

Project 2 - Anomaly Detection in caltech101

Group 5: Yin Yin Thu, Wiwat Pholsomboon, Wassim Mecheri

Introduction

- Detect non-airplane (anomaly) images in Caltech101 dataset
- Vanilla Autoencoder will be used to detect anomalies (non-airplane images) from normal class (airplane)
- VAE and β -VAE will be test as alternative model to improve the performance of the detection.

>Data Preprocessing and Exploration

Import libraries

```
In [1]: # Import libraries
import tensorflow_datasets as tfds
import tensorflow as tf
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import time
from tensorflow import keras
from tensorflow.keras import layers
from tensorflow.keras.layers import Conv2D, MaxPooling2D, UpSampling2D, concatenate, Input
from tensorflow.keras.models import Model
from tabulate import tabulate
from sklearn.metrics import confusion_matrix, precision_score, recall_score, f1_score, accuracy_score, roc_
BATCH_SIZE = 64
RANDOM_SEED = 555
IMAGE_SIZE = 96
LATENT_DIM = 512
THRESHOLD = 0.03
```

```
VAE_THRESHOLD = 0.05
BVAE_THRESHOLD = 0.03
UNET_THRESHOLD = 0.0009
VERBOSE = 0

np.random.seed(RANDOM_SEED)
tf.random.set_seed(RANDOM_SEED)
```

Load the Caltech101 dataset

- Images dataset which contain 102 classes.
- Visualize the distribution of classes in the dataset

```
In [2]: # Load and combine datasets
ds = tfds.load('caltech101', split='train+test', as_supervised=True)

# Count images per class
class_counts = {}
for _, label in ds:
    label_id = label.numpy()
    class_counts[label_id] = class_counts.get(label_id, 0) + 1

# Get class names
class_names = tfds.builder('caltech101').info.features['label'].names

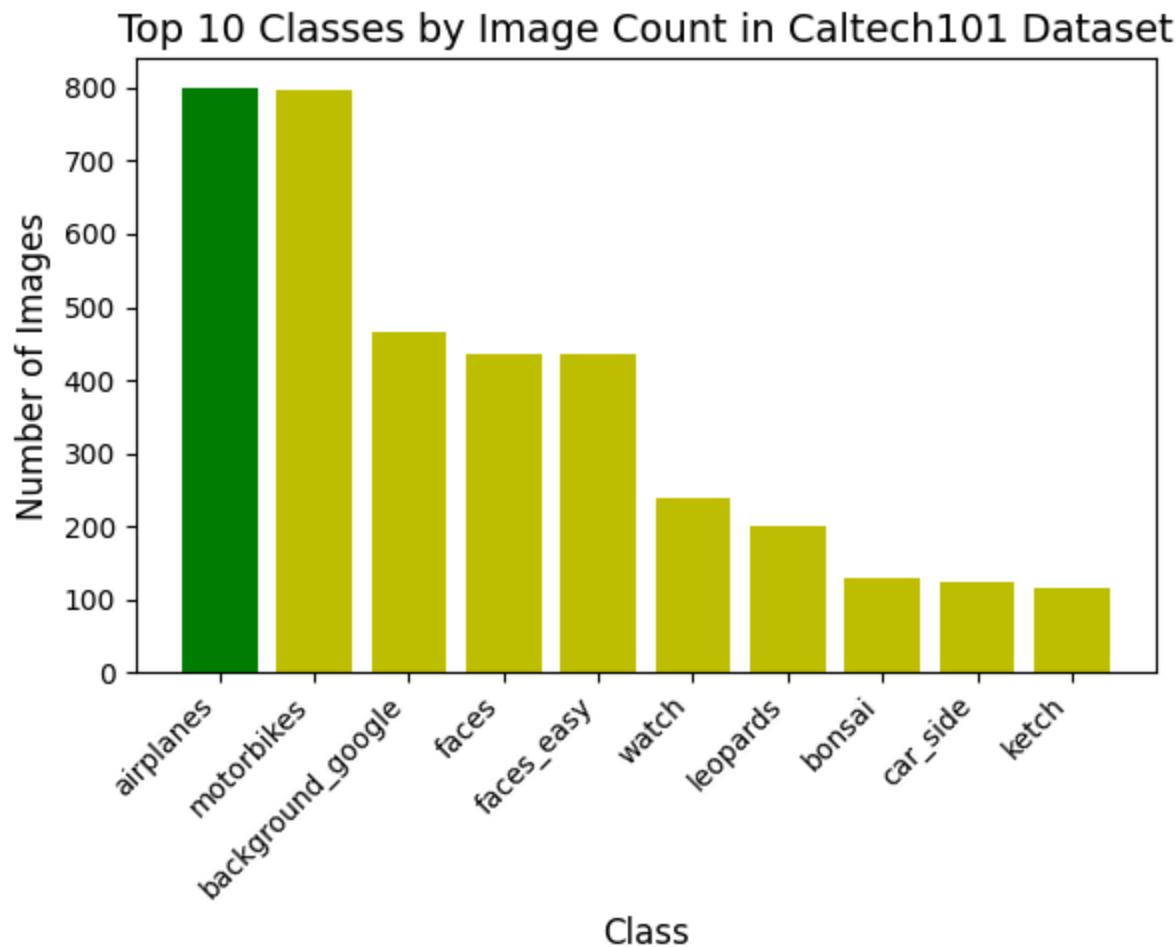
# Create and sort DataFrame
df = pd.DataFrame({
    'class_name': [class_names[i] for i in class_counts.keys()],
    'count': list(class_counts.values())
})
top10_df = df.nlargest(10, 'count')

# Plot
plt.figure(figsize=(6, 5))
plt.bar(top10_df['class_name'], top10_df['count'], color=['g', 'y', 'y', 'y', 'y', 'y', 'y', 'y', 'y', 'y'])
plt.xlabel('Class', fontsize=12)
plt.ylabel('Number of Images', fontsize=12)
plt.title('Top 10 Classes by Image Count in Caltech101 Dataset', fontsize=14)
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
```

```
plt.show()

# Print summary
print(f"Total classes in dataset: {len(class_names)}")
print(f"Total images in dataset: {sum(class_counts.values())}")
```

```
2025-04-01 12:37:27.083301: I metal_plugin/src/device/metal_device.cc:1154] Metal device set to: Apple M3 Max
2025-04-01 12:37:27.083331: I metal_plugin/src/device/metal_device.cc:296] systemMemory: 36.00 GB
2025-04-01 12:37:27.083335: I metal_plugin/src/device/metal_device.cc:313] maxCacheSize: 13.50 GB
WARNING: All log messages before absl::InitializeLog() is called are written to STDERR
I0000 00:00:1743525447.083346 42320482 pluggable_device_factory.cc:305] Could not identify NUMA node of platform GPU ID 0, defaulting to 0. Your kernel may not have been built with NUMA support.
I0000 00:00:1743525447.083365 42320482 pluggable_device_factory.cc:271] Created TensorFlow device (/job:localhost/replica:0/task:0/device:GPU:0 with 0 MB memory) -> physical PluggableDevice (device: 0, name: METAL, pci bus id: <undefined>)
2025-04-01 12:37:27.141600: I tensorflow/core/kernels/data/tf_record_dataset_op.cc:376] The default buffer size is 262144, which is overridden by the user specified `buffer_size` of 8388608
2025-04-01 12:37:27.752137: I tensorflow/core/framework/local_rendezvous.cc:405] Local rendezvous is aborting with status: OUT_OF_RANGE: End of sequence
```



Total classes in dataset: 102

Total images in dataset: 9144

Preprocessing

- Resize the image to 96x96 and preserve the aspect ratio
- Center crop the image to 96x96
- Normalize the image pixel values to 0-1
- Split the dataset into normal, anomaly images and evaluation dataset
- Random anomaly images 800 images from other classes to prevent imbalance dataset

```
In [3]: ### Split data into normal, anomaly sand evauiondataset
def scale_resize_image(image, label):
    image = tf.cast(image, tf.float32) / 255.0
    # Get center crop with aspect ratio preserved
    min_dim = tf.minimum(tf.shape(image)[0], tf.shape(image)[1])
    image = tf.image.crop_to_bounding_box(
        image,
        (tf.shape(image)[0] - min_dim) // 2,
        (tf.shape(image)[1] - min_dim) // 2,
        min_dim,
        min_dim
    )

    image = tf.image.resize(image, [IMAGE_SIZE, IMAGE_SIZE], ) # Resizing the image to 224x224 dimention
    return (image, label)

# Resize and center crop the image to 96x96 and convert to numpy array for visualization
combined_ds = (ds.map(scale_resize_image)).cache()
combined_images = [(image, label == 1) for image, label in combined_ds]

# Split 3 dataset
# For training and visualization
normal_images = np.array([image for image, airplane in combined_images if airplane == True ])
anomaly_images = np.array([image for image, airplane in combined_images if airplane == False ])

# Sample 800 images from anomaly dataset
anomaly_indices = np.random.choice(len(anomaly_images), size=800, replace=False)
anomaly_images = anomaly_images[anomaly_indices]

# For evaluation
test_X = np.concatenate([normal_images, anomaly_images])
test_y = np.concatenate([np.zeros(len(normal_images)), np.ones(len(anomaly_images))])
```

2025-04-01 12:37:28.559492: I tensorflow/core/framework/local_rendezvous.cc:405] Local rendezvous is aborting with status: OUT_OF_RANGE: End of sequence

Exploration

- Show sample images from normal (Airplane) and anomaly (Non-Airplane) classes

```
In [4]: # Display sample images from specific classes
def show_class_samples(dataset, class_name, num_samples=8):

    # Take the specified number of samples
    samples = dataset[:num_samples]

    # Create a figure to display the images
    fig, axes = plt.subplots(1, num_samples, figsize=(10, 2))
    fig.suptitle(f'Sample images from class: {class_name}', fontsize=14)

    # Display each sample
    for i, image in enumerate(samples):
        if i < num_samples:
            axes[i].imshow(image)
            axes[i].axis('off')

    plt.tight_layout()
    plt.show()

# Show samples from class 1 (airplane)
show_class_samples(normal_images, 'Normal Images')
show_class_samples(anomaly_images, 'Anomaly Images')
```

Sample images from class: Normal Images



Sample images from class: Anomaly Images



Unsupervised Learning Model Development

Vanilla Autoencoder

- Encoder contains 3 convolutional layers and 1 fully connected layer convert to latent vector
- Decoder contains 3 convolutional transpose layers and 1 convolutional layer to reconstruct the image

In [5]:

```
# Build the Autoencoder model
class Autoencoder(keras.Model):
    def __init__(self, img_shape, latent_dim, num_hidden_layer, num_filter, **kwargs):
        super(Autoencoder, self).__init__(**kwargs)

        encoder = tf.keras.Sequential()
        for _ in range(num_hidden_layer):
            encoder.add(layers.Conv2D(num_filter, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.he_normal()))
            num_filter *= 2
        encoder.add(layers.Flatten())
        encoder.add(layers.Dense(num_filter * (2 ** num_hidden_layer), activation="relu", kernel_initializer=tf.keras.initializers.he_normal()))
        encoder.add(layers.Dense(latent_dim, name="latent_vector", kernel_initializer=tf.keras.initializers.he_normal()))

        h, w, c = img_shape
        decoder_starting_dims = (h // (2 ** num_hidden_layer), w // (2 ** num_hidden_layer), 128)

        decoder = tf.keras.Sequential()
        decoder.add(layers.Dense(np.prod(decoder_starting_dims), activation="relu", kernel_initializer=tf.keras.initializers.he_normal()))
        decoder.add(layers.Reshape(decoder_starting_dims))
        for _ in range(num_hidden_layer):
            decoder.add(layers.Conv2DTranspose(num_filter, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.he_normal()))
```

```

        num_filter //= 2
        decoder.add(layers.Conv2D(c, 3, activation="sigmoid", padding="same", kernel_initializer=tf.keras.layers.
        self.encoder = encoder
        self.decoder = decoder

    def call(self, inputs):
        latent = self.encoder(inputs)
        return self.decoder(latent)

# Main training function
def train_autoencoder(dataset, img_shape=(224, 224, 3), latent_dim=64, epochs=20, learning_rate=0.001, num_
# Create Autoencoder
autoencoder = Autoencoder(img_shape, latent_dim, num_hidden_layer, num_filter)
autoencoder.compile(optimizer=keras.optimizers.Adam(learning_rate=learning_rate), loss="mse", metrics=
# Train the model
history = autoencoder.fit(dataset, dataset, epochs=epochs, batch_size=BATCH_SIZE, verbose=VERBOSE)

return autoencoder, history

# Generate reconstructions from the autoencoder
def generate_reconstructions(autoencoder, test_images, n=8):
    # Get test images
    test_sample = test_images[:n]
    # Get reconstructions
    reconstructed = autoencoder.predict(test_sample, verbose=VERBOSE)

    return test_sample, reconstructed

# Visualize original vs reconstructed images
def visualize_reconstructions(originals, reconstructions, n=8):
    import matplotlib.pyplot as plt

    plt.figure(figsize=(14, 4))
    for i in range(n):
        # Original
        ax = plt.subplot(2, n, i + 1)
        plt.imshow(originals[i])
        plt.title("Original")
        plt.axis("off")

        # Reconstruction

```

```

        ax = plt.subplot(2, n, i + 1 + n)
        plt.imshow(reconstructions[i])
        plt.title("Reconstructed")
        plt.axis("off")

    plt.tight_layout()
    plt.show()

# Calculate ROC AUC score for both models
def calculate_roc_auc(y_true, scores):
    # Normalize scores to 0-1 range
    scores_normalized = (scores - np.min(scores)) / (np.max(scores) - np.min(scores))
    fpr, tpr, _ = roc_curve(y_true, scores_normalized)
    return fpr, tpr, auc(fpr, tpr)

```

Train the Autoencoder and generate reconstructions

- Train the autoencoder with normal images (Airplane)
- Generate reconstructions from the autoencoder
- Loss function is Mean Squared Error (MSE)
- Optimizer is Adam with learning rate 0.001

In [6]:

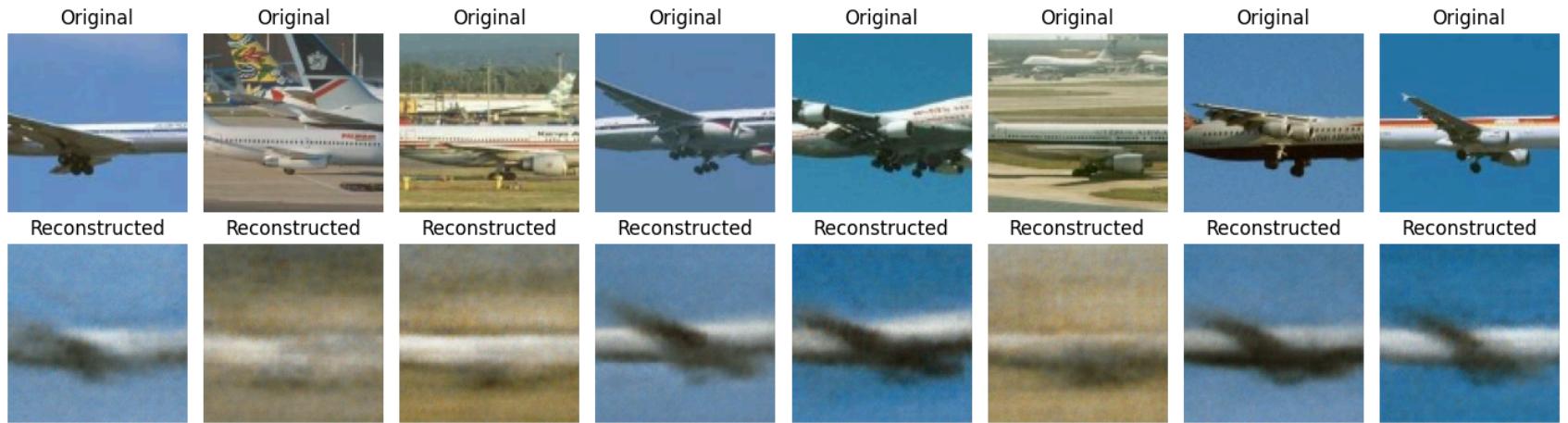
```

# Train the autoencoder
start_time = time.time()
autoencoder, history = train_autoencoder(normal_images, img_shape=(IMAGE_SIZE, IMAGE_SIZE, 3), latent_dim=latent_dim)
ae_training_time = time.time() - start_time

# Generate reconstructions
original, reconstructed = generate_reconstructions(autoencoder, normal_images, n=8)
visualize_reconstructions(original, reconstructed, n=8)

```

2025-04-01 12:37:29.991640: I tensorflow/core/grappler/optimizers/custom_graph_optimizer_registry.cc:117] Plugin optimizer for device_type GPU is enabled.



Anomaly Detection

- Detect anomalies (image from other classes) with trained autoencoder.
- Visualize the best and the worst images from anomaly classes and normal classes.
- Threshold is 0.03, come from hyperparameter tuning result.

```
In [7]: # Function to detect anomalies using autoencoder
def detect_anomalies_autoencoder(autoencoder, images, threshold):
    reconstructions = autoencoder.predict(images, verbose=VERBOSE)
    # Calculate MSE for each image
    mse = np.mean(np.square(images - reconstructions), axis=(1, 2, 3))
    # Determine if each image is an anomaly based on threshold
    anomalies = mse > threshold
    return mse, anomalies, reconstructions

# Function to visualize anomalies
def visualize_anomalies(original_images, reconstructions, scores, anomalies, model_name, n=8, is_normal_class=False):
    plt.figure(figsize=(14, 5))
    # Get indices with lowest and highest scores
    display_indices = np.concatenate([np.argsort(scores) [:n//2], np.argsort(scores) [:-1] [:n//2]])

    # Set colors based on class type
    anomaly_color = 'red' if is_normal_class else 'green'
    normal_color = 'green' if is_normal_class else 'red'
```

```

for i in range(min(n, len(original_images))):
    idx = display_indices[i]
    # Original image
    plt.subplot(2, n, i + 1)
    plt.imshow(original_images[idx])
    title = f"Original\nScore: {scores[idx]:.4f}"
    title += "\nANOMALY" if anomalies[idx] else "\n"
    plt.title(title, color=anomaly_color if anomalies[idx] else normal_color)
    plt.axis("off")

    # Reconstruction
    plt.subplot(2, n, i + n + 1)
    plt.imshow(reconstructions[idx])
    plt.title("Reconstructed")
    plt.axis("off")

plt.suptitle(f"Anomaly Detection with {model_name} (Threshold: {threshold})", fontsize=16)
plt.tight_layout()
plt.show()

# Detect anomalies with regular autoencoder
ae_scores, ae_anomalies, ae_reconstructions = detect_anomalies_autoencoder(autoencoder, anomaly_images, threshold)
visualize_anomalies(anomaly_images, ae_reconstructions, ae_scores, ae_anomalies, "Vanilla Autoencoder(Anomaly)")

print(f"Vanilla Autoencoder detected {np.sum(ae_anomalies)} anomalies from {len(anomaly_images)} images")

# Detect anomalies with regular autoencoder
normal_scores, normal_anomalies, normal_reconstructions = detect_anomalies_autoencoder(autoencoder, normal_images, threshold)
visualize_anomalies(normal_images, normal_reconstructions, normal_scores, normal_anomalies, "Vanilla Autoencoder(Normal)")

print(f"Vanilla Autoencoder disclassified {np.sum(normal_anomalies)} normal images from {len(normal_images)} images")

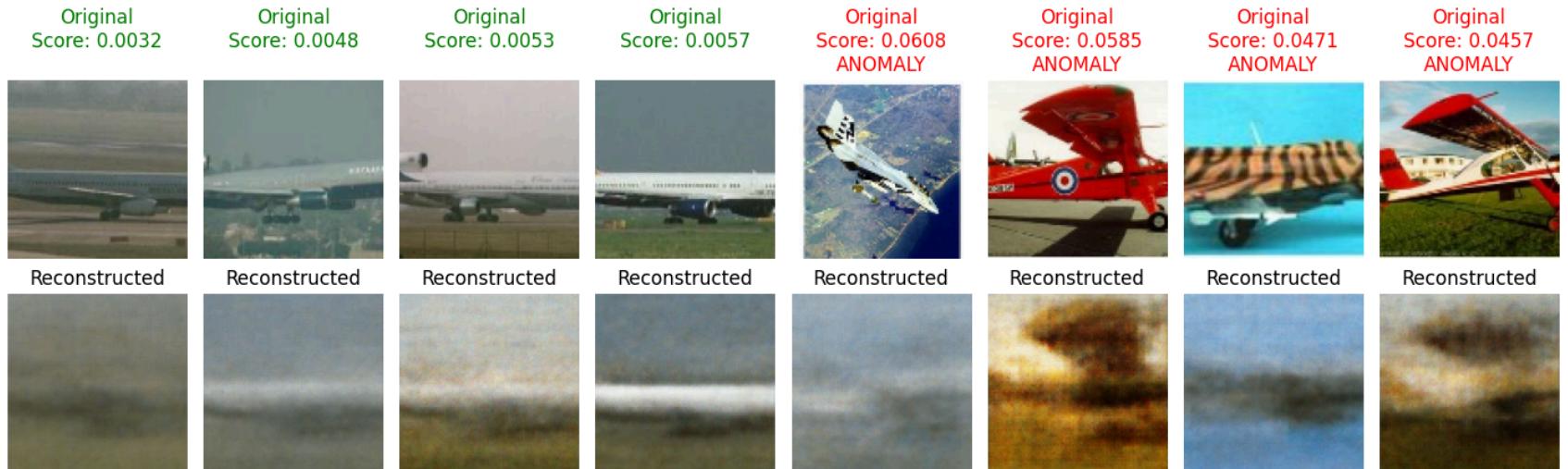
```

Anomaly Detection with Vanilla Autoencoder(Anomaly) (Threshold: 0.03)



Vanilla Autoencoder detected 639 anomalies from 800 images

Anomaly Detection with Vanilla Autoencoder(Normal) (Threshold: 0.03)



Vanilla Autoencoder disclassified 43 normal images from 800 images

👉 Actionable Recommendations

Feature Analysis

- Feature analysis below show the portion of the image that contribute to the anomaly detection.

1. Anomaly Item detected as Normal(Airplane)

- Has background color as grey or blue.
- Has object on the center of the image.
- Majority of the image is background.

2. Normal Item detected as Anomaly(Non-Airplane)

- Plane is in red color and cover more than 50% of the image.
- Plane which station on the ground and the background is not sky.

Actionable Recommendations

1. Increase complexity of the model to detect more patterns.
2. Increase size of dataset use to train the model. Especially on red airplane and airplane station on the ground.
3. Increase latent dimension to capture more features.

```
In [8]: # Actionable Recommendations (4 marks)

## Feature Importance Analysis
# Function to visualize the difference between original and reconstructed images
def visualize_difference_maps(originals, reconstructions, n=5):
    plt.figure(figsize=(15, 6))
    for i in range(min(n, len(originals))):
        # Original
        plt.subplot(3, n, i + 1)
        plt.imshow(originals[i])
        plt.title("Original")
        plt.axis("off")

        # Reconstruction
        plt.subplot(3, n, i + 1 + n)
        plt.imshow(reconstructions[i])
        plt.title("Reconstructed")
        plt.axis("off")
```

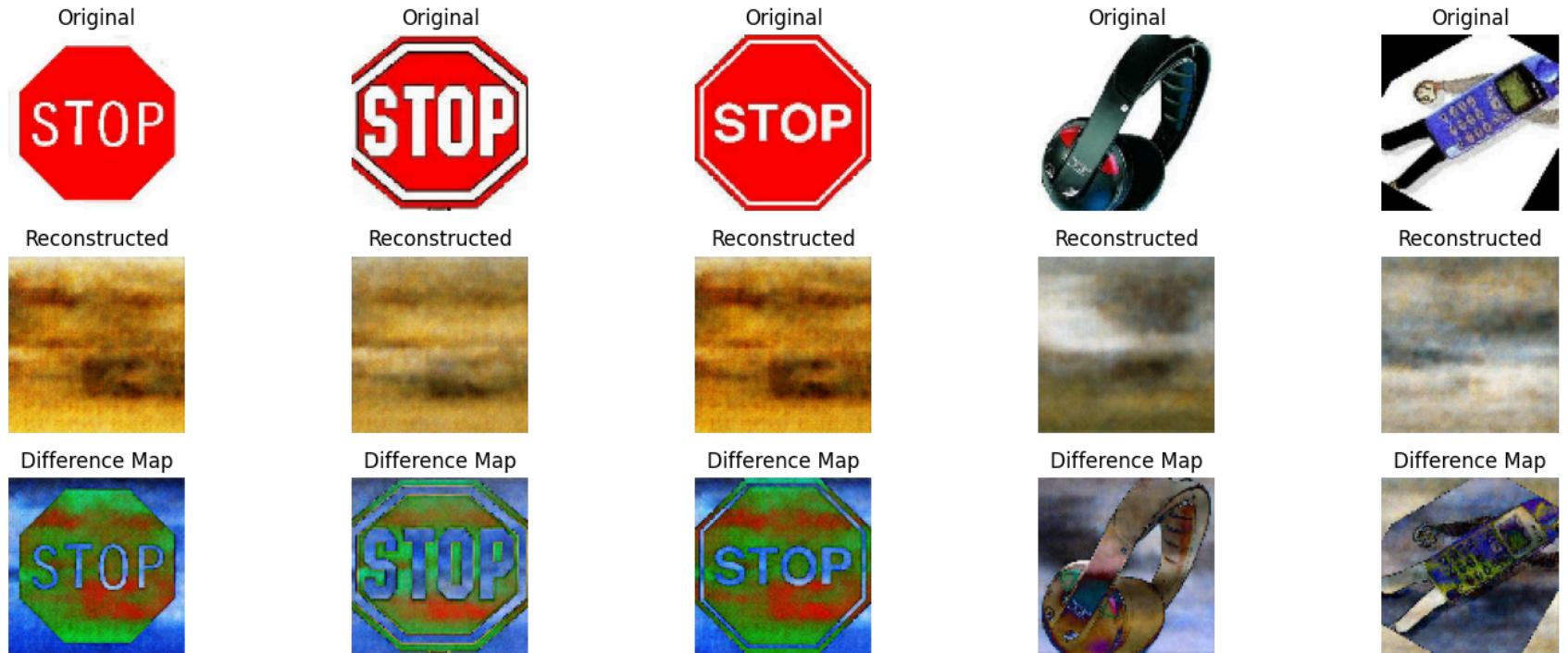
```
# Difference map
diff = np.abs(originals[i] - reconstructions[i])
# Normalize for better visualization
diff = diff / np.max(diff) if np.max(diff) > 0 else diff

plt.subplot(3, n, i + 1 + 2*n)
plt.imshow(diff, cmap='hot')
plt.title("Difference Map")
plt.axis("off")

plt.suptitle("Feature Contribution to Anomaly Detection", fontsize=16)
plt.tight_layout()
plt.show()

# Select some high-anomaly score examples
anomaly_indices = np.argsort(-ae_scores)[:5] # Top 5 highest anomaly scores
high_anomaly_imgs = anomaly_images[anomaly_indices]
high_anomaly_recon = ae_reconstructions[anomaly_indices]
visualize_difference_maps(high_anomaly_imgs, high_anomaly_recon)
```

Feature Contribution to Anomaly Detection



Bright color in Difference Map indicated area that is different between original and reconstructed images.

🧐 Evaluation and Visualization

- Calculate evaluation metrics for vanilla autoencoder.
- Plot ROC/AUC curve vanilla autoencoder.

```
In [9]: def calculate_metrics(y_true, scores, threshold):  
    # Convert scores to predictions based on threshold  
    y_pred = (scores > threshold).astype(int)  
  
    # Calculate basic metrics  
    accuracy = accuracy_score(y_true, y_pred)  
    precision = precision_score(y_true, y_pred)
```

```

recall = recall_score(y_true, y_pred)
f1 = f1_score(y_true, y_pred)

# Calculate confusion matrix
tn, fp, fn, tp = confusion_matrix(y_true, y_pred).ravel()

return [
    ["Metric", "Value"],
    ["Accuracy", f"{accuracy:.4f}"],
    ["Precision", f"{precision:.4f}"],
    ["Recall", f"{recall:.4f}"],
    ["F1 Score", f"{f1:.4f}"],
    ["✓ True Positives", tp],
    ["✗ False Positives", fp],
    ["✓ True Negatives", tn],
    ["✗ False Negatives", fn]
]

# Calculate metrics for both models on the test set
print("Calculating performance metrics for Vanilla Autoencoder...")
ae_scores, ae_anomalies, ae_reconstructions = detect_anomalies_autoencoder(autoencoder, test_X, THRESHOLD)
ae_metrics = calculate_metrics(test_y, ae_scores, THRESHOLD)

print("\nAutoencoder Performance Metrics:")
print(tabulate(ae_metrics, headers="firstrow", tablefmt="grid", numalign="right"))

# Calculate ROC AUC score
fpr_ae, tpr_ae, auc_ae = calculate_roc_auc(test_y, ae_scores)

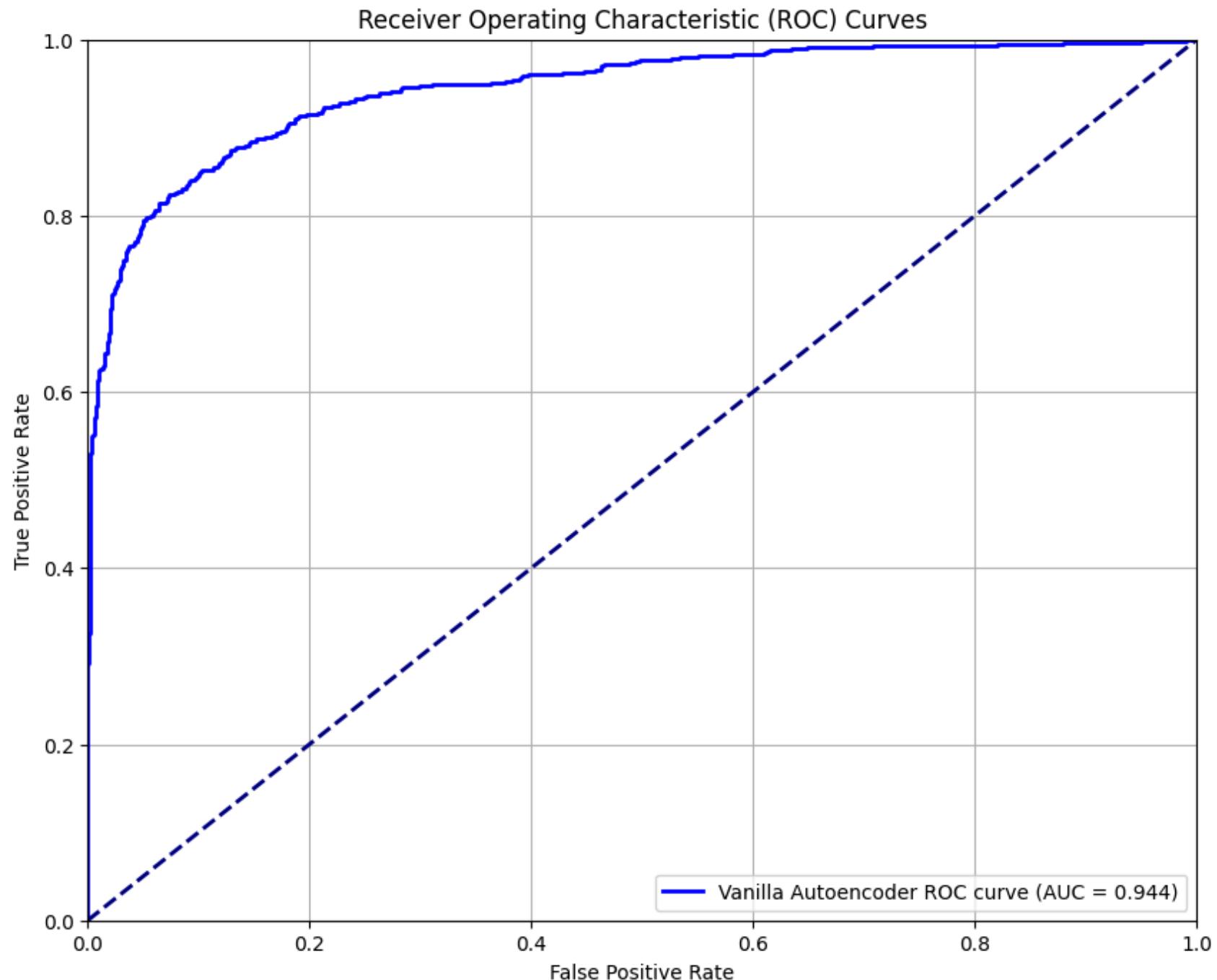
# Plot ROC curves
plt.figure(figsize=(10, 8))
plt.plot(fpr_ae, tpr_ae, color='blue', lw=2, label=f'Vanilla Autoencoder ROC curve (AUC = {auc_ae:.3f})')
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curves')
plt.legend(loc="lower right")
plt.grid(True)
plt.show()

```

Calculating performance metrics for Vanilla Autoencoder...

Autoencoder Performance Metrics:

Metric	Value
Accuracy	0.8725
Precision	0.937
Recall	0.7987
F1 Score	0.8623
✓ True Positives	639
✗ False Positives	43
✓ True Negatives	757
✗ False Negatives	161



Model Refinement and Optimization

List of hyperparameter tuning and alternative model finding:

1. Hyperparameter tuning on Vanilla Autoencoder in various threshold.
2. Hyperparameter tuning on Vanilla Autoencoder in various latent dimension.
3. Hyperparameter tuning on Vanilla Autoencoder in various number of hidden layer.
4. Try to use VAE to improve the performance of the detection.
5. Try to use β -VAE to improve the performance of the detection.
6. Try to use U-Net AE to improve the performance of the detection.

1. Hyperparameter Tuning on Vanilla Autoencoder in various threshold

- Threshold `0.01` has the highest f1 score and accuracy indicating best performance but it miss a lot on normal class.
- Threshold `0.03` show balanced performance on both normal and anomaly class.
- By the way, in real life scenario, we might need to consider the cost of false alarm and false negative on Airplane detection which might make `0.04` or `0.05` the best threshold if Airplant miss detection is costly.

```
In [10]: thresholds = [0.005, 0.01, 0.02, 0.025, 0.03, 0.035, 0.04, 0.05]

# Store metrics for each threshold
all_metrics = {}

for threshold in thresholds:
    # Get metrics for current threshold
    metrics = calculate_metrics(test_y, ae_scores, threshold)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]

    # Initialize dictionary on first run
    if not all_metrics:
        all_metrics = {row[0]: [] for row in metric_data}

    # Add values for this threshold
    for row in metric_data:
```

```

        all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"{t}" for t in thresholds]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("Vanilla Autoencoder Performance Across Thresholds:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))

```

Vanilla Autoencoder Performance Across Thresholds:

Metric	0.005	0.01	0.02	0.025	0.03	0.035	0.04	0.05
Accuracy	0.5012	0.5913	0.8462	0.8694	0.8725	0.8462	0.8075	0.7438
Precision	0.5006	0.5506	0.7966	0.8764	0.937	0.9663	0.9767	0.9949
Recall	1	0.9938	0.93	0.86	0.7987	0.7175	0.63	0.49
F1 Score	0.6672	0.7086	0.8581	0.8681	0.8623	0.8235	0.766	0.6566
✓ True Positives	800	795	744	688	639	574	504	392
✗ False Positives	798	649	190	97	43	20	12	2
✓ True Negatives	2	151	610	703	757	780	788	798
✗ False Negatives	0	5	56	112	161	226	296	408

2. Hyperparameter tuning on Vanilla Autoencoder in various latent dimension.

- Latent dimension 1280 show the best performance of F1-score balanced between precision and recall.
- Latent dimension 1024 show best performance on normal class make it most suitable for real life scenario.

```
In [11]: latent_dims = [512, 768, 1024, 1280]

# Store metrics for each threshold
all_metrics = {}

for latent_dim in latent_dims:
    # Get metrics for current threshold
    autoencoder, history = train_autoencoder(normal_images, img_shape=(IMAGE_SIZE, IMAGE_SIZE, 3), latent_d
    ae_scores, ae_anomalies, ae_reconstructions = detect_anomalies_autoencoder(autoencoder, test_X, THRESHOLD)
    metrics = calculate_metrics(test_y, ae_scores, THRESHOLD)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]

    # Initialize dictionary on first run
    if not all_metrics:
        all_metrics = {row[0]: [] for row in metric_data}

    # Add values for this threshold
    for row in metric_data:
        all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"{t}" for t in latent_dims]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("Vanilla Autoencoder Performance Across Latent Dimension:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))
```

Vanilla Autoencoder Performance Across Latent Dimension:

Metric	512	768	1024	1280
Accuracy	0.8725	0.8756	0.865	0.8588
Precision	0.937	0.9532	0.9665	0.9031
Recall	0.7987	0.79	0.7562	0.8037
F1 Score	0.8623	0.864	0.8485	0.8505
✓ True Positives	639	632	605	643
✗ False Positives	43	31	21	69
✓ True Negatives	757	769	779	731
✗ False Negatives	161	168	195	157

3. Hyperparameter tuning on Vanilla Autoencoder in various number of hidden layer.

- Number of hidden layer 5 show the best performance of F1-score balanced between precision and recall.
- Number of hidden layer 3 show best performance on normal class make it most suitable for real life scenario.

```
In [12]: hidden_layers = [3, 4, 5]

# Store metrics for each threshold
all_metrics = {}

for hidden_layer in hidden_layers:
    # Get metrics for current threshold
    autoencoder, history = train_autoencoder(normal_images, img_shape=(IMAGE_SIZE, IMAGE_SIZE, 3), num_hidden_layers=hidden_layer)
    ae_scores, ae_anomalies, ae_reconstructions = detect_anomalies_autoencoder(autoencoder, test_X, THRESHOLD)
    metrics = calculate_metrics(test_y, ae_scores, THRESHOLD)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]
```

```

# Initialize dictionary on first run
if not all_metrics:
    all_metrics = {row[0]: [] for row in metric_data}

# Add values for this threshold
for row in metric_data:
    all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"{t}" for t in hidden_layers]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("Vanilla Autoencoder Performance Across Hidden Layer:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))

```

Vanilla Autoencoder Performance Across Hidden Layer:

Metric	3	4	5
Accuracy	0.8662	0.8625	0.8506
Precision	0.9636	0.9167	0.8367
Recall	0.7612	0.7975	0.8712
F1 Score	0.8506	0.8529	0.8536
✓ True Positives	609	638	697
✗ False Positives	23	58	136
✓ True Negatives	777	742	664
✗ False Negatives	191	162	103

4. Variational Autoencoder

- define function to train, visualize, and detect anomalies for VAE

```
In [13]: # Build the VAE model
class VAE(keras.Model):
    def __init__(self, img_shape, latent_dim, **kwargs):
        super(VAE, self).__init__(**kwargs)
        self.encoder = keras.Sequential(
            [
                layers.Conv2D(32, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2D(64, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2D(128, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Flatten(),
                layers.Dense(256, activation="relu", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Dense(latent_dim * 2, name="latent_vector", kernel_initializer=tf.keras.initializers.GlorotUniform())
            ]
        )

        h, w, c = img_shape
        encoder_conv_layers = 3
        decoder_starting_dims = (h // (2 ** encoder_conv_layers), w // (2 ** encoder_conv_layers), 128)

        self.decoder = keras.Sequential(
            [
                layers.Dense(np.prod(decoder_starting_dims), activation="relu", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Reshape(decoder_starting_dims),
                layers.Conv2DTranspose(128, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2DTranspose(64, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2DTranspose(32, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2D(c, 3, activation="sigmoid", padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform())
            ]
        )

    def encode(self, x):
        mean_log_var = self.encoder(x)
        mean, log_var = tf.split(mean_log_var, num_or_size_splits=2, axis=1)
        return mean, log_var

    def reparameterize(self, mean, log_var):
        eps = tf.random.normal(shape=tf.shape(mean))
        return mean + eps * tf.exp(log_var * 0.5)
```

```

    def decode(self, z):
        return self.decoder(z)

    def call(self, x):
        mean, log_var = self.encode(x)
        z = self.reparameterize(mean, log_var)
        reconstructed = self.decode(z)

        # Add KL divergence loss as a model metric
        kl_loss = -0.5 * tf.reduce_mean(
            tf.reduce_sum(1 + log_var - tf.square(mean) - tf.exp(log_var), axis=1)
        )
        self.add_loss(kl_loss)

    return reconstructed

# Main training function
def train_vae(dataset, img_shape=(224, 224, 3), latent_dim=64, epochs=20, learning_rate=0.001):
    # Create VAE
    vae = VAE(img_shape, latent_dim)

    vae.compile(optimizer=keras.optimizers.Adam(learning_rate=learning_rate), loss="mse", metrics=["mse"])

    # Train the model
    history = vae.fit(dataset, dataset, batch_size=BATCH_SIZE, epochs=epochs, verbose=VERBOSE)

    return vae, history

# Generate reconstructions from the VAE
def generate_vae_reconstructions(vae, test_images, n=8):
    # Get test images
    test_sample = test_images[:n]
    reconstructed = vae.predict(test_sample, verbose=VERBOSE)
    return test_sample, reconstructed

def visualize_vae_reconstructions(originals, reconstructions, n=8):
    plt.figure(figsize=(14, 4))
    for i in range(min(n, len(originals))):
        # Original
        ax = plt.subplot(2, n, i + 1)
        plt.imshow(originals[i])

```

```

plt.title("Original")
plt.axis("off")

# Reconstruction
ax = plt.subplot(2, n, i + 1 + n)
plt.imshow(reconstructions[i])
plt.title("Reconstructed")
plt.axis("off")

plt.tight_layout()
plt.show()

def combine_metrics(*args, metric_names=None):
    all_metrics = {}

    # Process each experiment's metrics
    for _, metrics_dict in enumerate(args):
        metric_data = metrics_dict[1:] # Skip header row

        # Initialize dictionary on first run
        if not all_metrics:
            all_metrics = {row[0]: [] for row in metric_data}

        # Add values for this experiment
        for row in metric_data:
            all_metrics[row[0]].append(row[1])

    # Build the table
    header_row = ["Metric"] + [i for i in metric_names]
    combined_table = [header_row]

    # Add each metric row
    for metric, values in all_metrics.items():
        combined_table.append([metric] + values)

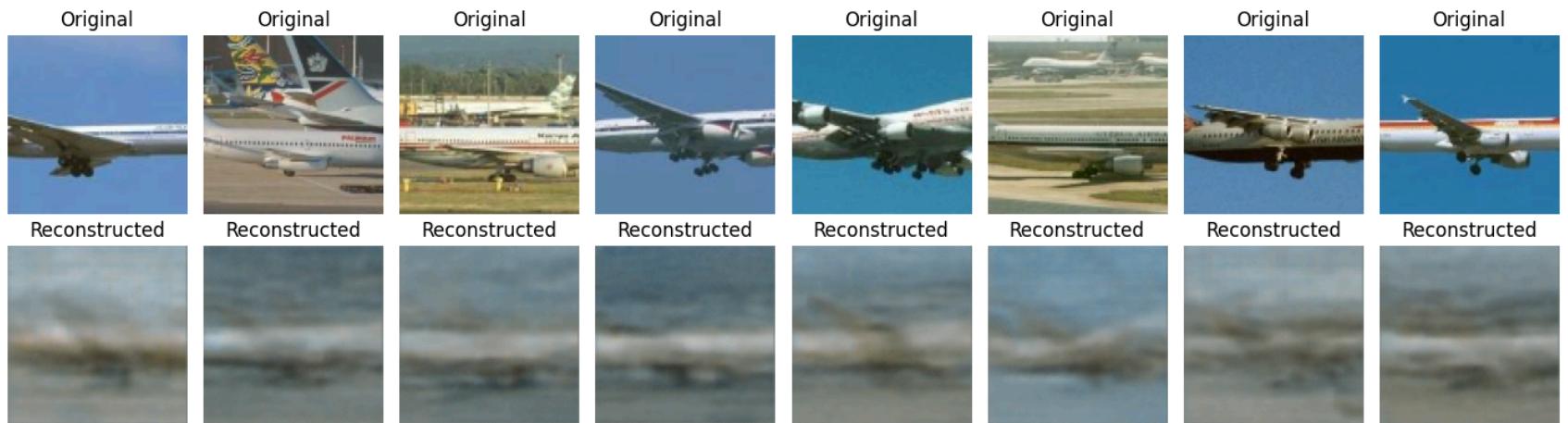
    return combined_table

```

Train the VAE model

```
In [14]: start_time = time.time()
vae, history = train_vae(normal_images, img_shape=(IMAGE_SIZE, IMAGE_SIZE, 3), latent_dim=LATENT_DIM, epochs=10)
vae_training_time = time.time() - start_time

original, reconstructed = generate_vae_reconstructions(vae, normal_images, n=8)
visualize_vae_reconstructions(original, reconstructed, n=8)
```



Visualize result of VAE

```
In [15]: def detect_anomalies_vae(vae, images, threshold):
    reconstructed = vae.predict(images, verbose=VERBOSE)

    # Calculate MSE for each image
    mse = np.mean(np.square(images - reconstructed), axis=(1, 2, 3))

    # Determine if each image is an anomaly based on threshold
    anomalies = mse > threshold
    return mse, anomalies, reconstructed

# Detect anomalies with VAE
vae_scores, vae_anomalies, vae_reconstructions = detect_anomalies_vae(vae, anomaly_images, threshold=VAE_THRESHOLD)
visualize_anomalies(anomaly_images, vae_reconstructions, vae_scores, vae_anomalies, "Variational Autoencoder")
print(f"VAE detected {np.sum(vae_anomalies)} anomalies from {len(anomaly_images)} images")

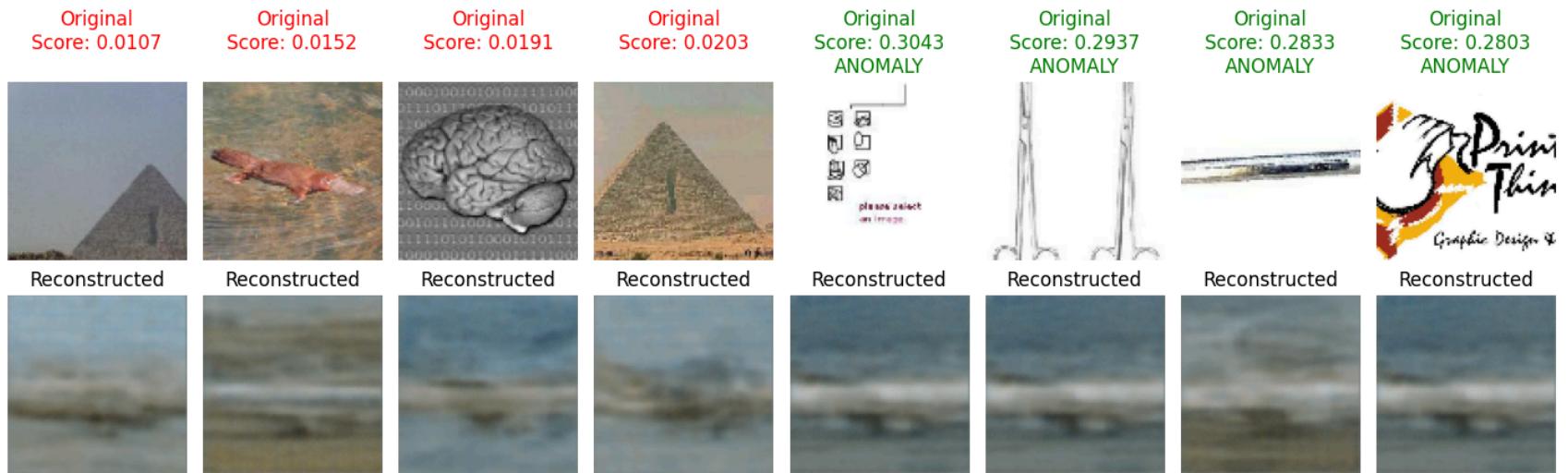
normal_scores, normal_anomalies, normal_reconstructions = detect_anomalies_vae(vae, normal_images, threshold=VAE_THRESHOLD)
visualize_anomalies(normal_images, normal_reconstructions, normal_scores, normal_anomalies, "Variational Autoencoder")
print(f"VAE detected {np.sum(normal_anomalies)} anomalies from {len(normal_images)} images")
```

```

visualize_anomalies(normal_images, normal_reconstructions, normal_scores, normal_anomalies, "Variational Autoencoder")
print(f"VAE disclassified {np.sum(normal_anomalies)} normal images from {len(normal_images)} images")

```

Anomaly Detection with Variational Autoencoder (Threshold: 0.05)



VAE detected 715 anomalies from 800 images

Anomaly Detection with Variational Autoencoder (Threshold: 0.05)



VAE disclassified 367 normal images from 800 images

Hyperparameter tuning on VAE Threshold

In [16]:

```
# Calculate metrics for VAE
vae_scores, _, _ = detect_anomalies_vae(vae, test_X, VAE_THRESHOLD)
vae_metrics = calculate_metrics(test_y, vae_scores, VAE_THRESHOLD)

thresholds = [0.005, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07]

# Store metrics for each threshold
all_metrics = {}

for threshold in thresholds:
    # Get metrics for current threshold
    metrics = calculate_metrics(test_y, vae_scores, threshold)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]

    # Initialize dictionary on first run
    if not all_metrics:
        all_metrics = {row[0]: [] for row in metric_data}

    # Add values for this threshold
    for row in metric_data:
        all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"t{t}" for t in thresholds]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("Vanilla Autoencoder Performance Across Thresholds:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))
```

Vanilla Autoencoder Performance Across Thresholds:

Metric	0.005	0.01	0.02	0.03	0.04	0.05	0.06	0.07
Accuracy	0.5	0.5006	0.5062	0.5619	0.6512	0.7175	0.755	0.7494
Precision	0.5	0.5003	0.5032	0.5336	0.5945	0.6608	0.7237	0.7574
Recall	1	1	0.9962	0.9825	0.9513	0.8938	0.825	0.7338
F1 Score	0.6667	0.6669	0.6686	0.6916	0.7317	0.7598	0.771	0.7454
✓ True Positives	800	800	797	786	761	715	660	587
✗ False Positives	800	799	787	687	519	367	252	188
✓ True Negatives	0	1	13	113	281	433	548	612
✗ False Negatives	0	0	3	14	39	85	140	213

Plot performance metrics of VAE

- Bases on the performance metrics, VAE perform worse than Vanilla Autoencoder in anomaly detection on airplane class in this dataset.

```
In [17]: print("\nAE vs. VAE Performance Metrics:")
combined = combine_metrics(ae_metrics, vae_metrics, metric_names=['AE', 'VAE'])
print(tabulate(combined, headers="firstrow", tablefmt="grid", numalign="right"))

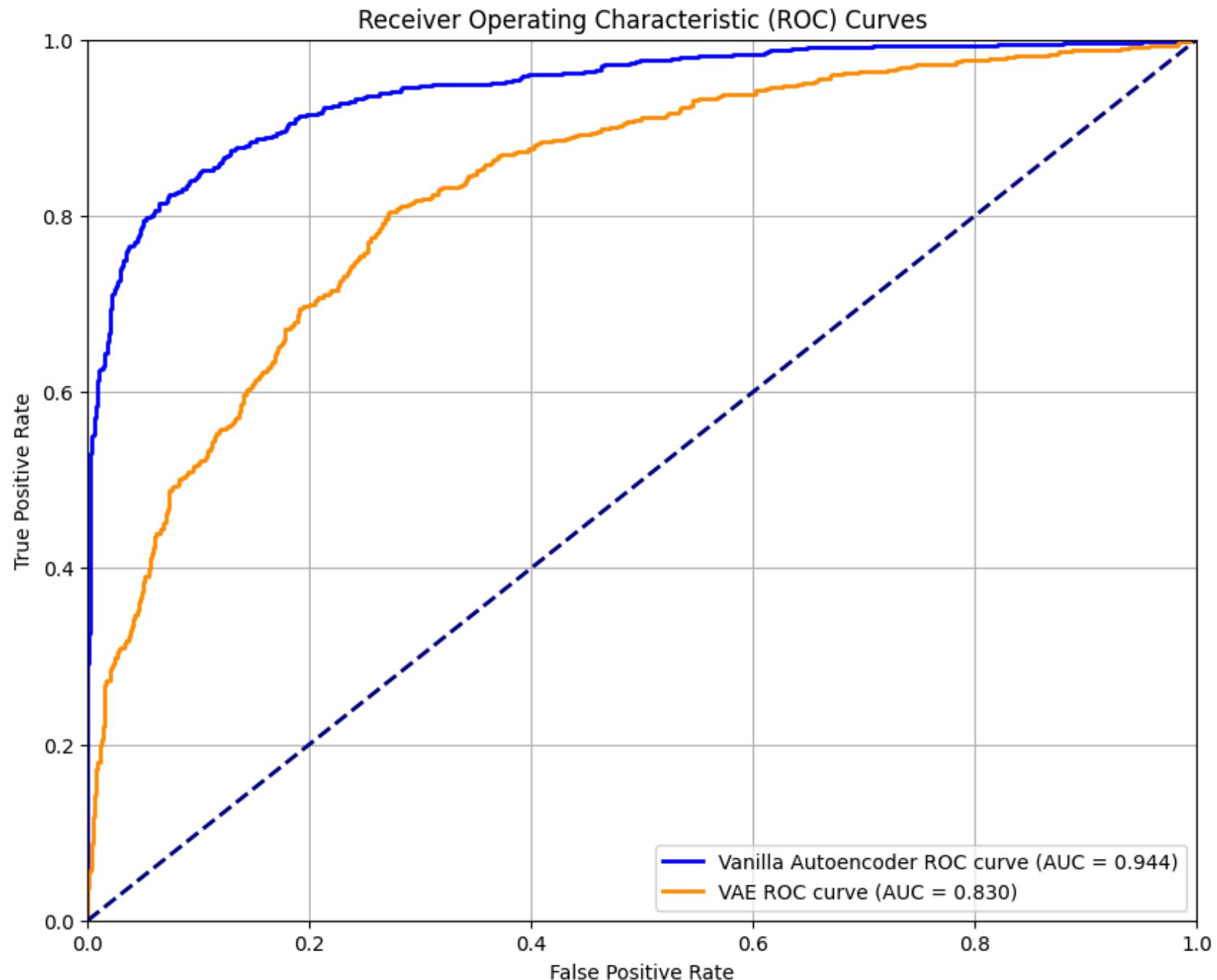
# Calculate ROC AUC score
fpr_vae, tpr_vae, auc_vae = calculate_roc_auc(test_y, vae_scores)

# Plot ROC curves
plt.figure(figsize=(10, 8))
plt.plot(fpr_ae, tpr_ae, color='blue', lw=2, label=f'Vanilla Autoencoder ROC curve (AUC = {auc_ae:.3f})')
plt.plot(fpr_vae, tpr_vae, color='darkorange', lw=2, label=f'VAE ROC curve (AUC = {auc_vae:.3f})')
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
```

```
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curves')
plt.legend(loc="lower right")
plt.grid(True)
plt.show()
```

AE vs. VAE Performance Metrics:

Metric	AE	VAE
Accuracy	0.8725	0.7175
Precision	0.937	0.6608
Recall	0.7987	0.8938
F1 Score	0.8623	0.7598
✓ True Positives	639	715
✗ False Positives	43	367
✓ True Negatives	757	433
✗ False Negatives	161	85



5. β-VAE

As we can see, VAE is unsuccesful at reconstructing the pictures. The output are grey images and it's impossible to recognize the planes.

A reason for that might be that the KL divergence overpowering the reconstruction loss and making the images like that, to deal with this we can add a β (beta) term to the KL divergence to reduce it's impact on the function.

In [18]:

```
# Build the VAE model
class BVAE(keras.Model):
    def __init__(self, img_shape, latent_dim, **kwargs):
        super(BVAE, self).__init__(**kwargs)
        self.encoder = keras.Sequential(
            [
                layers.Conv2D(32, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2D(64, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2D(128, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Flatten(),
                layers.Dense(256, activation="relu", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Dense(latent_dim * 2, name="latent_vector", kernel_initializer=tf.keras.initializers.GlorotUniform())
            ]
        )

        h, w, c = img_shape
        encoder_conv_layers = 3
        decoder_starting_dims = (h // (2 ** encoder_conv_layers), w // (2 ** encoder_conv_layers), 128)

        self.decoder = keras.Sequential(
            [
                layers.Dense(np.prod(decoder_starting_dims), activation="relu", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Reshape(decoder_starting_dims),
                layers.Conv2DTranspose(128, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2DTranspose(64, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2DTranspose(32, 3, activation="relu", strides=2, padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform()),
                layers.Conv2D(c, 3, activation="sigmoid", padding="same", kernel_initializer=tf.keras.initializers.GlorotUniform())
            ]
        )

    def encode(self, x):
```

```

        mean_log_var = self.encoder(x)
        mean, log_var = tf.split(mean_log_var, num_or_size_splits=2, axis=1)
        return mean, log_var

    def reparameterize(self, mean, log_var):
        eps = tf.random.normal(shape=tf.shape(mean))
        return mean + eps * tf.exp(log_var * 0.5)

    def decode(self, z):
        return self.decoder(z)

    def call(self, x):
        mean, log_var = self.encode(x)
        z = self.reparameterize(mean, log_var)
        reconstructed = self.decode(z)

        # Add KL divergence loss as a model metric
        beta = 0.0001 # add beta
        kl_loss = beta * (-0.5 * tf.reduce_mean(
            tf.reduce_sum(1 + log_var - tf.square(mean) - tf.exp(log_var), axis=1)
        ))
        self.add_loss(kl_loss)

        return reconstructed

# Main training function
def train_bvae(dataset, img_shape=(224, 224, 3), latent_dim=64, epochs=20, learning_rate=0.001):
    # Create VAE
    bvae = BVAE(img_shape, latent_dim)
    bvae.compile(optimizer=keras.optimizers.Adam(learning_rate=learning_rate), loss="mse", metrics=["mse"])

    # Train the model
    history = bvae.fit(dataset, dataset, batch_size=BATCH_SIZE, epochs=epochs, verbose=VERBOSE)

    return bvae, history

# Generate reconstructions from the VAE
def generate_bvae_reconstructions(bvae, test_images, n=8):
    # Get test images
    test_sample = test_images[:n]

    reconstructed = bvae.predict(test_sample, verbose=VERBOSE)

```

```
    return test_sample, reconstructed

def visualize_bvae_reconstructions(originals, reconstructions, n=8):
    plt.figure(figsize=(14, 4))
    for i in range(min(n, len(originals))):
        # Original
        ax = plt.subplot(2, n, i + 1)
        plt.imshow(originals[i])
        plt.title("Original")
        plt.axis("off")

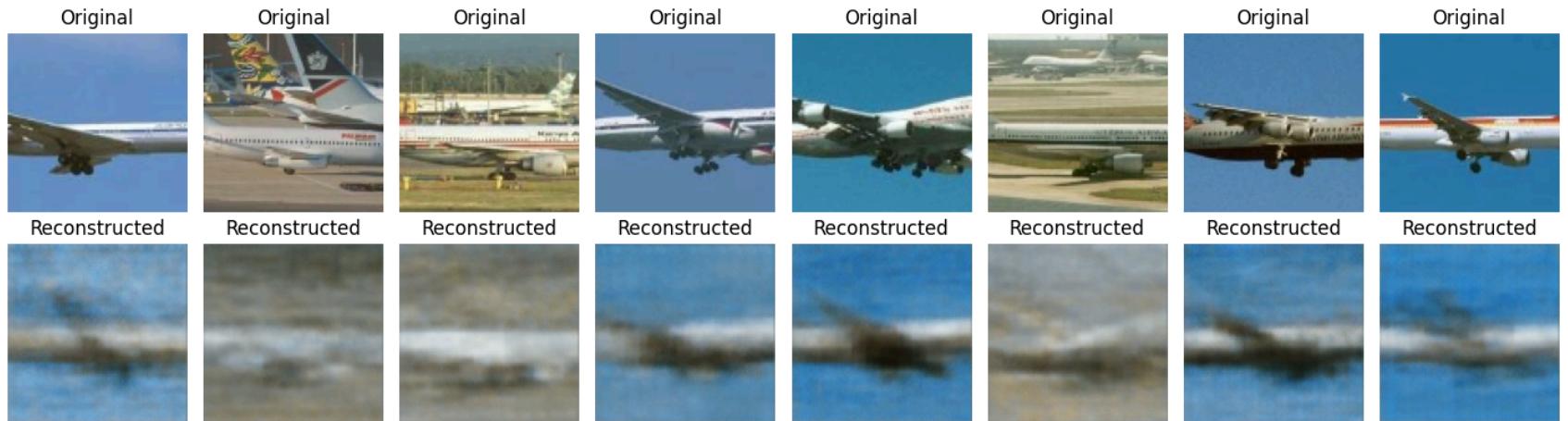
        # Reconstruction
        ax = plt.subplot(2, n, i + 1 + n)
        plt.imshow(reconstructions[i])
        plt.title("Reconstructed")
        plt.axis("off")

    plt.tight_layout()
    plt.show()
```

Train the β -VAE model

```
In [19]: start_time = time.time()
bvae, history = train_bvae(normal_images, img_shape=(IMAGE_SIZE, IMAGE_SIZE, 3), latent_dim=LATENT_DIM, epochs=100)
bvae_training_time = time.time() - start_time

original, reconstructed = generate_bvae_reconstructions(bvae, normal_images, n=8)
visualize_bvae_reconstructions(original, reconstructed, n=8)
```



Visualize result of β -VAE

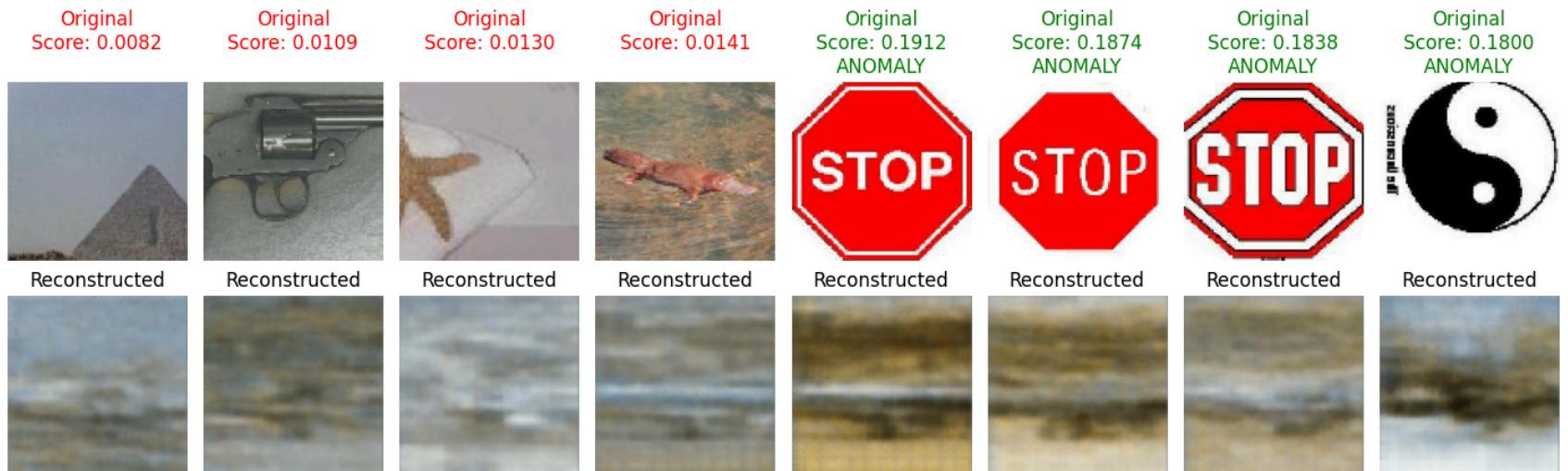
```
In [20]: def detect_anomalies_bvae(bvae, images, threshold):
    reconstructed = bvae.predict(images, verbose=VERBOSE)
    # Calculate MSE for each image
    mse = np.mean(np.square(images - reconstructed), axis=(1, 2, 3))
    # Determine if each image is an anomaly based on threshold
    anomalies = mse > threshold

    return mse, anomalies, reconstructed

# Detect anomalies with VAE
bvae_scores, bvae_anomalies, bvae_reconstructions = detect_anomalies_bvae(bvae, anomaly_images, threshold=threshold)
visualize_anomalies(anomaly_images, bvae_reconstructions, bvae_scores, bvae_anomalies, "β-VAE", is_normal=False)
print(f"β-VAE detected {np.sum(bvae_anomalies)} anomalies from {len(anomaly_images)} images")

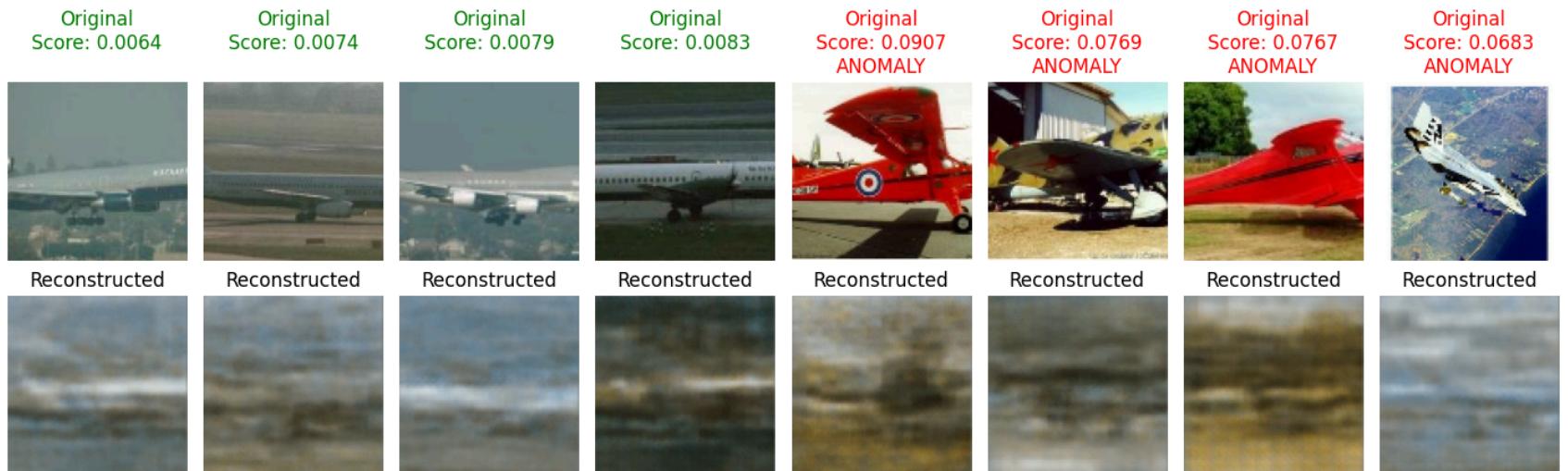
normal_scores, normal_anomalies, normal_reconstructions = detect_anomalies_bvae(bvae, normal_images, threshold=threshold)
visualize_anomalies(normal_images, normal_reconstructions, normal_scores, normal_anomalies, "β-VAE", is_normal=True)
print(f"β-VAE misclassified {np.sum(normal_anomalies)} normal images from {len(normal_images)} images")
```

Anomaly Detection with β -VAE (Threshold: 0.03)



β -VAE detected 699 anomalies from 800 images

Anomaly Detection with β -VAE (Threshold: 0.03)



β -VAE disclassified 154 normal images from 800 images

Hyperparameter tuning on β -VAE Threshold

```
In [21]: # Calculate metrics for VAE
bvae_scores, _, _ = detect_anomalies_bvae(bvae, test_X, BVAE_THRESHOLD)
```

```
bvae_metrics = calculate_metrics(test_y, bvae_scores, BVAE_THRESHOLD)

thresholds = [0.005, 0.01, 0.02, 0.025, 0.03, 0.035, 0.04, 0.05]

# Store metrics for each threshold
all_metrics = {}

for threshold in thresholds:
    # Get metrics for current threshold
    metrics = calculate_metrics(test_y, bvae_scores, threshold)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]

    # Initialize dictionary on first run
    if not all_metrics:
        all_metrics = {row[0]: [] for row in metric_data}

    # Add values for this threshold
    for row in metric_data:
        all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"{t}" for t in thresholds]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("β-VAE Performance Across Thresholds:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))
```

β -VAE Performance Across Thresholds:

Metric	0.005	0.01	0.02	0.025	0.03	0.035	0.04	0.05
Accuracy	0.5	0.5081	0.7244	0.8	0.8406	0.8438	0.83	0.7837
Precision	0.5	0.5041	0.6502	0.7386	0.8195	0.8618	0.894	0.9416
Recall	1	0.9988	0.9712	0.9287	0.8738	0.8187	0.7488	0.605
F1 Score	0.6667	0.67	0.7789	0.8228	0.8457	0.8397	0.815	0.7367
✓ True Positives	800	799	777	743	699	655	599	484
✗ False Positives	800	786	418	263	154	105	71	30
✓ True Negatives	0	14	382	537	646	695	729	770
✗ False Negatives	0	1	23	57	101	145	201	316

Plot performance metrics of β -VAE

```
In [22]: combined = combine_metrics(ae_metrics, vae_metrics, bvae_metrics, metric_names=['AE', 'VAE', ' $\beta$ -VAE'])
# Display metrics in a table format
print("\nAE vs. VAE vs.  $\beta$ -VAE Performance Metrics:")
print(tabulate(combined, headers="firstrow", tablefmt="grid", numalign="right"))

# Calculate ROC AUC score
fpr_bvae, tpr_bvae, thresholds = roc_curve(test_y, bvae_scores)

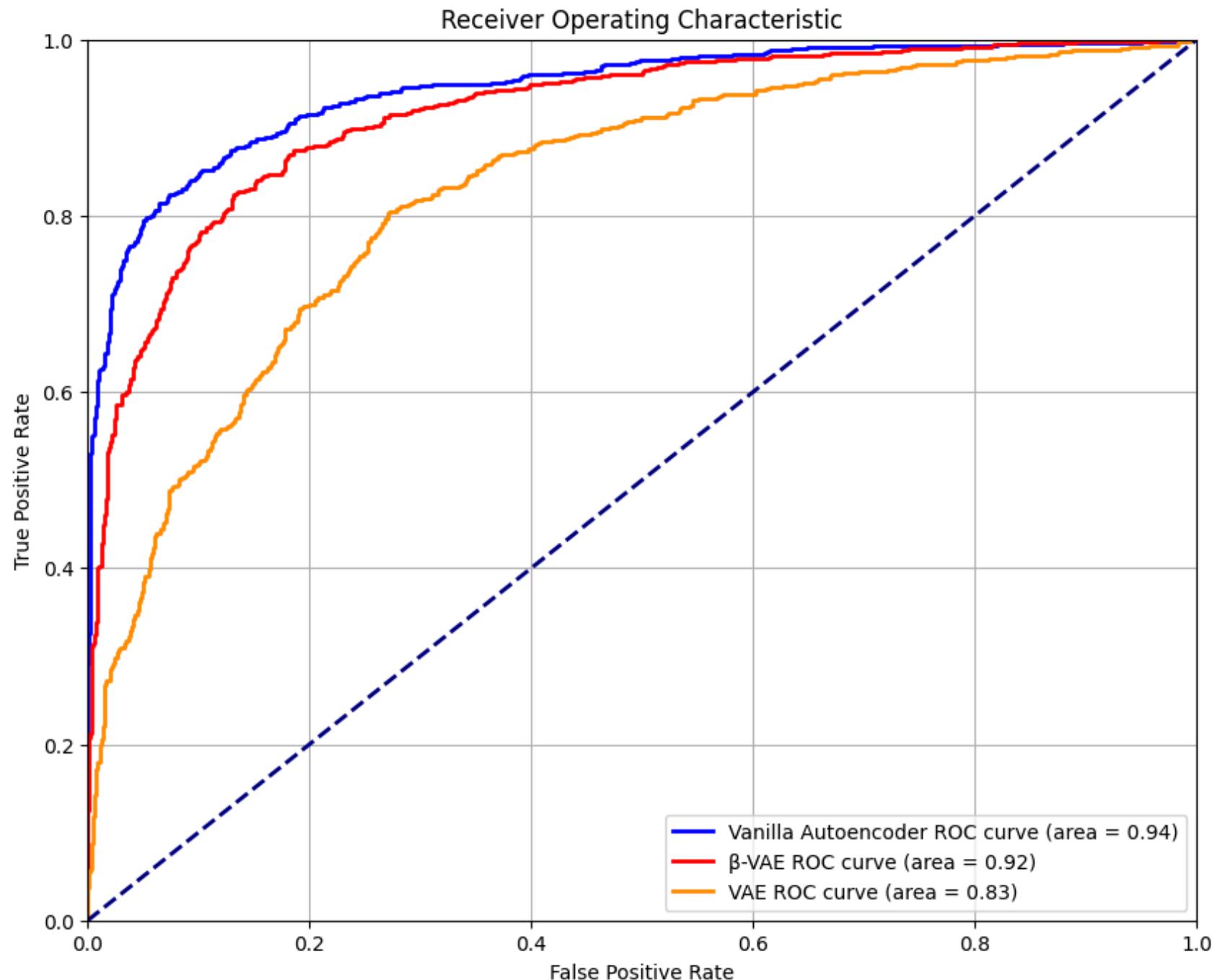
# Calculate Area Under Curve (AUC)
auc_bvae = auc(fpr_bvae, tpr_bvae)

# Plot ROC curve
plt.figure(figsize=(10, 8))
plt.plot(fpr_ae, tpr_ae, color='blue', lw=2, label=f'Vanilla Autoencoder ROC curve (area = {auc_ae:.2f})')
plt.plot(fpr_bvae, tpr_bvae, color='red', lw=2, label=f' $\beta$ -VAE ROC curve (area = {auc_bvae:.2f})')
plt.plot(fpr_vae, tpr_vae, color='darkorange', lw=2, label=f'VAE ROC curve (area = {auc_vae:.2f})')
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
```

```
plt.ylim([0.0, 1.0])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic')
plt.legend(loc="lower right")
plt.grid(True)
plt.show()
```

AE vs. VAE vs. β -VAE Performance Metrics:

Metric	AE	VAE	β -VAE
Accuracy	0.8725	0.7175	0.8406
Precision	0.937	0.6608	0.8195
Recall	0.7987	0.8938	0.8738
F1 Score	0.8623	0.7598	0.8457
✓ True Positives	639	715	699
✗ False Positives	43	367	154
✓ True Negatives	757	433	646
✗ False Negatives	161	85	101



6. U-Net Autoencoder

Build the U-Net Autoencoder model

```
In [23]: # Build the U-Net Autoencoder model as a subclass of keras.Model
class UnetAutoencoder(tf.keras.Model):
    def __init__(self, input_shape, **kwargs):
        super(UnetAutoencoder, self).__init__(**kwargs)
        # Define the functional model using the Keras Functional API
        inputs = Input(shape=input_shape)

        # --- Encoder ---
        # Block 1
        c1 = Conv2D(32, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        c1 = Conv2D(32, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        p1 = MaxPooling2D((2, 2))(c1)

        # Block 2
        c2 = Conv2D(64, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        c2 = Conv2D(64, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        p2 = MaxPooling2D((2, 2))(c2)

        # Block 3
        c3 = Conv2D(128, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        c3 = Conv2D(128, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        p3 = MaxPooling2D((2, 2))(c3)

        # --- Bottleneck ---
        b = Conv2D(256, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        b = Conv2D(256, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.

        # --- Decoder ---
        # Up Block 1
        u3 = UpSampling2D((2, 2))(b)
        u3 = concatenate([u3, c3])
        c4 = Conv2D(128, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.
        c4 = Conv2D(128, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializer.

        # Up Block 2
        u2 = UpSampling2D((2, 2))(c4)
```

```

u2 = concatenate([u2, c2])
c5 = Conv2D(64, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializers.he_normal())
c5 = Conv2D(64, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializers.he_normal())

# Up Block 3
u1 = UpSampling2D((2, 2))(c5)
u1 = concatenate([u1, c1])
c6 = Conv2D(32, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializers.he_normal())
c6 = Conv2D(32, (3, 3), activation='relu', padding='same', kernel_initializer=tf.keras.initializers.he_normal()

# Output layer: reconstruct the image with pixel values in [0, 1]
outputs = Conv2D(input_shape[-1], (1, 1), activation='sigmoid', kernel_initializer=tf.keras.initializers.he_normal())

# Create the underlying model
self.model = Model(inputs, outputs)
self.model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.001), loss='mse')

def call(self, x):
    return self.model(x)

def train_unet_autoencoder(dataset, input_shape, epochs=20, batch_size=64):
    """
    Train the U-Net autoencoder on the provided dataset.

    Args:
        dataset: Training dataset (e.g., normal_images).
        input_shape: Shape of the input images, e.g., (96, 96, 3).
        epochs: Number of training epochs.
        batch_size: Batch size for training.

    Returns:
        unet: Trained U-Net autoencoder.
        history: Training history.
    """
    unet = UnetAutoencoder(input_shape)
    history = unet.model.fit(dataset, dataset, epochs=epochs, batch_size=batch_size, verbose=0)
    return unet, history

# Example usage:
# unet_autoencoder, unet_history = train_unet_autoencoder(normal_images, input_shape=(96, 96, 3))

```

```
def generate_unet_reconstructions(unet, test_images, n=8):
    """
    Generate reconstructions from the U-Net autoencoder.

    Args:
        unet: Trained U-Net autoencoder.
        test_images: Images to reconstruct.
        n: Number of images to display.

    Returns:
        test_sample: Original images.
        reconstructed: Reconstructed images.
    """
    test_sample = test_images[:n]
    reconstructed = unet.model.predict(test_sample, verbose=0)
    return test_sample, reconstructed

def visualize_unet_reconstructions(originals, reconstructions, n=8):
    """
    Visualize original vs. reconstructed images for the U-Net autoencoder.
    """
    import matplotlib.pyplot as plt
    plt.figure(figsize=(14, 4))
    for i in range(n):
        # Original image
        ax = plt.subplot(2, n, i + 1)
        plt.imshow(originals[i])
        plt.title("Original")
        plt.axis("off")

        # Reconstructed image
        ax = plt.subplot(2, n, i + 1 + n)
        plt.imshow(reconstructions[i])
        plt.title("Reconstructed")
        plt.axis("off")

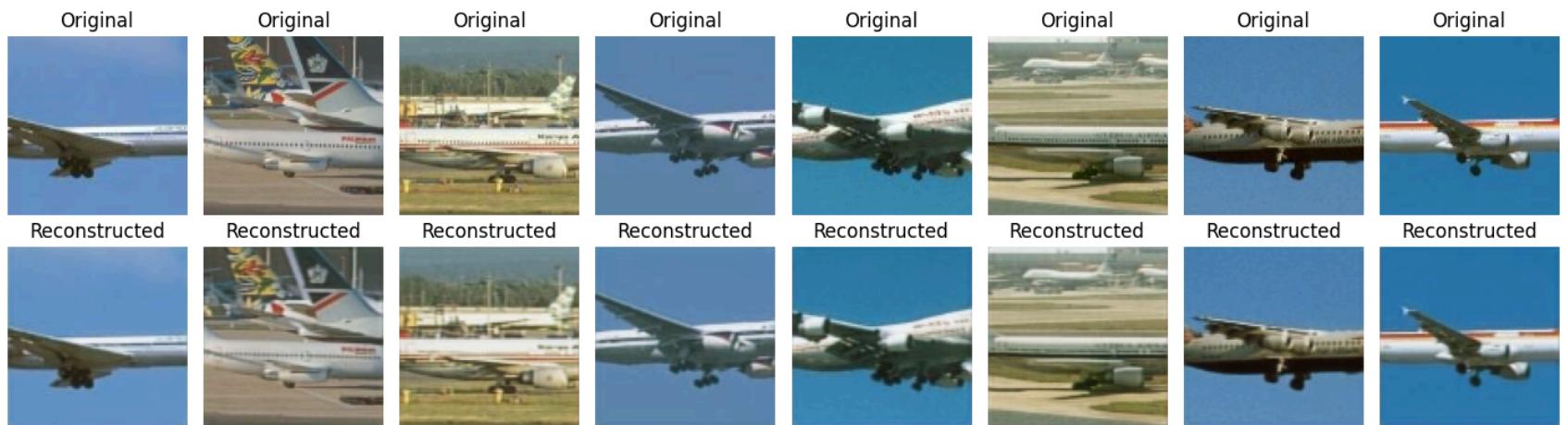
    plt.tight_layout()
    plt.show()

# Example usage:
# original, unet_reconstructed = generate_unet_reconstructions(unet_autoencoder, normal_images, n=8)
# visualize_unet_reconstructions(original, unet_reconstructed, n=8)
```

Train the U-Net Autoencoder model

```
In [24]: # Train U-Net Autoencoder
start_time = time.time()
unet_autoencoder, unet_history = train_unet_autoencoder(normal_images, input_shape=(IMAGE_SIZE, IMAGE_SIZE)
unet_training_time = time.time() - start_time

unet_original, unet_reconstructions = generate_unet_reconstructions(unet_autoencoder, normal_images, n=8)
visualize_unet_reconstructions(unet_original, unet_reconstructions, n=8)
```



Detect anomalies with U-Net Autoencoder

```
In [25]: def detect_anomalies_unet(unet, images, threshold):
    """
    Detect anomalies using the U-Net autoencoder.

    Args:
        unet: Trained U-Net autoencoder.
        images: Images to evaluate.
        threshold: Reconstruction error threshold.

    Returns:
        mse: Reconstruction error for each image.
        anomalies: Boolean array indicating anomalies.
        reconstructions: Reconstructed images.
    """

    # ... (implementation details)
```

```

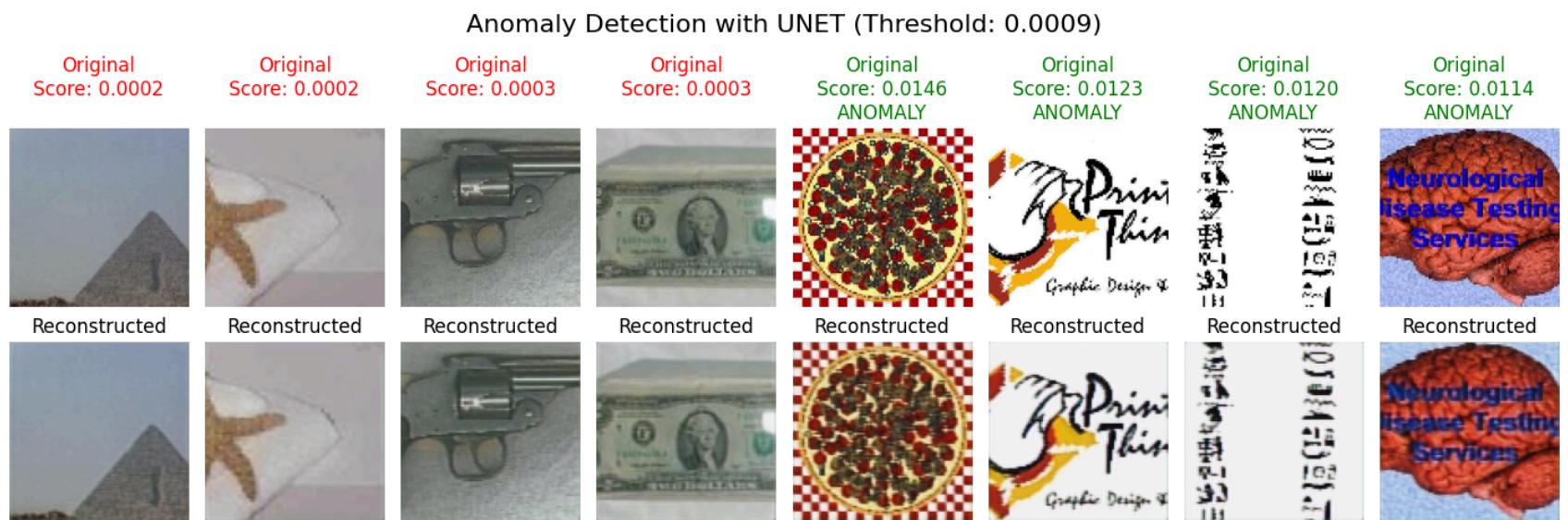
    .....
    reconstructions = unet.model.predict(images, verbose=0)
    mse = tf.reduce_mean(tf.square(images - reconstructions), axis=[1, 2, 3]).numpy()
    anomalies = mse > threshold
    return mse, anomalies, reconstructions

# Example usage:
# unet_scores, unet_anomalies, unet_reconstructions = detect_anomalies_unet(unet_autoencoder, test_X, THRESHOLD)

# Detect anomalies with Unet
unet_scores, unet_anomalies, unet_reconstructions = detect_anomalies_unet(unet_autoencoder, anomaly_images)
visualize_anomalies(anomaly_images, unet_reconstructions, unet_scores, unet_anomalies, "UNET", is_normal_class=False)
print(f"UNET detected {np.sum(unet_anomalies)} anomalies from {len(anomaly_images)} images")

normal_scores, normal_anomalies, normal_reconstructions = detect_anomalies_unet(unet_autoencoder, normal_images)
visualize_anomalies(normal_images, normal_reconstructions, normal_scores, normal_anomalies, "UNET", is_normal_class=True)
print(f"UNET disclassified {np.sum(normal_anomalies)} normal images from {len(normal_images)} images")

```



UNET detected 695 anomalies from 800 images

Anomaly Detection with UNET (Threshold: 0.0009)



UNET disclassified 102 normal images from 800 images

✖ Hyperparameter Tuning on U-Net Autoencoder in various threshold

```
In [26]: # Detect anomalies using the U-Net autoencoder
unet_scores, _, _ = detect_anomalies_unet(unet_autoencoder, test_X, UNET_THRESHOLD)
unet_metrics = calculate_metrics(test_y, unet_scores, UNET_THRESHOLD)

thresholds = [0.03, 0.0025, 0.0026, 0.0027, 0.0028, 0.0029, 0.0030]

# Store metrics for each threshold
all_metrics = {}

for threshold in thresholds:
    # Get metrics for current threshold
    metrics = calculate_metrics(test_y, unet_scores, threshold)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]

    # Initialize dictionary on first run
    if not all_metrics:
        all_metrics = {row[0]: [] for row in metric_data}
```

```

# Add values for this threshold
for row in metric_data:
    all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"{t}" for t in thresholds]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("U-Net Autoencoder Performance Across Thresholds:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))

```

U-Net Autoencoder Performance Across Thresholds:

Metric	0.03	0.0025	0.0026	0.0027	0.0028	0.0029	0.003
Accuracy	0.5	0.6775	0.6694	0.655	0.6388	0.6256	0.6169
Precision	0	0.9965	0.9963	0.996	0.9955	0.9951	0.9947
Recall	0	0.3563	0.34	0.3113	0.2787	0.2525	0.235
F1 Score	0	0.5249	0.507	0.4743	0.4355	0.4028	0.3802
✓ True Positives	0	285	272	249	223	202	188
✗ False Positives	0	1	1	1	1	1	1
✓ True Negatives	800	799	799	799	799	799	799
✗ False Negatives	800	515	528	551	577	598	612

/opt/anaconda3/envs/tensorflow/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, f'{metric.capitalize()} is', len(result))

```
In [27]: thresholds = [0.03, 0.001, 0.0015, 0.0020, 0.0021, 0.0022, 0.0023, 0.0024]

# Store metrics for each threshold
all_metrics = {}

for threshold in thresholds:
    # Get metrics for current threshold
    metrics = calculate_metrics(test_y, unet_scores, threshold)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]

    # Initialize dictionary on first run
    if not all_metrics:
        all_metrics = {row[0]: [] for row in metric_data}

    # Add values for this threshold
    for row in metric_data:
        all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"{t}" for t in thresholds]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("U-Net Autoencoder Performance Across Thresholds:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))
```

U-Net Autoencoder Performance Across Thresholds:

Metric	0.03	0.001	0.0015	0.002	0.0021	0.0022	0.0023	0.0024
Accuracy	0.5	0.8712	0.8087	0.7469	0.7281	0.7131	0.7025	0.6887
Precision	0	0.8992	0.9714	0.995	0.9946	0.9942	0.9969	0.9967
Recall	0	0.8363	0.6362	0.4963	0.4587	0.4288	0.4062	0.3787
F1 Score	0	0.8666	0.7689	0.6622	0.6279	0.5991	0.5773	0.5489
✓ True Positives	0	669	509	397	367	343	325	303
✗ False Positives	0	75	15	2	2	2	1	1
✓ True Negatives	800	725	785	798	798	798	799	799
✗ False Negatives	800	131	291	403	433	457	475	497

```
/opt/anaconda3/envs/tensorflow/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.
```

```
_warn_prf(average, modifier, f'{metric.capitalize()} is', len(result))
```

In [28]: thresholds = [0.03, 0.0001, 0.0002, 0.0003, 0.0004, 0.0005, 0.0006, 0.0007, 0.0008, 0.0009, 0.001]

```
# Store metrics for each threshold
all_metrics = {}

for threshold in thresholds:
    # Get metrics for current threshold
    metrics = calculate_metrics(test_y, unet_scores, threshold)

    # Extract metric names and values (skip header row)
    metric_data = metrics[1:]

    # Initialize dictionary on first run
    if not all_metrics:
        all_metrics = {row[0]: [] for row in metric_data}

    # Add values for this threshold
    for row in metric_data:
        all_metrics[row[0]].append(row[1])
```

```
for row in metric_data:
    all_metrics[row[0]].append(row[1])

# Build the table
header_row = ["Metric"] + [f"{{t}}" for t in thresholds]
combined_table = [header_row]

# Add each metric row
for metric, values in all_metrics.items():
    combined_table.append([metric] + values)

# Display results
print("U-Net Autoencoder Performance Across Thresholds:")
print(tabulate(combined_table, headers="firstrow", tablefmt="grid", numalign="right"))
```

U-Net Autoencoder Performance Across Thresholds:

Metric	0.03	0.0001	0.0002	0.0003	0.0004	0.0005	0.0006	0.0007
Accuracy	0.5	0.5	0.5038	0.5456	0.6569	0.7681	0.8269	0.8538
Precision	0	0.5	0.5019	0.5241	0.5941	0.6883	0.7581	0.805
Recall	0	1	1	0.9938	0.99	0.98	0.96	0.9337
F1 Score	0	0.6667	0.6683	0.6862	0.7426	0.8087	0.8472	0.8646
✓ True Positives	0	800	800	795	792	784	768	747
x False Positives	0	800	794	722	541	355	245	181
✓ True Negatives	800	0	6	78	259	445	555	619
x False Negatives	800	0	0	5	8	16	32	53

```
/opt/anaconda3/envs/tensorflow/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.  
    _warn_prf(average, modifier, f"{{metric.capitalize()}} is", len(result))
```

```
In [29]: import numpy as np  
from sklearn.metrics import confusion_matrix  
  
def calculate_roc_metrics(y_true, scores, thresholds):  
    """  
    Calculate ROC metrics for a list of thresholds.  
  
    Args:  
        y_true (array-like): True binary labels (0 for normal, 1 for anomaly).  
        scores (array-like): Reconstruction errors or model scores.  
        thresholds (list or array): Thresholds to evaluate.  
  
    Returns:  
        list of dicts: Each dict contains:  
            - 'threshold': Threshold value.  
            - 'TPR': True Positive Rate.  
            - 'FPR': False Positive Rate.  
            - 'YoudenIndex': TPR - FPR.  
            - 'Accuracy': Overall accuracy.  
            - 'Precision': Precision.  
            - 'TP': True Positives.  
            - 'FP': False Positives.  
            - 'TN': True Negatives.  
            - 'FN': False Negatives.  
    """  
    results = []  
    y_true = np.array(y_true)  
    scores = np.array(scores)  
  
    for thr in thresholds:  
        # Create binary predictions based on the threshold  
        y_pred = (scores > thr).astype(int)  
        tn, fp, fn, tp = confusion_matrix(y_true, y_pred).ravel()  
  
        TPR = tp / (tp + fn) if (tp + fn) > 0 else 0  
        FPR = fp / (fp + tn) if (fp + tn) > 0 else 0  
        YoudenIndex = TPR - FPR
```

```

Accuracy = (tp + tn) / len(y_true)
Precision = tp / (tp + fp) if (tp + fp) > 0 else 0

results.append({
    'threshold': thr,
    'TPR': TPR,
    'FPR': FPR,
    'YoudenIndex': YoudenIndex,
    'Accuracy': Accuracy,
    'Precision': Precision,
    'TP': tp,
    'FP': fp,
    'TN': tn,
    'FN': fn
})
return results

# Example usage:
# y_true: true labels from your test dataset.
# scores: reconstruction errors from your autoencoder (e.g., U-Net model).
# thresholds: a list of thresholds you want to evaluate.
thresholds = [0.03, 0.0001, 0.0002, 0.0003, 0.0004, 0.0005, 0.0006, 0.0007, 0.0008, 0.0009, 0.001]
roc_metrics = calculate_roc_metrics(test_y, unet_scores, thresholds)

# To print the results in a tabular format:
from tabulate import tabulate
print(tabulate(roc_metrics, headers="keys", tablefmt="grid"))

```

threshold	TPR	FPR	YoudenIndex	Accuracy	Precision	TP	FP	TN	FN
0.03	0	0	0	0.5	0	0	0	800	800
0.0001	1	1	0	0.5	0.5	800	800	0	0
0.0002	1	0.9925	0.0075	0.50375	0.501882	800	794	6	0
0.0003	0.99375	0.9025	0.09125	0.545625	0.524061	795	722	78	5
0.0004	0.99	0.67625	0.31375	0.656875	0.594149	792	541	259	8
0.0005	0.98	0.44375	0.53625	0.768125	0.688323	784	355	445	16
0.0006	0.96	0.30625	0.65375	0.826875	0.758144	768	245	555	32
0.0007	0.93375	0.22625	0.7075	0.85375	0.804957	747	181	619	53
0.0008	0.89625	0.16875	0.7275	0.86375	0.841549	717	135	665	83
0.0009	0.86875	0.1275	0.74125	0.870625	0.87202	695	102	698	105
0.001	0.83625	0.09375	0.7425	0.87125	0.899194	669	75	725	131

```
In [30]: # Detect anomalies with Unet
unet_scores, unet_anomalies, unet_reconstructions = detect_anomalies_unet(unet_autoencoder, anomaly_images)
visualize_anomalies(anomaly_images, unet_reconstructions, unet_scores, unet_anomalies, "UNET", is_normal_class=False)
print(f"UNET detected {np.sum(unet_anomalies)} anomalies from {len(anomaly_images)} images")

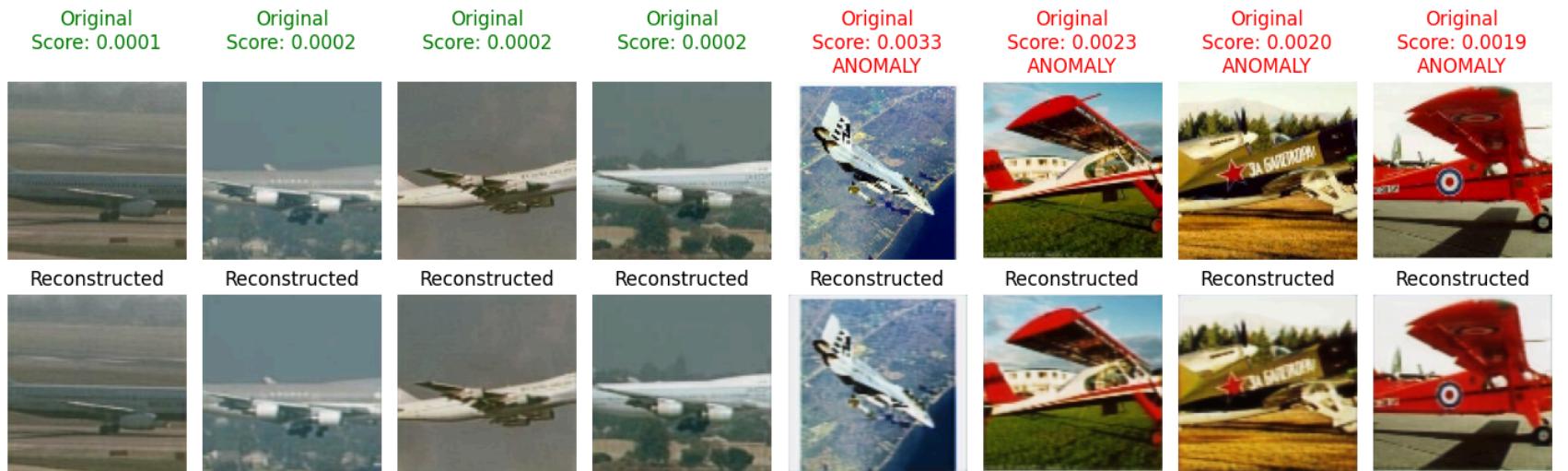
normal_scores, normal_anomalies, normal_reconstructions = detect_anomalies_unet(unet_autoencoder, normal_images)
visualize_anomalies(normal_images, normal_reconstructions, normal_scores, normal_anomalies, "UNET", is_normal_class=True)
print(f"UNET disclassified {np.sum(normal_anomalies)} normal images from {len(normal_images)} images")
```

Anomaly Detection with UNET (Threshold: 0.0009)



UNET detected 695 anomalies from 800 images

Anomaly Detection with UNET (Threshold: 0.0009)



UNET disclassified 102 normal images from 800 images

```
In [31]: # Detect anomalies using the U-Net autoencoder
unet_scores, _, _ = detect_anomalies_unet(unet_autoencoder, test_X, UNET_THRESHOLD)
unet_metrics = calculate_metrics(test_y, unet_scores, UNET_THRESHOLD)

# Combine metrics for all four models
```

```

combined = combine_metrics(ae_metrics, vae_metrics, bvae_metrics, unet_metrics,
                           metric_names=['AE', 'VAE', 'β-VAE', 'U-Net AE'])
print("\nAE vs. VAE vs. β-VAE vs. U-Net AE Performance Metrics:")
print(tabulate(combined, headers="firstrow", tablefmt="grid", numalign="right"))

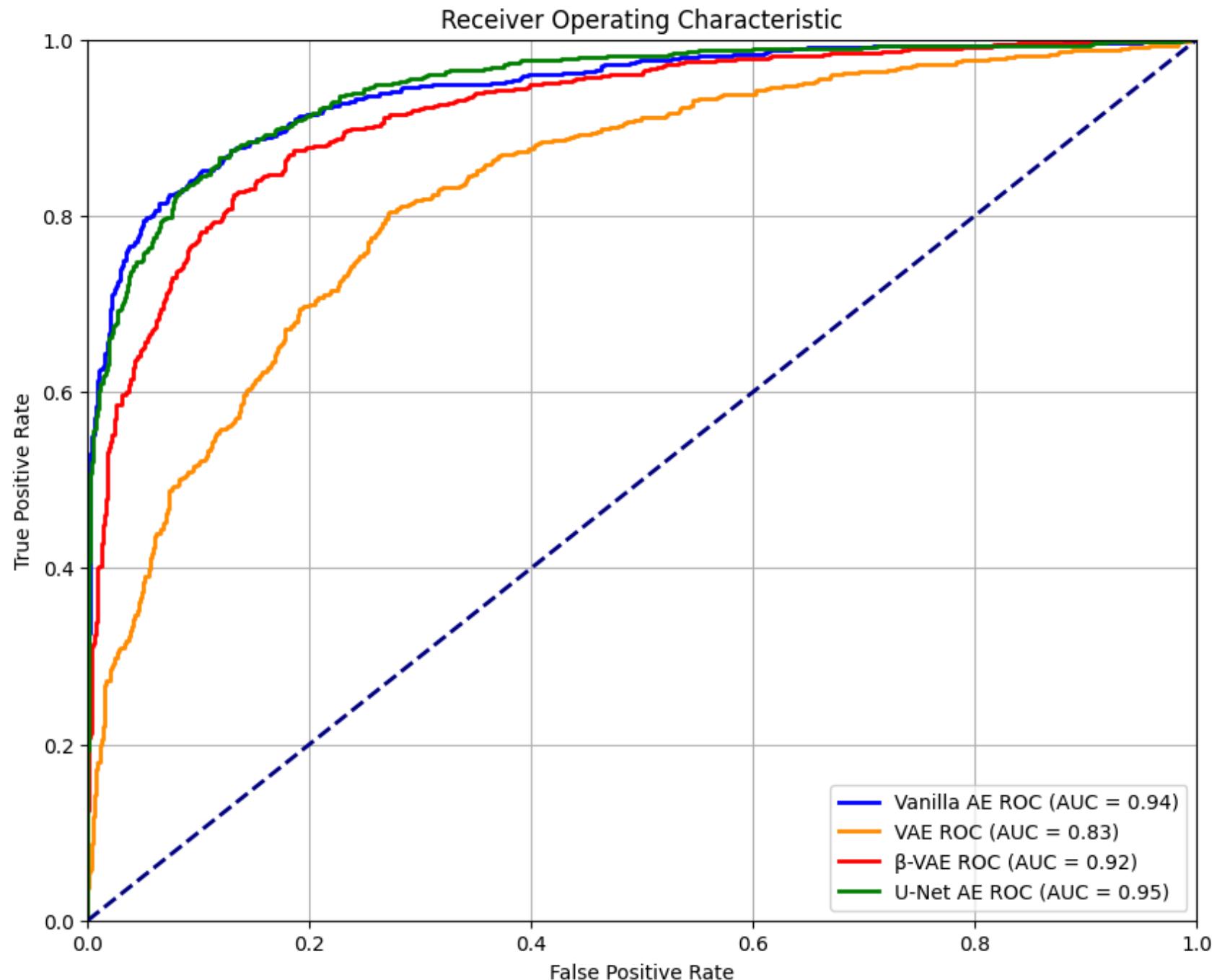
# Compute ROC curve for U-Net autoencoder
fpr_unet, tpr_unet, _ = roc_curve(test_y, unet_scores)
auc_unet = auc(fpr_unet, tpr_unet)

# Plot ROC curves for all models
plt.figure(figsize=(10, 8))
plt.plot(fpr_ae, tpr_ae, color='blue', lw=2, label=f'Vanilla AE ROC (AUC = {auc_ae:.2f})')
plt.plot(fpr_vae, tpr_vae, color='darkorange', lw=2, label=f'VAE ROC (AUC = {auc_vae:.2f})')
plt.plot(fpr_bvae, tpr_bvae, color='red', lw=2, label=f'β-VAE ROC (AUC = {auc_bvae:.2f})')
plt.plot(fpr_unet, tpr_unet, color='green', lw=2, label=f'U-Net AE ROC (AUC = {auc_unet:.2f})')
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic')
plt.legend(loc="lower right")
plt.grid(True)
plt.show()

```

AE vs. VAE vs. β -VAE vs. U-Net AE Performance Metrics:

Metric	AE	VAE	β -VAE	U-Net AE
Accuracy	0.8725	0.7175	0.8406	0.8706
Precision	0.937	0.6608	0.8195	0.872
Recall	0.7987	0.8938	0.8738	0.8688
F1 Score	0.8623	0.7598	0.8457	0.8704
✓ True Positives	639	715	699	695
✗ False Positives	43	367	154	102
✓ True Negatives	757	433	646	698
✗ False Negatives	161	85	101	105

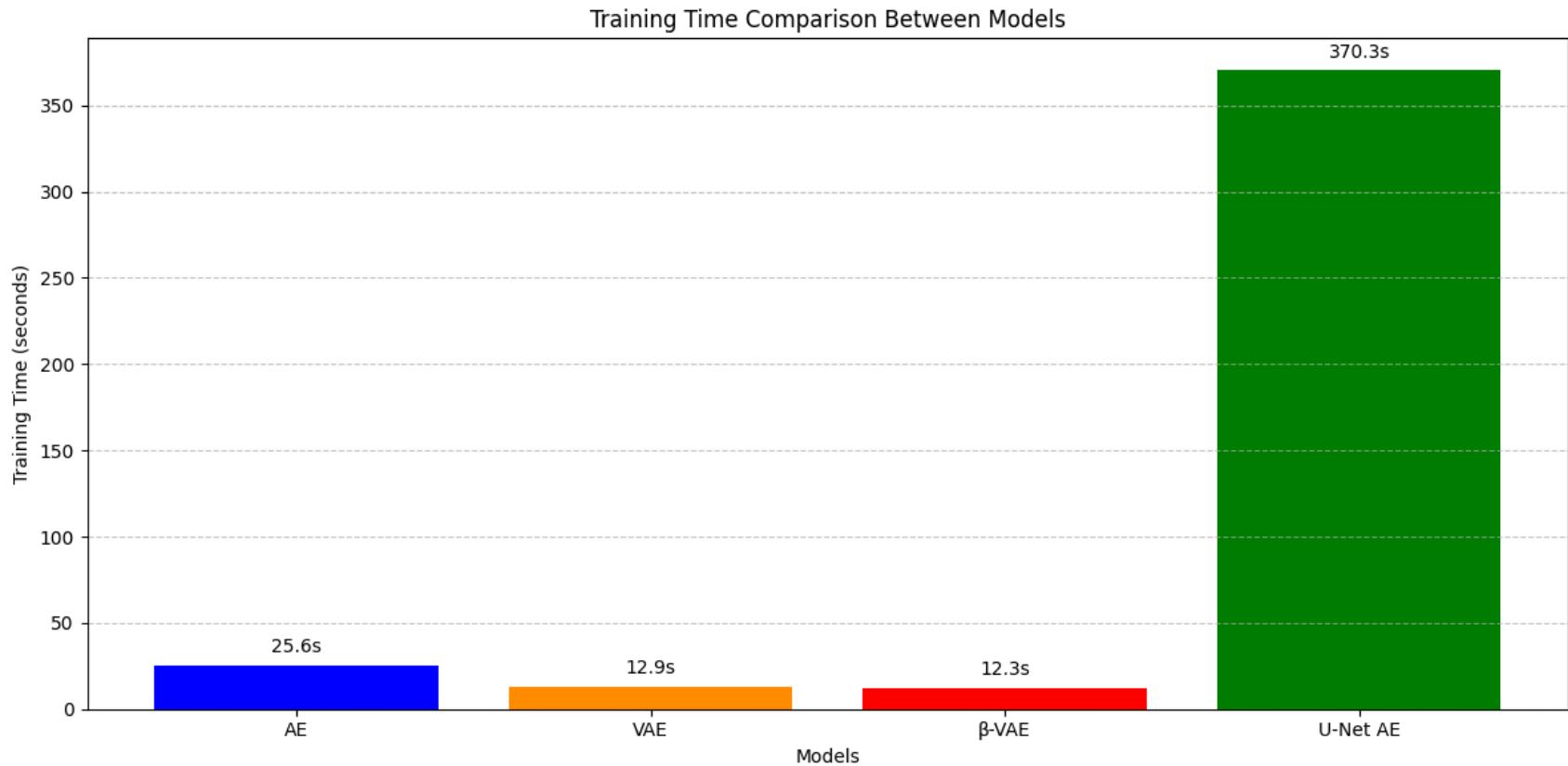


```
In [32]: # Compare training time for each model
plt.figure(figsize=(12, 6))
models = ['AE', 'VAE', 'β-VAE', 'U-Net AE']
# Replace these values with the actual training times from your experiments
training_times = [
    ae_training_time,
    vae_training_time,
    bvae_training_time,
    unet_training_time
]

# Create bar chart
bars = plt.bar(models, training_times, color=['blue', 'darkorange', 'red', 'green'])

# Add values on top of bars
for bar in bars:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height + 5,
             f'{height:.1f}s', ha='center', va='bottom')

plt.title('Training Time Comparison Between Models')
plt.xlabel('Models')
plt.ylabel('Training Time (seconds)')
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```



Conclusion

1. U-NET AE is the best model for anomaly detection on airplane class in this dataset with F-1 score of **0.8704**.
2. Vanilla Autoencoder is the second best model for anomaly detection on airplane class in this dataset with F-1 score of **0.8623**.
3. β -VAE and VAE performance poorly on anomaly detection bases on their architecture which add KL divergence to loss function, smooth the latent space distribution and reconstruct the image make it has high reconstruction error.
4. U-Net AE is very good at reconstructing the image because of its skip connection architecture which is able to preserve the fine details of the image. Result in very low reconstruction error threshold at **0.0009** comparing to other Autoencoder between **0.03 - 0.05**.
5. Regarding training time, other Autoencoders took less than 26 seconds, but U-Net AE took 370.3 seconds (~6mins)

6. U-NET AE has balanced between precision and recall, while Vanilla Autoencoder has higher Precision but lower Recall making it more suitable if miscalssify normal image has highly cost.