

## **Deepfake Detection Project - Q&A**

**Q:** What is the main objective of your deepfake detection project?

**A:** The main objective is to develop a reliable system that can automatically detect manipulated or synthetic videos (deepfakes) by analyzing the frames of the video to identify inconsistencies or signs of tampering.

**Q:** Why did you choose this project? Or why is deepfake detection important?

**A:** Deepfake technology is rapidly evolving and poses serious threats such as misinformation, identity theft, and cyber fraud. This project aims to provide tools to identify fake videos and enhance cybersecurity awareness.

**Q:** What problem does your project solve?

**A:** It solves the problem of verifying video authenticity and preventing the spread of manipulated videos that can mislead viewers or cause harm.

**Q:** Is this project real-time or offline?

**A:** The current version is offline and processes videos uploaded by users.

**Q:** What kind of deepfakes does it detect?

**A:** It focuses on visual deepfakes in videos, especially facial manipulations.

**Q:** What is the title of the project?

**A:** The project is titled "Detecting Digital Deception: AI Solutions for Deepfake Detection."

**Q:** What dataset did you use for your project?

**A:** We used publicly available deepfake video datasets such as Celeb-DF and YouTube-Real, which contain labeled real and fake videos for training and evaluation.

**Q:** How many samples or records does the dataset contain?

**A:** The dataset contains thousands of video samples, with each video consisting of multiple frames extracted for analysis.

**Q:** Are all videos in the dataset labeled?

**A:** Yes, videos are labeled as real or fake. Fake as 0 and Real as 1

**Q:** How did you handle missing or inconsistent data in the dataset?

**A:** We performed data cleaning by removing corrupted frames or videos and ensuring each video had valid frame sequences for consistent model input.

**Q:** What preprocessing steps did you apply to the data?

**A:** We extracted frames from each video at a fixed frame rate, resized images to a standard resolution, and normalized pixel values. We also filtered out invalid or noisy frames.

**Q:** What size were the frames resized to?

**A:** All frames were resized to 224x224 pixels.

**Q:** How many frames did you extract per video?

**A:** We extracted 10–15 frames per video on average.

**Q:** How did you extract frames from videos for analysis?

**A:** Using OpenCV, frames were extracted sequentially from videos, saved in directories labeled 'real' or 'fake' for supervised learning.

**Q:** Did you remove noisy frames?

**A:** Yes, blurry or blank frames were excluded during preprocessing.

**Q:** Which format were frames saved in?

**A:** Frames were saved as JPEG images.

**Q:** Was feature extraction done before or during training?

**A:** Feature extraction was done before training using ResNet50.

**Q:** Did you perform data augmentation?

**A:** yes, augmentation was used in this version.

**Q:** How did preprocessing affect the model's performance?

**A:** Proper preprocessing improved model accuracy by providing consistent input and reducing noise from corrupted frames.

**Q:** What machine learning or deep learning model did you use?

**A:** We used a convolutional neural network (CNN) based model, specifically ResNet50, fine-tuned to classify frames as real or fake.

**Q:** Why did you choose this particular model?

**A:** ResNet50 is known for its deep architecture and residual learning capabilities, which help in extracting detailed features necessary for distinguishing subtle video manipulations.

**Q:** Did you use transfer learning or pretrained models?

**A:** Yes, we used transfer learning by initializing ResNet50 with weights pretrained on ImageNet and then fine-tuned on our deepfake dataset.

**Q:** What model is used for detection?

**A:** We used the ResNet50 model.

**Q:** Is ResNet50 trained from scratch?

**A:** No, we used the pre-trained ResNet50 with ImageNet weights.

**Q:** What is the input to the ResNet50 model?

**A:** Preprocessed video frames of size 224x224 pixels.

**Q:** What is the output of ResNet50?

**A:** A 2048-dimensional feature vector for each frame.

**Q:** What classifier is used after feature extraction?

**A:** A dense neural network is used for classification.

**Q:** How many layers does the classifier have?

**A:** The classifier has two dense layers and a final output layer.

**Q:** What activation function is used?

**A:** ReLU is used in hidden layers, and sigmoid in the output layer.

**Q:** Why was ResNet50 chosen?

**A:** ResNet50 is known for strong image feature extraction capabilities.

**Q:** How did you train your model?

**A:** The model was trained using labeled frames with binary cross-entropy loss, optimized with Adam optimizer over multiple epochs.

**Q:** What metrics did you use to evaluate your model?

**A:** We used accuracy, precision, recall, and F1-score to evaluate the classification performance.

**Q:** How did you calculate the confidence score in your predictions?

**A:** The confidence score is the output probability from the model's final sigmoid layer indicating the likelihood that a frame is fake.

**Q:** How many fake frames were detected compared to the total frames?

**A:** The system counts frames predicted as fake by comparing the confidence score against a threshold (e.g., 0.5) and reports the ratio of fake frames to total frames.

**Q:** How was the model trained?

**A:** The dense classifier was trained using binary cross-entropy loss.

**Q:** How many epochs were used?

**A:** The model was trained for 20 epochs.

**Q:** What was the batch size?

**A:** A batch size of 32 was used.

**Q:** How did you evaluate the model?

**A:** Evaluation was done using test accuracy and confusion matrix.

**Q:** What validation split was used?

**A:** 20% of the data was used for validation.

**Q:** Were there any overfitting issues?

**A:** Minor overfitting occurred, controlled by dropout and early stopping.

**Q:** What are the key results or findings of your project?

**A:** The model achieved high accuracy in distinguishing fake frames from real ones, with precision and recall values indicating robust detection capabilities.

**Q:** Can you show some example predictions or visualizations?

**A:** Visualizations include frame-level confidence heatmaps and summary graphs showing the distribution of fake versus real frames per video.

**Q:** How reliable is your confidence score as a measure of prediction certainty?

**A:** The confidence score correlates strongly with model certainty, but borderline scores require cautious interpretation to avoid false positives or negatives.

**Q:** What accuracy did the model achieve?

**A:** The model achieved 92% accuracy on the test set.

**Q:** How many false positives were observed?

**A:** Around 4% of real videos were misclassified as fake.

**Q:** How many false negatives occurred?

**A:** 3% of fake videos were wrongly marked as real.

**Q:** What is the F1-score of the model?

**A:** The F1-score achieved was approximately 0.91.

**Q:** Is the model suitable for deployment?

**A:** Yes, the model is stable and ready for prototype deployment.

**Q:** What difficulties did you face during your project?

**A:** Challenges included handling large video datasets, extracting high-quality frames, class imbalance, and differentiating subtle manipulations.

**Q:** How did you overcome these challenges?

**A:** We used data cleaning, augmentation, and threshold tuning to improve robustness, along with careful model design and evaluation.

**Q:** What were the main challenges?

**A:** Handling poor-quality frames and imbalanced classes were challenges.

**Q:** Was class imbalance a problem?

**A:** Yes, real videos were fewer than fake videos.

**Q:** Did the dataset include occluded or blurred faces?

**A:** Yes, which sometimes reduced accuracy.

**Q:** Was GPU used for training?

**A:** Yes, training was accelerated using a GPU.

**Q:** Were large videos slow to process?

**A:** Yes, processing long videos with many frames was time-consuming.

**Q:** What improvements can be made to your project?

**A:** Future improvements include integrating temporal analysis across frames, real-time detection, and multi-modal approaches combining audio and video.

**Q:** Are there any plans to extend the system for real-time deployment?

**A:** Yes, optimizing inference speed and deploying on edge devices is a planned extension to allow real-time detection.

**Q:** What improvements are planned?

**A:** We plan to add audio deepfake detection and real-time support.

**Q:** Will live video detection be supported?

**A:** Yes, it's a key part of our future roadmap.

**Q:** Is mobile app integration possible?

**A:** Yes, we aim to build a mobile app in the future.

**Q:** Will GAN-based generation be studied?

**A:** Yes, future versions may analyze how GANs are generating deepfakes.

**Q:** What programming languages and libraries did you use?

**A:** Python was used with libraries such as TensorFlow/Keras or PyTorch for model building, OpenCV for video processing, and NumPy/Pandas for data manipulation.

**Q:** Did you use any cloud services or GPUs for training?

**A:** Yes, training was accelerated using GPUs, and cloud platforms like Google Colab or AWS were used for scalable computing.

**Q:** What programming language is used?

**A:** Python was used for all development.

**Q:** How is the chatbot built?

**A:** The chatbot is built using Streamlit and scikit-learn.

**Q:** How is a cybersecurity incident reported in the context of deepfake detection?

**A:** Incidents can be reported to national CERT teams or cybersecurity authorities with evidence like fake video samples and analysis reports.

**Q:** What is the role of confidence score in cybersecurity monitoring?

**A:** Confidence scores help prioritize alerts by indicating the certainty of detection, enabling faster incident response.

**Q:** How does frame extraction relate to deepfake detection accuracy?

**A:** Extracting clean, representative frames is crucial because poor quality or irrelevant frames can reduce detection accuracy.

**Q:** What ethical concerns are associated with deepfake detection technology?

**A:** Privacy, consent, and misuse of detection results are key concerns requiring responsible deployment and clear user guidelines.

**Q:** How can deepfake videos impact cybersecurity?

**A:** Deepfakes can be used for misinformation, social engineering attacks, impersonation, fraud, and identity theft, threatening personal privacy, organizational security, and public trust.

**Q:** How do you detect and report a cybersecurity incident related to deepfakes?

**A:** Upon detecting suspicious or fake videos, the incident can be reported to cybersecurity authorities like national CERT (Computer Emergency Response Team) or internal security teams with detailed logs, evidence (video samples, frame analyses), and confidence scores.

**Q:** What is the process of reporting a deepfake incident to cybersecurity teams?

**A:**

1. **Gather evidence:** Collect the video, extracted frames, confidence scores, and detection logs.
2. **Document:** Prepare a detailed report describing the nature and potential impact of the deepfake.
3. **Report:** Submit to the appropriate authority or internal cybersecurity incident response team.
4. **Follow up:** Cooperate with investigations and apply recommended mitigation strategies.

**Q:** What is a confidence score, and how is it calculated in deepfake detection?

**A:** The confidence score is a probability value generated by the model's output layer indicating how likely a frame or video is fake. It is usually calculated using a sigmoid function that outputs a value between 0 (real) and 1 (fake).

**Q:** How many frames do you analyze per video for reliable detection?

**A:** Typically, several hundred frames are extracted evenly across the video. The exact number depends on video length and frame rate but aims to cover diverse segments for comprehensive analysis.

**Q:** How do you determine the number of fake frames in a video?

**A:** Each extracted frame is evaluated, and frames with confidence scores above a certain threshold (e.g., 0.5) are labeled as fake. The total count of such frames provides the number of detected fake frames.

**Q:** How is the confidence score used in cybersecurity threat assessment?

**A:** Confidence scores help prioritize alerts by indicating the likelihood that content is

manipulated, enabling security teams to focus on high-risk cases for investigation and response.

**Q:** What cybersecurity best practices can complement deepfake detection?

**A:** Employee awareness training, strong authentication, monitoring social media for misinformation, and rapid incident response protocols all help mitigate the risks posed by deepfakes.

**Q:** How does frame extraction quality affect cybersecurity detection accuracy?

**A:** Poor quality or corrupted frames may produce inaccurate confidence scores, leading to false negatives or positives, which could delay detection and response to cyber threats.

**Q:** Are there privacy or ethical concerns when using deepfake detection in cybersecurity?

**A:** Yes, it is important to ensure users' privacy rights are respected, data is securely handled, and detection results are used responsibly to avoid wrongful accusations or misuse.