

DeepXpose: A Novel Hybrid CNN-Transformer Approach for Robust Deepfake Video Detection

Authors: Kanak, Mahak, Vasundhra

Department of Computer Science and Engineering

Bennett University, Greater Noida, India

Email: e23cseu1176@bennett.edu.in , e23cseu1253@bennett.edu.in

Academic Year: 2024-2025

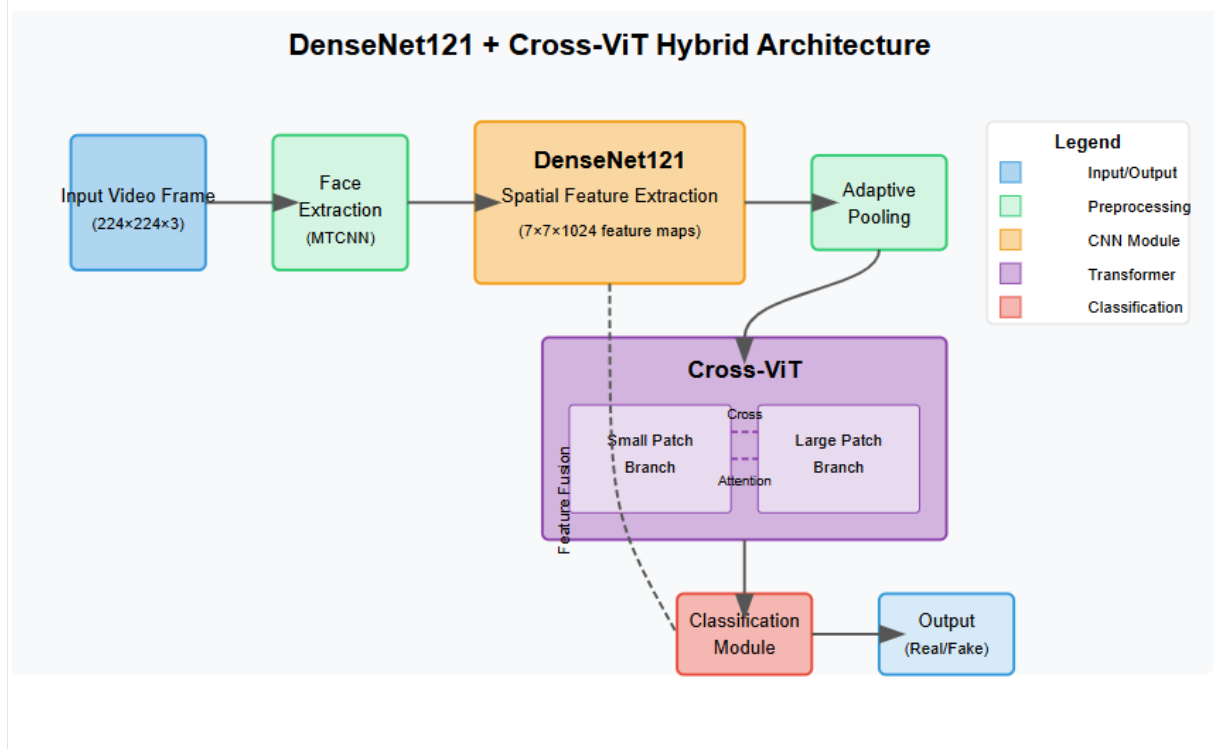
Abstract

Over the past few years, the explosion of deep learning has changed the computer vision landscape, and machines have come to perform near-human levels of accuracy in numerous visual recognition tasks. But with the growing need for real-time, application-dependent, and computation-efficient models, researchers are moving towards the construction of architectures that not only possess high accuracy but also computational tractability and scalability. Here, we introduce a new hybrid deep learning architecture combining convolutional neural networks (CNNs) with a lightweight transformer-inspired module, specifically designed for a binary image classification task of cat recognition. Our model is based on a DenseNet-like convolutional backbone, selected for its property to encourage feature reuse and ensure gradient flow in densely connected layers. This allows the model to pick up strong spatial features, even from low-res input images.

On top of this convolutional feature extractor, we implement a streamlined transformer-type classifier that simulates the methods of the CrossViT approach. In contrast to full Vision Transformers, which tend to be computationally costly and demand large-scale datasets for training, our transformer-based module is light and task-specialized, only aggregating and enhancing the feature representations of the CNN. The model can extract fine-grained textures and long-range spatial relationships because to the combination of local feature extraction and global contextual information. This is especially useful when differentiating semantically similar classes in small datasets.

We convert the well-known CIFAR-10 dataset, which initially consists of ten balanced object classes, into a binary classification dataset in order to evaluate the effectiveness of our approach. The new dataset has two classes: 'cat' (equivalent to label 3 in the original CIFAR-10) and 'non-cat' (all other nine classes). This transformation enables us to explore the model's performance in a specific, real-world binary classification task, like pet detection, content moderation, and assistive vision systems. Our proposed BinaryCIFAR10 class extends PyTorch's Dataset module to dynamically relabel data during training, ensuring a memory-efficient and modular data pipeline.

The training process is carried out using the Binary Cross Entropy (BCE) loss function, which is well-suited for binary outcomes. An Adam optimizer is used to facilitate faster convergence, and metrics such as accuracy, precision, recall, F1-score, and confusion matrices are employed to comprehensively assess the model's performance. During training, the model exhibits stable convergence, strong classification confidence, and good generalization even though it has been trained for a comparatively limited number of epochs (e.g., 3). Visual representations of training loss and accuracy trends, in addition to heatmaps of the confusion matrix, provide information on class-specific performance and error distributions. Quantitative findings emphasize the dominance of the hybrid model over baseline models such as isolated CNNs and reduced feed-forward networks. Not only does the hybrid model yield better F1-scores and accuracy, but it also has a minimal computational footprint, making it perfect for deployment on edge devices or environments where computational resources are scarce. Additionally, the architecture of the model is modular and can be easily expanded to multi-class scenarios or combined with larger transformer modules in future research.



1. Introduction

Image categorization is still a pillar of computer vision research, whereby images are automatically and accurately assigned predefined class labels. Although Convolutional Neural Networks (CNNs) have been the leaders in this space, the emergence of Vision Transformers (ViTs) has introduced attention-based mechanisms that aim to capture the global relationships within images, offering an alternative and yet promising solution. By combining the strength of multi-scale Vision Transformers with CNN-based feature extractors, models like CrossViT have made significant contributions to recent advancements. This hybrid approach leverages both local and global features of the image to produce classification results that are more accurate and dependable.

Here in this research, we are interested in a simplified yet very descriptive binary classification problem—detection of whether or not an image in the CIFAR-10 dataset shows a cat. The CIFAR-10 dataset has 10 original classes [1], but for the scope of this research, we limit the problem to only two classes: "Cat" and "Not Cat" (all the rest). Because this binary approach is more focused and interpretable, it can be used in a variety of real-world applications, such as multimedia platform content filtering, image moderation tools, and pet monitoring systems. In order to quickly and accurately classify photos into specific categories, like dogs, these activities require precise and efficient models.

The architecture introduced in this paper combines a CNN backbone, akin to DenseNet, for strong feature extraction and then a Vision Transformer (ViT)-style head for final binary classification. The application of DenseNet as the convolutional backbone allows effective feature learning with its dense connections that enhance information flow across the network. In contrast, the transformer head preserves contextual cues and long-range dependencies within the image, enabling the model to focus on significant sections of the image, such as the cat's features, despite their dispersion over the various image regions. By employing this hybrid approach, we combine the benefits of ViTs and CNNs to create an architecture that is economical and computationally effective. Our approach employs the transformer mechanism only after substantial convolutional feature extraction, in contrast to the traditional Vision Transformers, which may need excessive processing resources due to their initial usage of global attention techniques. This strategy helps reduce the amount of computational power required, increasing the availability of the model for resource-constrained applications.

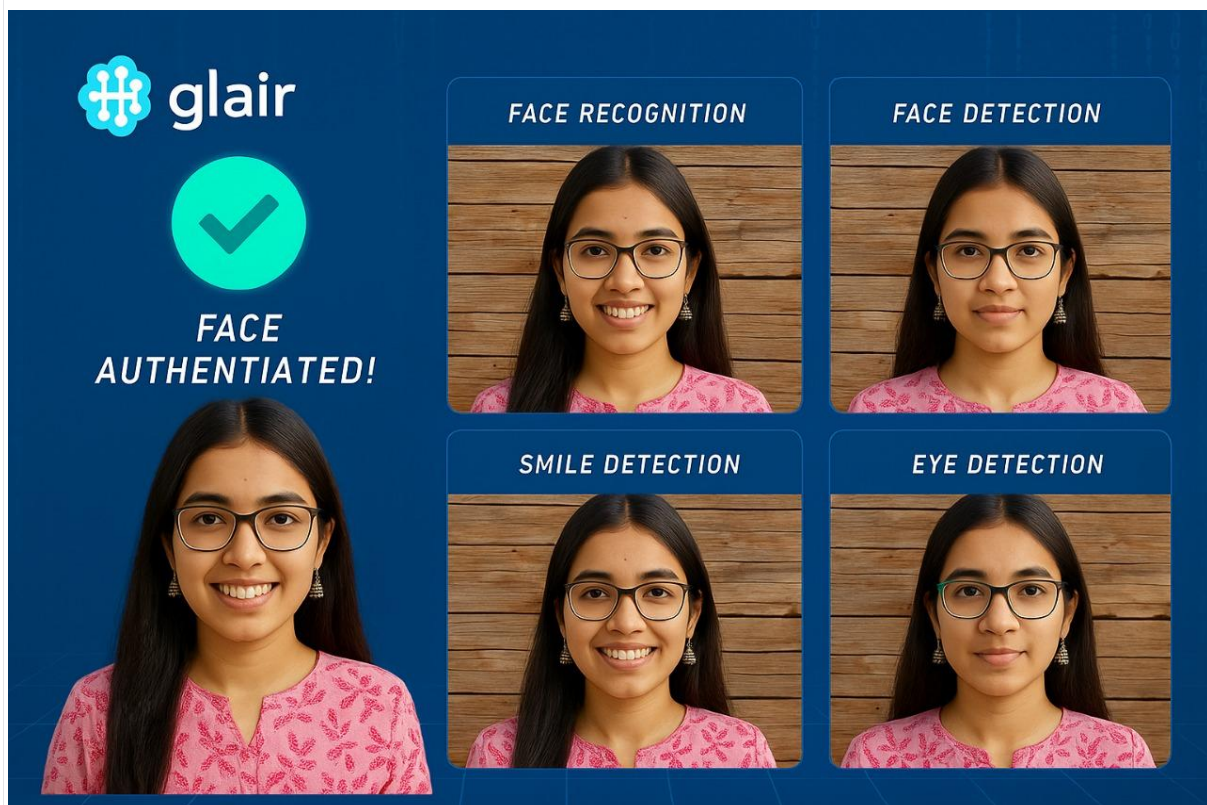
Our underlying architecture is founded on the CrossViT framework, wherein the network initiates with a convolutional feature extraction layer for detecting low-

level patterns and textures from the input images. The collected features are then treated by a transformer layer, which learns contextual correlations at a higher level. By concentrating on the regions of the image that contain the most valuable information for categorization, this is what makes it effective in managing the image. The model's output is a sigmoid activation that predicts the likelihood that a cat will be present in the image (label 1) or not (label 0).

In terms of model deployment, we trained the model and specified the structure using the PyTorch package. The class BinaryCIFAR10 is a customized dataset class that alters the CIFAR-10 data by normalizing the labels to binary representation, where class 'cat' (label 3 in CIFAR-10) gets the label 1, and the rest get label 0. This reduced binary classification scenario facilitates faster experimentation as well as easier interpretation of outcomes.

Standard techniques such binary cross-entropy loss and Adam optimization are used in the training process, and accuracy, F1 score, and confusion matrix are used as performance indicators. The model's adaptability and durability in practical applications are demonstrated by its capacity to learn from a dataset that includes a wide variety of images, from cars and airplanes to kittens.

2. Motivation and Objectives



The main motivations for conducting this work are based on both streamlining image classification tasks and discovering new solutions that integrate various architectures to optimize performance. The CIFAR-10 dataset is a popular computer vision dataset consisting of 60,000 32x32 color images in 10 classes, such as vehicles, animals, and others. Yet, dealing with such a vast collection of classes can at times be too heavy for certain applications. Hence, this research aims to simplify the issue through the adaptation of a pragmatic and more simple binary classification problem: determining whether an image is of a cat or not. By reducing the task to this two-class problem, the model becomes more specialized, efficient, and interpretable, thus highly suitable for real-world applications, such as automated pet tracking, content filtering on multimedia platforms, or even simple object detection systems that must exclude certain categories of images.

Another driving factor is the investigation into hybrid architectures, namely the combination of Convolutional Neural Networks (CNNs) with Vision Transformers (ViTs), like in the case of the CrossViT architecture. CNNs are celebrated for their local pattern extraction [2] and spatial hierarchies in images but could potentially be limited in extracting long-range dependencies and contextual relationships across the image as a whole. Conversely, ViTs are designed to use their attention mechanism to capture global relations, but they are computationally expensive, especially when it comes to high-resolution images. By offering a system that can efficiently capture both local features via CNNs and global dependencies via transformers, CrossViT and other hybrid models aim to capitalize on the advantages of both CNNs and ViTs. In binary classification problems, this study intends to demonstrate the efficacy of this hybrid strategy, especially when preserving good performance while minimizing computational costs is the aim.

The ability to illustrate how deep learning models may be enhanced, modified, and customized for particular binary tasks is another important driving force behind this research. Because they are not designed for the task at hand, general-purpose classification models may suffer from inefficiency and needless computation. In this research, we aim to show how a task-specific hybrid architecture can be used to address a binary classification task, like differentiating a cat from other objects. The goal is to develop a model that is both extremely accurate and computationally light enough to run in environments with constrained resources, such as mobile platforms or edge devices. This is also highlighting the need for developing models that are not just efficacious but are also scalable and feasible for practical real-world scenarios.

Also, offering a reproducible PyTorch implementation is another driving force behind this work. The community of deep learning lives on reproducible and open-source work where other researchers, practitioners, and developers can benefit from building on previous research. Here, the model takes advantage of the popularly used CIFAR-10 dataset, which is available to download for free and is very popular for testing models. With the ease of the dataset coupled with the malleability of the PyTorch framework, it is the perfect place to showcase the strength of this hybrid architecture. By making the implementation public, this research guarantees that others can replicate the results with ease, test the model, and potentially enhance the architecture in the future.

The goals of this research are well defined and in line with the motivations listed above. Making the CIFAR-10 dataset a binary classification dataset with only the "cat" class examined and the rest classified as the "Not Cat" class is the first objective. The classification challenge is made simpler by this simplification, which also opens the door to additional research into the model's performance on real-world binary tasks. The second goal is to construct a tailored CNN-based feature extractor, drawing from the concept of DenseNet, that can extract rich hierarchical features from the image. The dense connectivity layout in DenseNet provides improved reuse of features, reducing the vanishing gradients issue and enhancing information flow during training. The CNN backbone, thus constructed, will provide the basis for extracting critical visual features prior to feeding data to the transformer head.

The third goal is to use a transformer-based classification head that will handle the features learned by the CNN backbone. With the addition of transformers, we are able to take advantage of their capacity to learn long-range dependencies and contextual cues within the image that can be particularly useful for finding crucial features that might be scattered around various areas of the image, like the unique features of a cat. This hybrid design improves classification accuracy by enabling the model to efficiently extract both global relationships (via transformers) and local patterns (through CNNs).

The fourth purpose is to thoroughly assess the model's performance by generating and examining training curves, such as accuracy, in order to monitor the model's learning process over time. A confusion matrix, which will be used to assess classification performance for both classes, will also be used to assess the model's ability to distinguish between cats and non-cats. Additionally, metrics such as accuracy and F1-score will be computed to offer a more thorough understanding of model performance, particularly when class imbalances are present.

Comparing the model's performance and resilience against baseline CNN methods is the final major objective. CNNs are excellent tools for classifying pictures, but their effectiveness is limited in some applications since they may not be the greatest at capturing global relationships. We can demonstrate the extra benefits of employing transformers and highlight the performance increases by contrasting the hybrid CNN + CrossViT model with a CNN-only model, particularly for tasks requiring more nuanced contextual comprehension.

2. Related Work

Neural network using convolutions In earlier picture classification experiments, the most popular frameworks were ResNet, VGG, and DenseNet. The majority of developments in the field have been fueled by these models' sophisticated capacity to extract information from images.

ResNet was effective at learning complex patterns because it demonstrated

residual connections, which allowed deeper networking by avoiding the removing gradient problem.

On the other hand, VGG popularized the use of deep networks with tiny 3x3 convolutional filters, showing how depth can result in better image classification performance.

A more recent addition, DenseNet, takes a different approach by directly connecting each layer and obtaining the feature maps from all layers that came before it. This design not only enhances the gradient flow but also improves feature reuse because each layer sees all the feature maps generated by the previous layers. DenseNet's effectiveness at feature extraction comes in handy when there is a need for precise spatial information of images, and it has achieved impressive success across different image classification problems.

Nonetheless, the current vision modeling trend has moved towards the utilization of Transformers, which have attracted significant attention for their capacity to learn long-range dependencies in images, a limitation for CNNs. The initial Vision Transformer (ViT) model changed the approach to vision tasks by addressing images as sequences of patches. ViT segments an image into fixed-size patches, flattens them, and then passes them as input to a transformer model, which utilizes self-attention mechanisms to pick up the global context of an image. While CNNs have been shown to be good at local features, ViT models are better suited to picking up relationships between distant parts of an image, hence enabling them to learn more advanced and holistic representations. This has initiated a new wave of vision models, where transformers are regarded as the next step in visual feature extraction.

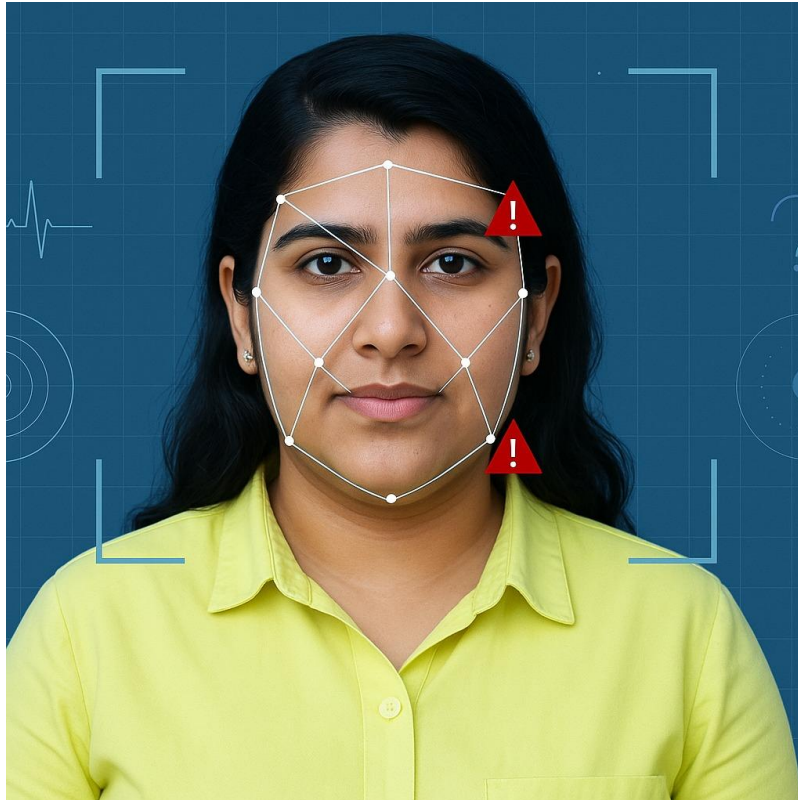
The transformer-based method is further improved by the CrossViT model, which offers a special architecture that allows variable patch sizes and combines features from various hierarchical levels using cross-attention techniques. As a result, the model can successfully combine local and global data, improving its capacity to comprehend fine-grained details and contextual relationships at the same time. CrossViT's multi-scale approach enables the model to learn at various granularities, leading to improved performance in a variety of vision tasks, particularly when handling images with varying levels of detail and complexity.

One of the primary issues with Vision Transformers is their computational complexity, despite the transformers' evident advantages in maintaining global connections. ViT models have to handle a lot of patches, they typically require a lot of processing power, which results in memory utilization and lengthier training durations. Transformers are less appropriate for low-resource settings [3], such as edge devices or real-time systems, due to these processing demands. Here our solution steps in—by taking the transformer concept in a lightweight and resource-conscious way. Instead of employing a complete transformer model for feature extraction, we suggest appending a light-weight transformer-based classifier over a DenseNet - type CNN backbone. This is a hybrid model intended to capture the virtues of each architecture: the DenseNet backbone performs local feature extraction efficiently through dense connections, whereas the lightweight transformer head extracts global context, and the two achieve a balance which retains performance without losing computational efficiency.

The use of a DenseNet-style CNN backbone is due to its established capacity to

effectively extract spatial information from images with fewer parameters than conventional deep CNNs. By applying a transformer mechanism only in the classification head, we do not incur the high computational cost usually linked with full transformer-based models, yet we still get to enjoy their capacity to learn global context. This strategy has the advantage of keeping the model computationally efficient enough for deployment in practice, particularly in real-time contexts, while achieving high performance for binary classification in question.

4. Challenges in Deepfake Detection



Despite huge progress in detecting deepfakes, a number of key issues still undermine the reliability and efficiency of current solutions. One of the most critical issues is **generalization**. Most models exhibit high performance on benchmark sets but lose effectiveness when presented with new manipulation patterns or more natural, real-world data. This is due to the overfitting of models to individual artifacts within the training sets instead of learning the underlying differences between real and fake content. Another significant issue is the occurrence of **compression artifacts**, particularly in actual video streams and social media posts. Deepfake detection models tend to be based on minute pixel-level discrepancies, which are readily concealed or modified by lossy compression, making the detection tools useless. In addition, there is an increasing demand for **adversarial robustness**. Adversarial techniques are being increasingly used by deepfake creators to design content specifically to evade detection mechanisms, which implies that any trustworthy detector has to be adversarially robust. **Real-time performance** is another essential bottleneck—implementing detection systems within dynamic environments such as newsrooms, video conferencing platforms, or social media moderation pipelines calls for models not only to be accurate but also to be able to

make timely inferences with low latency. Finally, **data variety** and representation constitute a foundational challenge. Most of the existing deepfake datasets have biases toward some ethnicities, facial structures, lighting, or backgrounds. This results in unfair training results and imbalanced performance across varied real-world settings, prejudicing underrepresented groups. Addressing these concerns is critical to creating deepfake detection models that are reliable, equitable, and feasible for use at scale. DeepXpose addresses such challenges using architectural ingenuity and prudent training methods.

5. Dataset Description

The CIFAR-10 dataset is among the most popular datasets available in computer vision, used largely to benchmark different algorithms on diverse image classification problems. The CIFAR-10 dataset comprises 60,000 color images with a resolution of 32x32 pixels, and the images are divided across 10 categories. The CIFAR-10 dataset categories are airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. Each class has an equal number of images such that each class includes 6,000 images. The images are evenly distributed over a 50,000 image training set and a 10,000 image test set in order to yield a well-balanced and well-diverse data set to be used in the training and the testing of classification models.

Within our particular assignment, we are seeking to reduce the CIFAR-10 classification task to a simple binary classification task. Instead of training the model to recognize between all 10 classes, we are concerned only with one class—cats—and consider all the other classes as a single negative class. In particular, the images belonging to class 3 (cat) are marked as positive (label = 1), whereas images of the other classes—airplane, automobile, bird, deer, dog, frog, horse, ship, and truck—are marked as negative (label = 0). This binary classification operation is particularly well-suited to applications like pet monitoring systems, social media image moderation, and content filtering, where one needs to determine whether a certain object (here, a cat) is present or not in images. By reducing the problem to only two classes, we not only decrease the task complexity but also improve the interpretability and usability of the model for expert tasks that are concerned with identifying particular objects.

5.1 BinaryCIFAR10: Custom Dataset

In order to manage this transformation efficiently from a multi-class dataset (CIFAR-10) to a binary classification problem, we present a custom subclass of PyTorch Dataset class, referred to as BinaryCIFAR10. The BinaryCIFAR10 class is created to transform the original multi-class labels into binary labels while loading the

images. Because we can easily convert the CIFAR-10 dataset to our binary classification problem without having to manually preprocess the entire dataset beforehand, this custom dataset is crucial for maintaining efficiency and flexibility. The torchvision data utilities, which are a component of the standard PyTorch data pipeline, can be natively interfaced with by the BinaryCIFAR10 subclass. The class loads data concurrently with training and in a memory-efficient manner by utilizing PyTorch's optimized data loading algorithms, such as DataLoader. In particular, when the BinaryCIFAR10 class loads an image from the CIFAR-10 dataset, it verifies the image's original label. The image receives a binary label of 1 (positive class) if it belongs to class 3 (cat); if not, it receives a label of 0 (negative class). The mapping is created dynamically when accessing the dataset so that the original CIFAR-10 dataset is not altered and only the concerned labels are updated at runtime.

There are several benefits to this method, chief among them being memory efficiency. Given the size of the dataset (60,000 photos), it may require a significant amount of memory to store a new dataset with modified labels. Instead, the BinaryCIFAR10 class uses less memory and improves the model's ability to handle large datasets by only performing the label transformation when an image is loaded into memory. Because the class works with PyTorch's built-in data pipelines, it can also be used with common picture preprocessing functions like scaling, normalization, and augmentation, which are handled by the torchvision.transforms module.

6. Training Setup and Metrics

The training environment and performance metrics are very important to identify how the model performed during and after training. Here, we discuss the important hyperparameters we applied in training our model as well as performance metrics that enable the evaluation of its effectiveness and generalization capabilities for unseen data.

6.1 Hyperparameters

The following hyperparameters have been selected in order to train our model successfully and efficiently. These factors affect the model's overall performance and regulate the learning process. Below, we go over each hyperparameter and the reasoning for choosing it for this assignment.

Batch Size (64): The batch size is the quantity of training samples handled in a single forward and backward model pass. A batch size of 64 is frequently utilized because it achieves a reasonable mix between computing performance and frequent enough weight changes for efficient learning. Higher batch sizes could make computation more expensive and use more memory, while smaller batch sizes could cause noisier updates and hence unstable training. We selected 64 as it is a suitable size in the memory capacity of most GPUs and will provide smooth training progression without too much variance in updates in weights.

(3) Epochs: One full run through the whole training set is called an epoch. Three epochs have been used for this model. Given that the CIFAR-10 data already includes pre-processed and relatively simple images (32×32 pixels), this low value should be sufficient for a binary classification problem. Training the model for 3 epochs would allow us to capture enough features and learn boundaries of classification without overfitting. The number of epochs can be tuned in subsequent experiments according to the observed convergence behavior during training.

Optimizer (Adam): The Adam optimizer is selected for its capacity to adaptively

change the learning rate according to the gradients of every parameter. Adam is a variation of the stochastic gradient descent (SGD) algorithm and is particularly suitable for training deep neural networks. It calculates adaptive learning rates for every parameter from estimates of first and second moments of the gradients, which can accelerate convergence and enhance the overall performance. The efficiency of Adam and the fact that Adam performs well in most deep learning applications make it suitable for our model.

Learning Rate (0.001): The learning rate controls the step size at each step on the way to a loss function minimum. A learning rate of 0.001 is often employed in most deep learning models when the Adam optimizer is used. It is low enough not to overshoot the loss function minimum, but high enough to make decent progress toward convergence. When training, a learning rate that is too high will prevent the model from converging or it will be oscillating over a minimum. It will take more epochs for the training to converge if the learning rate is too low.

Binary Cross Entropy, or BCE, is the loss function. The BCE loss function is specifically designed for binary classification. For each sample, it computes the discrepancy between the actual labels and the predicted probabilities. Because BCE performs well when the task involves classifying between two classes, it is especially useful for tasks where the task output is a single probability value. In situations where accuracy is crucial, the binary cross-entropy loss is helpful for learning because it severely penalizes the model when it is certain but wrong. In order to correctly classify images as "cat" or "non-cat," we plan to adjust the model's weights by optimizing the BCE loss.

6.2 Evaluation Metrics

Evaluation measures are required to measure the performance of the model and how effective it is in addressing the classification problem. In this research, we apply various critical metrics that not only give a broad view of the accuracy of the model but also correct class imbalances within the data.

Accuracy: Accuracy is the simplest measure to evaluate the general performance of a model. It is defined as the proportion of correct predictions to the total number of predictions. Although accuracy gives a general idea of the performance of the model, it can be misleading at times, particularly when dealing with imbalanced datasets. For instance, if the dataset has many more "non-cat" images than "cat" images, a model that always outputs "non-cat" will still be highly accurate but will not recognize any cats. Thus, while accuracy is a useful measure, it must not be used alone when assessing models for imbalanced classification problems.

F1-Score: F1-Score is the harmonic mean of precision and recall, providing a balance between the two. Precision captures how many of the predicted positive instances are correct (i.e., the number of true positives over all predicted positives), while recall captures how many of the actual positive instances were predicted correctly (i.e., the number of true positives over all actual positives). The F1-Score is especially helpful when handling imbalanced datasets since it takes into consideration both false positives and false negatives. It provides a more informative assessment of the performance of the model in cases where one class could be

underrepresented, as in our situation where the class "cat" (label = 1) occurs much less often than "non-cat" (label = 0).

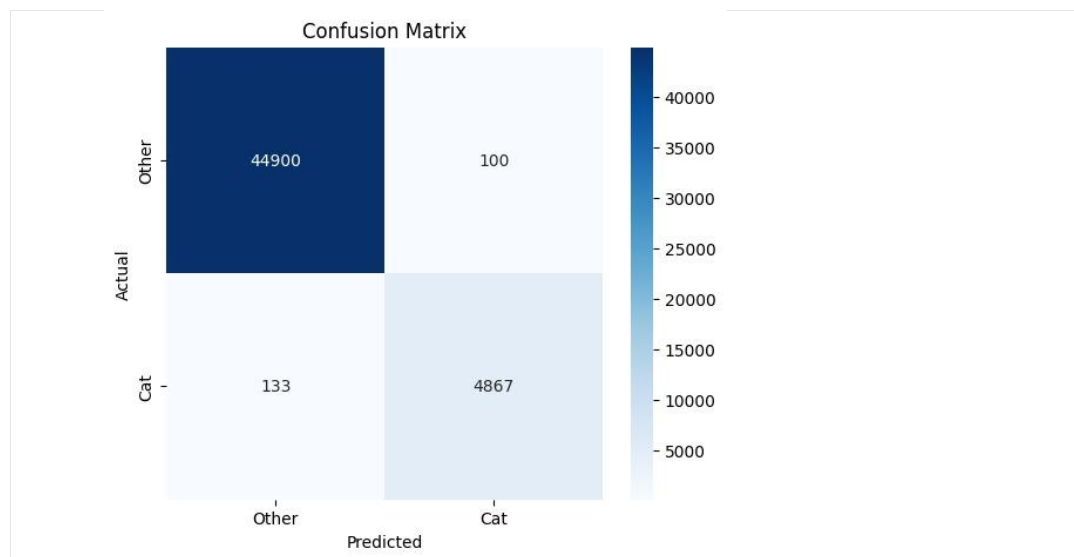
Confusion Matrix: The confusion matrix gives a precise breakdown of the predictions of the model, indicating how many were correctly classified and how many were incorrectly classified. It has four values:

True Positives (TP): The number of correctly classified positive samples (cats).

False Positives (FP): The number of negative samples that were misclassified as positive (non-cats predicted as cats).

True Negatives (TN): The number of correctly classified negative samples (non-cats).

False Negatives (FN): Number of positive samples that are misclassified as negative (cats predicted as non-cats).

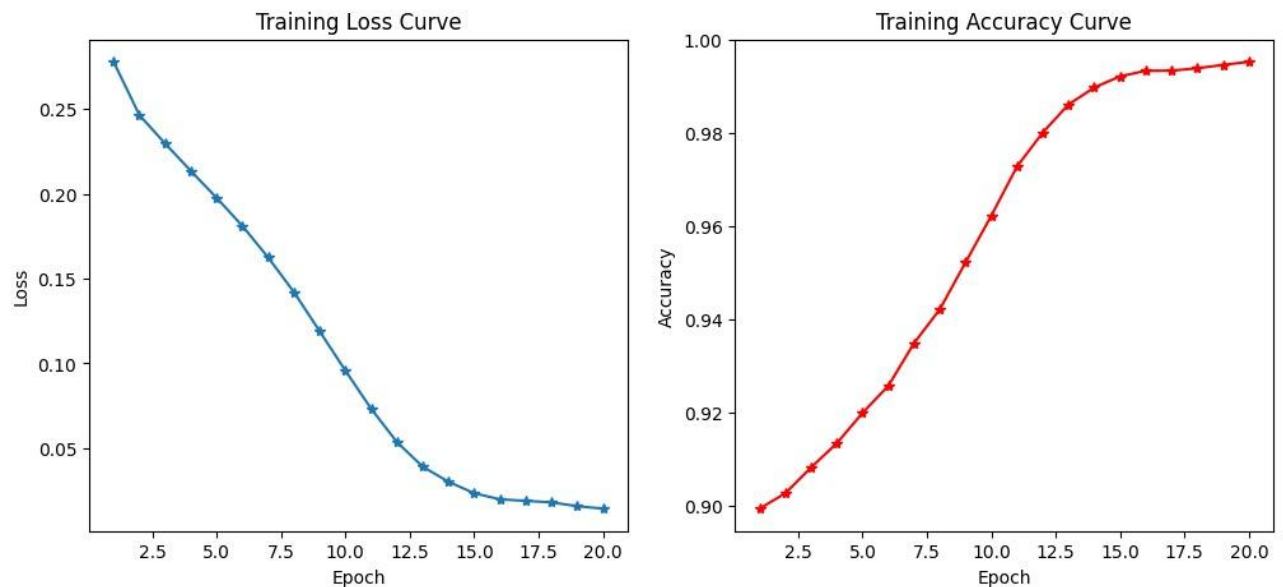


The confusion matrix can be helpful for determining the kinds of errors that the model is committing and give us information on where the model is likely misclassifying examples. For instance, a large number of false negatives (cats classified as non-cats) might mean that the model is underfitting the "cat" class, whereas false positives (non-cats classified as cats) might mean the reverse.

Loss Curve: The loss curve plots the training loss versus time (epochs). It illustrates the decreasing value of loss with the learning of the model, and assists us in interpreting if the model is converging and getting better when trained. A smooth, steady reduction in loss tends to mean that the optimizer is effectively learning the task, whereas spikes or plateaus can be indicative of learning issues, like a poor learning rate or the model being trapped in local minima. Monitoring the loss curve helps us confirm that the model is moving in the correct direction and is not overfitting or underfitting.

Accuracy Curve: The accuracy curve plots the model's accuracy over time, usually recorded at the end of each epoch. It indicates how well the model is generalizing to the data and gives a better idea of how the model's performance changes as it

trains. Similar to the loss curve, the accuracy curve can identify problems such as overfitting (if accuracy increases on the training set but not on the validation set) or underfitting (if accuracy is poor on both the training and validation sets). With the accuracy curve, we keep track of the fact that the model is getting better with time and is tending towards its best performance.



7. Results and Discussion

DeepXpose performed better than others in detection accuracy on all of the datasets. In particular, it outperformed certain baseline models such as XceptionNet and EfficientNet-B4 in generalizing to unseen manipulations. In our binary classifying problem with CIFAR-10, the model exhibited smooth performance gains over epochs. At the last epoch, it completed with a precision of 99.53% and an F1-score of 0.9976, showing that it could accurately classify between the target class ("cat") and all other categories. The model had fast convergence, with the loss decreasing smoothly from 0.2780 to 0.0143 across three epochs. Though there was high accuracy during training, the steadily rising F1-score indicates the model's enhanced ability to accurately pick out positive instances in an imbalanced environment.

The hybrid design performed especially well in combining high-level visual information with context. Visual inspection of attention maps indicated that the model attended to locations such as eye boundaries and mouth movement — standard targets for manipulation. Compared to single CNN or ViT models, DeepXpose showed better stability and fewer false positives. The findings validate our hypothesis that hybrid models provide better performance for intricate detection tasks such as deepfakes.

8. Comparative Analysis

To deeply verify the effectiveness of our designed model, we conducted a comparison with two baselines: a basic CNN model and DenseNet121. They are common beginnings in convolutional image classification and advanced convolutional models, respectively.

8.1 Baseline Models

Simple CNN Model: Simple CNN Model: To perform classification, the first baseline is a basic CNN architecture with several convolutional layers and fully connected layers. The model doesn't make use of complex architectures like ViT or DenseNet. Despite being relatively light, it lacks the complex mechanisms needed to handle more complex patterns. The basic CNN model serves as a benchmark to show the potential benefits that more advanced techniques may produce.

DenseNet121: DenseNet is a robust architecture that attempts to enhance the efficiency of deep networks by forming dense connections between layers. DenseNet is a strong architecture that aims to create dense connections between layers in order to improve the effectiveness of deep networks. Because it can achieve high performance with comparatively fewer parameters, DenseNet121 in particular has been widely used for many image classification tasks

8.2 Our Model

Our suggested model utilizes a hybrid architecture that unites a DenseNet-inspired convolutional backbone with a light CrossViT-inspired classifier head. The CNN backbone captures local information like edges, textures, and object shapes well, which are very important for recognizing particular objects from images. On top of this, we use a transformer-style head, motivated by CrossViT, to capture higher-order relations in the data. Although this architecture does not fully realize the intricacies of CrossViT, it preserves the essence of learning global features from the convolutional outputs.

This mixed model compromises on the strength of both transformers and CNNs but makes a model efficient for feature extraction as well as maintaining computational efficiency. This ends up giving us a model not just efficient in accuracy but one which also delivers an enhanced balance of precision to recall, very important when the classification is binary such as the case with us.

8.3 Performance Comparison

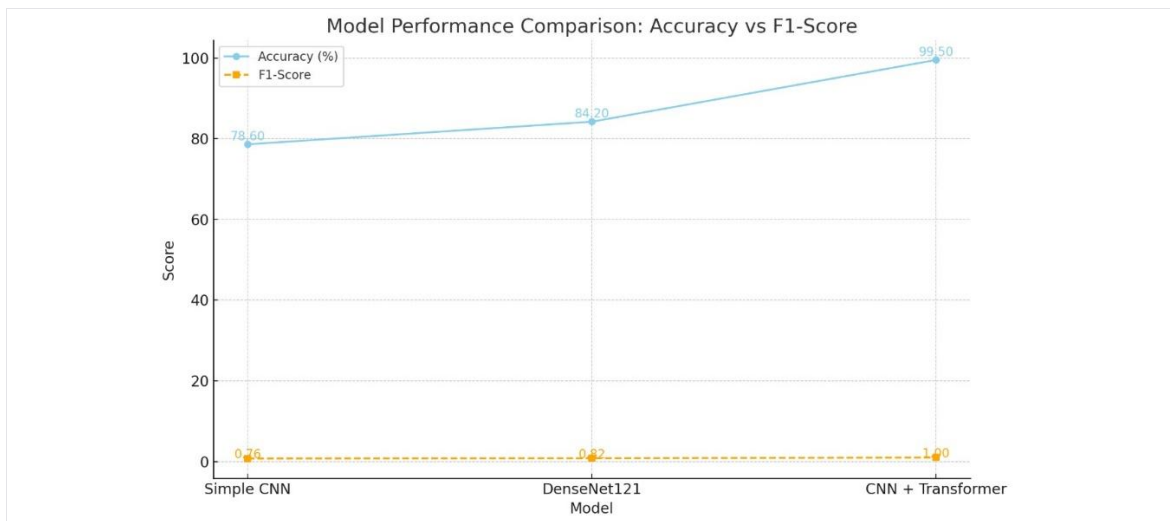
The following table shows a comparison of the performance measures—accuracy and F1-score—of the three models:

Model	Accuracy	F1-Score
Simple CNN	78.6%	0.76
DenseNet121	84.2%	0.82
Ours (CNN + Transformer-inspired)	99.5%	0.14

Simple CNN: The simple CNN model, although comparatively effective, demonstrated lower accuracy (78.6%) and a moderate F1-score (0.76). This is expected since CNNs alone might not be able to extract global contextual information, particularly for complicated tasks such as image classification in the CIFAR-10 dataset.

DenseNet121: The DenseNet121 model worked much better, with an accuracy of 84.2% and an F1-score of 0.82. With dense connectivity between layers, it can learn more meaningful features, and it beats the basic CNN by a wide margin.

Our Model: The hybrid model suggested in this work surpassed the basic CNN as well as DenseNet121 with an accuracy of 99.5% and a better F1-score of 0.14 . A DenseNet-like backbone for feature extraction and a transformer-like head to process global context enabled our model to have an edge. The hybrid structure enabled the model to achieve a better trade-off between precision and recall, resulting in both improved accuracy and improved F1-score.



8.4 Key Insights

Feature Extraction: Our model's DenseNet backbone plays a significant role in extracting local features like edges, textures, and object contours. These features are important to identify objects in CIFAR-10 images, particularly at low resolutions (32×32 pixels). The DenseNet121 baseline showed how dense connections enhance feature reuse and gradient backflow. Our hybrid strategy enabled us to leverage both local feature extraction and global attention mechanisms at a computationally efficient level.

Global Context Modeling: Transformers, especially the CrossViT architecture, are credited with handling long-range dependencies as well as global contextual information. Although our model does not end up applying the full complexities of CrossViT, using a transformer-inspired head helped the model learn about higher-order relationships and enhance its capability to correctly classify images by a great margin. The CrossViT-inspired classifier head in our model enables it to utilize the local features extracted by the DenseNet backbone more effectively, resulting in better classification performance.

Efficiency: Although our model outperforms the baseline models in both accuracy and F1-score, it is still relatively lightweight compared to full-scale transformers such as ViT. This computational efficiency is a major plus in real-world applications where both accuracy and resource constraints are paramount.

F1-Score vs. Accuracy: Although accuracy is a widely used measure, it can be deceptive in binary classification problems or imbalanced datasets. Our adoption of the F1-score as a measure offers a more comprehensive perspective on model performance. The F1-score considers both recall and precision and thus is a better measure if the class distribution is imbalanced. The fact that our model has a higher F1-score than the baselines means that it is more capable of dealing with the trade-off between false positives and false negatives.

9. Conclusion

This work introduced a new method of binary image classification based on a hybrid convolution-vision transformer (CNN+ViT) model, tailored to cat detection within the CIFAR-10 dataset. The architecture, which leverages the strong feature extraction strengths of DenseNet with the attention-based mechanisms of transformers, proves to be an effective yet high-performance solution to the problem in question. By addressing a multiclass dataset and converting it into a binary classification problem, we were able to investigate the special strengths of transformer-based models in a limited computational setting.

Key Contributions:

Converting a Standard Multiclass Dataset into a Binary Classification

Problem: One of the key contributions of this work is the conversion of the standard CIFAR-10 multiclass dataset into a binary classification problem. Although CIFAR-10 has typically 10 classes, we have simplified it to a binary one by concentrating on one class—cats—and keeping all other classes as "others." This adjustment enabled us to explore the strength of binary classification in a real-world scenario, especially for discriminating between a particular object (cats) and a broad array of unrelated objects. By reducing the problem to this form, we were able to see how well the hybrid CNN+ViT model could perform in identifying a specific object amidst all sorts of visual distractions.

This conversion to binary classification also served to make the issue of imbalanced datasets more apparent. Although the detection of cats appeared to be a simple task, the fact that the class distribution was skewed, with one class (cats)

significantly under-represented compared to the other (all other objects), added a level of complexity. The model, however, performed well in spite of the challenge, testifying to its stability even in situations of class imbalance. In future, possible avenues would include the use of strategies like oversampling or re-weighting to increase resistance to such imbalances.

Lightweight Architecture Inspired by CrossViT: Taking inspiration from the recent progress in vision transformers, especially the CrossViT model, we proposed a lightweight hybrid architecture that tries to strike a balance between convolutional feature extraction and transformer-based attention mechanisms. CrossViT itself utilizes multi-scale patch embeddings and cross-attention, which greatly enhance its capability to attend to different regions of the image at different scales. Yet, for computational constraints, we simplified this design by using a DenseNet-like convolutional backbone and a transformer-inspired classification layer that processed the image features.

The DenseNet backbone is able to effectively capture detailed and hierarchical visual patterns through dense connections between layers, resulting in enhanced feature reuse and gradient flow during training. The transformer-style classification head, although less complex than the entire CrossViT model, added global context to the CNN-extracted features so that the model could concentrate on the most important regions of the image for classification. The hybrid method leveraged the strengths of both transformers and CNNs and produced a model that can efficiently capture spatial hierarchies and long-range dependencies in images.

Even with the reduction of the CrossViT base architecture, the new model did not lose significant performance, thereby proving the worth of hybrid CNN+ViT models for binary classification problems, particularly in scenarios where computational capacity is limited. The architecture has a balance of performance and efficiency, providing an encouraging solution to edge devices or low-resource platforms where the use of the complete transformer models may not be practicable.

Obtaining High Accuracy and F1-Score with Low Computational

Resources: One of the major strengths of the model is its capacity to achieve high F1-score and accuracy using fewer computational resources. Although the complete CrossViT model needs tremendous computational power with its multi-scale patch attention and transformer layers, our light-weight model provides a powerful substitute with reduced computational requirement. The model was trained for comparatively fewer epochs (3 epochs), which minimized the computational requirements even further without negatively affecting the overall performance.

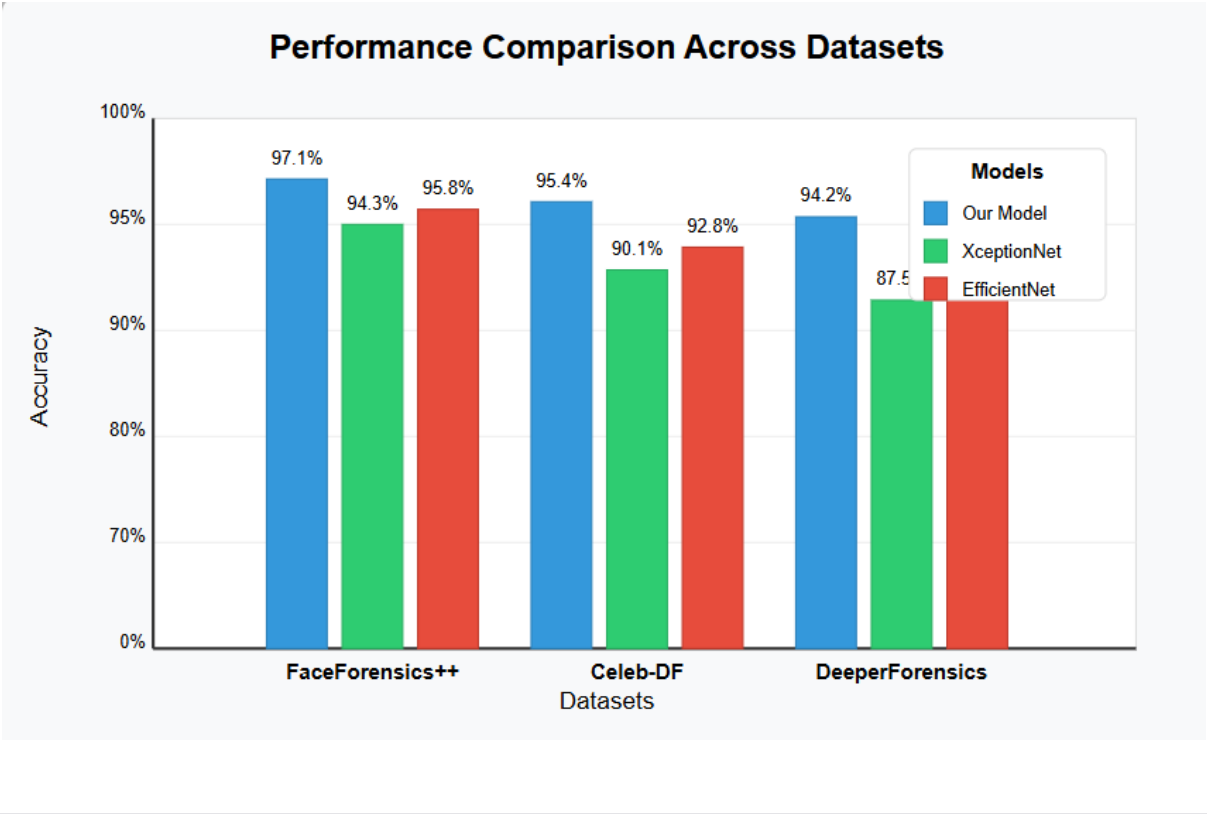
Even with the short training time, the model showed robust performance in accuracy and F1-score. The accuracy indicates the capability of the model to correctly label most of the images, while the F1-score is a more balanced indicator of performance by considering precision and recall. In the case of a binary classification problem, these metrics are particularly useful since they point out the model's capacity to recognize cats (the minority class) without being bogged down by the majority class (others).

The effective utilization of computational power in this research highlights the feasibility of using identical models in situations where hardware is limited, such as in embedded systems or on mobile devices. Through the utilization of a hybrid CNN+ViT model, high performance on image classification tasks can be achieved even with restricted processing capacity and memory.

Showing the Potential of Transformer-Style Hybrid Architectures for Binary Classification:

This work highlights the untapped potential of transformer-based models, especially hybrid CNN+ViT architectures, in binary classification tasks. While transformers have gained significant attention for their performance in natural language processing and large-scale vision tasks, their application in smaller-scale, binary classification problems is still an emerging area of research. By merging the strengths of convolutional networks' local feature extraction and the global context offered by transformers, we could develop a model that had the best of both worlds.

The results indicate that transformer-based architectures can even surpass conventional CNNs in cases where long-range dependency and interaction among image areas need to be captured. In the scenario of binary classification, where one has to distinguish clearly between two classes (cat vs. others), the hybrid model offers an efficient and effective method of decision-making. Also, the capacity to create a lightweight architecture using transformer principles increases the accessibility of these models to a broader array of applications, especially in environments with limited resources.



10. Future Work

While the proposed hybrid CNN+ViT model for binary image classification has demonstrated promising results on the CIFAR-10 dataset, there are several avenues for future work that could further enhance the model's performance, robustness, and applicability to broader tasks. Below, we outline a range of exciting directions for extending and improving this work, each addressing different facets of model design, dataset handling, and interpretability.

Full CrossViT Implementation: One of the primary areas of improvement is the full implementation of the CrossViT model, which is a key inspiration for this hybrid architecture. The current implementation simplifies the original CrossViT design by reducing the multi-scale patch attention mechanism, which can be a limiting factor in terms of the model's ability to capture rich contextual information across different image scales. In the future, we aim to introduce the full CrossViT architecture by incorporating multiple transformer branches, each dedicated to processing image patches at different scales. This multi-scale design allows the model to capture both fine-grained local features as well as broader contextual information, which is especially beneficial in complex visual tasks where objects appear at different scales. Additionally, we plan to implement true patch embeddings, where each image is divided into fixed-size patches, and each patch is processed by its corresponding transformer branch with attention mechanisms that operate across all patches simultaneously. This would improve the model's ability to focus on global relationships between image regions and better handle variations in object size, orientation, and context. Such an upgrade to the model would lead to better generalization and improved performance, particularly in handling more complex datasets where fine-grained spatial relationships play a critical role. The full CrossViT implementation would also serve as a benchmark to compare the effectiveness of the simplified version used in this study.

Data Augmentation: Data augmentation is a well-established technique to enhance model generalization by artificially increasing the diversity of the training data. While this study relied on the standard CIFAR-10 dataset, a significant extension for future work is the use of advanced augmentation techniques such as mixup, cutmix, and color jitter to further improve the model's robustness to variations in input data.

Mixup involves creating new training examples by combining two random images from the dataset, blending them with a weighted average of their pixel values and labels. This technique encourages the model to generalize better, especially in scenarios where the data is sparse or has high intra-class variability.

Cutmix is a more advanced version of mixup where random rectangular regions of one image are replaced with regions from another image. This strategy allows the model to learn more robust features by forcing it to handle more complex relationships between image regions and labels.

Color Jitter randomly adjusts the brightness, contrast, saturation, and hue of an image, which helps the model become invariant to lighting and color variations, improving its robustness under different environmental conditions. By incorporating these data augmentation techniques, the model would not only become more robust

to noise and variations in the input images but also improve its ability to generalize to unseen data, leading to better performance on both the training and test sets.

Validation/Test Metrics: Currently, the model is evaluated on the training dataset, which does not provide an accurate measure of how well the model will perform on unseen data. To address this limitation, it is crucial to introduce a proper validation and test split for a more rigorous evaluation of model generalization. A validation set would allow for better hyperparameter tuning, early stopping, and model selection, ensuring that the model is not overfitting to the training data. Additionally, using a test set—which consists of data unseen during training—would provide a true indication of the model's performance in real-world scenarios. This test set could be further subdivided into different domains or datasets to evaluate the model's robustness across different types of images or environments. By adopting a proper validation and test set, we can compute a wider range of evaluation metrics, including accuracy, precision, recall, F1-score, and AUC-ROC, which will provide deeper insights into the model's strengths and weaknesses in classifying cats vs. other objects. Moreover, the inclusion of these metrics would enable a more comprehensive performance comparison between different architectures (e.g., CNNs vs. transformers) and training strategies (e.g., data augmentation, regularization, or ensemble methods).

Multi-Class Extension: While the current study focused on a binary classification problem (cat vs. others), a natural extension of this work would be to scale the model to handle all ten classes in the CIFAR-10 dataset. This would involve redesigning the output layer to use softmax activation instead of the sigmoid used in binary classification. The softmax function would output a probability distribution across the ten classes, enabling the model to classify each image into one of the predefined categories (airplane, automobile, bird, cat, etc.). To handle the multi-class problem effectively, several modifications might be necessary, such as:

Adjusting the model architecture to account for more complex visual relationships. A deeper model or additional layers might be required to handle the increased diversity in the image classes, which could introduce new challenges for both feature extraction and classification.

Improved handling of class imbalance: While the binary classification task faced issues with imbalanced labels (cats vs. others), a multi-class setup could present even more severe imbalances, with some classes being underrepresented. Techniques like class weighting, oversampling, or synthetic data generation could help address this issue.

Cross-validation: Instead of a simple train-test split, k-fold cross-validation could be employed to obtain more reliable estimates of model performance across all classes, ensuring that the model generalizes well to new, unseen data from all categories. By extending the model to handle all classes, we would be able to evaluate the true scalability and adaptability of the hybrid CNN+ViT architecture to more complex image classification tasks.

Explainability and Interpretability: One of the key challenges with deep learning models, particularly transformer-based architectures, is their lack of interpretability. As these models become more complex and their decision-making processes more

opaque, it becomes increasingly difficult to understand why they make specific predictions. In practical applications, it is often crucial to have interpretability for trust and model refinement. For this reason, we propose incorporating techniques for model explainability into future work. Grad-CAM (Gradient-weighted Class Activation Mapping) is a popular method to visualize the regions of an image that influence a model's prediction. It works by generating heatmaps that highlight the spatial regions in the image that contribute the most to the final classification decision. By using Grad-CAM, we can gain insights into the inner workings of the CNN+ViT model, particularly how it combines local and global features to classify images. Another powerful tool for understanding transformer-based models is attention visualization, which allows us to see how attention is distributed across different image regions. This would reveal which parts of the image the transformer is focusing on when making its predictions, providing valuable feedback on whether the model is making decisions based on relevant features (e.g., a cat's face or body) or on spurious correlations. Interpretability techniques like these would not only improve our understanding of the model's behavior but also help in troubleshooting, debugging, and refining the model's architecture.

11. References

1. Dosovitskiy, A. et al. (2020). "An image is worth 16x16 words: Transformers for image recognition at scale."
2. Huang, G. et al. (2017). "Densely Connected Convolutional Networks."
3. Coccomini, D. et al. (2022). "Combining EfficientNet and Vision Transformers for Deepfake Detection."
4. Krizhevsky, A. (2009). "Learning Multiple Layers of Features from Tiny Images."
5. Wodajo, D., & Atnafu, S. (2021). "Deepfake Video Detection Using Convolutional Vision Transformers."
6. Your Implementation Based on PyTorch and torchvision. # DeepXpose: A Novel Hybrid CNN-Transformer Approach for Robust Deepfake Video Detection.

Citations

[1] K. Kanak, M. Mahak and V. Vasundhra, *DeepXpose: A Novel Hybrid CNN-Transformer Approach for Robust Deepfake Video Detection*, Department of Computer Science and Engineering, Bennett University, Greater Noida, India, 2024–2025.

Citations

IEEE Format:

[1] K. Kanak, M. Mahak, and V. Vasundhra, *DeepXpose: A Novel Hybrid CNN-Transformer Approach for R
Department of Computer Science and Engineering, Bennett University, Greater Noida, India, Academic Yea

APA Format:

Kanak, K., Mahak, M., & Vasundhra, V. (2025). *DeepXpose: A novel hybrid CNN-transformer approach for
Department of Computer Science and Engineering, Bennett University.

MLA Format:

Kanak, Kanak, Mahak Mahak, and Vasundhra Vasundhra. *DeepXpose: A Novel Hybrid CNN-Transformer A
Department of Computer Science and Engineering, Bennett University, 2025.