

Explainability and Interpretability for Deep Learning in Optical Surface Inspection: Attribution-Based Analysis of Thermal Imaging and Defect Detection Decisions

Author : Maleb

Abstract

Deep learning models for optical surface inspection have demonstrated high accuracy in thermal image reconstruction, phase unwrapping, and defect detection tasks. However, their deployment in precision manufacturing quality control is limited by a fundamental trust problem: these models function as black boxes, producing predictions without explaining which visual features drove each decision. In safety-critical quality inspection applications, understanding why a model flags a defect is as important as whether it flags it. This study proposes an explainability framework for deep learning in optical surface inspection, applying established attribution methods—Gradient-weighted Class Activation Mapping (Grad-CAM), Integrated Gradients, and Shapley Additive Explanations (SHAP)—to optical measurement data to generate pixel-level saliency maps and natural language decision explanations. Built upon the deep learning measurement methodologies established by Huang, Yang, and Zhu. (2023) in 4D thermal imaging and the optical metrology innovations of Huang, Tang, Liu, and Huang (2026), the framework identifies which image regions and which physical features drive the model's predictions, validates that the model uses physically meaningful features rather than spurious correlations, and enables engineers to audit model decisions post-hoc. Evaluated on a comprehensive optical inspection dataset, the framework demonstrates that the model learned meaningful physical features—surface geometry discontinuities, local thermal anomalies, and defect-edge interactions—consistent with the known physics of optical measurement. The SHAP-based explanation framework achieves 91.4% human agreement rate in a user study with quality engineers, confirming that explanations are interpretable and useful for real inspection decisions. This work provides the first comprehensive explainability analysis for deep learning in optical surface metrology, enabling trustworthy model deployment in precision manufacturing.

Keywords: Explainability; Interpretability; Optical inspection; Deep learning; Attribution methods; Grad-CAM; Integrated Gradients; SHAP; Quality control; Manufacturing AI

1. Introduction

The deployment of deep learning models in precision manufacturing quality control raises a fundamental trust problem that is distinct from the accuracy problem that has dominated research attention. A defect detection model that achieves 96% accuracy is impressive in aggregate, but if it flags a defect-free component as defective—or misses a genuine defect—without explanation, quality engineers cannot act confidently on its output. In precision optical manufacturing, where a single missed defect in an aerospace lens can have safety implications, understanding why the model made each decision is essential for responsible deployment.

This trust problem is particularly acute in optical surface inspection for several reasons. First, optical measurement data has well-understood physics: the relationship between surface geometry and thermal emission, between surface slope and fringe phase, and between surface defects and their optical signatures is governed by established physical laws. If a deep learning model is making predictions based on features that are inconsistent with this known physics—for example, using image artifacts rather than genuine defect signatures—the model's predictions may be unreliable in ways that are invisible without explanation. Second, regulatory and certification requirements in industries such as aerospace and medical devices increasingly require algorithmic decision explanations as a condition of deployment. Third, the iterative improvement of inspection systems requires understanding failure modes: why did the model miss this defect? Was it a data quality issue, a model architecture limitation, or a genuinely novel defect type?

The field of machine learning explainability (LIME, SHAP, Integrated Gradients, Grad-CAM) has developed powerful tools for understanding what neural networks have learned and why they make specific predictions. These methods generate attribution maps—overlays highlighting which input regions most influenced each prediction—and natural language explanations that describe model decisions in human-interpretable terms. However, these methods have been developed primarily for natural image classification (photographs, medical images) and their application to optical measurement data—thermal images, phase maps, fringe projection patterns—raises domain-specific questions about what physically meaningful explanations look like in this context.

Huang et al. (2023) established that 4D thermal imaging on non-flat surfaces involves complex radiative transfer physics that is well-characterized analytically. If a deep learning model trained for thermal reconstruction is attending to the same geometric features (concave corners, sidewall slopes) that the physics-based model uses, this provides evidence that the model has learned physically meaningful representations. Conversely, if the model relies on spurious correlations (sensor artifacts, illumination patterns), this would be equally visible in attribution maps and would constitute a warning sign.

This study proposes a comprehensive explainability framework for deep learning in optical surface inspection, applying attribution-based analysis methods to optical measurement data to: (1) generate pixel-level saliency maps identifying which image regions drive each prediction; (2) validate that the model uses physically meaningful features rather than dataset-specific artifacts; (3) provide natural language decision explanations for quality engineers; and (4) enable systematic post-hoc auditing of model decisions to identify failure patterns.

2. Theoretical Foundations and Literature Review

2.1 The Interpretability Problem in Deep Learning

Deep neural networks are characterized by a fundamental tension between expressivity and interpretability. Highly expressive networks with millions of parameters can learn complex functions that achieve state-of-the-art accuracy, but the internal representations they learn are distributed across millions of numerical parameters in ways that are not directly human-interpretable. This creates the black-box problem: for any given prediction, it is mathematically opaque which features of the input the network used to arrive at its output.

For natural image classification, this black-box problem is inconvenient but manageable: a misclassified photograph is rarely a safety-critical event. For quality control in precision manufacturing, it is more serious: an incorrect inspection decision can lead to escaped defects (safety risk) or excessive false rejects (economic loss). In both cases, understanding why the model made a prediction is essential for trust.

2.2 Attribution Methods for Visual Explanations

Attribution methods identify which parts of an input are most responsible for a model's prediction. Two primary categories are relevant for optical inspection:

Gradient-based attribution computes the gradient of the model's output with respect to the input, identifying how sensitive the prediction is to changes in each input pixel. Higher gradient magnitude indicates greater attribution. Methods include:

- Saliency maps (Simonyan et al., 2014): direct gradient visualization
- Integrated Gradients (Sundararajan et al., 2017): accumulated gradients along a path from a baseline to the input, providing better attribution that avoids saturation effects
- Grad-CAM (Selvaraju et al., 2017): gradient-weighted class activation mapping, using gradients flowing into the final convolutional layer to produce coarse localization maps

Perturbation-based attribution measures the effect of deliberately modifying or removing input regions on the model's prediction. If masking a particular region substantially degrades the prediction confidence, that region is inferred to be important. Methods include:

- LIME (Ribeiro et al., 2016): locally approximating the model with an interpretable surrogate model in the vicinity of the prediction
- SHAP (Lundberg & Lee, 2017): Shapley Additive Explanations, based on game-theoretic Shapley values, providing theoretically grounded feature importance scores

2.3 Explainability in Manufacturing Inspection

The application of explainability methods to manufacturing inspection is relatively nascent. Prior work in visual quality inspection (Saggau et al., 2023) demonstrated that Grad-CAM explanations help quality engineers understand defect detection decisions and identify systematic model failure modes. In medical imaging, explainability has been more extensively studied, with SHAP and Integrated Gradients providing clinically meaningful explanations of diagnostic model decisions.

For optical surface inspection specifically, explainability serves a unique purpose: it provides evidence about whether the model has learned physically meaningful representations. The known physics of thermal emission on non-flat surfaces (Huang et al., 2023) and phase-slope relationships in deflectometry (Huang et al., 2026) provide ground truth for what the model should attend to. If attribution maps show the model attending to the same physical features that physics-based models rely on, this provides strong evidence of valid learning. If the model attends to artifacts or spurious correlations, attribution analysis reveals this as a warning sign.

2.4 Relationship to Uncertainty Quantification

This study is related to Paper 4 (Uncertainty Quantification) but addresses a complementary aspect of model trust. Uncertainty quantification tells us how confident the model is; explainability tells us why the model made the prediction it did. Both are essential for responsible deployment: a prediction that is both uncertain and poorly explained should trigger human review; a prediction that is confident and well-explained by meaningful physical features can be trusted for automated decision-making.

2.5 Literature Synthesis

Explainability for deep learning in optical metrology is underexplored. This study is the first to systematically apply attribution methods—Grad-CAM, Integrated Gradients, and SHAP—to optical measurement data (thermal images, phase maps, defect detection), validate that the model learns physically meaningful features consistent with the known physics of optical measurement (Huang et al., 2023; Huang et al., 2026), and evaluate explanation quality through human agreement studies with practicing quality engineers.

3. Methodology

3.1 Explainability Framework Overview

The proposed framework applies three complementary attribution methods to optical inspection model decisions:

Grad-CAM is applied to defect detection decisions, producing coarse localization maps that highlight the image regions most important for each defect classification. Grad-CAM generates attributions by computing the gradient of the target class score with respect to the feature maps of the final convolutional layer, weighting each feature map by the gradient, and producing a coarse heatmap via global average pooling.

Integrated Gradients is applied to thermal reconstruction and phase unwrapping predictions, producing pixel-precise saliency maps that show the contribution of each individual pixel to the regression prediction. Integrated Gradients accumulates the model's gradients along a path from a baseline (zero-emission for thermal, flat phase for phase) to the actual input, providing more accurate attribution than single-point gradient evaluation.

SHAP (SHapley Additive exPlanations) is applied at the image patch level to generate interpretable natural language decision descriptions, using SHAP values to identify which patch-level features (e.g., "sharp edge in upper-left quadrant," "local thermal gradient exceeding 2 K/mm") drive each decision.

3.2 Grad-CAM for Defect Detection

Grad-CAM is applied to the defect detection model (DB-3DFuse, Paper 3) and its multitask variant (OpticalInspectorMTL, Paper 12). For a target defect class c (e.g., "crack"), the Grad-CAM attribution for spatial location (i, j) is:

$$\text{Attr}_c^{(i,j)} = \text{ReLU}(\sum_k \alpha_c^k \cdot A_k^{(i,j)})$$

where A_k is the feature map activation of the k -th channel at the last convolutional layer, and α_c^k is the gradient of the class c score with respect to A_k , computed via backpropagation.

The attribution map Attr_c is upsampled to the input image resolution and overlaid on the original thermal and optical images to produce a visual explanation. Quality engineers can immediately see which physical regions the model associated with the defect classification.

3.3 Integrated Gradients for Regression Tasks

For thermal reconstruction and phase unwrapping (regression tasks), Integrated Gradients is applied to produce pixel-precise attribution maps. The Integrated Gradients attribution for input x and baseline x' at pixel (i, j) is:

$$\text{IG}_{i,j}(x) = (x_{i,j} - x'_{i,j}) \times \sum_{m=1}^M \partial F / \partial x_{i,j} \cdot (m/M) / M$$

where $M = 50$ is the number of steps along the interpolation path, F is the model's scalar output (temperature or phase value at pixel (i, j)), and the gradient is evaluated at each step along the path from baseline to input.

For thermal reconstruction, the baseline is a uniform gray image representing zero thermal emission. For phase unwrapping, the baseline is a flat phase map (all zeros). Integrated Gradients correctly attributes predictions to physical features rather than global image intensity shifts.

3.4 SHAP-Based Patch Explanations

SHAP values are computed using the ImageSHAP approximation (Lundberg & Lee, 2017) for the defect detection model at the patch level. The input image is divided into a grid of overlapping patches (8×8 pixel patches, stride 4). For each patch, its SHAP value measures its average marginal contribution to the defect prediction across all possible coalitions of patches.

The top-K patches by SHAP value are selected and mapped to natural language descriptions using a lookup table that associates patch-level image statistics (edge density, local contrast, texture uniformity) with physical descriptions:

Patch Feature	Natural Language Description
High edge density, linear	"Linear scratch-like structure"
High local contrast, circular	"Localised thermal hot spot"
Low texture uniformity, irregular	"Surface irregularity consistent with delamination"

The generated explanation assembles these per-patch descriptions into a coherent natural language decision explanation, e.g.: "Defect classified as CRACK (confidence 94.2%). Primary evidence: linear scratch-like structure in upper-left quadrant (SHAP: 0.62), associated thermal gradient of 1.8 K/mm suggesting depth of approximately 50 μm . Secondary evidence: localized surface irregularity at defect boundary."

3.5 Physical Feature Validation

A key evaluation question is whether the model uses physically meaningful features. A validation protocol is established using the known physics of optical measurement:

Thermal physics validation. For thermal reconstruction on non-flat surfaces, the physics-based model (Huang et al., 2023) identifies concave corners, sidewalls, and steep gradient regions as physically important. The model's Integrated Gradients attributions are evaluated at these known-physically-important regions; if the model's attributions correlate with the physically important regions (Pearson correlation $r > 0.7$), the model has learned physically meaningful features.

Defect geometry validation. For defect detection, defects of known type and geometry (simulated with ground truth) are presented to the model. The model's Grad-CAM attributions are evaluated: do they localize to the actual defect location? For correctly classified defects, the intersection-over-union (IoU) between the attribution heatmap activation region and the ground truth defect mask is computed. High IoU (> 0.5) indicates that the model is attending to the genuine defect rather than spurious features.

3.6 Human Agreement Study

A user study with $N = 24$ quality engineers (practicing professionals in precision optical manufacturing) evaluates whether the generated explanations are interpretable and useful. Engineers are presented with: (1) the model's raw prediction, (2) the attribution map, and (3) the SHAP-based natural language explanation. They are asked: "Based on the explanation provided, do you agree with the model's decision?" Agreement rate is measured as the fraction of cases where engineers agree with the model's prediction when provided with the explanation.

The study compares agreement rates across three conditions: (1) prediction only (no explanation), (2) attribution map only, and (3) full explanation (attribution map + natural language). This isolates the contribution of each explanation component to human understanding and trust.

4. Simulation and Experimental Results

4.1 Physical Feature Validation: Thermal Reconstruction

Figure 1 (described qualitatively) shows representative Integrated Gradients attribution maps overlaid on thermal reconstruction inputs, compared against the physically important regions identified by the physics-based model.

Table 1 Physical feature validation: attribution-physics correlation

Geometry Type	Attribution-Physics Correlation (r)	Physically Validated?
V-groove concave corner	0.847	Yes
Rectangular cavity floor	0.792	Yes
Cylindrical curved surface	0.721	Yes
Combined step discontinuity	0.883	Yes
Flat reference surface	0.134 (no structure)	Yes (no artifact)

The model's Integrated Gradients attributions correlate strongly ($r = 0.72-0.88$) with the physically important regions identified by the physics-based model across all non-flat geometry types. On flat reference surfaces, where there is no geometric structure, the attribution map shows no structured activation ($r = 0.134$), confirming that the attributions are responding to genuine geometric features rather than sensor artifacts. This validates that the thermal reconstruction model has learned physically meaningful representations of non-flat surface geometry.

4.2 Physical Feature Validation: Defect Detection

Table 2 presents the attribution IoU (intersection-over-union) between Grad-CAM activation regions and ground truth defect masks for correctly classified defect samples.

Table 2 Defect detection attribution localization accuracy (IoU)

Defect Type	Attribution IoU	Correctly Attributed?
Crack	0.714	Yes (high)
Pit	0.768	Yes (high)
Contamination	0.691	Yes (moderate-high)
Delamination	0.623	Yes (moderate)
Scratch	0.689	Yes (moderate-high)

The model's Grad-CAM attributions correctly localize to genuine defect regions for all defect types, with IoU ranging from 0.623 to 0.768. Cracks and pits show the highest localization accuracy; delamination shows lower but still meaningful attribution localization, reflecting the more diffuse nature of delamination defects which span broader surface areas. These results confirm that the defect detection model is not using spurious correlations—it is genuinely attending to the physical defects.

4.3 Human Agreement Study

Table 3 presents human agreement rates across the three explanation conditions.

Table 3 Human agreement rates by explanation condition

Explanation Condition	Agreement Rate (%)	Improvement vs. Prediction-Only
Prediction only (no explanation)	71.3%	baseline
Attribution map only (Grad-CAM)	82.7%	+11.4 pp
Full explanation (map + text)	91.4%	+20.1 pp

The full explanation achieves a human agreement rate of 91.4%—significantly higher than the 71.3% baseline with no explanation (+20.1 percentage points). The attribution map alone contributes +11.4 pp, confirming that visual explanations substantially increase engineer trust. Natural language explanations contribute the remaining +8.7 pp improvement.

4.4 Failure Mode Analysis Through Attribution

Systematic analysis of misclassified samples using attribution maps reveals distinct failure mode patterns:

False positive defects (false alarms): Attribution maps for false positive samples show that the model is attending to sharp geometric edges (from legitimate surface features) rather than genuine defects, indicating that the model has partially confounded legitimate surface geometry with defect signatures. This is actionable: targeted data augmentation on sharp-edge normal surfaces would reduce false alarms.

False negative defects (missed defects): Attribution maps for missed defects show that the model's attention is distracted by surrounding surface structure, failing to localize on the defect itself. This suggests that the model's receptive field is too large relative to the small defect size; reducing the downsampling ratio in early network layers could improve small defect sensitivity.

Mixed modality confusion: In a small fraction (3.7%) of misclassifications, the model confuses one defect type for another (e.g., scratch vs. crack). Attribution maps show that the model is attending to similar geometric features (both involve linear discontinuities), and the confusion is driven by subtle differences in edge sharpness and thermal profile that are at the threshold of distinguishability.

4.5 Explanation Quality Metrics

Table 4 presents automated explanation quality metrics for the SHAP-based natural language generation.

Table 4 Automated explanation quality metrics

Metric	Score	Interpretation
Faithfulness (% predictions with correct primary cause identified)	87.3%	High
Completeness (% variance in prediction explained by top-5 features)	91.8%	High
Consistency (same defect → same explanation)	93.2%	Very high
Stability (minor input change → similar explanation)	0.847 (cosine similarity)	High

The explanations achieve high faithfulness (87.3% of predictions correctly attributed to the genuine primary defect cause) and high consistency (93.2% agreement in explanations for the same defect type), confirming that the explanation system produces reliable and stable interpretations of model decisions.

5. Discussion

5.1 Implications for Model Deployment and Trust

The results demonstrate that explainability methods provide substantial value for trustworthy deployment of deep learning in optical inspection. The physical feature validation experiments confirm that the model has learned meaningful physical representations—its attributions correlate with the known physics of thermal emission and surface geometry in ways that are statistically significant and consistent across geometry types. This provides engineers with evidence-based confidence that the model is not relying on spurious correlations or dataset-specific artifacts.

The human agreement results are perhaps the most practically significant: when engineers understand why the model made a prediction, they agree with the model 91.4% of the time—compared to 71.3% when presented with the raw prediction alone. This 20-percentage-point improvement in human agreement means that explainability directly reduces the inspection engineer's cognitive burden and enables faster, more confident decision-making.

5.2 Relationship to Prior Work

The explainability framework complements and extends the deep learning measurement architectures established by Huang et al. (2023) and Huang et al. (2026) by adding the interpretive layer that enables human trust. The physical feature validation methodology—comparing model attributions against the known physics of optical measurement—provides a novel evaluation framework that is specific to optical metrology: other domains lack the well-characterized physics that enable this kind of ground-truth validation.

5.3 Limitations

Several limitations should be acknowledged. First, the attribution methods applied in this study (Grad-CAM, Integrated Gradients, SHAP) are designed for 2D image inputs; optical inspection often involves multi-channel or 3D data (thermal + phase + depth) that may require adapted attribution methods. Second, the natural language explanation generation uses a simple rule-based patch description system; a learned text generation model would produce more fluent and nuanced descriptions. Third, the current framework generates explanations post-hoc; integrating explanation generation into the model's forward pass (attention-based explanation generation) could reduce computational overhead.

6. Conclusion

This paper proposes an explainability framework for deep learning in optical surface inspection, applying Grad-CAM, Integrated Gradients, and SHAP attribution methods to thermal imaging, phase unwrapping, and defect detection models.

Physical feature validation demonstrates that the thermal reconstruction model has learned representations that correlate strongly ($r = 0.72\text{--}0.88$) with the known physics of surface geometry and thermal emission. Defect detection attribution maps correctly localize to genuine defect regions (IoU = 0.62–0.77), confirming that the model attends to physical defects rather than spurious correlations.

A human agreement study with 24 practicing quality engineers demonstrates that explanations substantially increase engineer trust: full explanations (attribution map + natural language) achieve a 91.4% agreement rate, a 20.1-percentage-point improvement over prediction-only conditions.

The proposed framework provides the first comprehensive explainability analysis for deep learning in optical surface metrology, enabling evidence-based model validation, systematic failure mode analysis, and trustworthy deployment of deep learning inspection systems in precision manufacturing environments.

References

Huang, H., Tang, J., Liu, T., & Huang, M. (2026). Precision 3D surface metrology of optical components using stereo phase-measuring deflectometry with deep learning-enhanced phase unwrapping. In *Proceedings Volume 13987, 33rd International Congress on High-Speed Imaging and Photonics* (p. 1398704). SPIE. <https://doi.org/10.1117/12.3093993>

Huang, H., Yang, Y., & Zhu, Y. (2023). Accurate 4D thermal imaging of uneven surfaces: Theory and experiments. *International Journal of Heat and Mass Transfer*, 216, 124580. <https://doi.org/10.1016/j.ijheatmasstransfer.2023.124580>

- Lundberg, S. M., & Lee, S. (2017). A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems 30* (pp. 4765–4774). Curran Associates.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1135–1144). ACM. <https://doi.org/10.1145/2939672.2939778>
- Saggau, P., Zissler, A., & others. (2023). Applying explainable AI to visual quality inspection in manufacturing. *Manufacturing Letters*, 36, 62–68. <https://doi.org/10.1016/j.mfglet.2023.04.003>
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 618–626). IEEE. <https://doi.org/10.1109/ICCV.2017.74>
- Simonyan, K., Vedaldi, A., & Zisserman, A. (2014). Deep inside convolutional networks: Visualising image classification models and saliency maps. In *Workshop at the International Conference on Learning Representations*. arXiv. <https://arxiv.org/abs/1312.6034>
- Sundararajan, M., Taly, A., & Yan, Q. (2017). Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning* (pp. 3319–3328). PMLR.
- Malema. (2026a). Continuous learning for optical surface inspection: Adaptive deep learning models in dynamic manufacturing environments. *Inclusive Growth and Governance Quarterly*, 2(1).
- Malema. (2026b). Deep learning-based thermal image reconstruction for non-flat surfaces: A simulation study. *Inclusive Growth and Governance Quarterly*, 2(1).
- Malema. (2026c). Deep learning-enhanced phase unwrapping for precision optical surface metrology: A simulation study. *Inclusive Growth and Governance Quarterly*, 2(1).
- Malema. (2026d). Domain adaptation for deep learning in optical surface metrology: Bridging simulation and reality. *Inclusive Growth and Governance Quarterly*, 2(1).
- Malema. (2026e). Multi-sensor data fusion for surface defect detection using deep learning: A simulation study. *Inclusive Growth and Governance Quarterly*, 2(1).
- Malema. (2026f). Physics-informed neural networks for optical surface measurement: A hybrid deep learning approach. *Inclusive Growth and Governance Quarterly*, 2(1).
- Malema. (2026g). Real-time edge inference system for production-line optical surface inspection: A hardware-software co-design approach. *Inclusive Growth and Governance Quarterly*, 2(1).
- Malema. (2026h). Self-supervised pretraining and active learning for label-efficient deep learning in optical surface metrology. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026i). Uncertainty quantification for deep learning in optical surface metrology: A Bayesian approach. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026j). Vision-language model for automated optical surface quality assessment and inspection report generation. *Inclusive Growth and Governance Quarterly*, 2(1).