# An Improved VGG16 Model Based on Complex Invariant Descriptors for Medical Images Classification

Mohamed Amine Mezghich, Dorsaf Hmida, Taha Mustapha Nahdi and Faouzi Ghorbel

*GRIFT Research Group, CRISTAL laboratory, ENSI, Tunisia*

Keywords: Complex Moments, Invariant Descriptors, Stability, Completeness, Deep Learning, Classification.

Abstract: In this paper, we intent to present an improved VGG16 deep learning model based on an invariant and complete set of descriptors constructed by a linear combination of complex moments. First, the invariant features are studied to highlight it's stability and completeness properties over rigid transformations, noise and non rigid transformations. Then our proposed method to inject this family to the well know deep leaning VGG16 model is presented. Experimental results are satisfactory and the model accuracy is improved.

## 1 INTRODUCTION

Recently, images have captured great importance because they are the main source of information. It's used in many aspects to share information between different recipients in various fields such as military operations, medical imaging, artificial intelligence, social media and many other applications. Feature extraction technique consists in computing object descriptors with some interesting properties such as invariant to some geometric transformations, completeness which means discrimination between objects if they are different and finally stability over noise and shape distortion. It's the most crucial stage in every kind of classification system because the extracted features significantly influence the system performance. It is sometimes necessary to calculate descriptors directly from gray-scale objects when the external profile of the objects is not sufficiently discriminating. This problem is more complex than the previous one because the description must take into account information both on the contour and the internal texture of objects.

Many shape representation and description techniques have been developed in the past. Early works for computing invariant shape features were based on Fourier descriptors of the curve, e.g. (Zahn and Roskies, 1972) , (Persoon and Fu, 1977). (Ghorbel, 1998) introduced a joint topology and harmonic analysis formulation for the extraction of global shape descriptors which are invariant under a given group of geometrical transformations. The second family of object descriptors are based on geometric and complex moments of the image which take into account the gray-scale level of pixels. In (Hu, 1962), Hu introduced seven invariant descriptors based on a combination of centred and normalised geometric moments. This set is not complete which means that we can't reconstruct the original object from it's moments. Khotanzad and Hong introduced in (Khotanzad and Hong, 1990) Zernike moments which are invariant and complete set to rigid transformations. Fourier–Merlin moments were introduced in (Sheng and Shen, 1994) to extract invariant gray-scale features. In (Ghorbel, 1994), The author proposed the Analytical Fourier-Mellin Transform (AFMT) for invariant and complete shape description. A set of complete and stable invariants was then proposed and in (Derrode and Ghorbel, 2001), the authors proposed three methods to compute this set based on polar, log-polar and cartesian approximations. Flusser et al. proposed respectively in (Flusser, 2002) and (Suk and Flusser, 2003) a new set of descriptors based on complex moments which are invariant to rotation. In 2006, Ghorbel et al. (Ghorbel et al., 2006) proposed a systematic method to extract a complete set of similarity invariants (translation, rotation and scale), by means of some linear combinations of complex moments. Also the authors proposed an original method to reconstruct the image from its complex moments derived from the relationship between the Discrete Fourier Transform (DFT) and complex moments of the image. In recent years, many efforts have been made to apply convolutional neural networks (CNNs) on the medical imaging modality. In fact CNNs have already shown impressive perfor-

mance on the classification of natural images. The fusion of CNN with different techniques like local and global features have proven encourages and best results. (Sharma and Mehra, 2020) proposed two machine learning approaches. The first approach is based on handcrafted features which are extracted using Hu's moments, color histogram, and Haralick textures. The extracted features are then utilized to train the conventional classifiers, while the second approach is based on transfer learning where the pre-existing networks (VGG16, VGG19, and ResNet50) are utilized as feature extractor and as a baseline model. (Cao et al., 2020) proposed a ship recognition method based on morphological watershed image segmentation and Zernike moment to solve the problem of ship recognition in video images. The feature of ship image is extracted based on deep learning convolution neural network (CNN) and Zernike moments method. (Wang et al., 2020) developed a deep learning-based model for Drug-Target Interactions (DTIs) prediction. The proteins evolutionary features are extracted via Position Specific Scoring Matrix (PSSM) and Legendre Moment (LM). (Moung et al., 2021) proposed a fusion of a moment invariant (MI) method and a Deep Learning (DL) algorithm for feature extraction to address the instabilities in the existing COVID-19 classification models. (Ren et al., 2021) proposed a ship recognition approach based on Hu invariant moments and Convolutional Neural Network (CNN) to solve the issue with automatic recognition of ship images in video surveillance system.

In this work, an improved VGG16 model based on complex invariant descriptors is proposed for medical images classification. The proposed invariant features are added to the ones extracted by the considered CNN deep model in order to enhance the accuracy of classification. Our paper is organized as follows: In section 2, we introduce the set of invariant complex descriptors. By section 3, we recall the used VGG16 deep learning model then we present our approach to improve the accuracy of this model. Section 4 highlight experimental results. Finally in section 5, we summarize our work and provides brief insight into our planned future work.

## 2 INVARIANT FEATURES

In this paper, we are interested in invariant descriptors computed from image complex moments. It's a mathematical tools used in image processing to describe the intensity distribution in an image and it's a generalization of standard geometric moments used to extract features that are invariant under certain trans-

formations such as translation, rotation and scaling. First, we recall some basic terms that will be used to construct the invariant set. Let $f(x,y)$ be an image function having a bounded support and a finite nonzero integral. The complex moments of $f(x,y)$ are defined as:

$$c_f(p,q) = \iint_{\mathbb{R}^2} (x+\mathrm{i}y)^p (x-\mathrm{i}y)^q f(x,y)\mathrm{d}x\,\mathrm{d}y \quad (1)$$

where $p,q \in \mathbb{N}$. Eq.1 in polar coordinates becomes:

$$c_f(p,q) = \int_0^\infty \int_0^{2\pi} r^{p+q+1} \mathrm{e}^{\mathrm{i}(p-q)\theta} f(r,\theta)\mathrm{d}r\,\mathrm{d}\theta \quad (2)$$

In (Flusser, 2002) and (Flusser, 2006), a set of complete and independent set of rotation invariant descriptors is given as follows :

$$\forall p,q \in \mathbb{N}, \quad \Phi_f(p,q) = c_f(p_0-1,p_0)^{p-q} c_f(p,q) \quad (3)$$

And in (Ghorbel et al., 2006), the authors proposed a systematic method to extract a complete set of similarity invariant features to translation, rotation and scale by means of some linear combinations of complex moments. The complete formulation $I_f$ is defined as follows:

$$\forall p,q \in \mathbb{N}, \quad I_f(p,q) = \Gamma_f^{-(p+q+2)} \mathrm{e}^{-\mathrm{i}(p-q)\Theta_f} c_f(p,q) \quad (4)$$

with $\Theta_f = \arg(c_f(1,0))$ and $\Gamma_f = \sqrt{c_f(0,0)}$.

The completeness property is easy obtained by this formula:

$$\forall p,q \in \mathbb{N}, \quad c_f(p,q) = \Gamma_f^{p+q+2} \mathrm{e}^{\mathrm{i}(p-q)\Theta_f} I_f(p,q) \quad (5)$$

We will be based on this family to study first it's stability and robustness, then we will inject it to CNN based deep learning model in order to improve its accuracy.

## 3 AN IMPORVED VGG16 MODEL

In this section, we start by describing the CNN architecture which is utilised in this work , Sect. 3.1 then, we detail our proposed method in Sect. 3.2.

### 3.1 Deep Learning VGG16 Model

VGG16 is a convolutional neural network (CNN) architecture designed for image classification. It was introduced by (Simonyan and Zisserman, 2014)

which participated in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2014 and achieved notable success. It can more accurately express the characteristics of the data set when identifying and classifying images. This model has a 16 weight layers, which include 13 convolutional layers and 3 fully connected layers. It represented an advance on previous models, offering convolution layers with smaller convolution kernels (3×3) than had previously been the case.

During model training, the input for the first convolution layer is an RGB image of size 224 x 224. For all convolution layers, the convolution kernel is 3×3. These convolution layers are accompanied by Max-Pooling layers, each 2×2 in size, to reduce the size of the filters during learning.

In the proposed method, the model, as shown in 1, is inspired by (Moung et al., 2021). So, the last max-pooling layer was removed and replaced with an average pooling layer. this will introduce some feature generalization because this operation will take the average of pixels to retain fine-grain details in the final convolutions. During the classification phase, the output features from the average pooling layer of the VGG16 model are fed into the new fully-connected layers. The new classifier part of the VGG16 model consists of one flattened layer and two dense layers, with each generating 100 and 2 outputs, respectively. The first dense layer has a Rectified Linear Unit (ReLU) activation function with 0.5 dropout. The output layer, which is the last dense layer, has a Sigmoid activation function. The feature map size of the generated VGG16 features used in this work is $25{,}088 \times 1$ dimension per input image.

In this work, the proposed approach is built using a convolutional neural network (CNN) with batch normalization (BN) which is a popular and effective technique that consistently accelerates the convergence of deep networks.

## 3.2 Proposed Method

Our proposed method consists of 3 phases: feature extraction, feature concatenation and classification. Figure 2 illustrates the architecture of the proposed method. We will detail each phase in the next section.

### 3.2.1 Feature Extraction

In the feature extraction step, two types of features are extracted: those associated to DL (Deep Learning) and those relative to MI (moment invariant). For DL feature extraction, VGG16 is utilised in this

work. Meanwhile, the MI-based features are extracted using the (Ghorbel et al., 2006) moment invariant method. Then, these two features (DL-based and MI-based features) are concatenated using the join fusion. Lastly, classification based on the fused features is performed by using the fully connected (FC) layers.

### VGG16 Network Architecture

In this work, we use a CNN model VGG16 (see Section 3.1 to extract essential features. This model was composed of two parts: convolutional base and classifier. The convolutional base comprises convolutional and pooling layers to generate features, whereas the classifier categorises the image based on extracted features. We initialized the model weights using Kaiming Initialization for convolutions and normal initialization for fully connected layers. Additionally, we did not employ any transfer learning in our experiments. The VGG16 output features of shape 77512 at the final max-poling layer.

### Complete Complex Invariant Descriptors

Based on the work of (Ghorbel et al., 2006), the MI-based features of each CT image are extracted using complex invariant moments (see Section 2) to translation, rotation and scale, by means of some linear combinations of complex moments. We use both the absolute and phase values of the obtained complex invariant features, so no data are lost. More details are presented in the experimental results.

### 3.2.2 Feature Concatenation

In (Huang et al., 2020), authors described all the different fusion strategies using deep learning. In this work, we use the joint fusion as shown in Fig 3. It's the process of joining learned feature representations from intermediate layers of neural networks with features from other modalities as input to a final model. Thus, the features of each batches of images from the VGG16 model are concatenated with the Ghorbel moments features into a composite vector. The features extracted from VGG16 were multi-dimensional vectors of shape (7x7x512). Thus, it is flattened to be (25088x1) and the extracted Ghorbel complex moment invariant values were presented as a one-dimensional row vector. To be on the same scale as the deep learning feature, which is a vector in the $\mathbb{R}^{(25088x1)}$, we extracted the module and phase information from the Ghorbel complex invariant feature. Then we did the fusion. Fig3 illustrates the joint fusion process. The size of the final vector will vary
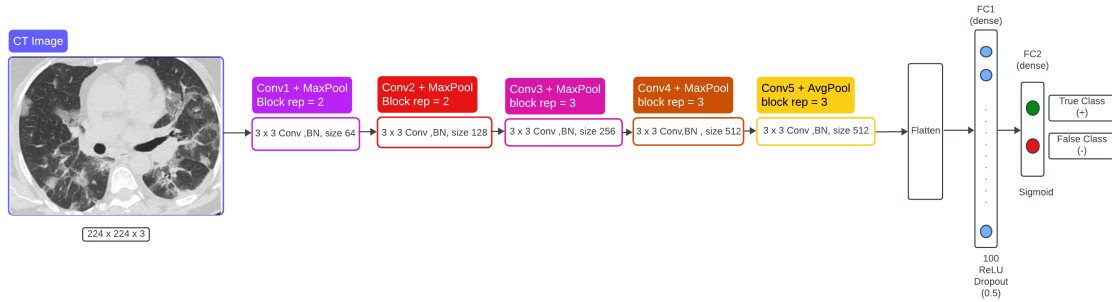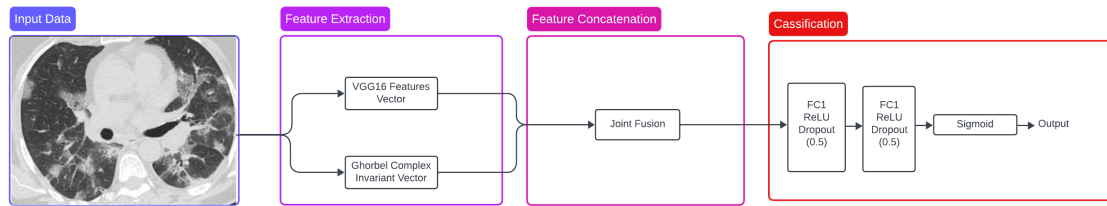
Figure 1: VGG16-based architecture with classifer block.



Figure 2: Overview of proposed method.

depending on the values chosen for $p$ and $q$ to calculate the MI features. This aspect will be discussed in the experimental section.

### 3.2.3 Classification Phase

For the classification phase, we followed the same method used in (Moung et al., 2021). We use a fully connected network which is built with two fully connected (FC) layers. On model VGG16, the first FC (dense) layers generates 100 output units with a ReLU activation function and a dropout of 50% and the second FC (dense) layers generates two output units with a sigmoid activation function.

## 4 EXPERIMENTAL RESULTS

In order to prove the robustness of our model, first, a set of experimental results where made in order to study invariant complex moments under rigid transformations, noise and non rigid ones performed on columbia coil database(Nene et al., 1996) which is available on its official website :https://www.col umbia.edu/CAVE/software/softlib/coil-20.php . Then, the proposed deep model is validated on SARS-COV-2 Ct-Scan Dataset.This dataset has been made public by the authors on 22 April 2020 and can be accessed via the Kaggle website at the following link: https://www.kaggle.com/datasets/plamenedua rdo/sarscov2-ctscan-dataset. The dataset contain 1252 CT scans that are positive for SARS-CoV-2 in-

fection (COVID-19) and 1230 CT scans for patients non-infected by SARS-CoV-2, 2482 CT scans in total. The selected data set was manually divided into training, validation and test sets with a ratio of 82.8 %, 9.1% and 8.1% respectively.

In this paper, we evaluate the performance of our predictor by calculating accuracy (ACC) for SARS-CoV-2 infection (COVID-19) .

$$\text{Accuracy} = \frac{\text{Number of Correct predictions}}{\text{Total no of predictions made}}$$

### 4.1 Stability of the Used Descriptors

In this section, we will study four families of invariant based moments which are Hu (Hu, 1962), Zernike (Khotanzad and Hong, 1990), Flusser (Flusser, 2002) and finally the Ghorbel et al. ones (Ghorbel et al., 2006). For 2D matrix moments (especially Flusser, Zernike and Ghorbel et al.), we used the zigzag technique described in (Ghorbel et al., 2006) to obtain 1D vector and keep the neighborhood. The comparison is performed according to the relative error as follows :

$$E_{I_1,I_2}(K) = |\frac{I_1(k) - I_2(k)}{I_1(k)}| \qquad (6)$$

We'll start by applying several rigid transformations to Lena's image to demonstrate the stability of these invariants. Figure 4 shows the image of Lena on which we'll be experimenting. For this first set of experiments, we fix $(p = q = 3)$ for Flusser and Ghorbel,(radios = digres = 4) for Zernike.
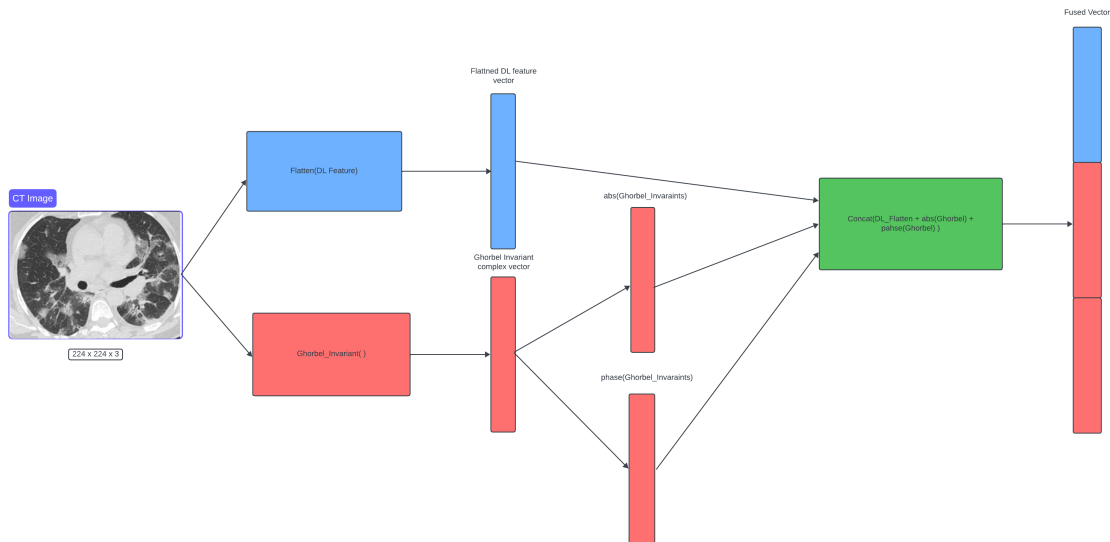
Figure 3: Feature Concatenation Process Between DL and Ghorbel Invariants Feature.



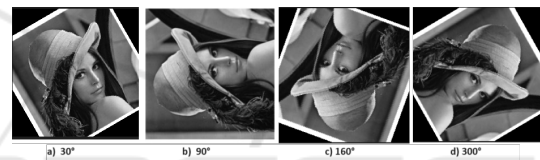Figure 4: Original image of Lena with size 224x224.



Figure 5: Several rotations of the original image.

## Invariance Against Rotation

To evaluate rotation invariance, we applied rotations to the image of Lena presented in Fig 4, rotating it by angles of 30°, 90°, 160° and 300°, as shown in Fig.5. An experimental demonstration of the stability between the invariants of (Flusser, 2006), (Hu, 1962), (Khotanzad and Hong, 1990) and (Ghorbel et al., 2006) is well illustrated in Fig.9 to Fig. 12. The results confirm that (Ghorbel et al., 2006) is invariant to rotation whatever the angle of rotation, compared with other families which have a higher relative error. We illustrate by Fig.14 and Fig.13 a summary of the evolution of relative error under two rotation angles (30° and 60°). It's clear that the descriptors we use have low, stable error values, whereas Flusser's descriptors give high error values. We also observe an oscillation in error values for the two families of Zernike and Hu moments.

## Invariance to Scale

To evaluate the robustness against scale factor, we have scaled the Lena image to scaling factors of 50%,

80%, 120% and 150%, as shown in Fig. 6. The plots in Fig.15 to Fig.18 show the relative errors between vector of invariants for the original image and those of the scaled ones. Relative error values for the Ghorbel et al. descriptors range from [0.1], while those for the Zernike descriptors range from [0.4] and those for the Hu'moments from [0.2]. For Flusser's descriptors, the error is huge, in fact Flusser'moments are only invariant under rotation. Also, we observe certain stability for the error only for Ghorbel et al. descriptors. This result confirms the robustness of (Ghorbel et al., 2006) against scale factor.



Figure 6: Several scaled images of the original image.

## Invariance Against Noise

To evaluate the robustness against noise, we introduced white Gaussian noise into the Lena image, in-

corporating a mean ($\mu$) of 0 and variances ($\sigma^2$) of 13, 20, 28, and 34, as shown in Fig.7. The plots in Figs.19 to 22 illustrate the relative errors between the invariant vector of the original image and the invariant vectors of the noisy images.
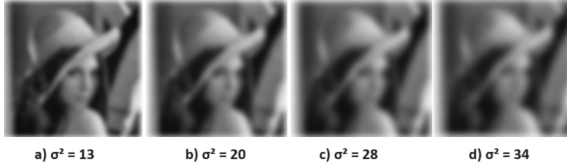


a) $\sigma^2 = 13$   b) $\sigma^2 = 20$   c) $\sigma^2 = 28$   d) $\sigma^2 = 34$

Figure 7: Several noisy images of Lena.

**Invariance to Non-Rigid Transformations**

In this section, we extend our experiment to non-rigid transformations with the images of ducks in Fig. 8 from Columbia Coil data base. The results in graphs 23 and 24 verified that complex invariant moments are stable to the non-rigid transformation using the (COIL-20) dataset.



a       b       c

Figure 8: Three test images used in the pattern experiments (original size 128×128). Fig. (a), (b) and (c) represent the same duck but with different orientation.

## 4.2 The Improved VGG16 Model

The performance of the VGG16 model with and without MI is summarized in table 1. In the first part of the experiments, the CT images are given only basic pre-processing steps before the training and test phases, which includes rescaling into 224 × 224 dimension using the resize function and we performed the normalization. All images that were fed into the VGG16 model for feature extraction retained their three-dimensional color channels format. However, the images that were fed into Ghorbel's invariant moment method (Ghorbel et al., 2006) were transformed from color channels to a single channel by taking the pixel-wise average across channels. We implemented the model using PyTorch.

We set a learning rate of 0.0004 and the number of epochs of 100 with an Adam optimiser based on the work (Moung et al., 2021). The accuracy that we obtained with this architecture is 97.512%. The concate-

Table 1: Performance of VGG16 model with and without MI on the testing dataset.

| Deep learning model | Accuracy (%) |
| --- | --- |
| VGG16 | 93.034 |
| VGG16 + 7 Hu'moments | 94.527 |
| VGG16 + Zernike descriptors | 93.532 |
| VGG16 + Flusser descriptors | Not available |
| **VGG16 + Ghorbel descriptors** | **97.512** |

nation of the invariant based complex moments in the VGG16 features improved the informativeness and discrimination capabilities of the flattened vector. The model's recognition ability improves with an increase in the number of features it learns. Furthermore, during training, the VGG16 model was exposed to complex moments, extent and solidity values in addition to the deep-learning feature vector. The absolute and phase values of the complex moments influenced the activations in the network. Feature(VGG16) + Abs(Ghorbel) + Phase(Ghorbel) showed a significant increase in accuracy over VGG16 (93.034%), VGG16 + Hu'moments (94.527%), VGG16 + Zernike descriptors (93.532%) and VGG16 + Flusser descriptors (Not available). In this case the Flusser invariant was so large that it resulted in numerical overflow and required more processing to be mapped to the VGG16 features, so we didn't train the model with Flusser.

## 5 CONCLUSION

In this research, we proposed an improved VGG16 deep learning model based on the fusion of CNN features and a complete and invariant set of descriptors based on the image's complex moments. Several experimental results were performed in order to demonstrate the stability of the proposed set of shape descriptors. Compared to other works based on Hu moments, Zernike and Flusser invariants, the proposed method achieves the highest accuracy. As future perspectives, we plan to study our approach on other databases with more sophisticated CNN architectures such as VGG19.

Figure 9: Stability study of Flusser invariants under rotation.



Figure 10: Stability study of Hu invariant under rotation.



Figure 11: Stability study of Zernike invariants under rotation.



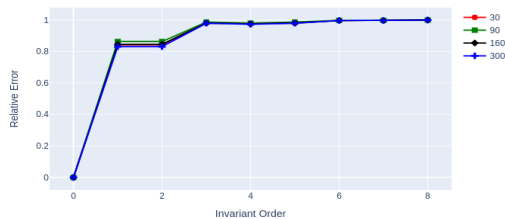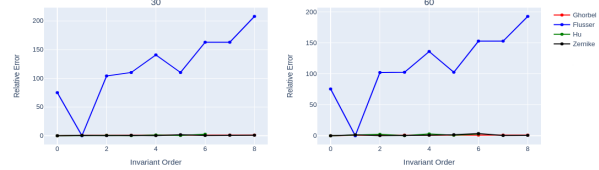Figure 12: Stability study of Ghorbel invariant under rotation.



Figure 13: Relative error against rotation, for the Hu's, Flusser's, Zernike's and Ghorbel's set of similarity invariants.



Figure 14: Relative error against rotation, for the Hu's, Zernike's and Ghorbel's set of similarity invariants (logarithmic scale except for the image).



Figure 15: Stability study of Flusser invariant moments under scale factor.



Figure 16: Stability study of Hu invariant moments under scale factor.

Figure 17: Stability study of Zernike invariant moments under scale factor.
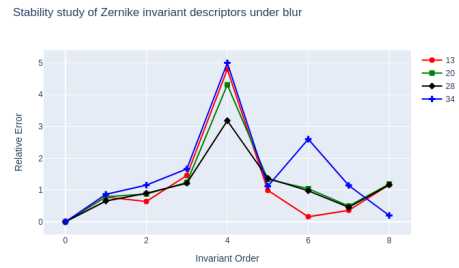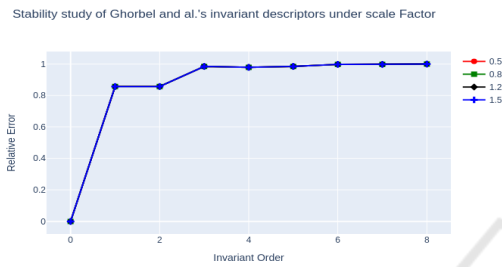


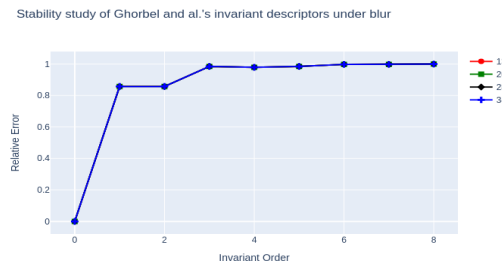Figure 18: Stability study of Ghorbel invariant moments under scale factor.



Figure 19: Stability study of Flusser invariant moments under blur.



Figure 20: Stability study of Hu invariant moments under blur.



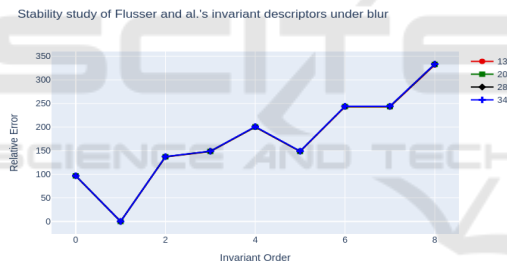Figure 21: Stability study of Zernike invariant moments under blur.



Figure 22: Stability study of Ghorbel invariant moments under blur.
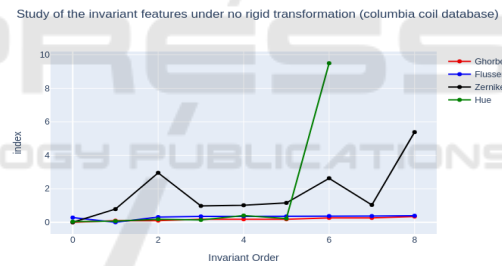


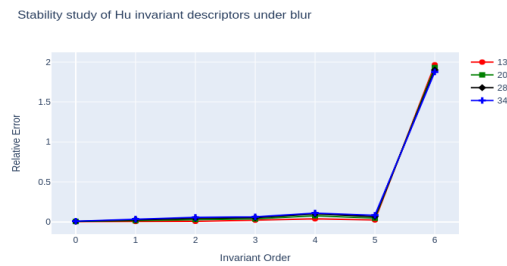Figure 23: Study of the invariant features under non-rigid transformation (duck image) between (a) and (b).
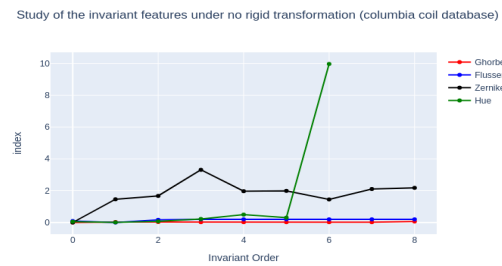


Figure 24: Study of the invariant features under non-rigid transformation (duck image) between (a) and (c).

# REFERENCES

Cao, X., Gao, S., Chen, L., and Wang, Y. (2020). Ship recognition method combined with image segmentation and deep learning feature extraction in video surveillance. *Multimedia Tools and Applications*, 79:9177–9192.

Derrode, S. and Ghorbel, F. (2001). Robust and efficient fourier-mellin transform approximations for invariant grey-level image description and reconstruction. *Computer Vision and Image Understanding*, 83(1):57–78.

Flusser, J. (2002). On the inverse problem of rotation moment invariants. *Pattern Recognition*, 35(12):3015–3017.

Flusser, J. (2006). Moment invariants in image analysis. In *proceedings of world academy of science, engineering and technology*, volume 11, pages 196–201. Citeseer.

Ghorbel, F. (1994). A complete invariant description for gray-level images by the harmonic analysis approach. *Pattern Recognit. Lett.*, 15(10):1043–1051.

Ghorbel, F. (1998). Towards a unitary formulation for invariant image description: application to image coding. *Annals of telecommunications*, 53(5):242–260.

Ghorbel, F., Derrode, S., Mezhoud, R., Bannour, T., and Dhahbi, S. (2006). Image reconstruction from a complete set of similarity invariants extracted from complex moments. *Pattern recognition letters*, 27(12):1361–1369.

Hu, M.-K. (1962). Visual pattern recognition by moment invariants. *IRE transactions on information theory*, 8(2):179–187.

Huang, S.-C., Pareek, A., Seyyedi, S., Banerjee, I., and Lungren, M. P. (2020). Fusion of medical imaging and electronic health records using deep learning: a systematic review and implementation guidelines. *NPJ digital medicine*, 3(1):136.

Khotanzad, A. and Hong, Y. H. (1990). Invariant image recognition by zernike moments. *IEEE Transactions on pattern analysis and machine intelligence*, 12(5):489–497.

Moung, E. G., Hou, C. J., Sufian, M. M., Hijazi, M. H. A., Dargham, J. A., and Omatu, S. (2021). Fusion of moment invariant method and deep learning algorithm for covid-19 classification. *Big Data and Cognitive Computing*, 5(4):74.

Nene, S. A., Nayar, S. K., Murase, H., et al. (1996). Columbia object image library (coil-20).

Persoon, E. and Fu, K.-S. (1977). Shape discrimination using fourier descriptors. *IEEE Transactions on systems, man, and cybernetics*, 7(3):170–179.

Ren, Y., Yang, J., Zhang, Q., and Guo, Z. (2021). Ship recognition based on hu invariant moments and convolutional neural network for video surveillance. *Multimedia Tools and Applications*, 80:1343–1373.

Sharma, S. and Mehra, R. (2020). Conventional machine learning and deep learning approach for multi-classification of breast cancer histopathology images—a comparative insight. *Journal of digital imaging*, 33:632–654.

Sheng, Y. and Shen, L. (1994). Orthogonal fourier–mellin moments for invariant pattern recognition. *JOSA A*, 11(6):1748–1757.

Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.

Suk, T. and Flusser, J. (2003). Combined blur and affine moment invariants and their use in pattern recognition. *Pattern Recognition*, 36(12):2895–2907.

Wang, Y.-B., You, Z.-H., Yang, S., Yi, H.-C., Chen, Z.-H., and Zheng, K. (2020). A deep learning-based method for drug-target interaction prediction based on long short-term memory neural network. *BMC medical informatics and decision making*, 20(2):1–9.

Zahn, C. T. and Roskies, R. Z. (1972). Fourier descriptors for plane closed curves. *IEEE Transactions on computers*, 100(3):269–281.