

Enhancing Gastric Cancer Diagnosis through Deep Learning and Ensemble Techniques: A Comprehensive Analysis of Sub-database Performance

Angati Kalyan Kumar¹, Gangadhara Rao Kancharla²

¹Research Scholar, Department of Computer Science & Engineering
University College of Sciences, Acharya Nagarjuna University
Nagarjuna Nagar, Guntur, Andhra Pradesh, India.
e-mail:kalyank442@gmail.com

²Professor, Department of Computer Science & Engineering
University College of Sciences, Acharya Nagarjuna University
Nagarjuna Nagar, Guntur, Andhra Pradesh, India.
e-mail:kancharla123@gmail.com

Abstract—Gastric cancer is the fifth most prevalent cancer worldwide and the fourth most fatal. Early detection is critical for effective treatment. Histopathological examination is the established diagnostic method for gastric cancer. Recent advancements in computer technology have accelerated the use of digital tools to aid pathologists in diagnosing gastric cancer from pathological images. Ensemble learning is employed to enhance algorithm precision, involving the integration of multiple complementary learning models. The experimental platform focused on three subdatabases within GasHisSDB. Four deep learning classifiers, specifically VGG19, Inception-V4, ResNeXt, and ResNet152, were employed for classification experiments on the GasHisSDB database. The evaluation encompassed performance metrics, including accuracy, precision, recall, specificity, and F1 score. In the examination of the 160 x 160 pixel sub-database, ResNet152 stood out by delivering exceptional results in both categories. It achieved a remarkable accuracy of 98.02% in the "Normal" category and 87.9% in the "Abnormal" category. In the 120 x 120 pixel sub-database, ResNeXt displayed strong performance with a 96.98% accuracy in the "Normal" category, while its accuracy dropped to 89.09% in the "Abnormal" category. Notably, in the 80 x 80 pixel sub-database, ResNet152 emerged as the top performer with a remarkable 98.67% accuracy in the "Normal" category and 95.12% in the "Abnormal" category. Across these diverse sub-databases, ResNet152 consistently outperformed other models, maintaining high accuracy and precision while ensuring balanced performance in both categories.

Keywords—Image Classification; Cancer; Deep Learning; Performance Metrics.

I. INTRODUCTION

Gastric cancer is a prevalent and life-threatening malignancy that poses a significant global health challenge. Recent global cancer statistics indicate that gastric cancer ranks as the fifth most frequently diagnosed cancer and is responsible for a substantial portion of cancer-related deaths, accounting for 18.0% of the total [1]. Pathologists traditionally identify diseased regions by visually examining pathological slides with the naked eye, followed by further scrutiny using low-power microscopes [2]. Although the prevailing approach involves analyzing whole-slide images, practical constraints, such as limited computer performance, often necessitate subdividing these images into smaller components for analysis. To address this need, a sub-size image database is essential to assess the effectiveness of various medical image classification techniques in the face of real-world challenges.

The rapid advancements in computer vision, particularly in the realm of medical image classification, offer the potential to efficiently and swiftly scrutinize every microscopic image [3]. This technological progress presents an opportunity to address the complexities involved in gastric cancer diagnosis. Image classification techniques have ushered in a new era of progress, enabling the differentiation between benign and malignant cancers, the identification of various tumor differentiation stages, and the categorization of distinct tumor subtypes. These techniques provide valuable support to pathologists during the diagnostic process. Furthermore, the evolving landscape of image classification technology [4] is primarily focused on enhancing the precision of classification algorithms and fortifying their resilience against interference. In this context, ensemble learning has emerged as an effective solution. It is paramount to identify multiple efficient classification

algorithms that exhibit complementary properties to maximize classification accuracy.

This paper introduces the Gastric Histopathology Sub-size Image Database (GasHisSDB) [5], which is publicly available and a valuable resource. This database comprises a substantial collection of 245,196 sub-size pathological images related to gastric cancer and features three distinct sub-size labels. Each image within this database is associated with the calculation of three distinct features. The evaluation results encompass various classification methodologies based on both image features and raw images. Deep learning approaches, including VGG19, Inception V4, ResNeXt, and ResNet152, are employed to showcase the discriminative capabilities of each classifier.

In summary, gastric cancer is a formidable global health challenge. The reliance on traditional diagnostic methods, such as visual inspection and microscopic analysis, is being complemented by advanced computer vision technology. This technology enables the rapid and accurate analysis of pathological images, offering crucial support to pathologists. The introduction of the gastric histopathology sub-size image database is a significant milestone, providing researchers with a comprehensive dataset for evaluating and advancing medical image classification techniques. Additionally, the deployment of state-of-the-art deep learning models demonstrates the potential for these technologies to enhance the precision and efficiency of gastric cancer diagnosis.

II. LITERATURE REVIEW

A deep learning framework designed for the analysis of histopathological images is presented in [6]. Deep-Hipo leverages multi-scale receptive fields to enhance the accuracy and effectiveness of analyzing such complex images. This research contributes to the growing field of medical image analysis and underscores the importance of deep learning techniques for improving the diagnosis and understanding of histopathological images. The results of a Delphi survey regarding histopathologic tumor regression grading in patients with gastric carcinoma who had received neoadjuvant treatment were observed in [7]. The study also highlighted the outcomes of this survey, shedding light on tumor regression grading methodologies, particularly when applied to patients with gastric carcinoma who had undergone preoperative treatment. A study was conducted by Cruz-Roa et al [8] for the purpose of detecting invasive breast cancer in histopathological images. High-throughput adaptive sampling for whole-slide histopathology image analysis (HASHI) using convolutional neural networks was employed in this research. The research aimed to improve the efficiency and accuracy of breast cancer detection through automated image analysis. The utilization of cutting-edge deep learning techniques showed great potential for enhancing the diagnostic process for invasive breast cancer,

making it more efficient and precise. The utilization of deep learning for the analysis of histopathological images, with a specific focus on incorporating explanation methods was presented in [9]. The primary objective was to enhance the interpretability and transparency of deep learning models when applied to histopathological tasks. A novel dataset comprising gastric histopathology images [10] was developed. This dataset was developed to support computer-aided diagnostic applications related to gastric cancer, introduced a valuable resource for the research and development of AI systems in the field of gastric cancer diagnosis. With the dataset's availability, it is anticipated that advances in computer-aided diagnosis for gastric cancer will be facilitated, potentially enhancing early detection and patient care. [11] presents a significant contribution to computer vision. It discusses a method for training deep neural networks on a large-scale dataset for general visual representation learning. The authors propose a novel architecture and training scheme that leverages large-scale labeled and unlabeled data to achieve remarkable performance in a wide range of computer vision tasks, including object recognition and classification. This paper has had a notable impact on the field of computer vision and deep learning, demonstrating the importance of pre-trained representations for various vision tasks. It is particularly relevant in the context of transfer learning and can serve as a foundational reference for researchers and practitioners in the field. The work in [12] delves into techniques for unraveling the inner workings of deep neural networks. It offers insights into the field of neural network interpretability, addressing the challenge of understanding how these complex models make predictions. The paper explores various methods and tools for interpreting and visualizing the learned representations and decision-making processes within deep neural networks. Understanding deep learning models is crucial for their application in fields like computer vision, natural language processing, and healthcare. Architectures of CNNs, dataset characteristics, and the application of transfer learning were discussed in [13]. By investigating these key aspects, the paper provides valuable insights into the development of robust and effective CNN models for medical image analysis, particularly in the context of aiding diagnosis and detection in healthcare. This work contributes significantly to the intersection of deep learning and medical imaging. The authors in [14] proposed a novel architecture that utilized residual blocks to enable the training of very deep networks. A profound impact on the field of computer vision and deep learning, and ResNets have become a fundamental building block for various applications involving image and feature recognition. The study in [15] provides a comprehensive overview of the applications of deep learning in the field of medical image analysis, discusses how deep learning techniques, including convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have

been applied to tasks such as image segmentation, disease classification, and object detection in various medical imaging modalities. The study highlights the potential for deep learning to revolutionize healthcare by improving diagnostic accuracy and automating image interpretation. Ensemble learning involves combining predictions from multiple models to enhance classification accuracy. By leveraging the diversity of different models, ensemble learning mitigates errors and improves performance, particularly in medical image classification, where precision is crucial. The GasHisSDB dataset serves as the testing ground for evaluating how these deep learning models can work together to improve the accuracy and efficacy of medical image classification. The study explores their complementarity and potential in ensemble learning.

III. MATERIALS AND METHODS

A. Dataset: GasHisSDB

This study employs the GasHisSDB dataset to evaluate the performance of deep learning models, including VGG19, Inception V4, ResNeXt, and ResNet152. The goal is to explore the complementarity of these models and their effectiveness in ensemble learning for medical image classification. GasHisSDB is a public dataset containing three sub-datasets with a total of 245,196 images. These sub-datasets consist of images in three different sizes: 160x160, 120x120, and 80x80 pixels. The images are categorized into "normal" and "abnormal" classes. "Abnormal" images represent pathological slices with cancerous areas exceeding 50% of the image, while "normal" images depict typical pathological tissue slices without abnormalities. The study capitalizes on the diversity within the GasHisSDB dataset, taking into account the varying image dimensions. This approach recognizes that medical image classification tasks can vary in complexity, and the choice of image dimensions significantly impacts algorithm performance. The deep learning models under examination, VGG19, Inception V4, ResNeXt, and ResNet152, have a track record of excellence in image classification. They excel in capturing fine details and patterns, making them valuable for medical image analysis.

B. Deep learning models

This section focuses on the utilization of deep learning models for the classification of gastric cancer pathology images. Initially, the model undergoes training using training and validation sets, which are derived from three distinct sub-datasets within GasHisSDB. Following the training phase, the experiment utilizes a test set to assess the performance of these models. To ascertain the potential complementarity of these classifiers in the domain of deep learning, a thorough comparative analysis of the obtained classification results is conducted. This analysis involves the examination of various

evaluation metrics. The study employs three prominent deep learning models, specifically VGG19, Inception-V4, ResNeXt, and ResNet152, to carry out the classification task. These models are renowned for their efficacy in image classification and have been chosen for their ability to capture intricate features within the pathology images.

VGG [16], a convolutional neural network (CNN), represents an improvement over AlexNet and was jointly developed by the Visual Geometry Group and Google DeepMind in 2014. VGG19 is the most frequently employed variant of this architecture in image classification tasks. In the same year, Google introduced InceptionNet [17] at the ILSVRC competition. This marked the inception of a series of InceptionNet versions, with InceptionV4 standing out as a prominent member of this extensive family. Within the domain of image classification, ResNeXt [18] and ResNet152 [18] has gained widespread recognition and is often the model of choice for researchers and practitioners.

IV. EXPERIMENTAL ANALYSIS

The primary phase of the complementarity experiments is carried out on a laptop running the Windows 11 operating system. This computer boasts 16 gigabytes of RAM and is furnished with a 4-gigabyte NVIDIA Quadro RTX 4000. Coding for the experiment was done using Python 3.6 with pytorch version 1.7.1 deep learning framework. A total of 100 experimental epochs are executed to observe the classification outcomes of this dataset using various models. The assessment encompasses an in-depth analysis and evaluation of the classification deep learning models. The visual representation of the experimental process is depicted in Figure 1.

A. Evaluation metrics

The choice of evaluation metrics holds great significance in papers that seek to make complementarity comparisons. In the experiments conducted within this thesis, Accuracy takes precedence as the most crucial metric. In addition to Accuracy, Precision, Recall, Specificity, and F1-score (F1) are also included as selected metrics for evaluation as calculated in (1-5). These metrics are widely employed in comparative studies to assess the performance of classifiers. Their utilization aids in the comprehensive analysis of classifier performance and plays a pivotal role in identifying complementarities, thereby contributing to the enhancement and refinement of ensemble learning techniques.

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

$$\text{Precision} = \frac{TP}{TP+FP} \quad (2)$$

$$\text{Recall} = \frac{TP}{TP+FN} \quad (3)$$

$$\text{Specificity} = \frac{TN}{TN+FP} \quad (4)$$

$$\text{F1-score} = \frac{2 \times (\text{Precision} \times \text{Recall})}{(\text{Precision} + \text{Recall})} \quad (5)$$

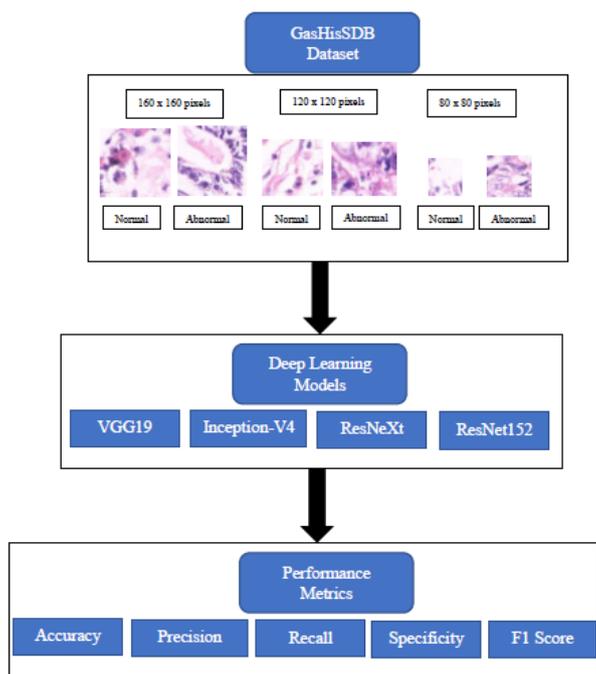


Figure 1. Experimental process

V. EVALUATION OF RESULTS

A. Evaluation of deep learning methods on 160x160 pixel subdatabase

This section deals with evaluation of deep learning methods such as VGG19, Inception-V4, ResNeXt and ResNet152 on 160x 160 pixel subdatabase. According to the Table I, on 160 x 160 pixel database, four deep learning models, namely VGG19, Inception-V4, ResNeXt, and ResNet152, were evaluated for their performance after 100 training epochs as shown in Table. I. The results of these experiments reveal interesting insights and offer a basis for model selection in various medical image classification

tasks. First, in terms of accuracy, VGG19 achieved an impressive accuracy rate of 95.9% in classifying "Normal" instances, but it slightly lagged behind in the "Abnormal" category with an accuracy of 93.78%. On the other hand, Inception-V4 displayed consistent accuracy rates in both categories, with a slightly lower overall accuracy of 94.57%. ResNeXt outperformed the other models with an accuracy of 96.09% in the "Normal" category, although it faced challenges in classifying "Abnormal" instances with an accuracy of 86.89%. Remarkably, ResNet152 excelled in both categories, achieving an accuracy of 98.02% in the "Normal" category and 87.9% in the "Abnormal" category. In terms of precision, recall, and specificity, ResNet152 consistently outperformed the other models in both the "Normal" and "Abnormal" categories. It achieved the highest precision, recall, and specificity values in the "Normal" category, emphasizing its capability to accurately identify "Normal" instances. Notably, in the "Abnormal" category, while ResNeXt achieved the highest precision, ResNet152 still exhibited competitive values, demonstrating a well-balanced performance across all categories. These results highlight the trade-offs between different models and their ability to perform effectively in classifying "Normal" and "Abnormal" instances within medical image datasets. While VGG19 demonstrated high accuracy in the "Normal" category, its performance in the "Abnormal" category was slightly inferior. Inception-V4 showed consistent but slightly lower accuracy overall. ResNeXt excelled in the "Normal" category but encountered challenges in classifying "Abnormal" instances. ResNet152 showcased a remarkable ability to maintain high accuracy, precision, recall, and specificity across both categories, making it a robust choice for medical image classification tasks as shown in Fig. 2 and Fig. 3.

TABLE I. DEEP LEARNING METHODS ON 160 X 160 PIXEL SUBDATABASE

Sub-database Size	Model	Quantity of epoch	Model size(MB)	Best Epoch	Training time(s)	Accuracy	Category	Precision	Recall	Specificity	F1 Score
160 x 160 pixels	VGG19	100	272.14	98	13654	95.9	Normal	96.56	93.87	96.89	96.03
							Abnormal	93.78	91.45	94.89	93.98
	Inception-V4	100	87.24	94	9894	94.57	Normal	94.78	96.89	92.76	95.87
							Abnormal	94.87	95.12	94.89	95.87
	ResNeXt	100	82.9	81	10941	96.09	Normal	94.43	94.87	95.87	96.78
							Abnormal	86.89	88.67	84.54	83.89
	ResNet152	100	32.18	98	2456	98.02	Normal	97.54	96.98	96.89	97.65
							Abnormal	87.9	82.98	89.09	92.9

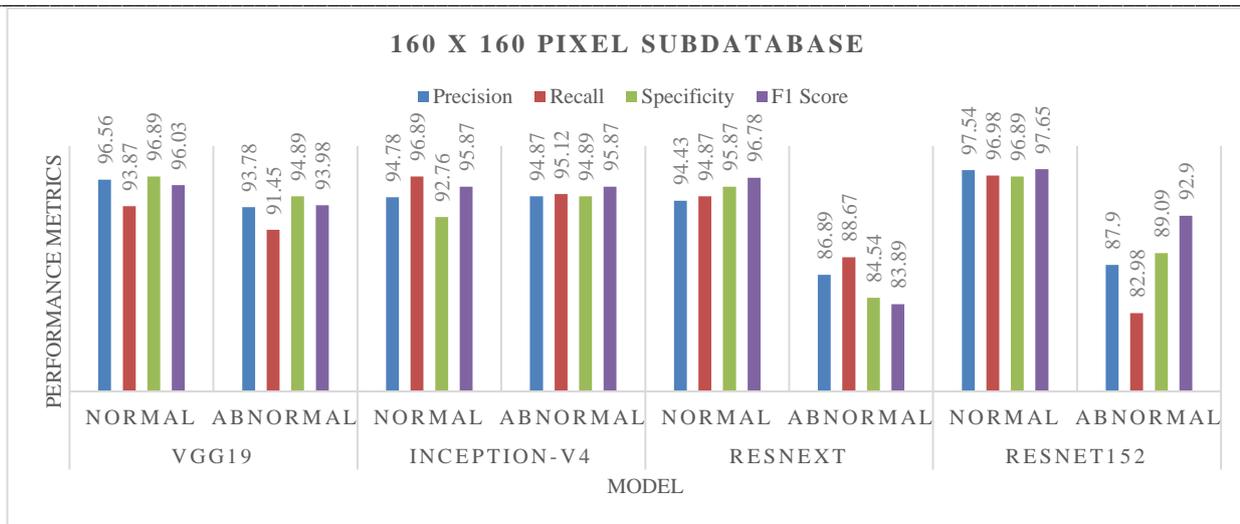


Figure 2. Performance metrics for 160 x 160 pixel subdatabase

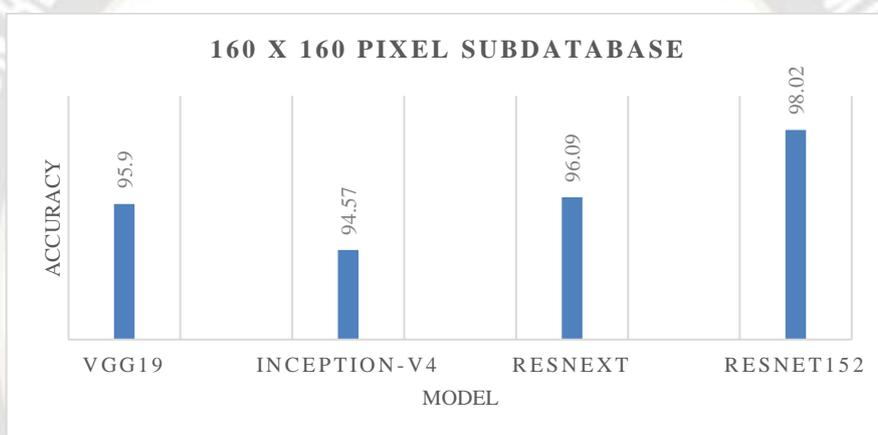


Figure 3. Accuracy for 160 x 160 pixel subdatabase

B. Evaluation of deep learning methods on 120 x120 pixel subdatabase

This section deals with evaluation of deep learning methods such as VGG19, Inception-V4, ResNeXt and ResNet152 on 120 x 120 pixelsubdatabase. According to the Table II, on 120 x 120 pixel database, four deep learning models (VGG19, Inception-V4, ResNeXt, and ResNet152) were evaluated after 100 training epochs as shown in Table. II. These results provide a comprehensive view of their performance in classifying "Normal" and "Abnormal" instances within medical image datasets. VGG19 demonstrated solid accuracy, achieving 92.9% in the "Normal" category and 95.87% in the "Abnormal" category. Notably, its precision and specificity were quite high in both categories, indicating its ability to accurately identify instances. Inception-V4 maintained a competitive accuracy rate of 94.12% overall. Its precision, recall, and specificity values were relatively balanced in both categories. ResNeXt showed strong performance with an accuracy of 96.98% in the "Normal" category, although its accuracy dropped to 89.09% in the

"Abnormal" category. The precision and specificity remained high in the "Normal" category but were less favorable in the "Abnormal" category. ResNet152 emerged as the top performer, achieving an impressive accuracy of 97.98% in the "Normal" category and 94.9% in the "Abnormal" category. It demonstrated high precision, recall, and specificity in both categories, underscoring its ability to consistently classify instances accurately. These results highlight the differences in model performance concerning the sub-database with 120 x 120-pixel images. VGG19 exhibited solid performance with high accuracy, particularly in the "Abnormal" category, where it excelled in precision and specificity. Inception-V4 maintained competitive accuracy and balanced precision and recall values across both categories. ResNeXt showcased strong performance in the "Normal" category but faced challenges with "Abnormal" instances. ResNet152 consistently outperformed the other models, achieving high accuracy, precision, recall, and specificity in both categories as shown in Fig. 4 and Fig. 5.

TABLE II. DEEP LEARNING METHODS ON 120 X 120 PIXEL SUBDATABASE

Sub-database Size	Model	Quantity of epoch	Model size(MB)	Best Eopch	Training time(s)	Accuracy	Category	Precision	Recall	Specificity	F1 Score
120 x 120 pixels	VGG19	100	272.14	100	13453	92.9	Normal	96.78	98.87	93.98	98.09
							Abnormal	95.87	95.87	97.98	94.12
	Inception-V4	100	87.24	96	9467	94.12	Normal	95.89	96.89	93.78	97.54
							Abnormal	95.98	92.7	98.02	95.87
	ResNeXt	100	82.9	95	10356	96.98	Normal	96.34	98.98	94.78	97.9
							Abnormal	89.09	85.89	91.87	84.78
	ResNet152	100	32.18	98	2671	97.98	Normal	91.94	91.89	85.89	91.12
							Abnormal	89.01	94.9	93.9	93.87

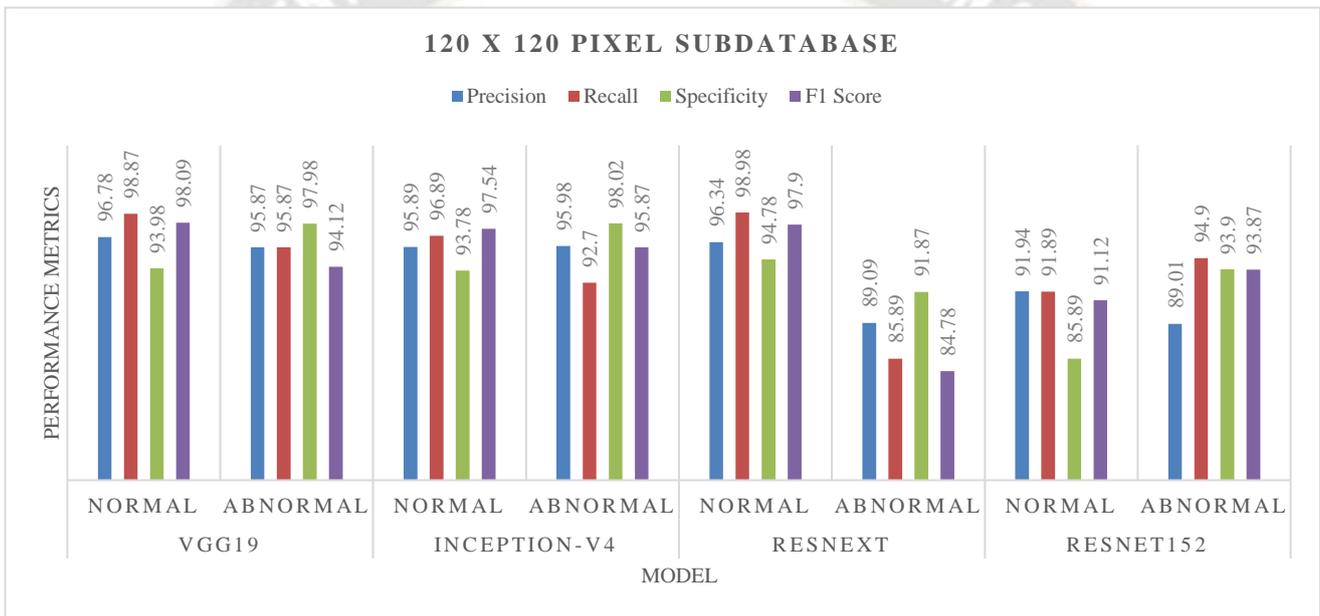


Figure 4. Performance metrics for 120 x 120 pixel subdatabase

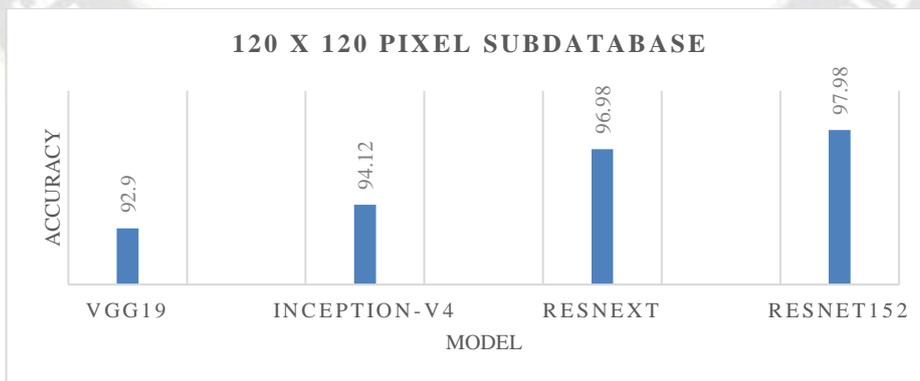


Figure 5. Accuracy for 120 x 120 pixel subdatabase

C. Evaluation of deep learning methods on 80 x 80 pixel subdatabase

This section deals with evaluation of deep learning methods such as VGG19, Inception-V4, ResNeXt and

ResNet152 on 80 x 80 pixelsubdatabase. According to the Table III, on 80 x 80 pixel database, four deep learning models (VGG19, Inception-V4, ResNeXt, and ResNet152) was evaluated after 100 training epochs as shown in Table. III. These results offer insights into the models' performance in

classifying "Normal" and "Abnormal" instances within medical image datasets. VGG19 exhibited a commendable accuracy of 93.76%, with particularly high precision, recall, and specificity in the "Normal" category. It also performed well in the "Abnormal" category, with precision and specificity values surpassing 95%. Inception-V4 achieved an impressive accuracy of 95.76%. Its precision, recall, and specificity values were well-balanced in both categories, making it a robust performer. ResNeXt displayed robust performance with an accuracy of 96.12% in the "Normal" category, though its accuracy was slightly lower in the "Abnormal" category at 88.09%. It maintained high precision and specificity in the "Normal" category but faced challenges with "Abnormal" instances. ResNet152 emerged as the top performer, achieving a remarkable accuracy of 98.67% in the "Normal" category and 95.12% in the "Abnormal" category. It

consistently demonstrated high precision, recall, and specificity in both categories, underscoring its ability to accurately classify instances. These results underscore the variations in model performance within the sub-database with 80 x 80-pixel images. VGG19 exhibited solid accuracy, especially in the "Normal" category, where it excelled in precision and specificity. Inception-V4 displayed a balanced performance across both categories, achieving high overall accuracy. ResNeXt performed well in the "Normal" category but faced challenges in the "Abnormal" category. ResNet152 consistently outperformed the other models, demonstrating high accuracy and precision while maintaining balance between categories as shown in Fig. 6 and Fig. 7.

TABLE III. DEEP LEARNING METHODS ON 80 X 80 PIXEL SUBDATABASE

SUB-DATABASE SIZE	Model	Quantity of epoch	Model size(MB)	Best Epoch	Training time(s)	Accuracy	Category	Precision	Recall	Specificity	F1 Score
80 x 80 pixels	VGG19	100	272.14	98	14345	93.76	Normal	97.87	98.09	94.76	97.89
							Abnormal	95.78	94.78	98.9	93.67
	Inception-V4	100	87.24	98	9656	95.76	Normal	96.12	97.98	92.35	98.78
							Abnormal	95.78	92.67	97.9	96.23
	ResNeXt	100	82.9	94	11432	96.12	Normal	96.12	97.34	94.78	98.9
							Abnormal	88.09	84.78	91.76	87.87
	ResNet152	100	32.18	100	2356	98.67	Normal	90.21	90.87	94.87	97.9
							Abnormal	88.78	95.12	96.89	96.9

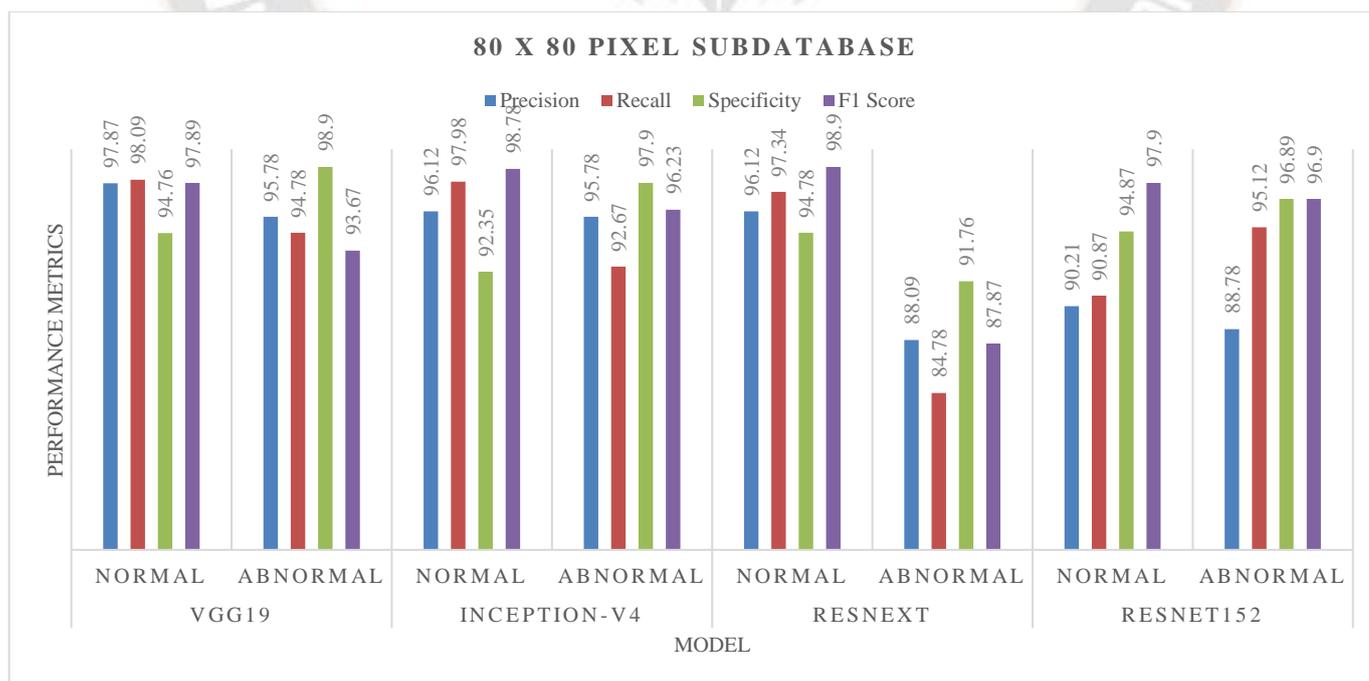


Figure 6. Performance metrics for 80 x 80 pixel subdatabase

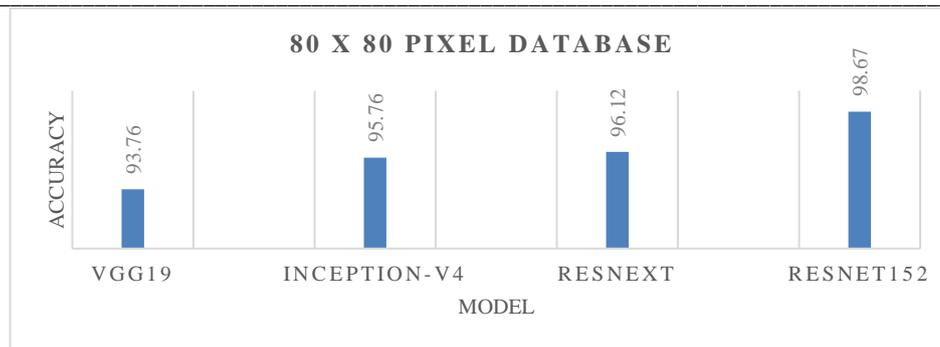


Figure 7. Accuracy for 80 x 80 pixel subdatabase

VI. CONCLUSION AND FURTHER DIRECTIONS

In conclusion, gastric cancer poses a significant global health challenge, being one of the most prevalent and deadliest forms of cancer. Timely diagnosis is crucial for effective treatment, and histopathological examination has long been the gold standard for diagnosing gastric cancer. With the advancement of computer technology, the field of medical image analysis has benefited greatly from digital tools that assist pathologists in making accurate diagnoses from pathological images. Ensemble learning has emerged as a valuable technique for improving the precision of diagnostic algorithms. By combining the strengths of multiple complementary learning models, ensemble learning holds the potential to enhance the accuracy and robustness of gastric cancer diagnosis. The experiments conducted in this study represent a significant step forward in the evaluation of deep learning classifiers for gastric cancer diagnosis. The use of three distinct sub-databases within GasHisSDB, each with varying image sizes, allowed for a comprehensive assessment of model performance in different scenarios. The four deep learning classifiers, VGG19, Inception-V4, ResNeXt, and ResNet152, were evaluated using key performance metrics such as accuracy, precision, recall, specificity, and F1 score. These metrics provide a well-rounded view of each model's capability to distinguish between "Normal" and "Abnormal" instances within gastric cancer images. Notably, ResNet152 consistently demonstrated outstanding performance across all sub-databases, with remarkable accuracy and precision. Its ability to maintain a high level of accuracy while ensuring a balanced performance in both categories makes it a promising choice for gastric cancer diagnosis. The findings from these experiments contribute to the field of medical image analysis and gastric cancer diagnosis, providing valuable insights for researchers and practitioners. The continued development of advanced diagnostic tools and techniques is essential in the ongoing fight against gastric cancer, with the goal of achieving earlier detection and more effective treatment. Future directions focus on expanding the dataset and exploring additional deep learning models. Moreover, research can delve into transfer learning techniques and real-world clinical applications for

gastric cancer diagnosis, ultimately striving for more accurate and efficient diagnostic tools to improve patient outcomes.

ACKNOWLEDGMENT

The authors are thankful to the management of Acharya Nagarjuna University, Nagarjuna Nagar, Guntur, Andhra Pradesh, India that greatly assist research.

REFERENCES

- [1] Hyuna Sung, Jacques Ferlay, Rebecca L. Siegel, Mathieu Laversanne, Isabelle Soerjomataram, Ahmedin Jemal, and Freddie Bray. Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians*, 0(0):1–41, 2021.
- [2] Jing Wang and Xiuping Liu. Medical image recognition and segmentation of pathological slices of gastric cancer based on deeplab v3+neural network. *Computer Methods and Programs in Biomedicine*, page 106210, 2021.
- [3] Harshita Sharma, Norman Zerbe, Daniel Heim, Stephan Wienert, Hans-Michael Behrens, Olaf Hellwich, and Peter Hufnagel. A multiresolution approach for combining visual information using nuclei segmentation and classification in histopathological images. In *VISAPP (3)*, pages 37–46, 2015.
- [4] Peng Jin, Xiaoyan Ji, Wenzhe Kang, Yang Li, Hao Liu, Fuhai Ma, Shuai Ma, Haitao Hu, Weikun Li, and Yantao Tian. Artificial intelligence in gastric cancer: A systematic review. *Journal of Cancer Research and Clinical Oncology*, pages 1–12, 2020.
- [5] Sai Chandra Kosaraju, Jie Hao, Hyun Min Koh, and Mingon Kang. Deep-hipo: Multi-scale receptive field deep learning for histopathological image analysis. *Methods*, 179:3–13, 2020.
- [6] Jun Zhang, Zhiyuan Hua, Kezhou Yan, Kuan Tian, Jianhua Yao, Eryun Liu, Mingxia Liu, and Xiao Han. Joint fully convolutional and graph convolutional networks for weakly-supervised segmentation of pathology images. *Medical image analysis*, 73:102183, 2021.
- [7] Andrianos Tsekrekos, Sönke Detlefsen, Robert Riddell, James Conner, Luca Mastracci, Kieran Sheahan, Jayant Shetye, Lars Lundell, and Michael Vieth. Histopathologic tumor regression grading in patients with gastric carcinoma submitted to neoadjuvant treatment: results of a delphi survey. *Human pathology*, 84:26–34, 2019.

- [8] Cruz-Roa, A., Gilmore, H., Basavanthally, A., Feldman, M., Ganesan, S., Shih, N., Tomaszewski, J., Madabhushi, A., Gonzalez, F., 2018. High-throughput adaptive sampling for whole-slide histopathology image analysis (HASHI) via convolutional neural networks: Application to invasive breast cancer detection. *PLOS ONE* 13, e0196828. doi:10.1371/journal.pone.0196828.
- [9] Hagele, M., Seegerer, P., Lapuschkin, S., Bockmayr, M., Samek, W., Klauschen, F., Muller, K.R., Binder, A., 2020. Resolving challenges in deep learning-based analyses of histopathological images using explanation methods. *Scientific Reports* 10, 6423. doi:10.1038/s41598-020-62724-2.
- [10] Hu, W., Li, C., Li, X., Rahaman, M.M., Ma, J., Zhang, Y., Chen, H., Liu, W., Sun, C., Yao, Y., Sun, H., Grzegorzec, M., 2022. GasHisSDB: A new gastric histopathology image dataset for computer aided diagnosis of gastric cancer. *Computers in Biology and Medicine* 142, 105207. doi:10.1016/j.combiomed.2021.105207.
- [11] Kolesnikov, A., Beyer, L., Zhai, X., Puigcerver, J., Yung, J., Gelly, S., Houlsby, N., 2020. Big transfer (bit): General visual representation learning, in: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (Eds.), *Computer Vision – ECCV 2020*, Springer International Publishing, Cham. pp. 491–507.
- [12] Montavon, G., Samek, W., Muller, K.R., 2018. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing* 73, 1–15. doi:10.1016/j.dsp.2017.10.011.
- [13] Hoo-Chang Shin, Holger R Roth, Mingchen Gao, Le Lu, Ziyue Xu, Isabella Nogues, Jianhua Yao, Daniel Mollura, and Ronald M Summers. Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning. *IEEE transactions on medical imaging*, 35(5):1285–1298, 2016. doi:https://doi.org/10.1109/TMI.2016.2528162.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [15] Justin Ker, Lipo Wang, Jai Rao, and Tchoyoson Lim. Deep learning applications in medical image analysis. *Ieee Access*, 6:9375–9389, 2017. doi:https://doi.org/10.1109/ACCESS.2017.2788044.
- [16] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [17] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. doi:https://doi.org/10.1109/CVPR.2016.308.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. doi:https://doi.org/10.1109/CVPR.2016