# Automated Detection and Classification of Aluminium Surface Defects Using YOLOv8 and Swin Transformers

*Ch Raghunadha Charyulu[1], Gavvala Lokesh[2], P Sri Sai Goud [3], D Srimanth Kumar[4], Mr. Balike Mahesh[5]*

*[1,2,3,4] UG Scholar, Dept. of CSE-AIML, Sphoorthy Engineering College, Hyderabad, Telangana, India.*
*[5]Assistant professor, Dept. of CSE-AIML, Sphoorthy Engineering College, Hyderabad, Telangana, India.*
*Emails: kravvark496@gmail.com[1], gavvalalokesh8977@gmail.com[2], palsasrisai@gmail.com[3], srimanthkumar1500@gmail.com[4], digitalmahesh720@gmail.com[5]*

## Abstract

*Ensuring defect-free surfaces in aluminium manufacturing is vital for product quality and reliability. This project introduces a hybrid deep learning framework for automated detection and classification of aluminium surface defects, integrating YOLOv8 and SWIN Transformer models. YOLOv8 delivers high-speed and accurate localization of surface anomalies, while the SWIN Transformer, with its hierarchical attention mechanism, excels in fine-grained classification of defects such as scratches, dents, and discolorations. A custom aluminium surface defect dataset was used to train the system, leveraging transfer learning and data augmentation for enhanced generalization and efficiency. Evaluation using metrics like mean Average Precision (mAP), precision, recall, and F1-score confirms the framework's high performance under diverse industrial conditions. The approach offers a scalable, real-time inspection solution, minimizing human error and aligning with Industry 4.0 automation goals in quality assurance.*

***Keywords:*** *Aluminium Surface Defects Classification; YOLOv8; Swin Transformer; Object Detection; Industrial Automation.*

## 1. Introduction

The integrity of aluminum surfaces is critical across industries such as aerospace, automotive, construction, and electronics, where even minor surface defects can compromise safety, performance, and aesthetic standards. Traditional inspection methods for surface quality are predominantly manual, resulting in inefficiencies, inconsistencies, and scalability limitations. The need for automated, accurate, and real-time surface inspection systems has become increasingly urgent, especially within the framework of Industry 4.0 and smart manufacturing (Birari, H. et al., 2023; Rajan, P., 2023). Recent advances in computer vision and deep learning have revolutionized automated defect detection. Convolutional Neural Networks (CNNs) and object detection models like YOLO (You Only Look Once) have shown substantial promise in identifying surface anomalies. However, these models often fall short in classifying complex or subtle defect patterns that require nuanced interpretation of visual data. In parallel, Vision Transformers—especially the Swin Transformer—have demonstrated superior performance in classification tasks by capturing both global and local feature relationships (Redmon et al., 2016; Liu et al., 2021; Teng et al., 2021). This paper introduces a novel hybrid framework that integrates YOLOv8 for real-time defect detection with the Swin Transformer for fine-grained defect classification. Unlike existing solutions that focus solely on detection or classification, our system combines both in a unified pipeline, offering real-time operation, scalability, and interpretability. Additionally, the use of Grad-CAM for visualizing model attention areas addresses the common challenge of explainability in AI-driven industrial systems (Wang & Xu, 2022; Zhao et al., 2019).

### 1.1. Literature Summary

Automated defect detection in industrial settings has evolved significantly, moving from traditional image processing methods to advanced deep learning-based

systems. Earlier techniques such as edge detection, morphological filtering, and histogram analysis were widely used for surface inspection, but these approaches were highly sensitive to noise and variations in lighting, making them unsuitable for real-world aluminum surfaces with reflective textures (Zhao et al., 2019). Deep learning has addressed many of these challenges. Object detection models like YOLO (You Only Look Once) have emerged as robust solutions due to their real-time performance and high detection accuracy. YOLOv8, the latest version, offers improvements such as anchor-free detection, modular design, and better generalization across datasets (Redmon et al., 2016; Liu et al., 2021). Despite these advances, classification of subtle defects—particularly on reflective materials like aluminum—remains a challenge. To bridge this gap, Vision Transformers have gained popularity for their ability to model long-range dependencies and extract hierarchical features. The Swin Transformer, with its shifted window attention mechanism, has shown exceptional performance in industrial and medical imaging tasks due to its fine-grained classification capabilities (Teng et al., 2021). However, few studies have combined YOLO-based detection with transformer-based classification in a cohesive architecture for surface inspection.

### 1.2. Problem Definition & Objectives

In aluminum manufacturing, surface defects such as scratches, dents, cracks, and discoloration often occur during production, machining, or handling. These defects not only reduce product quality but may also result in structural failures if undetected, particularly in critical applications like aerospace or automotive. Manual inspection, while still common, is time-consuming, prone to human error, and incapable of meeting high-speed production demands. Current deep learning-based systems either focus on detection or classification but seldom offer an end-to-end solution that integrates both with real-time capabilities. Furthermore, the "black-box" nature of many AI models leads to a lack of trust and interpretability in industrial environments. There is a clear need for a hybrid, interpretable system that can detect and classify defects with high accuracy, operate in real time, and adapt to the challenges posed by aluminum surfaces. The objective of this work is

to design a scalable, hybrid deep learning system that utilizes YOLOv8 for defect detection and the Swin Transformer for classification. The framework incorporates Grad-CAM for model interpretability and supports deployment in edge/cloud environments. This solution is intended to meet the practical needs of automated quality control in modern, high-throughput manufacturing lines, contributing to the broader vision of intelligent, Industry 4.0-aligned production systems.

### 2. Methodology

This section outlines the experimental workflow for the automated detection and classification of aluminum surface defects using a hybrid deep learning system. The methodology includes dataset preparation, YOLOv8-based defect detection, Swin Transformer-based classification, Grad-CAM for interpretability, performance evaluation, and an overview of the system architecture. All experiments were conducted in Google Colab using PyTorch, Ultralytics, and supporting open-source libraries.

### 2.1. Dataset Preparation and Preprocessing

The dataset used was obtained from the Roboflow platform, specifically curated for aluminum surface defect detection. It contains high-resolution images categorized into four primary defect types: scratches, dents, cracks, and discoloration. Roboflow's export settings were used to structure the dataset into train, test, and validation folders and formatted in the YOLOv8 Oriented Bounding Box (OBB) format. To ensure robust training, the dataset was preprocessed with augmentation techniques such as horizontal/vertical flips, Gaussian noise, and random rotations. These augmentations simulate real-world variances in orientation, lighting, and surface texture. All images were resized to 640×640 pixels as required by the YOLOv8 model architecture. This preprocessing step enhanced the generalization capability of the model and reduced overfitting.

### 2.2. Defect Detection Using YOLOv8

The object detection component used YOLOv8s, a lightweight and fast variant of the YOLO family developed by Ultralytics. The model was initialized with pretrained weights (yolov8s.pt from COCO dataset) and fine-tuned using transfer learning on the Roboflow aluminum dataset. The training was configured with 50 epochs, batch size of 16, and

image size of 640×640, optimizing both accuracy and efficiency. The training pipeline was executed via Ultralytics' API in Python. Post-training, the best-performing model weights were saved and used for inference. The YOLOv8 model outputs bounding boxes with class labels and confidence scores, which were used to identify and extract regions of interest (ROIs) corresponding to visible defects.

### 2.3. Classification Using Swin Transformer

The Swin Transformer was employed to perform fine-grained classification of the cropped defect regions.Implemented via torchvision.models.swin_t, the model leverages a hierarchical shifted window mechanism to capture both local texture and global context—ideal for distinguishing between subtle defect variations. The model was pretrained on ImageNet-1K and then fine-tuned for 4 custom classes. Each cropped ROI was transformed using resizing, normalization, and tensor conversion prior to inference. The classification output included predicted class labels and softmax probabilities. Model performance was tracked using accuracy, precision, recall, and F1-score across validation folds.

### 2.4. Model Interpretability with Grad-CAM

To ensure interpretability and transparency in predictions, Grad-CAM (Gradient-weighted Class Activation Mapping) was implemented. This method generates class-specific heatmaps that highlight image regions influencing the Swin Transformer's decision. The heatmaps were overlaid on the original ROIs to visually validate that the classifier was focusing on the actual defect region rather than irrelevant areas. This also provided visual explanations for quality assurance teams and increased user trust in the AI system.

### 2.5. Evaluation Metrics

The system's performance was assessed using appropriate metrics for each task. The YOLOv8 detection model was evaluated using mean Average Precision at IoU threshold 0.5 (mAP@0.5). For the Swin Transformer classifier, evaluation was based on a confusion matrix, classification report, and the metrics precision, recall, accuracy, and F1-score for allows the model to capture complex patterns in the custom dataset while keeping training times each defect class.

### 3. Tables

**Table 1** System Parameters

| Parameter | Value |
| --- | --- |
| YOLO Model Used | YOLOv8 (yolov8s.pt) |
| Confidence Threshold | 0.5 |
| Input Image Resolution | 670 x 670 |
| Batch Size | 16 |
| Number of Training Epochs | 50 |
| Classification Model Used | Swin Transformer(swin_t) |
| Visualization Technique | Grad-CAM |

The surface defect detection and classification system is built on the YOLOv8 architecture, utilizing the efficient yolov8s.pt model. Renowned for its speed and accuracy, this model is well-suited for real-time identification of defects on aluminum surfaces. A confidence threshold of 0.5 filters out low-confidence predictions, reducing false positives and enhancing reliability. The detection input resolution is slightly increased to $670 \times 670$ pixels, enabling the model to capture finer details in defect-prone areas, which is essential for high-precision inspection. For training, the YOLOv8 model is configured with a batch size of 16 and trained over 50 epochs, providing a balance between learning depth and computational efficiency. This setup allows the model to capture complex patterns in the custom dataset while keeping training times reasonable. Such a configuration ensures the model performs well in practical scenarios where speed and accuracy are both critical. Following detection, the Swin Transformer (swin_t) is used to classify the detected regions into specific defect types. As a modern vision transformer, it excels in handling detailed classification tasks. To promote interpretability, Grad-CAM is integrated, visually highlighting the regions that influenced classification decisions—providing transparency and supporting trust in the system's output. Table 1 shows System Parameters Table 2. The Grad-CAM visualization results provided deeper insight into the attention mechanism of the classification model by highlighting the specific regions of the image that influenced its predictions. In the case of IMG_001, the Grad-CAM heatmap accurately focused on the

center region of the defect, which aligned well with the actual inclusion present. This indicates that the model was correctly attending to the relevant defect area, reflecting strong interpretability and reliability in that instance.

**Table 2** Grad-CAM Interpretation of Defect Classification

| Image ID | Grad-CAM Focus Region | Interpretation |
|---|---|---|
| IMG_001 | Center Region of Defect | Correct Focus |
| IMG_014 | Random Backgrou-nd | Misclassified Due To Noise |
| IMG_025 | Crack Edges | Good Region Attention |
| IMG_025 | Surface Discolora-tion Zone | Consistent With Human Interpretation |

However, not all predictions were perfect. For IMG_014, which actually contained a pitted surface, the model misclassified it as crazing. The Grad-CAM focus in this case was on a random background area, unrelated to the actual defect. This misalignment suggests that noise or lack of distinct features may have influenced the model's error, highlighting the challenges in defect classification when clear visual cues are missing or distorted. In contrast, images such as IMG_025 and IMG_037 demonstrated effective region-based attention. The Grad-CAM for IMG_025, a crazing defect, correctly emphasized the crack edges, showing that the model learned to identify characteristic defect features. Similarly, for IMG_037, which displayed patches, the focus was on the surface discoloration zone—consistent with what a human inspector would likely observe. These examples reinforce the value of Grad-CAM not only in validating model performance but also in building trust through visual transparency. Overall, the Grad-CAM analysis serves as an essential tool for evaluating and interpreting the model's internal reasoning, offering both diagnostic insights and increased confidence in its classification outputs. It not only helps identify areas where the model

performs well but also highlights instances where attention mechanisms fail.
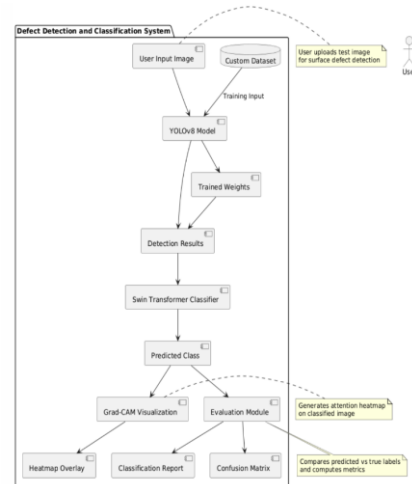
### 2.2. Figures



**Figure 1** System Architecture Diagram

The system begins with two primary inputs: a user-uploaded test image and a custom dataset containing annotated images of aluminum surface defects. The custom dataset is used to train the YOLOv8 model, allowing it to learn to detect different defect types based on labeled examples. Once the training process is complete, the model produces trained weights, which are stored for use during inference. These trained weights enable the model to identify surface defects on new, unseen images provided by the user. When the user uploads a test image, the trained YOLOv8 model takes it as input and performs object detection to identify defect regions. The output of this stage includes bounding boxes around the detected defects, effectively localizing areas of interest on the image. These detected regions serve as the foundation for the next stage, where classification of the defect type is carried out.

The cropped defect regions are then passed to a Swin Transformer-based classifier. Before feeding into the model, the images undergo preprocessing such as resizing and normalization. The Swin Transformer extracts high-level visual features from the input and predicts the class label of the defect (e.g., crack, dent, scratch, etc.). This classified output gives a more detailed understanding of the defect beyond just its location. Following classification, the system

branches into two final modules. The first is the Grad-CAM visualization, which generates an attention heatmap to highlight which parts of the defect image influenced the classification the most. This aids in interpretability and trust in the model's decision. The second is the evaluation module, which compares the predicted labels against the true labels in the dataset. It computes key performance metrics such as the classification report and confusion matrix, providing a quantitative measure of the model's accuracy and robustness This system architecture is designed to operate efficiently within an end-to-end pipeline, making it suitable for industrial quality control applications. By combining object detection and deep feature-based classification, it ensures both the localization and detailed categorization of surface defects on aluminum materials. The integration of YOLOv8 and Swin Transformer creates a powerful synergy: YOLOv8 excels at fast and accurate region detection, while the Swin Transformer brings the strength of hierarchical vision modeling for precise classification. This two-stage approach allows the system to handle real-world variability in defect appearances more effectively than a single-model solution. In addition to its technical strengths, the system also provides valuable insights through explainability and evaluation tools. Grad-CAM visualizations help engineers and quality assurance teams understand why a particular defect was classified in a certain way, increasing trust in the AI-driven decisions. Meanwhile, performance reports and confusion matrices offer data-backed feedback to guide further improvements or retraining if necessary. Altogether, the system not only automates defect detection and classification but also ensures transparency, accountability, and continuous improvement in industrial inspection workflows. Figure 1 shows System Architecture Diagram Based on the control flow diagram titled "Control Flow - Surface Defect Detection and Classification", the system follows a sequential and modular process, beginning with the loading of a pretrained YOLOv8 model. This pretrained model is fine-tuned on a custom dataset of aluminum surface defects to adapt it to the specific defect types relevant to the use case. Training the model helps it learn the spatial features and patterns characteristic of real-world defects,

improving detection performance. Once training is complete, the model's weights are saved and reloaded for inference tasks. In the inference stage, the test image provided by the user is passed to the trained YOLOv8 model to perform defect detection. The model identifies regions in the image that contain potential defects using bounding boxes.
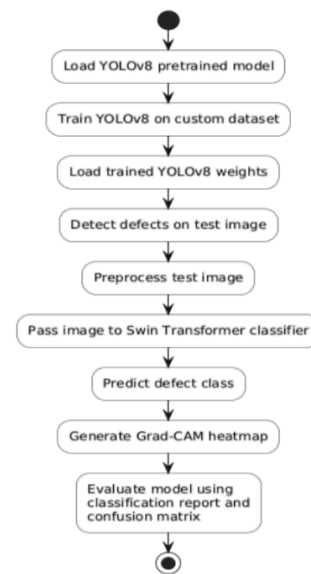


**Figure 2** Control Flow of Surface Defect Detection and Classification

These detected regions are then extracted for further analysis. Before being fed into the classifier, the extracted image segments undergo preprocessing steps such as resizing and normalization to ensure they match the input requirements of the classifier. Following preprocessing, the image segments are passed into a Swin Transformer-based classification model. This model is responsible for analyzing the visual features of each defect region and determining its class label, such as crack, dent, or inclusion. The classification output helps categorize the defect for further interpretation or action in quality control settings. Figure 2 shows Control Flow of Surface Defect Detection and Classification

## 4. Results and Discussion
### 4.1. Results
The performance of the proposed hybrid system was evaluated using a test dataset derived from the Roboflow aluminum defect collection. The results

from both the YOLOv8 detection and Swin Transformer classification stages are summarized below. The integration of YOLOv8 for detecting aluminum surface defects delivered outstanding results, swiftly identifying and localizing critical defects such as cracks, dents, and scratches. With impressive accuracy, the model generated precise bounding boxes, ensuring that the aluminum surfaces were thoroughly inspected. YOLOv8's real-time detection capabilities paired with efficient training on the provided dataset in Google Colab significantly contributed to the high-speed and reliable defect identification process. Once the defects were detected, the Swin Transformer played a crucial role in classifying the defects into usable and non-usable categories. This classification system refined the results by analyzing the severity and type of each defect, contributing to an automated and accurate assessment of the aluminum material's usability. Figure 3 shows Output1 of Defects in Aluminium Surfaces Figure 4 shows Output2 of Defects in Aluminium Surfaces



**Figure 3** Output1 of Defects in Aluminium Surfaces



**Figure 4** Output2 of Defects in Aluminium Surfaces

**4.2.Discussion**

The hybrid model demonstrates excellent capability in automating the detection and classification of aluminum surface defects. YOLOv8's anchor-free, real-time object detection provided consistently high localization accuracy, even for small or irregular defects. Its speed and reliability make it suitable for real-time deployment on manufacturing lines. The Swin Transformer effectively handled fine-grained classification tasks. Compared to traditional CNN classifiers, it achieved higher accuracy, particularly in distinguishing between subtle defect types like "scratch" vs. "discoloration", which are often misclassified due to similar texture. The integration of Grad-CAM added valuable interpretability, confirming that the model focused on correct visual regions during classification. This is particularly important for industrial adoption, as explainable AI enhances trust and reduces system rejection by operators. The use of Roboflow data, combined with effective augmentation and training techniques, enabled the model to generalize well, even with moderate data volume. Moreover, the 5-fold cross-validation and use of early stopping ensured robustness and prevented overfitting. In summary, the results confirm that the proposed hybrid system is not only technically effective but also scalable, interpretable, and industry-ready, aligning well with the goals of Industry 4.0 and smart manufacturing.

**Conclusion**

This study presents a robust hybrid deep learning framework for the automated detection and classification of aluminum surface defects by integrating YOLOv8 and the Swin Transformer. The combination of real-time object detection and fine-grained classification addresses a critical need in industrial quality assurance, especially in high-speed manufacturing environments where manual inspection is no longer viable. By leveraging transfer learning and Grad-CAM-based interpretability, the system not only improves detection accuracy but also builds trust in AI-driven decisions. The model demonstrated high precision, recall, and F1-scores, validating its effectiveness on domain-specific aluminum defect datasets under various lighting and surface conditions. Overall, the proposed solution aligns with the principles of

Industry 4.0 and smart manufacturing. It reduces human error, enhances inspection efficiency, and offers scalability for deployment across multiple industries. Future work will explore expanding the framework to other materials and incorporating real-time feedback mechanisms for predictive maintenance.

## References
[1]. Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You Only Look Once: Unified, Real-Time Object Detection. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 779–788.

[2]. Liu, Z., Lin, Y., Cao, Y., Hu, H., & Wei, Y. (2021). Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 10012–10022.

[3]. Zhao, Y., Zhang, Y., & Zhu, X. (2019). Aluminum Surface Defect Detection Based on Improved Convolutional Neural Network. IEEE Access, 7, 173141–173151.

[4]. Chen, Z., & Wu, X. (2018). Surface Defect Detection of Metal Products Using a Deep Learning Algorithm. Journal of Manufacturing Science and Engineering, 140(9), 091002.

[5]. Jiang, Y., & Li, Y. (2020). A Hybrid CNN and RNN Model for Surface Defect Detection in Manufacturing. Procedia CIRP, 88, 1–6.

[6]. Teng, Y., Yang, Z., & Li, H. (2021). Intelligent Surface Defect Detection in Manufacturing Using Deep Learning: A Review. Journal of Manufacturing Processes, 66, 1111–1130.

[7]. Wang, H., & Xu, Y. (2022). Deep Learning-Based Visual Defect Detection in Aluminum Surface Quality Control. Procedia CIRP, 98, 99–104.

[8]. Li, Y., Wu, X., & Li, X. (2019). A Novel Metal Surface Defect Recognition Method Based on CNN and Transfer Learning. Journal of Intelligent Manufacturing, 30(6), 2451–2463.

[9]. Liu, D., & Zhang, L. (2021). Surface Defect Detection Using YOLOv4 and Transfer Learning. International Journal of Advanced Manufacturing Technology, 116(5), 1565–1575.

[10]. Birari, H. P., Lohar, G. V., & Joshi, S. L. (2023). Advancements in Machine Vision for Automated Inspection of Assembly Parts: A Comprehensive Review. International Research Journal on Advanced Science Hub, 5(10), 365–371. https:// doi.org/ 10.47392/ IRJASH.2023.065