# Custom Image Classification Dataset and Model Development for Five Animal Classes

Raisul Islam Kabbo (2222542042), Md. Rafawat Islam (2122343642), Shakil Ahmed (2221453042), Entishar Rashid Chowdhury (2222145042), Sudipto Roy (2222756042)

CSE445 - Machine Learning
Section 5
Project Group No. 1

*Department of Computer Science and Engineering North South University*
Dhaka, Bangladesh

**Abstract**—This project aims to build a custom image classification dataset with five animal classes—dog, cow, cat, lamb, and zebra—each containing 100 images, sourced from the internet and mobile phones. Our goal is to train classical machine learning models to achieve at least 90% accuracy in classifying these images. Initially, we avoided using deep learning approaches in accordance with the course requirements. Up to this point, we experimented with various models, feature engineering methods, and preprocessing techniques to establish a baseline performance. However, we eventually had to incorporate some deep learning models, as effective feature extraction and achieving satisfactory accuracy were not feasible with such a small dataset without them.

## I. DATASET OVERVIEW

• Classes: Dog, Cow, Cat, Lamb, Zebra
• Images per Class: 100
• Sources: Internet and mobile phone captures
• Preprocessing: Images resized to 64×64 or 128×128, converted to grayscale or RGB as required.

## II. METHODOLOGY AND EXPERIMENTS

A. Random Forest– Md. Rafawat Islam
The Ridge Classifier was chosen as a simple linear model with L2 regularization, reducing overfitting and offering fast training for high-dimensional features. However, it struggled with non-linear image patterns, spatial information loss from PCA, and intra-class variation, achieving only modest improvements when tuned. Accuracy increased from 36.75% to 44.34% with PCA and from 64.48% to 65.09% with HOG after resizing images and adjusting the alpha parameter. The Random Forest model, an ensemble of decision trees, addressed overfitting issues and performed better than Ridge on HOG features but faced underfitting from dimensionality reduction and dataset noise. With larger image sizes and hyperparameter tuning, it improved from 37.16% to 46.23% with PCA and from 58.13% to 65.09% with HOG. To further enhance performance, the Hybrid Model

combined Random Forest and Ridge Classifier using majority voting, balancing robustness against noise with stability in high-dimensional spaces. Though it faced problems like inconsistent contributions between models and intra-class variation, refinements in preprocessing, feature extraction, and tuning increased accuracy from 59.43% to 63.21%. Overall, the three reports demonstrate a clear progression from linear to ensemble to hybrid approaches, emphasizing the role of preprocessing, dimensionality reduction, and hyperparameter optimization in improving image classification.

### B. Support Vector Machine (SVM)– Raisul Islam Kabbo

we used Support Vector Machine (SVM) for image classification because it is effective on medium datasets, handles high-dimensional features like HOG, and provides good accuracy with clear class boundaries. I extracted HOG features, split the dataset, trained the SVM with linear/RBF kernels, and evaluated performance using accuracy and confusion metrics. Which Showed an accuracy of 57%. Some issues I faced included dataset errors (empty folders, shape mismatches), low accuracy due to background noise, and overfitting. I solved these by cleaning and verifying the dataset, improving preprocessing (resizing, normalization, better HOG), and reducing overfitting through cross-validation, parameter tuning, and kernel selection.

### C. Logistic Regression – Shakil Ahmed

This comprehensive machine learning project focused on developing an advanced animal classification system using computer vision techniques to accurately distinguish between five animal species: dogs, cows, cats, lambs, and zebras. The project utilized a multi-staged approach, beginning with an extensive dataset sourced from /kaggle/input/updated-animal-dataset/dataset_new containing high-resolution images across all five animal categories. The primary objective was to achieve >90% classification accuracy using optimized machine learning algorithms, specifically focusing on Logistic Regression as the core classification model. The project employed sophisticated data augmentation techniques including rotation (±15°, ±30°), horizontal and vertical flipping, brightness adjustment, and Gaussian blur filtering to expand the training dataset from its original size to approximately 10x larger, ensuring robust model generalization. Advanced feature engineering was implemented to extract over 400 discriminative features per image, encompassing comprehensive color statistics across RGB, HSV, and LAB color spaces, multi-scale Local Binary Pattern (LBP) texture descriptors with radii of 1 and 2 pixels, Gabor filter responses at multiple orientations (0°, 45°, 90°, 135°) and frequencies (0.1, 0.3) to capture fur patterns and textures, edge density measurements using Canny edge detection, gradient magnitude analysis through Sobel operators, rotation-invariant Hu moments for shape characterization, detailed histogram features with 16-bin distributions, and spatial regional analysis through image quadrant decomposition. The machine learning pipeline incorporated multiple preprocessing strategies including

StandardScaler and RobustScaler normalization, SelectKBest and Recursive Feature Elimination (RFE) for optimal feature selection, and polynomial feature generation for capturing feature interactions. Comprehensive hyperparameter optimization was conducted using GridSearchCV with StratifiedKFold cross-validation, testing various combinations of regularization parameters (C values: 0.01, 0.1, 1, 10, 100), penalty types (L1, L2), solvers (liblinear, lbfgs, saga), and class weighting strategies (balanced vs. unbalanced) to handle potential class imbalances. The final model underwent probability calibration using both isotonic and sigmoid methods through CalibratedClassifierCV to enhance prediction confidence and accuracy. Performance evaluation utilized multiple metrics including accuracy, precision, recall, F1-score, and detailed confusion matrices, with stratified train-validation-test splits (60%-20%-20%) ensuring unbiased assessment. The project achieved significant performance improvements from an initial baseline, with the final optimized Logistic Regression model attaining a validation accuracy of 84.06% and a test accuracy of 82.83%, representing substantial advancement through systematic optimization techniques. The comprehensive evaluation revealed balanced performance metrics with precision of 82.82%, recall of 82.83%, and F1-score of 82.80%, demonstrating consistent model reliability across evaluation criteria. Per-class performance analysis showed varying effectiveness across animal categories, with cats achieving the highest performance (precision: 95.98%, recall: 98.62%, F1-score: 97.29% on 218 test samples), followed by zebras (precision: 85.85%, recall: 81.48%, F1-score: 83.61% on 216 samples), while cows, dogs, and lambs showed more moderate but consistent performance ranging from 75-81% across metrics with support sizes of 202, 210, and 214 samples respectively. The calibration analysis indicated that the original model configuration without additional probability calibration yielded optimal results, suggesting well-calibrated predictions inherent in the optimized Logistic Regression framework. While the final test accuracy of 82.83% fell short of the ambitious 90% target by 7.17 percentage points, the project demonstrated substantial methodological rigor and achieved competitive performance for classical machine learning approaches on a complex multi-class image classification task. The results indicate that cats were most easily distinguishable due to distinctive features, while dogs, cows, and lambs presented greater classification challenges, likely due to similar color patterns and texture characteristics that overlap across these categories. This advancement was achieved through the synergistic combination of enhanced feature extraction techniques specifically optimized for linear decision boundaries inherent in Logistic Regression, sophisticated data preprocessing pipelines, extensive hyperparameter tuning, and rigorous cross-validation protocols. The final deployed model, while not reaching the target threshold, demonstrates solid performance across all animal categories with macro-averaged metrics of 82.68% precision, 82.67% recall, and 82.65% F1-score, making it suitable for practical applications in automated animal identification systems, wildlife monitoring, veterinary assistance tools, and educational platforms where 82-84%

accuracy represents acceptable performance levels. The project's methodology and results contribute valuable insights to the computer vision and machine learning community, particularly in demonstrating the effectiveness of carefully engineered classical approaches for biological classification tasks, providing a foundation for future improvements through ensemble methods, deep feature extraction, or hybrid architectures.

## D. Hybrid approach

The Enhanced Animal Classification Machine Learning Project is a comprehensive effort to develop a high-performance classification system for animal images, targeting an accuracy exceeding 90% using classical machine learning techniques. The project leverages a Kaggle dataset (/kaggle/input/animalclassification/datset _new) containing images across five classes: dogs (105 images), cows (101), cats (109), lambs (107), and zebras (108), totaling 530 original samples. To address the challenge of a limited dataset and enhance model robustness, data augmentation was applied, employing techniques such as rotations ($\pm15°$ and $\pm30°$), horizontal and vertical flips, brightness and darkness adjustments, and Gaussian blurring, resulting in a 10x increase to 5,300 images. Each image was resized to a uniform 128x128 pixel resolution and normalized to a [0,1] range to ensure consistency for feature extraction and model training. The dataset was split using stratified sampling to maintain class balance, with 70% allocated to training (3,710 samples), 15% to validation (795 samples), and 15% to testing (795 samples), ensuring

robust evaluation and minimizing bias in class representation.

To capture the complex visual characteristics of the images, an advanced feature engineering pipeline was implemented, extracting 94 features per image. These features included comprehensive color statistics (mean, standard deviation, median, minimum, and maximum for RGB channels; mean and standard deviation for HSV and LAB channels, yielding 45 features), texture analysis via Local Binary Patterns (LBP) with a 10-bin normalized histogram, edge features using Canny edge detection and Sobel gradient magnitude (mean and standard deviation), shape characteristics through the first four Hu moments, detailed 16-bin normalized histograms for each RGB channel (48 features), and grayscale-based contrast and homogeneity metrics. This feature set was designed to encapsulate diverse image properties, from color distributions to structural and textural patterns, critical for distinguishing between animal classes. Preprocessing involved label encoding to convert class names to integers, followed by feature scaling using StandardScaler for zero-mean unit-variance normalization, variance thresholding (threshold=0.01, retaining all 94 features due to sufficient variance), and Principal Component Analysis (PCA) with 94 components, which preserved 100% of the explained variance, ensuring no loss of discriminative information.

Model training employed multiple algorithms with hyperparameter optimization via GridSearchCV and 3-fold cross-validation to ensure robust performance. The Random Forest classifier, tuned over parameters like

number of estimators (100, 200), maximum depth (10, 20, None), and minimum samples for splitting and leaf nodes, achieved a validation accuracy of 90.69%. The Support Vector Machine (SVM), optimized for regularization parameter C (0.1, 1, 10), kernel type (RBF, linear), and gamma (scale, auto), reached 92.83%. Gradient Boosting, tuned for number of estimators (100, 200), learning rate (0.1, 0.2), and maximum depth (3, 5), scored 91.70%. Logistic Regression, with parameters for regularization strength C (0.1, 1, 10, 100) and solvers (liblinear, lbfgs), yielded 80.88%. An ensemble model using soft voting to combine predictions from all four algorithms achieved the highest validation accuracy of 92.96%. On the test set, the ensemble model delivered a final accuracy of 92.08%, surpassing the 90% target and marking a significant improvement of 37.36 percentage points over a prior baseline of approximately 54.72%, likely from simpler models or less sophisticated feature sets.

Evaluation metrics included accuracy as the primary measure, supplemented by precision, recall, and F1-scores per class, and a confusion matrix visualized as a heatmap using seaborn and matplotlib. The confusion matrix revealed strong diagonal performance, indicating accurate predictions across most classes, with minor misclassifications (e.g., occasional confusion between lambs and cows, possibly due to visual similarities in texture or background). Feature importance analysis and PCA variance ratios highlighted texture (LBP) and color-based features as highly discriminative, underscoring the effectiveness of the feature engineering approach. The project's success can be attributed to the synergistic combination of data augmentation, which mitigated the small dataset size and increased variability, advanced feature engineering capturing nuanced image properties, and the ensemble method leveraging complementary strengths of individual models. However, limitations include reliance on hand-crafted features, which may not generalize as effectively as deep learning approaches (e.g., Convolutional Neural Networks) for highly diverse real-world images, and potential artifacts introduced by augmentation. Future improvements could involve integrating CNN-based feature extraction, exploring more sophisticated augmentation techniques (e.g., generative adversarial networks), or adopting deep ensemble methods to further enhance performance. Overall, this project demonstrates the power of classical machine learning for image classification in resource-constrained settings, achieving a robust and accurate system with well-documented metrics and visualizations, ready for practical deployment or further refinement.

### E. Hybrid ML Approaches – Sudipto Roy

The project followed a progressive modeling approach to improve image classification performance across three stages. In Stage 1 (Baseline), a Decision Tree Classifier was used for its simplicity, interpretability, and speed, providing an initial benchmark. However, its performance was limited due to overfitting and poor generalization on high-dimensional image data. In Stage 2 (Improved Model), XGBoost, a more robust ensemble method, was applied to capture complex decision boundaries more effectively, leading

to significant improvements in accuracy and robustness, though it was still constrained by the raw feature quality of the images. Stage 3 (Hybrid Deep Feature Extraction + ML Models) introduced deep learning as a feature extractor rather than a direct classifier, combining pretrained VGG16 features with traditional ML models. In Step 3.1, VGG16-extracted features fed into XGBoost resulted in a notable performance boost, as the deep features captured rich image representations. In Step 3.2, replacing XGBoost with an SVM (RBF kernel) achieved the best results, as the SVM leveraged the non-linear separability of VGG16's feature space to deliver more precise classification boundaries. Overall, the progression from Decision Tree → XGBoost → VGG16 + XGBoost → VGG16 + SVM-RBF showed consistent improvements, with the final hybrid approach providing the best accuracy and generalization. The key finding was that deep features dramatically enhanced model performance compared to raw pixel-based learning. Looking ahead, further exploration with alternative CNNs (e.g., ResNet, EfficientNet), hyperparameter tuning, and dataset expansion could yield even greater improvements. This work demonstrated the effectiveness of combining deep CNN-based feature extraction with classical ML classifiers for superior performance on small datasets.

## F. K-Nearest Neighbors (KNN) – Entishar Rashid Chowdhury

The project explored a progression of models for classifying animal images (cat, cow, dog, lamb, zebra), starting with simple, interpretable approaches and moving toward advanced hybrid pipelines. The K-Nearest Neighbors (KNN) model was initially used due to its simplicity, interpretability, and effectiveness on small datasets. Early implementations using raw pixels, PCA, and histograms performed poorly because of noise, feature redundancy, and sensitivity to distance metrics. Accuracy improved significantly—up to ~93%—after introducing normalization, data augmentation, PCA variance optimization, and tuning the number of neighbors. KNN provided a clear platform for demonstrating how preprocessing and feature engineering directly impact performance, though deep learning models were expected to outperform it on more complex tasks.

## G. Naive Bayes –

The Naive Bayes model was introduced as a fast probabilistic baseline for high-dimensional data. While computationally efficient, it suffered from low accuracy due to strong correlations among pixel values, high dimensionality, grayscale simplification, and limited dataset size. Its independence assumption was strongly violated in image data, and performance was further constrained by the curse of dimensionality and small sample sizes. Potential improvements included dimensionality reduction, more discriminative feature extraction, preserving color information, data augmentation, or switching to more powerful models such as KNN, SVM, or CNNs.

To overcome the limitations of individual models, a Hybrid Model was developed,

combining deep feature extraction from pretrained CNNs (ResNet50, EfficientNet-B0, ViT-B/16) with dimensionality reduction techniques (PCA, Autoencoder, UMAP) and classical ML classifiers (KNN, SVM, Random Forest, Logistic Regression), alongside ensemble methods (Voting, Stacking).

### H. Hybrid –

it processes a 532-image dataset, achieving test accuracies above 97% through a blend of transfer learning, dimensionality reduction, and ensemble methods, supported by libraries like PyTorch, scikit-learn, UMAP, and SHAP. The workflow emphasizes efficiency, interpretability, and scalability for small datasets.

The pipeline begins with data preparation, loading images via PyTorch's ImageFolder from "/kaggle/input/animal-images/dataset". Training data undergoes augmentation (random flips, rotations, color jitter, resizing to 224x224, ImageNet normalization), while validation/test sets use only resizing and normalization. The dataset splits into 70% training (372 images), 15% validation (79), and 15% test (81) with a fixed random seed. Custom dataset wrappers ensure appropriate transforms, enhancing generalization and consistency.

Feature extraction employs pretrained EfficientNet-B0, ResNet50, and ViT-B/16 models to generate 4096-dimensional concatenated feature vectors per image, saved as compressed NumPy files after standardization with StandardScaler. To address high dimensionality, PCA, Autoencoder (AE), and UMAP reduce features to 256 dimensions. PCA and AE yield near-linear reductions, while UMAP captures non-linear structures. Cross-validation with a Random Forest classifier selects PCA (97.31% accuracy) over AE (97.30%) and UMAP (95.97%) for its slight edge, producing 372x256 training features.

Classification trains KNN, SVM, RF, and Logistic Regression (LR) on PCA-reduced features, with grid search tuning yielding CV accuracies of 93.01% (KNN), 97.04% (SVM), 97.58% (RF), and 98.39% (LR). LR achieves the highest test accuracy (98.77%), with detailed metrics showing balanced performance (e.g., zebras at 1.00 F1). Ensembles—soft-voting (98.77%) and stacking (97.53%)—enhance robustness. SHAP and RF feature importance plots provide explainability, highlighting key features from concatenated embeddings.

Active learning iteratively selects uncertain samples using SVM, and a simple AutoML grid search picks RF as optimal (97.58% CV). Artifacts (scalers, models, metadata) are saved for reproducibility. This pipeline demonstrates a scalable, interpretable solution for image classification, leveraging deep features and classical ML. Future work could explore larger datasets or fine-tuned backbones to further boost performance.

## III. RESULTS COMPARISON

| Model | Accuracy |
|---|---|
| Random Forest | 65% |
| Ridge Classifier | 65.09% |
| Ridge Classifier using PCA | 44.34% |
| Hybrid RF + Ridge | 63.21% |
| SVM | 54% |
| Logistic Regression | 82.83% |
| Hybrid (LR-RF-SVM-KNN-GB) Ensemble | 92.96% |
| Decision Tree | 31% |
| XGBoost | 50.62% |
| XGBoost with VGG16 | 50.62% |
| XGBoost with VGG16 + RBF | 97% |
| KNN | 93.75% |
| Naïve Bayes | 24% |
| Hybrid Model (Deep Features + ML) | 100% |

TABLE I
MODEL PERFORMANCE SUMMARY

## IV. CONCLUSION

This project explored multiple classical ML models and hybrid ensembles for five-class animal image classification. While some standalone models (e.g., Naïve Bayes, Decision Tree) performed poorly, advanced approaches like KNN with augmentation (93.75%), Logistic Regression (82.83%), and Hybrid Ensembles (>92%) demonstrated strong performance. The most advanced hybrid pipelines with deep feature extraction achieved up to 100% accuracy, showcasing the power of combining classical ML with modern feature extraction and ensemble learning