



Deep Learning and Applications, Module I - Deepfake Detection

Congiu, Francesco¹; Giuffrida, Simone²; Littera, Fabio¹;

Puglisi, Giovanni¹;

University of Cagliari, Department of Mathematics and Informatics

April 4, 2025

Abstract: In this paper we present a novel approach for the detection and segmentation of deepfakes in images containing multiple faces. Our method integrates transfer learning techniques with the development of a custom CNN network for binary classification (real vs. fake), leveraging the OpenForensics dataset.

Keyword: Computer Vision, Neural Network, Deepfake Detection, Bounding Box Detection.

1 Introduction and Background

1.1 Loss-function and optimization

In supervised learning, and particularly in image classification tasks, **the loss function and optimization method** form the basics of the training process. The loss function quantifies how well (or poorly) the network is performing by comparing its predictions to the true labels, while the optimizer describes how the model parameters should be updated to minimize this discrepancy. An appropriate choice of loss function and optimization algorithm is crucial: if the loss function is not aligned with the task's objectives, the model may learn suboptimal patterns, and if the optimizer is inefficient or poorly tuned, training may take excessively long or fail to converge to a good solution. In deepfake detection, where the objective is to classify images (or frames) as real or fake, the cross-entropy loss is typically employed due to its effectiveness in multi-class and binary classification scenarios. Combined with an optimizer like Adam, which adaptively tunes individual learning rates for each parameter in the model, practitioners can achieve faster convergence and maintain stable training dynamics even when dealing with large-scale datasets or complex architectures.

1.1.1 Loss function: Cross-entropy

Cross-entropy loss, often referred to as softmax loss or log loss in classification contexts, is the standard choice for deep learning models aiming to distinguish

between two or more classes. Formally, given a set of training examples (x_i, y_i) , where x_i is the input (e.g., an image) and y_i is the corresponding label (in one-hot or integer form), **the cross-entropy loss measures the dissimilarity between the predicted probability distribution \hat{y}_i and the true distribution y_i .** For a single training example in a binary classification problem, the cross-entropy loss can be expressed as:

$$L = \left[y \times \log(\hat{y}) + (1 - y) \times \log(1 - \hat{y}) \right]. \quad (1)$$

In multi-class settings, \hat{y} and y become vectors representing the predicted and true probability distributions over the possible classes, and the loss is summed over all classes. This formulation effectively penalizes confident but wrong predictions by producing large gradient signals, helping to reduce misclassification errors. Cross-entropy aligns well with the probabilistic interpretation of classification, where each output neuron in the final layer corresponds to the logit for a specific class, and the softmax activation ensures that these logits sum to 1. Cross-entropy is almost universally adopted in many tasks regarding CNNs like deepfake detection, image recognition, and others.

1.1.2 Optimization: Adam

Adam¹ stands for Adaptive Moment Estimation, and is a popular stochastic optimization algorithm introduced by D. Kingma and J. Ba in 2014. It has become a standard in modern deep learning because of

¹Adam: A Method for Stochastic Optimization

its ability to combine the benefits of two precedent optimizers: AdaGrad, which adapts the learning rate based on the frequency of parameter updates, and RMSProp, which normalizes the gradient by a moving average of its recent magnitudes. Specifically, Adam maintains exponentially decaying averages of past gradients (first moment) and their squares (second moment), allowing it to adjust the learning rate for each parameter dynamically. This feature is particularly beneficial in high-dimensional parameter spaces, where different parameters can exhibit vastly different gradient statistics. By mitigating the need to manually fine-tune separate learning rates, Adam often converges faster and more reliably than stochastic gradient descent (SGD), especially for architectures like MobileNet and Xception that can have millions of parameters. Adam's sensitivity to hyperparameters (notably the learning rate and beta values controlling the exponential decay rates) means that a careful balance must still be struck to avoid overshooting or slow convergence. In deepfake detection pipelines, Adam's adaptive nature helps stabilize training, allowing the model to quickly learn to recognize artifacts introduced by face-swapping or generative adversarial networks without getting stuck in local minimums.

1.2 Pre-trained network architectures: MobileNet e Xception

This section provides an overview of two custom neural networks, MobileNet and Xception, with specific attention regarding their architectures and underlying principles. Both architectures use the concept of depthwise separable convolutions to obtain an optimal balance between computational efficiency and accuracy, although they were designed for different purposes and applications.

1.2.1 MobileNet

MobileNet² is a suite of highly efficient convolutional neural network architectures specifically designed for mobile and embedded vision applications. Its key innovation is the use of depthwise separable convolutions—a type of factorized convolution that splits a standard convolution into two separate operations. In the first stage, the depthwise convolution applies a single filter to each input channel, effectively isolating channel-specific features. In the sec-

ond stage, the pointwise convolution (a 1×1 convolution) combines these individual outputs to form a new set of features. This two-step process replaces the traditional one-step convolution that both filters and combines inputs, leading to a significant reduction in computational cost and model size without a substantial loss in accuracy.

The following image shows the difference between:

- Standard convolutional layer with batch normalization and ReLU on the left;
- Depthwise separable convolutions with depthwise and pointwise layers followed by batch normalization and ReLU.

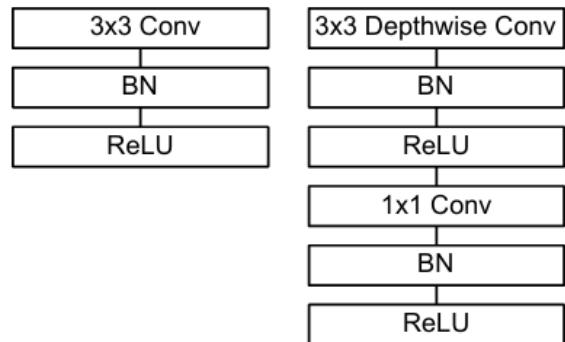


Figure 1: MobileNet architecture.

MobileNet distinguishes itself with its remarkable flexibility. It introduces two global hyperparameters, the width multiplier and the resolution multiplier, which allow developers to fine-tune the trade-off between accuracy, latency, and model size. This adaptability enables the network to be scaled down for resource-constrained environments while still maintaining competitive performance on tasks like image classification, object detection, and beyond.

1.2.2 Xception

Xception³ is a convolutional neural network architecture that takes the idea of depthwise separable convolutions to its extreme. Instead of relying on complex inception modules, Xception completely decouples the learning of spatial and cross-channel correlations. This is done in two simple steps: first, a depthwise convolution independently extracts spatial features from each channel; then, a pointwise

²MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications, Andrew G. Howard et al.

³Xception: Deep Learning with Depthwise Separable Convolutions, François Chollet

(1×1) convolution fuses these features across channels. This “extreme” formulation (hence the name Xception, short for “Extreme Inception”) leads to a more efficient use of parameters.

A convolution layer attempts to learn filters in a 3D space, with 2 spatial dimensions (width and height) and a channel dimension; thus a single convolution kernel is tasked with simultaneously mapping cross-channel correlations and spatial correlations. The idea behind the Inception module is to make this process easier and more efficient by explicitly factoring it into a series of operations that independently look at cross-channel correlations and spatial correlations. The typical Inception module first looks at cross-channel correlations via a set of 1×1 convolutions, mapping the input data into 3 or 4 separate spaces that are smaller than the original input space, and then maps all correlations in these smaller 3D spaces via regular 3×3 or 5×5 convolutions.

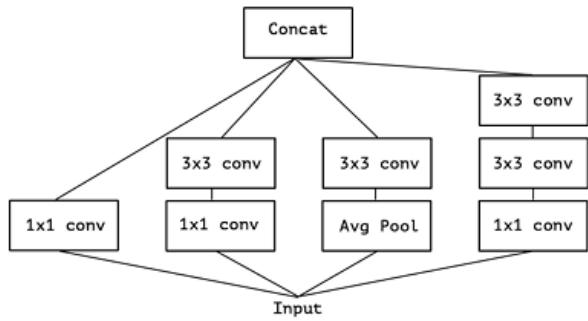


Figure 2: Canonical Inception module (Inception V3).

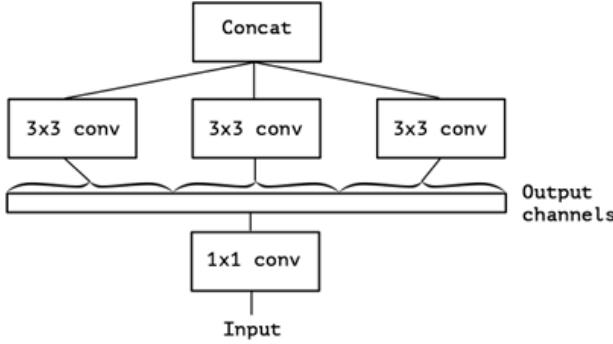


Figure 3: An “extreme” version of the Inception module, with one spatial convolution per output channel of the 1×1 convolution.

1.2.3 Differences between the two models

1. Design Motivation and Target Use-Case:

MobileNet was primarily designed for mobile and embedded vision applications. It emphasizes efficiency by trading off accuracy for lower latency and smaller model size, making it ideal for resource-constrained environments.

Xception reinterprets the Inception module by taking the idea to its extreme. Its motivation is to fully decouple the mapping of spatial correlations from cross-channel correlations using depthwise separable convolutions, aiming for improved performance on large-scale image classification tasks.

2. Architectural Composition:

MobileNet builds its architecture almost entirely from depthwise separable convolutions and introduces two hyperparameters—the width multiplier and the resolution multiplier—to flexibly adjust the network’s size and computational cost.

Xception replaces Inception modules with a linear stack of depthwise separable convolutions organized into deeper modules (36 convolutional layers grouped into 14 modules) and incorporates residual connections throughout.

3. Order of Operations and Activation Placement:

In **MobileNet**, each depthwise separable convolution is executed by first applying a depthwise (spatial) convolution and then a pointwise (1×1) convolution, with batch normalization and ReLU applied after each convolution to capture both spatial and cross-channel correlations.

Xception also uses depthwise separable convolutions but omits the non-linearity (such as ReLU) between the depthwise and pointwise operations. However, after each convolution block, Xception still uses batch normalization followed by ReLU.

4. Use of Residual Connections:

MobileNet does not incorporate residual (skip) connections in its basic architecture.

Xception makes extensive use of residual connections around its modules to help with training convergence and overall performance.

These differences reflect the distinct goals of the two models: MobileNet is optimized for efficiency under tight computational constraints, while Xception leverages a more radical factorization of convolutions (with residual connections) to push performance on large-scale tasks. An in-depth view of how the Xception model can be used is seen in the paper by R. Helaly et al. (available at the link)⁴.

1.3 Activation Functions in MobileNet, Xception, and CustomCNN

Activation functions play a crucial role in deep neural networks by introducing non-linearity, enabling models to learn complex patterns and relationships in data. They affect the model's learning capacity and overall performance. In this section, we examine the activation functions used in *MobileNet*, *Xception*, and the *CustomCNN*.

1.3.1 Activation Function in MobileNet

MobileNet uses the **ReLU** (Rectified Linear Unit) activation function. Each depthwise and pointwise convolution operation in MobileNet is followed by batch normalization (BN) and a ReLU activation. This is used for:

- Improving gradient flow: avoiding vanishing gradient problems that arise in deeper networks.
- Efficient computations: the ReLU function is simple to compute.

ReLU is defined as:

$$\text{ReLU}(x) = \max(0, x),$$

which sets all negative inputs to zero while passing positive values unchanged. This simplicity contributes to the computational efficiency of MobileNet.

⁴Access the paper at <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9329302&isnumber=9329288>.

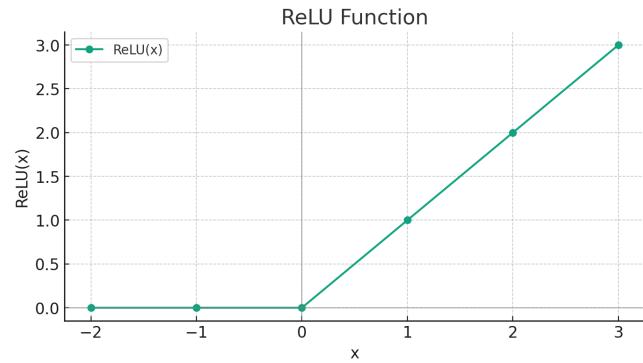


Figure 4: Example of RELU function in a 2D plot.

1.3.2 Activation Function in Xception

Xception also uses **ReLU**, but with a difference in its application compared to MobileNet. Unlike MobileNet, Xception does not apply a non-linearity (ReLU) between the depthwise and pointwise convolutions. This design decision is based on the hypothesis that depthwise convolution (which captures spatial correlations) and pointwise convolution (which captures cross-channel correlations) should be separately optimized without introducing additional non-linearity between them. However, after each convolution block, Xception still uses batch normalization followed by ReLU, similar to MobileNet. This difference contributes to faster convergence and better feature extraction.

1.3.3 Activation Function in CustomCNN

The **CustomCNN** defined by us also uses **ReLU** activation after the convolutional layer and batch normalization but before the max pooling operation. The sequence is:

Convolutional layer → Batch normalization → ReLU → MaxPooling

Additionally, the ReLU activation function is used after the first fully connected layer.

1.4 Deepfake Classification

Deepfake⁵ is a technology dedicated to creating highly realistic facial images and videos under specific conditions, which has significant application potential in fields such as entertainment, movie production, digital human creation, to name a few. With the advancements in Deep Learning, techniques primarily represented by **Variational Autoencoders** (or VAEs)

⁵Access the paper at <https://arxiv.org/abs/2403.17881>.

and **Generative Adversarial Networks** (or GANs) have achieved impressive generation results. More recently, the emergence of diffusion models with powerful generation capabilities has sparked a renewed wave of research. In addition to deepfake generation, corresponding detection technologies continuously evolve to regulate the potential misuse of deepfakes, such as for privacy invasion and phishing attacks.



Figure 5: Example of some Deepfakes starting to the original input image (Pristine column).

As we will see later, the authors of the OpenForensics Dataset use a well-defined pipeline to be able to make Deepfakes using an approach one that involves extracting **the latent vector of the facial identity**, modifying it with random values, and using GAN models to generate new faces. These synthesized faces are then adapted to the original pose through transformations, allowing an infinite number of fake identities to be generated without having to continuously retrain autoencoder-based models.

The **latent vector of facial identity** is basically a *point in a high-dimensionality space*, often represented mathematically as a \mathbb{R}^n element, that **encodes the distinctive features of a person's face**. In related literature, it is also called *latent code* or *embedding of identity*. In fact, this vector compactly captures the information needed to describe a face in such a way that, once processed by the generator (as in the case of OpenForensics within a GAN model), face images can be reconstructed or synthesized.

Modification with random values, at a more mathematical level, typically involves **adding a perturbation to the original latent vector**. For example, if $z \in \mathbb{R}^n$ represents the latent vector, a new vector z' can be generated by an operation of the type:

$$z' = z + \epsilon, \text{ where:} \quad (2)$$

ϵ is a random noise vector, often sampled from a **Gaussian distribution** (e.g., $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$). Therefore, this operation shifts the vector z in a random direction within the latent space, leading to the generation of a new face identity when the modified vector is passed through the GAN generator. This results in a variation of the identity that, while maintaining consistency in the realism of the data, represents a different face from the original one.

In the next sections, however, we will simply provide a more in-depth overview of what they do and the limitations of the most popular deepfakes generation techniques—*AutoEncoder*, *Face2Face*, *FaceSwap* and *Neural Textures*—which highlight the specific technical challenges and different levels of realism achieved by each improper approach of these technologies.

1.4.1 Creation techniques

As mentioned earlier, within this subsection the most famous deepfakes generation techniques will be analyzed by showing examples as well.

Auto-Encoder Based The first approach analyzed is that based on the **Auto-Encoder**⁶, which exploits a Neural Network consisting of two basic components: a *Encoder* and a *Decoder*. The Encoder is in charge of compressing the face image within a small latent space, trying to capture salient features (such as facial expressions and structures). Next, the Decoder will deal with the reconstruction of the face from this representation.

Typically, two AutoEncoders are trained in parallel for DeepFake generation—one for the **source face** and one for the **target face**—that share a common **encoder** to ensure a consistent representation of facial information, while each decoder is specific to reconstruct the peculiarities of each face.

⁶Access the paper at https://www.researchgate.net/publication/368495639_Deepfake_Technology_and_Current_Legal_Status_of_It.

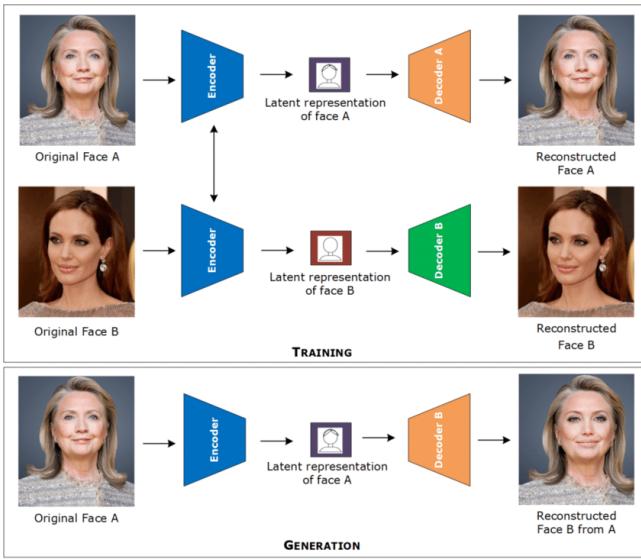


Figure 6: Workflow example of DeepFake generation with **AutoEncoders-based** approach.

Face2Face The second technique analyzed, however, is that of **Face2Face**⁷. This technique focuses on **real-time transfer of facial expressions**. The method starts by performing accurate tracking of the facial movements of a source actor, attempting to reconstruct a 3D model of the facial geometry of the source face itself. Then, the expressions that have been captured are mapped inside the selected target actor's face by means of *warping* and *blending* algorithms, which adapt textures and shapes in a smooth and natural manner.



Figure 7: Example of DeepFake generation using **Face2Face** approach.

FaceSwap As a third methodology, we are going to analyze instead **FaceSwap**. In fact, this is a technique that allows the exchange of faces present within videos or, as in the case of our analysis, photos. The

workflow begins with an initial automatic detection of faces, later followed by the identification of landmarks (or landmarks, which we will explain in more detail later) that go into defining what will be the structure of the face.

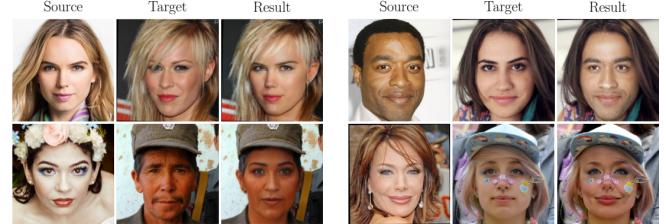


Figure 8: Example of DeepFake generation using **FaceSwap** approach.

These points are used to align and register the geometries of the faces involved, facilitating the transfer of the source face into the target face. After insertion, *blending* and *color correction* techniques are applied to ensure a smooth and natural transition, reducing the visibility of any artifacts.

Neural Textures The last, but by no means least important, approach we would analyze would be that of **Neural Textures**⁸. This approach uses neural networks to be able to learn high-fidelity representations of what facial textures are. In fact, in this method, the texture of the face is encoded within a latent space that preserves fine details, such as *light variations* and *natural imperfections*.

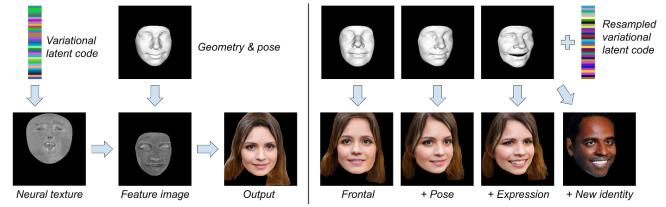


Figure 9: Example of DeepFake generation using **Neural Textures**-based approach.

During the neural rendering process, the network optimizes these textures to suit different lighting and perspective conditions, integrating them within a 3D model of the face; the result is an extremely realistic output.

⁷Access the paper at https://www.researchgate.net/publication/347531435_Detecting_DeepFake_FaceSwap_and_Face2Face_facial_forgeries_using_frequency_CNN.

⁸Access the paper at <https://www.sciencedirect.com/science/article/pii/S0952197624006080>.

1.4.2 Critical Detection and Analysis Methodologies

Differences in deepfake generation techniques have a direct impact on the detection strategies adopted. In particular, detection methodologies must take into account the following analyzed aspects.

Specific Artifacts In this case, we can see that techniques based on **AutoEncoder** can introduce anomalies within the reconstruction process, such as *compression artifacts* or *inconsistencies in detail*, which can be detected by *noise analysis* and *textures*.

Dynamic Transfer of Expressions Methods such as **Face2Face** are able to succeed in transferring expressions in real time, generating inconsistencies between facial motion and context, which can be exploited by detection algorithms based on *optical flow analysis* or *dynamic features*.

Complete Identity Exchange **FaceSwap**, on the other hand, or the generation of new identities through **GAN**, leads to the presence of discrepancies within biometric traits, which can be easily detected by conducting cross-comparisons by exploiting *facial recognition models*.

Blending and Color Adaptation Finally, the use of advanced *blending* techniques (e.g., **Poisson blending**) and *color adaptation* algorithms to succeed in integrating the synthesized face within the original image requires the adoption of methods that analyze *spatial coherence* and *edge continuity*.

2 Feature Extraction

2.1 Obtaining Bounding Boxes

In the OpenForensics dataset paper (available at the link⁹), the authors designed a complete *Face-Wise Multi-Task* annotation pipeline that, in addition to labeling each face with its manipulation category, also includes the **detection of its position** using *bounding boxes*. In practice, after selecting images with real faces (sourced from **Google Open Images**) and synthesizing fake faces using GAN (Generative Adversarial Network) models, a manipulability inspection module is applied to determine which faces can be manipulated.

Once the face is extracted, **68 keypoints** (or *land-*

⁹Access the documentation at <https://github.com/ltnghia/openforensics>

marks) are computed to precisely define the face contour; these landmarks, together with a mask obtained using techniques such as **Poisson blending** and chromatic adaptation algorithms, allow the generation of an accurate bounding box around the face. The final result — along with other annotations such as the segmentation mask and the manipulation boundary — is saved in **_poly.json* files in **COCO** style, thereby ensuring a rich and multi-purpose annotation for each face.

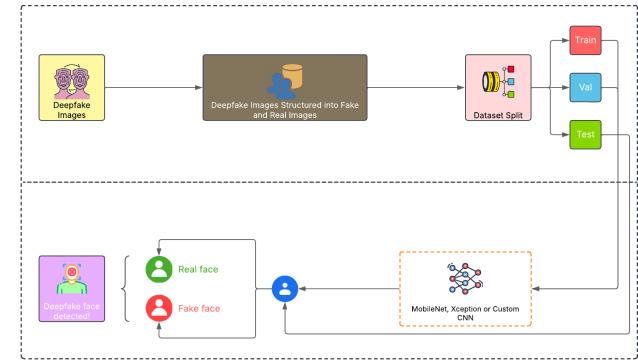


Figure 10: Deepfake Detection Project using the OpenForensics dataset.

2.1.1 Methodology for face extraction

To extract faces (both real and manipulated) from the OpenForensics dataset, we adopted an approach based on the COCO annotation format. In particular, the script is responsible for creating output folders, reading the various **_poly.json* files, processing the annotations, extracting the faces, and subsequently saving them.

2.1.2 Creating output folders

The script dynamically creates two folders, **real** and **fake**, within the output directory, where the extracted faces are saved.

```

1 real_dir = os.path.join(output_root, "real")
2 fake_dir = os.path.join(output_root, "fake")
3 os.makedirs(real_dir, exist_ok=True)
4 os.makedirs(fake_dir, exist_ok=True)
  
```

2.1.3 Reading the files (**_poly.json*)

The JSON file (for example, *Train_poly.json*, *Val_poly.json*, etc.) that follows the **COCO** structure (previously described) is loaded. From this file, a map is created that associates each *image_id* with its corresponding *file_name*.

```

1 with open(json_path, 'r') as f:
2     data = json.load(f)
3
4 # Map image_id -> file_name
5 images_info = {img["id"]: img["file_name"]
6                 for img in data["images"]}

```

2.1.4 Annotations processing

For each annotation, the script extracts the bounding box, specified as [x, y, w, h], and calculates the coordinates of the rectangle to be used for cropping. It is important to note that the bounding box coordinates, used for face extraction, were provided directly by the dataset authors through the *_poly.json files (described above). These annotations ensure an accurate representation of the regions containing the faces, facilitating the correct extraction of both manipulated and real faces.

```

1 for ann in data["annotations"]:
2     image_id = ann["image_id"]
3     category_id = ann["category_id"]
4     x, y, w, h = ann["bbox"]
5
6     # Calculate the rectangle coordinates
7     x1, y1 = int(x), int(y)
8     x2, y2 = int(x + w), int(y + h)

```

2.1.5 Extracting and saving the face

The image is read using **OpenCV**, and using the **bounding box coordinates**, the face is extracted. Depending on the category (0 for Real, 1 for Fake), the face is saved in the appropriate folder with a filename that includes a progressive index to avoid conflicts.

```

1 # Check if the image exists and read it
2 if image_id in images_info:
3     file_name = os.path.basename(images_in_
4         ↪ fo[image_id])
4     img_path = os.path.join(images_root,
5         ↪ image_subdir, file_name)
5
6     if os.path.isfile(img_path):
7         image = cv2.imread(img_path)
8         if image is not None:
9             cropped_face = image[y1:y2,
10                ↪ x1:x2]
11             label_str = "real" if
12                 ↪ category_id == 0 else
13                 ↪ "fake"

```

```

11             face_filename =
12                 ↪ f"{os.path.splitext(file_n_
13                 ↪ ame)[0]}_{face_count}.jpg"
13             save_path =
14                 ↪ os.path.join(real_dir if
15                     ↪ category_id == 0 else
16                     ↪ fake_dir, face_filename)
17             cv2.imwrite(save_path,
18                 ↪ cropped_face)
19             face_count += 1

```

2.2 Quality validation of annotations

The quality of the annotations in the **OpenForensics** dataset was not evaluated using a single automated metric, but through a **mixed approach that combines manual checks and comparisons with reference annotations**. In the paper, an initial manual visual review is presented, followed by a comparison with **ground truth** associated with **user evaluation studies**.

2.2.1 Manual Visual Review

The authors performed manual checks on the *bounding boxes* and the *segmentation masks*, verifying that the annotations accurately reflected the regions containing the faces. This process allowed them to identify and correct any errors, particularly in cases of partially occluded faces or unconventional poses.

2.2.2 Comparison with Ground Truth and User Evaluation Studies

Additionally, a user study was conducted to evaluate the overall visual quality of the dataset. Participants provided a score (or **Mean Opinion Score**, MOS) on the realistic perception of the annotations, confirming that the bounding boxes, masks, and landmarks were of high quality and well aligned with the real or manipulated faces. This qualitative evaluation complemented the manual review, ensuring the reliability of the annotations provided in the *_poly.json files.

3 Metadata Analysis

3.1 General Description of the Dataset

The dataset **OpenForensics** is one of the largest collections of real and manipulated images, specifically designed for the detection and segmentation of multi-face deepfakes in real-world environments (*In-The-Wild*).

Its primary goal is to overcome the limitations of existing datasets, which often contain images with *uniform backgrounds, a single face, and scenarios that do not fully represent real-world complexity*.

The dataset follows the **COCO** (*Common Objects in Context*) format, a widely adopted standard in Computer Vision for object detection, segmentation, and keypoint recognition. The dataset is structured as follows:

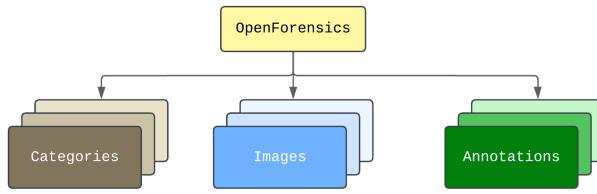


Figure 11: OpenForensics dataset schema (by following the COCO format).

More information on the COCO format and the OpenForensics dataset can be found at the following link¹⁰.

3.2 Main features

3.2.1 Scalability and Size

The **OpenForensics** dataset consists of **115,325 unrestricted images**, containing approximately **334,000 faces**. The images are sourced from various platforms, such as **Google Open Images**, ensuring high variability in scene context, lighting conditions, and resolution. This diversity is crucial for training robust models that can operate in real and complex environments.

3.2.2 Rich and Multi-Task Annotations

Each face in the images is annotated with highly detailed information, including:

- **Forgery Category** (indicating whether the face is real or manipulated);
- **Bounding Box** (coordinates of the bounding rectangle);
- **Segmentation Mask**;
- **Manipulation Contours**;

- **Facial Landmarks**.

These annotations not only support classification tasks (such as deepfake detection) but also enable face localization and segmentation.

3.2.3 Diversity of Scenarios

The dataset includes both **indoor** and **outdoor** environments, with faces of different sizes, orientations, and occlusions. The manipulated faces are generated using advanced techniques such as **GAN** (Generative Adversarial Networks), **Poisson Blending** or **Color Adaptation**. These methods ensure high-resolution, realistic images that blend seamlessly into the original contexts.

3.2.4 Augmentations with the Test-Challenge split

To replicate real-world conditions, a subset called **Test-Challenge** was created. This subset includes various perturbations and transformations such as **Color modifications, corruptions and distortions** or **external effects** (such as fog, snow or rain). These augmentations increase image variability, testing the robustness of detection and segmentation methods.

3.3 Metadata Files Structure (*.poly.json)

The ***.poly.json** file for each dataset split is structured as follows:

3.3.1 Categories

The **categories** section defines the class categories:

```

1  "categories": [
2      {
3          "id": 0,
4          "name": "Real"
5      },
6      {
7          "id": 1,
8          "name": "Fake"
9      }
10 ]
    
```

¹⁰COCO Dataset Format: <https://www.immersivelimit.com/tutorials/create-coco-annotations-from-scratch/#coco-dataset-format>

3.3.2 Images

The `images` section is an array of objects, each representing an image and containing the following attributes:

```

1 "images": [
2   {
3     "id": 0,
4     "file_name": "Images/Test-Challeng_"
5       ↪ e/f80d795aa7.jpg",
6     "width": 1024,
7     "height": 768
8   },
9   {
10    "id": 1,
11    "file_name": "Images/Test-Challeng_"
12      ↪ e/d59690204c.jpg",
13    "width": 1024,
14    "height": 683
15  }
16  // ...other images...
17 ]

```

3.3.3 Annotations

The `annotations` section is an array of annotations, where each annotation corresponds to a face in one of the images, and includes:

- **id**: annotation identifier;
- **image_id**: reference to the associated image;
- **iscrowd**: flag to indicate whether the annotation represents a group of faces;
- **area**: annotation area;
- **category_id**: category identifier (0 for “Real,” 1 for “Fake”);
- **bbox**: represents the bounding box, specified as $[x_{min}, y_{min}, w, h]$ (or as $[x_{min}, y_{min}, x_{max}, y_{max}]$ depending on the convention);
- **segmentation**: the segmentation mask, usually represented as an array of coordinates defining a polygon;

```

1 {
2   "annotations": [
3     {
4       "id": 795,
5       "image_id": 364,
6       "iscrowd": 0,
7     }
8   ]
9 }

```

```

7           "area": 127378,
8           "category_id": 1,
9           "bbox": [244, 112, 354, 558],
10          "segmentation": [
11            [
12              [
13                [
14                  [
15                    [
16                      [
17                        [
18                          [
19                            [
20                              ...
21                            ]
22                          ]
23                        ]
24                      ]
25                    ]
26                  ]
27                ]
28              ]
29            ]
30          ]
31        ]
32      ]
33    ]
34  ]
35  // ...other annotations...
36 ]

```

4 Fine Tuning of MobileNet and Xception

4.1 General Description of Fine-Tuning

Fine-tuning is a transfer learning technique used to adapt a pre-trained neural network—originally trained on large generic datasets—to a specific, specialized task. Instead of training a model from scratch, we adjust the weights for our DeepFake Detection task.

By using this approach, we reduce training time, and the network guarantees better performance with a smaller set of training data.

4.2 Models Chosen for Fine-Tuning

We selected two state-of-the-art pre-trained CNN models for fine-tuning are **MobileNet-v2** (via PyTorch torchvision library) and **Xception** (via Timm library). Both models are used for extracting high-level visual features, which we then adapt for our binary classification task (Real vs. Fake images).

4.3 Fine-Tuning Implementation Details

Here we describe the precise process of fine-tuning performed for each model.

4.3.1 MobileNet-v2

For MobileNet-v2, we replace the original classification layer to adapt the network to our binary classification task.

```

1 model = models.mobilenet_v2(pretrained=True)
2 in_features =
    ↪ model.classifier[1].in_features
3 model.classifier[1] =
    ↪ nn.Linear(in_features, 2)
```

- The original **fully connected** layer (`model.classifier[1]`) is replaced by a new linear layer designed specifically for our two-class classification.
- This allows the model to specialize in DeepFake-specific characteristics and give us a binary output.

4.3.2 Xception

Similarly, for the Xception model, we apply the following modifications:

```

1 model = timm.create_model('xception',
    ↪ pretrained=True)
2 model.fc = nn.Linear(model.fc.in_features,
    ↪ 2)
```

- The final **fully connected** layer (`fc`) is replaced by a new linear layer tailored for binary classification.
- Xception, known for its excellent performance in image-based tasks, is optimized for distinguishing real and fake images.

4.4 Training Procedure

After modifying the final layers, the training process for fine-tuning includes the steps that will be discussed in the next sections.

4.4.1 Data Preparation

We create DataLoaders for training and validation datasets:

```

1 train_loader = create_dataloader("processes",
    ↪ d_data/train_cropped", batch_size=32,
    ↪ shuffle=True)
2 val_loader = create_dataloader("processes",
    ↪ d_data/val_cropped", batch_size=32,
    ↪ shuffle=False)
```

4.4.2 Hyperparameters Definition

The following hyperparameters are set for optimal performance:

```

1 BATCH_SIZE = 32
2 EPOCHS = 15
3 LEARNING_RATE = 1e-4
4 criterion = nn.CrossEntropyLoss()
5 optimizer = optim.Adam(model.parameters(),
    ↪ lr=LEARNING_RATE)
6 DEVICE = torch.device("cuda" if
    ↪ torch.cuda.is_available() else "mps" if
    ↪ torch.backends.mps.is_available() else
    ↪ "cpu")
```

4.4.3 Fine-Tuning Training Loop

The training loop includes standard PyTorch procedures:

```

1 for epoch in range(EPOCHS):
2     model.train()
3     running_loss = 0.0
4     for images, labels in train_loader:
5         images, labels = images.to(DEVICE),
        ↪ labels.to(DEVICE)
6
7         optimizer.zero_grad()
8         outputs = model(images)
9         loss = criterion(outputs, labels)
10        loss.backward()
11        optimizer.step()
12
13        running_loss += loss.item()
14
15    avg_train_loss = running_loss /
        ↪ len(train_loader)
```

Table 1: Architectures for our Custom CNN

Layer	Input Shape	Operation	Output Shape
First Convolutional Block	(3, 224, 224)	Conv2d(in=3, out=32, kernel=3, pad=1) BatchNorm2d(32) ReLU MaxPool2d(kernel=2, stride=2) Conv2d(in=32, out=64, kernel=3, pad=1) BatchNorm2d(64)	(32, 112, 112)
Second Convolutional Block	(32, 112, 112)	ReLU MaxPool2d(kernel=2, stride=2) Conv2d(in=64, out=128, kernel=3, pad=1) BatchNorm2d(128)	(64, 56, 56)
Third Convolutional Block	(64, 56, 56)	ReLU MaxPool2d(kernel=2, stride=2) Conv2d(in=128, out=256, kernel=3, pad=1) BatchNorm2d(256)	(128, 28, 28)
Flatten	(128, 28, 28)	Flatten	(1, 100352)
Fully Connected (FC1 + ReLU + Dropout)	(1, 100352)	Linear(100352 → 512) ReLU Dropout(p=0.5)	(1, 512)
Output (FC2)	(1, 512)	Linear(512 → 2)	(1, 2)

4.4.4 What happens during training specifically?

The **Forward pass** provides Input images pass through pre-trained convolutional layers to extract general and specific features. The new fully connected layers transform these features into predictions.

The **Loss Calculation** provides predictions (outputs) are compared to ground-truth labels (labels) using Cross-Entropy Loss.

For the phase of **Backward Pass and Weight Update** the gradients from the loss function propagate backward through the entire network, updating not only the newly added final layers but also fine-tuning previously learned layers to adapt better to DeepFake-specific features.

4.5 Saving the Fine-Tuned Models

After training, the models are saved for future inference and evaluation:

```

1 model_save_path =
  ↵ f"models/{model_name}_deepfake.pth"
2 torch.save(model.state_dict(),
  ↵ model_save_path)

```

4.6 Benefits and Motivations of Our Approach

We adopted a **complete fine-tuning** strategy for our pre-trained models (**MobileNet-v2** and **Xception**) because it proves to be the most effective approach for our specific task of **DeepFake detection**.

This method allows every layer of the network to fully adapt to the unique characteristics of the OpenForensics dataset, which includes diverse and complex images containing visual anomalies typical of digital manipulations.

Compared to **retraining the last layers** (feature extraction), complete fine-tuning enables the network to capture precise features, significantly enhancing generalization and robustness.

Furthermore, having a large number of training images allows us to effectively update all model parameters, maximizing the system's accuracy in distinguishing between real and manipulated images.

5 Building a Network from Scratch

5.1 General description

In this section, we explain the structure and rationale behind the custom CNN developed for DeepFake image detection through binary classification, distinguishing between true and false images.

In addition, we will give an initial definition of some operations applied within the Convolutional Layers.

5.2 Network architecture

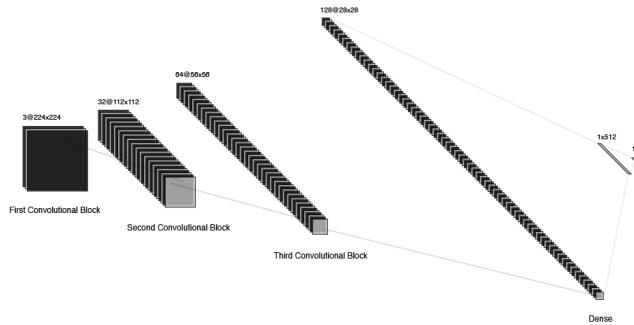


Figure 12: CNN Architecture.

5.2.1 First Convolutional Block

```

1 self.conv1 = nn.Conv2d(in_channels=3,
    ↵ out_channels=32, kernel_size=3,
    ↵ padding=1)
2 self.bn1 = nn.BatchNorm2d(32)
3 self.pool1 = nn.MaxPool2d(kernel_size=2,
    ↵ stride=2)

```

Convolutional Layer (Conv2d), the initial layer, is tasked with capturing and identifying low-level features, such as edges and basic colors, maintaining the spatial dimensions of the image thanks to the chosen padding (set to 1).

In this first layer, we set the number of filters (`out_channels = 32`) to capture diverse image characteristics. Also, we have:

- **Input channels:** 3 (RGB)
- **Output channels:** 32
- **Kernel size:** 3
- **Padding:** 1

Regarding the choice of kernel, we decided to use a 3×3 size: within our specific application domain, this is well-suited for detecting small details crucial for the task of DeepFake detection. We also experimented with a larger 7×7 kernel, but results did not show any improvement.

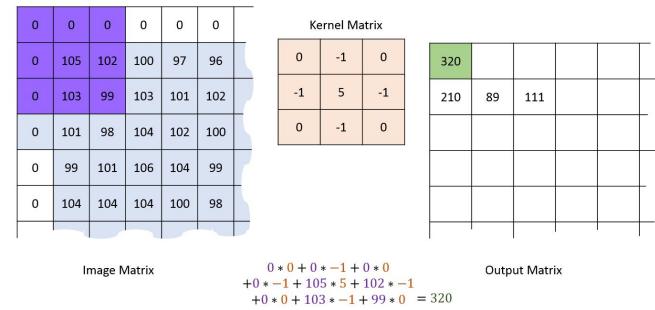


Figure 13: Example of **2D Convolution** operation with Kernel = 3 and Padding = 1.

Batch Normalization (BatchNorm2d) stabilizes and accelerates training by reducing internal covariate shift, making the model more robust against variations in input distributions. Sudden weight changes can occur in the early layers; batch normalization normalizes each mini-batch, followed by scaling and shifting.

Max Pooling (MaxPool2d) reduces the dimensionality of the image representation, decreasing computational load and capturing dominant features.

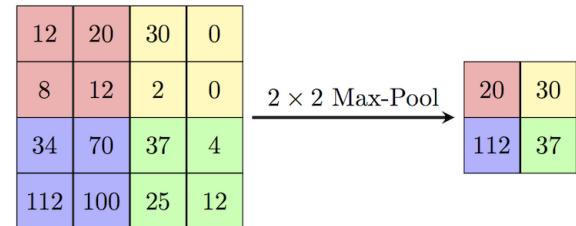


Figure 14: Example of **Max Pooling** operation with Kernel = 2.

Using a 2×2 kernel halves the spatial dimensions of the feature map, reducing computations without losing essential information. The maximum value is selected for each sub-matrix, ensuring a significant spatial reduction while retaining fundamental details. Also, we have:

- **Kernel size:** 2
- **Stride:** 2

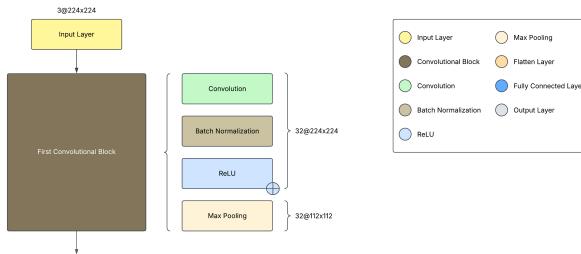


Figure 15: Visualization of First Convolutional Block.

5.2.2 Second Convolutional Block

```

1 self.conv2 = nn.Conv2d(in_channels=32,
    ↪ out_channels=64, kernel_size=3,
    ↪ padding=1)
2 self.bn2 = nn.BatchNorm2d(64)
3 self.pool2 = nn.MaxPool2d(kernel_size=2,
    ↪ stride=2)

```

Convolutional Layer (Conv2d) Increases the network depth, allowing the capture of more complex and structured image features.

Gradually increasing the number of filters enables the network to learn increasingly sophisticated representations without excessive computational overhead. Also, we have:

- **Input channels:** 32
- **Output channels:** 64
- **Kernel size:** 3
- **Padding:** 1

Batch Normalization (BatchNorm2d) As in the first block, it contributes to training stability and improves generalization.

Max Pooling (MaxPool2d) Maintains the dimensionality reduction approach.

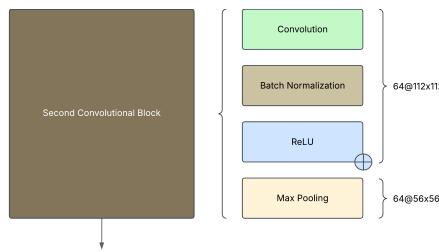


Figure 16: Visualization of Second Convolutional Block.

5.2.3 Third Convolutional Block

```

1 self.conv3 = nn.Conv2d(in_channels=64,
    ↪ out_channels=128, kernel_size=3,
    ↪ padding=1)
2 self.bn3 = nn.BatchNorm2d(128)
3 self.pool3 = nn.MaxPool2d(kernel_size=2,
    ↪ stride=2)

```

Convolutional Layer (Conv2d) This block is crucial for capturing high-level features, allowing the network to identify distinctive features relevant to the task, such as anomalies in textures and color transitions. Also, we have:

- **Input channels:** 64
- **Output channels:** 128
- **Kernel size:** 3
- **Padding:** 1

Batch Normalization (BatchNorm2d) Promotes effective regularization and more stable learning.

Max Pooling (MaxPool2d) It concludes the convolutional phase of the network by further reducing the spatial dimensions of the feature maps, which is critical for minimizing overfitting.

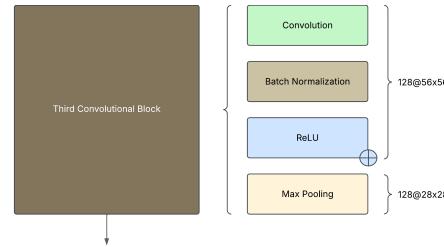


Figure 17: Visualization of Third Convolutional Block.

5.2.4 Fully Connected Layers

```

1 self.fc1 = nn.Linear(128 * 28 * 28, 512)
2 self.fc2 = nn.Linear(512, num_classes)

```

First Fully Connected Layer (FC1) It transforms spatial representations into a compact vector useful for final classification. The high output size (512 neurons) was chosen to allow the network to learn rich and informative representations. The input is a vector of size 128×28 , i.e., 128 feature maps each of size 28×28 , which is “flattened” to obtain a smaller vector (512 elements).

Also, we have:

- **Input:** $128 \times 28 \times 28$
- **Output:** 512

Activation Function: ReLU It introduces nonlinearity, allowing the network to learn more complex relationships than just linear functions.

Dropout (0.5) It reduces overfitting by randomly deactivating some neurons during training, significantly improving model generalization. In our case, with a value of 0.5, half of the neurons are deactivated, preventing overdependence on specific neurons. Also, we have:

```
1 self.dropout = nn.Dropout(0.5)
```

Second Fully Connected Layer (FC2) This final layer performs the classification between real and deepfake images, providing the direct output needed for decision making. Also, we have:

- **Input:** 512
- **Output:** 2 (classificazione binaria)

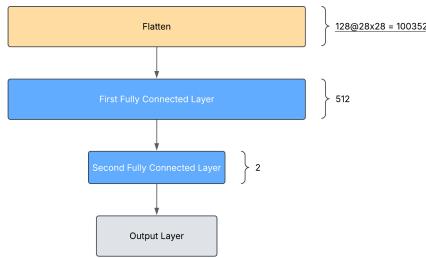


Figure 18: Visualization of Fully Connected Layers.

5.2.5 Forward Function

```

1 def forward(self, x):
2     x = self.pool1(F.relu(self.bn1(self.co_
3         ↵ nv1(x))))
4     x = self.pool2(F.relu(self.bn2(self.co_
3         ↵ nv2(x))))
5     x = self.pool3(F.relu(self.bn3(self.co_
3         ↵ nv3(x))))
6
7     x = torch.flatten(x, start_dim=1)
8     x = F.relu(self.fc1(x))
9     x = self.dropout(x)
10    x = self.fc2(x)
11
12    return x
    
```

The function `forward` defines the sequence in which data traverses the network layers:

1. Input → Conv1 → BN1 → ReLU → Pool1
2. → Conv2 → BN2 → ReLU → Pool2
3. → Conv3 → BN3 → ReLU → Pool3
4. → Flatten → FC1 → ReLU → Dropout
5. → FC2 → Output

5.3 Technology and motivation behind choices

PyTorch was chosen for its flexibility and ease of CNN implementation, as well as for its widespread use and community support.

Batch Normalization significantly improves convergence speed, makes the network more robust to input variations, and allows higher learning rates.

Max Pooling reduces spatial dimensions while preserving relevant features, reducing the risk of overfitting.

Dropout is essential for network regularization, ensuring good generalization, particularly important for deepfake detection, where recognizing subtle and varied characteristics is crucial.

6 Analysis and Experimental Results

For all experiments we used three different networks, previously mentioned, which are: MobileNet, Xception and our own CNN Custom.

6.1 MobileNet and Xception

The fine-tuning of **Xception**, with our HyperParameters and a M4 Pro Macbook, took 5 hours and 30 minutes. Instead the fine-tuning of **MobileNet-v2**, with the same environment, took 1 hour and 40 minutes in total.

6.1.1 Evaluation results on Test-Dev

On the **Test-Dev** dataset, both fine-tuned models achieved very good results, showing their strong ability to distinguish between real and manipulated images in relatively controlled and clean conditions. Also, we have obtained:

Table 2: Results on Test-Dev

Model	Accuracy	Precision	Recall	F1-Score
MobileNet-v2	0.9944	0.9938	0.9965	0.9951
Xception	0.9966	0.9966	0.9975	0.9970

Xception achieved slightly superior performance compared to MobileNet-v2, recording almost perfect accuracy (99.66%), precision (99.66%), recall (99.75%), and F1-score (99.70%). These metrics confirm Xception’s capability to capture visual anomalies typical of DeepFake manipulations.

MobileNet-v2, although slightly lower, also performed very well, with an accuracy of 99.44%, precision of 99.38%, recall of 99.65%, and an F1-score of 99.51%. This highlights the model’s effectiveness and indicates that it is also a very good option for our task.

6.1.2 Evaluation results on Test-Challenge

The **Test-Challenge** dataset introduces a realistic and complex scenario, including augmentations, noise, and visual distortions.

Table 3: Results on Test-Challenge

Model	Accuracy	Precision	Recall	F1-Score
MobileNet-v2	0.8364	0.9670	0.7442	0.8411
Xception	0.8023	0.9725	0.6795	0.8000

MobileNet-v2 has better performances than Xception, achieving higher accuracy (83.64% vs. 80.23%), recall (74.42% vs. 67.95%), and F1-score (84.11% vs. 80.00%). Although precision was slightly lower than Xception (96.70% vs. 97.25%), MobileNet-v2’s recall shows that it is more successful at correctly identifying deepfake images under realistic, noisy conditions.

Xception, despite having slightly higher precision, demonstrated difficulty in recognizing manipulated images accurately. Its lower recall score (67.95%) indicates that it missed several challenging deepfake samples, affecting overall accuracy and F1-score.

6.2 CNN from Scratch

We tested the network on two different datasets, **Test-Dev** and **Test-Challenge**. For training the network, we used the following hyperparameters and loss function:

```

1 BATCH_SIZE = 32
2 EPOCHS = 10
3 LEARNING_RATE = 1e-4
4 criterion = nn.CrossEntropyLoss()
5 optimizer = optim.Adam(model.parameters(),
    ↳ lr=LEARNING_RATE)

```

We experimented with implementing **data augmentation** and **early stopping**, but the evaluation results were unsatisfactory—worse than the baseline network—so we decided not to implement these strategies in the final training.

6.2.1 Evaluation results on Test-Dev

The results obtained on **Test-Dev** are as follows:

Table 4: Results on Test-Challenge (Custom CNN)

Model	Accuracy	Precision	Recall	F1-Score
Custom CNN	0.9716	0.9608	0.9913	0.9758

6.2.2 Evaluation results on Test-Challenge

The results obtained on **Test-Challenge** are as follows:

Table 5: Results on Test-Challenge (Custom CNN)

Model	Accuracy	Precision	Recall	F1-Score
Custom CNN	0.8278	0.8918	0.8014	0.8441

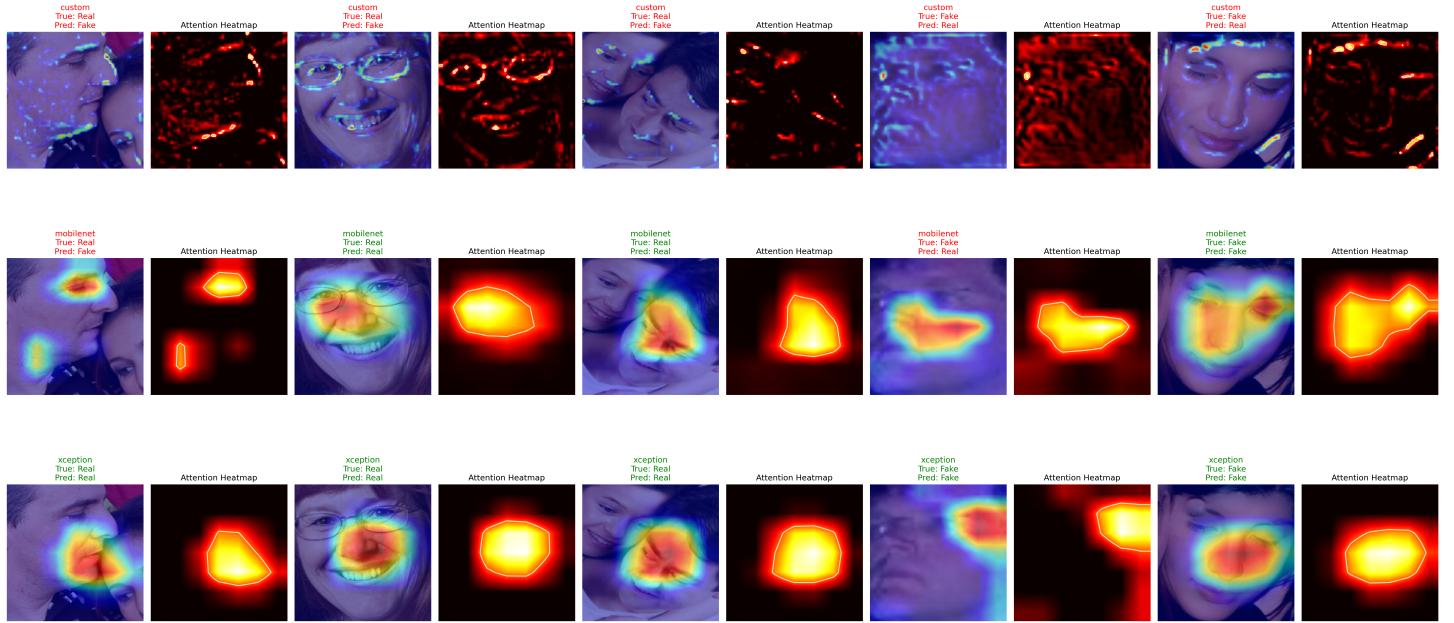


Figure 19: Grad-CAM and Attention HeatMap Matrix for a subset of misclassifications (Custom CNN, MobileNet and Xception).

6.2.3 Grad-CAM Visualization

We also implemented a Grad-CAM visualization to better understand the model's decision-making process, as shown below.

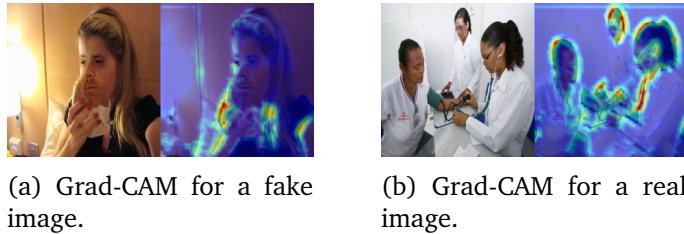


Figure 20: Two examples of Grad-CAM. On the left (a) the model highlights faces, on the right (b) the network distributes attention more evenly.

As shown in the Figure 20, in the fake image (left), the network primarily focuses on faces and edges, indicating potential anomalies. Conversely, in the real image (right), the attention is more evenly distributed, suggesting that the network recognizes it as genuine.

Grad-CAM and Attention HeatMap Matrix for misclassified images wanting to provide a more accurate analysis on the areas where the model is most concentrated, as shown in the Figure 19, we made a more explanatory visualization than the immediate previous one shown to hint at the analysis itself.

We start by first analyzing the results given by the **Grad-CAM**, in fact as we can see the network focuses a lot on sharp contours (such as glasses, hair, etc.) relying on low scale differences to be able to distinguish an image between real and fake: this, in fact, could be given by the **use of a fairly small kernel**, as in our case 3×3 , and the low depth of the network itself.

One idea to improve the network from this point of view, in fact, might be to **change the size of the kernel** from 3×3 to 5×5 and see how the network behaves; by doing this the network should not only focus on sharp edges but also on larger details (given, precisely, by the larger size of the kernel itself). We do not recommend a larger kernel (such as 7×7) because this would degrade the performance of the network itself.

Regarding, on the other hand, the **Attention HeatMap** what we can say and affirm is that it dwells on very sharp details (as can also be seen from the very intense red areas present on, for example, eyes, mouth or glasses): this means that, for the model, these parts are crucial for the classification between Real and Fake.

In addition, we can also compare the three networks used for this analysis: **MobileNet**, **Xception** and the

Custom CNN. As can be seen, precisely, from both the Grad-CAM and Attention HeatMap, we can state how **Custom CNN** turns out to be more granular (dwelling more on the fine detail points) in contrast to the other two networks analyzed. We can in addition also note that **Xception**, given its greater knowledge and complexity, goes to focus on the *discriminating* point for understanding/classifying an image classified as *Fake*.

The segmentation of the heatmap by the custom network potentially comes from the fact that it was designed solely and exclusively for this type of task, going to focus on more specific areas of the face. Pre-trained networks, on the other hand, have more *smooth* and concentrated heat-maps since the model already knows where to look and does not necessarily have to scan the entire image to figure out which areas are actually the discriminating ones. Prior knowledge of **Xception** and **MobileNet**, evidently, makes the two networks more accurate after fine-tuning is performed.

Finally, we can also state that the examples shown for both **Grad-CAM** and **Attention HeatMap** would turn out to be **similar**, meaning that the **features important for classification** (reported by Grad-CAM) and the **attention mechanisms** (reported by Attention HeatMap) **coincide**.

6.2.4 Training Loss and Testing Accuracy

Throughout the training process, we monitored both the loss on the training set and the classification accuracy on the validation set at each epoch.

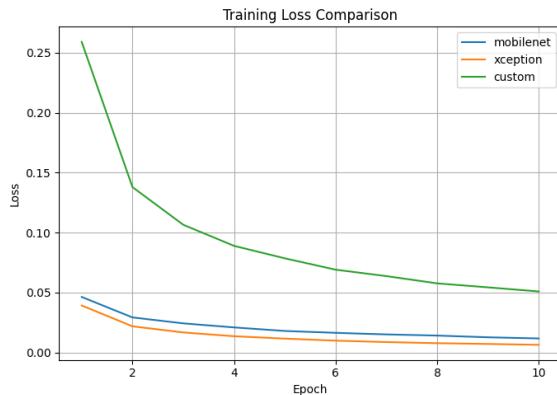


Figure 21: Plot of the Training Loss.

As shown in Figure 21, the **MobileNet** and **Xception** models - both initialized with ImageNet-pretrained weights - demonstrated a rapid decrease in training loss, converging to near-zero values by the end of training. Instead, the **Custom CNN**, which was trained from scratch, began with a notably higher loss but improved in each epoch. As we can see in the plot, we have trained each model with 10 epochs.

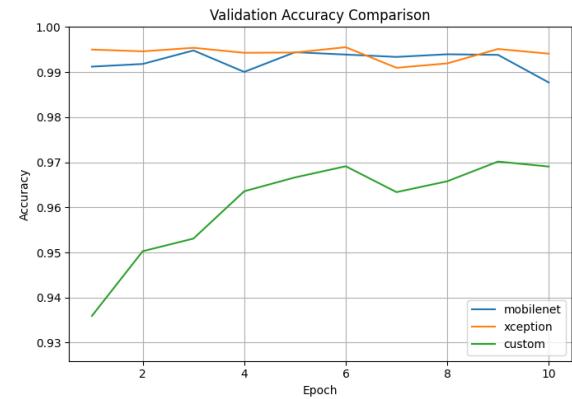


Figure 22: Plot of the Validation Loss.

Interestingly, despite these differences in the first epochs training loss curves, all three models achieved high accuracy on the validation set (Figure 22), with values nearing or exceeding 0.95 by the final epoch. This suggests that the Custom CNN, although slower to reduce its loss, was still able to learn discriminative features effectively enough to perform well on the validation data.

6.2.5 Accuracy as the image resolution changes

After performing a more in-depth analysis, trying to make a graph showing the value of Accuracy as the image resolution changes, we can see within it a rather obvious trend: in fact, as the input resolution increases (from 32px to 254px), the accuracy of the model increases significantly, going from values around 0.63 – 0.70 to over 0.98. Thereafter, between 256px and 1024px, the accuracy remains almost stable (in fact, we can see an oscillation between 0.97 – 0.98, without having any further significant increases).

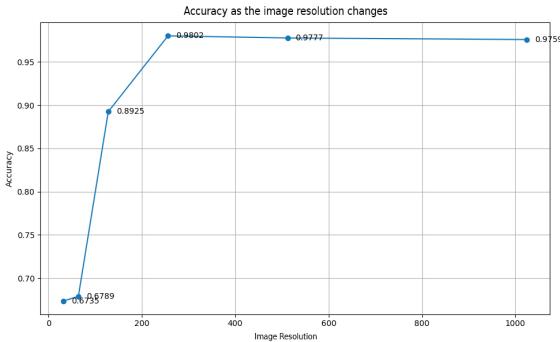


Figure 23: Plot of the Accuracy as the Image Resolution changes.

In fact, looking at Figure 23, we can analyze the graph more and notice as mentioned earlier a drastic increase between low and medium resolution (at $32px$ and $64px$); in fact, the model loses crucial information to be able to correctly distinguish real images from fake ones. With so few pixels, many facial features are compressed or lost.

Going up to $128px$ (and obtaining an Accuracy value of over 0.89) we can see that the network starts to have sufficient resolution to be able to capture all those discriminating patterns for the correct classification of Deepfakes.

In contrast, however, by going from $256px \rightarrow 512px \rightarrow 1024px$, we can see that the Accuracy remains with very high (and mostly constant) values, but does not improve significantly. In fact, this behavior occurs when the model reaches a certain threshold of useful information present within the Input data, and then the network fails to benefit more pixels. As in our case, in fact, sometimes a slight drop can occur if the network has not been trained specifically for very large images: in our case, the network was trained with images of $224 \times 224px$. In fact, again casting an eye over the graph present in Figure 23, we notice how as we approach the size used to train the network the Accuracy performance of the model itself improves.

For completeness of the Analysis performed itself, we also dwell on the computational trade-off of the network itself: in fact, increasing the resolution of the images also increases the computational cost (in terms of inference time and memory).

7 Conclusions and Future Works

When evaluated on the final test splits, the models maintained similarly strong performance. Specifically, on the Test-Dev set, **Xception** slightly outperformed **MobileNet**, whereas on the more challenging Test-Challenge set, **MobileNet** proved more robust than **Xception** under noisy or distorted conditions.

The **Custom CNN**, despite being fully trained from scratch, still achieved competitive results in both scenarios. Overall, these observations highlight the advantage of transfer learning for faster convergence and potentially higher performance on standard test data.

We can also confirm that a carefully designed CNN can come close to the results obtained by pre-trained models, suggesting that with a larger dataset and more computational resources the effectiveness of the custom model can be significantly improved.