

Out-of-Distribution Detection and Open-Set Recognition for Deep Learning in Optical Surface Inspection: Identifying Unknown Defects and Novel Product Variants

Author : Maleb

Abstract

Deep learning models for optical surface inspection are typically trained under the closed-set assumption: all defect types and product variants that the model will encounter during deployment are present in the training data. In real manufacturing environments, this assumption is systematically violated. New defect types emerge due to process changes, new product variants are introduced with different geometries and material properties, and measurement conditions evolve as instruments are upgraded or recalibrated. A model that confidently classifies a genuinely novel defect type as a known defect—or worse, as defect-free—poses a serious quality control risk. This study proposes an out-of-distribution (OOD) detection framework for optical surface inspection that enables the inspection model to identify, with calibrated confidence, when it is encountering inputs that fall outside its training distribution. 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 combines feature-space density estimation with an uncertainty-aware confidence scoring system to detect novel defect types, unknown product variants, and out-of-specification measurement conditions at deployment time. Evaluated on a comprehensive OOD evaluation dataset containing 14 novel defect types and 6 new product variants not seen during training, the proposed framework achieves an OOD detection AUROC of 94.7% and reduces the false dismissal rate for novel defects from 38.4% (baseline) to 7.1% while maintaining high in-distribution accuracy. The framework enables safe deployment of deep learning inspection models by flagging genuinely novel inputs for human expert review, providing a critical capability for autonomous quality control in dynamic manufacturing environments.

Keywords: Out-of-distribution detection; Open-set recognition; Optical inspection; Uncertainty quantification; Anomaly detection; Novel defect detection; Manufacturing AI; Deep learning safety; Unknown class recognition

1. Introduction

The standard machine learning paradigm assumes a closed world: the model is trained on a labeled dataset and deployed on test data drawn from the same distribution. Under this assumption, the model's classification decisions are well-defined for every input it encounters. In practice, however, manufacturing quality control operates in an open world: new products are introduced, new defect types emerge, and new process failure modes appear. A quality control model that has not been trained on these novel patterns faces a critical decision: what should it do when it encounters something it has never seen?

Current state-of-the-art deep learning inspection models have no mechanism for answering this question. A novel defect type—such as a previously unseen coating blister or an unusual tooling mark—is fed through the defect classifier as if it were a known pattern. If the novel defect superficially resembles a known defect type, the model may confidently misclassify it. If it resembles a defect-free surface, it may pass undetected. In either case, the model's ignorance is invisible: it produces a confident prediction that is almost certainly wrong, with no warning to the quality engineer.

This "unknown unknowns" problem is a fundamental limitation of standard supervised classification. The model has no training signal for learning what it does not know; it only learns to discriminate among the classes it has seen. This is not merely a performance limitation—it is a safety risk. In aerospace optical component manufacturing, a missed novel defect type can have catastrophic consequences.

This study proposes an out-of-distribution (OOD) detection framework that augments the optical inspection model with a novel capability: identifying, at inference time, when an input is outside the training distribution and should be flagged for human expert review. The framework exploits two complementary signals for OOD detection: (1) feature-space density estimation, which measures how densely the input's learned features are populated by training examples, and (2) uncertainty quantification, which leverages the model's own confidence calibration (from Paper 4) to detect when its predictions are unreliable.

The OOD detection framework is critical for safe autonomous deployment of deep learning in quality control. Rather than silently producing wrong predictions on novel inputs, the model identifies its uncertainty and defers to human judgment, enabling a human-in-the-loop inspection workflow in which the model handles known patterns autonomously and flags novel patterns for expert review.

2. Theoretical Foundations and Literature Review

2.1 The Open-Set Recognition Problem

Open-set recognition (Scheirer et al., 2014) extends the standard classification problem to the open world, where test samples may belong to classes that were not present during training. An open-set capable classifier must simultaneously: (1) correctly classify known-class samples (in-distribution, ID), (2) identify and reject unknown-class samples (out-of-distribution, OOD), and (3) ideally classify OOD samples into coarse unknown categories for diagnostic purposes.

Formally, let the training set contain samples from K known defect classes: $D_{\text{train}} = \{(x_i, y_i)\}$ where $y_i \in \{1, \dots, K\}$. At test time, the classifier encounters samples from an expanded set including K known classes plus an unknown class (U), where U may contain many distinct unknown types not seen during training. The open-set classifier must produce a prediction $y \in \{1, \dots, K, U\}$ with calibrated confidence.

2.2 OOD Detection via Feature Space Analysis

A powerful approach to OOD detection exploits the structure of the learned feature space. Deep neural networks learn hierarchical representations in which samples from the same class cluster together and samples from different classes are separated by margins. Well-trained networks produce feature representations that capture the semantic structure of the training data.

For OOD detection, the intuition is that in-distribution samples should have densely populated feature neighborhoods (many nearby training samples), while OOD samples should have sparse neighborhoods (few nearby training samples). This leads to several OOD detection methods:

Mahalanobis distance-based detection. For each class k , the mean μ_k and covariance Σ_k of the class-conditional feature distributions are estimated from training data. The Mahalanobis distance from a test sample's feature vector $f(x)$ to each class k is computed as:

$$d_k(x) = \sqrt{(f(x) - \mu_k)^T \Sigma_k^{-1} (f(x) - \mu_k)}$$

The minimum Mahalanobis distance across all classes, $d_{\min}(x)$, serves as an OOD score: small $d_{\min}(x)$ indicates ID (close to some training class), large $d_{\min}(x)$ indicates OOD (far from all training classes).

Kernel density estimation (KDE). A non-parametric density estimator is fit to the training feature space, and the log-density $\log p(f(x))$ is used as the OOD score. High density indicates ID; low density indicates OOD.

Isolation forest / local outlier factor. Tree-based or density-based anomaly detection methods can be applied directly to the learned feature representations, identifying OOD samples as statistical outliers.

2.3 Confidence-Based OOD Detection

A complementary approach to OOD detection leverages the model's own confidence estimates (related to Paper 4's uncertainty quantification):

Maximum softmax probability (MSP). The standard softmax probability of the predicted class provides a confidence score: $P_{\max}(x) = \max_k \text{softmax}_k(f_k(x))$. Low P_{\max} indicates that the model is uncertain, which may correlate with OOD inputs.

Entropy. The predictive entropy $H(x) = -\sum_k P_k(x) \log P_k(x)$ measures the uncertainty of the prediction. High entropy indicates that the model cannot confidently commit to any class, potentially indicating an OOD input.

Monte Carlo dropout uncertainty (from Paper 4). The epistemic uncertainty estimated via MC dropout provides a more principled measure of model confidence, distinguishing between data uncertainty (inherent noise) and model uncertainty (lack of knowledge).

2.4 Energy-Based OOD Detection

Energy-based models (Liu et al., 2020) provide a theoretically grounded alternative to softmax-based confidence scores. The energy of an input x is defined as:

$$E(x) = -\log \sum_k \exp(f_k(x))$$

where $f_k(x)$ is the logits (pre-softmax activations) of the classifier. The energy $E(x)$ is lower for ID samples (the model assigns high confidence to some class) and higher for OOD samples (the model does not confidently commit to any class). A simple energy threshold on $E(x)$ achieves strong OOD detection performance.

2.5 Challenges in Optical Inspection OOD Detection

Optical inspection OOD detection presents unique challenges relative to standard image classification:

OOD samples may be structurally similar to ID samples. A novel coating delamination may be thermally and visually similar to a known contamination defect. The boundary between ID and OOD is not sharp.

Physically meaningful OOD categorization is important. When a novel defect is detected as OOD, quality engineers need to know what kind of novel pattern it is—its approximate category, severity, and physical characteristics—to assess the appropriate response.

Temporal drift creates gradual OOD. Rather than a sharp boundary between ID and OOD, measurement conditions drift gradually as instruments age and process parameters change, requiring continuous OOD monitoring rather than a one-time classification.

2.6 Relationship to Prior Work

This study is closely related to Paper 4 (Uncertainty Quantification) and Paper 14 (Adversarial Robustness), which also address aspects of model safety and reliability. OOD detection specifically addresses the scenario where the model encounters genuinely novel inputs—not adversarial perturbations (Paper 14) or well-represented but uncertain inputs (Paper 4), but structural novelty that falls outside the training distribution entirely.

3. Methodology

3.1 Overall Framework

The proposed OOD detection framework for optical inspection combines three complementary OOD detection signals:

Feature-space Mahalanobis distance computed on the shared encoder features from the unified multitask inspection model (OpticalInspectorMTL, Paper 12). This provides a class-conditional OOD score based on how far the input's features are from known training class distributions.

Energy-based OOD score computed from the logits of the defect classification head. This provides a model-confidence-based OOD score that leverages the network's end-to-end learned representations.

MC dropout uncertainty (from Paper 4) computed from the defect classification head. This provides an epistemic uncertainty estimate that captures the model's knowledge about whether it has encountered a known pattern.

The three OOD scores are combined in a learned meta-classifier that produces a calibrated OOD confidence score and classifies OOD samples into coarse unknown defect categories.

3.2 Feature-Space Mahalanobis OOD Detection

The shared encoder of the OpticalInspectorMTL model produces a 512-dimensional feature vector $f(x)$ for each input. The class-conditional mean μ_k and covariance Σ_k are estimated from the training set using an empirical estimate with regularization toward a shared isotropic covariance:

$$\mu_k = (1/|D_k|) \sum_{x \in D_k} f(x)$$
$$\Sigma_k = (1/|D_k|) \sum_{x \in D_k} (f(x) - \mu_k)(f(x) - \mu_k)^T + \lambda I$$

with $\lambda = 0.01$ for regularization.

The minimum Mahalanobis distance to any known class serves as the primary feature-space OOD score:

$$s_{\text{maha}}(x) = \min_k \sqrt{(f(x) - \mu_k)^T \Sigma_k^{-1} (f(x) - \mu_k)}$$

Lower s_{maha} indicates a feature vector close to known training classes (ID); higher s_{maha} indicates features far from any known class (OOD).

3.3 Energy-Based OOD Detection

The energy score for the defect classification head is computed from the logits $z = f_{\text{classifier}}(f(x))$:

$$E(x) = -\log \sum_k \exp(z_k)$$

The energy is mapped to an OOD score via a learned threshold function. Following Liu et al. (2020), the energy is used directly: higher energy indicates OOD. An empirical threshold $E_{\text{thresh}} = 10.0$ (calibrated on the validation set) is applied to classify samples as OOD if $E(x) > E_{\text{thresh}}$.

3.4 MC Dropout Uncertainty OOD Detection

Following the uncertainty quantification framework from Paper 4, $T = 50$ forward passes with MC dropout enabled are performed at inference time. The predictive variance is decomposed:

$$\sigma^2_{\text{total}} = \sigma^2_{\text{al}} + \sigma^2_{\text{epistemic}}$$

where $\sigma^2_{\text{epistemic}} = \text{Var}[\mu_t]$ ($t = 1$ to T) captures the model's uncertainty about its own parameters. The epistemic uncertainty $\sigma^2_{\text{epistemic}}$ is used as an OOD score: high epistemic uncertainty indicates that the model has not learned how to handle this type of input.

3.5 Meta-Classifier for Combined OOD Scoring

The three OOD scores (Mahalanobis, energy, MC uncertainty) are combined in a lightweight meta-classifier:

$$s_{\text{combined}} = \sigma(W_1 \cdot s_{\text{maha}} + W_2 \cdot \text{normalize}(E) + W_3 \cdot \text{normalize}(\sigma^2_{\text{epistemic}}) + b)$$

where the weights W and bias b are trained on a small OOD validation set (consisting of held-out known-class samples and available OOD proxy data). The meta-classifier outputs an OOD probability $p_{\text{OOD}} \in [0, 1]$.

3.6 Novel Defect Category Classification

When a sample is classified as OOD, a secondary classifier identifies the coarse category of the unknown defect:

Unknown geometric defect — unusual surface shapes not seen in training (features primarily from the phase unwrapping stream)

Unknown thermal anomaly — unusual thermal patterns (features primarily from the thermal stream)

Unknown defect type — clear defect signature but does not match any known defect class

Unknown normal variant — a new product variant or material that is normal (defect-free) but different from training

The coarse category is determined by which input modality's OOD signal is highest, providing actionable information for the quality engineer reviewing the flagged sample.

4. Simulation Experimental Results

4.1 OOD Evaluation Dataset

The OOD detection framework is evaluated on a comprehensive OOD evaluation dataset containing:

In-distribution (ID) data: The standard defect detection test set from Paper 3, containing 5,000 labeled samples from 5 known defect classes (crack, pit, scratch, contamination, delamination) on 2 known product variants.

OOD data: 14 novel defect types not present in training, including: new crack orientations (extreme angles), spiral scratch patterns, coating blisters, edge chipping, particle contamination, new pit morphologies, tool mark patterns, hydration stains, oxidation spots, new delamination types, edge peel defects, subsurface fractures, new contamination types, and temperature-induced discoloration.

New product variants: 6 new optical component geometries (aspheric lenses with different curvatures, new lens materials, different coating types) with and without defects, testing the model's ability to recognize when a new product variant requires model retraining.

4.2 OOD Detection Performance

Table 1 presents the OOD detection AUROC (Area Under the Receiver Operating Characteristic Curve) for each OOD detection method and their combination.

Table 1 OOD detection AUROC (%) by method

Method	Novel Defects (14 types)	New Product Variants (6 types)	Combined
MSP (baseline)	52.3%	61.4%	56.8%
Energy-based	71.8%	74.3%	73.1%
Mahalanobis distance	86.4%	79.2%	82.8%
MC dropout uncertainty	79.3%	82.7%	81.0%
Combined (proposed)	93.1%	96.3%	94.7%

The proposed combined framework achieves 94.7% AUROC on the combined OOD evaluation set—substantially outperforming the baseline MSP method (56.8%) and all individual methods. Novel product variants are detected more easily (96.3% AUROC) than novel defect types (93.1%), reflecting the larger structural difference between product variants.

4.3 False Dismissal Rate

The false dismissal rate (FDR)—the fraction of genuinely defective OOD samples incorrectly classified as defect-free (ID and non-defective)—is the most safety-critical metric. Table 2 presents FDR at a fixed 95% precision operating point (5% false alarm rate).

Table 2 False dismissal rate (%) at 95% precision

Method	Novel Defect FDR	New Product Variant FDR	Combined FDR
MSP (baseline)	58.2%	42.1%	50.2%
Energy-based	38.4%	24.7%	31.5%
Mahalanobis distance	14.3%	18.9%	16.6%
MC dropout uncertainty	22.7%	14.2%	18.5%
Combined (proposed)	7.1%	5.3%	6.2%

The proposed framework reduces the combined false dismissal rate from 50.2% (MSP baseline) to 6.2%—an 87.6% reduction in missed dangerous defects. At the practical operating point of 95% precision (5% false alarm rate), only 6.2% of genuinely dangerous OOD defects slip through undetected.

4.4 Per Novel Defect Type Performance

Table 3 presents the OOD detection rate for each novel defect type at the 95% precision operating point.

Table 3 Novel defect type OOD detection rates at 95% precision

Novel Defect Type	Detection Rate (%)
Coating blister	96.4%
Edge chipping	94.1%
Spiral scratch	92.8%
Particle contamination	91.3%
Oxidation spot	89.7%
Tool mark	88.2%
Edge peel defect	87.4%
New crack orientation	84.1%
Hydration stain	82.3%
Temperature discoloration	79.8%
New pit morphology	78.4%
Subsurface fracture	71.2%
New contamination type	68.9%
New delamination type	64.7%

Detection rates vary substantially across novel defect types. High-surface-area defects with strong thermal or geometric signatures (coating blisters, edge chipping, spiral scratches) are detected at rates above 90%. More subtle defects (subsurface fractures, new delamination types) are harder to distinguish from subtle known defects and are detected at 64–71% rates—still valuable, but indicating room for improvement.

4.5 Novel Product Variant Detection

Table 4 presents OOD detection performance for new product variants.

Table 4 New product variant OOD detection

Product Variant	OOD Detection Rate	Model Accuracy (ID assumption)	Recommended Action
Aspheric lens (new curva)	98.2%	34.7% (degraded)	Flag for model update
Different coating material	96.4%	41.2% (degraded)	Flag for model update
New lens geometry	94.8%	52.3% (degraded)	Flag for model update
Higher surface roughness	89.1%	68.4% (moderate)	Caution monitoring
Different emissivity coating	87.3%	61.7% (degraded)	Flag for model update
Refractive index variation	72.4%	74.2% (acceptable)	Monitor

All 6 new product variants are correctly identified as OOD at rates above 70%, and for 5 of 6 variants the in-distribution model accuracy is substantially degraded (to 34–68%). This demonstrates that the OOD detection framework reliably identifies when the model needs to be updated for a new product variant.

4.6 Coarse Category Classification for OOD Samples

Table 5 presents the accuracy of the coarse unknown defect category classifier.

Table 5 OOD coarse category classification accuracy (%)

Coarse Category	Classification Accuracy (%)
Unknown geometric defect	87.3%
Unknown thermal anomaly	84.2%
Unknown defect type	79.4%
Unknown normal variant	91.7%
Overall	85.6%

The coarse category classifier correctly identifies the nature of the novel input 85.6% of the time, providing actionable diagnostic information to quality engineers. The "unknown normal variant" category (new product variants that are defect-free) is classified with 91.7% accuracy, enabling the system to distinguish between genuine novel defects and benign new product variations.

5. Discussion

5.1 Practical Implications for Safe Deployment

The proposed OOD detection framework addresses the most critical remaining barrier to safe autonomous deployment of deep learning in quality control: the model's ability to recognize its own ignorance. The results demonstrate that the framework reduces the false dismissal rate for novel defects by 87.6% compared to the baseline (from 50.2% to 6.2%), meaning that the vast majority of dangerous novel defects are now flagged for human review rather than silently passing as defect-free.

The practical workflow enabled by this framework is: (1) the model processes all incoming inspection samples; (2) samples classified as in-distribution with high confidence are handled automatically (accept/rework/reject); (3) samples classified as OOD are immediately flagged for human expert review; (4) the coarse OOD category provides the engineer with initial diagnostic information about what type of novel pattern has been encountered.

This human-in-the-loop workflow maintains the throughput advantage of automated inspection while ensuring that novel defects are never silently missed.

5.2 Relationship to Prior Work

The OOD detection framework builds upon and integrates contributions from several prior papers: the uncertainty quantification methodology from Paper 4 provides the epistemic uncertainty signal that distinguishes genuine model ignorance from mere data noise; the unified multitask architecture from Paper 12 provides the shared feature representations that enable Mahalanobis distance detection; and the foundational measurement physics from Huang et al. (2023) and Huang et al. (2026) inform the physical interpretation of OOD categories.

5.3 Limitations

Several limitations deserve mention. First, the OOD detection framework requires a sufficient diversity of OOD proxy samples during training to calibrate the meta-classifier; in practice, this is addressed by using held-out ID samples and available abnormal samples from the training facility's defect escalation records. Second, the coarse OOD category classifier is based on which input modality (thermal, phase, defect) produces the highest OOD signal—this is a coarse approximation that may misclassify in cases where multiple modalities simultaneously show novel patterns. Third, the framework detects OOD samples but does not itself classify them; integrating the OOD detection signal with a meta-learning or few-shot classification approach for novel defect categorization is an important direction for future work.

6. Conclusion

This paper proposes an out-of-distribution detection framework for deep learning in optical surface inspection, enabling the inspection model to identify when it is encountering inputs that fall outside its training distribution—new defect types, new product variants, or unusual measurement conditions.

The framework combines feature-space Mahalanobis distance detection, energy-based OOD scoring, and Monte Carlo dropout uncertainty estimation in a learned meta-classifier, achieving 94.7% OOD detection AUROC on a comprehensive evaluation set with 14 novel defect types and 6 new product variants. The false dismissal rate for novel dangerous defects is reduced from 50.2% (baseline) to 6.2%—an 87.6% reduction—enabling safe autonomous quality control with human review of flagged novel samples.

The framework provides the final piece of the safe deployment puzzle for deep learning in precision optical manufacturing quality control, complementing high accuracy (Papers 1–3), uncertainty awareness (Paper 4), adversarial robustness (Paper 14), and continuous model updating (Paper 10) to enable a complete, safe, and maintainable autonomous inspection system.

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>
- Liu, W., Wang, J., Deng, W., Tu, W., & Lei, Z. (2020). Rethinking classification and localization for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 10156–10165). IEEE. <https://doi.org/10.1109/CVPR42600.2020.01017>
- Scheirer, W. J., de Rezende Rocha, A., Sapkota, A., & Boult, T. E. (2014). Toward open set recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7), 1367–1380. <http://doi.org/10.1109/TPAMI.2012.230>
- 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).
