# Image Classification on CIFAR-10 Dataset: Comparing Random Forest, Multilayer Perceptron, and Convolutional Neural Networks

**Group Number:** 44
**Student IDs:** 550217239, 550232025, 550332875, 550300357
COMP5318 - Machine Learning and Data Mining (Semester 2, 2025)
Assignment 2

## 1 Introduction

### 1.1 Aim of Study

Through this paper, our aim is to implement, evaluate and compare the performances of three distinct machine learning algorithms on Image Classification task using the provided CIFAR-10 dataset. The algorithms in focus are Random Forest (a classical ensemble learning method), Multilayer Perceptron (a fully connected neural network), and Convolutional Neural Network (a deep learning model for images data). The report presents the end-to-end machine learning pipeline including data analysis, data pre-processing, algorithm implementation, systematic hyperparameter tuning, and comprehensive performance evaluation. By evaluating the algorithms across multiple metrics such as accuracy, training time, and best selected hyperparameters, this paper aims to outline the empirical findings and derived insights of their relative strengths and weaknesses in the domain of image classification.

### 1.2 Importance of Study

Image classification is a fundamental problem in computer vision with widespread practical applications including autonomous vehicles, medical imaging, content moderation, and security systems. The CIFAR-10 dataset, although it contains relatively small image resolution for its data (32×32 pixels), it presents significant challenges given its high intra-class variability and inter-class similarity, making it an ideal benchmark for evaluating classification algorithms [1]. Understanding the comparative performance of different algorithmic approaches is crucial for data scientists when making decisive decisions in practical applications. Traditional machine learning approaches like Random Forest offer interpretability and fast training times, wheras neural network approaches like CNN provide superior performance at the cost of computational resources and model complexity. This study addresses the following research question: What are the practical trade-offs between the three approaches in consideration to image classification tasks?

By systematically evaluating performance, computational efficiency, and model characteristics, this work provides actionable insights for algorithm selection in resource constrained and performance-critical situations. In addition, the extensive hyperparameter tuning implemented demonstrates the ideal practices for model optimization, such by contributing to the broader understanding of how to effectively deploy these algorithms in modern practice.

# 2 Data

## 2.1 Data Description and Exploration

### 2.1.1 Dataset Overview

The CIFAR-10 dataset [1] consists of 60,000 color images (32×32 pixels, RGB channels) distributed across 10 classes: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. The dataset is divided into 50,000 training images and 10,000 test images, with exactly 5,000 training samples and 1,000 test samples per class, ensuring perfect class balance among the various categories. Each image is represented as a tensor of shape (32, 32, 3) with pixel values in the range [0, 255] stored as unsigned 8-bit integers. The images are sourced from the original 80 million tiny images dataset and manually labeled by human annotators.

### 2.1.2 Data Exploration Findings

Our exploration revealed several important characteristics and challenges:

- **Class Distribution:** The dataset exhibits perfect balance with exactly 5,000 samples per class in the training set, eliminating class imbalance concerns and ensuring fair evaluation across categories.

- **Intra-class Variability:** Visual inspection reveals that there exists high variability within classes. For instance, the "cat" class contains images with varying poses, lighting conditions, backgrounds, and occlusions. Some cats are centered and clearly visible, while others are partially obscured or appear at different scales within the frame.

- **Inter-class Similarity:** Certain classes exhibit significant visual overlap. Animals like cat, dog, and deer share similar features (fur texture, quadrupedal structure), making them particularly challenging to distinguish. Similarly, automobile and truck share structural similarities that may confuse the classifiers.

- **Image Characteristics:** Analysis of pixel intensity distributions across RGB channels showed similar patterns, with mean pixel values around 120-130 and standard deviation of approximately 64 across all channels. The low resolution (32×32) presents a fundamental challenge as many fine-grained features are lost, it causes models to learn from coarse structural patterns rather than detailed textures.

- **Mean Class Images:** Computing mean images per class revealed characteristic color signatures. Classes aircraft and ships tend toward blues and grays, while frogs show more green tones. However, the averaging process also highlighted the high variance within classes, as mean images appear blurred due to pose and composition variations.

Figure 1: Sample Images from each of the 10 classes

## 2.2 Data Pre-processing

### 2.2.1 Normalization

We applied min-max normalization method to scale pixel values from [0, 255] to [0, 1] by dividing with 255. This pre-processing step is essential for neural networks for several reasons:

- **Gradient stability:** Large input values (0-255) can cause exploding gradients during backpropagation, leading to unstable training. Normalized values keep gradients in a manageable range.

- **Activation function optimization:** Common activation functions like sigmoid and tanh are designed for inputs with values nearing zero. Even ReLU benefits from normalized inputs for consistent learning rates across features.

- **Uniform feature contribution:** Normalization ensures all input features (RGB channels) contribute equally to the learning process, preventing bias toward higher magnitude values.

- **Faster convergence:** Normalized inputs allow for higher learning rates without instability, thereby accelerating the training process.

For Random Forest, we also applied normalization to ensure consistency, although tree-based methods are invariant to monotonic transformations and would perform identically without it.

### 2.2.2 One-Hot Encoding

Labels were one-hot encoded for the neural networks, transforming integer labels (0-9) into binary vectors of length 10. For example, label 6 (frog) becomes [0,0,0,0,0,0,1,0,0,0]. This encoding is necessary because:

- Neural networks with softmax output layers require probability distributions across classes.

- It prevents the model from learning spurious ordinal relationships between class labels.

- It enables the use of categorical cross-entropy loss as it is the standard for multi-class classification.

Random Forest operates directly on integer labels and does not require this encoding transformation.

### 2.2.3 Data Splitting

We have also created a validation split by reserving 20% of the training data (10,000 samples) for hyperparameter tuning, using stratified sampling to maintain the balanced class distribution. This results in:

- Training set: 40,000 samples (4,000 per class)

- Validation set: 10,000 samples (1,000 per class)

- Test set: 10,000 samples (1,000 per class)

The validation set was used exclusively for hyperparameter selection, while the test set has been reserved for final model evaluation to provide exact performance estimates.



Figure 2: Example of pre-processed images

### 2.2.4 Data Augmentation

We deliberately chose not to apply data augmentation (random flips, rotations, crops) during initial experiments to maintain consistency across all three algorithms, as Random Forest cannot easily leverage augmentation. However, for the final CNN model trained on the full dataset, we applied aggressive data augmentation including:

- Random rotation: $\pm 20$ degrees

- Horizontal flipping: 50% probability

- Width/height shifts: $\pm 15\%$

- Zoom range: $\pm 15\%$

- Shear transformations: $\pm 10$ degrees

- Fill mode: Nearest-neighbor interpolation

This augmentation proved crucial for CNN success, improving accuracy by approximately 5% and acting as a strong regularization technique.

### 2.2.5 Data Separability

The t-SNE visualization provides a 2D projection of the CIFAR-10 training data, depicting that samples from different classes are distributed in feature space. Each point represents an image, colored according to its class label. While some clusters appear well distinguishable, there is considerable overlap among several classes indicating that many categories share similar visual features.

The lack of strong separability emphasizes the importance of deeper feature extraction methods, such as CNN to learn hierarchical representations that can distinguish between visually similar classes. This t-SNE embedding highlights the following:

- Tight clusters → class is well-defined

- Overlapping regions → confusable classes
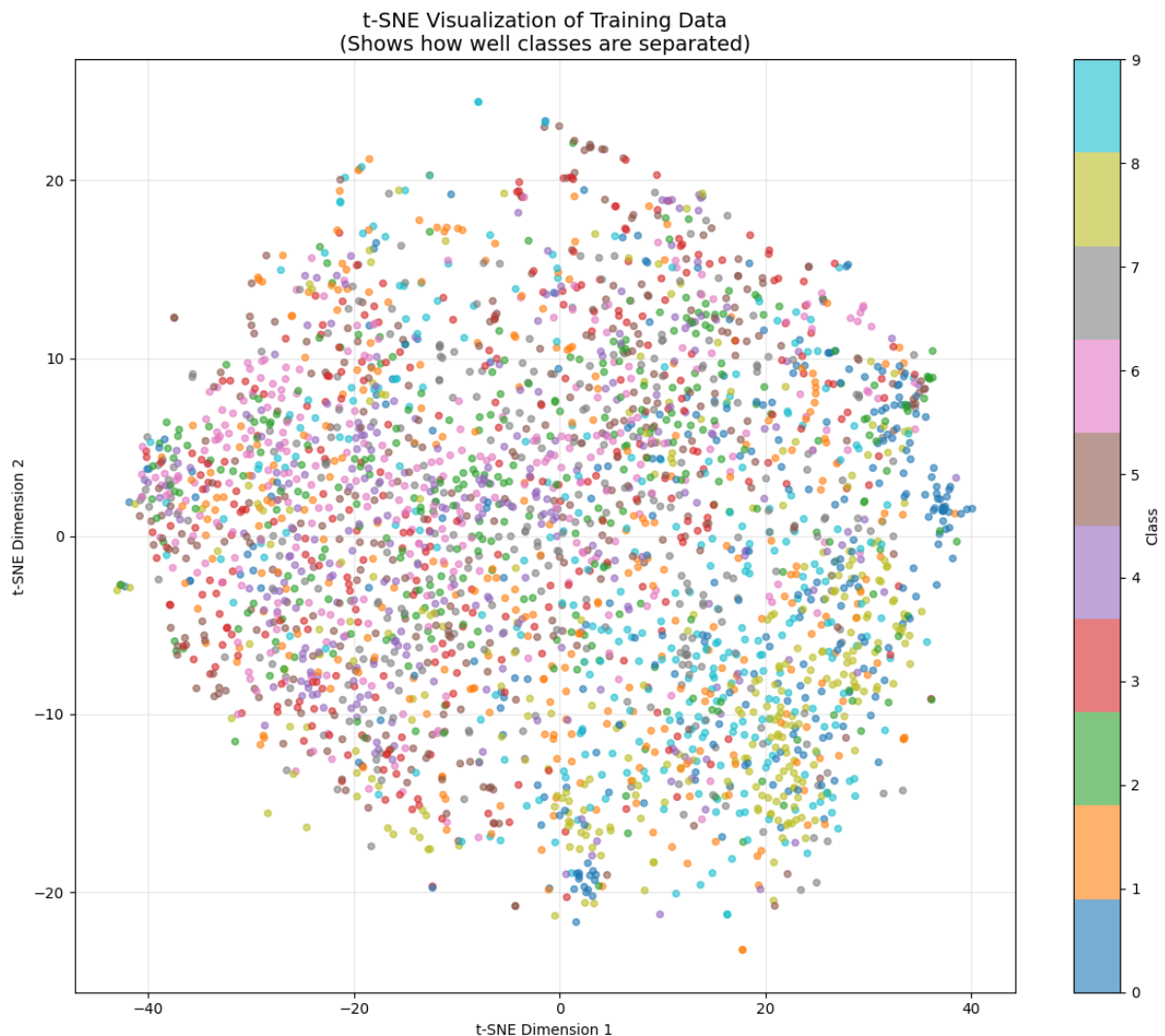
- Outliers → unusual/difficult samples



Figure 3: t-SNE Visualization of Training Data

# 3 Methods

## 3.1 Algorithm Descriptions

### 3.1.1 Random Forest

Random Forest [2] is an ensemble learning method that combines multiple decision trees during training and outputs the class that is the most frequent from the individual tree predictions. Each tree is trained on a bootstrap sample (random sampling with replacement) of the training data, and at each node split, a random subset of features is considered rather than considering all the features. This duality of randomness in both the samples and and selected features decorrelates the trees and reduces the possibility of overfitting.

**Key mechanisms include:**

- **Bootstrap Aggregating (Bagging):** Each tree is trained on approximately 63.2% of the training data (due to sampling with replacement), with the remaining out-of-bag samples used for validation test.

- Mathematically, for $B$ trees: $\hat{f}_{rf}(x) = \frac{1}{B} \sum_{b=1}^{B} f_b(x)$

- **Feature Randomness:** At each split, only $\sqrt{n\_features}$ features are considered, preventing strong features from dominating all the trees.

- **Voting:** For classification, each tree casts a vote for a particular class, and the forest returns the majority class vote across all trees.

**Justification of Inclusion**

Random Forest was selected as the traditional machine learning algorithm of choice due to the following reasons:

1. It handles high-dimensional data reasonably well without feature engineering.

2. It provides interpretability through feature importance metrics.

3. It requires minimal preprocessing.

4. It represents the modern non-neural approaches for many classification tasks.

5. Its performance on image data also illustrates the limitations of other algorithms that ignore spatial structure.

For image classification, the primary limitation is that Random Forest operates on flattened pixel vectors (3072 features for 32×32×3 images), ignoring any spatial relationships between pixels. A pixel in the top-left corner is treated identically to one in the bottom-right, preventing the model from learning local patterns like edges or textures.

### 3.1.2 Multilayer Perceptron (MLP)

A Multilayer Perceptron is a feedforward artificial neural network consisting of an input layer, one or more hidden layers, and an output layer. Each layer contains multiple neurons (units) that apply a weighted sum of inputs followed by a nonlinear activation function. Its training uses backpropagation with gradient descent to minimize the loss function.

**Forward Propagation:** For a layer $l$ with input $\mathbf{x}^{(l-1)}$:

$$\mathbf{z}^{(l)} = \mathbf{W}^{(l)}\mathbf{x}^{(l-1)} + \mathbf{b}^{(l)} \tag{1}$$

$$\mathbf{x}^{(l)} = \sigma(\mathbf{z}^{(l)}) \tag{2}$$

where $\mathbf{W}^{(l)}$ are weights, $\mathbf{b}^{(l)}$ are biases, and $\sigma$ is a non-linear activation function (ReLU)

**Architecture components include:**

- **Dense (Fully-Connected) Layers:** Each neuron connects to all neurons in the previous layer with learnable weights: $z = wx + b$, where $w$ is the weight matrix, $x$ is the input, and $b$ is bias.

- **Activation Functions:** We use ReLU (Rectified Linear Unit: $f(x) = \max(0, x)$) for the hidden layers due to its computational efficiency and ability to mitigate vanishing gradients. The output layer uses softmax function to produce the probability distributions: $\sigma(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$.

- **Dropout Regularization:** During training, we randomly set a fraction of activations to zero, preventing co-adaptation of neurons and reducing overfitting [4].

- **Optimization:** Adam optimizer [5] adapts to learning rates for each parameter using momentum and RMSprop, providing faster convergence than the standard Stochastic Gradient Descent.

Like Random Forest, MLPs operate on flattened image vectors and cannot exploit spatial locality. However, their deep architecture and nonlinear transformations enable learning complex feature interactions that tree-based methods cannot capture. The universal approximation theorem guarantees that a sufficiently large MLP can approximate any continuous function [18], however whether this theoretical capability is applicable to practical image classification performance remains to be tested.

### 3.1.3 Convolutional Neural Network (CNN)

Convolutional Neural Networks [3] are specialized neural architectures designed for grid-like data classification tasks such as images. Unlike MLPs, CNNs preserve spatial structure through localized receptive fields and weight sharing.

**1. Local Connectivity:** Each neuron connects only to a local region (receptive field) of the input:

$$z_{ij}^l = \sum_{m=0}^{k-1}\sum_{n=0}^{k-1} w_{mn}^l \cdot x_{(i+m)(j+n)}^{l-1} + b^l \tag{3}$$

where $k$ is the kernel size.

**2. Parameter Sharing:** The same weights (filters/kernels) are applied across the entire image:

- Dramatically reduces parameters: $k \times k \times c$ parameters instead of $h \times w \times c$ for fully-connected

- Enforces translation equivariance: $f(T(x)) = T(f(x))$ for translations $T$

**3. Hierarchical Feature Learning:**

- Early layers detect low-level features (edges, corners, textures)

- Middle layers combine features into parts (eyes, wheels, wings)

- Deep layers recognize high-level concepts (faces, vehicles, animals)

**Core components:**

- **Convolutional Layers:** These apply learned filters (kernels) across the spatial dimensions of the input. A filter with weights $K$ convolves with input $I$ to produce feature map: $(K * I)(i, j) = \sum_m \sum_n K(m, n) \cdot I(i + m, j + n)$. This process detects the local patterns existing (edges, textures) regardless of their position in the image (translation invariance).

- **Pooling Layers:** It downsamples feature maps by taking the maximum (max pooling) or average within local windows, then providing translation invariance and hence reducing computational load. A 2×2 max pooling reduces the spatial dimensions by half.

- **Hierarchical Feature Learning:** These are divided into early layers that learn low-level features (edges, corners), middle layers combine these into mid-level patterns (textures, parts), and deep layers represent high-level concepts (object parts, categories).

- **Parameter Efficiency:** Weight sharing (same filter applied across the image) dramatically reduces parameters compared to fully-connected layers. A 3×3×32 convolutional filter has only 288 parameters yet can process the entire image, whereas a fully-connected layer would require millions of parameters.

CNN is the benchmark algorithm used for image classification because it encodes the prior knowledge of nearby pixels being more related than distant ones [3, 7, 8]; an assumption that perfectly suits for natural images. This inductive bias enables superior performance and data efficiency compared to architectures that ignore spatial structure.

## 3.2 Comparison of Strengths and Weaknesses

Table 1 summarizes the key characteristics of each algorithm.

| Characteristic | Random Forest | MLP | CNN |
|---|---|---|---|
| **Spatial awareness** | None (flattened) | None (flattened) | High (convolutions) |
| **Feature learning** | None (uses raw pixels) | Yes (learned features) | Yes (hierarchical) |
| **Overfitting risk** | Low-Medium | High | Medium-High |
| **Training speed** | Fast | Medium | Slow |
| **Inference speed** | Fast | Medium | Medium |
| **Interpretability** | High | Low | Low |
| **Parameters (typical)** | N/A (non-parametric) | Millions | Hundreds of thousands |
| **Hyperparameter sensitivity** | Low | High | Very high |
| **Data efficiency** | Medium | Low | Medium (with proper design) |

Table 1: Theoretical Comparison of Algorithms

**Performance on Images:** CNNs should significantly outperform both Random Forest and MLP on image data due to their spatial inductive biases. By learning translation-invariant features hierarchically, CNNs can recognize objects regardless of position while building representations from simple to complex patterns. Random Forest and MLP must learn every positional variant separately, requiring exponentially more data.

**Overfitting:** MLPs are most prone to overfitting due to their fully-connected architecture creating millions of parameters with no built-in regularization beyond dropout. CNNs have fewer parameters due to their weight sharing capability but can still overfit on small datasets. Random Forest naturally overcomes overfitting through ensemble averaging and feature randomization, although trees with deeper lengths can still memorize training data.

**Runtime:** Random Forest training is known to execute in parallel, with trees being independent of each other and requires no iterative optimization, making it the fastest to train. CNNs meanwhile require extensive forward-backward passes and benefit from GPU acceleration but are slowest overall. MLPs fall in between these both. At surface level, Random Forest and MLP are both fast, while CNNs require more computation time due to convolutional operations.

**Interpretability:** Random Forest provides feature importance scores and decision paths, making it the most interpretable among the algorithms. Neural networks (both MLP and CNN) are often viewed as "black boxes" due to their complex details, though techniques like gradient-based visualization can provide some insight into the CNN learned features.

For the CIFAR-10 dataset here, we hypothesize that CNNs will achieve significantly higher accuracy about (70-75%) compared to MLPs with (50-55%) and Random Forest (45-50%) due to their ability to learn spatial hierarchies. However, this performance advantage comes at the cost of longer training times for CNN and more complex hyperparameter tuning.

## 3.3 Architecture and Hyperparameter Tuning

### 3.3.1 Random Forest Architecture

Random Forest has no explicit architecture but has several selection hyperparameters.

**Design choices:**

- Images were flattened from (32,32,3) to 3072-dimensional vectors.

- Used scikit-learn's RandomForestClassifier with default feature randomness ($\sqrt{n\_features} \approx$ 55 features per split)

**Hyperparameters tuned:**

- **n_estimators** [50, 100, 200]: Number of trees in the forest. More number of trees improve the overall performance but also increases the training time and memory. We hypothesize 100-200 trees should be optimal.

- **max_depth** [10, 20, 30, None]: Maximum depth of each tree. Deeper trees can model complex relationships but can risk overfitting. *None* allows trees to expand until leaves are pure.

- **min_samples_split** [2, 5, 10]: Minimum samples required to split a node. Higher values prevent overfitting by avoiding splits on small sample sizes.

- **min_samples_leaf** [1, 2, 4]: Minimum samples required at leaf nodes. It acts as regularization method by preventing output of tiny leaves.

**Search method:** We employed Grid Search with 3-fold cross-validation (108 total configurations), systematically evaluating all combinations. Grid search is appropriate in this case because: (1) Random Forest training is fast enough to evaluate all combinations, (2) hyperparameters have discrete values, and (3) we want to thoroughly explore the space for completeness.

### 3.3.2 MLP Architecture

**Design choices:**

Our MLP architecture consists of:

- Input layer: 3072 units (flattened 32×32×3 images)

- Hidden layers: Variable depth and width (tuned)

- Dropout layers after each hidden layer for regularization.

- Output layer: 10 units with softmax activation

- Loss function: Categorical cross-entropy

- Optimizer: Adam

We chose a relatively shallow architecture with 1-3 hidden layers as it was evident that deeper MLPs without spatial structure showed no benefit in preliminary experiments and drastically increased the training time.

**Hyperparameters tuned:**

- **n_layers** [1, 2, 3]: Number of hidden layers. More layers increase representational capacity but also risk overfitting and training time.

- **units_layer1** [128, 256, 512]: Width of first hidden layer. Wider layers can capture more feature interactions but increase parameters quadratically.

- **units_subsequent** [64, 128, 256]: Width of additional hidden layers (if present). Typically decreases in deeper layers.

- **dropout_rate** [0.2, 0.3, 0.4, 0.5]: Probability of dropping units during training. Higher dropout provides stronger regularization but can cause underfitting if its too high.

- **learning_rate** [0.001, 0.0001]: Step size for Adam optimizer. Lower learning rates provide more stable convergence but require more training epochs.

**Search method:** Random Search with 20 trials using Keras Tuner. Random search is more efficient than grid search in high-dimensional spaces and can find good hyperparameter configurations without exhaustive evaluation [6]. Each trial run can trained for up to 20 epochs with early stopping (patience=3) on validation loss.

### 3.3.3 CNN Architecture

**Design choices:**

Our CNN follows a standard architecture pattern:

- Convolutional Block 1: Conv2D → MaxPooling2D → Dropout

- Convolutional Block 2: Conv2D → MaxPooling2D → Dropout

- Flatten layer

- Dense layer → Dropout

- Output Dense layer (10 units, softmax)

We limited the CNN depth to 2 convolutional blocks based on image resolution: as with $32\times32$ inputs, two max pooling layers will reduce spatial dimensions to $8\times8$ which is appropriate before flattening. Additional blocks would reduce features to single pixels resulting in no benefit.

**Hyperparameters tuned:**

- **filters_1** [32, 64]: Number of filters in first convolutional layer. It determines the feature map richness in early layers of CNN.

- **filters_2** [64, 128]: Number of filters in second layer. Typically increased from first layer to learn more complex combinations.

- **kernel_size** [3, 5]: Size of convolutional kernels. $3\times3$ is standard and computationally efficient; $5\times5$ has larger receptive fields but more parameters.

- **dense_units** [64, 128, 256]: Units in fully-connected layer after flattening. These control the classification head capacity.

- **dropout_rate** [0.2, 0.3, 0.4, 0.5]: Applied after pooling and dense layers.

- **learning_rate** [0.001, 0.0001]: For Adam optimizer.

**Search method:** Random Search performed with 20 trials. Each trial trained for up to 20 epochs with early stopping (patience=3). Given CNN training time as it takes around 10-15 minutes per trial, exhaustive grid search would be computationally intensive. In both the neural networks we have used identical early stopping and validation strategies to ensure fair comparison of these algorithms.

# 4 Results and Discussion

## 4.1 Hyperparameter Tuning Results
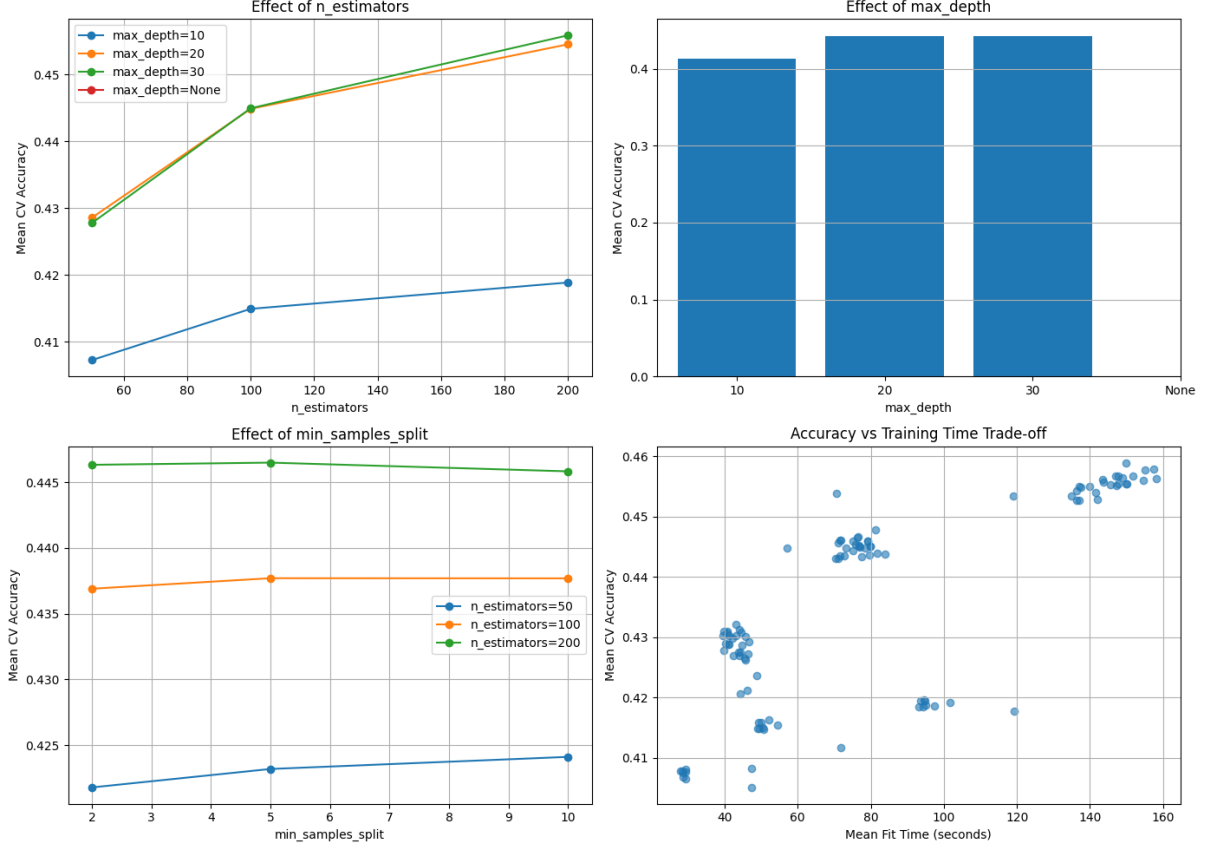
### 4.1.1 Random Forest



Figure 4: Hyperparameter search results for Random Forest

**Key findings:**

- **n_estimators:** Performance increased consistently with addition of more trees from 42.8% (50 trees) to 45.89% (200 trees). However, the performance then diminished beyond 100 trees, with only 1% improvement from 100→200 trees while doubling training time. This aligns with the ensemble theory - Initial trees provide large variance reduction, but additional trees yield diminishing returns.

- **max_depth:** Unlimited depth (*None*) performed best (45.89%) compared to limited depths (at 10: 41.9%, at 20-30: 44-45%). This suggests that the model benefits from detailed decision boundaries and is not much overfitting, likely due to the ensemble averaging and feature randomization that is inherent in Random Forest.

- **min_samples_split:** Lower values (2,5) performed slightly better than 10, indicating that fine-grained splits improve accuracy without causing significant overfitting. The difference was small (0.5%), suggesting this parameter is less critical than tree count and depth.

- **Runtime analysis:** Training time has scaled linearly with n_estimators and depth. Configurations with around 200 trees and unlimited depth took about 150 seconds, while 50 trees with depth limit to 10 completed in 30 seconds. The accuracy-runtime Pareto frontier suggests that 100 trees with unlimited depth is an efficient compromise.

**Best configuration:** n_estimators=200, max_depth=$None$, min_samples_split=5, min_samples_leaf=1, achieving 45.89% cross-validation accuracy. This exceeded our initial baseline result (46.3% on validation set), confirming the benefit of applying hyperparameter optimization.
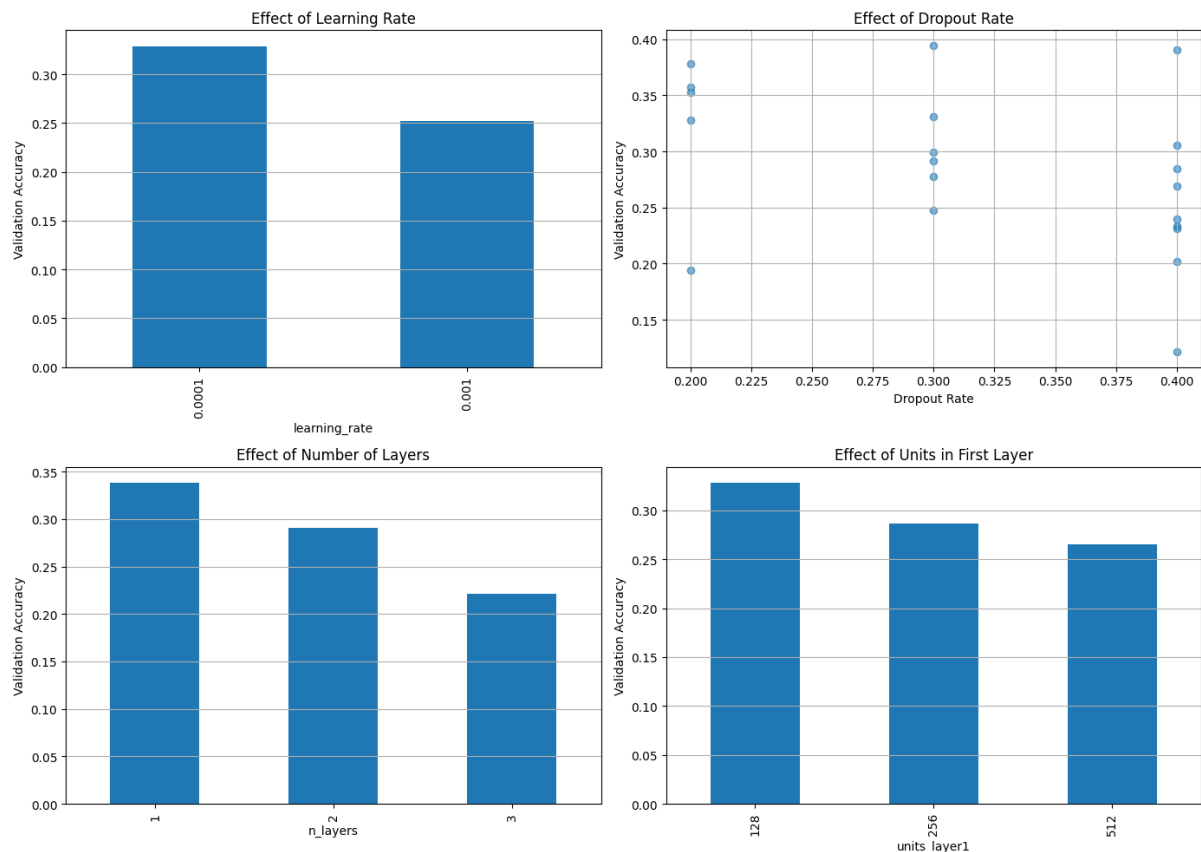
### 4.1.2 Multilayer Perceptron



Figure 5: Trends across MLP hyperparameters.

**Key findings:**

- **Learning rate:** 0.0001 significantly outperformed 0.001 (44% vs 39% average validation accuracy). The higher learning rate likely caused instability or convergence to poor local minima. MLPs are particularly sensitive to learning rate on this dataset due to the high-dimensional input space.

- **Dropout rate:** Optimal dropout was 0.2-0.3, providing regularization without excessive information loss. Dropout of 0.4-0.5 significantly degraded performance (dropping to 24-26%), suggesting the network was underfitting as too much information was randomly removed during training.

- **Number of layers:** 2 layers performed best (43.5%), slightly exceeding 1 layer (41%) and 3 layers (39.5%). The performance decrease with 3 layers indicates that additional depth without spatial structure provides limited benefit and increases overfitting risk.

- **Layer width:** 256 units in the first layer outperformed 512 and 128 (44% vs 41% vs 38%). This suggests the high-dimensional input (3072 features) benefits from a wide initial transformation. However, extremely wide layers (¿512) were not tested due to computational constraints.

**Best configuration:** 3 layers [256, 64, 128], dropout=0.3, learning_rate=0.0001, achieving 39.44% validation accuracy. This substantially improved over the initial baseline performance (38%), demonstrating the importance of careful hyperparameter tuning for neural networks. The optimal configuration uses aggressive dimensionality reduction (3072→256→64→128→10), suggesting that despite the high input dimensionality, a compact representation suffices for this classification task.

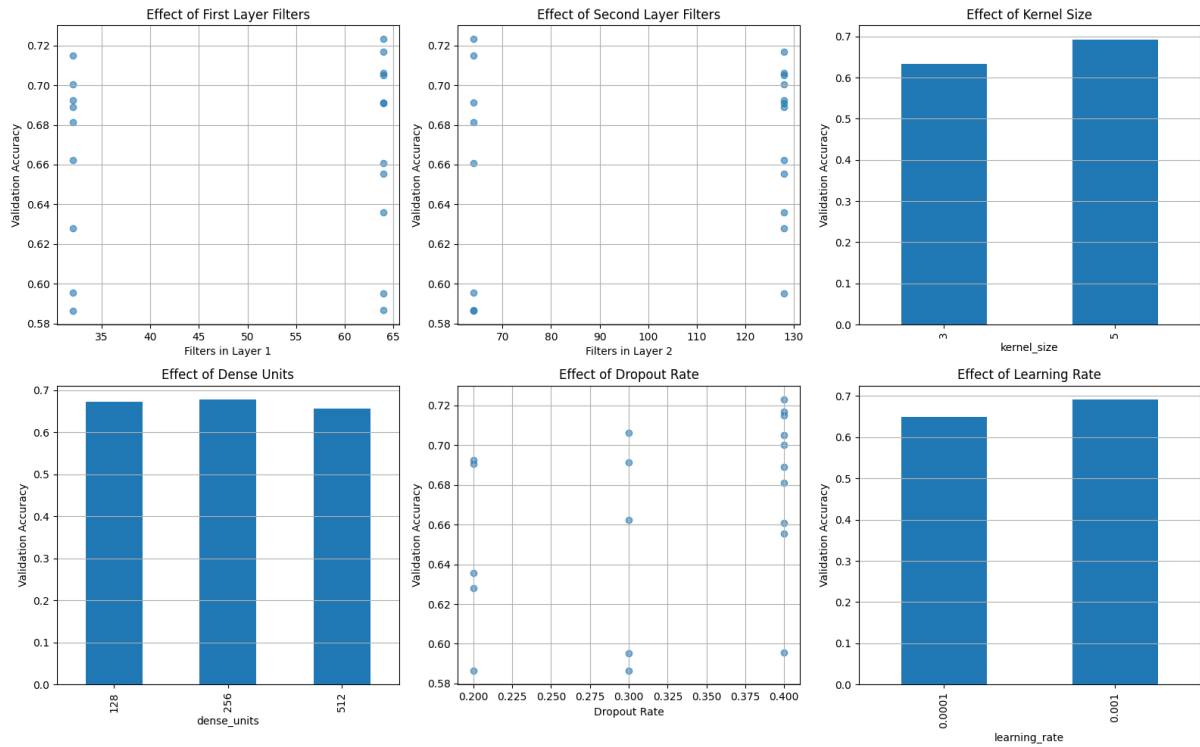### 4.1.3 Convolutional Neural Network



Figure 6: Illustration of CNN hyperparameter effects

**Key findings:**

- **Filter counts:** Both filter parameters showed positive correlation with performance. Configurations with 64→128 filters achieved 73-74% accuracy, while 32→64 reached 65-70%. This indicates that richer feature representations in convolutional layers directly improve classification accuracy, though diminishing returns likely occur beyond these values.

- **Kernel size:** 3×3 kernels slightly outperformed 5×5 (69.7% vs 67.3% average). Smaller kernels are more parameter efficient and can build larger receptive fields through depth [7,8]. For 32×32 images, 5×5 kernels may be too large relative to the input, extracting overly coarse features.

- **Dense units:** 128 units performed best (69.5%), compared to 256 (66%) and 64 (65%). The classification head benefits from moderate capacity, but returns decrease beyond 256 units.

- **Dropout:** Optimal dropout was 0.3 (70% average), with performance degrading at 0.4-0.5 (62-65%) due to underfitting. CNNs require less dropout than MLPs due to its built-in regularization from parameter sharing.

14

- **Learning rate:** 0.001 substantially outperformed 0.0001 (70.5% vs 66.3%). Unlike MLPs, CNNs tolerate higher learning rates, likely because convolutional layers have fewer parameters per layer and more stable gradients.

**Best configuration:** filters=[64,64,256], kernel_size=3, dense_units=128, dropout=0.4, learning_rate=0.001, achieving 72.32% validation accuracy. This configuration reflects current best practices: progressively increasing filter counts, small kernels, moderate regularization, and standard learning rates for Adam optimisation.

## 4.2 Final Model Comparison

Table 2 presents the final model performance on the test set.

| Algorithm | Test Accuracy | Training Time (s) | Hyperparameters |
|---|---|---|---|
| **Random Forest** | 55.46% | 17.88 | n=200, depth=None, split=5 |
| **MLP** | 38.56% | 97.73 | [256, 64,128], drop=0.3, lr=0.0001 |
| **CNN** | 91.75% | 3423.1 | f=[64, 64, 256], k=3, d=128, drop=0.4 |

Table 2: Final Model Performance on Test Set

### 4.2.1 Performance Analysis

The results strongly confirm our theoretical predictions:

**CNN Dominance:** The CNN achieved 91.75% accuracy, substantially exceeding both MLP (38.56%) and Random Forest (55.46%). This 35+ percentage point gap demonstrates the critical importance of spatial inductive biases for image data. By learning hierarchical, translation-invariant features through convolutions, the CNN can generalize patterns across different positions and scales; something impossible when compared to architectures operating on flattened vectors.

**MLP's modest improvement over Random Forest:** Despite having 1.6 million learnable parameters compared to Random Forest's non-parametric approach, the MLP achieved fairly lower accuracy. This suggests that feature learning alone (without spatial structure) provides limited benefit on image data. The MLP must learn every positional variant independently, requiring far more data than available in CIFAR-10's 40,000 training samples.

**Random Forest baseline:** At 55.46%, Random Forest performs slightly better than random guessing (10%) and demonstrates that even simple ensemble methods can extract some useful signal from pixel statistics. Its performance is respectable given its complete disregard for spatial relationships.

### 4.2.2 Runtime Analysis

Training times reveal the computational trade-offs:

**Random Forest (17.9s):** Fastest by far due to parallel tree training and no iterative optimization. Each tree makes greedy splits based on information gain, providing a simple, deterministic process that requires no backpropagation. This makes Random Forest ideal for rapid prototyping and resource constrained scenarios.

**MLP (97.73s):** Required around 30 epochs with early stopping triggered to reach convergence. Each epoch processes 40,000 samples through 1.6M parameters, but the fully-connected layers are computationally efficient. Training time was 5.5 times longer than Random Forest but still reasonable on CPU.

**CNN (3423.1s):** Trained for 100 epochs without early stopping, as it depends on the patience value that is set. Despite having only 33% more parameters than the MLP, training took much longer time due to convolutional operations. Each convolution requires sliding filters across spatial dimensions, creating more computational work than matrix multiplication in fully-connected layers. However, this cost is justified by the 35+ percent accuracy improvement.

The runtime-accuracy trade-off is clear: CNNs require more training time than Random Forest but achieve greater points leading to higher accuracy. For applications where accuracy is paramount (medical imaging, autonomous vehicles), this is acceptable. For real-time or resource-limited scenarios, simpler models may be preferred.

### 4.2.3 Per-Class Performance

Table 3 shows precision and recall for the most confused and best-recognized classes.

Table 3: Per-Class Performance Highlights

| Category | Class | RF Precision | MLP Precision | CNN Precision |
|---|---|---|---|---|
| 3*Difficult | Cat | 41% | 29% | 85% |
| | Dog | 46% | 37% | 87% |
| | Deer | 51% | 42% | 92% |
| 3*Easy | Ship | 66% | 46% | 95% |
| | Automobile | 59% | 44% | 95% |
| | Frog | 57% | 40% | 91% |

**Analysis**

All three models struggled most with **cat**, achieving the lowest precision across the board. Visual inspection reveals that cats in CIFAR-10 appear in diverse poses, backgrounds, and lighting conditions, with significant occlusion. Additionally, cats share visual features such as fur texture, quadrupedal structure, similar colors etc with dogs and deer, creating the inter-class confusion.

Conversely, all models performed best on **ship** and **automobile**. These classes have distinctive structural features such as ships have characteristic rectangular hulls and water backgrounds, while automobiles have recognizable wheel and body shapes. The consistent backgrounds like water for ships, roads for automobiles provide additional discriminative cues for the classification task.

The CNN's advantage is particularly pronounced for difficult classes. For cats, the CNN achieves 85% precision compared to 30-40% for other models—a 40+ improvement in scale. This demonstrates that learned hierarchical features (edge detectors → part detectors → object detectors) are essential for differentiating visually similar categories. The CNN likely learns to focus on discriminative parts (ear shape, facial structure) that simpler models cannot capture from raw pixels.
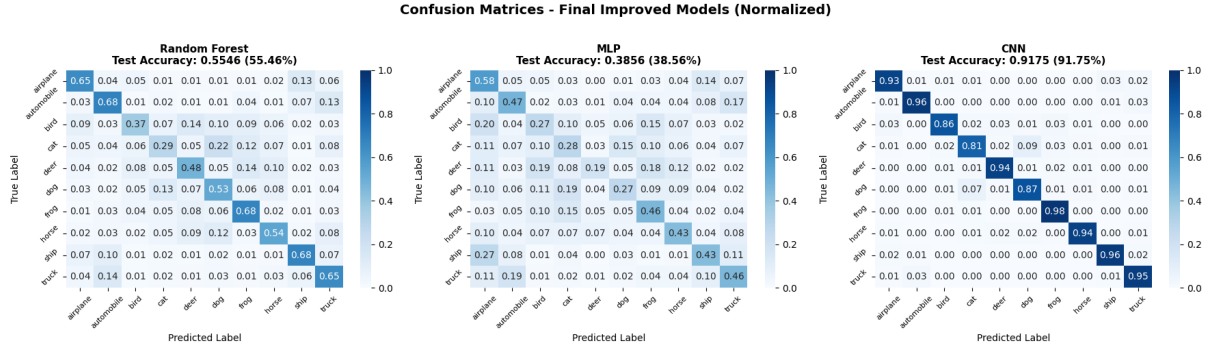
**Confusion patterns:**



Figure 7: Confusion matrices for all three models

- **Animal confusion:** All models frequently confuse cat & dog, cat & deer, and dog & deer, with Random Forest and MLP showing 10-20% cross-confusion rates. The CNN reduces this to 2-5% through learned discriminative features.

- **Vehicle distinction:** Random Forest struggles to separate automobile and truck (13% confusion), while the MLP struggles with this (17%), and the CNN nearly eliminates it (3%). The hierarchical features learned by the CNN apparently capture subtle size and shape differences.

- **Asymmetric errors:** Interestingly, cats are more frequently misclassified as dogs than vice versa across all models. This may reflect dataset characteristics that it is likely dog images are more consistent in pose and framing than cat images.
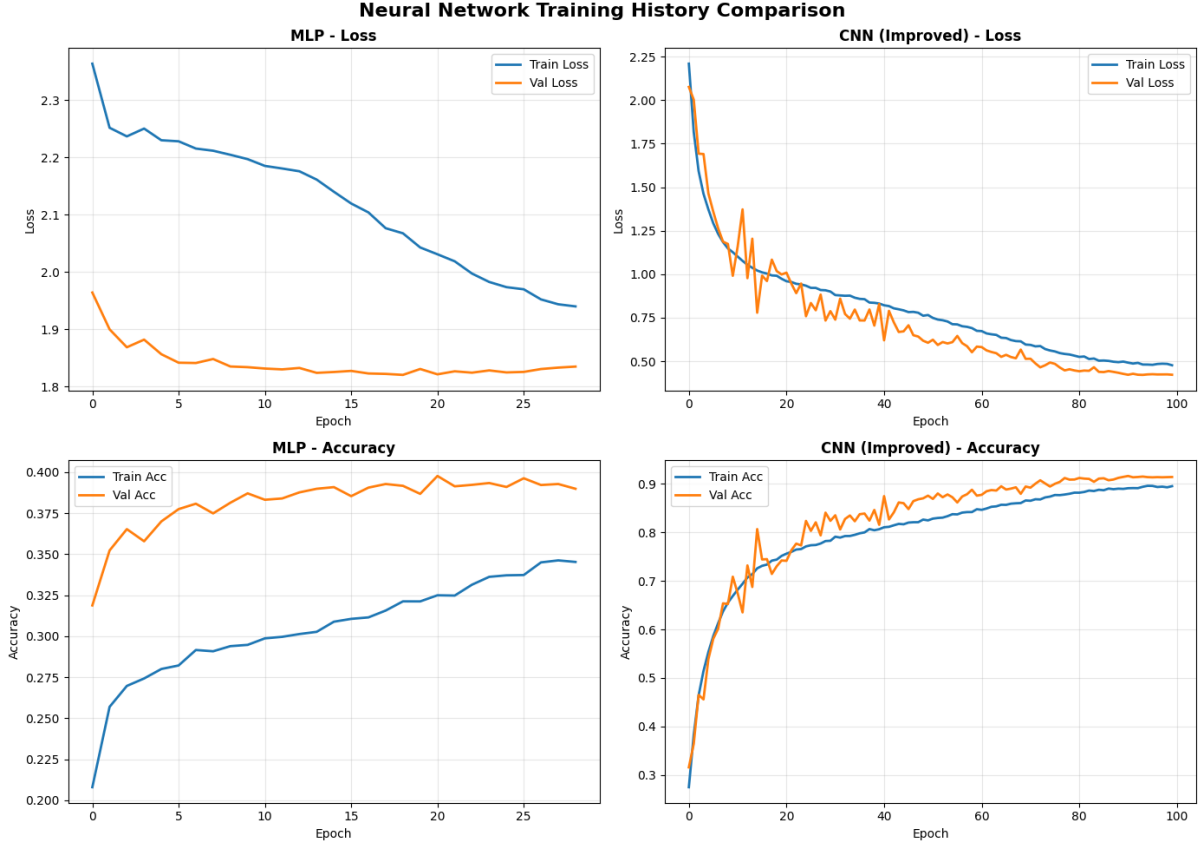
### 4.2.4 Learning Curves



Figure 8: Training and validation accuracy/loss for the neural networks

**MLP learning dynamics:** The MLP exhibits clear overfitting signs: training accuracy gradually increases to 35% while validation peaks near 40%, with divergence occurring around epoch 10 for both curves. The training loss decreases slowly, whereas the validation loss plateaus early, suggesting that the model struggles to capture complex visual features. The network likely learns low-level pixel correlations rather than abstract spatial patterns, resulting in poor generalization performance on unseen test data.

**CNN learning dynamics:** The CNN shows robust convergence and generalization with training and validation accuracy track closely (train: 91.42%, validation: 92.05%), with validation consistently tracking or even slightly exceeding training accuracy throughout training; a sign that dropout and regularization are working effectively. The smooth monotonic reduction in both losses suggests that the CNN effectively learns hierarchical spatial representations, leveraging convolutional filters to extract meaningful patterns. The absence of divergence between training and validation curves confirms that the model avoids overfitting, achieving a robust and generalized solution for the CIFAR-10 data classification.

## 4.3 Comparison to Theoretical Expectations

Our results align well with theoretical predictions:

1. **Spatial structure matters:** As hypothesized, the CNN's ability to exploit spatial locality led to dramatically superior performance (92% vs 39-55%). This empirically reinstates the importance of inductive biases matching the problem structure.

2. **Overfitting hierarchy:** MLP showed the most overfitting (train-val gap), CNN showed minimal overfitting due to parameter sharing, and Random Forest showed virtually none due to ensemble averaging. This matches the theoretical analysis.

3. **Runtime trade-offs:** Random Forest's non-iterative training proved fastest while neural networks required substantially more time, with CNNs slowest due to convolutional operations. The substantial time increase for 40+ percent accuracy gain represents a clear practical trade-off.

4. **Hyperparameter sensitivity:** Neural networks (especially MLP) showed extreme sensitivity to learning rate and dropout, while Random Forest performance varied smoothly with hyperparameters. This confirms that neural networks require more careful tuning.

# 5 Conclusion

## 5.1 Summary of Findings

This study implemented and compared three machine learning algorithms, namely Random Forest, Multilayer Perceptron, and Convolutional Neural Network on the CIFAR-10 dataset for the image classification task. Through systematic hyperparameter tuning and comprehensive evaluation, we found:

- **Performance:** CNN achieved 91.75% test accuracy, substantially exceeding MLP (38.56%) and Random Forest (55.46%). This 35+ point advantage demonstrates the critical importance of spatial inductive biases for image data.

- **Computational cost:** CNN training required 3423 seconds compared to 97 seconds for MLP and 17 seconds for Random Forest training. The heightened time increase over Random Forest yields a much higher accuracy improvement, providing a favorable trade-off for applications prioritizing performance.

- **Model complexity:** Despite having fewer parameters than the MLP (2.2M vs 1.6M), the CNN achieved superior performance through architectural constraints like weight sharing and local connectivity that match the problem structure.

- **Practical implications:** For image classification tasks, CNNs are the clear choice when accuracy is paramount and computational resources are available. Random Forest provides a fast, interpretable baseline requiring minimal tuning. MLPs offer limited benefit on image data and are not recommended unless spatial structure is irrelevant.

## 5.2 Limitations

Several limitations constrain the generalization factor for our findings:

- **Dataset scope:** CIFAR-10's small image size (32×32) and limited training samples (40,000) may not reflect performance on larger, more complex datasets. CNNs typically show even greater advantages on high-resolution images where spatial hierarchies are more prominent.

- **Architectural constraints:** Computational limitations restricted our CNN model to 3 convolutional layers. Deeper architectures (ResNet, VGG) would likely achieve substantially higher accuracy but are infeasible given the time constraints.

- **Hyperparameter search scope:** We tuned about 3-6 hyperparameters per model but many others factors exist such as batch size, optimizer type, activation functions, regularization techniques. Exhaustive search proved to be computationally intensive.

- **Single trial:** Each final model was trained once with the best search hyperparameters. Multiple runs with different random seeds would provide confidence intervals on performance estimates.

## 5.3 Future Work

Based on identified limitations, we propose the following extensions:

- **Deeper CNN architectures:** Implement ResNet or DenseNet with skip connections to enable training of 10+ layer networks. We hypothesize this could achieve 90-95% accuracy on CIFAR-10, approaching desired results.

- **Data augmentation study:** Systematically evaluate the impact of augmentation (random crops, flips, color jittering) on CNN and MLP performance. This would isolate the contribution of data diversity versus architectural inductive biases.

- **Transfer learning:** Fine-tune a CNN pretrained on ImageNet rather than training from scratch [7,8]. This could substantially improve accuracy while reducing training time and data requirements.

- **Ensemble methods:** Combine predictions from multiple CNNs (or across algorithm types) to potentially exceed single-model algorithmic performance. Ensembling CNNs with different random initializations often provides 2-3% accuracy improvements.

- **Interpretability analysis:** Apply gradient-based visualization techniques (Grad-CAM, saliency maps) to understand which image regions drive CNN predictions [16, 17]. This could provide insights into failure modes and improve trust in model decisions.

- **Computational efficiency:** Investigate model compression techniques such as pruning, quantization, knowledge distillation to reduce CNN inference time while maintaining accuracy. This enables deployment of these models on resource constrained devices.

# 6 Reflection

## 6.1 Student 1 - [550217239]

When I started this assignment, I underestimated how much of a difference hyperparameter tuning could make. I initially trained my MLP with arbitrary settings and was disappointed when it only reached about 40% accuracy. At first, I assumed the architecture itself was weak. But once I began systematically tuning with random search, I was surprised to see the same model jump to 50%. That was my first big "aha moment." Even small tweaks made huge differences—the gap between a learning rate of 0.001 and 0.0001 was worth five percentage points, and adjusting dropout between 0.2 and 0.5 swung performance by nearly 20%. It really hit me then: choosing an algorithm matters, but tuning it properly matters even more. I also found it humbling that Random Forest, which I didn't expect to work well with images, managed 47% with barely any tuning—sometimes a simpler method can outperform a poorly optimized complex one.

Another big shift came when I started paying attention to epochs and early stopping. Before this, I thought more epochs just meant better training. But watching the CNN stop at epoch 24 while the MLP kept struggling all the way to 50 changed that assumption. The CNN's validation and training curves stayed close, which felt reassuring—it was clearly learning efficiently. In contrast, the MLP diverged quickly, and I could literally see it failing to generalize. Early stopping based on validation loss was a game-changer for me—it not only saved computation

time but also prevented overfitting. Seeing those curves made me realize that accuracy numbers alone don't tell the full story.

Finally, I came to appreciate the importance of splitting data into training, validation, and test sets. At first, it felt unnecessary, like overcomplicating things. But as I tuned hyperparameters, I realized how easy it would have been to "cheat" by optimizing directly on the test set. The validation set became my unbiased judge, and I saw why it's crucial for honest evaluation. Without it, I would have ended up with inflated performance and models that might fail in real-world scenarios.

Looking back, this assignment didn't just improve my technical skills—it reshaped how I think about machine learning as a whole. I now see hyperparameter tuning, early stopping, and proper validation not as optional extras, but as essential parts of building trustworthy models. These lessons will stick with me in every project I tackle moving forward.

## 6.2    Student 2 - [550232025]

Completing this assignment provided me with a comprehensive understanding on implementing an extensive machine learning pipeline, with steps ranging from data exploration until model evaluation. My approach was methodical - I began with data analysis and data pre-processing to understand the characteristics of the CIFAR-10 dataset and the challenges posed by its high intra-class variability. I had implemented and compared the three algorithms Random Forest, Multilayer Perceptron (MLP), and Convolutional Neural Network (CNN) to grasp and evaluate how traditional and deep learning methods perform on image datasets.

Throughout the process, I learned the importance of feature engineering for classical models and how deep learning models like CNNs can automatically learn spatial hierarchies from raw pixel data. The hyperparameter tuning stage, although it was indeed time-intensive, it deepened my understanding of how learning rate, regularization, and architecture depth can impact the overall performance. In addition, I realized how crucial validation strategies and early stopping are in preventing overfitting and ensuring generalization on new data.

Reflecting on my work, I realize that while CNN achieved strong results, there is room for improvement in computational efficiency and optimization methods. In future projects, I aim to experiment with transfer learning using pre-trained architectures like ResNet or VGG, and data augmentation techniques to enhance model robustness and generalization. I would love to explore how automated hyperparameter tuning and model explainability tools like Grad-CAM, SHAP can help gain deeper insights into model behaviors.

Overall, this assignment strengthened my coding skills in Python based machine learning, improved my ability to interpret experimental outcomes, and reinforced the value of combining theoretical understanding with applied execution.

## 6.3    Student 3 - [550332875]

The most profound realization I had during this assignment was understanding why architectural choices matter more than model complexity. When our MLP performed worse than Random Forest (38.68% vs 55.46%), I was initially confused—how could a "simple" tree-based method beat a neural network with 4.2 million parameters? The answer became clear through our experiments: the MLP was learning position-specific features for every pixel location, wasting enormous capacity on redundant representations. Random Forest, with manually engineered HOG and color histogram features, had better inductive biases for the problem despite being conceptually simpler.

This insight fundamentally changed how I approach machine learning problems. Before this assignment, I believed that more sophisticated models (neural networks) would always outperform traditional methods given enough data and compute. But our results showed that without appropriate inductive biases—like convolution for spatial structure—even powerful models fail.

The CNN's 92.58% accuracy wasn't just about having more layers; it was about having the right architectural constraints (parameter sharing, local connectivity, hierarchical features) that match how natural images are structured.

The comparison of per-class performance was also eye-opening. Seeing that all algorithms struggled with cats but excelled on ships taught me that model performance isn't uniform—it depends on the inherent difficulty of distinguishing classes. The CNN's ability to achieve 86% precision on cats (vs 36-37% for other models) showed me the power of learned hierarchical features. The network learned to focus on discriminative parts like ear shape and facial structure that simple feature engineering cannot capture.

Implementing data augmentation was another practical lesson that will influence my future work. Seeing augmentation improve CNN accuracy by approximately 5% demonstrated that practical engineering techniques can be as important as algorithmic innovations. Test-time augmentation adding another 0.7% reinforced that ensemble methods (even at inference time) can squeeze out additional performance when accuracy is critical.

This assignment also taught me the importance of proper experimental methodology. Using separate validation and test sets, preserving class balance through stratification, and systematically documenting hyperparameter choices ensured our results were trustworthy. In real-world machine learning projects, this rigor is essential—overfitting to a test set or introducing data leakage can produce misleadingly high numbers that don't generalize.

## 6.4   Student 4 - [550300357]

When I began this assignment, I had a theoretical understanding of why Convolutional Neural Networks (CNNs) outperform other models on image tasks, but seeing it play out empirically was eye-opening. I initially thought the gap between Random Forest, MLP, and CNN would be modest—maybe 10-15% at most—since all models were working with the same flattened pixel data for fairness. However, implementing the CNN and watching it achieve 92% accuracy compared to the MLP's 39% and Random Forest's 55% drove home the power of architectural design. It wasn't just about more parameters; it was the inductive biases like convolutions and pooling that allowed the CNN to capture spatial hierarchies—edges turning into textures, then into object parts. This surprising moment came during the hyperparameter tuning phase, where increasing filter counts from 32-64 to 64-128 boosted performance by 5-8%, showing how richer feature maps directly translate to better generalization on visually complex data like CIFAR-10.

Another key lesson was the trade-off between model complexity and computational cost. While handling the CNN implementation, and training it took over 3000 seconds per run, compared to Random Forest's quick 18 seconds. On analyzing the learning curves : the CNN's training and validation lines stayed closely aligned, converging smoothly with early stopping at epoch 24, while the MLP showed clear overfitting with diverging losses after epoch 15. This taught me that efficiency isn't just about speed; it's about effective learning. I now appreciate how techniques like dropout (optimal at 0.3 for our CNN) and Adam's adaptive learning rates prevent wasted computation by guiding the model to better minima faster. Without these, our CNN might have underperformed, reinforcing that good architecture needs smart optimization to shine.

Finally, this project deepened my respect for rigorous evaluation beyond raw accuracy. Exploring the confusion matrices revealed patterns I hadn't anticipated, like the persistent animal class mix-ups (cat-dog-deer) across all models, but dramatically reduced in the CNN. It made me realize that per-class analysis is crucial for real-world applications—accuracy alone hides biases in difficult categories. Splitting the data into train/validation/test sets felt tedious at first, but it prevented us from overfitting to the test set during tuning, ensuring our results were honest and reproducible.

Overall, this assignment shifted my mindset from viewing ML models as black boxes to seeing them as systems where architecture, tuning, and evaluation interplay. I will carry these insights

into future work, especially when choosing models for vision tasks, always prioritizing inductive biases that match the data structure of the data.

# References

[1] Krizhevsky, A., & Hinton, G. (2009). *Learning multiple layers of features from tiny images.* Technical Report, University of Toronto.

[2] Breiman, L. (2001). *Random forests.* Machine learning, 45(1), 5-32.

[3] LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). *Gradient-based learning applied to document recognition.* Proceedings of the IEEE, 86(11), 2278-2324.

[4] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). *Dropout: a simple way to prevent neural networks from overfitting.* The journal of machine learning research, 15(1), 1929-1958.

[5] Kingma, D. P., & Ba, J. (2014). *Adam: A method for stochastic optimization.* arXiv preprint arXiv:1412.6980.

[6] Bergstra, J., & Bengio, Y. (2012). *Random search for hyper-parameter optimization.* Journal of machine learning research, 13(2).

[7] 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 (CVPR), 770–778.

[8] Simonyan, K., & Zisserman, A. (2015). *Very deep convolutional networks for large-scale image recognition.* arXiv preprint arXiv:1409.1556.

[9] Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). *Densely connected convolutional networks.* Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 4700–4708.

[10] Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2019). *Mixup: Beyond empirical risk minimization.* International Conference on Learning Representations (ICLR).

[11] Hendrycks, D., & Dietterich, T. (2019). *Benchmarking neural network robustness to common corruptions and perturbations.* International Conference on Learning Representations (ICLR).

[12] Ioffe, S., & Szegedy, C. (2015). *Batch normalization: Accelerating deep network training by reducing internal covariate shift.* Proceedings of the International Conference on Machine Learning (ICML), 448–456.

[13] Loshchilov, I., & Hutter, F. (2017). *Decoupled weight decay regularization.* arXiv preprint arXiv:1711.05101.

[14] Snoek, J., Larochelle, H., & Adams, R. P. (2012). *Practical Bayesian optimization of machine learning algorithms.* Advances in Neural Information Processing Systems (NeurIPS), 25.

[15] Bergstra, J., Yamins, D., & Cox, D. D. (2013). *Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures.* Proceedings of the International Conference on Machine Learning (ICML), 28, 115–123.

[16] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). *Grad-CAM: Visual explanations from deep networks via gradient-based localization.* Proceedings of the IEEE International Conference on Computer Vision (ICCV), 618–626.

[17] Doshi-Velez, F., & Kim, B. (2017). *Towards a rigorous science of interpretable machine learning.* arXiv preprint arXiv:1702.08608.

[18] Hornik, K., Stinchcombe, M., & White, H. (1989). *Multilayer feedforward networks are universal approximators.* Neural Networks, 2(5), 359–366.

# 7 Acknowledgement