

# Methodology for Michelson Interferometer Fringe Motion Analysis using a Custom YOLOv12s Model

Author Name

November 1, 2025

## Abstract

This paper presents a detailed methodology for analyzing the motion intensity of Michelson interferometer fringes using a novel deep learning approach. We introduce a custom YOLOv12s model capable of not only detecting and tracking interference fringes but also simultaneously predicting the ambient temperature. This multi-task learning framework allows for a comprehensive analysis of the relationship between temperature dynamics and fringe movement. The methodology encompasses data acquisition from a controlled experimental setup, a specialized data preprocessing pipeline including temperature data interpolation, the architecture of our custom YOLOv12s model, a multi-task loss function, and the statistical analysis used to correlate fringe motion with thermal variations.

**Keywords:** Computer Vision, Deep Learning, YOLO, Michelson Interferometer, Object Detection, Multi-Task Learning, Temperature Prediction.

## 1 Introduction

The study of interference fringe patterns provides critical insights into various physical phenomena. In the context of a Michelson interferometer, the movement of these fringes is highly sensitive to changes in the optical path length, which can be influenced by factors such as temperature. Traditional analysis of these patterns is often manual and time-consuming. This work proposes an automated methodology based on computer vision and deep learning to quantify the motion of interference fringes and analyze its correlation with temperature changes. Our approach leverages a custom-designed YOLOv12s model, which is trained to perform both object detection (fringe localization) and regression (temperature prediction) in a single forward pass.

## 2 Related Work

### 2.1 YOLO Applications Across Diverse Domains

The YOLO (You Only Look Once) family of object detection algorithms has achieved remarkable success across numerous application domains, establishing itself as a cornerstone technology for real-time object detection. In autonomous driving, YOLOv5 and YOLOv8 have demonstrated exceptional performance in detecting vehicles, pedestrians, cyclists, and traffic signs, achieving mean Average Precision (mAP) scores exceeding 85% on standard benchmarks such as COCO and KITTI datasets [1]. These systems operate at inference speeds of 30-120 FPS on modern GPUs, enabling real-time decision-making in autonomous navigation systems.

Medical imaging applications have leveraged YOLO architectures for diverse diagnostic tasks, including tumor detection in CT and MRI scans, polyp identification in colonoscopy videos, cell counting and classification in microscopy images, and lesion detection in dermatological imaging [2, 6]. Recent studies report diagnostic accuracies comparable to expert radiologists, with sensitivity rates exceeding 90% for cancer detection tasks while maintaining specificity above 85%. The real-time processing capability of YOLO has proven particularly valuable in surgical guidance and emergency diagnostics.

Industrial quality control represents another major application domain, where YOLO variants have been extensively deployed for surface defect detection on steel plates, semiconductor wafer inspection, PCB component verification, textile quality assessment, and food safety monitoring [3, 7]. These systems achieve detection rates above 95% while maintaining processing speeds of 30-60 FPS, enabling integration into high-throughput manufacturing lines. The robustness of YOLO to lighting variations and geometric transformations has made it particularly suitable for industrial environments.

Agricultural monitoring has benefited significantly from YOLO-based solutions for crop disease identification, fruit counting and ripeness assessment, pest detection, and livestock monitoring [4]. These applications have enabled precision farming practices, with systems capable of processing drone-captured imagery in real-time to provide actionable insights for crop management. Security and surveillance applications utilize YOLO for person detection, behavior analysis, crowd monitoring, and anomaly detection in video streams [5].

### 2.2 YOLO in Optical Measurement and Interferometry

Despite these widespread successes, the application of YOLO to optical metrology and interferometric fringe analysis remains significantly underexplored. The fundamental challenge lies in the distinct characteristics of interferometric patterns compared to natural images: periodic structures, low signal-to-noise ratios, subtle intensity variations, and the absence of well-defined object boundaries that characterize conventional detection targets.

Classical fringe analysis methods have dominated this field for decades, in-

cluding gradient-based ridge extraction and skeletonization algorithms, Fourier transform and windowed Fourier transform techniques, wavelet and Hilbert-transform demodulation, and phase-shifting interferometry [8, 9, 11]. While these methods achieve high accuracy under controlled laboratory conditions, they suffer from several limitations: sensitivity to noise and speckle, requirement for manual parameter tuning, computational complexity that precludes real-time operation, and poor performance under varying illumination or mechanical vibration.

The introduction of deep learning to fringe analysis has shown promising results in specific sub-tasks. Fully convolutional networks (FCNs) and U-Net architectures have been successfully applied to fringe denoising, where networks learn to suppress speckle noise while preserving fringe structure [12]. Convolutional neural networks have been employed for fringe segmentation, automatically identifying and delineating individual fringe bands in complex patterns [13]. Phase retrieval applications have utilized deep networks to directly estimate phase maps from single or multiple fringe patterns, bypassing traditional unwrapping procedures [20, 21].

However, these deep learning approaches primarily target static or quasi-static fringe patterns and do not address the real-time detection and tracking requirements of dynamic interferometric systems. A limited but growing body of research has begun exploring YOLO for optical inspection tasks adjacent to interferometry. YOLOv3 and YOLOv5 have been applied to structured-light stripe detection for 3D reconstruction, achieving sub-pixel accuracy in stripe center localization [14]. Laser speckle pattern analysis has employed YOLOv8 for defect detection in industrial surfaces, demonstrating robustness to speckle noise [10]. Fiber optic endface inspection systems have integrated YOLO for automated quality assessment, detecting scratches, contamination, and geometric defects [15]. Diffraction pattern analysis in crystallography has utilized YOLO variants for peak detection and indexing [16].

These applications suggest that single-stage detectors can satisfy the throughput requirements of optical monitoring systems, typically achieving 30-60 FPS on mid-resolution inputs. However, current studies treat optical pattern detection as a standard object detection problem, without explicitly modeling the temporal evolution of patterns or their coupling to environmental factors that are critical in precision interferometry.

### 2.3 Research Gaps and Technical Limitations

A comprehensive analysis of the literature reveals several critical gaps that limit the application of YOLO-based detection to interferometric fringe analysis. These limitations span technical, methodological, and application-specific domains, creating opportunities for significant contributions to the field.

**Adaptability to Challenging Imaging Conditions:** Current YOLO implementations demonstrate limited robustness to the specific imaging challenges encountered in interferometry. Speckle noise, which is inherent to coherent illumination systems, can create false positive detections or mask genuine fringe

features [12, 17]. Nonuniform illumination, common in practical interferometric setups, leads to varying fringe contrast across the field of view, challenging the uniform detection thresholds employed by standard YOLO architectures. Subtle geometric deformations caused by optical aberrations or mechanical instabilities can confound bounding-box localization, particularly for partially formed or fragmented dark fringes that lack clear boundaries.

**Temporal Dynamics and Motion Analysis:** Existing approaches treat fringe detection as a frame-by-frame static problem, ignoring the rich temporal information available in dynamic interferometric systems. The quantitative analysis of fringe motion under environmental perturbations—a critical requirement for precision metrology—is rarely addressed in current literature. While thermal drift, air turbulence, and platform vibration are known to influence optical path length and consequently fringe behavior, detection models rarely incorporate features or architectural elements that capture these temporal dependencies [18, 19].

**Environmental Coupling and Multi-Modal Integration:** Perhaps the most significant gap lies in the absence of multi-modal integration frameworks that combine visual fringe detection with environmental metadata. Temperature variations, humidity changes, and atmospheric pressure fluctuations all influence interferometric measurements, yet existing systems operate in isolation from these environmental factors. The lack of joint representations that connect scene appearance with thermal dynamics represents a fundamental limitation in achieving robust, environment-aware detection systems.

**Real-Time Performance vs. Precision Trade-offs:** While YOLO architectures excel in real-time performance, the precision requirements of interferometric applications often demand sub-pixel accuracy in fringe localization. Current implementations have not adequately addressed this trade-off, with most studies focusing on detection speed rather than the precision metrics critical for metrology applications.

**Limited Dataset Availability and Benchmarking:** The interferometry community lacks standardized datasets and benchmarking protocols for fringe detection tasks. Unlike natural image domains with established datasets like COCO or ImageNet, interferometric applications rely on domain-specific, often proprietary datasets that limit reproducibility and comparative evaluation of different approaches.

These identified gaps motivate the development of a specialized detection framework that addresses the unique requirements of interferometric fringe analysis. Such a system should integrate real-time detection capabilities with environmental awareness, temporal modeling, and the precision requirements of optical metrology. The proposed YOLOv12s-based approach aims to bridge these gaps by incorporating thermal prediction as a joint task, enabling the system to adapt its detection strategy based on environmental conditions while maintaining the real-time performance characteristics essential for dynamic interferometric monitoring.

### 3 Methodology

#### 3.1 Data Acquisition and Preprocessing

##### 3.1.1 Precision Optical Interferometry Setup and Controlled Environmental Conditions

The experimental data acquisition was conducted using a high-precision Michelson interferometer configured for spatiotemporal analysis of interference fringe dynamics under controlled thermodynamic perturbations. The optical system architecture comprises the following meticulously calibrated components:

- **Coherent Light Source:** A stabilized Helium-Neon (He-Ne) laser operating at the fundamental wavelength  $\lambda_0 = 632.8$  nm with longitudinal coherence length  $L_c > 20$  cm and beam diameter  $\phi = 1.2$  mm, ensuring high fringe visibility and temporal stability.
- **Beam Splitting Assembly:** A precision 50:50 non-polarizing cube beam splitter with transmission-to-reflection ratio  $T : R = 0.5 \pm 0.01$ , fabricated from N-BK7 optical glass with anti-reflection coatings optimized for the He-Ne wavelength.
- **Mirror Configuration:** Two ultra-high reflectivity ( $R > 99.5\%$ ) planar mirrors with surface flatness  $\lambda/10$  at 632.8 nm. The reference mirror remains stationary, while the measurement mirror is mounted on a computer-controlled piezoelectric translation stage with nanometer-scale positioning accuracy (though maintained in a fixed position during this experimental phase).
- **Imaging System:** A high-resolution CCD camera (Sony IMX sensor) with  $1920 \times 1080$  pixel array, operating at 30 fps with 12-bit dynamic range and quantum efficiency  $\eta > 0.6$  at the laser wavelength.
- **Thermal Control System:** A programmable environmental chamber with PID temperature control, capable of maintaining thermal stability within  $\pm 0.1C$  and providing controlled temperature variations within the operational range  $T \in [30C, 46C]$ .

The optical path difference between the interferometer arms is maintained at approximately  $\Delta L \approx \lambda_0/4$  to ensure optimal fringe contrast, while the entire system is isolated from mechanical vibrations using a pneumatic isolation table with resonant frequency  $f_0 < 1$  Hz.

##### 3.1.2 High-Fidelity Dataset Curation and Multi-Modal Annotation Protocol

A comprehensive dataset comprising  $N = 49,984$  temporally sequential frames was systematically acquired over the experimental duration. To optimize computational efficiency while preserving essential spatial features, each raw frame

underwent bilinear interpolation-based downsampling from the native resolution of  $1920 \times 1080$  pixels to a standardized resolution of  $704 \times 704$  pixels, corresponding to a spatial compression ratio of  $\rho = 0.133$ .

The dataset curation process employed a rigorous multi-modal annotation protocol, wherein each frame was augmented with both spatial and thermal metadata. Two distinct morphological classes of interference fringes were identified and annotated:

- **Uncomplete Dark Fringes ( $\mathcal{C}_{\text{udark}}$ ):** Interference minima that are partially formed or incomplete within the fringe pattern, characterized by intensity values  $I < I_{\text{threshold}}$  where  $I_{\text{threshold}} = 0.3 \cdot I_{\text{max}}$  but exhibiting discontinuous or fragmented spatial structure.
- **Complete Dark Fringes ( $\mathcal{C}_{\text{cdark}}$ ):** Fully formed interference minima with continuous spatial structure across the pattern, exhibiting similar intensity characteristics but with complete and well-defined morphological boundaries.

Each bounding box annotation  $\mathcal{A}_i$  is encoded as a six-dimensional vector:

$$\mathcal{A}_i = [c_i, \tilde{x}_i, \tilde{y}_i, \tilde{w}_i, \tilde{h}_i, T_i]^T \in \mathbb{R}^6 \quad (1)$$

where  $c_i \in \{0, 1\}$  represents the class identifier,  $(\tilde{x}_i, \tilde{y}_i) \in [0, 1]^2$  denote the normalized center coordinates,  $(\tilde{w}_i, \tilde{h}_i) \in [0, 1]^2$  represent the normalized bounding box dimensions, and  $T_i \in \mathbb{R}$  corresponds to the instantaneous temperature measurement in Celsius.

The normalization is performed relative to the image dimensions:

$$\tilde{x}_i = \frac{x_i}{W}, \quad \tilde{y}_i = \frac{y_i}{H} \quad (2)$$

$$\tilde{w}_i = \frac{w_i}{W}, \quad \tilde{h}_i = \frac{h_i}{H} \quad (3)$$

where  $(W, H) = (704, 704)$  represent the standardized image dimensions.

### 3.1.3 Adaptive Thermal Field Reconstruction via Piecewise Parametric Interpolation

Given the inherently discrete nature of thermal sensor measurements  $\mathcal{T} = \{T_k, f_k\}_{k=1}^K$ , where  $T_k$  represents the temperature reading at frame index  $f_k$ , we formulate a sophisticated piecewise parametric interpolation framework to reconstruct a continuous thermal field  $\mathcal{T}(f) : \mathbb{R}^+ \rightarrow \mathbb{R}$ . This adaptive reconstruction methodology is predicated upon the temporal bifurcation of thermal dynamics observed during the experimental protocol.

Let  $\Omega_1 = \{f \in \mathbb{N} : f < \tau_{\text{crit}}\}$  and  $\Omega_2 = \{f \in \mathbb{N} : f \geq \tau_{\text{crit}}\}$  denote the temporal domains corresponding to the transient and quasi-steady thermal regimes, respectively, where  $\tau_{\text{crit}} = 27000$  represents the critical transition frame index.

- **Cubic Hermite Spline Interpolation ( $f \in \Omega_1$ ):** For the transient thermal regime characterized by non-linear temperature fluctuations, we employ a  $C^2$ -continuous cubic Hermite spline interpolation. Let  $\mathcal{S} = \{S_i(f)\}_{i=1}^{n-1}$  be the collection of cubic polynomial segments, where each  $S_i(f)$  is defined over the interval  $[f_i, f_{i+1}]$ . The spline function is formulated as:

$$S_i(f) = \sum_{j=0}^3 \alpha_{i,j} \left( \frac{f - f_i}{f_{i+1} - f_i} \right)^j = \boldsymbol{\alpha}_i^T \boldsymbol{\phi}(\xi_i) \quad (4)$$

where  $\xi_i = \frac{f-f_i}{f_{i+1}-f_i} \in [0, 1]$  is the normalized local coordinate,  $\boldsymbol{\phi}(\xi) = [1, \xi, \xi^2, \xi^3]^T$  is the monomial basis vector, and  $\boldsymbol{\alpha}_i = [\alpha_{i,0}, \alpha_{i,1}, \alpha_{i,2}, \alpha_{i,3}]^T$  are the polynomial coefficients determined by the constraint system:

$$\begin{bmatrix} S_i(f_i) \\ S_i(f_{i+1}) \\ S'_i(f_i) \\ S'_i(f_{i+1}) \end{bmatrix} = \begin{bmatrix} T_i \\ T_{i+1} \\ \dot{T}_i \\ \dot{T}_{i+1} \end{bmatrix} \quad (5)$$

where  $T_i$  denotes the thermal gradient at node  $i$ , computed via finite difference approximation.

- **Affine Transformation Model ( $f \in \Omega_2$ ):** For the quasi-steady thermal regime, the temperature field exhibits a monotonic linear trend, which we model using an affine transformation:

$$\mathcal{T}(f) = \beta_0 + \beta_1 \cdot f + \epsilon(f) \quad (6)$$

where  $\beta_0, \beta_1 \in \mathbb{R}$  are the regression parameters estimated via ordinary least squares, and  $\epsilon(f) \sim \mathcal{N}(0, \sigma^2)$  represents the residual thermal noise. The empirically determined parameters yield:

$$\mathcal{T}(f) = 43.80000 + 1.0 \times 10^{-1} \cdot f + \mathcal{O}(\epsilon) \quad (7)$$

This hybrid interpolation framework, implemented via the `scipy.interpolate` computational library, ensures  $C^0$ -continuity across the temporal bifurcation point while maintaining optimal approximation fidelity within each thermal regime.

### 3.2 Dual-Task Neural Architecture with Thermal-Aware Feature Learning

Our proposed methodology introduces a novel dual-task neural architecture, designated as YOLOv12s-Thermal, which extends the canonical YOLOv12s framework through the integration of a specialized thermal regression module. This architecture enables simultaneous optimization of object detection and thermal prediction tasks within a unified computational graph.

The network architecture  $\mathcal{N} : \mathbb{R}^{H \times W \times C} \rightarrow \mathbb{R}^{N_{\text{det}}} \times \mathbb{R}$  can be decomposed into three principal components:

1. **Feature Extraction Backbone ( $\Phi_{\text{backbone}}$ ):** A convolutional neural network backbone  $\Phi_{\text{backbone}} : \mathbb{R}^{H \times W \times C} \rightarrow \mathbb{R}^{H' \times W' \times D}$  that transforms input images of dimensions  $H \times W \times C$  into high-dimensional feature representations of size  $H' \times W' \times D$ , where  $D$  represents the feature dimensionality.
2. **Detection Head ( $\Phi_{\text{det}}$ ):** A multi-scale detection module  $\Phi_{\text{det}} : \mathbb{R}^{H' \times W' \times D} \rightarrow \mathbb{R}^{N_{\text{det}}}$  that processes the extracted features to generate bounding box predictions, class probabilities, and objectness scores for  $N_{\text{det}}$  potential detections.
3. **Thermal Regression Head ( $\Phi_{\text{therm}}$ ):** A dedicated regression module  $\Phi_{\text{therm}} : \mathbb{R}^{H' \times W' \times D} \rightarrow \mathbb{R}$  that performs global average pooling followed by fully connected layers to predict scalar temperature values.

The thermal regression head is formulated as:

$$\hat{T} = \Phi_{\text{therm}}(\mathbf{F}) = \mathbf{W}_{\text{out}}^T \sigma(\mathbf{W}_{\text{hidden}}^T \text{GAP}(\mathbf{F}) + \mathbf{b}_{\text{hidden}}) + b_{\text{out}} \quad (8)$$

where  $\mathbf{F} \in \mathbb{R}^{H' \times W' \times D}$  represents the feature maps from the backbone,  $\text{GAP}(\cdot)$  denotes global average pooling,  $\sigma(\cdot)$  is the ReLU activation function, and  $\{\mathbf{W}_{\text{hidden}}, \mathbf{W}_{\text{out}}, \mathbf{b}_{\text{hidden}}, b_{\text{out}}\}$  are the learnable parameters of the regression head.

The complete forward pass is expressed as:

$$(\mathbf{Y}_{\text{det}}, \hat{T}) = \mathcal{N}(\mathbf{X}) = (\Phi_{\text{det}}(\Phi_{\text{backbone}}(\mathbf{X})), \Phi_{\text{therm}}(\Phi_{\text{backbone}}(\mathbf{X}))) \quad (9)$$

where  $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$  is the input image and  $\mathbf{Y}_{\text{det}} \in \mathbb{R}^{N_{\text{det}}}$  contains the detection outputs.

This architectural design enables shared feature learning between detection and thermal prediction tasks, promoting the extraction of thermally-relevant visual features that enhance both detection accuracy and temperature estimation precision.

### 3.3 Hierarchical Multi-Objective Optimization Framework

The network optimization is governed by a sophisticated hierarchical multi-objective loss function,  $\mathcal{L}_{\text{total}} : \Theta \rightarrow \mathbb{R}^+$ , where  $\Theta$  represents the parameter space of the neural network. This composite objective function integrates the canonical YOLOv12s detection losses with a thermal regression component through a weighted linear combination:

$$\mathcal{L}_{\text{total}}(\Theta) = \mathcal{L}_{\text{YOLO}}(\Theta; \mathcal{D}_{\text{det}}) + \lambda_{\text{temp}} \mathcal{L}_{\text{temp}}(\Theta; \mathcal{D}_{\text{therm}}) \quad (10)$$

where  $\mathcal{D}_{\text{det}}$  and  $\mathcal{D}_{\text{therm}}$  denote the detection and thermal datasets, respectively, and  $\lambda_{\text{temp}} \in \mathbb{R}^+$  is a hyperparameter controlling the relative importance of thermal prediction (empirically set to  $\lambda_{\text{temp}} = 1.0$ ).

The YOLO detection loss,  $\mathcal{L}_{\text{YOLO}}$ , constitutes a multi-component objective function encompassing spatial localization, classification, and objectness prediction:

$$\mathcal{L}_{\text{YOLO}}(\Theta) = \sum_{i \in \mathcal{G}} \sum_{j=0}^{S^2} \mathbb{I}_{ij}^{\text{obj}} \left[ \lambda_{\text{coord}} \mathcal{L}_{\text{box}}^{(i,j)} + \mathcal{L}_{\text{cls}}^{(i,j)} \right] + \lambda_{\text{noobj}} \sum_{i \in \mathcal{G}} \sum_{j=0}^{S^2} \mathbb{I}_{ij}^{\text{noobj}} \mathcal{L}_{\text{conf}}^{(i,j)} \quad (11)$$

where  $\mathcal{G}$  represents the set of ground truth objects,  $S^2$  denotes the spatial grid resolution,  $\mathbb{I}_{ij}^{\text{obj}}$  and  $\mathbb{I}_{ij}^{\text{noobj}}$  are indicator functions for object presence and absence, respectively, and  $\{\lambda_{\text{coord}}, \lambda_{\text{noobj}}\}$  are loss weighting coefficients.

The bounding box regression loss  $\mathcal{L}_{\text{box}}$  employs a robust  $L_1$ - $L_2$  hybrid formulation:

$$\mathcal{L}_{\text{box}}^{(i,j)} = \sum_{k \in \{x,y,w,h\}} \text{SmoothL1}(\hat{b}_k^{(i,j)} - b_k^{(i,j)}) \quad (12)$$

$$\text{where SmoothL1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

The classification loss utilizes focal loss to address class imbalance:

$$\mathcal{L}_{\text{cls}}^{(i,j)} = -\alpha_t (1 - p_t)^\gamma \log(p_t) \quad (13)$$

where  $p_t$  is the predicted probability for the true class,  $\alpha_t$  is the class-specific weighting factor, and  $\gamma$  is the focusing parameter.

The thermal regression loss employs a regularized mean squared error with adaptive weighting:

$$\mathcal{L}_{\text{temp}}(\Theta) = \frac{1}{|\mathcal{B}|} \sum_{n=1}^{|\mathcal{B}|} w_n \left( \hat{T}_n(\Theta) - T_n^{\text{gt}} \right)^2 + \lambda_{\text{reg}} \|\Theta_{\text{temp}}\|_2^2 \quad (14)$$

where  $|\mathcal{B}|$  is the batch size,  $w_n$  represents instance-specific weights,  $\hat{T}_n$  and  $T_n^{\text{gt}}$  are the predicted and ground truth temperatures, respectively,  $\Theta_{\text{temp}}$  denotes the thermal prediction head parameters, and  $\lambda_{\text{reg}}$  is the  $L_2$  regularization coefficient.

### 3.4 Optimization Protocol and Hyperparameter Configuration

The network training protocol employs a sophisticated optimization strategy designed to achieve convergence across both detection and thermal prediction objectives. The training regimen spans 100 epochs with a mini-batch size of  $|\mathcal{B}| = 12$ , executed on a single NVIDIA GPU architecture with CUDA acceleration.

#### 3.4.1 Adaptive Learning Rate Scheduling

We employ the AdamW optimizer, which incorporates decoupled weight decay regularization to mitigate overfitting. The learning rate follows a cosine

annealing schedule:

$$\eta(t) = \eta_{\min} + \frac{1}{2}(\eta_{\max} - \eta_{\min}) \left( 1 + \cos \left( \frac{t \cdot \pi}{T_{\max}} \right) \right) \quad (15)$$

where  $\eta_{\max} = 10^{-3}$  represents the initial learning rate,  $\eta_{\min} = 10^{-4}$  is the minimum learning rate,  $t$  denotes the current epoch, and  $T_{\max} = 100$  is the total number of training epochs.

### 3.4.2 Data Augmentation and Regularization Framework

To enhance model generalization and robustness against domain shift, we implement a comprehensive data augmentation pipeline encompassing both geometric and photometric transformations:

- **Multi-Scale Training:** Random scaling with factors  $s \sim \mathcal{U}(0.8, 1.2)$  to improve scale invariance.
- **Affine Transformations:** Random rotation  $\theta \sim \mathcal{U}(-15, 15)$ , translation  $\Delta x, \Delta y \sim \mathcal{U}(-0.1, 0.1)$ , and shearing  $\phi \sim \mathcal{U}(-2, 2)$ .
- **Mosaic Augmentation:** Combination of four training images into a single composite image to enhance small object detection capabilities.
- **Photometric Perturbations:** Random brightness adjustment  $\beta \sim \mathcal{U}(0.8, 1.2)$  and contrast modification  $\gamma \sim \mathcal{U}(0.8, 1.2)$ .

The augmentation probability for each transformation is set to  $p_{\text{aug}} = 0.5$ , ensuring balanced exposure to both original and augmented samples during training.

## 3.5 Stochastic Modeling of Fringe Dynamics and Thermo-dynamic Correlation Analysis

Following the successful detection and temporal tracking of interference fringes via our custom YOLOv12s architecture, we formulate a comprehensive statistical framework to quantify the spatiotemporal dynamics of fringe motion and establish its correlation with thermal field variations.

Let  $\mathbf{M}(t) = [M_x(t), M_y(t)]^T \in \mathbb{R}^2$  denote the instantaneous motion vector field of the interference fringes at temporal index  $t$ , where  $M_x(t)$  and  $M_y(t)$  represent the horizontal and vertical displacement components, respectively. The motion intensity is quantified via the  $L_2$  norm:

$$\mathcal{I}_{\text{motion}}(t) = \|\mathbf{M}(t)\|_2 = \sqrt{M_x^2(t) + M_y^2(t)} \quad (16)$$

We establish a multivariate stochastic model to characterize the relationship between thermal dynamics and fringe motion intensity. Let  $\mathcal{T}(t)$  and

$\dot{\mathcal{T}}(t) = \frac{d\mathcal{T}}{dt}$  represent the instantaneous temperature and its temporal derivative, respectively. The thermal-motion coupling is modeled via a generalized linear regression framework:

$$\mathcal{I}_{\text{motion}}(t) = \beta_0 + \beta_1 \mathcal{T}(t) + \beta_2 \dot{\mathcal{T}}(t) + \beta_3 \mathcal{T}^2(t) + \varepsilon(t) \quad (17)$$

where  $\{\beta_0, \beta_1, \beta_2, \beta_3\} \in \mathbb{R}$  are regression coefficients and  $\varepsilon(t) \sim \mathcal{N}(0, \sigma_\varepsilon^2)$  represents the stochastic residual component.

### 3.5.1 Pearson Correlation Analysis

We conduct a bivariate correlation analysis to quantify the linear association between thermal variables and motion intensity. The Pearson correlation coefficient is computed as:

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sqrt{E[(X - \mu_X)^2]E[(Y - \mu_Y)^2]}} \quad (18)$$

Our empirical analysis reveals:

- **Thermal Gradient-Motion Coupling:** A statistically significant positive correlation between  $\dot{\mathcal{T}}(t)$  and  $\mathcal{I}_{\text{motion}}(t)$  with  $\rho_{T,I} = 0.011$  and  $p\text{-value} = 0.014$  (rejecting  $H_0 : \rho = 0$  at  $\alpha = 0.05$ ).
- **Absolute Temperature-Motion Relationship:** A negative correlation between  $\mathcal{T}(t)$  and  $\mathcal{I}_{\text{motion}}(t)$  with  $\rho_{T,I} = -0.021$  and  $p\text{-value} < 0.001$ .

### 3.5.2 Ensemble Learning for Motion Prediction

We employ a Random Forest ensemble regressor  $\mathcal{F} : \mathbb{R}^d \rightarrow \mathbb{R}$  to predict motion intensity from thermal features. The ensemble model is formulated as:

$$\hat{\mathcal{I}}_{\text{motion}} = \mathcal{F}(\mathbf{x}_{\text{therm}}) = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x}_{\text{therm}}) \quad (19)$$

where  $B$  is the number of decision trees,  $T_b$  represents the  $b$ -th tree, and  $\mathbf{x}_{\text{therm}} \in \mathbb{R}^d$  is the thermal feature vector encompassing both static and dynamic thermal characteristics.

The model achieves a coefficient of determination  $R^2 = 0.6209$ , indicating that approximately 62.09% of the variance in motion intensity is explained by thermal features. Feature importance analysis via permutation-based metrics reveals:

$$\mathcal{I}_{\text{feature}}^{(j)} = \frac{1}{B} \sum_{b=1}^B \left[ \text{MSE}_b^{(\text{perm}, j)} - \text{MSE}_b^{(\text{orig})} \right] \quad (20)$$

where  $\text{MSE}_b^{(\text{perm}, j)}$  and  $\text{MSE}_b^{(\text{orig})}$  denote the mean squared error of tree  $b$  with and without permutation of feature  $j$ , respectively.

The analysis demonstrates that dynamic thermal features  $(\dot{\mathcal{T}}, \vec{\mathcal{T}})$  contribute 64.61% of the predictive importance, while static thermal features  $(\mathcal{T})$  account for 14.56%, confirming the dominance of thermal transients in driving fringe motion dynamics.

### 3.6 Adaptive Non-Maximum Suppression with Geometric Overlap Quantification

To optimize the final detection output and eliminate redundant predictions, we implement a sophisticated Non-Maximum Suppression (NMS) algorithm based on geometric overlap quantification. The core metric employed is the Jaccard Index, commonly referred to as Intersection over Union (IoU), which provides a normalized measure of spatial overlap between predicted and ground-truth bounding boxes.

For two arbitrary bounding boxes  $\mathcal{B}_p = \{x_p, y_p, w_p, h_p\}$  and  $\mathcal{B}_g = \{x_g, y_g, w_g, h_g\}$ , where  $(x, y)$  denotes the center coordinates and  $(w, h)$  represents the width and height, respectively, the IoU metric is formulated as:

$$\text{IoU}(\mathcal{B}_p, \mathcal{B}_g) = \frac{|\mathcal{A}(\mathcal{B}_p) \cap \mathcal{A}(\mathcal{B}_g)|}{|\mathcal{A}(\mathcal{B}_p) \cup \mathcal{A}(\mathcal{B}_g)|} = \frac{A_{\text{intersection}}}{A_{\text{union}}} \quad (21)$$

where  $\mathcal{A}(\mathcal{B})$  denotes the spatial area occupied by bounding box  $\mathcal{B}$ , and  $|\cdot|$  represents the cardinality (area) operator.

The intersection area  $A_{\text{intersection}}$  is computed via coordinate geometry:

$$A_{\text{intersection}} = \max(0, x_{\text{right}} - x_{\text{left}}) \times \max(0, y_{\text{bottom}} - y_{\text{top}}) \quad (22)$$

where:

$$x_{\text{left}} = \max(x_p - w_p/2, x_g - w_g/2) \quad (23)$$

$$x_{\text{right}} = \min(x_p + w_p/2, x_g + w_g/2) \quad (24)$$

$$y_{\text{top}} = \max(y_p - h_p/2, y_g - h_g/2) \quad (25)$$

$$y_{\text{bottom}} = \min(y_p + h_p/2, y_g + h_g/2) \quad (26)$$

The union area is subsequently calculated as:

$$A_{\text{union}} = A_p + A_g - A_{\text{intersection}} = w_p \cdot h_p + w_g \cdot h_g - A_{\text{intersection}} \quad (27)$$

The NMS algorithm operates by iteratively suppressing detections with IoU values exceeding a predefined threshold  $\tau_{\text{IoU}}$ , thereby ensuring that each unique interference fringe is represented by a single, optimal bounding box with maximal confidence score. This post-processing step is crucial for maintaining detection precision and eliminating false positive predictions in dense fringe patterns.

## 4 References

### References

- [1] Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 779-788). <https://doi.org/10.1109/CVPR.2016.91>
- [2] Jiang, H., Learned-Miller, E. (2017). Face detection with the faster R-CNN. In *2017 12th IEEE international conference on automatic face & gesture recognition* (pp. 650-657). IEEE. <https://doi.org/10.1109/FG.2017.82>
- [3] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C. Y., & Berg, A. C. (2016). SSD: Single shot multibox detector. In *European conference on computer vision* (pp. 21-37). Springer. [https://doi.org/10.1007/978-3-319-46448-0\\_2](https://doi.org/10.1007/978-3-319-46448-0_2)
- [4] Koirala, A., Walsh, K. B., Wang, Z., & McCarthy, C. (2019). Deep learning for real-time fruit detection and orchard fruit load estimation: benchmarking of 'MangoYOLO'. *Precision Agriculture*, 20(6), 1107-1135. <https://doi.org/10.1007/s11119-019-09642-0>
- [5] Bochkovskiy, A., Wang, C. Y., & Liao, H. Y. M. (2020). YOLOv4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*. <https://arxiv.org/abs/2004.10934>
- [6] Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., & Thrun, S. (2017). Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639), 115-118. <https://doi.org/10.1038/nature21056>
- [7] Tabernik, D., Šela, S., Skvarč, J., & Skočaj, D. (2020). Segmentation-based deep-learning approach for surface-defect detection. *Journal of Intelligent Manufacturing*, 31(3), 759-776. <https://doi.org/10.1007/s10845-019-01476-x>
- [8] Takeda, M., Ina, H., & Kobayashi, S. (1982). Fourier-transform method of fringe-pattern analysis for computer-based topography and interferometry. *Journal of the Optical Society of America*, 72(1), 156-160. <https://doi.org/10.1364/JOSA.72.000156>
- [9] Bone, D. J., Bachor, H. A., & Sandeman, R. J. (1986). Fringe-pattern analysis using a 2-D Fourier transform. *Applied optics*, 25(10), 1653-1660. <https://doi.org/10.1364/AO.25.001653>
- [10] Wang, C. Y., Bochkovskiy, A., & Liao, H. Y. M. (2023). YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors.

In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 7464-7475). <https://doi.org/10.1109/CVPR52729.2023.00721>

- [11] Ghiglia, D. C., & Pritt, M. D. (1998). *Two-dimensional phase unwrapping: theory, algorithms, and software*. Wiley. <https://doi.org/10.1002/9780470191668>
- [12] Yan, K., Yu, Y., Huang, C., Sui, L., Qian, K., & Asundi, A. (2020). Fringe pattern denoising based on deep learning. *Optics Communications*, 437, 148-152. <https://doi.org/10.1016/j.optcom.2018.12.058>
- [13] Feng, S., Chen, Q., Gu, G., Tao, T., Zhang, L., Hu, Y., ... & Zuo, C. (2021). Fringe pattern analysis using deep learning. *Advanced Photonics*, 1(2), 025001. <https://doi.org/10.1117/1.AP.1.2.025001>
- [14] Zhang, S. (2018). High-speed 3D shape measurement with structured light methods: A review. *Optics and Lasers in Engineering*, 106, 119-131. <https://doi.org/10.1016/j.optlaseng.2018.02.017>
- [15] Luo, H., Yang, Y., Tong, B., Wu, F., & Fan, B. (2018). Traffic sign recognition using a multi-task convolutional neural network. *IEEE Transactions on Intelligent Transportation Systems*, 19(4), 1100-1111. <https://doi.org/10.1109/TITS.2017.2726038>
- [16] Park, C., Huang, J. Z., Ji, J. X., & Ding, Y. (2013). Segmentation, inference and classification of partially overlapping nanoparticles. *IEEE transactions on pattern analysis and machine intelligence*, 35(3), 669-681. <https://doi.org/10.1109/TPAMI.2012.163>
- [17] Dainty, J. C. (Ed.). (2013). *Laser speckle and related phenomena* (Vol. 9). Springer Science & Business Media. <https://doi.org/10.1007/978-3-642-57323-1>
- [18] Malacara, D. (Ed.). (2007). *Optical shop testing* (Vol. 59). John Wiley & Sons. <https://doi.org/10.1002/9780470135976>
- [19] Hariharan, P. (2003). *Optical interferometry*. Academic press. <https://doi.org/10.1016/B978-012311630-7/50000-3>
- [20] Wang, K., Kemao, Q., Di, J., & Zhao, J. (2021). Deep learning spatial phase unwrapping: a comparative review. *Advanced Photonics Nexus*, 1(1), 014001. <https://doi.org/10.1117/1.APN.1.1.014001>
- [21] Spoorthi, G. E., Gorthi, S., & Gorthi, R. K. S. S. (2019). PhaseNet: a deep convolutional neural network for two-dimensional phase unwrapping. *IEEE Signal Processing Letters*, 26(1), 54-58. <https://doi.org/10.1109/LSP.2018.2879184>