

# SIFT And CNN in Objects Classification

COMP61342&41342 Cognitive Robotics And Computer Vision Assignment Report

Yuyi Yang  
*yuyi.yang@postgrad.manchester.ac.uk*  
*p68393yy*

## I. INTRODUCTION

In recent years, equipping robots with reliable perception capabilities has become a foundation for autonomous manipulation, navigation, and human–robot interaction.

Traditional computer-vision pipelines, based on hand-crafted local features (e.g. SIFT) and classifiers (e.g. SVM, K-NN), offer interpretability and low inference latency; however, they may falter when faced with dynamic and complex environment and large intra-class variability. Conversely, convolutional neural networks (CNNs) have revolutionized vision performance, automatically learning hierarchical features, but at the cost of greater computational demand and complex hyperparameter tuning. Despite extensive research in object recognition, systematic comparisons between traditional handcrafted methods and CNN-based approaches remain limited. Thus, this work aims to fill this gap by systematically benchmarking and analyzing these paradigms, providing practical insights for researchers and practitioners in robotic perception.

This work systematically compares CNN models and Local Feature methods (i.e., SIFT+BoVW) on the *CIFAR-10* [1] and *iCubWorld 1.0* [2] benchmark, researching how choice of feature extractor, classifier, network architecture, and training strategy affect performance of classification accuracy and inference speed, and model robustness to environmental changes.

## II. LITERATURE REVIEW: STATE OF THE ART

Computer vision technology plays a crucial role in robotic autonomy tasks, including target recognition, manipulation, navigation, and human-robot interaction. Early robotic vision methods typically relied on hand-crafted localized features, such as by SIFT [3], [4], SURF [5] and ORB [6] extracting from the image feature points with interpretability, as well as certain invariance (to illumination, rotation and scaling) to achieve simple recognition and classification of objects. However, these methods rely too much on pre-defined features and lack the ability to represent more abstract features, making them difficult to extract effective features when facing more dynamic and variable complex real-world environments [7].

And over the past 15 years, deep learning, especially convolutional neural networks (CNNs), has dramatically reshaped computer vision research [8]. Its ability to automatically learn hierarchical feature representations of

images through multilayer nonlinear transformations has led to significant advances in computer vision tasks, enabling a leapfrog over traditional methods [9]. However, as the depth of the network increases, the problem of gradient vanishing or performance degradation in CNNs during training becomes more and more serious. In this context, the proposal of ResNet [10] provides a new idea for deep network training, which uses a residual learning mechanism to alleviate the degradation problem in deep network training by introducing skip connections to allow shallower features and the gradients flow directly to the deeper network, thus effectively alleviating the problem of gradient vanishing or explosion in deep network training.

And in recent years, the Attention mechanism and the Transformer model [11] have made breakthroughs in the field of computer vision, which is able to capture long-distance dependencies in data through its efficient Self-Attention mechanism and demonstrates excellent scalability to ultra-large capacity networks and large datasets. These advantages have led to exciting advances in tasks such as image classification, target detection, and semantic segmentation using Transformer networks [12]. Vision Transformer [13] is the first to enable the use of Transformer, which “completely” replace standard convolution in deep neural networks on large-scale image datasets and demonstrated performance that rivaled or even bettered CNNs. In addition, Swin Transformer [14] has further improved the computational efficiency of the Transformer model by introducing a window self-attention mechanism, achieving leading results in several visual tasks. Furthermore, recent researchers have proposed “Norm-free” Transformer [15], NoProp [16], and Model-Seeking [17] as an alternative to Transformer model (with no normalization) and learning mechanisms such as back/forward-propagation have been improved to provide new ideas for research in the field of Computer Vision.

## III. METHODOLOGY

### A. Dataset

**CIFAR-10** and ***iCubWorld 1.0*** are two datasets used in this project for object classification tasks. **CIFAR-10** is a well-known dataset containing 60,000 32x32 color images (10,000 for test) in 10 classes, with 6,000 images per class. While, ***iCubWorld 1.0*** is a dataset that objects are manually picked up for human and captured by iCub

robot. It includes 10 categories, each category contains 3 instances in the training set. Both datasets are famous benchmark in the field of computer vision, which is very suitable to compare the performance of local feature and deep learning methods. For each dataset, I applied corresponding data pre-processing methods separately, trying to improve the accuracy and robustness of the model.

In addition, *iCubWorld 1.0* contains 4 different testsets to fit different scenarios, including: CATEGORIZATION, DEMONSTRATOR, ROBOT and BACKGROUND. In order to ensure the objectivity and precision of the experiments and to improve the efficiency of the experiments, only the CATEGORIZATION (marked as *CATE*) and DEMONSTRATOR (marked as *DEMO*) testsets are used in this project to evaluate the performance of the model. The *CATE* testset is used to evaluate the generalization ability of the model when facing new instances, while the *DEMO* testset is used to evaluate the generalization ability of the model when facing same instances with different demonstrators (background). Thus, a comprehensive consideration based on these two test sets is sufficient to assess the performance and generalization ability of the models.

Additionally, I apply different resizing approaches, since the image size variance. For *CIFAR-10*, the original image size is only 32x32, which is definitely too small for SIFT method (patch size is 16x16). Therefore, I resized the images to 128x128 to ensure that it can get enough SIFT descriptors from the images. For *iCubWorld 1.0*, the original image size is 640x480, and all of the training images are pre-cropped to 160x160. In this case, I also applied resizing method to resize the image to 224x224, which better fits the input size of ResNet models.

### B. SIFT+BoVW

In this section, the brief workflow of SIFT+BoVW method will be introduced.

1) *Keypoints Detection*: The first step is to detect interest points in the image, which are more likely to be stable and distinctive, making them suitable for matching across different images. In order to detect these interest points for the following SIFT descriptor extraction, this project tries two different methods: Difference of Gaussian (**DoG**) and **Harris-Laplace**.

In computer vision, scale invariance is one of the key requirements for feature detection. And DoG is an efficient scale-space keypoint detection method, which approximates the Laplace-Gaussian (LoG) operator by making a subtraction between different Gaussian smoothing scales and octaves, and searches for spatial-scale extreme points in the generated DoG pyramid, obtaining points of interest with invariance to scale, which in turn significantly reduces the amount of computation.

While, for the Harris-Laplace method, its central idea lies in effectively combining the spatially accurate Harris corner detector with the scale-invariant properties of the scale-normalized Laplacian-of-Gaussian (LoG) operator.

This ensures that detected interest points not only have high spatial precision but are also stable across scales. Its brief workflow is as follows:

1. Apply the Harris corner detection independently to the input image at several scales (constructed via Gaussian smoothing) to provide precise spatial localization of corner-like structures across each individual scale.
2. For each corner candidate detected, the LoG response is computed across scales. The scale at which the candidate corner exhibits the strongest (local extrema) response is selected as its characteristic scale, ensuring robust scale invariance.

By this way, the Harris-Laplace method effectively emphasizes corner-type interest points, which typically convey richer structural information compared to edge points, while simultaneously preserving the desirable property of scale invariance inherent in LoG-based approaches. In addition, the traditional **Harris** algorithm is also tested, and their results are compared in the Experiments section.

2) *SIFT Descriptor Extraction*: After gaining the interest points in Section III-B1, the project manually implements the various steps of **SIFT** descriptor extraction. The following will briefly describe the SIFT workflow in the project:

1. Gradient magnitude and direction are calculated at each keypoint neighborhood through a 16\*16 region. Then, an orientation histogram with 36 (or 8) bins (each covering 10°) is constructed. The peak direction from this histogram is selected as the keypoint's main orientation, ensuring rotation invariance.
2. After aligning (rotating) the region with the keypoint's main orientation, the region is divided into a  $4 \times 4$  grid of smaller cells. Within each of these 16 cells, gradient orientations are quantized into an 8-bin histogram covering 360°. These histograms are then concatenated, forming a  $4 \times 4 \times 8 = 128$ -dimensional vector.
3. Applying  $L_2$  normalization to the 128-D vector to obtain a final descriptor that is robust to make them independent of illumination and contrast variations, thus enhancing robustness.

3) *Build Bag of Visual Words*: After extracting SIFT descriptors (the 128-D vectors) for all the training images, all the descriptors are clustered by **kMeans** method to get the cluster center (i.e. visual word) of each cluster. And the dimensionality of the final BoVW image coding vector is determined by the number of target clusters, presenting the expressive capacity of the BoVW.

For each image, the descriptor calls `kmeans.predict()` to assign each descriptor to the nearest visual word; then use `np.histogram` to count the number of descriptors that fall in each cluster, and obtain a one-dimensional histogram.

After that, the histogram is normalized to convert the frequency into a probability distribution to form the final BoVW feature vector. And these BoVW vectors are uti-

lized as feature inputs to a classifier (e.g. SVM) to perform the image classification task (target of this project).

Finally, this project test different classifiers on the BoVW feature vectors, including: K-Nearest Neighbors (**KNN**), Support Vector Machine (**SVM**), and **Softmax** Regression. The performance of these classifiers will be compared with each other to explore the best classifier for the BoVW feature vectors.

### C. CNN

CNN is selected as another approach to achieve object classification, which is more state of the art and widely used in the field of computer vision.

1) *Network Architecture*: CNN Network Architecture is a key factor that affects the performance of CNN, which may influence the accuracy, speed, and robustness of the model. In this project, three different network architectures are tested: **Self-Designed Simple CNN**, **ResNet-18** and **ResNet-50** (with *He initialization* [18]). All of them will be implemented by self-coding, and the performance of them will be compared with each other.

The following Figure 1 is the brief architecture of the self-designed simple CNN model. It contains 2 stage, each stage contains 1 convolutional layers, 1 max pooling layer, and 1 ReLU activation layer. Then, the output of the last stage is flattened and passed to a fully connected layer. Its

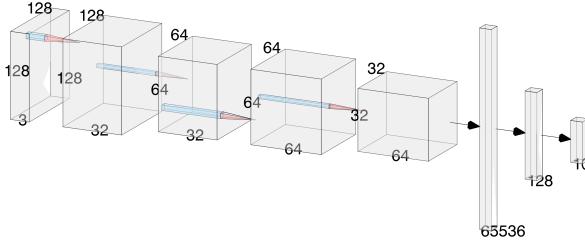


Fig. 1: Architecture of the Self-Designed Simple CNN

character is to play as the baseline model, and compare with the other two ResNet architectures (mainly ResNet-50) to explore the performance of different network.

Then, Table I shows the architecture of ResNet-18 and ResNet-50, which are two popular architectures in the field of computer vision [10]. Compared with the self-designed simple CNNs, ResNets are more complex and deeper, consisting of a “stem” module followed by four sequential stages that progressively reduce spatial resolution while increasing feature depth.

In addition, in order to alleviate the degradation problem that arises in very deep networks, ResNets introduce skip (identity) connections in each residual block. Rather than directly learning an underlying mapping  $H(x)$ , each block learns a residual function  $F(x) = H(x) - x$ ; its output is then  $H(x) = F(x) + x$ . This identity shortcut not only makes it easier for the optimizer to learn small perturbations around the identity (improving convergence), but also provides a direct pathway for both activations

layer name	output size	18-layer	50-layer
conv1	112×112	7×7, 64, stride 2	
		3×3 max pool, stride 2	
conv2_x	56×56	$\left[ \begin{array}{c} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$
		$\left[ \begin{array}{c} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$
conv3_x	28×28	$\left[ \begin{array}{c} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$
		$\left[ \begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$
conv4_x	14×14	$\left[ \begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$
		average pool, 1000-d fc, softmax	
FLOPs		$1.8 \times 10^9$	$3.8 \times 10^9$

TABLE I: Architectures for ResNet (adapted from [10])

and gradients between shallow and deep layers, effectively mitigating vanishing or exploding gradients in very deep architectures.

2) *Coordinate Attention*: In deep learning, attention mechanisms have been widely used to enhance the model’s performance. This project also applies the **Coordinate Attention** (CA) mechanism [19] to the ResNets. CA is a lightweight attention mechanism designed to enhance the CNNs’ ability to model spatial information while maintaining high efficiency. Unlike traditional channel attention mechanisms (e.g. SE [20]) that focus on “channel dimensions”, CA encodes both “spatial location information (coordinates)” and “channel dependencies”. The specific schematic comparison in shown in Figure 2.

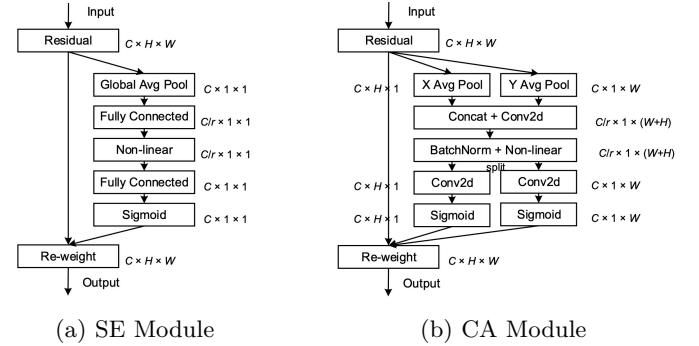


Fig. 2: Schematic comparison of the proposed CA block (b) to the classic SE block (a). (Adapted from [19])

## IV. EXPERIMENTS AND RESULTS

In the experimental section, this project will experiment the performance of each component, present the model accuracy (i.e. *F1-score*) and show the final models with the best results.

### A. Local Feature Methods

This section will cover the experiments of grid search, classifiers, keypoint detection methods, and the final model selection.

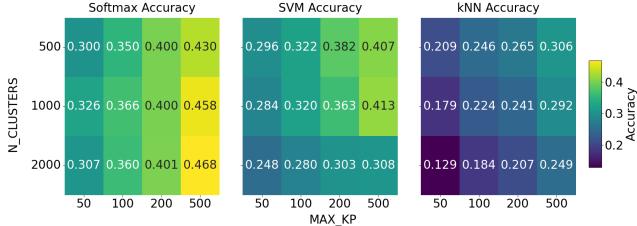


Fig. 3: Grid Search Results of SIFT on Different Classifiers on *CIFAR-10* dataset (From left to right: Softmax, SVM and KNN)

1) *Grid Search*: For Grid Search in the Local Feature section, this experiment mainly tests the following basic parameters:

- **MAX\_KP**: The maximum number of keypoints in each image. The test range is: [50, 100, 200, 500].
- **n\_CLUSTERS**: The number of clusters in the kMeans clustering algorithm, determining the size of the BoVW. The test range is: [500, 1000, 2000].

This grid search is performed on both *CIFAR-10* and *iCubWorld 1.0* datasets, and the best performance is selected as the basic hyperparameters for the local feature methods. In this process, I downsampled both datasets to ensure efficient grid search. A 30% downsampling was performed for the *CIFAR-10* dataset and a 70% downsampling was performed for *iCubWorld 1.0*. In this case, there are 15000 images in *CIFAR-10* and 4200 images in *iCubWorld 1.0* being used for the grid search, which is representative. In addition, 80% of the data will be used as train set, and the rest of the data will be used as validation set to measure the performance.

Figure 3 shows the grid search results of SIFT+BoVW on *CIFAR-10* dataset with different classifiers. The three classifiers show different tendency with the hyperparameters; while, the same phenomenon occurs on the *iCubWorld 1.0* dataset. For the SVM and kNN classifiers, the accuracy of the model grows as *MAX\_KP* grows, but this is not the case for *n\_CLUSTERS*, where the best performance of both classifiers occurring at *n\_CLUSTERS*=1000, and larger or smaller *n\_CLUSTERS* both lead to a decrease in model performance. As for the Softmax classifier, when *MAX\_KP* is small (e.g., 50 and 100), it also conforms to the above law; and when *MAX\_KP* is large, the larger the *n\_CLUSTERS* is, the more accurate the model is.

The possible reasons for this phenomenon could be the sensitivity of the three classifiers to the number of keypoints and BoVW size differs. First, as *MAX\_KP* increases, the number of SIFT descriptors increases significantly, and the statistics of each visual word in the BoVW histogram is more stable, which enables both kNN based

on the distance metric and SVM based on the maximal spacing to obtain more accurate and robust classification boundaries. Second, for the *n\_CLUSTERS*, both SVM and kNN degrade in performance when the dictionary is too small (underfitting) or too large (high-dimensional sparse/overfitting), but they reach a performance balance at a medium size (1000 words). Finally, Softmax classifier is with L2 regularization, which may help the model to suppress noisy words to take full advantage of the finer visual word segmentation. Therefore, for the comprehensive performance consideration of multiple classifiers, [*MAX\_KP*=500, *n\_CLUSTERS*=1000] will be chosen as the benchmark hyperparameters for Local Feature in the subsequent experiments.

2) *Classifiers*: Figure 3 also shows the performance difference between the three classifiers. Overall, the softmax classifier achieves the best performance, while the kNN performs the worst. And the SVM classifier is in the middle. In addition, the three classifiers differ significantly in terms of training time, with kNN being the most efficient on the same device (taking only 0.4s), softmax next (taking 1.8s for 10 epochs), and SVM (using LinearSVC) being the slowest (taking 11 mins). Therefore, for the combined consideration of model performance and training efficiency, Softmax was used as the default classifier in all subsequent tests in this project.

3) *Keypoints Detection*: In the SIFT pipeline, I tested three different keypoint detection methods: Difference of Gaussian (DoG), Harris-Laplace, and Harris. The results of these methods will be compared with each other to explore the best keypoint detection method for the SIFT pipeline.

Dataset	Harris	DoG	Harris-Laplace
<i>CIFAR-10</i>	45.4	40.2	46.5
<i>CATE</i>	23.8	21.6	25.7
<i>DEMO</i>	22.6	18.9	24.4

TABLE II: Accuracy (%) with Different Keypoint Detection Methods

In Table II, the experimental results demonstrate some differences in the performance of the three keypoint detection methods on different datasets. Harris-Laplace has the highest accuracy (achieves 46.5% on *CIFAR-10*, and 25.7% on *CATE*), followed by Harris (achieves 45.4% on *CIFAR-10*, and 23.8% on *CATE*), while DoG is the worst (achieves 40.2% on *CIFAR-10*, and 21.6% on *CATE*).

This performance is not expected. The performance ranking of the three was envisioned to be Harris < DoG < Harris-Laplace, but in fact Harris outperforms DoG on all test sets and closes to Harris-Laplace. Presumably, the possible reasons are, firstly, that the image resolution of *CIFAR-10* is low, and thus there is limited spatial difference between different scales of blur images, making it difficult for DoG to find a truly stable blob structure;

and secondly, the background of the shooting environment in the *iCubWorld 1.0* dataset is relatively simple and consistent, and the surface of the object seldom shows obvious blob structures, but instead, it is dominated by edges and corner points. This makes Harris more able to obtain more stable and recognizable keypoints on both datasets, while DoG and Harris-Laplace have difficulties in obtaining better performance because of the difficulty in obtaining blobs.

Finally, Harris-Laplace was chosen as the keypoint detection method for this project in terms of overall accuracy. Therefore, the final SIFT+BoVW model is [ $MAX\_KP=500$ ,  $n\_CLUSTERS=1000$ , Softmax classifier, Harris-Laplace], whose F1-score can achieve 46.5% on *CIFAR-10* and 25.7% on *CATE*, and 24.4% on *DEMO* (has been shown in Table II).

## B. CNN

1) *Grid Search*: For Grid Search in the Local Feature section, this experiment mainly tests the following basic parameters:

- **Learning Rate**: The test range is: [1e-3, 5e-4, 1e-4].
- **Batch Size**: The test range is: [32, 64, 128].
- **Epoch Number**: The test range is: [10, 20, 30].

Similarly, the two datasets used in the grid search were downsampled with the same downsampling strategy as the Local Feature grid search.

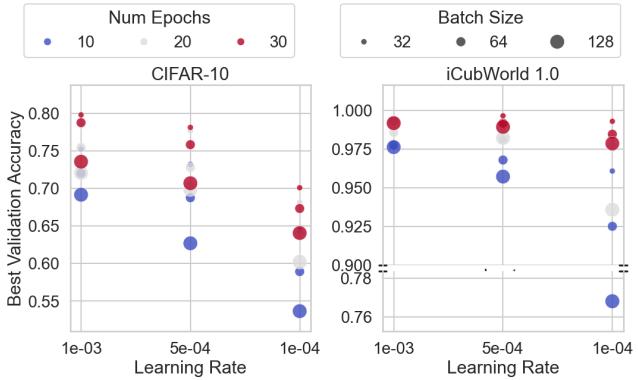


Fig. 4: Grid Search Results of CNN on *CIFAR-10* and *iCubWorld 1.0*

Figure 4 shows the results of the grid search on the CNN. In the *CIFAR-10*, the accuracy distribution is more regular basically following the law that the smaller the batch size, the larger the learning rate, the larger the epoch number, and ultimately the larger the accuracy. Therefore, [batch size=32, learning rate=1e-3, epoch number=30] is chosen as the basic hyperparameter combination for the subsequent experiments.

And on the *iCubWorld 1.0*, all parameter combinations achieve a high level of fit, except for the case of epoch=10 where underfitting may occur. At a learning rate of 1e-3, all combinations achieve 97.5% accuracy. Also, it can

be noticed that the batch size does not have a significant effect on this dataset (about 1%). Therefore, due to the concern of overfitting and comprehensive performance, [batch size=128, learning rate=5e-4, epoch number=30] is chosen as the basic hyperparameter combination on the *iCubWorld 1.0* dataset.

2) *Data Pre-Processing*: This project applies various transformations to the original training images to synthesize training samples with higher diversity, thus expanding the data distribution and reducing the risk of model overfitting.

This experiment tests three different data augmentation strategies, each new strategy is based on the previous one for further processing as follows:

- **Strategy 1 - No Augmentation**: No data augmentation is applied to the original images.
- **Strategy 2 - Rotation and Cropping**: The images are rotated by 10 degrees and centeral cropped to 224x224 pixels.
- **Strategy 3 - ColorJitter and RandomPerspective**: The images are randomly adjusted in brightness, contrast, saturation, and hue, and the perspective is randomly changed.

Model Performance Under Different Augmentation Strategies

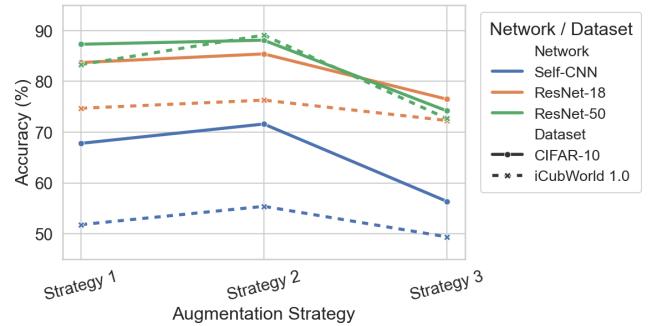


Fig. 5: Accuracy (%) with Different Augmentation Strategies

Figure 5 shows the performance of the model with different data augmentation strategies on both datasets. By comparing the accuracy line, it can be found that the model with data augmentation strategy 2 achieves the best performance (average increase of about 4.1% compared to Strategy 1); on the other hand, Strategy 3 shows a negative impact in all conditions, with an average decrease in accuracy of 10.2% compared to Strategy 1. Therefore, the data augmentation strategy 2 is selected as the final data augmentation strategy for this project.

3) *Model Architecture*: As mentioned before, this project tests three different network architectures: *self-designed simple CNN*, ResNet-18, and ResNet-50. The following Table III shows their performance on the both datasets.

It can be seen that on CIAFR-10, the performance of all models increases from 71% to 88% with model com-

Dataset	Self-CNN	ResNet-18	ResNet-50
CIFAR-10	71.6	85.4	88.1
CATE	56.3	75.4	75.7
DEMO	55.4	76.3	89.1

**TABLE III:** Accuracy (%) of Models on *CIFAR-10* and *iCubWorld 1.0* test sets

plexity, reflecting the positive correlation between model complexity and test set accuracy under the current hyperparameters. This phenomenon is even more pronounced for *iCubWorld 1.0*, where ResNet-18 generally improves accuracy by 10% compared to Self-CNN, but ResNet-50 does not improve as much compared to ResNet-18.

4) *Attention Mechanism*: Attention mechanisms are used as an important component of present-band machine learning, and this implementation conducts comparative experiments with several of the previously mentioned mechanisms. In Table IV, I present part of the experiment results about the attention mechanism on ResNet-50. The results show that the attention mechanism can significantly improve the performance of the model. And in the training processing, it can be seen that the model with the attention mechanism tends to have a better fitting efficiency in the early stages of training, especially when the training with smaller learning rate. And compared with the SE attention mechanisms, CA can achieve a better performance on both datasets. For example, on *CIFAR-10*, the accuracy of the model with CA is 88.1%, while the accuracy of the model with SE are 88.0%, respectively. And on *iCubWorld 1.0*, the accuracy of the model with CA is 75.7% and 89.1%, while the accuracy of the model with SE are 70.9% and 84.7%, respectively.

Dataset	Original	SE	CA
CIFAR-10	85.7	88.0	88.1
CATE	64.4	70.9	75.7
DEMO	81.0	84.7	89.1

**TABLE IV:** Accuracy (%) with Different Attention Mechanisms

5) *Optimizer*: In the training process of CNN, the choice of optimizer is also a key factor that affects the performance of the model. This project mainly tests two different optimizers: *AdamW* [21], and *Ranger* [22]. By comparing the performance of the two optimizers on individual models and training sets, it was found that the final accuracy between the two was similar (typically between 1 and 2%), but the Ranger optimizer had a better smoother fit, while the AdamW optimizer's fitting process had larger fluctuations.

Therefore, the Ranger optimizer was chosen as the default option to ensure a smoother training process and facilitate the tuning of other parameters on both datasets, and ResNet-50 with CA module is selected as the final

model with Strategy 2 data augmentation. As shown in Table IV, on *CIFAR-10*, the basic hyperparameters are [batch size=32, learning rate=1e-3, epoch number=30], and the final F1-score is 88.1%. On *iCubWorld 1.0*, the basic hyperparameters are [batch size=128, learning rate=5e-4, epoch number=30], and the final F1-score is 75.7% on *CATE* and 89.1% on *DEMO*.

## V. DISCUSSION

This section will discuss the difference between the two methods from different perspectives, including the model performance, training time, and the interpretability of the model.

Overall, in terms of performance, the accuracy of the ResNet-50 model (with CA module) is significantly better than the Local Feature (SIFT+BoVW) method on either test set. This proves that the CNN model can still show good generalization performance on the image classification task, even with a smaller data volume and a test set that requires higher generalization ability of the model, such as *iCubWorld 1.0*. As for the SIFT method, its accuracy is backward, and it can only achieve an accuracy of 46.5% even on a simple dataset like *CIFAR-10* with sufficient sample size, which illustrates the limitation of the SIFT method on the image classification task.

However, the performance comparison still reveals the essential difference in generalization performance between the two methods due to the different underlying logic. On the *iCubWorld 1.0* dataset, the accuracy of SIFT+BoVW on *CATE* is higher than that on *DEMO*, while the accuracy of ResNet-50 is higher on *DEMO* than on *CATE*, which suggests that the SIFT+BoVW method relies more on low-level features such as color and texture of the image for feature extraction, and it can classify the same kind of different images in the *CATE* test set by learning the “local structure” of the image, and the model performs better when facing the same category of different instances in the *CATE* test set; however, in the *DEMO* test set, the performance of the model is degraded due to the drastic change of the background information by the change of the demonstrator. ResNet-50, on the other hand, achieves object classification by learning abstract semantic features of the image, which makes it more robust to changing backgrounds; however, since the CNN model is a Data-Driven approach, it is more dependent on sample size, which causes it to perform more poorly when confronted with small sample sizes of unseen examples in the *CATE* test set.

There is also a significant difference in the computational cost between the two methods. the computational time of the SIFT+BoVW method is mainly focused on feature extraction and BoVW dictionary construction, and the number of images that can be processed per second ranges from about 5it/s to 30it/s (using CPU only) after limiting the number of available keypoints. In contrast, the computational time of ResNet-50 (with a reference number

of about  $2.55 \times 10^6$ ) is mainly focused on model training, with its training time on the two datasets being 1.5 hours and 1 hour respectively (under the condition of using GPU acceleration). This reflects the huge difference in computation between the two and explains why the Local Feature approach has always been active on edge devices and has not been replaced by CNNs. Also the interpretability of the SIFT+BoVW method is significantly better than that of ResNet-50. Since the SIFT+BoVW method classifies images based on their local features, its classification results can be explained by visualizing the features in the BoVW dictionary; whereas, the training process of ResNet-50 is a black-box process, which makes it difficult to explain the reason. Therefore, the SIFT+BoVW method is easier to analyze and debug, while ResNet-50 can only rely on the loss function during the training process to evaluate the performance of the model.

## VI. LIMITATIONS AND FUTURE WORK

Although the project has achieved good results, there are still some limitations and much work remains to be done in the future. For example, in terms of hyperparameter tuning, this project only performed grid searches for a few basic hyperparameters, which may have caused the model to miss many potentially more available hyperparameter combinations. In this case, more automatical tuning methods such as Optuna can be used to find the optimal hyperparameters with more detailed search space in the further work.

Moreover, although this project involves a large number of models and methods, there is a lack of innovation for them, rather than a preference for reproduction and application. To solve this, the techniques mentioned in the Section II, e.g. Vision Transformer, can be applied to further improve the performance of the model.

## VII. CONCLUSION

To summarize, this project has successfully implemented two different image classification methods: Local Feature (SIFT+BoVW) and CNN (ResNet-50). The experimental results show that the CNN model outperforms the Local Feature method in terms of accuracy and generalization performance. The project also explores the impact of various hyperparameters, data augmentation strategies, and model architectures on the performance of both methods. The practice and discovery of this project has deepened my understanding and inspired my interest in the field, and I look forward to further continuous improvement of the program and further study of the field afterwards.

## REFERENCES

- [1] A. Krizhevsky, G. Hinton *et al.*, “Learning multiple layers of features from tiny images.(2009),” 2009.
- [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] D. G. Lowe, “Object recognition from local scale-invariant features,” in *Proceedings of the seventh IEEE international conference on computer vision*, vol. 2. Ieee, 1999, pp. 1150–1157.
- [4] ———, “Distinctive image features from scale-invariant keypoints,” *International journal of computer vision*, vol. 60, pp. 91–110, 2004.
- [5] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, “Speeded-up robust features (surf),” *Computer vision and image understanding*, vol. 110, no. 3, pp. 346–359, 2008.
- [6] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, “Orb: An efficient alternative to sift or surf,” in *2011 International conference on computer vision*. Ieee, 2011, pp. 2564–2571.
- [7] T. Tuytelaars, K. Mikolajczyk *et al.*, “Local invariant feature detectors: a survey,” *Foundations and trends® in computer graphics and vision*, vol. 3, no. 3, pp. 177–280, 2008.
- [8] N. Robinson, B. Tidd, D. Campbell, D. Kulić, and P. Corke, “Robotic vision for human-robot interaction and collaboration: A survey and systematic review,” *ACM Transactions on Human-Robot Interaction*, vol. 12, no. 1, pp. 1–66, 2023.
- [9] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [10] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” 2015. [Online]. Available: <https://arxiv.org/abs/1512.03385>
- [11] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [12] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah, “Transformers in vision: A survey,” *ACM computing surveys (CSUR)*, vol. 54, no. 10s, pp. 1–41, 2022.
- [13] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [14] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 10 012–10 022.
- [15] J. Zhu, X. Chen, K. He, Y. LeCun, and Z. Liu, “Transformers without normalization,” *arXiv preprint arXiv:2503.10622*, 2025.
- [16] Q. Li, Y. W. Teh, and R. Pascanu, “Noprop: Training neural networks without back-propagation or forward-propagation,” *arXiv preprint arXiv:2503.24322*, 2025.
- [17] K. Sargent, K. Hsu, J. Johnson, L. Fei-Fei, and J. Wu, “Flow to the mode: Mode-seeking diffusion autoencoders for state-of-the-art image tokenization,” *arXiv preprint arXiv:2503.11056*, 2025.
- [18] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” 2015. [Online]. Available: <https://arxiv.org/abs/1502.01852>
- [19] Q. Hou, D. Zhou, and J. Feng, “Coordinate attention for efficient mobile network design,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 13 713–13 722.
- [20] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [21] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [22] L. Wright, “Ranger - a synergistic optimizer.” <https://github.com/lessw2020/Ranger-Deep-Learning-Optimizer>, 2019.

## APPENDIX

```

# --- 1. Interest Point Detection ---
def detect_interest_points(gray, method='DoG', threshold=5):
    if method == 'DoG':
        # Define multiple sigma values for Gaussian blur
        num_octaves=3 # the level of the image pyramid
        scales_per_octave=3 # the number of scales per octave
        sigma0=1.6 # otiginal sigma value

        keypoints = []
        k = 2 ** (1 / scales_per_octave) # the scale factor for each scale

        # form the 0th octave to num_octaves-1 octaves
        for octave in range(num_octaves):
            # create the base image for the current octave
            if octave == 0:
                base = gray.copy()
            else:
                # downsample the previous base image
                base = cv2.pyrDown(prev_base)

            # save the current base for the next octave
            prev_base = base
            # calculate the size of the base image
            h, w = base.shape

            # create the Gaussian pyramid
            # scales_per_octave + 3 scales, to produce scales_per_octave+2 DoG
            # layers
            sigmas = [sigma0 * (k ** (i + octave*scales_per_octave)) for i in
                      range(scales_per_octave + 3)]
            gaussians = [cv2.GaussianBlur(base, (0,0), sigmaX=s) for s in
                         sigmas]
            # calculate DoG
            dogs = []
            for i in range(len(gaussians)-1):
                diff = gaussians[i+1].astype(np.float32) - gaussians[i].astype
                (np.float32)
                dogs.append(diff)

            # find the scale-space extrema
            for i in range(1, len(dogs)-1):
                # the three layers to compare
                prev_d, curr_d, next_d = dogs[i-1], dogs[i], dogs[i+1]

                # dilation is used to find the maximum
                max_spat = cv2.dilate(curr_d, np.ones((3,3),np.uint8))
                # erosion is used to find the minimum
                min_spat = cv2.erode( curr_d, np.ones((3,3),np.uint8))


```

```

        mask_max = (curr_d == max_spat) & (curr_d > prev_d) & (curr_d
        > next_d) & (np.abs(curr_d) > threshold)
        mask_min = (curr_d == min_spat) & (curr_d < prev_d) & (curr_d
        < next_d) & (np.abs(curr_d) > threshold)

        # Find all coordinates of the points that satisfy the
        # conditions
        coords = np.vstack((np.argwhere(mask_max), np.argwhere
        (mask_min)))

        # map the coord back to the original image
        for y, x in coords:
            # make downsampling back
            scale = (2 ** octave)
            kp = cv2.KeyPoint(float(x*scale), float(y*scale),
            size=PATCH_SIZE * scale)
            keypoints.append(kp)

    keypoints = sorted(keypoints, key=lambda kp: kp.response if hasattr
    (kp, 'response') else 1.0, reverse=True)
    return keypoints[:MAX_KP]

elif method == 'Harris':
    # calculate Harris corner response
    H = cv2.cornerHarris(
        gray.astype(np.float32),
        blockSize=2,
        ksize=3,
        k=0.04
    )

    # normalize the response to [0, 1]
    Hn = cv2.normalize(H, None, 0, 1, cv2.NORM_MINMAX)

    # get local maximum
    dil = cv2.dilate(Hn, np.ones((3,3), np.uint8))
    nms_mask = (Hn == dil)

    # thresholding
    mask = nms_mask & (Hn > 0.01*threshold)

    # get kps
    ys, xs = np.where(mask)
    keypoints = []
    for y, x in zip(ys, xs):
        kp = cv2.KeyPoint(
            x=float(x), y=float(y),
            size=float(PATCH_SIZE),
            response=float(Hn[y, x])
        )
        keypoints.append(kp)

    # order keypoints by response
    keypoints.sort(key=lambda kp: kp.response, reverse=True)
    return keypoints[:MAX_KP]

```

```

# --- Interest Point Detector: Harris-Laplace ---
def detect_interest_points(gray,
                           method=INTEREST_METHOD,
                           num_octaves=3,
                           scales_per_octave=3,
                           sigma0=1.6,
                           harris_block=2,
                           harris_ksize=3,
                           harris_k=0.04,
                           harris_thresh=0.01):
    keypoints = []
    k = 2 ** (1 / scales_per_octave)
    # 1) Build Gaussian pyramid bases
    gp = []
    for octave in range(num_octaves):
        base = gray if octave == 0 else cv2.pyrDown(gp[-1])
        gp.append(base)

    # 2) Build Gaussian blurred images per octave
    gauss_pyr = [] # list of lists: gauss_pyr[octave][scale]
    for o, base in enumerate(gp):
        sigmas = [sigma0 * (k ** (i + o*scales_per_octave))
                  for i in range(scales_per_octave + 3)]
        gaussians = [cv2.GaussianBlur(base, (0,0), sigmaX=s) for s in
                     sigmas]
        gauss_pyr.append(gaussians)

    # 3) Harris detection per scale (spatial localization)
    harris_pyr = [] # harris_pyr[octave][scale] = boolean mask
    for o, gaussians in enumerate(gauss_pyr):
        masks = []
        for img in gaussians:
            H = cv2.cornerHarris(np.float32(img), blockSize=harris_block,
                                 ksize=harris_ksize, k=harris_k)
            Hn = cv2.normalize(H, None, 0, 1, cv2.NORM_MINMAX)
            # spatial Non-Max Suppression
            max_spat = cv2.dilate(Hn, np.ones((3,3), np.uint8))
            mask = (Hn == max_spat) & (Hn > harris_thresh)
            masks.append(mask)
        harris_pyr.append(masks)

```

```

# 4) Compute Laplacian-of-Gaussian responses per scale for scale
selection
log_pyr = [] # log_pyr[octave][scale]
for gaussians in gauss_pyr:
    logs = [cv2.Laplacian(img, cv2.CV_32F, ksize=3) for img in
            gaussians]
    log_pyr.append(logs)

```

```

# 5) Combine: for each Harris candidate, pick scale with max |LoG|
seen = set()
for o, masks in enumerate(harris_pyr):
    for s_idx, mask in enumerate(masks):
        ys, xs = np.where(mask)
        for y, x in zip(ys, xs):
            # scale selection
            responses = [abs(log_pyr[o][i][y, x]) for i in range(len(log_pyr[o]))]
            best_i = int(np.argmax(responses))
            # map to original image coordinates
            scale_factor = 2 ** o
            X = int(x * scale_factor)
            Y = int(y * scale_factor)
            if (X, Y) in seen:
                continue
            seen.add((X, Y))
            # size proportional to sigma
            best_sigma = sigma0 * (k ** (best_i + o*scales_per_octave))
            size = PATCH_SIZE * best_sigma
            kp = cv2.KeyPoint(float(X), float(Y), size)
            kp.response = responses[best_i]
            keypoints.append(kp)

# 6) sort by response and limit
keypoints.sort(key=lambda kp: kp.response, reverse=True)
return keypoints[:MAX_KP]

```

The above shows the DoG, Harris, Harris-Laplace module details.

```

# --- 2. Compute Custom SIFT-like Descriptor ---
def compute_custom_sift_descriptor(gray, keypoints, patch_size=PATCH_SIZE):
    # Compute gradients once
    grad_x = cv2.Sobel(gray, cv2.CV_32F, 1, 0, ksize=3)
    grad_y = cv2.Sobel(gray, cv2.CV_32F, 0, 1, ksize=3)
    mag = cv2.magnitude(grad_x, grad_y)
    ori = cv2.phase(grad_x, grad_y, angleInDegrees=True)

    descriptors = []
    half = patch_size // 2
    h, w = gray.shape

    for kp in keypoints:
        x, y = int(kp.pt[0]), int(kp.pt[1])
        # boundary check
        if x-half < 0 or y-half < 0 or x+half >= w or y+half >= h:
            continue
        # extract full patch
        mag_full = mag[y-half:y+half, x-half:x+half]
        ori_full = ori[y-half:y+half, x-half:x+half]
        # 1) Orientation assignment: build 8-bin histogram over full patch
        hist8, bin_edges = np.histogram(ori_full.ravel(), bins=8, range=(0, 360), weights=mag_full.ravel())

```

```

# main orientation is center of max bin
max_idx = np.argmax(hist8)
main_ori = (bin_edges[max_idx] + bin_edges[max_idx+1]) / 2.0
# normalize full patch orientations relative to main orientation
ori_full = (ori_full - main_ori + 360) % 360
# slice into subregions and build descriptor
desc = []
for i in range(4):
    for j in range(4):
        sub_mag = mag_full[i*4:(i+1)*4, j*4:(j+1)*4].ravel()
        sub_ori = ori_full[i*4:(i+1)*4, j*4:(j+1)*4].ravel()
        # 2) histogram on aligned orientations
        hsub, _ = np.histogram(sub_ori, bins=8, range=(0,360),
                               weights=sub_mag)
        desc.append(hsub)
desc = np.array(desc, dtype=np.float32)
desc /= (np.linalg.norm(desc) + 1e-7)
descriptors.append(desc)

return np.array(descriptors) if descriptors else None

```

The above shows how SIFT works

```

all_desc = []
img_desc_idx = []
train_labels = []
print('Extracting training descriptors...')
for img, label in tqdm(trainset, desc='Train images'):
    img_np = (img.numpy().transpose(1,2,0) * 255).astype(np.uint8)
    gray = cv2.cvtColor(img_np, cv2.COLOR_RGB2GRAY)
    kps = detect_interest_points(gray, method=INTEREST_METHOD)
    des = compute_custom_sift_descriptor(gray, kps)
    if des is not None:
        all_desc.append(des)
        img_desc_idx.append(len(all_desc))
        train_labels.append(label)
all_desc = np.vstack(all_desc)

```

The above shows how to get descriptor from train images

```

# Build visual vocabulary
kmeans = MiniBatchKMeans(n_clusters=N_CLUSTERS, batch_size=1000, verbose=1)
kmeans.fit(all_desc)

```

The above shows how use K-Means to find visual words

```

# Build BoVW features for train
def build_bovw(descriptors):
    words = kmeans.predict(descriptors)
    hist, _ = np.histogram(words, bins=np.arange(N_CLUSTERS+1))
    return hist.astype(float) / (hist.sum() + 1e-7)

train_feats = []
start = 0
# for idx in img_desc_idx:
for idx in tqdm(img_desc_idx, desc='Building BoVW features'):
    end = idx
    if end - start > 0:
        train_feats.append(build_bovw(all_desc[start:end]))
    else:
        train_feats.append(np.zeros(N_CLUSTERS, dtype=float))
    start = end
train_feats = np.array(train_feats)
train_labels = np.array(train_labels)

```

The above shows how to build BoVW

```

# Extract BoVW features for test set
test_feats = []
test_labels = []
print('Extracting test BoVW features...')
for img, label in tqdm(testset, desc='Test images'):
    img_np = (img.numpy().transpose(1,2,0) * 255).astype(np.uint8)
    gray = cv2.cvtColor(img_np, cv2.COLOR_RGB2GRAY)
    kps = detect_interest_points(gray, method=INTEREST_METHOD)
    des = compute_custom_sift_descriptor(gray, kps)
    if des is not None:
        hist = build_bovw(des)
    else:
        hist = np.zeros(N_CLUSTERS, dtype=float)
    test_feats.append(hist)
    test_labels.append(label)

test_feats = np.array(test_feats)
test_labels = np.array(test_labels)

# SVM Pipeline
pipe = Pipeline([
    ('scaler', StandardScaler()),
    ('svc', LinearSVC(dual=False, max_iter=500, random_state=42,
                      verbose=2))
])

# TODO: further grid search
param_grid = {
    'svc__C': [1],
    'svc__tol': [1e-4]
}

```

```

# K-Fold val
grid = GridSearchCV(pipe,param_grid, cv=5, n_jobs=-1, verbose=2,
scoring='accuracy')

grid.fit(train_feats, train_labels)
print("Best parameters:", grid.best_params_)
print("Best CV accuracy:", grid.best_score_)

# test
best_model = grid.best_estimator_
test_acc = best_model.score(test_feats, test_labels)
print(f"Test accuracy with best SVM: {test_acc:.4f}")

```

```

clf_knn = KNeighborsClassifier(n_neighbors=500, metric='euclidean',
n_jobs=-1)
clf_knn.fit(train_feats, train_labels)
preds_knn = clf_knn.predict(test_feats)
acc_knn = accuracy_score(test_labels, preds_knn)
print(f"KNN Test Accuracy: {acc_knn:.4f}")

```

```

# get train and val
full_ds = TensorDataset(
    torch.from_numpy(train_feats).float(),
    torch.from_numpy(train_labels).long()
)
train_n = int(0.9 * len(full_ds))
val_n   = len(full_ds) - train_n
train_ds, val_ds = random_split(full_ds, [train_n, val_n])

# DataLoader
batch_size = 256
train_loader = DataLoader(train_ds, batch_size=batch_size, shuffle=True)
val_loader   = DataLoader(val_ds,   batch_size=batch_size)

# model, loss, optimizer
model = nn.Sequential(
    nn.BatchNorm1d(N_CLUSTERS),
    nn.Linear(N_CLUSTERS, 10)
).to(device)
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=LR, weight_decay=1e-4)

# train
best_val_acc = 0.0
best_path     = 'best_model.pt'

```

The above shows how to get descriptor for test images, and the cores of LinearSVC, KNN,

Softmax classifiers.

```
# Self-CNN Architecture
class SimpleCNN(nn.Module):
    def __init__(self, num_classes=10):
        super(SimpleCNN, self).__init__()
        self.features = nn.Sequential(
            # con, relu, maxpool, repeat twice
            nn.Conv2d(3, 32, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(2),
            nn.Conv2d(32, 64, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(2)
        )
        self.classifier = nn.Sequential(
            nn.Flatten(),
            nn.Linear(64 * 32 * 32, 128),
            nn.ReLU(inplace=True),
            nn.Dropout(0.5),
            nn.Linear(128, num_classes)
    )

    def forward(self, x):
        x = self.features(x)
        return self.classifier(x)
```

The above shows the Self-CNN architecture.

```
class Bottleneck(nn.Module):
    expansion = 4
    def __init__(self, in_planes, planes, stride=1, downsample=None,
                 ca_reduction=32):
        super().__init__()
        # 1x1 reduce dim
        self.conv1 = nn.Conv2d(in_planes, planes, kernel_size=1, bias=False)
        self.bn1 = nn.BatchNorm2d(planes)
        # 3x3 conv
        self.conv2 = nn.Conv2d(planes, planes, kernel_size=3,
                           stride=stride, padding=1, bias=False)
        self.bn2 = nn.BatchNorm2d(planes)
        # 1x1 restore dim
        self.conv3 = nn.Conv2d(planes, planes * self.expansion,
                           kernel_size=1, bias=False)
        self.bn3 = nn.BatchNorm2d(planes * self.expansion)
```

```

# CoordAtt
self.ca    = CoordAtt(planes * self.expansion,
                      planes * self.expansion,
                      reduction=ca_reduction)
# self.se = SELayer(planes * self.expansion, reduction=16)

self.relu = nn.ReLU(inplace=True)
self.downsample = downsample

def forward(self, x):
    identity = x

    out = self.relu(self.bn1(self.conv1(x)))
    out = self.relu(self.bn2(self.conv2(out)))
    out = self.bn3(self.conv3(out))

    # apple Coordinate Attention here
    out = self.ca(out)
    # out = self.se(out)

    if self.downsample is not None:
        identity = self.downsample(x)
    out += identity
    return self.relu(out)

```

```

class SimpleResNetLike(nn.Module):
    def __init__(self, num_classes=10):
        super().__init__()

        # stem of resnet50
        self.stem = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=7, stride=2, padding=3,
                     bias=False),
            nn.BatchNorm2d(64),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2, padding=1),
        )
        # four stages of resnet50
        self.layer1 = self.make_layer(64, 64, blocks=3, stride=1) # 256
        self.layer2 = self.make_layer(256, 128, blocks=4, stride=2) # 512
        self.layer3 = self.make_layer(512, 256, blocks=6, stride=2) # 1024
        self.layer4 = self.make_layer(1024, 512, blocks=3, stride=2) # 2048

        # avgpool and fc
        self.avgpool = nn.AdaptiveAvgPool2d((1,1))
        self.fc = nn.Linear(512 * Bottleneck.expansion, num_classes)

        self._init_weights()

```

```

def make_layer(self, in_planes, planes, blocks, stride=1):
    downsample = None
    out_planes = planes * Bottleneck.expansion
    if stride != 1 or in_planes != out_planes:
        # 1x1 conv to downsample
        downsample = nn.Sequential(
            nn.Conv2d(in_planes, out_planes,
                      kernel_size=1, stride=stride, bias=False),
            nn.BatchNorm2d(out_planes),
        )
    layers = [Bottleneck(in_planes, planes, stride, downsample)]
    for _ in range(1, blocks):
        layers.append(Bottleneck(out_planes, planes))
    return nn.Sequential(*layers)

```

```

def _init_weights(self):
    for m in self.modules():
        if isinstance(m, nn.Conv2d):
            # He init
            nn.init.kaiming_normal_(m.weight, mode='fan_out',
                                   nonlinearity='relu')
            if m.bias is not None:
                nn.init.zeros_(m.bias)
        elif isinstance(m, nn.Linear):
            nn.init.kaiming_uniform_(m.weight, mode='fan_in',
                                   nonlinearity='linear')
            if m.bias is not None:
                nn.init.zeros_(m.bias)
        elif isinstance(m, nn.BatchNorm2d):
            nn.init.ones_(m.weight)
            nn.init.zeros_(m.bias)

def forward(self, x):
    x = self.stem(x)
    x = self.layer1(x)
    x = self.layer2(x)
    x = self.layer3(x)
    x = self.layer4(x)
    x = self.avgpool(x)
    x = torch.flatten(x, 1)
    return self.fc(x)

```

The above shows the ResNet-50 with CA module and He init.  
The below shows the several data augmentation strategies in CNN part.

```
def get_data_augmentation(train_indices, val_indices):
    # data augmentation
    transform_train = transforms.Compose([
        # Resize
        transforms.Resize(128, interpolation=transforms.InterpolationMode.BICUBIC),
        # Random crop and horizontal flip
        transforms.RandomCrop(128, padding=4),
        transforms.RandomHorizontalFlip(),

        transforms.ToTensor(),
        transforms.Normalize((0.4914, 0.4822, 0.4465), (0.2470, 0.2430, 0.2610)) # norm the train data
    ])
    transform_test = transforms.Compose([
        # Only apple resize
        transforms.Resize(128, interpolation=transforms.InterpolationMode.BICUBIC),

        transforms.ToTensor(),
        transforms.Normalize((0.4914, 0.4822, 0.4465), (0.2470, 0.2430, 0.2610)) # norm the test data
    ])

def get_data_augmentation(train_indices, val_indices):
    transform_train = transforms.Compose([
        # resize
        transforms.Resize(256, interpolation=transforms.InterpolationMode.BICUBIC),
        #random crop and flip
        transforms.RandomResizedCrop(224, scale=(0.8, 1.0)),
        transforms.RandomHorizontalFlip(),
        # color jitter, perspective change and affine
        transforms.ColorJitter(brightness=0.2, contrast=0.2),
        transforms.RandomAffine(degrees=10, translate=(0.1, 0.1), scale=(0.9, 1.1)),
        transforms.RandomPerspective(distortion_scale=0.1, p=0.25),

        transforms.ToTensor(),
        # transforms.RandomErasing(p=0.5, scale=(0.02, 0.2), ratio=(0.3, 3.3)),
        transforms.Normalize(mean, std),
    ])

    transform_test = transforms.Compose([
        # only resize and norm
        transforms.Resize(256, interpolation=transforms.InterpolationMode.BICUBIC),
        transforms.CenterCrop(224),
        transforms.ToTensor(),
        transforms.Normalize(mean, std),
    ])
```

If you want to find more details about the code, you can visit my GitHub repository: [https://github.com/SchwarzerYV/CRCV\\_Assignment.git](https://github.com/SchwarzerYV/CRCV_Assignment.git).