provides flexibility for different clinical scenarios (screening vs. diagnosis).

Next Steps:

1. Deploy pilot system in partner clinic for prospective validation
2. Implement Grad-CAM explainability for clinician trust
3. Expand dataset with more diverse populations and rare pathologies
4. Pursue regulatory approval pathway for commercial deployment

This work demonstrates that AI-driven ocular disease screening is not merely a proof-of-concept but a deployable solution ready for real-world clinical integration with appropriate human oversight.

APPENDIX A: GLOSSARY

AUC (Area Under Curve): Measures model's ability to discriminate between positive and negative classes (0.5=random, 1.0=perfect)

Binary Crossentropy: Loss function for multi-label classification where each class is independent

Class Imbalance: When some classes have far fewer samples than others (e.g., 203 Hypertension vs. 4,291 Normal)

DenseNet: Convolutional architecture where each layer receives inputs from all previous layers

F1-Score: Harmonic mean of precision and recall (balances false positives and false negatives)

Fundus Image: Color photograph of the back of the eye (retina)

Macro Average: Average of per-class metrics (treats all classes equally regardless of size)

Multi-Label Classification: Each sample can have multiple classes simultaneously (vs. multi-class with one label)

Precision: Of all positive predictions, what fraction were correct ($TP / (TP + FP)$)

Recall (Sensitivity): Of all actual positives, what fraction were detected ($TP / (TP + FN)$)

Transfer Learning: Using a model pre-trained on one task (ImageNet) and adapting it to another (fundus images)

Weighted Average: Average of per-class metrics weighted by class support (larger classes contribute more)

APPENDIX B: DATASET SOURCES

ODIR-5K (Ocular Disease Intelligent Recognition):

- Source: Peking University
- URL:
<https://www.kaggle.com/datasets/andrewmvd/ocular-disease-recognition-odir5k>
- Images: 6,392 fundus photographs
- Patients: 3,358 (both eyes captured)

External Dataset 1:

- Source: Kaggle augmented collection
- Images: 7,986 (augmented)

External Dataset 2:

- Source: Kaggle preprocessed collection
- Images: 7