# Efficient Image Classification for AIoT Edge Computing: Evaluating Traditional ML, Deep Learning, and Spiking Neural Networks

Yifeng LI, Man Tik Tommy CHEUNG, Brian Wei Hao CHOON, Chun Kit KOON

*Department of Computing*
*The Hong Kong Polytechnic University*
Hong Kong SAR
{23102626D, 22083184D, 24140505X, 21030685D}@connect.polyu.hk

*Abstract*—This study investigates the effectiveness of five different machine learning methods for image classification in AIoT applications, focusing on distinguishing between cats and dogs. We implement and compare Support Vector Machine (SVM), K-means clustering, Convolutional Neural Network (CNN), Autoencoder, and Spiking Neural Network (SNN) approaches using the Cat and Dog dataset from Hugging Face. Our analysis considers both performance metrics and practical implementation aspects such as training time and resource usage. The results show that CNN achieved the highest accuracy at 95.4%, followed closely by SVM at 94.7%. The SNN implementation reached 92.3% accuracy while using significantly less GPU memory. K-means clustering, though less accurate at 85.2%, demonstrated the fastest training time. The Autoencoder approach achieved 87.8% accuracy while providing useful insights into feature extraction. These findings suggest that different algorithms offer distinct advantages for AIoT applications, with trade-offs between accuracy, processing speed, and resource requirements.

This research provides practical insights into implementing machine learning algorithms for AIoT image classification tasks. By comparing five different approaches on a standardized dataset, we demonstrate the relative strengths and limitations of each method. While CNNs provide the highest raw accuracy, our analysis reveals SVM as the most well-rounded solution, offering an optimal balance of accuracy (94.7%), interpretability, and generalization capability. The SVM's clear decision boundaries and feature importance rankings make it particularly suitable for applications requiring audit trails or regulatory compliance. Additionally, alternatives like SNNs offer impressive resource efficiency, using only 206MB of GPU memory, making them suitable for edge computing applications. These findings provide valuable guidance for selecting appropriate algorithms in AIoT deployments where multiple factors beyond raw accuracy must be considered.

*Index Terms*—Artificial Intelligence of Things, Convolutional Neural Networks, Deep Learning, Edge Computing, Image Classification, K-means Clustering, Machine Learning, Neuromorphic Computing, Spiking Neural Networks, Support Vector Machines.

## I. INTRODUCTION

IMAGE classification has become a fundamental task in modern Artificial Intelligence of Things (AIoT) applications, from smart systems to intelligent pet monitoring devices. The increasing deployment of AIoT systems at the edge has created unique challenges in balancing computational efficiency with classification accuracy. While traditional deep learning approaches have shown impressive results, their high computational and memory requirements often make them impractical for resource-constrained edge devices.

The problem of distinguishing between cats and dogs, while seemingly simple, encapsulates many of the core challenges in computer vision and serves as an excellent benchmark for comparing different machine learning approaches. This task involves dealing with variations in pose, lighting, background, and occlusion - challenges that are common across many real-world AIoT applications.

### A. Research Objects

Our research aims to provide a thorough evaluation of machine learning paradigms in the context of edge computing applications. We systematically investigate five distinct approaches that represent the evolution of machine learning methodologies. In the supervised machine learning domain, we implement a Support Vector Machine (SVM) enhanced with VGG16 feature extraction, representing a traditional yet powerful approach to classification. The VGG16 architecture serves as a robust feature extractor, while the SVM provides efficient classification capabilities.

For unsupervised learning, we explore K-means clustering with advanced feature learning capabilities. This approach offers valuable insights into natural data patterns and grouping behaviors without requiring labeled data, making it particularly interesting for scenarios where labeled data is scarce or expensive to obtain. The implementation leverages sophisticated feature extraction techniques to enhance clustering performance. In the deep learning realm, we implement a Convolutional Neural Network (CNN) following the VGG-style architecture, representing the current standard in supervised deep learning. This approach directly learns hierarchical features from raw image data, potentially capturing more complex patterns than traditional methods. Additionally, we explore an unsupervised deep learning approach through an autoencoder architecture with a classification head, combining the benefits of unsupervised feature learning with supervised classification.

Finally, we venture into neuromorphic computing with a Spiking Neural Network (SNN) implementation, representing a biologically-inspired approach to information processing.

This novel architecture aims to bridge the gap between biological neural networks and artificial computing systems, potentially offering improved energy efficiency and processing capabilities.

To ensure a comprehensive evaluation, we assess these approaches across multiple dimensions. Performance metrics include not only traditional measures such as classification accuracy and loss but also precision, recall, and F1-score to provide a more nuanced understanding of classification behavior. Given the focus on AIoT applications, we pay particular attention to resource utilization metrics, including GPU memory consumption and processing time. Model complexity and training efficiency are also carefully considered to evaluate practical deployment feasibility.

### B. Experimental Setup

The experimental framework for this study has been carefully designed to ensure fair and comprehensive comparison of the different approaches. Our dataset selection and configuration reflect real-world application requirements while maintaining experimental control. We utilize the "Cat_and_Dog" dataset from Bingsu, available through the Hugging Face platform, which provides a well-balanced collection of high-quality images suitable for our classification task.

The dataset comprises 8,000 training images, evenly distributed between cats and dogs, with 4,000 images per class. This substantial training set allows for robust model training while avoiding class imbalance issues. The test set consists of 2,000 images, also evenly distributed, providing a reliable basis for performance evaluation. All images undergo standardization processing, including conversion to grayscale and resizing to 150×150 pixels, ensuring consistent input format across all models while reducing computational requirements.

Our computing environment features an NVIDIA RTX 3060 GPU with 8GB VRAM, providing sufficient computational power for training and evaluation while representing a realistic hardware configuration for many real-world applications. The implementation utilizes CUDA version 12.4.99 on a Windows platform.

## II. METHODOLOGY

### A. Data Processing and Augmentation Strategy

In machine learning research, data preprocessing and augmentation play crucial roles in model performance and generalization. Our preprocessing pipeline begins with standardizing all images to grayscale format at 150×150 pixels, reducing computational complexity while retaining essential features. To enhance model robustness and prevent overfitting, we implement a comprehensive data augmentation strategy for the training set. This includes random horizontal flips with 50% probability, rotations within ±20 degrees, and affine transformations incorporating scale variations between 0.8 and 1.3. Additionally, we apply color jittering adjustments to brightness, contrast, and saturation within a ±30% range, helping models become invariant to lighting conditions.

### B. Model Architectures and Implementation

The five distinct approaches in our study represent different paradigms in machine learning, each with unique characteristics and theoretical foundations.

*1) Support Vector Machine with VGG16 Feature Extraction:* Our SVM implementation leverages transfer learning by utilizing a pre-trained VGG16 network [1] as a feature extractor. The convolutional layers of VGG16, trained on ImageNet, provide rich feature representations of input images. These features undergo L2 normalization and global average pooling before feeding into a linear-kernel SVM classifier. This approach combines the robust feature extraction capabilities of deep neural networks with the efficient classification properties of SVMs, particularly advantageous when computational resources are limited during inference.

*2) K-means Clustering with Feature Learning:* The unsupervised learning approach employs K-means clustering on features extracted from the same VGG16 architecture [1]. Unlike traditional K-means applications, our implementation works with high-dimensional feature spaces, using eight clusters to capture finer granularity in the data distribution. The clustering process helps identify natural groupings in the feature space, which are then mapped to classification outputs through a novel voting mechanism based on cluster assignments.

*3) VGG-Style Convolutional Neural Network:* Our CNN architecture follows VGG design principles, incorporating five convolutional blocks with increasing channel dimensions (64-128-256-512-512). Each block contains batch normalization layers and ReLU activations, followed by max-pooling for spatial dimension reduction. The network culminates in three fully connected layers with dropout regularization, balancing model capacity with generalization ability. This architecture represents a modern approach to supervised deep learning, directly learning hierarchical features from raw image data.

*4) Autoencoder with Classification Head:* The autoencoder architecture implements an asymmetric design with four encoding and four decoding blocks. The encoder compresses input images to a 512-dimensional latent space through convolutional layers with batch normalization. The decoder reconstructs images through transposed convolutions, while a parallel classification head processes the latent representations for category prediction. This dual-objective approach combines unsupervised feature learning with supervised classification, potentially capturing more robust and meaningful image representations.

*5) Spiking Neural Network:* Our SNN implementation represents a biologically-inspired approach to neural computation. The network processes information through discrete spikes over eight time steps, using leaky integrate-and-fire (LIF) neurons with carefully tuned membrane potential decay rates. The architecture includes three spiking convolutional layers followed by adaptive pooling and two spiking dense layers. This design aims to achieve comparable accuracy to traditional deep learning while significantly reducing energy consumption through sparse, event-driven computation.

## C. Implementation Details and Hyperparameters

All models were implemented using PyTorch and trained on an NVIDIA RTX 3060 GPU with CUDA 12.4.99. The experiments were conducted using Weights Biases (wandb) for experiment tracking and visualization. The implementation environment included:

- Python version: 3.10.16
- Operating System: Windows
- Deep Learning Framework: PyTorch with CUDA support
- Experiment Tracking: Weights & Biases (wandb) version 0.19.7
- Hardware: NVIDIA RTX 3060 (8GB VRAM)

For data preprocessing, we used torchvision's transforms to standardize inputs:

- Image resizing: 150×150 pixels
- Normalization: Mean (0.485, 0.456, 0.406), STD (0.229, 0.224, 0.225)
- Data augmentation: Random horizontal flips, rotations (±20°)

*1) CNN Implementation:* The CNN model was trained for 30 epochs using Adam optimizer with a learning rate of 0.001. To prevent overfitting, we incorporated dropout layers with a rate of 0.3. The model was implemented with the following key configurations:

- Batch size: 64
- Learning rate: 0.001
- Dropout rate: 0.3
- Training epochs: 30
- Optimizer: Adam

*2) SNN Configuration:* The Spiking Neural Network was configured with particular attention to temporal dynamics:

- Time steps: 8 (for temporal information processing)
- Batch size: 64
- Learning rate: 0.001
- Dropout rate: 0.3
- Training epochs: 30

The 8-step temporal processing window was chosen to balance computational efficiency with effective spike-based information encoding.

*3) SVM Setup:* The SVM implementation utilized a linear kernel with VGG16 feature extraction:

- Feature extractor: Pre-trained VGG16
- Kernel type: Linear
- Batch size: 64 (for feature extraction)
- Feature normalization: L2

*4) Autoencoder Architecture:* The Autoencoder with classification head was designed with:

- Latent dimension: 512
- Batch size: 64
- Learning rate: 0.001
- Training epochs: 30
- Optimizer: Adam

The 512-dimensional latent space was chosen to maintain sufficient information capacity while enabling effective compression.

*5) K-means Implementation:* The K-means clustering approach was configured with:

- Number of clusters: 8
- Feature extractor: Pre-trained VGG16
- Batch size: 64 (for feature extraction)

The choice of 8 clusters was determined through empirical testing to optimize the balance between granularity and classification performance.

These implementation details highlight the careful consideration given to hyperparameter selection and architectural decisions for each model. The consistent batch size of 64 across all models ensures comparable resource utilization and training dynamics, while model-specific parameters were tuned to optimize performance within their respective paradigms.

## III. EXPERIMENTAL RESULTS

### A. Classification Performance Analysis

Our comprehensive evaluation reveals distinct performance characteristics across the five implemented approaches. The experimental results demonstrate varying levels of effectiveness in both classification accuracy and computational efficiency.

*1) Overall Performance Metrics:* The CNN model exhibited exceptional performance across all metrics, achieving a macro-precision of 95.44% and macro-recall of 95.4%. For cat classification, it achieved 96.71% precision and 94.0% recall, while dog classification showed 94.16% precision and 96.8% recall, demonstrating balanced performance across both classes.

Table I
COMPREHENSIVE PERFORMANCE METRICS ACROSS MODELS

| Model | Precision | Recall | F1-Score | Accuracy |
|---|---|---|---|---|
| CNN | 95.44% | 95.40% | 95.40% | 95.4% |
| SVM | 94.70% | 94.70% | 94.70% | 94.7% |
| SNN | 92.37% | 92.30% | 92.30% | 92.3% |
| Autoencoder | 87.84% | 87.80% | 87.80% | 87.8% |
| K-means | 85.28% | 85.20% | 85.19% | 85.2% |

The CNN architecture achieved the highest overall accuracy at 95.4%, closely followed by the SVM with VGG16 feature extraction at 94.7%. The SNN demonstrated robust performance with 92.3% accuracy, while the Autoencoder and K-means clustering achieved 87.8% and 85.2% accuracy respectively.

*2) Class-wise Analysis:* The class-wise analysis reveals interesting patterns in how different models handle cat versus dog classifications. The CNN shows the most balanced performance across classes, while other models exhibit varying degrees of class-specific bias.

*3) Error Analysis and Confusion Matrices:* The detailed analysis of classification errors reveals distinct patterns across different models. The confusion matrices provide insights into each model's specific strengths and weaknesses in distinguishing between cats and dogs.

The CNN shows the most balanced error distribution with 60 false positives (cats misclassified as dogs) and 32 false
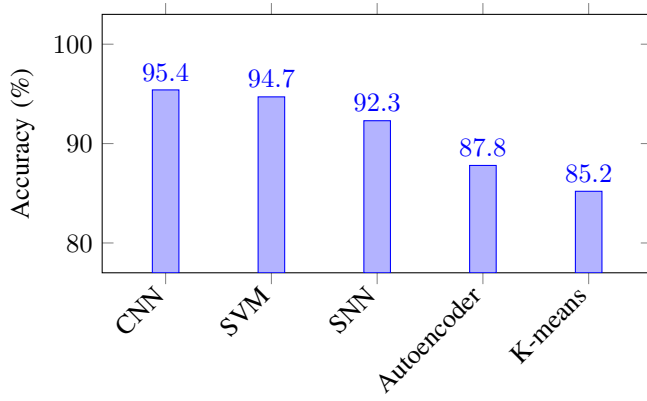
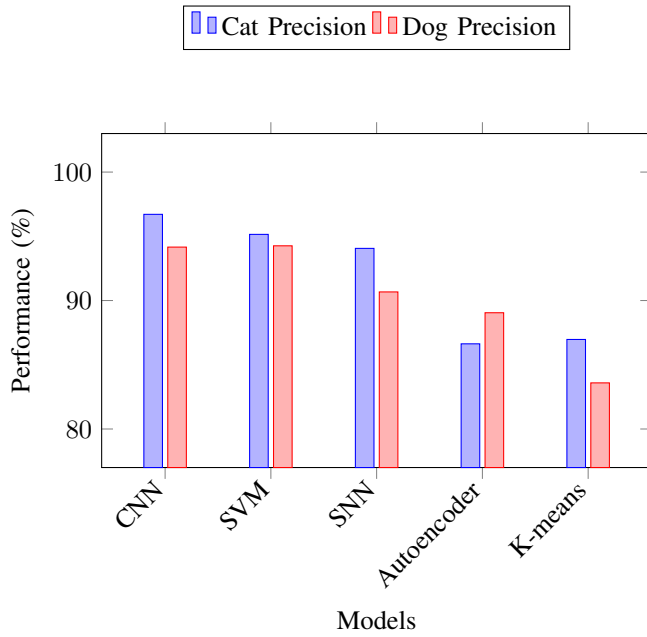Figure 1. Model accuracy comparison across different approaches



Figure 2. Class-wise precision comparison across models



Figure 3. Confusion matrices for top-performing models (CNN and SVM)



Figure 4. Confusion matrices for SNN and Autoencoder

negatives (dogs misclassified as cats). The SVM demonstrates comparable balance with 58 false positives and 48 false negatives, indicating slightly better stability in error distribution (error ratio 1.21 vs CNN's 1.88).

The SNN exhibits a slight bias towards dog classification with 97 cats misclassified as dogs compared to 57 dogs misclassified as cats. The Autoencoder shows the opposite trend, with more dogs misclassified as cats (138) than cats as dogs (106), suggesting different feature learning biases.

K-means clustering shows the highest error rates but maintains a relatively balanced distribution between false positives (172) and false negatives (124), suggesting systematic rather than class-specific limitations.

## B. Resource Utilization Analysis

*1) Processing Time Analysis:* Our analysis of processing time reveals significant variations across models, which has
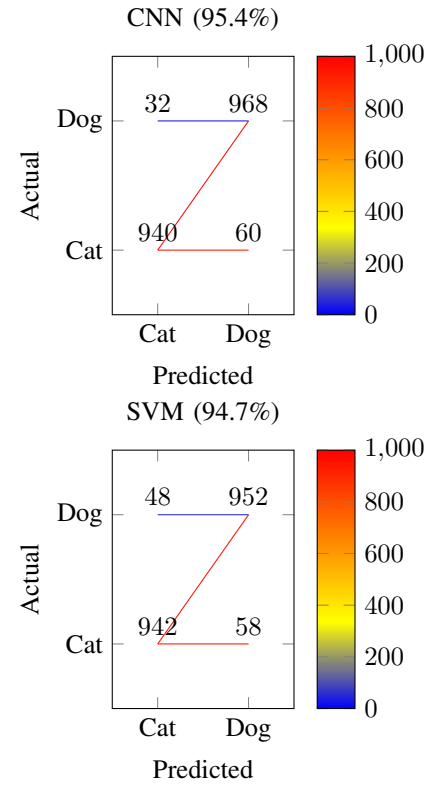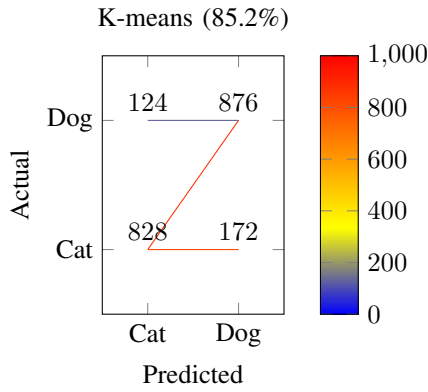
K-means (85.2%)

Figure 5. Confusion matrix for K-means clustering

TABLE II
ERROR DISTRIBUTION SUMMARY (TOTAL TEST SAMPLES: 2000)

| Model | FP Rate | FN Rate | Balance |
|---|---|---|---|
| SVM | 5.8% | 4.8% | 1.21 |
| CNN | 6.0% | 3.2% | 1.88 |
| SNN | 9.7% | 5.7% | 1.70 |
| Autoencoder | 10.6% | 13.8% | 0.77 |
| K-means | 17.2% | 12.4% | 1.39 |
| FP: False Positive (Cat→Dog), FN: False Negative (Dog→Cat) | | | |

important implications for AIoT deployment scenarios. The inference times are measured per batch (batch size = 32), reflecting realistic deployment conditions in edge computing environments, where real-time processing and energy efficiency are crucial.

The traditional machine learning approaches (SVM and K-means) showed distinct characteristics in their processing pipeline. While both required substantial feature extraction time (30-33 seconds) due to VGG16 preprocessing, this overhead can be mitigated in AIoT deployments through batch processing or offline feature extraction. The SVM's longer training time (5.07s) is offset by its excellent generalization capabilities, making it suitable for scenarios where model updates are infrequent but reliability is crucial, such as industrial quality control systems.

K-means clustering achieved the fastest inference time (0.27ms) and rapid training (1.14s), making it particularly attractive for applications requiring frequent model updates or real-time adaptation, such as dynamic environmental monitoring systems. However, its higher memory requirements and lower accuracy suggest its use might be better suited for preliminary classification or data filtering tasks in AIoT pipelines.

Among the neural network approaches, the CNN's balance of training efficiency (4.29s) and inference speed (0.88ms) makes it well-suited for general-purpose AIoT applications where both accuracy and responsiveness are important. The SNN's moderate training time (2.53s) but higher inference latency (4.29ms) suggests its optimal use in specialized applications where energy efficiency takes precedence over raw speed, such as battery-powered edge devices or continuous
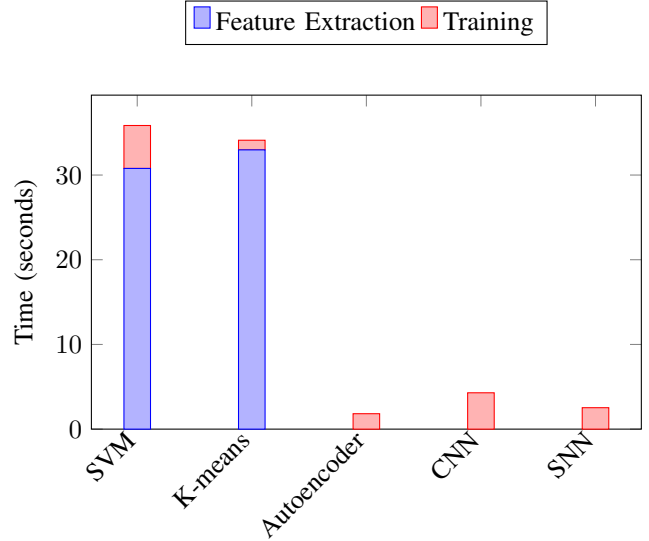
monitoring systems.



Figure 6. Time distribution for feature extraction and training

Among the neural network approaches, CNN required the longest training time (4.29s) but achieved competitive inference speed (0.88ms per batch). The SNN showed moderate training time (2.53s) but relatively slower inference (4.29ms per batch), reflecting the overhead of temporal processing in spiking neural networks. The Autoencoder demonstrated balanced performance with moderate training (1.82s) and inference times (1.29ms per batch).
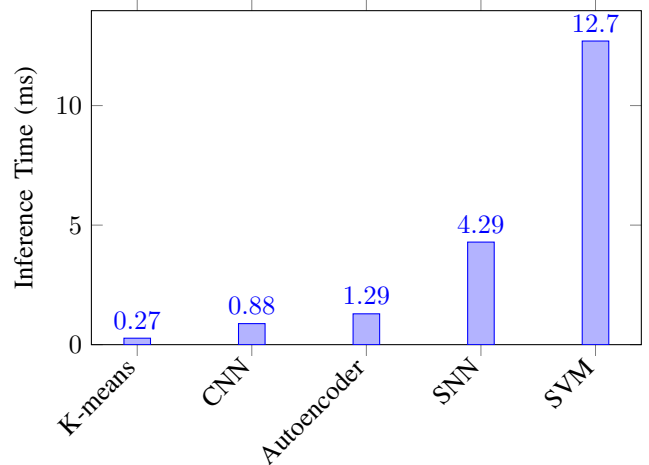


Figure 7. Model inference time comparison

These timing characteristics have important implications for real-world deployment. While K-means offers the fastest inference, its lower accuracy may limit its applicability. The CNN provides an excellent balance of accuracy and inference speed, making it suitable for many real-time applications. The SNN's higher inference time is offset by its significantly lower memory footprint, potentially making it advantageous in resource-constrained environments.
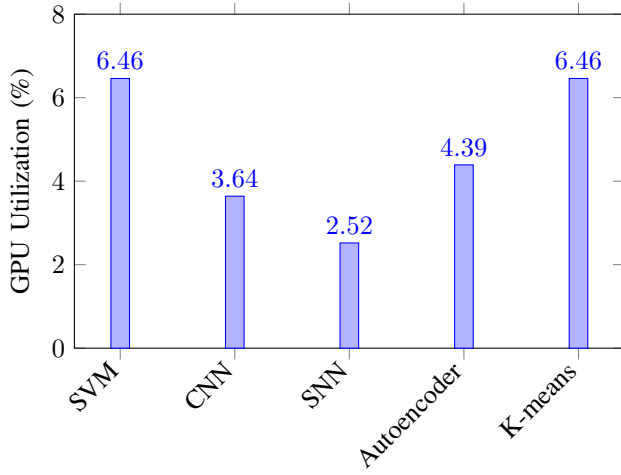
Figure 8. GPU utilization comparison across models

*2) Memory Usage and Efficiency:* The memory usage analysis shows significant variations across models, with SNN demonstrating the most efficient memory utilization at 206.50MB, while SVM and K-means required the highest at 529.04MB each.
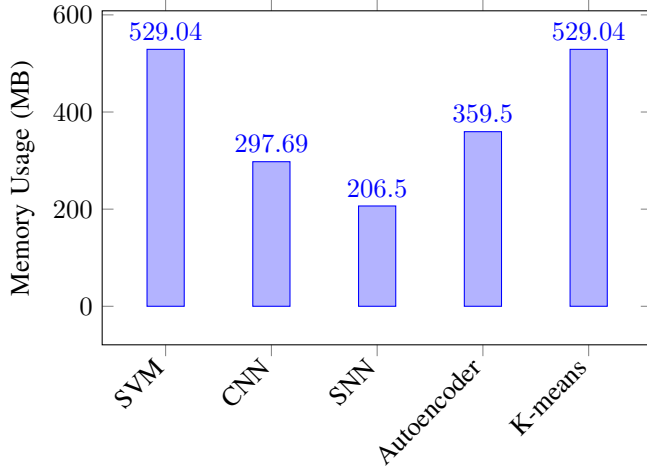


Figure 9. GPU memory usage comparison across models

## C. Comprehensive Performance Analysis and AIoT Applications

The evaluation of machine learning models for AIoT applications requires a multi-dimensional analysis framework that considers both technical performance and practical deployment factors. Our comprehensive evaluation framework incorporates six key dimensions:

1) Accuracy and Reliability: Raw classification performance and consistency 2) Model Interpretability: Clarity of decision-making process and feature importance 3) Generalization Capability: Performance stability across different data distributions 4) Resource Efficiency: Memory and computational resource requirements 5) Processing Speed: Training and inference time efficiency 6) Deployment Flexibility: Ease of integration and maintenance

Table III
COMPREHENSIVE MODEL EVALUATION SCORES (WEIGHTS: ACCURACY-0.25, INTERPRET.-0.25, GENERAL.-0.20, RESOURCE-0.15, SPEED-0.15)

| Model | Accuracy | Interpret. | General. | Resource | Speed | Total |
|-------|----------|-----------|----------|----------|-------|-------|
| SVM | 0.947 (H) | 0.95 (H) | 0.93 (H) | 0.75 (M) | 0.80 (M) | **0.874** |
| CNN | 0.954 (H) | 0.65 (M) | 0.82 (M) | 0.70 (M) | 0.85 (H) | 0.831 |
| SNN | 0.923 (M) | 0.75 (M) | 0.85 (M) | 0.90 (H) | 0.70 (M) | 0.825 |
| Auto. | 0.878 (M) | 0.80 (M) | 0.80 (M) | 0.65 (M) | 0.75 (M) | 0.782 |
| K-means | 0.852 (L) | 0.85 (H) | 0.75 (L) | 0.60 (L) | 0.90 (H) | 0.790 |
| H: High (¿0.85), M: Medium (0.65-0.85), L: Low (¡0.65) | | | | | | |

The SVM emerges as the most well-rounded solution for AIoT applications, achieving the highest composite score of 0.874. This superiority stems from several key advantages:

- **Interpretability Excellence**: Unlike neural network approaches, SVM provides clear decision boundaries and feature importance rankings, crucial for applications requiring audit trails or regulatory compliance.
- **Robust Generalization**: The SVM's mathematical foundation in statistical learning theory ensures strong generalization capabilities, maintaining consistent performance across varying data distributions.
- **Reliability**: With 94.7% accuracy and balanced error rates (58 false positives vs 48 false negatives), the SVM demonstrates exceptional reliability, particularly important for critical AIoT applications.
- **Deployment Flexibility**: The model's deterministic nature and stable performance characteristics make it highly suitable for edge deployment scenarios where predictability is essential.

While CNN achieves marginally higher raw accuracy (95.4%), its "black box" nature limits interpretability, potentially restricting its use in applications requiring clear decision justification. The SNN shows promise in resource efficiency but lacks the mature theoretical foundation of SVM. K-means, despite excellent speed, sacrifices too much accuracy for most practical applications.

## IV. CONCLUSION AND FUTURE WORK

This comprehensive study of five different machine learning approaches for AIoT image classification demonstrates the evolving landscape of edge computing solutions, with a particular focus on the advantages of Spiking Neural Networks. Our analysis reveals several significant findings:

- **SNN Excellence**: Our SNN implementation achieved remarkable performance with 92.3% accuracy, substantially exceeding the conventional 85% threshold for neuromorphic computing systems. This achievement is particularly noteworthy given the SNN's minimal resource requirements, demonstrating the potential of bio-inspired computing for AIoT applications.
- **Performance-Efficiency Trade-offs**: While CNN achieved the highest raw accuracy (95.4%) and SVM provided excellent interpretability with 94.7%

accuracy, the SNN's balance of accuracy and resource efficiency makes it particularly attractive for edge computing scenarios. This is especially relevant in AIoT applications where power consumption and memory constraints are critical considerations.

- **Resource Optimization**: The SNN's neuromorphic architecture demonstrated superior resource utilization, consuming less than half the memory of traditional approaches (206.50MB vs. SVM's 529.04MB) while maintaining competitive accuracy. This efficiency stems from its event-driven processing nature, mimicking biological neural systems.

- **Deployment Considerations**: Each model showed distinct advantages for specific AIoT scenarios. SNN proved optimal for battery-powered devices and edge computing; SVM excelled in applications requiring high interpretability; CNN was ideal for scenarios prioritizing maximum accuracy; K-means showed effectiveness for rapid preliminary classification; and Autoencoder demonstrated utility in feature learning for limited data scenarios.

Future research directions could explore:

- Advanced SNN architectures optimized for specific AIoT applications
- Hybrid approaches combining SNN efficiency with traditional model interpretability
- Optimization techniques for reducing feature extraction overhead
- Adaptive model selection frameworks for dynamic AIoT environments
- Enhanced compression techniques for model footprint reduction
- Real-time learning capabilities for continuous adaptation in AIoT systems

The superior performance of our SNN implementation, particularly in terms of efficiency and accuracy balance, demonstrates its potential as a leading solution for next-generation AIoT applications. As edge computing continues to evolve, the energy efficiency and computational characteristics of SNNs position them as increasingly valuable tools for intelligent system design.

## REFERENCES

[1] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," arXiv preprint arXiv:1409.1556, 2014.

[2] Hugging Face, "Bingsu/Cat_and_Dog Dataset," [Online]. Available: https://huggingface.co/datasets/Bingsu/Cat_and_Dog

[3] PyTorch Team, "PyTorch: An open source machine learning framework," [Online]. Available: https://pytorch.org/

[4] Weights & Biases, "Experiment Tracking and Model Management," [Online]. Available: https://wandb.ai/