

Veriforge: An Image Forgery Detection Model

RISHIKESH RAVI¹, HAMZA RANGWALA², AYZAF AFZAL³ and TRUSHIT PATEL⁴,

¹Wilfrid Laurier University, Waterloo ON Canada (email: rish8180@mylaurier.ca)

²Wilfrid Laurier University, Waterloo ON Canada (email: rang8720@mylaurier.ca)

³Wilfrid Laurier University, Waterloo ON Canada (email: afza9880@mylaurier.ca)

⁴Wilfrid Laurier University, Waterloo ON Canada (email: pate9410@mylaurier.ca)

ABSTRACT This paper proposes a new method of tampered image detection by combining the CNN with a backbone of ResNet50 with ELA. Our method is designed for digital forensics, which is a very important domain for the detection and assessment of image modifications in many other domains, such as cybersecurity, journalism, and criminal investigations. This procedure comprises image pre-processing to ascertain altered areas by comparing variations in compression levels to find compression artifacts using ELA. It serves as the foundation for reliable feature extraction: the ResNet50 model, pre-trained on ImageNet, has its basic layers frozen to preserve learned weights. For customization with regard to binary classification, some custom layers are appended on top: global average pooling, batch normalization, dense layers with ReLU, and dropout regularization. The architecture here comes up with a generator-based approach so that large datasets can be handled with better memory consumption, enabling real-time data preprocessing during training of models. Our methodology is assured to optimize resource utilization and enhance generalization, which has been evidenced with extensive experiments. The obtained model offers high accuracy and robustness in detecting tampered images; thus, it presents a scalable solution for real-world forensic applications. This integration of ELA and CNN provides the backbone of fast and accurate analysis of digital evidence to help in better decision-making in sensitive situations.

INDEX TERMS Tampered Image Detection, Error Level Analysis (ELA), Convolutional Neural Network (CNN), ResNet50, Digital Forensics, Image Manipulation Detection, Cybersecurity, Journalism, Criminal Investigations, Feature Extraction, Binary Classification, Data Preprocessing, Dropout Regularization, Compression Artifacts, Real-time Analysis, Scalable Forensic Applications.

I. INTRODUCTION

IN this digital era, where information is usually imparted through the means of digitization, a tampered image might start a wildfire of misinformation, controversies, and disasters in a lot of sectors. The more use of social media escalates the scale of it, which calls for this challenge to be taken seriously. The instances of tampered images being used for spreading fake news, malicing individuals or organizations, and manipulating public opinion have started becoming alarmingly common. A serious issue of this aspect, our research focuses on VeriForge, which is a model designed for the detection and exposure of image tampering using the powerful deep learning and Error Level Analysis [1]. This innovative approach is targeted at offering an effective solution in finding discrepancies in digital images and building trust in visual content. But accessibility to sophisticated image editing software has put the art of doctoring images within everybody's reach, irrespective of any kind of technical expertise. For that reason, it has become an increasingly demanding

task to differentiate between original and tampered images by commoners, but also by professionals in vital fields like law and order or journalism. But available techniques of image tampering detection, though at one's fingertips, are not short-circuited by numerous limitations. Most of the known methods are outdated, slowing down, and cannot find a complex type of manipulation. Nowadays, there is a serious demand to create modern and effective methods which could cope with complications arising in image forgery techniques [2].

VeriForge fills this gap by using the synergy of traditional forensic methods and deep learning. ELA, for example, is a well-established method in forensics: it underlines inconsistencies between image compression levels, hence highlighting the image anomalies caused by editing [3]. VeriForge couples ELA with deep learning to extract meaningful features from these discrepancies in a way that has allowed the attainment of high accuracy enabled by advanced neural network architecture. It allows VeriForge to find the tampered images with quite an intelligent attack.

For the training and verification of VeriForge, we utilize CASIA V2, a database of well-labeled images in their original and forged aspects. In this database, tampering is performed manually with several image editing tools; thus, it is one of the good sources for training a model to detect different types of tampering [4]. This is also ensured by the great variety of images in the dataset, thus the model generalizes well to unseen data, avoiding overfitting and improving its application in real-world scenarios.

The VGG-16 architecture was adopted for the first version of VeriForge. Although a promising result could be achieved in a preliminary way, we obtained only mediocre performance, given that after some point the model converged, resulting in an inability to increase its accuracy. Then, in order to capture the intricate image features, residual learning with the much deeper convolutional neural network architecture, ResNet50 [5], was chosen. This switch to ResNet50 allowed VeriForge to capture more complicated patterns of tampered images, thereby increasing its detection capability by a big margin.

The general objective of the work is to protect the integrity of digital content by offering a reliable, scalable, and efficient tool for image forgery detection. In this regard, VeriForge uses advanced techniques such as convolutional neural networks (CNNs) [6], transfer learning [7], and state-of-the-art architectures like ResNet50 [8]. Put together, these technologies ensure that VeriForge will find tampered images and continually improve by keeping pace with the changing world of image manipulation.

In a nutshell, VeriForge is a leap into the future of image tampering detection: it merges classical forensic techniques such as ELA with the latest powers of deep learning for an emerging need to have reliable tools safeguarding digital content in this world where visual communication is well on its way to dominating other forms.

II. PROPOSED METHOD

In this paper, we demonstrate that combining ELA with CNN model while freezing the base layers and adding custom layers to the base model which will provide a significant increase in the robustness of the final model and significantly improve performance of the classification task to detect images, which are tampered. We explore this idea in the context of digital forensics, where the goals are to analyze the images and detect suspicious activity in the field of crime, journalism etc. Since digital forensics helps uncover activities and patterns, determine the root causes of incidents, and establish a chain of evidence admissible in court, it becomes essential to identify, recover, analyze and present digital evidence from electronic devices and digital storage. Additionally, an ELA based model helps in identifying the difference in various layers of the images to help us uncover the underlying changes made to the image that is suspected to be tampered with. The model allows us to help various fields where criminal investigations, legal disputes and cybersecurity incidents occur and images need to be analyzed quickly to help user make quick informed

decisions. The following sections outline the details of each step in the process.

A. INITIALIZATION

The initialization of the framework is preprocessing the images based on the number of images available in the dataset, it corresponds to the total number of images which are categorized in two classes specifically: 1) Authentic and 2) Tampered. The classes with their labels are split initially to ensure we have a sizeable amount of image to work with, we ensure that the classes are not biased and are equally balanced to provide a robust model and reduce overfitting. The classes are initialized with ones and zeros, zero if the image is authentic and one if the image has been tampered. For images to get preprocessed, we need to handle a chunk of images at a time, which will require us to take images in batches to preprocess and analyze further.

After reshaping the images according to our input shape which will reduce the size of the image so that our next function would work effectively. Creating batches to handle memory overflow and preprocessing image before feeding them to the ELA function ensures that the system doesn't require a lot of resources and is more optimized to create a more robust model.

Further, we take each resized image and input it to our ELA function, this step is iterated over multiple images to get robust images with reduced quality to take forward for our ELA function.

B. ELA PROCESS WITH GENERATORS

For each image in our ELA function, we generate a buffer to save the original image to a temporary in-memory file with reduced quality, which helps us in reducing the overall memory utilization of the system. We load the compressed image from the buffer and calculate the difference between the original image and the one which was compressed to get the Pillow image object, this is a technique used in image forensics to detect tampering or inconsistencies in an image. We enhance the brightness of the image according to the extrema we get in the pillow image object.

This returns us the ELA image which is a new image that emphasizes regions with differing compression levels, useful for detecting tampering or manipulation. This image will highlight differences due to compression artifacts. To feed the image forward to the model for training, it requires the image to be converted to an array for numerical processing and is normalized making the data compatible with many machine learning models. We close the image further to free up resources and it is a good practice to avoid file locks or memory leaks, especially when processing multiple images in a loop. For the model to handle such large number of images at once would require a lot of resources, to tackle this issue we provide generators to the model.

These generators help in handle image preprocessing and feeding batches of data to the model during training. They are often used in deep learning workflows when the model

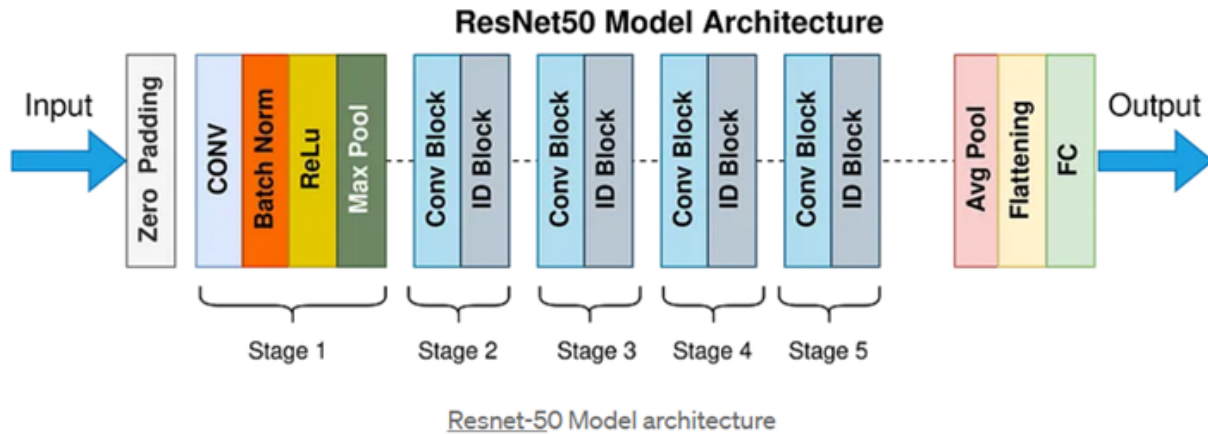


FIGURE 1. Architecture of the base model

requires to load images in batches and apply transformations or augmentation. Instead of loading all images into memory at once, generators load and process batches on-the-fly, reducing memory consumption. Generators often include on-the-fly preprocessing, which enhances the dataset and helps the model generalize better, increasing its robustness and generalizability for better performance. We also try to visualize the batch of images after the ELA function to ensure that the shape of the generator and preprocessed images are according to the requirements.

C. RESNET50 MODEL ARCHITECTURE WITH TRAINING EVOLUTION

In this study, we employ a Resnet50 model to classify the images into two categories (Authentic and Tampered), we use a transfer learning approach based on the ResNet50 architecture, pre-trained on ImageNet to solve our classification problem. This method capitalizes on the robust feature extraction capabilities of ResNet50, significantly reducing training time and computational resources while improving model accuracy. The model was customized to adapt the pre-trained base for the task by introducing additional trainable (Custom) layers specifically designed for binary classification. Figure 1: shows the architecture of the base model we have used featuring Resnet50 in a five-stage design to represent the overall function of the layers

The base of our architecture is the ResNet50 model, known for its depth and efficient handling of the vanishing gradient problem through residual connections. The model was loaded without the top layers, removing the full-connected top layers helped us retained only the convolutional layers which are responsible for extracting hierarchical features from the input images. The input shape was fixed to match the requirement of the model while ensuring compatibility with pre-trained

weights. The freezing of all layers in the base model during training aimed at preserving the learned weights and keeping the feature extraction capability unaltered by the new task-specific data. To this binary classification problem, we attached five different types of custom layers to the base pre-trained model in an attempt to make it adapt. First, there is a global average pooling layer that reduces the high-dimensional feature maps output by ResNet50 into a single vector for each feature map, reducing the number of trainable parameters and hence improving generalization. Unlike the traditional Flatten layer, this preserves spatial information and improves model performance when using ResNet-based architectures. Further, we added batch normalization to normalize the activations from the pooling layer, reducing internal covariate shift and accelerating convergence.

Once we are done with that, we added a series of two dense layers. The first layer consists of 512 neurons providing substantial capacity to capture complex patterns. The second contains 256 neurons for further refinement of feature extraction. Both are activated by the ReLU activation function, which introduces non-linearity in deep networks and provides output zero for negative inputs, creating sparsity in the activations that could further be helpful in efficient computations with a reduction in overfitting. While being computationally less expensive than sigmoid or tanh function. The custom layers then involved a dropout regularization technique which was employed after each dense layer to with a rate of 0.75 after the first dense layer to impose strong regularization and rate of 0.5 was applied after the second dense layer for moderate regularization. This helped in preventing overfitting of the model during training. After dropout regularization, there was a single neuron in the output layer with sigmoid activation to predict the probability of the positive class, which is suitable for binary classification tasks.

It achieves a very good trade-off among accuracy, training efficiency, and robustness, making it pretty effective in real-world applications. This model was trained for more than 100 epochs with a step size of about 145 while capturing essential evaluation metrics such as accuracy, the f1 score, and validation accuracy. This will make sure we minimize the gap between the training and validation accuracy of the model, which will be robust and very effective in real-world applications.

III. RESULTS AND ANALYSIS

A. DATASETS

The CASIA V2 dataset, created by the Institute of Automation at the Chinese Academy of Sciences, CASIA, is one of the more general resources for research in the area of image forgery detection. Designed to support studies on image tampering, it includes manipulative techniques such as splicing, copy-move forgery, and other methods. The Table 1 outlines the details of the dataset and gives us an overview of what the dataset contains. The dataset is a very diverse collection of 12,614 images of roughly equal amounts of authentic and tampered examples from various categories such as nature, animals, objects, and urban settings. Such diversity will enable models trained with CASIA V2 to generalize rather well across different scenarios. There are also ground truth masks provided for every tampered image in the dataset, highlighting manipulated regions and thus making the dataset indispensable for supervised learning tasks and performance evaluation. The resolution of the images also varies to reflect real-world conditions, as images naturally come from different devices and sources. CASIA V2 has a variety of applications, including training and validation of image forgery detection models, testing of pre-processing techniques such as Error Level Analysis, and benchmarking deep learning architectures like CNNs for forgery classification. Features such as these make it invaluable in the field of digital forensics and research on image authenticity.

Feature	Description
Dataset Name	CASIA TIDE v2.0
Purpose	Designed for research and development in image forgery detection
Released By	Chinese Academy of Sciences Institute of Automation (CASIA)
Image Type	Digital images, both authentic and tampered
Number of Images	12,614 images (7,491 authentic and 5,123 tampered)
Forgery Techniques	Splicing, copy-move, and other common image tampering techniques
Resolution	Varies (typically medium resolution)
File Format	JPEG (for compressed images)
Dataset Structure	Contains two folders: one for authentic images and one for tampered images
Applications	Used for image forgery detection, digital forensics, and tampered image localization research
Public Availability	Available for academic and non-commercial use upon request from CASIA

TABLE 1. Overview of the CASIA TIDE v2.0 Dataset

B. RESULTS EVALUATION

This section compiles the quantitative results and presents an analysis of the project outcomes. The performance of ResNet50 and VGG16 CNN models trained on CASIA V2 dataset by using traditional methods like ELA will be judged on Accuracy and F1 metrics. Accuracy measures the overall correctness of the model's predictions. At the same time, the F1 score represents the balance between precision and recall, which is suited for imbalanced datasets. Or in other words,

it shows the likelihood that the model correctly classifies an image with its true label.

Model Variant	Accuracy	F1 Score
VGG16 + ELA	83%	0.86
ResNet50 + Improved ELA	89%	0.92

TABLE 2. Quantitative Results

According to the qualitative results, RestNet50, combined with improved ELA, performed better than VGG16, which used traditional ELA.

C. ANALYSIS OF THE INITIAL VGG16 MODEL

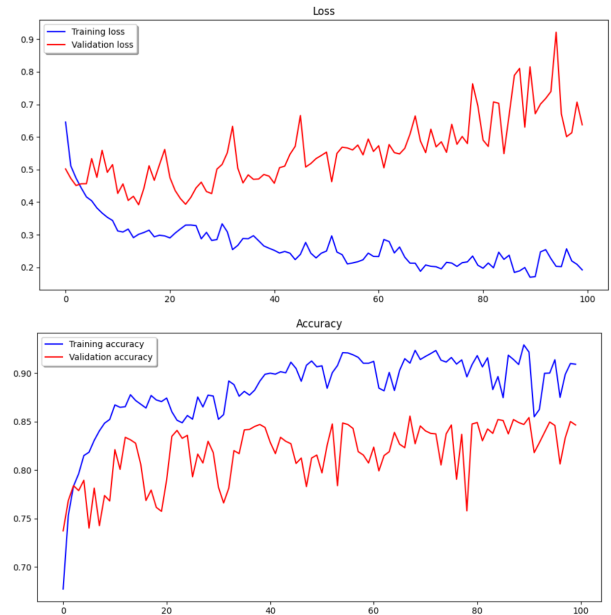


FIGURE 2. VGG16 Model: Loss and Accuracy Graphs

The first prototyping at the earlier stages of the project used a VGG16 CNN model combined with traditional ELA preprocessing. This approach focused on finding edges to detect forged regions on a grey-scale-converted image. Although this yielded a promising 83% in validation accuracy, the training and validation graphs were too inconsistent and abrupt, making the F1 score stand at a low of 0.86. That meant the balance between precision and recall had to be promoted, reducing the higher number of false positives and false negatives, which could be achieved by using more sophisticated preprocessing techniques. The model VGG16 has anomalies shown in Figure due to sudden changes and also huge differences between Training and Validation accuracy graphs. Because the model failed to generalise into a valid test set, it could be defined that the model went to overfit while in training, also justified with the low F1 value of the VGG16 model

Figure 2 reveals inconsistencies, with abrupt variations and a wide gap between the Training and Validation accuracy

graphs for the VGG16 model. As the model failed to generalise on the validation test set, it can be termed that the model went overfitting during the training, which is evident by the low F1 score of the VGG16 model.

D. RESNET50 MODEL AND IMPROVED ELA WITH EVALUATION METRICS

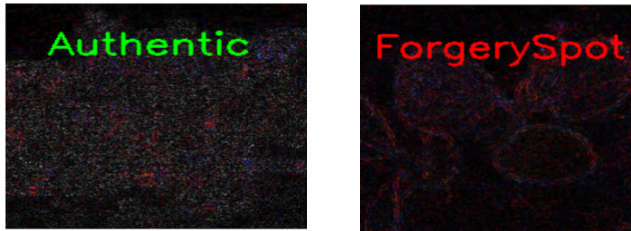


FIGURE 3. Authentic and Forgery Predicted Label images

Resnet50 with an enhanced ELA preprocessing was implemented to address these limitations. The Improved preprocessing technique modified the ELA to analyse the RGB channels instead of grayscale. This allowed the model to detect colour discrepancies and variations in JPEG compression artefacts which are more effective in detecting forgery compared to finding edges. The improved ELA led to the clear identification of tampered regions, as shown in Fig 3.

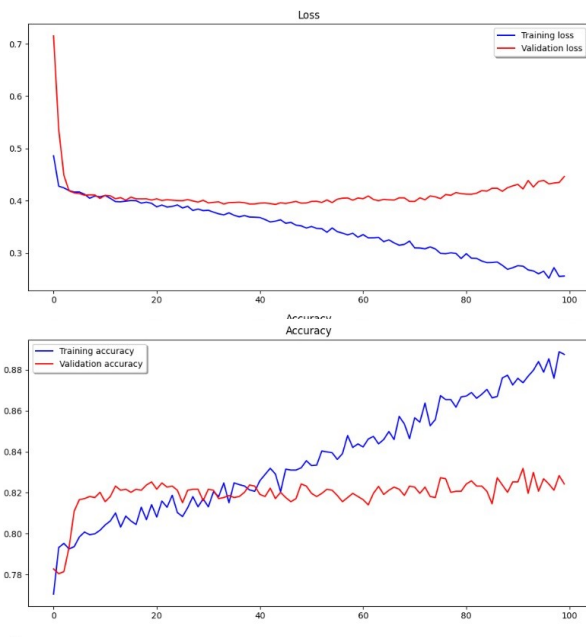


FIGURE 4. Resnet50 Model: Loss and Accuracy Graphs

Besides that, compared to VGG16, several advantages have been gained with the ResNet50 architecture. Instead of flattening the feature maps, global average pooling was used that reduces overfitting because of minimized trainable parameters. More than that, batch normalisation layers stabilize training by normalizing the inputs, and dropout with

various intensities added to the regularization helped reduce overfitting. With these changes, training-validation accuracy curves were smoother and more convergent, as depicted in Fig 4.

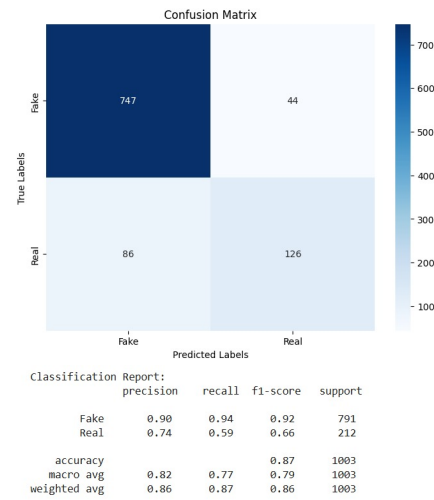


FIGURE 5. Confusion Matrix and Evaluation Chart

The improved ELA coupled with transitioning from VGG16 to Resnet50 significantly reduced the number of false positives. Smoothing of both the validation and training accuracy curves showed reduced overfitting compared to the curves of the VGG16 model. The resultant 89% accuracy and 0.92 weighted F1 score reflect the better generalization capability of the model compared to VGG16. A higher F1 score indicated a balanced recall and precision where false positives and false negatives are minimised as evident in Fig 5.

It is noteworthy to mention that data augmentation techniques – rotation, flip, zoom, width shift height shift etc – when applied did not have any effect on forged images. The best possible explanation for this phenomenon is that data augmentation does not present any variation in the forged region. After finalising the model configuration, training, and validation, it was programmatically uploaded to Hugging Face, from where it would be remotely instantiated in the backend service layer using the API provided by Hugging Face. Python Django was used to create a backend application to serve user requests, process the images present in the user requests, predict the labels, and respond with the predicted results.

IV. CONCLUSION

This study is essential in maintaining integrity is crucial in fields like Law enforcement, Cyber security, and Legal investigations. In this project, we aimed to design a robust system capable of determining whether an image is tampered with. The research compares the performance of two widely known CNN models, VGG16 and ResNet50, with traditional preprocessing methods like ELA and data augmentation. The project concludes with viable results that taking colour in-

consistencies and compression differences significantly improves model capability compared to conventional grayscale-based ELA methods. The system can be integrated into social media to moderate content in real-time, flagging misleading images and advancing platform integrity. In journalism, it makes sure that media is authentic by verifying images before publication; this helps in debunking fake news and fosters ethical reporting as a way of maintaining the public's trust in digital journalism. Although promising performance may be seen in the CASIA v2 dataset, it does not generalize upon exposure to other datasets or real-world scenarios. This needs to be handled in future work by enrichment of the dataset with diversity and updates for further strengthening the resistance and longevity of the model. Moreover, interpreting it with Explainable AI by using Grad-CAM will illustrate the regions of interest influential in the model's choices for transparent decisions and gaining the trust of the users. The dataset and code for further research are available [here] and [repository link].

ACKNOWLEDGMENT

We would like to express our sincere gratitude to Prof. ANK Zaman for his invaluable guidance and support throughout the course of this research. His insightful suggestions, expert advice, and continuous encouragement significantly contributed to the successful completion of this project. We also appreciate his assistance in refining our methodology and providing

constructive feedback, which greatly enhanced the quality of our work. This research would not have been possible without his mentorship.

REFERENCES

- [1] A. Gupta, R. Joshi, and R. Laban, "Detection of tool based edited images from error level analysis and convolutional neural network," *arXiv preprint arXiv:2204.09075*, 2022.
- [2] A. Kaur, D. Chahal, and L. Kharb, "Weak form efficiency of currency futures: Evidence from india," in *2019 International Conference on Computing, Communication, and Intelligent Systems (ICCCIS)*. IEEE, 2019, pp. 30–35.
- [3] A. Singh and J. Singh, "Image forgery detection using deep neural network," in *2021 8th International conference on signal processing and integrated networks (SPIN)*. IEEE, 2021, pp. 504–509.
- [4] N. B. A. Warif, M. Y. I. Idris, A. W. A. Wahab, and R. Salleh, "An evaluation of error level analysis in image forensics," in *2015 5th IEEE International Conference on System Engineering and Technology (ICSET)*, 2015, pp. 23–28.
- [5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2015. [Online]. Available: <https://arxiv.org/abs/1409.1556>
- [6] S. D. Thepade, S. Bhandari, C. Bagde, R. Chaware, and K. Lodha, "Image forgery detection using machine learning with fusion of global and local thepade's sbtc features," in *2021 International Conference on Disruptive Technologies for Multi-Disciplinary Research and Applications (CENT-CON)*, vol. 1. IEEE, 2021, pp. 234–238.
- [7] S. Han, W. Gao, Y. Wan, and Y. Wu, "Scene-unified image translation for visual localization," in *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2020, pp. 2266–2270.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.



FIRST A. AUTHOR received the B.S. and M.S. degrees in aerospace engineering from the University of Virginia, Charlottesville, in 2001 and the Ph.D. degree in mechanical engineering from Drexel University, Philadelphia, PA, in 2008.

From 2001 to 2004, he was a Research Assistant with the Princeton Plasma Physics Laboratory. Since 2009, he has been an Assistant Professor with the Mechanical Engineering Department, Texas A&M University, College Station. He is the

author of three books, more than 150 articles, and more than 70 inventions. His research interests include high-pressure and high-density nonthermal plasma discharge processes and applications, microscale plasma discharges, discharges in liquids, spectroscopic diagnostics, plasma propulsion, and innovation plasma applications. He is an Associate Editor of the journal *Earth, Moon, Planets*, and holds two patents.

Dr. Author was a recipient of the International Association of Geomagnetism and Aeronomy Young Scientist Award for Excellence in 2008, and the IEEE Electromagnetic Compatibility Society Best Symposium Paper Award in 2011.



THIRD C. AUTHOR, JR. (M'87) received the B.S. degree in mechanical engineering from National Chung Cheng University, Chiayi, Taiwan, in 2004 and the M.S. degree in mechanical engineering from National Tsing Hua University, Hsinchu, Taiwan, in 2006. He is currently pursuing the Ph.D. degree in mechanical engineering at Texas A&M University, College Station, TX, USA.

From 2008 to 2009, he was a Research Assistant with the Institute of Physics, Academia Sinica, Taipei, Taiwan. His research interest includes the development of surface processing and biological/medical treatment techniques using nonthermal atmospheric pressure plasmas, fundamental study of plasma sources, and fabrication of micro- or nanostructured surfaces.

Mr. Author's awards and honors include the Frew Fellowship (Australian Academy of Science), the I. I. Rabi Prize (APS), the European Frequency and Time Forum Award, the Carl Zeiss Research Award, the William F. Meggers Award and the Adolph Lomb Medal (OSA).

...



SECOND B. AUTHOR (M'76–SM'81–F'87) and all authors may include biographies. Biographies are often not included in conference-related papers. This author became a Member (M) of IEEE in 1976, a Senior Member (SM) in 1981, and a Fellow (F) in 1987. The first paragraph may contain a place and/or date of birth (list place, then date). Next, the author's educational background is listed. The degrees should be listed with type of degree in what field, which institution, city, state, and country, and year the degree was earned. The author's major field of study should be lower-cased.

The second paragraph uses the pronoun of the person (he or she) and not the author's last name. It lists military and work experience, including summer and fellowship jobs. Job titles are capitalized. The current job must have a location; previous positions may be listed without one. Information concerning previous publications may be included. Try not to list more than three books or published articles. The format for listing publishers of a book within the biography is: title of book (publisher name, year) similar to a reference. Current and previous research interests end the paragraph.

The third paragraph begins with the author's title and last name (e.g., Dr. Smith, Prof. Jones, Mr. Kajor, Ms. Hunter). List any memberships in professional societies other than the IEEE. Finally, list any awards and work for IEEE committees and publications. If a photograph is provided, it should be of good quality, and professional-looking. Following are two examples of an author's biography.