

# Machine Learning Operations Canvas (v1.1)

Product name:

AI/Real Image Classification

Designed by:

Group 75

Date:

09/01/2026

Iteration:

1

| Problem  | Data  | Model   | Operations   | Monitoring   | Risk   |
|--|---|---|--|--|--|
| <b>Background</b><br>   | <b>Data Collection</b><br><br>Data sources include the Kaggle competition dataset, featuring authentic images from Shutterstock paired with AI-generated counterparts using state-of-the-art models. Methods involve downloading the dataset. Volume: 79,950 training images and 19,986 test images. Labeling process: Binary labels (0 = Real, 1 = AI-generated) provided in train.csv and test.csv. Frequency: One-time static collection. | <b>Modelling</b><br><br>Use ResNet models (e.g., ResNet50) pretrained on ImageNet, fine-tuned for binary classification. Also evaluate DIRE using diffusion model reconstruction errors. Features: Resize to 224x224, normalize, augment. Transfer learning for ResNet; DIRE as non-neural complement.   | <b>Inference</b><br><br>Deployment process: Package models into a service using Docker for API-based inference. Include details on preprocessing pipelines and batch handling.<br><br>Infrastructure: Cloud-based with GPU support for efficient processing.  | <b>Feedback</b><br><br>Since this is an academic project without live users, feedback is collected via: automated validation on hold-out sets and cross-validation folds, error analysis on misclassified samples periodic benchmarking against new generators or updated test splits, and qualitative review of failure modes. How used: Insights drive iterative improvements. In a future production extension, this would expand to logged production misclassifications and user-reported flags. | <b>Fairness</b><br><br>Check biases across categories (humans vs. others); avoid errors from imbalanced data or biased generators. Use equalized odds metrics, augment data, perform audits.  |
| <b>Value Proposition</b><br><br>The solution reliably detects synthetic images, enhancing trust in visual content and aiding in quality assurance. Benefits: Supports content platforms in verification, helps users identify real vs. AI images, enables integration for automated checks. | <b>Metrics and Evaluation</b><br><br>Metrics: Accuracy, Precision, Recall, F1, AUC-ROC. Methods: 5-fold CV on train, test set eval. Effectiveness: Beat 50% baseline, aim >90% accuracy vs. benchmarks.  | <b>Data Verification and Governance</b><br><br>Data management policies: Ensure quality through checks for corrupt files, label accuracy, and class balance (50/50 real/AI). Privacy and compliance: Adhere to dataset licenses (Shutterstock/Kaggle); no personal data involved.<br><br>Mechanisms: Access controls via secure storage, automated quality scripts for validation, and compliance monitoring through version tracking. | <b>Decision</b><br><br>Predictions integrated: Output probability scores for AI-generated; threshold for binary decisions. Human oversight: In high-stakes scenarios, flag uncertain cases for manual review. Automated systems: Direct integration into workflows like content filters, with logging for audits. | <b>Lifetime</b><br><br>Outline after deployment: Continuous monitoring for drift via periodic re-evaluation on new data. Conditions for retraining: If metrics drop >5% or new generators emerge. For decommissioning: If accuracy falls irrecoverably or replaced by superior models, with graceful shutdown and data archival.  | <b>Explainability</b><br><br>Use Grad-CAM for ResNet visuals; error maps for DIRE. Add SHAP/LIME. Provide explanations like "AI due to textures" for trust.   |
| <b>Objectives</b><br><br>Goals: ≥95% test accuracy; compare ResNet variants and DIRE. Outcomes: Deployable robust classifier. Criteria: Outperform baselines, low FP/FN, via confusion matrices/ROC.  | <b>Model Governance</b><br><br>Process for managing versions: Use Weights & Biases for tracking experiments and artifacts. Conditions from staging to production: Pass validation thresholds. Procedures: Automated pipelines for retraining on new data, with rollback options for failed updates.  |   |  |  | <b>Security</b><br><br>Identify risks: Adversarial attacks, data breaches in storage. Adversarial attacks, system vulnerabilities: Model poisoning or inference-time exploits. Include measures: Adversarial training, input sanitization, secure APIs, and regular vulnerability scans to ensure robustness. |