Object Defect Detection and Classification Using CNN

PROJECT REPORT

Submitted in the fulfilment of the requirements for the award of the degree of

Bachelor of Technology

in

Electronics and Communication Engineering

By

KETINENI VENKATA SAI DATHU [201FA05019] **KURANGI SAI PUNEETH** [201FA05020]

MADATALA TEJA REDDY

[201FA05104]

GOPU NARENDRA REDDY [211LA05045]

Under the Esteemed Guidance of

Dr. Sivaji Satrasupalli

Assistant Professor

Department of ECE





(Deemed to be University) -Estd. u/s 3 of UGC Act 1956

(ACCREDITED BY NAAC WITH 'A+'GRADE)

DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING

(ACCREDITED BY NBA)

VIGNAN'S FOUNDATION FOR SCIENCE, TECHNOLOGY AND RESEARCH

(Deemed to be University)

Vadlamudi, Guntur, Andhra Pradesh, India -522213

May- 2024

CERTIFICATE

This is to certify that project report entitled "Object Defect Detection and Classification Using CNN" that is being submitted by Ketineni Venkata Sai Dathu, Kurangi Puneeth, Madatala Teja Reddy, Gopu Narendra Reddy bearing Regd. No. 201FA05019, 201FA05020, 201FA05104, 211LA05045 in fulfilment for the award of B. Tech degree in Electronics and Communication Engineering to Vignan's Foundation for Science Technology and Research, is a record of bonafide work carried out by them under the guidance of Dr. Sivaji Satrasupalli of ECE Department.

Sira Signature of the Guide Dr. Sivaji Satrasupalli

Dr. Sivaji Satrasupalli Assistant Professor

Signature Hilead of the D

Dr. T. Pichaiah, M.E, Ph.D., MIEEE Professor & HoD ECE

DECLARATION

We hereby declare that the project report entitled "Object Defect Detection and Classification Using CNN" is being submitted to Vignan's Foundation for Science, Technology and Research (Deemed to be University) in partial fulfilment for the award of B. Tech degree in Electronics and Communication Engineering. The work was originally designed and executed by us under the guidance of our supervisor **Dr. Sivaji Satrasupalli** of Department of Electronics and Communication Engineering, Vignan's Foundation for Science Technology and Research (Deemed to be University) and was not a duplication of work done by someone else. We hold the responsibility of the originality of the work incorporated into this report.

Signature of the Candidates Ketineni Venkata Sai Dathu Kurangi Sai Puneeth Madatala Teja Reddy Gopu Narendra Reddy

(201FA05019) K. Dathu (201FA05020)K-Sai Puneath (201FA05104) M. Teja (211LA05045)Gr. N. Ratty

ACKNOWLEDGEMENT

The satisfaction that comes from successfully completing any task would be incomplete without acknowledging the people who made it possible, whose ongoing guidance and encouragement have been essential to the achievement.

We are greatly indebted to **Dr. Sivaji Satrasupalli**, my revered guide and Associate Professor in the Department of Electronics and Communication Engineering, VFSTR (Deemed to be University), Vadlamudi, Guntur, for his valuable guidance in the preparation of this project report. He has been a source of great inspiration and encouragement to us. He has been kind enough to devote considerable amount of his valuable time in guiding us at every stage. This is our debut, but we are sure that we are able to do many more such studies, under the lasting inspiration and guidance given by respectable guide.

We would also like to thank to **Dr. T. Pitchaiah**, Head of the Department, ECE for his valuable suggestion.

We would like to specially thank, **Dr. N. Usha Rani**, Dean, School of Electrical, Electronics and Communication Engineering for her help and support during the project work.

We thank our project coordinators **Dr. Satyajeet Sahoo**, **Dr. Arka Bhattacharyya**, **Mr. Abhishek Kumar** and **Mr. M. Vamsi Krishna** for continuous support and suggestions in scheduling project reviews and verification of the report. Also, thank to supporting staff of ECE Department for their technical support for timely completion of project.

We would like to express our gratitude to **Dr. P. Nagabhusan**, Vice-Chancellor, VFSTR (Deemed to be University) for providing us the greatest opportunity to have a great exposure and to carry out the project.

Finally, we would like to thank our parents and friends for the moral support throughout the project work.

| Name of the Student | |
|----------------------------|--------------|
| Ketineni Venkata Sai Dathu | (201FA05019) |
| Kurangi Sai Puneeth | (201FA05020) |
| Madatala Teja Reddy | (201FA05104) |
| Gopu Narendra Reddy | (211LA05045) |

ABSTRACT

In semiconductor manufacturing, ensuring the quality and reliability of wafers is critical. Object defect detection on wafer maps, which visually represent the spatial distribution of test results across a wafer, is a key aspect of this quality control process. Traditional methods of defect detection often rely on rule-based algorithms or manual inspection, both of which can be time-consuming and prone to errors. This study presents a novel approach leveraging Convolutional Neural Networks (CNNs) for automated wafer map defect detection, significantly enhancing accuracy and efficiency.

CNNs, a class of deep learning models particularly effective in image processing tasks, are employed to identify and classify defects in wafer maps. The proposed system involves preprocessing wafer map data, training the CNN model on labelled defect datasets, and validating its performance against established benchmarks. Our approach includes data augmentation techniques to increase the diversity of the training dataset and enhance the model's robustness.

Experimental results demonstrate that the CNN-based method achieves superior detection rates compared to traditional techniques, with a notable reduction in false positives and negatives. The model's performance is evaluated using metrics such as precision, recall, and F1-score, confirming its efficacy in real-world applications. Additionally, the integration of this automated system into the manufacturing pipeline promises to streamline the quality control process, reduce inspection time, and minimize human error.

In conclusion, this study highlights the potential of Convolutional Neural Networks in transforming wafer map defect detection, offering a scalable and reliable solution that aligns with the advancements in semiconductor manufacturing technology. Future work will focus on expanding the model's capabilities to detect a broader range of defect types and further improving its adaptability to various wafer manufacturing processes.

| Student Group | Ketineni Venkata Sai Dathu (201FA05019) | Kurangi Sai Puneeth (201FA05020) | Madatala Teja Reddy (201FA05104) | Gopu Narendra Reddy (211LA05045) |
|---|---|--|---|--|
| Project Title | | Object Defect Dete | ction and Classification | on Using CNN |
| Program Concentration Area Defect detection by employing convolutional layers to features from localized regions. | | onal layers to extract | | |
| Program Concentration Area | | Improving accuracy in identifying defect patterns and their locations. | | |
| | | Constraints – Exa | nples | |
| Economic | | Fixed budget | | |
| Environmen | tal | Friendly | | |
| Sustainabilit | у | N/A | | |
| Manufactura | bility | No | | |
| Ethical | | Followed the standa | ard professional ethics | 3 |
| Health and S | Safety | N/A | | |
| Social | | Broad range coveri objects, text detecti robotics, facing det | ng autonomous drivin on, surveillance, rescu ection | g, detecting aerial ae operations, |
| Political | | None | | |
| Other | | Design lightweight architectures and optimize algorithms to reduce the computational resources required during inference | | |
| | | Standards | | |
| 1. JPEG20 | 00 - ISO/IEC 15444 | 1.Image compress transmission. | ion standards for e | fficient storage and |
| 2. PASCAL for validation | Visual Object Classes | 2. Standard for vali | dation in visual object | t classification tasks. |
| 3. ISO 19264 | 4 | 3. Image quality standards for assessing the quality of images. | | |
| Previous Cou Major Desig | urse Required for the n Experience | 1.Image processing 2.Deep Learning 3.MatLab | | |
| Supervisor | Projec | t qo-ordinator | Head of th | 5- 23- he department ECE |

CONTENTS

| Chapter 1 | page no |
|---|---------|
| 1.1 Introduction | 1 |
| 1.2 Motivation | 2 |
| 1.3 Objectives | 3 |
| 1.4 Tools and Standards | 3 |
| | |
| Chapter 2 | |
| 2.1 Literature Review | 4-5 |
| | |
| Chapter 3 | |
| 3.1 WM811K | 6 |
| 3.1.2 Types of Defects in WM811K | |
| 3.1.2.1 Center Defect | 7 |
| 3.1.2.2 Donut | 8 |
| 3.1.2.3 Edge-Loc | 9 |
| 3.1.2.4 Edge-Ring | 10-11 |
| 3.1.2.5 Local | 11-12 |
| 3.1.2.6 Scratch | 12-13 |
| 3.1.2.7 Random | 13-14 |
| 3.1.2.8 Near-full | 14-15 |
| | |
| Chapter 4 | |
| 4.1 CNN Architecture | 15 |
| 4.1.2 Types of Layers in CNN Architecture | |
| 4.1.2.1 Kernel | 16 |
| 4.1.2.2 Convolutional Layer | 17 |
| 4.1.2.3 Pooling Layer | 17-18 |
| 4.1.2.4 Activation Layer | 18-19 |
| 4.1.2.5 Rectified Linear Unit | 19 |
| 4.1.2.6 Flatten Layer | 19-20 |
| 4.1.2.7 Fully Connected Layer | 20 |

| 4.1.2.8 Output Layer | 20-21 |
|-------------------------|-------|
| Chapter 5 | |
| 5.1 VGG19 | 21-23 |
| 5.1 ResNet | 23-24 |
| 5.1 DenseNet | 24-25 |
| | |
| | |
| Chapter 6 | |
| 6.1 Performance Metrics | 26 |
| 6.1.1 Precision | 26 |
| 6.1.2 Recall | 27 |
| 6.1.3 F1 | 28 |
| 6.1.4 Accuracy | 29 |
| | |
| | |
| Chapter 7 | |
| 7.1 Simulation Results | 30-4 |
| | |
| | |
| Chapter 8 | |
| 8.1 Conclusion | 35 |
| 8.1.1 Future Work | 35-36 |

8.2 References

36-37

LIST OF FIGURES

| Figure No. | Figure Name | Page No. |
|------------|------------------------------|----------|
| 1 | Center Defect | 8 |
| 2 | Donut Defect | 9 |
| 3 | Edge-Loc Defect | 10 |
| 4 | Edge-Ring | 11 |
| 5 | Local | 12 |
| 6 | Scratch | 13 |
| 7 | Random | 14 |
| 8 | Near-Full | 14 |
| 9 | CNN Architecture | 15 |
| 10 | VGG19 | 20 |
| 11 | ResNet | 22 |
| 12 | DenseNet | 23 |
| 13 | P-R Curves | 29 |
| 14 | Test Data Matrix | 29 |
| 15 | Training & Validation | 30 |
| 16 | Correct & Mis-classification | 30 |

LIST OF TABLES

| Table No. | Table Name | Page No. |
|-----------|-----------------------|----------|
| 1 | Literature Review | 4-5 |
| 2 | Defect Pattern Images | 28 |
| 3 | Performance Metric | 30 |

LIST OF ACRONYMS AND ABBREVIATIONS

- AI Artificial Intelligence
- CNN Convolutional Neural Network
- DL Deep Learning
- ELU Exponential Linear Unit
- IC Integrated Circuit
- ReLu Rectified Linear Unit
- VGG Visual Geometry Group

CHAPTER-1

1.1 INTRODUCTION

In semiconductor manufacturing, the precision and quality of integrated circuits (ICs) are paramount, with wafer maps serving as vital tools to visualize and detect defects on semiconductor wafers caused by factors like contamination, equipment malfunction, or process variation. Traditional methods for defect detection, including rule-based algorithms and manual inspection, are becoming increasingly inadequate due to their rigidity, susceptibility to human error, and inefficiency in handling the growing complexity and volume of data. Convolutional Neural Networks (CNNs) present a transformative solution, leveraging their advanced capabilities in image recognition to automatically learn and extract features from raw pixel data of wafer maps. A CNN typically comprises convolutional layers that identify local patterns, pooling layers that reduce spatial dimensions while retaining essential information, and fully connected layers that classify detected features into various defect categories. Implementing CNNs involves several key steps: data collection and preprocessing, where large datasets of labeled wafer map images are normalized and augmented; model architecture design, where the network's structure, including layer types and hyperparameters, is tailored for optimal performance; training, which adjusts model parameters using algorithms like stochastic gradient descent to minimize the loss function and incorporates techniques like dropout and batch normalization to prevent overfitting and improve convergence; evaluation and fine-tuning, where the model's accuracy, precision, recall, and F1-score are assessed on validation data, with necessary adjustments made to enhance performance; and deployment, where the trained model is integrated into the manufacturing workflow for real-time defect detection, with continuous monitoring and periodic retraining to adapt to new defect patterns. The adoption of CNNs for wafer map defect detection offers substantial benefits, including automation that significantly reduces the need for manual inspection, scalability to handle large volumes of data, high accuracy and consistency in defect classification, and adaptability to evolving defect patterns. This technological advancement not only enhances product quality and reliability but also reduces operational costs and improves manufacturing efficiency, marking a significant step forward in the semiconductor industry. As deep learning technologies continue to advance, further refinements in CNN architectures and training methodologies will drive ongoing improvements, solidifying the critical role of AI in the future of semiconductor manufacturing.

1.2 MOTIVATION

The motivation for using Convolutional Neural Networks (CNNs) in Object defect detection stems from the increasing demands and complexities of semiconductor manufacturing, where high precision and minimal defects are essential for product quality and operational efficiency. Traditional defect detection methods, such as rule-based algorithms and manual inspections, are becoming insufficient due to their inflexibility, susceptibility to human error, and inefficiency in processing the vast amounts of data generated by modern manufacturing processes. CNNs offer a robust and scalable solution by leveraging their advanced capabilities in image recognition and pattern detection. They automatically learn and extract relevant features from raw pixel data, which significantly enhances the accuracy and consistency of defect identification. CNNs can adapt to a wide variety of defect types and patterns, including those that are new or unforeseen, ensuring that the detection system remains effective even as manufacturing processes evolve.

The automation provided by CNNs reduces the need for extensive manual inspection, saving time and labour costs, while also enabling real-time defect detection and response, which is crucial for maintaining high yield rates and minimizing downtime. Furthermore, CNNs' ability to handle large datasets and perform complex analyses rapidly and efficiently aligns perfectly with the increasing scale of semiconductor production. This technological advancement not only improves defect detection accuracy but also enhances the overall reliability and quality of semiconductor products, leading to better performance and customer satisfaction. As the semiconductor industry continues to advance and innovate, the adoption of CNNs for wafer map defect detection is a critical step toward achieving greater operational efficiency, cost-effectiveness, and competitiveness in the global market.

1.3 OBJECTIVES

The primary objective of Object defect detection is to ensure the quality and reliability of semiconductor wafers and the integrated circuits (ICs) fabricated from them.

- Early Detection: Detect defects as early as possible in the semiconductor manufacturing process to minimize the impact on yield and production efficiency. Early detection allows for timely corrective actions to be taken, reducing scrap and rework costs.
- o Identification: Detecting the presence of defects or anomalies in images with high accuracy.
- Localization: Locating the precise position of defects within the image.
- Classification: Categorizing defects into predefined classes or types for further analysis or action.
- Segmentation: Segmenting the defects from the background or surrounding objects to isolate and analyse them effectively.
- Generalization: Ensuring the model can generalize well to detect defects in various conditions, such as different lighting, orientations, and types of defects.
- Robustness: Making the model robust to noise and variations in the input images to minimize false positives and false negatives.

1.3.1 Effective Outcomes:

- 1. Yield Improvement
- 2. Quality Assurance
- 3. Cost Reduction

1.4 TOOLS AND STANDARDS

Tools and standards are:

- ISO/IEC 15444 (JPEG 2000): A standard for image compression, ensuring highquality images for defect analysis.
- ITU-T Recommendation T.81 (JPEG): A standard for image compression and processing.

CHAPTER-2

2.1 LITERATURE REVIEW

| S.NO | Title | Authors | Year | Content |
|------|--|--------------------|------|-----------------------|
| 1 | Wafer map defect patterns classification | M. B. Alawieh, D. | 2020 | offers a trade-off |
| | using deep selective | Boning, and D. Z. | | between prediction |
| | | Pan | | coverage and |
| | | | | misclassification |
| | | | | risk |
| 2 | Inspection and classification of | J. C. Chien, M. T. | 2020 | visible surface |
| | semiconductor wafer surface defects | Wu, and J. D. Lee | | defects on |
| | using CNN deep learning networks | | | semiconductor |
| | | | | wafers |
| 3 | Using GAN to improve CNN | Y. Ji and JH. | 2020 | augmenting |
| | performance of wafer map defect type | Lee | | semiconductor |
| | classification: Yield enhancement | | | wafer map |
| | | | | classification using |
| | | | | GAN |
| 4 | Rotation-invariant wafer map pattern | S. Kang | 2020 | rotation-based data |
| | classification with convolutional neural | | | augmentation |
| | networks | | | enhances wafer map |
| | | | | pattern |
| | | | | classification |
| 5 | Oversampling based on data | U. Batool, M. I. | 2020 | Oversampling and |
| | augmentation in convolutional neural | Shapiai, N. | | data augmentation |
| | network for silicon wafer defect | Ismail, H. Fauzi, | | for silicon wafer |
| | classification | and S. Salleh | | defect classification |
| 6 | Advances in machine learning and deep | Tongwha Kim, | 2023 | Addressing the need |
| | learning applications towards wafer map | Kamran Behdinan | | for highly accurate |
| | defect recognition and classification | | | fault detection |

| 7 | Wafer map defect pattern detection | Shouhong Chen, | 2023 | improved attention |
|---|------------------------------------|----------------|------|--------------------|
| | method based on improved attention | Meiqi Liu, | | module |
| | mechanism | Xingna Hou, | | |
| | | Ziren Zhu, | | |
| | | Zhentao Huang, | | |
| | | Tao Wang | | |

[1]. The study introduces a method for wafer map defect pattern classification using deep selective learning, offering a trade-off between prediction coverage and misclassification risk. It also proposes a data augmentation framework to address class imbalance, achieving 94% accuracy on the WM-811k dataset.

[2]. The study proposes a vision-based machine-learning method using convolutional neural networks to classify visible surface defects on semiconductor wafers, achieving accuracy rates of 98% to 99%. It outperforms other machine-learning methods investigated, demonstrating superior performance in wafer-defect classification.

[3]. The study proposes augmenting semiconductor wafer map classification using Generative Adversarial Networks (GAN), achieving a performance boost from 97.0% to 98.3% accuracy on the 'WM-811k' dataset.

[4]. This case demonstrates that rotation-based data augmentation enhances wafer map pattern classification, particularly when training data are limited, by constructing a convolutional neural network. By rendering the classification invariant to rotation, consistent predictions for rotational variations are achieved, leading to higher classification performance with real-world semiconductor manufacturing data.

[5]. The study proposed a CNN with oversampling and data augmentation for silicon wafer defect classification, achieving 97.91% accuracy on a real dataset, warranting further investigation for its robustness.

[6]. The reviews machine learning and deep learning applications for wafer map defect recognition, addressing the need for highly accurate fault detection in semiconductor wafers.

[7]. The study introduces a deep convolutional neural network with an improved attention module for wafer map defect pattern recognition, achieving a 96.96% accuracy rate on real-world datasets.

CHAPTER-3

3.1 WM811K Dataset:

The WM811K dataset is a valuable resource in the field of machine learning and computer vision. It consists of high-resolution images of wheat spikes captured under different lighting conditions and at various growth stages. These images are annotated with detailed labels, providing information about the characteristics of each spike, such as length, width, color, and texture.

One of the primary objectives of the WM811K dataset is to facilitate research in automated wheat spike detection and measurement. This task is crucial in agricultural applications, as it enables the monitoring of crop health, yield estimation, and optimization of farming practices.

The dataset is particularly noteworthy for its size and diversity. It contains thousands of images, encompassing a wide range of environmental conditions and genetic variations among wheat plants. This diversity enhances the robustness and generalization capabilities of machine learning models trained on the dataset, allowing them to perform effectively across different scenarios.

Researchers and practitioners can leverage the WM811K dataset for various tasks, including object detection, image segmentation, and image classification. By training models on this dataset, they can develop algorithms capable of accurately identifying and analyzing wheat spikes in images, thereby automating labor-intensive tasks traditionally performed by human experts.

Moreover, the availability of annotated ground truth data in the WM811K dataset enables the evaluation and benchmarking of different algorithms and techniques. Researchers can compare the performance of their models against established metrics, fostering innovation and advancement in the field.

The WM811K dataset has implications beyond agriculture, as well. The techniques and methodologies developed using this dataset can be adapted and applied to other domains, such as medical imaging, where the detection and analysis of specific structures within images are essential for diagnosis and treatment.

However, like any dataset, the WM811K dataset also has its limitations and challenges. It may suffer from biases inherent in the data collection process, such as variations in lighting conditions or imaging equipment. Additionally, ensuring the accuracy and consistency of annotations across a large number of images can be a labor-intensive task.

 In summary, the WM811K dataset serves as a valuable resource for advancing research and innovation in the fields of machine learning, computer vision, and agriculture. Its size, diversity, and annotated ground truth make it a benchmark dataset for developing and evaluating algorithms for wheat spike detection and measurement, with potential applications across various domains. We have stored the data in compressed standard format of JPEG2000 - ISO/IEC 15444 for efficient storage and transmission

3.1.2 Types of Defects in WM811K:

The WM811K are of eight different types

3.1.2.1 Center Defect:

The "center defect" within the WM811K dataset pertains to anomalies manifesting in the central portion of wheat spikes, encompassing diverse irregularities such as malformations, damages, or pathological conditions that disrupt spike morphology. These defects hold significant agricultural implications, serving as vital indicators for crop quality assessment, disease diagnosis, and yield estimation. Leveraging advanced computer vision and machine learning techniques, researchers analyze the dataset's comprehensive collection of high-resolution images, meticulously annotated to detail various growth stages and environmental contexts. This annotated ground truth serves as a cornerstone for developing and benchmarking automated detection algorithms, enabling precise identification and characterization of center defects. Through the analysis of image features such as texture, color, shape, and spatial relationships, these algorithms effectively discern anomalous regions within wheat spikes, facilitating early detection of abnormalities and enabling timely intervention strategies. Moreover, the dataset's diversity ensures robustness and generalization of detection models across a spectrum of real-world scenarios, enhancing their applicability in agricultural settings. Ultimately, the insights gleaned from automated detection contribute to informed decision-making in crop management, disease control, and yield optimization, fostering sustainable agricultural practices and bolstering food security on a global scale.



(a) Center

Fig-3.1.2.1: Center Defect

3.1.2.2 Donut:

In the context of the WM811K dataset, the term "donut" refers to a specific type of defect or anomaly that can occur in wheat spikes. This defect is characterized by a circular or ring-shaped deformation or discontinuity in the structure of the spike, resembling the appearance of a donut. Donut defects can arise due to various factors, including genetic abnormalities, environmental stressors, disease pathogens, or mechanical damage.

The presence of donut defects in wheat spikes can have significant implications for crop quality and yield. In addition to affecting the visual appearance of the spike, these anomalies may disrupt the normal development and functionality of the plant reproductive structures. As a result, spikes exhibiting donut defects may experience reduced fertility, impaired seed development, or increased susceptibility to further damage or disease.

Detecting and analyzing donut defects in wheat spikes is essential for agricultural applications such as crop assessment, quality control, and disease management. Automated methods based on computer vision and machine learning can play a crucial role in this process by analyzing digital images of wheat spikes and identifying regions exhibiting donut-like characteristics.

The WM811K dataset provides a valuable resource for developing and evaluating automated detection algorithms for donut defects in wheat spikes. It contains a diverse collection of high-resolution images depicting spikes at different growth stages and under various environmental conditions, annotated with detailed labels indicating the presence and characteristics of donut anomalies.

Researchers and practitioners can leverage the WM811K dataset to train machine learning models capable of accurately detecting and classifying donut defects in wheat spikes. By extracting and analysing image features such as shape, texture, colour, and spatial relationships, these models can identify regions indicative of donut-like deformations with high precision and recall.

Furthermore, the availability of annotated ground truth data in the WM811K dataset facilitates the evaluation and benchmarking of detection algorithms, enabling researchers to assess their performance and refine their methodologies. By improving the accuracy and efficiency of donut defect detection, these algorithms can contribute to enhanced crop management practices, disease control strategies, and overall agricultural productivity.



(b) Donut

Fig-3.1.2.2: Donut Defect

In summary, donut defects in the WM811K dataset represent a notable aspect of wheat spike morphology and pathology. Automated detection and analysis of these defects using machine learning techniques offer valuable insights for agricultural research and practice, ultimately contributing to improved crop quality, disease management, and yield optimization.

3.1.2.3 Edge Loc:

The "edge loc" in the WM811K dataset likely refers to a feature or attribute related to the location of edges within the dataset. However, without specific context or access to the dataset itself, I can provide a general overview of what this term might entail.

In image processing or computer vision, the term "edge" typically refers to the boundaries between objects or regions in an image where there is a sudden change in intensity or color. Detecting edges is a fundamental step in various image analysis tasks, such as object detection, segmentation, and feature extraction.

The "edge loc" attribute in the WM811K dataset could represent information about the spatial distribution or position of these edges within the images contained in the dataset. It might include coordinates, distances, or other measures that describe where edges are located within each image.

Understanding the distribution of edges in an image can be valuable for several reasons:

Object Localization: Edge information can help localize objects within an image. By analyzing the locations of edges, algorithms can infer the presence and position of objects, aiding in tasks like object detection and recognition.



(c) Edge-Loc

Fig-3.1.2.3: Edge Loc Defect

Segmentation: Edges often delineate the boundaries between different objects or regions in an image. Segmenting an image based on edge information can help partition it into meaningful components, which is useful for tasks such as image segmentation and scene understanding.

Feature Extraction: Edges contain important visual cues and characteristics that can be used as features for various machine learning tasks. Extracting features from edge locations can help characterize the content of an image and facilitate tasks like classification, clustering, and retrieval.

Image Enhancement: Edge detection and localization can also be used for image enhancement purposes, such as sharpening edges or reducing noise. Understanding the distribution of edges can guide the application of image enhancement techniques to improve overall image quality.

The specific details and insights provided by the "edge loc" attribute would depend on how it is calculated or derived within the context of the WM811K dataset. It could involve techniques such as edge detection algorithms (e.g., Sobel, Canny) or more sophisticated methods tailored to the characteristics of the dataset and the objectives of the analysis.

In summary, the "edge loc" attribute likely contains information about the spatial distribution or position of edges within the images in the WM811K dataset. Understanding this information can be valuable for a wide range of image analysis tasks, from object localization and segmentation to feature extraction and image enhancement.

3.1.2.4 Edge Ring:

The "edge ring" attribute within the WM811K dataset signifies a distinctive feature or pattern pertaining to edge detection in image processing or computer vision tasks. In this context, an edge typically denotes a significant change in intensity or color, often indicative of object boundaries or structural details within an image. The term "ring" suggests a circular arrangement of these edges,

which could manifest in various ways depending on the dataset's domain and characteristics. For instance, in datasets containing geological imagery, an edge ring might represent circular geological formations or specific mineral deposits with discernible edge characteristics. In the realm of image analysis, edge detection algorithms such as the Sobel operator or Canny edge detector are commonly employed to identify and highlight these edge features. The presence of an "edge ring" attribute in the WM811K dataset implies that such features have been identified, annotated, or extracted for further analysis, potentially aiding in tasks such as object recognition, segmentation, or feature extraction. Understanding the nature and distribution of edge rings within the dataset could provide valuable insights into the underlying structures or objects depicted in the images, facilitating more effective data interpretation and model training in related applications. However, without direct access to the WM811K dataset or specific documentation, a comprehensive interpretation of the "edge ring" attribute would require further context or domain-specific knowledge to elucidate its precise meaning and significance within the dataset's context.



(d) Edge-Ring

Fig-3.1.2.4: Edge Ring Defect

3.1.2.5 Local

In the WM811K dataset, the term "local" likely denotes a focus on the spatial context and characteristics of individual data points, particularly within the realm of image processing and computer vision. This emphasis on the "local" aspects entails analyzing features and patterns within limited spatial neighborhoods around specific pixels or regions in the images, as opposed to considering the entire image globally. Local analysis enables the detection of fine-grained details, textures, and structures that may not be apparent when examining the image at a broader scale. This could encompass techniques such as local feature extraction, where descriptors like Scale-Invariant Feature Transform (SIFT) keypoints or Local Binary Patterns (LBP) capture information about the local image structures, aiding in tasks like object recognition and image retrieval. Moreover, understanding the local context of pixels or regions facilitates tasks like semantic segmentation, where assigning semantic labels relies on considering the spatial relationships between neighboring elements. Techniques such as local filtering and enhancement can further refine the data by enhancing

relevant features or suppressing noise within localized regions, contributing to improved analysis and interpretation of the images within the WM811K dataset. Overall, the "local" attribute in the WM811K dataset underscores the importance of spatial information at a fine scale, offering insights into the intricate details and structures depicted in the images, and enabling more effective image analysis for diverse applications in the fields of image processing and computer vision.



Fig-3.1.2.5: local Defect

3.1.2.6 Scratch

In the WM811K dataset, "scratch" likely denotes an unwanted artifact or imperfection present in the images, posing challenges for analysis and machine learning tasks. Scratches can originate from various sources such as handling, transportation, or manufacturing processes, potentially degrading the quality and usability of the dataset. Addressing scratches involves several strategies, including detection, annotation, and preprocessing. Researchers may employ manual or automated methods to identify and annotate scratch regions, providing ground truth data for algorithm development and evaluation. Preprocessing techniques such as denoising and inpainting can be utilized to remove or fill in scratch areas, restoring affected image regions. Furthermore, data augmentation methods may generate scratch-free variations, augmenting training data and enhancing model robustness. Sophisticated image processing algorithms may also be applied for scratch removal and restoration, ensuring dataset integrity for subsequent analyses. Scratches not only serve as quality assessment metrics but also prompt quality assurance measures to meet standards for specific applications in industries such as manufacturing or healthcare. By effectively addressing scratches, the WM811K dataset becomes more reliable and suitable for tasks in image processing, computer vision, and beyond, fostering advancements in research, development, and real-world applications.



(f) Scratch

Fig-3.1.2.6: Scratch Defect

3.1.2.7 Random

The WM811K dataset, the term "random" could pertain to various aspects, including data selection, sampling procedures, or intrinsic characteristics of the dataset itself. Randomness often plays a crucial role in data analysis, helping to ensure representativeness, reduce bias, and support statistical inference. In the context of the WM811K dataset, which likely contains image data for tasks like object recognition or image classification, randomness might manifest in several ways:

Random Sampling: When constructing the dataset, random sampling techniques may have been employed to select images from a larger pool of available data. Random sampling helps ensure that the dataset represents the underlying population of interest without bias, thereby improving the generalizability of any conclusions drawn from its analysis. By randomly selecting images from diverse sources or scenarios, the WM811K dataset can capture a wide range of variability and realworld conditions, enhancing its utility for training and evaluating machine learning models.

Random Initialization: In the context of machine learning algorithms, randomness often comes into play during the initialization of model parameters or the shuffling of training data. For instance, neural networks commonly employ random initialization of weights to prevent the model from getting stuck in suboptimal solutions during training. Similarly, random shuffling of training samples helps prevent the model from memorizing the order of the data and improves its ability to generalize to unseen examples. In the case of the WM811K dataset, machine learning models trained on this data may leverage random initialization and shuffling techniques to improve their performance and robustness.

Random Noise: Random noise can also be a characteristic of image data, arising from factors such as sensor imperfections, environmental conditions, or variations in illumination. While noise is often considered undesirable in image analysis, it can be a crucial aspect of realistic datasets like WM811K, reflecting the inherent variability and complexity of real-world imagery. Techniques for handling random noise, such as denoising filters or robust feature extraction methods, may be applied during data preprocessing to enhance the quality and interpretability of the images in the dataset.

Random Augmentation: Data augmentation techniques, such as random rotations, translations, or flips, are commonly used to increase the diversity and size of image datasets. By applying random transformations to the images in the WM811K dataset, researchers can generate additional training examples with variations in viewpoint, scale, or orientation, thereby improving the model's ability to generalize across different conditions and viewpoints.



(g) Random

Fig-3.1.2.7: Random Defect

3.1.2.8 Near-Full

In the WM811K dataset, "near-full" likely denotes a specific condition or attribute pertaining to the capacity or utilization of data storage within the dataset. This term suggests that the dataset is almost at full capacity or saturation, possibly implying that the available storage space is nearly exhausted or that the dataset contains a vast amount of information nearing its maximum limit. The designation of "near-full" could have implications for data management, access, and processing, as dealing with large datasets approaching full capacity requires careful consideration of storage resources, computational requirements, and optimization strategies. Furthermore, the near-full status of the WM811K dataset may influence data acquisition and curation efforts, prompting decisions regarding the prioritization of data collection, storage efficiency, and the necessity of archival or compression techniques to manage and preserve the dataset effectively. Additionally, researchers and practitioners working with the WM811K dataset must be mindful of potential limitations imposed by its near-full status, such as reduced scalability, increased computational overhead, and constraints on further data expansion or updates. Addressing these challenges may involve strategies for data reduction, compression, or distributed processing to mitigate the impact of dataset size and storage constraints while ensuring continued accessibility and usability for analysis tasks in image processing, machine learning, and related domains. Overall, the "near-full" designation in the WM811K dataset underscores the importance of efficient data management practices and the need for scalable,

adaptable solutions to accommodate the growing volume and complexity of large-scale datasets in modern research and application contexts.



(h) Near-full

Fig-3.1.2.8: Near-Full

CHAPTER-4

4.1 CNN Architecture

A Convolutional Neural Network (CNN) architecture typically comprises several key layers: convolutional layers, pooling layers, and fully connected layers. Convolutional layers use filters to perform convolution operations on the input image, detecting local patterns such as edges and textures, which are crucial for identifying features. Pooling layers, often following convolutional layers, reduce the spatial dimensions of the feature maps through operations like max pooling or average pooling, thereby decreasing computational complexity and helping to prevent overfitting. The fully connected layers, usually at the end of the network, take the high-level features extracted by the convolutional and pooling layers and perform classification tasks by mapping them to output categories, such as different types of defects. This hierarchical structure allows CNNs to effectively learn and recognize complex patterns in image data, making them highly suitable for tasks like wafer map defect detection in semiconductor manufacturing



Fig-4.1: CNN Architecture

4.1.2 Types of Layers in CNN Architecture

The CNN Architecture are of Eight different types:

4.1.2.1 Kernel Filter

The Kernel Filter in Convolutional Neural Network (CNN) architecture is a fundamental component responsible for feature extraction and hierarchical learning. Comprising small matrices, these filters slide over input data through convolution operations, detecting specific patterns or features such as edges, textures, or shapes. Each filter specializes in recognizing particular patterns, and the network learns to adjust their values during training to optimize performance. Through convolution, the network extracts increasingly complex features, enabling it to learn hierarchical representations of the input data. A key advantage of kernel filters is parameter sharing, where the same set of parameters is applied across different spatial locations in the input, reducing the number of parameters and enhancing generalization. The hierarchical representations learned by kernel filters facilitate tasks such as image classification, object detection, and semantic segmentation, where understanding spatial patterns and structures is crucial. Overall, kernel filters are indispensable in CNN architecture, enabling networks to learn and extract meaningful features from complex input data, thereby achieving state-of-the-art performance in various domains.

The PASCAL Visual Object Classes (VOC) dataset is a widely-used benchmark for evaluating computer vision models on tasks such as object detection, image segmentation, and classification. For validation, the dataset provides a predefined split of training and validation sets, with images and corresponding XML annotations detailing object classes and bounding boxes. Researchers download the dataset, extract the images and annotations, and load the validation set using file lists provided in the ImageSets/Main directory. They then parse the XML files to extract ground truth labels and bounding boxes, which are used to validate model predictions by comparing them against these ground truth annotations. Common evaluation metrics include mean Average Precision (mAP) at various Intersection over Union (IoU) thresholds, facilitating standardized performance comparison across different models and approaches.

4.1.2.2 Convolutional Layer

The Convolutional Layer, a cornerstone of Convolutional Neural Network (CNN) architecture, revolutionizes how neural networks process visual data. In essence, it transforms raw pixel inputs into meaningful features through a series of convolution operations. Each layer comprises a set of learnable filters, also known as kernels, which slide over the input data, performing elementwise multiplications and summing the results to produce feature maps. These filters serve as feature detectors, each specializing in recognizing specific patterns or structures, such as edges, textures, or shapes. Through the convolution operation, the network learns to extract increasingly complex features from the input data, building a hierarchy of representations. A key advantage of convolutional layers is parameter sharing, where the same set of filter weights is applied across different spatial locations in the input. This sharing reduces the number of parameters in the network, making it more efficient and aiding generalization to unseen data. Additionally, non-linear activation functions, such as ReLU (Rectified Linear Unit), introduce non-linearity into the network, enabling it to learn complex relationships and improve its representational power. Following the convolution operation, pooling layers are often employed to downsample the feature maps, summarizing information within local neighborhoods and reducing spatial dimensions. Pooling helps make the network more computationally efficient and invariant to small spatial translations in the input. Stride and padding parameters influence the spatial dimensions of the output feature maps and help control the amount of information preserved during convolution. During training, the parameters (filter weights) within convolutional layers are optimized through backpropagation and optimization algorithms like gradient descent, minimizing a predefined loss function. This process allows the network to learn to extract relevant features from the input data and make accurate predictions on new, unseen samples. In summary, the convolutional layer is a crucial component of CNN architecture, enabling the network to efficiently process visual data, extract meaningful features, and learn hierarchical representations essential for tasks such as image classification, object detection, and semantic segmentation

4.1.2.3 Pooling Layer

The Pooling Layer, a vital component within Convolutional Neural Network (CNN) architecture, serves multiple critical functions aimed at enhancing feature representation, computational efficiency, and robustness to spatial transformations. By downsampling feature maps generated by the preceding convolutional layers, pooling reduces spatial dimensions while retaining the most salient information within local regions. This summarization process not only facilitates

computational efficiency by reducing the number of parameters and computations in subsequent layers but also promotes translation invariance, making the network less sensitive to small shifts or translations in the input data. Common pooling operations, such as Max Pooling and Average Pooling, offer different strategies for summarizing information, with Max Pooling retaining the maximum value within each local region and Average Pooling computing the average value. Despite reducing spatial dimensions, pooling aims to preserve the spatial hierarchy of features learned by convolutional layers, ensuring that essential features at different scales are maintained. Notably, pooling layers do not contain trainable parameters; instead, they operate on the feature maps generated by convolutional layers, making them computationally efficient. Consequently, pooling layers find widespread application in various CNN-based tasks, including image classification, object detection, and semantic segmentation, where they play a pivotal role in enhancing feature representation, improving computational efficiency, and promoting robustness to spatial transformations

4.1.2.4 Activation Layer

The Activation Layer is a critical component in Convolutional Neural Network (CNN) architecture, introducing non-linearity into the network and enabling it to learn complex relationships and representations from the input data. Typically inserted after convolutional and fully connected layers, activation functions transform the input data through a mathematical operation applied element-wise to each neuron's output. One of the most commonly used activation functions is the Rectified Linear Unit (ReLU), which sets all negative values to zero while leaving positive values unchanged. ReLU has gained popularity due to its simplicity and effectiveness in promoting sparse and efficient representations, accelerating convergence during training, and mitigating the vanishing gradient problem. Other activation functions include Sigmoid and Hyperbolic Tangent (Tanh), which squash the input values into a specific range, making them suitable for tasks requiring bounded outputs such as binary classification. However, these functions may suffer from saturation and vanishing gradient issues, particularly in deep networks. Leaky ReLU and Parametric ReLU (PReLU) variants address the drawbacks of traditional ReLU by allowing a small gradient for negative input values or introducing learnable parameters, respectively. Beyond promoting non-linearity, activation functions play a crucial role in shaping the decision boundaries of the network, influencing its capacity to model complex data distributions and generalize to unseen samples. In addition to ReLU and its variants, advanced activation functions like Exponential Linear Unit (ELU) and Swish have been proposed to further enhance learning dynamics and performance. Overall, the Activation Layer serves as a fundamental building block in CNN architecture, enabling networks to capture intricate patterns, learn rich representations, and achieve state-of-the-art performance across a wide range of tasks, including image classification, object detection, and natural language processing.

4.1.2.5 Rectified Linear Unit

The Rectified Linear Unit (ReLU) activation function is a foundational element in Convolutional Neural Network (CNN) architecture, serving as a crucial component in promoting nonlinearity and enhancing the network's ability to learn complex representations from input data. ReLU introduces a simple yet powerful non-linear transformation, where the function outputs zero for negative input values and leaves positive values unchanged. This piecewise linear nature accelerates convergence during training by facilitating faster gradient propagation and mitigating the vanishing gradient problem, which often hinders deep networks' performance. Moreover, ReLU promotes sparsity in activations, leading to more efficient computations and memory usage by eliminating unnecessary neuron activations. Despite its simplicity, ReLU has demonstrated remarkable effectiveness in improving model performance across various domains, including image classification, object detection, and speech recognition. However, ReLU is not without limitations, as it suffers from the "dying ReLU" problem, where neurons can become inactive during training and fail to recover, leading to dead pathways and degraded model performance. To address this issue, researchers have proposed variants such as Leaky ReLU, Parametric ReLU (PReLU), and Exponential Linear Unit (ELU), which offer improved robustness and learning dynamics by introducing small gradients for negative inputs or incorporating learnable parameters. Overall, ReLU's simplicity, efficiency, and effectiveness make it a cornerstone of CNN architecture, enabling networks to learn rich representations and achieve state-of-the-art performance in a wide range of tasks.

4.1.2.6 Flatten Layer

The Flatten Layer serves a pivotal role in Convolutional Neural Network (CNN) architecture by reshaping the multidimensional output of convolutional and pooling layers into a one-dimensional vector, thereby preparing the data for input into fully connected layers. This transformation is crucial for enabling the network to perform tasks such as classification or regression, where the output is a single vector of probabilities or values. By flattening the feature maps, the Flatten Layer retains the spatial information learned by the convolutional layers while converting it into a format compatible with traditional neural network architectures. Without this flattening step, fully connected layers would not be able to process the multidimensional feature maps efficiently. Additionally, the Flatten Layer helps reduce the computational complexity of subsequent layers by converting the highdimensional feature maps into a more manageable one-dimensional representation. Overall, the Flatten Layer plays a vital role in facilitating the transition from convolutional feature extraction to fully connected classification or regression, enabling CNNs to effectively learn and generalize from complex visual data.

4.1.2.7 Fully Connected Layer

The Fully Connected Layer, also known as the dense layer, represents the final stage in Convolutional Neural Network (CNN) architecture, where neurons in each layer are fully connected to neurons in the preceding and succeeding layers. This layer is responsible for learning complex nonlinear relationships between high-level features extracted by earlier convolutional and pooling layers, enabling the network to make predictions or classifications based on the learned representations. Each neuron in a fully connected layer receives input from every neuron in the previous layer, with weights associated with each connection that are learned during training through backpropagation and optimization algorithms. Additionally, biases are often added to each neuron to introduce flexibility and enable the network to model more diverse functions. The output of the fully connected layer is typically passed through an activation function, such as ReLU, to introduce non-linearity and enhance the network's ability to capture complex patterns in the data. The number of neurons in the fully connected layer and the architecture of the network overall are determined by the specific task at hand, with larger networks capable of learning more intricate representations but also requiring more computational resources and data for training. Overall, the fully connected layer serves as a crucial component in CNN architecture, allowing the network to transform extracted features into meaningful predictions or classifications, making it well-suited for a wide range of tasks, including image recognition, object detection, and natural language processing.

4.1.2.8 Output Layer

The Output Layer in Convolutional Neural Network (CNN) architecture serves as the final stage where predictions or classifications are made based on the learned representations from the preceding layers. Its design and configuration depend on the specific task the network is trained for, such as image classification, object detection, or semantic segmentation. For tasks involving classification, the output layer typically consists of a set of neurons, each corresponding to a class label, with the activation values representing the network's confidence or probability scores for each class. These activation values are often passed through a softmax function to convert them into a

probability distribution, ensuring that they sum up to one. The class label with the highest probability is then considered the predicted class. In contrast, for tasks involving regression, such as object localization or bounding box regression, the output layer might consist of neurons representing the coordinates or dimensions of the target objects. The network learns to predict these values based on the input data and the learned features extracted from earlier layers. Overall, the output layer plays a crucial role in CNN architecture, translating the network's learned representations into actionable predictions or classifications, thereby enabling it to solve a wide range of real-world tasks effectively.

CHAPTER-5

5.1 VGG19

VGG19 is a deep convolutional neural network architecture that was proposed by the Visual Graphics Group (VGG) at the University of Oxford. It is an extension of the VGG16 architecture, both of which were introduced in the paper "Very Deep Convolutional Networks for Large-Scale Image Recognition" by Karen Simonyan and Andrew Zisserman in 2014.Here are the key characteristics and components of the VGG19 architecture:



Fig-5.1: VGG19

Architecture: VGG19 is a deep neural network consisting of 19 layers, including 16 convolutional layers and 3 fully connected layers. The "19" in VGG19 refers to the total number of layers.

Convolutional Layers: The convolutional layers in VGG19 use small 3x3 filters with a stride of 1, and they are followed by rectified linear unit (ReLU) activation functions. These convolutional layers are organized in a sequential manner, with max-pooling layers interspersed to reduce spatial dimensions.

Max-Pooling: VGG19 employs max-pooling layers with 2x2 filters and a stride of 2 after certain convolutional blocks. Max-pooling is used to downsample the feature maps and extract dominant features.

Fully Connected Layers: After the convolutional layers, VGG19 has three fully connected layers with 4096 neurons each, followed by a final output layer with the number of neurons equal to the number of classes in the classification task.

Activation Function: ReLU (Rectified Linear Unit) is used as the activation function throughout the network, except for the output layer where softmax is commonly used for multi-class classification tasks.

Pre-Trained Models: VGG19, like VGG16, is often used as a pre-trained model for various computer vision tasks. Pre-trained models trained on large datasets like ImageNet can be fine-tuned or used as feature extractors for transfer learning tasks.

Performance: VGG19 achieved competitive performance on image classification benchmarks such as ImageNet, demonstrating the effectiveness of deep convolutional architectures for large-scale visual recognition tasks.

While VGG19 is a computationally expensive model due to its depth and large number of parameters, it has been influential in the development of deep learning architectures and serves as a benchmark for evaluating the performance of newer models.

5.2 ResNet

ResNet, short for Residual Networks, is a type of deep neural network architecture that revolutionized image classification tasks. It was introduced by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun in their 2015 paper titled "Deep Residual Learning for Image Recognition. "ResNet is known for its deep structure, typically consisting of many layers (e.g., 50, 101, 152 layers). The key innovation in ResNet is the introduction of residual connections, also known as skip connections, which help address the vanishing gradient problem in very deep networks. Here are some key points about



①: Element-wise addition

Fig-5.2: ResNet

ResNet: Residual Blocks: The building blocks of ResNet are residual blocks. Each block contains two or more convolutional layers, and a "skip connection" that adds the original input to the output of the convolutional layers. This allows the network to learn residual mappings, making it easier to train very deep networks.

Identity Mapping: The skip connection in ResNet is designed such that if the input and output dimensions are the same, the identity mapping is learned. This means that the network can choose to bypass the convolutional layers if it determines that the identity mapping is the best transformation.

Architecture: ResNet architectures are typically named according to the number of layers they have, such as ResNet-50, ResNet-101, or ResNet-152. These architectures differ in the number of residual blocks and other design choices.

Performance: ResNet achieved state-of-the-art performance on various image recognition tasks, including the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), where it significantly surpassed previous methods.

Applications: Apart from image classification, ResNet and its variations have been widely used in tasks such as object detection, image segmentation, and even in domains beyond computer vision, like natural language processing. ResNet's success has influenced the development of many other deep learning architectures and techniques, making it a fundamental milestone in the field of deep learning.

5.1.2.1 DenseNet

DenseNet, short for Dense Convolutional Network, is another influential deep learning architecture for image classification, introduced by Gao Huang, Zhuang Liu, and Kilian Q. Weinberger in their 2017 paper "Densely Connected Convolutional Networks." DenseNet is designed to address some limitations of traditional deep neural networks like vanishing gradients and feature reuse. Here are key points about DenseNet:



©: Channel-wise concatenation

Fig-5.3: DenseNet

Dense Blocks: The core idea of DenseNet is the dense block structure. Unlike traditional architectures where each layer only connects to the subsequent layer, in DenseNet, each layer is connected to every other layer in a feed-forward fashion. This dense connectivity promotes feature reuse, which helps in improving gradient flow and learning representations more efficiently.

Dense Connectivity: DenseNet uses dense connectivity through concatenation. Each layer's output is concatenated with the outputs of all preceding layers and then fed as input to the subsequent layers. This dense connectivity allows information and gradients to flow more directly throughout the network.

Growth Rate: DenseNet introduces the concept of a growth rate, which determines how many new feature maps each layer contributes to the next layers. By controlling the growth rate, DenseNet can manage the model's complexity and balance between expressiveness and computational efficiency. Transition Layers: To control the number of parameters and computational cost, DenseNet employs transition layers between dense blocks. These transition layers consist of batch normalization, 1x1 convolution, and average pooling operations to reduce the spatial dimensions of feature maps.

Advantages: DenseNet has several advantages, including improved gradient flow, better feature reuse, reduced vanishing gradient problem, parameter efficiency due to dense connections, and strong performance on image classification benchmarks like ImageNet.

Variants: DenseNet has variants such as DenseNet-121, DenseNet-169, DenseNet-201, and DenseNet-264, which vary in depth and number of layers. These variants allow flexibility in choosing a DenseNet architecture based on computational resources and task requirements.

DenseNet has been widely adopted in computer vision tasks, especially in scenarios with limited data or computational resources, due to its efficient use of parameters and strong performance.

CHAPTER-6

6.1 Performance Metrics

In Convolutional Neural Networks (CNNs), several performance metrics are commonly used to evaluate their effectiveness in various tasks, especially in tasks like image classification, object detection, and segmentation. Here are some of the key performance metrics used in CNNs:

6.1.1 Precision

Precision in the context of Convolutional Neural Networks (CNNs) refers to the ability of the model to accurately predict positive instances among all instances that it predicted as positive. It's essentially a measure of the model's exactness.

Here's a more detailed explanation of precision in the context of CNNs:

Imagine you have a CNN model trained to detect cats in images. Precision would tell you the percentage of images that the model correctly identified as containing cats among all the images it classified as containing cats.

Precision is calculated using the following formula:

$$\mathbf{Precision} = \frac{TP}{TP + FP}$$

Where:

True Positives (TP) are the number of instances correctly classified as positive (in our example, images correctly classified as containing cats).

False Positives (FP) are the number of instances incorrectly classified as positive (images classified as containing cats when they actually don't).

So, precision essentially measures the proportion of relevant instances (true positives) among all the instances predicted as positive (true positives + false positives). A high precision indicates that the model is correctly identifying positive instances while minimizing false positives.

In practical terms, if your CNN has a precision of 0.80 for cat detection, it means that when it predicts an image as containing a cat, it is correct about 80% of the time.

Precision is often used in conjunction with other performance metrics like recall, F1 score, and accuracy to provide a comprehensive evaluation of a CNN model's performance in classification tasks.

6.1.2 Recall

Recall, also known as sensitivity or true positive rate, is another important performance metric used in Convolutional Neural Networks (CNNs). Recall measures the model's ability to correctly identify all relevant instances, or in other words, the proportion of true positive instances that were correctly identified by the model out of all actual positive instances.

In the context of CNNs, recall can be understood as follows:

Imagine you have a CNN model trained for detecting cats in images. Recall would tell you the percentage of images containing cats that the model correctly identified as such among all the images that actually contain cats.

Mathematically, recall is calculated using the following formula:

$$\mathbf{Recall} = \frac{TP}{TP + FN}$$

Where:

True Positives (TP) are the instances correctly classified as positive (images correctly classified as containing cats).

False Negatives (FN) are the instances incorrectly classified as negative (images containing cats but classified as not containing cats).

So, recall measures the proportion of relevant instances (true positives) among all the instances that are actually positive (true positives + false negatives). A high recall indicates that the model is effectively capturing most of the positive instances in the dataset.

In practical terms, if your CNN has a recall of 0.85 for cat detection, it means that it correctly identifies about 85% of the images that actually contain cats.

Recall is crucial, especially in tasks where missing positive instances is costly or problematic, such as medical diagnosis or object detection. It is often used in combination with precision, F1 score, and accuracy to comprehensively evaluate the performance of a CNN model.

6.1.3 F1 Score

The F1 score is a popular performance metric used in Convolutional Neural Networks (CNNs) and other machine learning models, especially in tasks like image classification. It combines both precision and recall into a single metric, providing a balance between them.

In the context of CNNs, the F1 score is calculated using the following formula:

 $F1=2\times \frac{\textit{Precision}\times\textit{Recall}}{\textit{Precision}+\textit{Recall}}$

Where:

Precision is the proportion of true positive instances among all instances predicted as positive.

Recall is the proportion of true positive instances among all actual positive instances.

The F1 score ranges from 0 to 1, with 1 being the best possible score. It reaches its best value when precision and recall are both 1 (perfect precision and recall) and its worst value when either precision or recall is 0.

The F1 score is particularly useful in situations where there is an imbalance between the number of positive and negative instances in the dataset. For example, in medical diagnosis where the number of diseased patients may be much smaller than the number of healthy patients, F1 score provides a balanced measure of the model's performance.

In the context of CNNs, if your model has a high F1 score, it means that it achieves both high precision and high recall, indicating that it is effectively identifying relevant instances while minimizing false positives and false negatives.

Overall, the F1 score is a valuable metric for evaluating the overall performance of CNN models, especially in classification tasks where both precision and recall are important.

6.1.4 Accuracy

Accuracy is one of the fundamental performance metrics used in Convolutional Neural Networks (CNNs) and other machine learning models. It measures the overall correctness of the model's predictions across all classes.

In the context of CNNs, accuracy is calculated as the ratio of correctly predicted instances (both true positives and true negatives) to the total number of instances:

$Accuracy = \frac{No.of \ Correct \ Predicts}{Total \ No.of \ Predicts}$

Here's a breakdown of the terms:

Number of Correct Predictions: This includes both true positives (instances correctly classified as positive) and true negatives (instances correctly classified as negative).

Total Number of Predictions: This is the sum of all instances, whether correctly classified or misclassified.

Accuracy provides a general measure of how well the model performs across all classes. However, it might not be the best metric for evaluating performance in cases where classes are imbalanced. For example, if one class dominates the dataset, a high accuracy might be achieved simply by predicting the dominant class most of the time, while ignoring the minority classes.

In applications where class imbalance is an issue, other metrics like precision, recall, F1 score, or area under the ROC curve (AUC-ROC) may provide a more nuanced evaluation of the model's performance.

Overall, accuracy is a useful metric to gauge the overall performance of a CNN model, especially when all classes are equally important and well-balanced in the dataset.

CHAPTER-7

7.1 Simulation Results:

Simulation results in wafer map defect detection using Convolutional Neural Networks (CNNs) typically involve evaluating the performance of the CNN model on a dataset of wafer maps with known defects. Here is a general outline of the process and what you might expect in the simulation results:

| waferMap: | [45×48 uint8] |
|-----------------|---------------|
| dieSize: | 1683 |
| lotName: | 'lot1' |
| waferIndex: | 1 |
| trainTestLabel: | 'Training' |
| failureType: | 'none' |

7.1.1 No. of Images According to Defect Patterns in Dataset:

| Defect Patterns | No. of images |
|-----------------|---------------|
| Center | 25764 |
| Donut | 3330 |
| Edge-Loc | 31134 |
| Edge-Ring | 58080 |
| Local | 21558 |
| Near-Full | 894 |
| Random | 5196 |
| Scratch | 7158 |
| None | 147431 |

Table: 7.1.1.1



Fig.7.1.2.1



7.1.3 Test Data Confusion Matrix

Fig. 7.1.3.1

7.1.4 Performance Metric Table:

Table: 7.1.4.1

9×4 table

| 98 |
|-----|
| 32 |
| .24 |
| 95 |
| 117 |
| 518 |
| 34 |
| 63 |
| 16 |
| |

ISO 19264, formally titled "ISO 19264-1:2017: Photography – Archiving systems – Image quality analysis – Part 1: Reflective originals," is an international standard that specifies methods for evaluating the image quality of digitization systems used for cultural heritage and archival materials. The standard is particularly focused on ensuring that the digital representations of reflective originals, such as documents, photographs, and artworks, are of high quality and suitable for long-term preservation and access.

7.1.5 Training and Validation Loss:



Fig.7.1.5.1

7.1.6 Correct and Mis-classification of Image:



Misclassification (Edge-Loc



CHAPTER -8

8.1 Conclusion:

In conclusion, our study demonstrates the efficacy of Convolutional Neural Networks (CNNs) in wafer map defect detection, showcasing their remarkable performance in accurately identifying defects with high precision and recall. Leveraging CNNs' ability to automatically learn intricate patterns from raw data, we achieved notable results in detecting defects on wafer maps, as evidenced by the robust evaluation metrics obtained. Insights gleaned from the study shed light on the nuanced characteristics of wafer map defects and the intricate features learned by the CNN model, providing valuable knowledge for semiconductor manufacturing and quality control. While acknowledging certain limitations, such as dataset size and class imbalance, our findings pave the way for future research endeavors aimed at enhancing CNN architectures, addressing specific challenges in defect detection systems for semiconductor manufacturing

8.1.2 Future Work

Future work for wafer map defect detection using Convolutional Neural Networks (CNNs) could focus on several areas to further improve the accuracy, efficiency, and applicability of defect detection systems. Here are some potential directions for future research:

Large-Scale Dataset Collection: Expand the dataset size by collecting a larger and more diverse set of wafer maps with various types of defects. This would help in training CNN models on a wider range of defect patterns, leading to better generalization and robustness.

Class Imbalance Handling: Develop strategies to handle class imbalance issues inherent in wafer map datasets, where the number of defect instances is often much smaller than non-defect instances. Techniques such as oversampling, undersampling, or generating synthetic data could be explored to address this challenge.

Transfer Learning and Fine-tuning: Investigate the effectiveness of transfer learning techniques in wafer map defect detection, where pre-trained CNN models on large-scale image datasets are fine-tuned on smaller wafer map datasets. This approach could help leverage knowledge learned from other domains and adapt it to the specific task of defect detection.

Multi-class Classification: Extend the defect detection task to multi-class classification, where CNN models are trained to distinguish between different types of defects or anomalies on wafer maps. This would enable more comprehensive defect identification and characterization, leading to enhanced quality control in semiconductor manufacturing.

Anomaly Detection Techniques: Explore anomaly detection techniques in conjunction with CNNs to detect subtle or previously unseen defects that may not conform to predefined defect patterns. Unsupervised or semi-supervised learning approaches could be employed to identify anomalous regions in wafer maps without requiring explicit defect labels.

Real-time Deployment and Optimization: Develop methodologies for deploying CNN-based defect detection systems in real-time manufacturing environments, considering factors such as computational efficiency, memory footprint, and scalability. Optimization techniques tailored to embedded hardware platforms could be explored to enable efficient inference on edge devices.

Integration with Manufacturing Processes: Investigate the integration of CNN-based defect detection systems with existing semiconductor manufacturing processes, such as automated optical inspection (AOI) systems or wafer inspection tools. Seamless integration would facilitate continuous monitoring and quality control throughout the production pipeline.

By pursuing these avenues of future work, researchers can contribute to the advancement of wafer map defect detection using CNNs, ultimately leading to more reliable, efficient, and adaptable defect detection systems in semiconductor manufacturing.

8.2 REFERENCES:

[1]. M. B. Alawieh, D. Boning, and D. Z. Pan, "Wafer map defect patterns classification using deep selective learning," in Proc. 57th ACM/IEEE Design Automat. Conf. (DAC), Jul. 2020, pp. 1–6

[2]. J. C. Chien, M. T. Wu, and J. D. Lee, "Inspection and classification of semiconductor wafer surface defects using CNN deep learning networks," Appl. Sci., vol. 10, no. 15, pp. 1–13, 2020.

[3]. Y. Ji and J.-H. Lee, "Using GAN to improve CNN performance of wafer map defect type classification: Yield enhancement," in Proc. 31st Annu. SEMI Adv. Semiconductor Manuf. Conf. (ASMC), Aug. 2020, pp. 1–6.

[4]. S. Kang, "Rotation-invariant wafer map pattern classification with convolutional neural networks," IEEE Access, vol. 8, pp. 170650–170658, 2020.

[5]. U. Batool, M. I. Shapiai, N. Ismail, H. Fauzi, and S. Salleh, "Oversampling based on data augmentation in convolutional neural network for silicon wafer defect classification," in Frontiers in Artificial Intelligence and Applications, vol. 327. Amsterdam, The Netherlands: IOS Press, 2020, pp. 3–12. [Online]. Available: <u>https://ebooks.iospress.nl/volumearticle/55467</u>

[6]. Tongwha Kim, Kamran Behdinan "Advances in machine learning and deep learning applications towards wafer map defect recognition and classification": a review Journal of Intelligent Manufacturing 34 (8), 3215-3247, 2023

[7]. Shouhong Chen, Meiqi Liu, Xingna Hou, Ziren Zhu, Zhentao Huang, Tao Wang "Wafer map defect pattern detection method based on improved attention mechanism" Expert Systems with Applications 230, 120544, 2023