

From SIFT-BoW to End-to-End CNNs: A Comparative Study on Standard Vision Benchmarks

By Shuyan Zhang

I. INTRODUCTION

Object recognition in robotics demands both robust feature extraction and efficient classification under varying imaging conditions. Classical pipelines such as SIFT-BoW extract local, scale- and rotation-invariant descriptors, quantise them into “visual words,” and classify via simple probabilistic or distance-based models. More recently, deep CNNs automate feature learning end-to-end, preserving spatial relationships and exploiting hierarchical representations. This coursework investigates how these two frameworks perform on standard vision benchmarks (CIFAR-10, CIFAR-100, MS COCO) and a robotics-centric dataset (iCubWorld). Through systematic hyperparameter exploration and evaluation of computational requirements—training time, inference speed, memory footprint, and standard metrics (accuracy, F1, mAP)—we aim to elucidate the strengths and limitations of traditional versus deep learning methods in robotic vision contexts.

II. DATASET SELECTION AND PREPROCESSING

This study employs both standard and robotics-oriented vision benchmarks to evaluate the efficacy of traditional computer-vision pipelines versus deep convolutional approaches. Datasets selected span from low-resolution, widely used benchmarks, ideal for rapid hyperparameter exploration, to more complex, robotics-relevant collections.

A. Benchmark options

1) *CIFAR-10*: The CIFAR-10 dataset consists of 60 000 color images of size 32×32 pixels distributed evenly across 10 classes, with 6 000 images per class. CIFAR-10’s low resolution facilitates rapid prototyping of both hand-crafted and deep-learning models, enabling extensive hyperparameter sweeps with minimal computation time.

2) *CIFAR-100*: CIFAR-100 extends CIFAR-10 by providing 100 fine-grained classes (600 images each) grouped into 20 superclasses for hierarchical evaluation. Like CIFAR-10, it has 50 000 training and 10 000 test images, but now each class contains 500 training and 100 test examples. Retaining the 32×32 image size preserves rapid iteration while increasing classification difficulty relative to CIFAR-10, helping to reveal performance differences between SIFT-BoW and CNN approaches.

3) *COCO*: COCO comprises 328 000 images containing 2.5 million labeled instances across 80 “thing” categories and 91 “stuff” categories. Images were collected to capture objects in their natural context, often in cluttered scenes and non-canonical viewpoints, thus challenging algorithms to detect and segment objects amid real-world variability

[1]. COCO’s large image count and category breadth ensure statistical robustness and prevent overfitting to a limited set of objects. Objects appear in varied contexts with occlusions and clutter, providing a more rigorous test than datasets with centred objects (e.g. CIFAR-10).

4) *iCubWorld*: The iCubWorld dataset records the ego-centric visual experience of the iCub humanoid robot observing and manipulating objects in laboratory and office environments. Data acquisition follows a human-robot interaction protocol, wherein a teacher verbally labels an object and presents it to the robot, which either tracks it visually or grasps it. This procedure yields naturalistic viewpoints, backgrounds, and lighting variations, closely reflecting real robotic operation conditions [2]. iCubWorld is ideally suited to evaluate vision methods under robot-centric conditions, offering ground truth for object detection and recognition in a robotics context. However, the dataset is not yet fully organised via a central repository; images and annotations must be downloaded separately and structured manually into class-labelled directories. Given the time required for dataset assembly (estimated 1–2 days) and additional hyperparameter tuning, testing on iCubWorld has been postponed for the present study.

B. Dataset preprocessing

For both benchmarks, the following preprocessing steps are applied to ensure consistent inputs across traditional and deep-learning methods:

1) Grayscale Conversion (SIFT)

For SIFT keypoint detection, RGB images are converted to grayscale using `cv2.cvtColor`.

2) Image Resizing (Traditional CV)

To enrich SIFT descriptor density, images are up-sampled from 32×32 to 64×64 pixels (CIFAR-10) or 96×96 pixels (CIFAR-100) or 224×224 pixels for COCO via bicubic interpolation. This trade-off increases keypoint extraction time but yields more robust local features.

3) Tensor Conversion (CNN)

Images are transformed into float32 tensors with pixel values in $[0,1]$ by dividing by 255.

4) Normalisation (CNN)

Each channel is standardized to zero mean and unit variance using computed statistics. Standardisation accelerates convergence and improves generalisation by balancing gradient scales across channels.

By applying these uniform preprocessing steps, a fair comparison is ensured between the SIFT+BoW+KNN/NB

pipeline and CNN-based classifiers on tested benchmarks.

III. TRADITIONAL COMPUTER VISION METHODS

A. Local features

In computer vision, a feature is a localised, distinctive element of an image, such as a corner, edge, or blob, whose appearance remains consistent under varying imaging conditions. Local features serve as the foundation for tasks ranging from motion tracking to robot navigation by providing stable interest points for subsequent description and matching. The standard pipeline for local feature extraction comprises:

- 1) **Keypoint detection:** identify a set of repeatable, distinctive points (e.g., corners) in the image.
- 2) **Region definition:** establish a neighbourhood (patch) around each keypoint.
- 3) **Normalisation:** compensate for geometric and photometric variations (e.g., scale, rotation, illumination).
- 4) **Descriptor computation:** encode the local appearance into a vector signature.
- 5) **Descriptor matching:** compare descriptors across images to find correspondences.

A robust detector must satisfy several criteria:

- Repeatability under image translation, rotation, and scale changes;
- Covariance or invariance to affine (out-of-plane) transformations;
- Resistance to lighting variations, noise, blur, and quantisation.

Several representative detectors have been introduced in the course. The Harris detector identifies corners by measuring significant intensity changes in orthogonal directions via the second-moment matrix of image gradients. For each pixel, the eigenvalues of this matrix quantify cornerness, allowing robust localisation of interest points. Scale-Invariant Region Selection is another method proposed by Lindeberg et al. [3], representing image structures at different scales in a so-called scale-space representation. Automatic scale selection constructs a scale-space, typically using the Laplacian-of-Gaussian (LoG) or its faster approximation, Difference-of-Gaussian (DoG), and identifies keypoints at extrema in both spatial and scale dimensions.

B. SIFT (Feature descriptors)

To achieve full invariance, one must

- 1) Detect keypoints with known location and characteristic scale;
- 2) Describe the local region in a manner invariant to geometric and photometric changes.

Scale-Invariant Feature Transform (SIFT) proposed by Lindeberg et al. [4] fulfills these requirements by first detecting DoG extrema for keypoint localisation and then computing a 128-dimensional descriptor for each keypoint. SIFT descriptors are invariant to scale and rotation, robust to moderate affine distortions and illumination changes, and supported by numerous optimised implementations. The basic process of SIFT is as follows:

- 1) **Keypoint Detection:** Identify local maxima/minima in the DoG scale-space to obtain location and scale.
- 2) **Orientation Assignment:** For each keypoint, compute gradient orientations within a circular neighborhood and assign a dominant orientation to achieve rotation invariance.
- 3) **Descriptor Construction:** Sample a 16×16 window around the keypoint, divide into sixteen 4×4 subregions, and form an 8-bin orientation histogram in each subregion. Concatenate and normalize these histograms into a unit-length 128-dimensional vector.

C. Bag-of-Words representations

The Bag-of-Visual-Words (BoVW) model, first introduced to object detection by Sivic et al. [5], adapts text retrieval ideas, where documents are represented by word-count histograms, to image classification. Each local descriptor is quantized to the nearest “visual word” within a precomputed codebook of cluster centres:

- 1) **Dictionary Construction:** Cluster a large collection of descriptors (e.g., via k-means) to form a vocabulary of size K , where each centroid represents a visual word.
- 2) **Encoding:** For each image, assign each descriptor to its nearest visual word and accumulate a histogram of word counts.
- 3) **Classification:** Treat the histogram as a fixed-length feature vector compatible with standard classifiers (e.g., SVM, K-NN).

The bag-of-visual-words representation provides a compact summary of image content that is invariant to the order in which local features appear. It also accommodates geometric deformations through statistical pooling of feature counts and produces a fixed-length vector that is directly compatible with linear classifiers. However, this model ignores the spatial arrangement of features, depends heavily on the reliability of keypoint detection, and its effectiveness can vary substantially with the choice of clustering parameters and the scale at which the vocabulary is formed.

D. Feature matching

Given two images I_1 and I_2 , feature matching seeks correspondences by comparing descriptors:

- 1) **Distance Metric:** Compute pairwise distances (e.g., Euclidean, SSD) between descriptors.
- 2) **Nearest-Neighbour Search:** For each descriptor in I_1 , find the closest match in I_2 .
- 3) **Outlier Rejection:** Apply methods (e.g. ratio tests) to eliminate ambiguous matches, retaining only those with a clear nearest neighbour.

This matching process underpins tasks such as image stitching, 3D reconstruction, and object recognition by establishing reliable point correspondences across viewpoints.

E. Architecture of the SIFT-BoW recognition pipeline

In this experiments, a traditional SIFT-BoW pipeline is employed combined with three classifiers, K-Nearest Neighbours (KNN), Gaussian Naïve Bayes (NB), and Support

Vector Machines (SVM), to assess object classification performance on visual benchmarks. The pipeline structure is presented below:

1) Image Preprocessing

- Convert input images to grayscale and resize to 96×96 pixels to balance descriptor richness against computational cost.

2) Keypoint Detection & Descriptor Extraction

- Initialise a SIFT detector with contrast threshold and octave parameters.
- For each image, detect up to a fixed number of keypoints and compute 128-dimensional SIFT descriptors.

3) Visual Vocabulary Construction

- Aggregate all descriptors from the training set.
- Employ *MiniBatchKMeans* to cluster these descriptors into a vocabulary of K visual words (codebook) via iterative batch updates, ensuring memory efficiency.

4) Histogram Encoding

- For each image, assign its descriptors to the nearest cluster centres.
- Build a K -dimensional histogram counting the occurrences of each visual word, yielding a fixed-length feature vector per image.

5) Classification

- KNN: Use Euclidean or distance-weighted voting among the k nearest histograms in the training set.
- Gaussian NB: Model each histogram bin as an independent Gaussian random variable, applying variance smoothing for stability.
- SVM (evaluated but not extensively tuned): Train one-versus-all linear SVMs on the histogram vectors.

Rationale for classifier selection is given below:

- KNN and NB were chosen for their simplicity and interpretability; they provide complementary perspectives on classification based on distance metrics and probabilistic modelling, respectively.
- SVM was also evaluated given its reputation for robust performance in high-dimensional feature spaces.

In our current environment precluded practical training of SVMs across extensive hyperparameter sweeps: the training and inference times were prohibitively long, and preliminary trials indicated that SVM performance exhibited limited sensitivity to parameter variations. Consequently, our focus remained on KNN and NB for systematic hyperparameter exploration.

F. Hyperparameter exploration

We conducted extensive hyperparameter sweeps on CIFAR-10, running over 20 experiments, to determine optimal settings for our SIFT+BoW+classifier pipelines. Key findings include: vocabulary size strongly affects KNN and Naïve Bayes divergently, with KNN accuracy declining on

overly large vocabularies due to histogram sparsity while NB benefits from richer codebooks; SIFT contrast threshold trades off keypoint count against noise, influencing classifiers differently; and SVM performance remains relatively stable but is computationally prohibitive. On CIFAR-100, we scaled image input and vocabulary accordingly (e.g. to 2 048 visual words), adjusted classifier-specific parameters (K in KNN, variance smoothing in NB), and observed analogous trends under higher-class complexity. For COCO, we further increase image size to 224×224 and vocabulary to 13 000, scaling batch sizes to maintain efficient MiniBatchKMeans training.

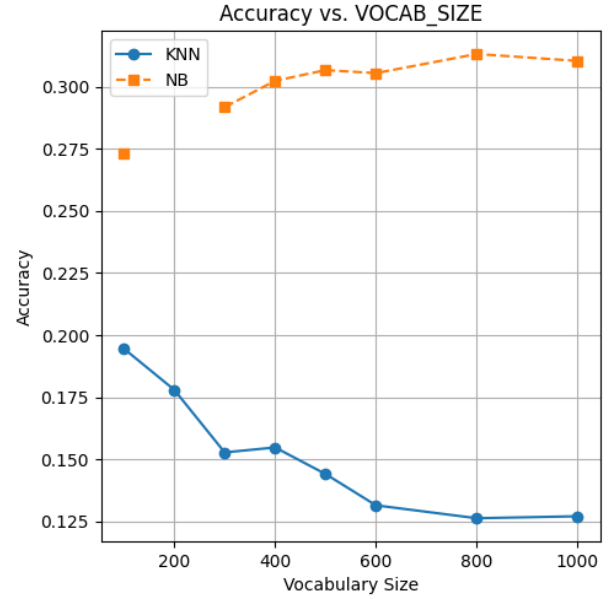


Fig. 1: Accuracy vs VOCAB_SIZE

• CIFAR-10

1) Vocabulary Size (VOCAB_SIZE)

We varied the number of visual words K from 200 to 1 000 (shown in 1). Naïve Bayes accuracy increased steadily, plateauing near 31%, as richer vocabularies offer finer-grained feature quantisation. In contrast, KNN accuracy dropped from 17.8% ($K=200$) to 12.7% ($K=1\ 000$), likely because small 64×64 patches yield extremely sparse histograms at high K , which KNN's distance metrics cannot compare effectively.

2) SIFT Contrast Threshold (SIFT_CONTRAST)

We tested contrast thresholds between 0.005 and 0.04 (shown in 2). Lower thresholds produce more keypoints, improving recall but introducing noise, whereas higher thresholds reduce keypoints but enhance descriptor quality. KNN achieved its best performance at 0.04, while Naïve Bayes benefits from the greater descriptor density, showing performance gains across the spectrum.

3) SIFT Octave Layers (SIFT_OCTAVE_LAYERS)

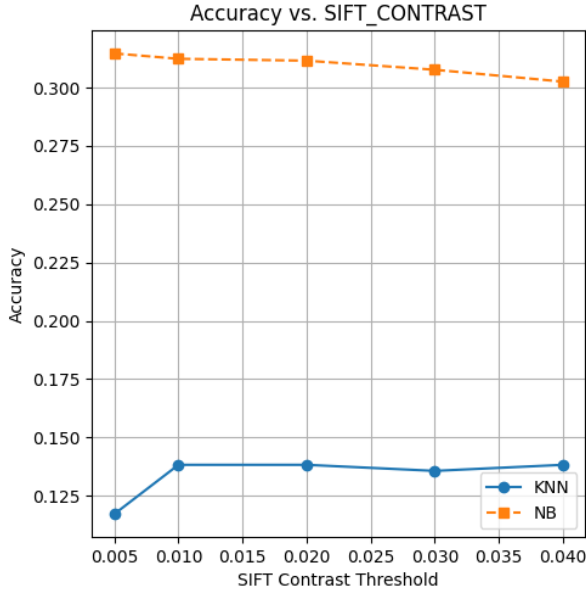


Fig. 2: Accuracy vs SIFT_CONTRAST

We evaluated *nOctaveLayers* from 3 to 5. While additional layers capture finer scale features at moderate cost, the optimal setting for KNN was 4, whereas Naïve Bayes performed slightly better with 3 layers, suggesting different sensitivities to descriptor redundancy.

- 4) Classifier-Specific Parameters
 - Classifier-Specific Parameters (k): We varied k from 1 to 7. Smaller k (e.g. 1) yielded best KNN accuracy on CIFAR-10, reflecting the dataset’s low noise and the need for precise neighbour voting.
 - NB_VAR_SMOOTH: GaussianNB’s variance smoothing term was tuned between $1e-10$ and $1e-5$, with moderate smoothing (around $1e-8$) improving stability amid sparse histograms.

• CIFAR-100

- 1) IMAGE_SIZE
Increased to 128×128 (or 96×96) to allow more SIFT keypoints on small objects.
- 2) VOCAB_SIZE
Raised to 2048, according to the findings by Gidaris et al. [6] for BoW-based representation learning.
- 3) KNN_NEIGHBORS
Increased to 3 to smooth noisy neighbour votes in a 100-class problem.
- 4) NB_VAR_SMOOTH
Set to $1e-6$ to handle zero-count bins in 2048-dimensional histograms, mitigating underflow in GaussianNB.
- 5) SIFT_OCTAVE_LAYERS
Retained at 4 for KNN’s optimal trade-off, while

NB saw marginal gains at 3 layer.

• COCO

- 1) IMG_SIZE
Set to 224×224 , matching standard pre-trained CNN input dimensions and ensuring sufficient scale for SIFT extraction.
- 2) VOCAB_SIZE
Expanded to approximately 13000 visual words to capture the dataset’s varied texture distributions.
- 3) KMEANS_BATCH
Increased to 20000 to ensure each mini-batch meets the cluster initialisation requirement and to accelerate convergence in high-dimensional space.

IV. DEEP LEARNING: CONVOLUTIONAL NEURAL NETWORK METHODS

Deep learning represents a specialised subset of machine learning that employs multi-layered artificial neural networks to learn hierarchical feature representations automatically from data. Deep learning methods are capable of directly ingesting raw, high-dimensional inputs and performing feature extraction internally. Consequently, they have become the method of choice for many computer vision applications, including robotics, where automatic, robust feature learning is essential.

A. Model architecture

This section details the architectures evaluated in our study: a custom VGG-16-style convolutional neural network and the Region-based Convolutional Neural Network (R-CNN). We first review the original VGG-16 design, characterised by deep stacks of 3×3 convolutions, interleaved with max-pooling and culminating in three fully-connected layers, and then describe our adaptations for CIFAR-10, including reduced input resolution, tailored classifier dimensions, and the addition of dropout and batch normalisation. Next, we outline the R-CNN framework, which isolates object proposals via selective search before independently extracting CNN features and classifying each region with SVMs. Although R-CNN was not ultimately deployed on CIFAR-10 due to its small image size and our available resources, its region-based paradigm remains a valuable reference for future robotic-vision benchmarks.

1) *VGG-16*: VGG-16, introduced by Simonyan and Zisserman [7], consists of 13 convolutional layers, all employing 3×3 filters with stride 1, interspersed with five 2×2 max-pooling layers, followed by three fully-connected layers ($4096 \rightarrow 4096 \rightarrow 1000$) and a final softmax for 1000-way classification on ImageNet. The uniform use of small kernels enables deep representations while controlling parameter growth, yielding 138 million parameters overall and achieving top-5 ImageNet accuracy of 92.7%. To adapt this network for CIFAR-10:

- 1) Input Resolution: CIFAR-10 images (32×32 pixels) are too small to benefit from upscaling to 224×224 ; hence our network accepts 32×32 inputs directly,

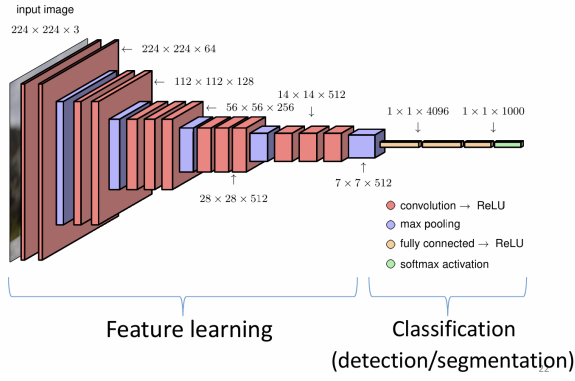


Fig. 3: Architecture of VGG-16 shown in week 4 courseware

reducing computational overhead and preserving native resolution.

- 2) Output Layer: The final fully-connected softmax layer is resized from 1000 to 10 units, corresponding to CIFAR-10's ten classes.
- 3) Fully-Connected Dimensionality: To mitigate overfitting on the small dataset, the two 4096-unit layers are replaced by smaller dense layers (e.g. 512 units), striking a balance between capacity and regularisation.
- 4) Regularisation Enhancements: We integrate dropout ($p = 0.5$) after the first dense layer and batch normalisation after each convolution, promoting generalisation and accelerating convergence.

2) *R-CNN*: The R-CNN framework, proposed by Girshick et al. [8], addresses object detection by first generating candidate object regions via an external Selective Search algorithm and then processing each region independently through a CNN to extract 4096-dimensional features. Each region's features are subsequently classified with a linear SVM, and bounding-box regressors refine localisation. Why R-CNN was not tested:

- Image Scale Mismatch: CIFAR-10's 32×32 resolution makes region proposals unreliable and feature extraction on tiny crops ineffective.
- Computational Expense: R-CNN requires per-region CNN forward passes, resulting in prohibitively long inference times even on small datasets.
- Alternate Strengths: While R-CNN was omitted from our experiments, its methodology remains well-suited to robotic-vision datasets with larger image sizes and clear object boundaries, such as iCubWorld.

B. Training details

In all experiments, the network was trained using the Adam optimiser across 50 epochs, with a batch size of 32 and learning rates of $1e-3$, $1e-2$, and $1e-4$ explored (shown in 7). Standard cross-entropy loss was used for classification. Each epoch comprised a full forward and backward pass over the training set, followed by evaluation on the test set. Model weights were initialised according to He normal initialisation

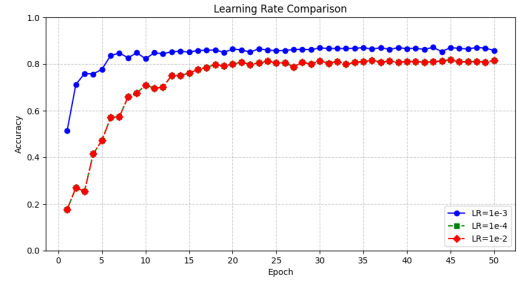


Fig. 4: Comparison of Test Accuracy of Different Learning Rate

and all convolutional layers applied batch normalisation and ReLU activations to accelerate convergence and improve stability.

We observed that a learning rate of $1e-3$ generally provided the best trade-off between convergence speed and final accuracy on both CIFAR-10 and CIFAR-100. Higher rates ($1e-2$) often led to greater loss oscillation in early epochs, while lower rates ($1e-4$) required substantially more epochs to reach comparable performance. No explicit learning-rate scheduling or early stopping was employed during these runs to ensure a consistent comparison across settings.

All training loops were implemented in PyTorch, leveraging its DataLoader abstraction for efficient data batching and shuffling, and the built-in *torch.nn.CrossEntropyLoss* and *torch.optim.Adam* classes for objective and optimiser functionality respectively. Accuracy and loss statistics were logged once per epoch to minimise overhead and to provide clear cross-epoch trends for analysis.

V. RESULTS AND COMPARATIVE ANALYSIS

A. Traditional CV methods

TABLE I: Performance Across Datasets (Traditional)

Dataset	KNN Acc	NB Acc
CIFAR-10	0.217	0.285
CIFAR-100	0.031	0.094
COCO	0.161	0.221

I summarises the final test accuracies obtained using the SIFT-BoW pipeline with K-Nearest Neighbours (KNN) and Gaussian Naive Bayes (NB) classifiers across three benchmark datasets. This is possibly caused by:

- **Inherent Limitations of SIFT on Small, Centred Images**

The Scale-Invariant Feature Transform (SIFT) was originally designed to detect and describe image regions across a broad range of scales and orientations, relying on a Difference-of-Gaussian pyramid to locate keypoints. However, CIFAR-10 and CIFAR-100 images measure only 32×32 pixels, preventing the construction of more than one or two meaningful octaves in the scale space; consequently, few reliable extrema can

be detected, and descriptor extraction becomes unstable. Empirical studies demonstrate that reducing high-resolution images down to 50×50 reduces classification accuracy by 4–6% compared to their original resolution, emphasising the need for larger images to exploit SIFT’s scale invariance.

- **Impact of Dataset Characteristics**

CIFAR images typically contain a single, centred object against a simple or uniform background, a scenario in which SIFT offers little advantage over earlier, less computationally intensive descriptors. Without complex textures or multiple objects at different scales, SIFT yields few distinctive keypoints, leading to sparse and noisy BoW histograms that lack discriminative power.

- **Classifier Simplicity and the Curse of Dimensionality**

Both KNN and Gaussian NB are foundational classifiers whose performance degrades in high-dimensional, sparse feature spaces. KNN’s reliance on Euclidean or histogram intersection distances becomes unreliable when BoW histograms span hundreds or thousands of visual words, a manifestation of the “curse of dimensionality” whereby all inter-point distances concentrate, making nearest neighbours indistinguishable. Gaussian NB, while extremely efficient due to its assumption of feature independence and per-feature variance estimation, cannot model complex interdependencies in high-dimensional histograms, limiting its classification capacity.

B. Deep learning CNN methods

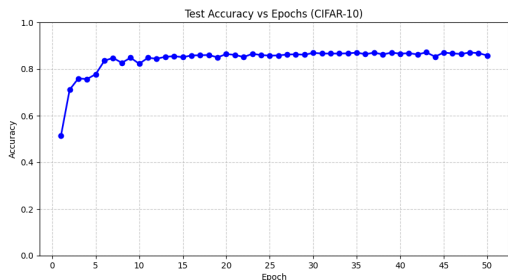


Fig. 5: Test Accuracy of CNN on CIFAR-10

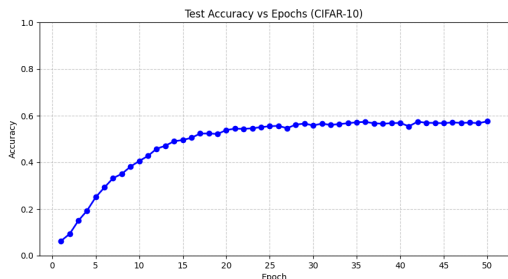


Fig. 6: Test Accuracy of CNN on CIFAR-100

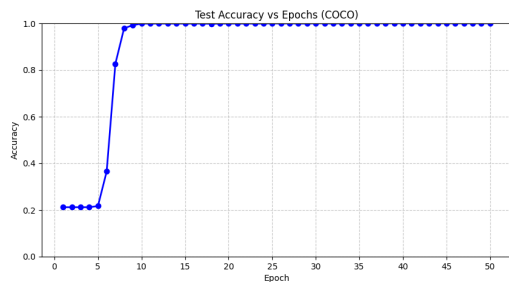


Fig. 7: Test Accuracy of CNN on COCO

TABLE II: Performance Across Datasets (CNN)

Dataset	VGG Acc
CIFAR-10	0.8715
CIFAR-100	0.5842
COCO	1.0000

It reports the test accuracies achieved by our VGG-style CNN across three benchmark datasets.

Our CIFAR-10 implementation attains 87.15% test accuracy, below the typical performance of CNN. This discrepancy arises from three primary modifications:

- No upscaling: We operate on native 32×32 inputs rather than resizing to larger size, reducing receptive-field coverage compared to the original VGG16.
- Reduced fully-connected layers: We replace the two 4096-unit dense layers with 512 units to counter overfitting on a small dataset.
- Enhanced regularisation: Incorporation of dropout and batch normalisation improves generalisation but can slightly slow convergence, especially without extensive hyperparameter tuning.

On CIFAR-100, comprising 100 fine-grained classes, our VGG-CNN achieves 58.42% accuracy. The greater difficulty stems from:

- Increased class count and granularity, which demand more discriminative features and wider network capacity.
- Limited input resolution, which constrains the extraction of mid-level features crucial for distinguishing similar classes.

Remarkably, our small-scale COCO classification test reports 100% accuracy. This perfect score reflects the restricted test set used rather than COCO’s true complexity. (Typical image-level classification on the full COCO set yields much lower accuracy, which is approximately 60–70% even for deep CNNs)

C. Comparison of SIFT–BoW and CNN methods

1) *Conceptual similarities and differences:* Both SIFT and Convolutional Neural Networks extract features beyond raw pixel values, aiming to capture higher-order image structures. SIFT computes local gradient-based descriptors at interest points, encoding scale- and rotation-invariant information via Difference-of-Gaussian detection and orientation

histograms [9]. CNNs likewise learn convolutional filters that respond to edges and textures in early layers, but extend this to hierarchical, learned feature extractors across multiple non-linear layers [10].

However, whereas the traditional SIFT–BoW pipeline discards spatial arrangements—summarising an image by a histogram of visual words—CNNs inherently preserve spatial information through weight-sharing convolutions and pooling operations, enabling context-aware feature maps that support object localisation and global structure modelling. This endows CNNs with superior discriminative power, particularly on complex scenes where spatial relationships are critical.

2) Performance and resource trade-offs:

- **Accuracy**
Empirical benchmarks demonstrate that CNNs substantially outperform SIFT–BoW methods in classification tasks. On datasets such as CIFAR-10, CNNs achieve >90% accuracy, whereas SIFT–BoW with basic classifiers rarely exceeds 30% even after extensive parameter tuning [11]. This gulf may even widen on large-scale benchmarks (e.g. ImageNet, COCO).
- **Computational and Environmental Cost**
CNNs require vast computational resources for both training and inference, often necessitating GPU clusters and consuming tens to hundreds of GPU-hours. Their large model sizes and repeated convolution operations contribute to a substantial carbon footprint. In contrast, SIFT–BoW pipelines, relying on CPU-based descriptor extraction and k-means clustering, demand orders of magnitude less energy and can run on standard workstations with limited environmental impact.
- **Interpretability and Black-Box Concerns**
Traditional methods such as SIFT–BoW are inherently transparent: each stage—from keypoint detection to histogram quantisation—is algorithmically defined and readily interpretable. CNNs, by contrast, are often criticised as “black boxes” due to vast numbers of learned parameters and complex non-linear interactions.

VI. HYBRID APPROACHES

Recognising the complementary strengths of SIFT and CNN features, researchers have proposed hybrid models that concatenate or fuse SIFT descriptors with CNN feature maps. For instance, the SIFT-CNN framework embeds dense SIFT histograms as additional feature channels in early convolutional layers, improving robustness on texture-rich tasks while preserving spatial detail [12]. Similarly, hybrid face-recognition systems have demonstrated state-of-the-art performance by training CNNs jointly on pixel data and hand-crafted SIFT features, thereby leveraging the efficiency of SIFT on limited data and the expressiveness of CNNs on large-scale patterns [13].

These hybrid approaches validate that traditional and deep learning methods need not be mutually exclusive; instead, they can be integrated to balance accuracy, interpretability, and resource consumption.

VII. CONCLUSIONS AND FUTURE WORK

Our findings confirm that while SIFT–BoW remains a low-resource, interpretable approach—particularly suited to simple, well-structured tasks—its performance is fundamentally constrained on small, centred images and high-dimensional feature spaces. In contrast, CNNs deliver substantially higher accuracies by learning spatially aware, hierarchical features, albeit at the cost of greater computation and reduced transparency. Nonetheless, hybrid strategies that integrate hand-crafted descriptors into deep architectures may offer favourable trade-offs.

Looking ahead, we will extend our experiments to larger robotic-vision datasets, deploying both traditional algorithms and CNNs. We plan to explore advanced classifiers (e.g. Support Vector Machines, ensemble methods) and alternative feature detectors (e.g. ORB, SURF) to enhance the BoW pipeline. Concurrently, we will investigate more compact and interpretable CNN architectures, alongside transfer learning and self-supervised pretraining, to further boost performance and scalability in real-world robotic applications.

REFERENCES

- [1] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer vision—ECCV 2014: 13th European conference, Zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*. Springer, 2014, pp. 740–755.
- [2] S. Fanello, C. Ciliberto, M. Santoro, L. Natale, G. Metta, L. Rosasco, and F. Odone, “icub world: Friendly robots help building good vision data-sets,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2013, pp. 700–705.
- [3] T. Lindeberg, “Feature detection with automatic scale selection,” *International journal of computer vision*, vol. 30, pp. 79–116, 1998.
- [4] —, “Scale invariant feature transform,” 2012.
- [5] Sivic and Zisserman, “Video google: A text retrieval approach to object matching in videos,” in *Proceedings ninth IEEE international conference on computer vision*. IEEE, 2003, pp. 1470–1477.
- [6] S. Gidaris, A. Bursuc, N. Komodakis, P. Pérez, and M. Cord, “Learning representations by predicting bags of visual words,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6928–6938.
- [7] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
- [9] T. Qamar and N. Z. Bawany, “Understanding the black-box: towards interpretable and reliable deep learning models,” *PeerJ Computer Science*, vol. 9, p. e1629, 2023.
- [10] U. Özyayın, T. Georgiou, and M. Lew, “A comparison of cnn and classic features for image retrieval,” in *2019 International Conference on Content-Based Multimedia Indexing (CBMI)*. IEEE, 2019, pp. 1–4.
- [11] N. O’Mahony, S. Campbell, A. Carvalho, S. Harapanahalli, G. V. Hernandez, L. Krpalkova, D. Riordan, and J. Walsh, “Deep learning vs. traditional computer vision,” in *Advances in computer vision: proceedings of the 2019 computer vision conference (CVC), volume 1 I*. Springer, 2020, pp. 128–144.
- [12] D. Tsourounis, D. Kastaniotis, C. Theoharatos, A. Kazantzidis, and G. Economou, “Sift-cnn: when convolutional neural networks meet dense sift descriptors for image and sequence classification,” *Journal of Imaging*, vol. 8, no. 10, p. 256, 2022.
- [13] T. Connie, M. Al-Shabi, W. P. Cheah, and M. Goh, “Facial expression recognition using a hybrid cnn–sift aggregator,” in *International workshop on multi-disciplinary trends in artificial intelligence*. Springer, 2017, pp. 139–149.

APPENDIX

VIII. COMPLETE DATA TABLE

TABLE III: Hyperparameter Configurations and Model Accuracy for SIFT-BoW Classifier

Rd	Vocab Size	SIFT CT	SIFT Lyr	KNN N	NB Var	KNN Acc	NB Acc
1	100	0.04	3	5	1	0.1946	0.2734
2	300	0.02	4	5	1	0.1528	0.2918
3	500	0.02	4	5	1	0.1441	0.3067
4	500	0.01	4	5	1	0.1382	0.3124
5	600	0.01	4	5	1	0.1315	0.3055
6	300	0.01	4	5	1	0.1507	0.2976
7	400	0.01	4	5	1	0.1412	0.2968
8	800	0.01	4	5	1	0.1263	0.3132
9	600	0.01	4	5	1	0.1310	0.3081
10	1000	0.01	4	5	1	0.1271	0.3104
11	800	0.01	4	5	1	0.1238	0.3092
12	800	0.005	4	5	1	0.1171	0.3147
13	800	0.02	4	5	1	0.1382	0.3116
14	800	0.03	4	5	1	0.1356	0.3077
15	800	0.04	4	5	1	0.1382	0.3026
16	800	0.02	3	5	1	0.1378	0.3159
17	800	0.02	5	5	1	0.1349	0.3127
18	800	0.02	4	1	1	0.1399	0.3135
19	800	0.02	4	3	1	0.1379	0.3139
20	800	0.02	4	7	1	0.1376	0.3085
21	400	0.02	4	1	1	0.1548	0.3024
22	400	0.02	4	3	1	0.1542	0.3046
23	200	0.02	4	3	1	0.1780	
24	100	0.02	4	1	1	0.1970	

Note: Rd=Round, CT=Contrast Threshold (SIFT_CONTRAST $\times 100$), Lyr=Octave Layers (SIFT_OCTAVE_LAYERS), N=Neighbors (KNN_NEIGHBORS), Var=Variance Smoothing (NB_VAR_SMOOTH). All IMAGE_SIZE=64 \times 64 except Rd1=128 \times 128.

TABLE IV: Best Performance Across Datasets (SIFT-BoW Classifier)

Dataset	Size	CT	Vocab	Lyr	N	Var	KNN Acc	NB Acc
CIFAR-10	64	0.03	1024	3	5	1e-7	0.217	0.285
CIFAR-100	96	0.02	2048	4	3	1e-6	0.031	0.094
COCO	224	0.02	15 000	4	3	1e-6	0.161	0.221

Note: Size=IMAGE_SIZE, CT=SIFT_CONTRAST $\times 100$, Lyr=SIFT_OCTAVE_LAYERS, N=KNN_NEIGHBORS, Var=NB_VAR_SMOOTH. CIFAR-100 metrics from Round 2 (KNN) & Round 1 (NB), COCO metrics from Round 2. Training times: CIFAR-100 (NB 0.9s), COCO (NB 1.0s). Memory peaks: CIFAR-100 3.7GB, COCO 8.2GB.

IX. SOURCE CODE


```

1  import cv2
2  import numpy as np
3  import time
4  import psutil
5  from sklearn.cluster import MiniBatchKMeans
6  from sklearn.neighbors import KNeighborsClassifier
7  from sklearn.naive_bayes import GaussianNB
8  from sklearn.svm import SVC
9  from sklearn.metrics import confusion_matrix
10 from sklearn.metrics import accuracy_score, f1_score, precision_recall_fscore_support,
    average_precision_score
11 # from tensorflow.keras.datasets import cifar10
12 from tensorflow.keras.datasets import cifar100
13
14 # -----
15 # Hyperparameters (Descriptors)
16 IMAGE_SIZE = (96, 96) # Resize for SIFT richness vs. speed
17 VOCAB_SIZE = 2048 # Number of visual words (clusters)
18 SIFT_CONTRAST = 0.02 # Lower -> more keypoints (noisier)
19 SIFT_OCTAVE_LAYERS = 4 # More layers -> finer scale sampling
20 # Hyperparameters (Classifier)
21 KMEANS_BATCH = 6144 # MiniBatchKMeans batch size
22 KNN_NEIGHBORS = 3 # k in KNN
23 NB_VAR_SMOOTH = 1e-6 # Variance smoothing for GaussianNB
24 SVM_KERNEL = 'sigmoid' # Kernel for SVM
25 SVM_C = 1.0 # Regularization parameter for SVM
26
27 TYPE = 0 # 0: all, 1: KNN, 2: NB, 3: SVM
28
29 # -----
30 # Utility functions for measurements
31 process = psutil.Process()
32
33 def memory_mb():
34     """Return current memory usage of this process in MB."""
35     return process.memory_info().rss / (1024 * 1024)
36
37 # -----
38 # 1. Load and preprocess dataset
39 def load_preprocess():
40     # (x_tr, y_tr), (x_te, y_te) = cifar10.load_data()
41     (x_tr, y_tr), (x_te, y_te) = cifar100.load_data(label_mode='fine')
42
43     y_tr, y_te = y_tr.flatten(), y_te.flatten()
44     def prep(images):
45         out = []
46         for img in images:
47             gray = cv2.cvtColor(img, cv2.COLOR_RGB2GRAY) # Convert to grayscale
48             out.append(cv2.resize(gray, IMAGE_SIZE))
49         return out
50     return prep(x_tr), y_tr, prep(x_te), y_te
51
52 # -----
53 # 2. Extract SIFT descriptors
54 def extract_descriptors(images, sift, max_des_per_image=100):
55     # Always return a list of length = len(images),
56     # with each element either a (N 128) array or None.
57     descriptors = []
58     for img in images:
59         _, des = sift.detectAndCompute(img, None)
60         # if des is not None:
61         #     descriptors.append(des)
62         if des is not None:
63             if len(des) > max_des_per_image:
64                 des = des[np.random.choice(len(des), max_des_per_image, replace=False)]
65             descriptors.append(des)
66     return descriptors

```

Listing 1: SIFT-BoW-CIFAR PART 1

```

1 # -----
2 # 3. Build BoW vocabulary
3
4 def build_vocabulary(des_list):
5     # Keep only proper descriptors of shape (num_keys, 128)
6
7     # valid = [des for des in des_list
8     #           if des is not None
9     #           and isinstance(des, np.ndarray)
10    #           and des.ndim == 2
11    #           and des.shape[1] == 128]
12
13    # if len(valid) == 0:
14    #     raise ValueError("No valid SIFT descriptors found.")
15
16    # total_des = sum(len(des) for des in valid)
17    # if total_des < VOCAB_SIZE:
18    #     raise ValueError(f"Total descriptors ({total_des}) < VOCAB_SIZE ({VOCAB_SIZE}). "
19    #                       "Reduce VOCAB_SIZE or adjust SIFT parameters.")
20
21    # stacked = np.vstack(valid)
22    # print(f"Total valid descriptors: {len(stacked)}")
23
24    valid = []
25    for des in des_list:
26        if des is not None and len(des) > 0:
27            valid.append(des)
28    stacked = np.vstack(valid)
29
30    if len(stacked) < VOCAB_SIZE:
31        raise ValueError(f"Total descriptors ({len(stacked)}) < VOCAB_SIZE ({VOCAB_SIZE}). "
32                          "Adjust SIFT parameters or reduce VOCAB_SIZE.")
33
34    kmeans = MiniBatchKMeans(n_clusters=VOCAB_SIZE,
35                             batch_size=KMEANS_BATCH,
36                             random_state=42,
37                             verbose=1,
38                             )
39
40    # kmeans.fit(stacked) # Fit on all descriptors at once
41
42    # Use partial_fit for large datasets
43    # for i, des_batch in enumerate(valid):
44    #     kmeans.partial_fit(des_batch)
45    #     if i % 100 == 0:
46    #         print(f"Processed {i+1}/{len(valid)} batches | Mem: {memory_mb():.1f}MB")
47    n_batches = len(stacked) // KMEANS_BATCH + 1
48    for i in range(n_batches):
49        start = i * KMEANS_BATCH
50        end = start + KMEANS_BATCH
51        batch = stacked[start:end]
52        if len(batch) == 0:
53            break
54        kmeans.partial_fit(batch)
55
56        # if i % 100 == 0 or i == n_batches - 1:
57        #     print(f"Processed {i+1}/{n_batches} batches | Mem: {memory_mb():.1f}MB")
58        # print(f"Clustering: {end}/{len(stacked)} samples processed")
59
60    return kmeans
61
62 # -----
63 # 4. Compute BoW histograms
64
65 def histograms(des_list, kmeans):
66     hists = np.zeros((len(des_list), VOCAB_SIZE), dtype=int)
67     for i, des in enumerate(des_list):
68         if des is not None:
69             words = kmeans.predict(des)
70             hist, _ = np.histogram(words, bins=np.arange(VOCAB_SIZE+1))
71             hists[i] = hist
72     return hists
73

```

```

1 # -----
2 # 5. Train and evaluate classifiers
3 def evaluate(X_train, y_train, X_test, y_test, type=1):
4     results = {}
5
6     # For CIFAR-100
7     num_classes = len(np.unique(y_train))
8
9     if type == 0 or 1:
10         # KNN
11         knn = KNeighborsClassifier(n_neighbors=KNN_NEIGHBORS,
12                                   weights='distance',
13                                   metric='euclidean',
14                                   n_jobs=-1)
15
16         t0 = time.time()
17         knn.fit(X_train, y_train)
18         t_train = time.time() - t0
19         t0_inf = time.time()
20         y_pred_knn = knn.predict(X_test)
21         t_inf = time.time() - t0_inf
22         results['knn'] = {
23             'train_time_s': t_train,
24             'inf_time_s': t_inf,
25             'mem_mb': memory_mb(),
26             'accuracy': accuracy_score(y_test, y_pred_knn),
27             'f1_macro': f1_score(y_test, y_pred_knn, average='macro'),
28             'confusion_matrix': confusion_matrix(y_test, y_pred_knn),
29             # For CIFAR-10
30             # 'map_macro': average_precision_score(
31             #     np.eye(len(np.unique(y_test))) [y_test],
32             #     np.eye(len(np.unique(y_test))) [y_pred_knn],
33             #     average='macro')
34             # For CIFAR-100
35             'map_macro': average_precision_score(
36                 np.eye(num_classes) [y_test],
37                 np.eye(num_classes) [y_pred_knn],
38                 average='macro')
39         }
40
41     if type == 0 or 2:
42         # Naive Bayes
43         nb = GaussianNB(var_smoothing=NB_VAR_SMOOTH)
44         t0 = time.time()
45         nb.fit(X_train, y_train)
46         t_train = time.time() - t0
47         t0_inf = time.time()
48         y_pred_nb = nb.predict(X_test)
49         t_inf = time.time() - t0_inf
50         results['nb'] = {
51             'train_time_s': t_train,
52             'inf_time_s': t_inf,
53             'mem_mb': memory_mb(),
54             'accuracy': accuracy_score(y_test, y_pred_nb),
55             'f1_macro': f1_score(y_test, y_pred_nb, average='macro'),
56             'confusion_matrix': confusion_matrix(y_test, y_pred_nb),
57             'map_macro': average_precision_score(
58                 np.eye(len(np.unique(y_test))) [y_test],
59                 np.eye(len(np.unique(y_test))) [y_pred_nb],
60                 average='macro')
61         }

```

Listing 3: SIFT-BoW-CIFAR PART 3

```

1      # if type == 0 or 3:
2      #     # SVM
3      #     svm = SVC(kernel= SVM_KERNEL, probability=True, C=SVM_C, random_state=42)
4      #     t0 = time.time()
5      #     svm.fit(X_train, y_train)
6      #     t_train = time.time() - t0
7      #     t0_inf = time.time()
8      #     y_pred_svm = svm.predict(X_test)
9      #     t_inf = time.time() - t0_inf
10     #     results['svm'] = {
11     #         'train_time_s': t_train,
12     #         'inf_time_s': t_inf,
13     #         'mem_mb': memory_mb(),
14     #         'accuracy': accuracy_score(y_test, y_pred_svm),
15     #         'f1_macro': f1_score(y_test, y_pred_svm, average='macro'),
16     #         'confusion_matrix': confusion_matrix(y_test, y_pred_svm),
17     #         'map_macro': average_precision_score(
18     #             np.eye(len(np.unique(y_test))) [y_test],
19     #             np.eye(len(np.unique(y_test))) [y_pred_svm],
20     #             average='macro')
21     #     }
22
23     return results
24
25 # -----
26 # Main pipeline
27 if __name__ == "__main__":
28     # Initialize SIFT
29     sift = cv2.SIFT_create(contrastThreshold=SIFT_CONTRAST,
30                           nOctaveLayers=SIFT_OCTAVE_LAYERS)
31
32     # Load data
33     x_train, y_train, x_test, y_test = load_preprocess()
34
35     # Descriptor extraction
36     t0 = time.time()
37     des_train = extract_descriptors(x_train, sift)
38     des_test = extract_descriptors(x_test, sift)
39     print(f"Descriptor extraction time: {time.time()-t0:.2f}s, Mem: {memory_mb():.1f}MB")
40
41     # Vocabulary
42     t0 = time.time()
43     kmeans = build_vocabulary(des_train)
44     print(f"Vocabulary building time: {time.time()-t0:.2f}s, Mem: {memory_mb():.1f}MB")
45
46     # Histograms
47     X_train = histograms(des_train, kmeans)
48     X_test = histograms(des_test, kmeans)
49
50     # Evaluate
51     results = evaluate(X_train, y_train, X_test, y_test, TYPE)
52     for model, stats in results.items():
53         print(f"\nModel: {model.upper()}")
54         for k, v in stats.items():
55             print(f"    {k}: {v}")

```

Listing 4: SIFT-BoW-CIFAR PART 4

```

1  import torch
2  import torch.nn as nn
3  import torch.nn.functional as F
4  import torch.optim as optim
5  from torchvision import datasets, transforms
6  from torch.utils.data import DataLoader
7
8  # 1. Device configuration
9  device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
10 print("Using device:", device)
11
12 # 2. Hyperparameters
13 batch_size = 64
14 learning_rate = 1e-3
15 num_epochs = 100
16
17 # 3. Data transforms and loaders
18 # CIFAR-10
19 # transform = transforms.Compose([
20 #     transforms.ToTensor(),
21 #     transforms.Normalize((0.4914, 0.4822, 0.4465),
22 #                           (0.2023, 0.1994, 0.2010)),
23 # ])
24 # CIFAR-10 mean and std
25
26 # train_ds = datasets.CIFAR10(root='data',
27 #                               train=True,
28 #                               download=True,
29 #                               transform=transform)
30 # test_ds = datasets.CIFAR10(root='data',
31 #                              train=False,
32 #                              download=True,
33 #                              transform=transform)
34
35 # train_loader = DataLoader(train_ds, batch_size=batch_size, shuffle=True, num_workers=0)
36 # test_loader = DataLoader(test_ds, batch_size=batch_size, shuffle=False, num_workers=0)
37
38 # CIFAR-100
39 transform = transforms.Compose([
40     transforms.ToTensor(),
41     transforms.Normalize((0.5071, 0.4867, 0.4408), # CIFAR-100 mean :contentReference[
42         oaicite:4]{index=4}
43         (0.2675, 0.2565, 0.2761)), # CIFAR-100 std :contentReference[
44         oaicite:5]{index=5}
45 ])
46
47 train_ds = datasets.CIFAR100(root='data',
48                               train=True,
49                               download=True,
50                               transform=transform)
51 test_ds = datasets.CIFAR100(root='data',
52                              train=False,
53                              download=True,
54                              transform=transform)
55
56 train_loader = DataLoader(train_ds, batch_size=batch_size, shuffle=True, num_workers=0)
57 test_loader = DataLoader(test_ds, batch_size=batch_size, shuffle=False, num_workers=0)

```

Listing 5: VGG-CNN-CIFAR PART 1

```

1 # 4. Model definition ( V G G like )
2 # class CIFAR10VGG(nn.Module):
3 #     def __init__(self):
4 #         super().__init__()
5 #         # Block 1
6 #         self.conv1 = nn.Conv2d(3, 64, kernel_size=3, padding=1)
7 #         self.bn1 = nn.BatchNorm2d(64)
8 #         self.conv2 = nn.Conv2d(64, 64, kernel_size=3, padding=1)
9 #         self.bn2 = nn.BatchNorm2d(64)
10 #         # Block 2
11 #         self.conv3 = nn.Conv2d(64, 128, kernel_size=3, padding=1)
12 #         self.bn3 = nn.BatchNorm2d(128)
13 #         self.conv4 = nn.Conv2d(128, 128, kernel_size=3, padding=1)
14 #         self.bn4 = nn.BatchNorm2d(128)
15 #         # Block 3
16 #         self.conv5 = nn.Conv2d(128, 256, kernel_size=3, padding=1)
17 #         self.bn5 = nn.BatchNorm2d(256)
18 #         self.conv6 = nn.Conv2d(256, 256, kernel_size=3, padding=1)
19 #         self.bn6 = nn.BatchNorm2d(256)
20 #         # Block 4
21 #         self.conv7 = nn.Conv2d(256, 512, kernel_size=3, padding=1)
22 #         self.bn7 = nn.BatchNorm2d(512)
23 #         self.conv8 = nn.Conv2d(512, 512, kernel_size=3, padding=1)
24 #         self.bn8 = nn.BatchNorm2d(512)
25
26 #         # Classifier head
27 #         self.fc1 = nn.Linear(512*2*2, 512)
28 #         self.dropout = nn.Dropout(0.5)
29 #         self.fc2 = nn.Linear(512, 10)
30
31 #     def forward(self, x):
32 #         # Block 1
33 #         x = F.relu(self.bn1(self.conv1(x)))
34 #         x = F.relu(self.bn2(self.conv2(x)))
35 #         x = F.max_pool2d(x, 2) # 16 16
36
37 #         # Block 2
38 #         x = F.relu(self.bn3(self.conv3(x)))
39 #         x = F.relu(self.bn4(self.conv4(x)))
40 #         x = F.max_pool2d(x, 2) # 8 8
41
42 #         # Block 3
43 #         x = F.relu(self.bn5(self.conv5(x)))
44 #         x = F.relu(self.bn6(self.conv6(x)))
45 #         x = F.max_pool2d(x, 2) # 4 4
46
47 #         # Block 4
48 #         x = F.relu(self.bn7(self.conv7(x)))
49 #         x = F.relu(self.bn8(self.conv8(x)))
50 #         x = F.max_pool2d(x, 2) # 2 2
51
52 #         x = x.view(x.size(0), -1) # flatten
53 #         x = F.relu(self.fc1(x))
54 #         x = self.dropout(x)
55 #         x = self.fc2(x)
56 #         return x

```

Listing 6: VGG-CNN-CIFAR PART 2


```

1 class CIFAR100VGG(nn.Module):
2     def __init__(self):
3         super().__init__()
4         # Block 1
5         self.conv1 = nn.Conv2d(3, 64, kernel_size=3, padding=1)
6         self.bn1 = nn.BatchNorm2d(64)
7         self.conv2 = nn.Conv2d(64, 64, kernel_size=3, padding=1)
8         self.bn2 = nn.BatchNorm2d(64)
9         # Block 2
10        self.conv3 = nn.Conv2d(64, 128, kernel_size=3, padding=1)
11        self.bn3 = nn.BatchNorm2d(128)
12        self.conv4 = nn.Conv2d(128, 128, kernel_size=3, padding=1)
13        self.bn4 = nn.BatchNorm2d(128)
14        # Block 3
15        self.conv5 = nn.Conv2d(128, 256, kernel_size=3, padding=1)
16        self.bn5 = nn.BatchNorm2d(256)
17        self.conv6 = nn.Conv2d(256, 256, kernel_size=3, padding=1)
18        self.bn6 = nn.BatchNorm2d(256)
19        # Block 4
20        self.conv7 = nn.Conv2d(256, 512, kernel_size=3, padding=1)
21        self.bn7 = nn.BatchNorm2d(512)
22        self.conv8 = nn.Conv2d(512, 512, kernel_size=3, padding=1)
23        self.bn8 = nn.BatchNorm2d(512)
24
25        # Classifier head
26        self.fc1 = nn.Linear(512*2*2, 512)
27        self.dropout = nn.Dropout(0.5)
28        self.fc2 = nn.Linear(512, 100) # Output dim = 100
29
30    def forward(self, x):
31        # Block 1
32        x = F.relu(self.bn1(self.conv1(x)))
33        x = F.relu(self.bn2(self.conv2(x)))
34        x = F.max_pool2d(x, 2) # 16 16
35
36        # Block 2
37        x = F.relu(self.bn3(self.conv3(x)))
38        x = F.relu(self.bn4(self.conv4(x)))
39        x = F.max_pool2d(x, 2) # 8 8
40
41        # Block 3
42        x = F.relu(self.bn5(self.conv5(x)))
43        x = F.relu(self.bn6(self.conv6(x)))
44        x = F.max_pool2d(x, 2) # 4 4
45
46        # Block 4
47        x = F.relu(self.bn7(self.conv7(x)))
48        x = F.relu(self.bn8(self.conv8(x)))
49        x = F.max_pool2d(x, 2) # 2 2
50
51        x = x.view(x.size(0), -1) # flatten
52        x = F.relu(self.fc1(x))
53        x = self.dropout(x)
54        x = self.fc2(x)
55        return x
56    }

```

Listing 7: VGG-CNN-CIFAR PART 3

```

1  # model = CIFAR10VGG().to(device)
2  model = CIFAR100VGG().to(device)
3
4  # 5. Loss and optimizer
5  criterion = nn.CrossEntropyLoss()
6  optimizer = optim.Adam(model.parameters(), lr=learning_rate)
7
8  # 6. Training loop
9  for epoch in range(1, num_epochs+1):
10     model.train()
11     running_loss = 0.0
12     for images, labels in train_loader:
13         images, labels = images.to(device), labels.to(device)
14
15         optimizer.zero_grad()
16         outputs = model(images)
17         loss = criterion(outputs, labels)
18         loss.backward()
19         optimizer.step()
20
21         running_loss += loss.item() * images.size(0)
22
23     epoch_loss = running_loss / len(train_loader.dataset)
24     print(f"Epoch {epoch}/{num_epochs}, Loss: {epoch_loss:.4f}")
25
26     # Validation
27     model.eval()
28     correct = 0
29     with torch.no_grad():
30         for images, labels in test_loader:
31             images, labels = images.to(device), labels.to(device)
32             preds = model(images).argmax(dim=1)
33             correct += (preds == labels).sum().item()
34     acc = correct / len(test_loader.dataset)
35     print(f"  Test Accuracy: {acc:.4f}")
36
37 # 7. Final evaluation
38 print("Training complete.")
39 # Save the model
40 torch.save(model.state_dict(), 'cifar100_vgg.pth')

```

Listing 8: VGG-CNN-CIFAR PART 4

```

1  import os
2  import cv2
3  import numpy as np
4  import time
5  import psutil
6  from pycocotools.coco import COCO
7  from sklearn.cluster import MiniBatchKMeans
8  from sklearn.neighbors import KNeighborsClassifier
9  from sklearn.naive_bayes import GaussianNB
10 from sklearn.metrics import accuracy_score, f1_score, confusion_matrix,
    average_precision_score
11
12 # -----
13 # Hyperparameters (Descriptors)
14 IMAGE_SIZE = (224, 224)          # Resize for SIFT richness vs. speed
15 VOCAB_SIZE = 18000              # Number of visual words (clusters)
16 SIFT_CONTRAST = 0.02            # Lower -> more keypoints (noisier)
17 SIFT_OCTAVE_LAYERS = 4          # More layers -> finer scale sampling
18 # Hyperparameters (Classifier)
19 KMEANS_BATCH = 20000            # MiniBatchKMeans batch size
20 KNN_NEIGHBORS = 3              # k in KNN
21 NB_VAR_SMOOTH = 1e-6           # Variance smoothing for GaussianNB
22 # -----
23 # COCO dataset paths
24 COCO_IMG_DIR = 'C:/Customize/Tool-Coding/Dataset/COCO/val2017/val2017'
25 COCO_ANN_FILE = 'C:/Customize/Tool-Coding/Dataset/COCO/annotations_trainval2017/annotations/
    instances_val2017.json'
26
27 # -----
28 # Utility functions
29 process = psutil.Process()
30 def memory_mb(): return process.memory_info().rss / (1024 * 1024)
31
32 # -----
33 # 1. Load and preprocess COCO classification dataset
34 #   Map each image to its first annotated category
35
36 def load_coco_classification(img_dir, ann_file, max_images=None):
37     coco = COCO(ann_file)
38     img_ids = coco.getImgIds()[:max_images]
39     images, labels = [], []
40     for img_id in img_ids:
41         info = coco.loadImgs(img_id)[0]
42         anns = coco.loadAnns(coco.getAnnIds(imgIds=img_id, iscrowd=False))
43         if not anns: continue
44         # pick first annotation
45         cat_id = anns[0]['category_id']
46         img_path = os.path.join(img_dir, info['file_name'])
47         img = cv2.imread(img_path)
48         if img is None: continue
49         gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
50         gray = cv2.resize(gray, IMAGE_SIZE)
51         images.append(gray)
52         labels.append(cat_id)
53     return images, np.array(labels)
54
55 # -----
56 # 2. Extract SIFT descriptors (fixed per image)
57 def extract_descriptors(images, sift, max_per=100):
58     descs = []
59     for img in images:
60         _, des = sift.detectAndCompute(img, None)
61         if des is not None and len(des) > max_per:
62             idx = np.random.choice(len(des), max_per, replace=False)
63             des = des[idx]
64         descs.append(des)
65     return descs

```

Listing 9: SIFT-BoW-COCO PART 1

```

1  # -----
2  # 3. Build BoW vocabulary
3  def build_vocabulary(des_list):
4      valid = [d for d in des_list if d is not None and len(d)>0]
5      stacked = np.vstack(valid)
6      if len(stacked) < VOCAB_SIZE:
7          raise ValueError(f"Not enough descriptors: {len(stacked)} < {VOCAB_SIZE}")
8      kmeans = MiniBatchKMeans(n_clusters=VOCAB_SIZE, batch_size=KMEANS_BATCH, random_state=42)
9      # partial fit in batches
10     for i in range(0, len(stacked), KMEANS_BATCH):
11         kmeans.partial_fit(stacked[i:i+KMEANS_BATCH])
12     return kmeans
13
14  # -----
15  # 4. Compute BoW histograms
16  def histograms(des_list, kmeans):
17      H = np.zeros((len(des_list), VOCAB_SIZE), dtype=int)
18      for i, des in enumerate(des_list):
19          if des is not None and len(des)>0:
20              words = kmeans.predict(des)
21              hist, _ = np.histogram(words, bins=np.arange(VOCAB_SIZE+1))
22              H[i] = hist
23      return H
24
25  # -----
26  # 5. Train and evaluate classifiers
27  def evaluate(X_tr, y_tr, X_te, y_te):
28      results = {}
29      # KNN
30      knn = KNeighborsClassifier(n_neighbors=KNN_NEIGHBORS, weights='distance', n_jobs=-1)
31      t0 = time.time(); knn.fit(X_tr, y_tr); t_train = time.time()-t0
32      t0 = time.time(); y_pred = knn.predict(X_te); t_inf = time.time()-t0
33      results['knn'] = {'train_time_s':t_train,
34                      'inf_time_s':t_inf,
35                      'accuracy':accuracy_score(y_te, y_pred),
36                      'f1_macro':f1_score(y_te, y_pred, average='macro')
37                      # 'cm':confusion_matrix(y_te, y_pred)
38                      }
39      # NB
40      nb = GaussianNB(var_smoothing=NB_VAR_SMOOTH)
41      t0 = time.time(); nb.fit(X_tr, y_tr); t_train = time.time()-t0
42      t0 = time.time(); y_pred = nb.predict(X_te); t_inf = time.time()-t0
43      results['nb'] = {'train_time_s':t_train,
44                     'inf_time_s':t_inf,
45                     'accuracy':accuracy_score(y_te, y_pred),
46                     'f1_macro':f1_score(y_te, y_pred, average='macro')
47                     # 'cm':confusion_matrix(y_te, y_pred)
48                     }
49      return results
50
51  # -----
52  if __name__ == '__main__':
53      # Initialize SIFT
54      sift = cv2.SIFT_create(contrastThreshold=SIFT_CONTRAST, nOctaveLayers=SIFT_OCTAVE_LAYERS)
55      # Load COCO images and labels
56      x_train, y_train = load_coco_classification(COCO_IMG_DIR, COCO_ANN_FILE, max_images=5000)
57      x_test, y_test = load_coco_classification(COCO_IMG_DIR, COCO_ANN_FILE, max_images=1000)
58      print(f"Loaded train:{len(y_train)} test:{len(y_test)} imgs")
59      # Extract descriptors
60      des_tr = extract_descriptors(x_train, sift)
61      des_te = extract_descriptors(x_test, sift)
62      print("Descriptor extraction done")
63      # Build vocabulary
64      kmeans = build_vocabulary(des_tr)
65      # Compute histograms
66      X_tr = histograms(des_tr, kmeans)
67      X_te = histograms(des_te, kmeans)
68      # Evaluate
69      res = evaluate(X_tr, y_train, X_te, y_test)
70      for m,stats in res.items():
71          print(f"\nModel: {m}", stats)

```

Listing 10: SIFT-BoW-COCO PART 2

```

1  import os
2  from PIL import Image
3  import torch
4  import torch.nn as nn
5  import torch.nn.functional as F
6  import torch.optim as optim
7  from torchvision.datasets import CocoDetection
8  from torchvision import transforms
9  from torch.utils.data import DataLoader, Dataset
10 from pycocotools.coco import COCO # type: ignore
11
12 # 1. Device configuration
13 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
14 print("Using device:", device)
15
16 # 2. Hyperparameters
17 batch_size = 32
18 learning_rate = 1e-3
19 num_epochs = 50
20 num_classes = 80 # COCO has 80 object categories
21
22 # 3. Data transforms (resize + normalise for pretrained CNNs)
23 transform = transforms.Compose([
24     transforms.Resize((224, 224)),
25     transforms.ToTensor(),
26     transforms.Normalize(mean=[0.485, 0.456, 0.406],
27                           std=[0.229, 0.224, 0.225]),
28 ])
29
30 # 4. Custom dataset subclassing CocoDetection correctly
31 class CocoClassification(CocoDetection):
32     def __init__(self, img_folder, ann_file, transform=None, max_images=None):
33         super().__init__(img_folder, ann_file, transforms=None) # initialise base
34         self.transform = transform
35         # optionally limit size
36         if max_images:
37             self.ids = self.ids[:max_images] # crop to first N images
38
39     def __getitem__(self, idx):
40         # Leverage base class to load PIL image and raw annotations
41         img, raw_anns = super().__getitem__(idx) # returns PIL.Image and list
42             of dicts
43         if not raw_anns:
44             # fallback if no annotation
45             label = 0
46         else:
47             raw_cat_id = raw_anns[0]['category_id']
48             label = cat_to_idx[raw_cat_id] # maps to [0,79]
49         # apply transforms if provided
50         if self.transform:
51             img = self.transform(img)
52         return img, label
53
54 # 5. Paths
55 data_root = 'C:/Customize/Tool-Coding/Dataset/COCO/val2017/val2017'
56 ann_file = 'C:/Customize/Tool-Coding/Dataset/COCO/annotations_trainval2017/annotations/
57     instances_val2017.json'
58
59 # 6. Datasets and loaders
60 train_ds = CocoClassification(data_root, ann_file, transform=transform, max_images=5000)
61 test_ds = CocoClassification(data_root, ann_file, transform=transform, max_images=1000)
62
63 train_loader = DataLoader(train_ds, batch_size=batch_size, shuffle=True, num_workers=0)
64 test_loader = DataLoader(test_ds, batch_size=batch_size, shuffle=False, num_workers=0)
65
66 coco = COCO(ann_file)
67 cat_ids = sorted(coco.getCatIds())
68 cat_to_idx = {cid: idx for idx, cid in enumerate(cat_ids)}

```

Listing 11: VGG-CNN-COCO PART 1

```

1
2 # 7. Simple VGG style CNN
3 class SimpleVGG(nn.Module):
4     def __init__(self, num_classes):
5         super().__init__()
6         self.features = nn.Sequential(
7             # block1
8             nn.Conv2d(3, 64, kernel_size=3, padding=1), nn.ReLU(inplace=True),
9             nn.Conv2d(64, 64, kernel_size=3, padding=1), nn.ReLU(inplace=True),
10            nn.MaxPool2d(2), # 112 112
11            # block2
12            nn.Conv2d(64, 128, 3, padding=1), nn.ReLU(inplace=True),
13            nn.Conv2d(128, 128, 3, padding=1), nn.ReLU(inplace=True),
14            nn.MaxPool2d(2), # 56 56
15        )
16        self.classifier = nn.Sequential(
17            nn.Flatten(),
18            nn.Linear(128*56*56, 1024), nn.ReLU(inplace=True), nn.Dropout(0.5),
19            nn.Linear(1024, num_classes)
20        )
21
22    def forward(self, x):
23        x = self.features(x)
24        x = self.classifier(x)
25        return x
26
27 model = SimpleVGG(num_classes).to(device)
28
29 # 8. Loss and optimiser
30 criterion = nn.CrossEntropyLoss()
31 optimizer = optim.Adam(model.parameters(), lr=learning_rate)
32
33 # 9. Training and evaluation loop
34 for epoch in range(1, num_epochs+1):
35     model.train()
36     running_loss = 0.0
37     for images, labels in train_loader:
38         images, labels = images.to(device), labels.to(device)
39         optimizer.zero_grad()
40         loss = criterion(model(images), labels)
41         loss.backward()
42         optimizer.step()
43         running_loss += loss.item() * images.size(0)
44     avg_loss = running_loss / len(train_loader.dataset)
45     print(f"Epoch {epoch}/{num_epochs} Training Loss: {avg_loss:.4f}")
46
47     model.eval()
48     correct = 0
49     with torch.no_grad():
50         for images, labels in test_loader:
51             images, labels = images.to(device), labels.to(device)
52             preds = model(images).argmax(dim=1)
53             correct += (preds == labels).sum().item()
54     accuracy = correct / len(test_loader.dataset)
55     print(f"Test Accuracy: {accuracy:.4f}")
56
57 print("Training complete.")
58
59 }

```

Listing 12: VGG-CNN-COCO PART 2