

---

# ADAPTING OPENAI'S CLIP MODEL FOR FEW-SHOT IMAGE INSPECTION IN MANUFACTURING QUALITY CONTROL: AN EXPOSITORY CASE STUDY WITH MULTIPLE APPLICATION EXAMPLES

---

A PREPRINT

Fadel M. Megahed<sup>1</sup>, Ying-Ju Chen<sup>2</sup>, Bianca Maria Colosimo<sup>3</sup>, Marco Luigi Giuseppe Grasso<sup>3</sup>, L. Allison Jones-Farmer<sup>1</sup>, Sven Knoth<sup>4</sup>, Hongyue Sun<sup>5</sup>, and Inez Zwetsloot<sup>6,\*</sup>

<sup>1</sup>Farmer School of Business, Miami University, Oxford OH, USA.

<sup>2</sup>College of Arts and Sciences, University of Dayton, Dayton OH, USA.

<sup>3</sup>Department of Mechanical Engineering, Politecnico di Milano, Milan, Italy.

<sup>4</sup>Mathematics & Statistics, Helmut Schmidt University, Hamburg, Germany.

<sup>5</sup>College of Engineering, University of Georgia, Athens, GA, USA.

<sup>6</sup>Amsterdam Business School, University of Amsterdam, Amsterdam, The Netherlands.

\*Corresponding author. Email: [I.M.Zwetsloot@uva.nl](mailto:I.M.Zwetsloot@uva.nl)

January 23, 2025

## ABSTRACT

This expository paper introduces a simplified approach to image-based quality inspection in manufacturing using OpenAI's CLIP (Contrastive Language-Image Pretraining) model adapted for few-shot learning. While CLIP has demonstrated impressive capabilities in general computer vision tasks, its direct application to manufacturing inspection presents challenges due to the domain gap between its training data and industrial applications. We evaluate CLIP's effectiveness through five case studies: metallic pan surface inspection, 3D printing extrusion profile analysis, stochastic textured surface evaluation, automotive assembly inspection, and microstructure image classification. Our results show that CLIP can achieve high classification accuracy with relatively small learning sets (50-100 examples per class) for single-component and texture-based applications. However, the performance degrades with complex multi-component scenes. We provide a practical implementation framework that enables quality engineers to quickly assess CLIP's suitability for their specific applications before pursuing more complex solutions. This work establishes CLIP-based few-shot learning as an effective baseline approach that balances implementation simplicity with robust performance, demonstrated in several manufacturing quality control applications.

*Key Words:* Computer vision; Industry 4.0; supervised fault detection; vision transformer; and visual inspection

## 1 Introduction

The monitoring and inspection of image data in manufacturing environments presents a fundamental dichotomy in statistical quality monitoring. Traditional approaches, primarily leveraging control charts, have provided rigorous theoretical frameworks for detecting process shifts. Concurrently, computational approaches ranging from traditional image processing to modern deep learning architectures have demonstrated remarkable empirical success. This dichotomy has led to two distinct schools of thought, each with inherent strengths and limitations.

The statistical process monitoring school, grounded in control chart theory, offers two compelling advantages. First, these methods require only in-control data for implementation, circumventing the often insurmountable challenge of collecting a representative sample of defective products. Second, they provide theoretical guarantees through their ability to aggregate information temporally, offering insights into shift location, timing, and magnitude. However, the practical implementation of such methods faces significant challenges. They typically require extensive in-control samples, verification of process stability, and specific mathematical formulations closely tied to particular quality data model assumptions, often resulting in methods tailored to specific applications with uncertain generalizability. Furthermore, the lack of accessible software packages poses a general challenge. When applied to image and video image data, traditional approaches typically begin with image preprocessing to extract specific features, followed by control charting to assess process stability over time. Stability can also be evaluated using functional or profile data, representing the in-control dynamics of the video image within the region of interest Colosimo (7), Colosimo et al. (8), Megahed et al. (22), Menafoglio et al. (23). Spatio-temporal modeling has proven effective for identifying out-of-control states, although it introduces additional complexity in the mathematical formulation Yan et al. (32).

The computational school, conversely, has evolved from basic image-processing techniques to sophisticated deep-learning architectures. Early works, such as Megahed and Camelio (18), demonstrated the utility of relatively simple image processing approaches. The field then witnessed a dramatic shift with the success of deep learning architectures, beginning with Krizhevsky et al. (14)'s AlexNet, followed by increasingly sophisticated networks such as VGG (29), ResNet (11), and EfficientNet (30). While powerful, these modern deep-learning methods often require substantial datasets for fine-tuning and optimization. Thus, such methods demand significant expertise in machine learning and domain-specific knowledge. In our estimation, the statistical community's skepticism toward these approaches stems from their empirical nature, extensive data requirements (for both in-control and out-of-control cases, which the only exception of methods based on the one-class classifier), and inability to accumulate evidence over time—though this limitation is primarily relevant when detection sensitivity is suboptimal. Note that the need for large datasets remains true for applications which combine statistical control charts with deep learning methods; for example, Kang et al. (13) used an augmented training set of 16,850 in-control and 7,410 out-of-control images.

Okhrin et al. (24) recently noted that most image process monitoring procedures rely on aggregated image characteristics, such as entropy or arithmetic averages, due to pixel-level analysis's computational and theoretical challenges. This observation parallels the fundamental design of vision transformer models, which inherently operate on aggregated patch-level features rather than individual pixels. This architectural similarity suggests that vision transformer models might naturally align with image quality control applications while addressing some limitations across inspection and monitoring applications.

The emergence of OpenAI's CLIP (Contrastive Language-Image Pretraining) model (27) presents an intriguing opportunity to address these challenges. Trained on 400 million image-text pairs, CLIP's dual-encoder architecture has demonstrated exceptional feature extraction capabilities across various domains. However, our empirical investigations reveal that while CLIP often fails as a zero-shot classifier in manufacturing applications (due to the domain gap between its training data and industrial applications), its performance as a few-shot classifier is remarkably robust.

This finding suggests that CLIP could serve as an initial benchmark for image inspection and monitoring applications, potentially offering:

- Reduced data requirements compared to traditional deep learning approaches
- More generalizable performance across different manufacturing contexts
- Natural handling of high-dimensional image data through its transformer architecture

These characteristics make CLIP appealing for quality engineering practitioners seeking robust, yet implementable, solutions for image-based inspection tasks.

This paper systematically evaluates CLIP’s utility in manufacturing quality control through few-shot learning. Following Megahed et al. (19)’s argument that proper benchmarking is essential before deploying complex methods, we position CLIP as a simple yet powerful baseline that should be evaluated before implementing more sophisticated approaches. Our proposition is straightforward: if CLIP’s few-shot learning capabilities prove sufficient for a given application, there may be no need for more complex methods, especially since our approach can be implemented with minimal code and computational overhead. Specifically, the objectives of this study are threefold:

- (1) To demonstrate how CLIP can be effectively adapted using few-shot learning for manufacturing quality control. We provide a practical framework that reduces implementation complexity while maintaining high accuracy. We emphasize that this framework is a baseline that should be tested before pursuing more complex solutions.
- (2) To investigate the relationship between learning set size and classification performance, and offer guidance on the minimal data requirements for effective implementation.
- (3) To examine the impact of vision transformer model selection (by comparing ViT-L/14 and ViT-B/32) on classification accuracy and computational efficiency.

Through these objectives, we aim to establish a systematic framework for evaluating whether this simplified approach can meet practitioners’ needs before investing in more complex solutions.

To examine the utility of our approach, we present five diverse case studies: metallic pan surface inspection (18), 3D printing extrusion profile analysis, stochastic textured surface evaluation (4, 5), automotive assembly inspection (6), and microstructure image classification for metal additive manufacturing (31). These cases represent a broad spectrum of manufacturing quality control challenges, from well-defined defect patterns to subtle variations in surface texture. Through these examples, we demonstrate that CLIP-based few-shot learning can offer a practical middle ground between traditional statistical process monitoring and modern computational approaches. Our examples indicate that CLIP-based few-shot learning can provide a new paradigm for image-based quality control in manufacturing environments due to its predictive performance and implementation simplicity.

## 2 Background

### 2.1 The CLIP Model Architecture

OpenAI’s CLIP (Contrastive Language-Image Pretraining) represents a significant advancement in multimodal learning. Pretrained on more than 400 million image-text pairs, CLIP employs a dual-encoder architecture that simultaneously processes visual and textual inputs. The uses of the CLIP model are varied and include image search and retrieval, content moderation, text-to-image and image-to-text generation, and zero-shot image classification (27). At its core, CLIP consists of an image encoder  $f_\theta$  and a text encoder  $g_\phi$ , where  $\theta$  and  $\phi$  represent their respective parameters. These encoders map images and text into a shared high-dimensional embedding space  $\mathbb{R}^d$ , where  $d$  varies by encoder model. Table 1 captures some popular image encoder models that can be used within the CLIP architecture.

Table 1: Popular image encoder models used with the CLIP model and their characteristics

Encoder Model	Architectural Details	Embedding Dim.	Compute Efficiency	Detail Capturing	Typical Use Case	Linked Model Card
ViT-B/32	Splits image into 32×32 pixel patches, trading spatial resolution for speed	512	High	Moderate	Simple tasks, small datasets, limited resources	<a href="#">ViT-B/32 (25)</a>
ViT-B/16	Uses 16×16 pixel patches, offering better spatial resolution than ViT-B/32	512	Moderate	High	Medium complexity tasks, more variability	<a href="#">ViT-B/16 (25)</a>
ViT-L/14	Uses 14×14 pixel patches and larger model capacity, maximizing spatial detail	768	Low	Very High	Complex tasks, large datasets, ample resources	<a href="#">ViT-L/14 (25)</a>
RN50	Standard 50-layer ResNet architecture with basic scaling	1024	High	Moderate	Legacy compatibility, fast inference	<a href="#">RN50 (25)</a>
RN50x16	50-layer ResNet with 16x wider layers, offering increased model capacity	768	Moderate	High	Detailed tasks, medium resources	<a href="#">RN50x16 (25)</a>

The CLIP architecture, depicted in Figure 1, operates through three distinct stages. During pretraining, both image and text encoders transform millions of image-text pairs into a shared embedding space. The image encoder first standardizes inputs to model-specific dimensions; typically  $224 \times 224$  pixels for models like ViT-B/32 and ViT-B/16, or  $336 \times 336$  pixels for ViT-L/14 to capture finer details. Simultaneously, the text encoder processes corresponding textual descriptions and creates embeddings that enable the model to learn semantic relationships between visual and textual content. Both the image and text encoders produce embeddings of the same dimensionality, ensuring that the outputs can be directly compared in the shared embedding space. This alignment enables the model to learn semantic relationships between visual and textual content effectively.

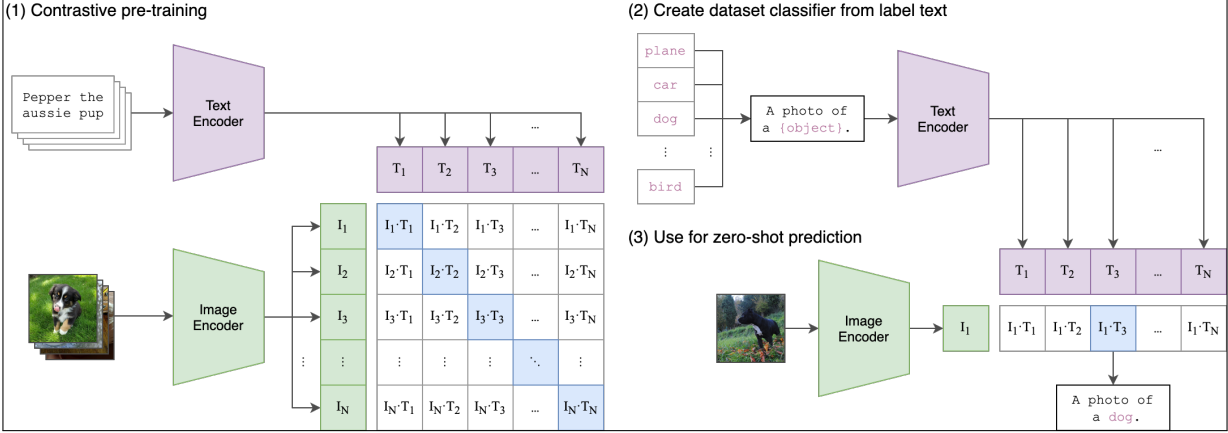


Figure 1: CLIP’s architecture and workflow: (1) the model learns to align image and text embeddings through contrastive learning during pretraining. (2) text templates are created from class labels for classification tasks. (3) Zero-shot prediction compares image embeddings with text embeddings of possible classes. Image source: OpenAI, provided through an MIT license at <https://github.com/openai/CLIP/blob/main/CLIP.png>.

## 2.2 Mathematical Foundations

Input images undergo critical preprocessing steps before entering the CLIP encoders. For an input image  $x \in \mathbb{R}^{H \times W \times 3}$  of arbitrary height  $H$  and width  $W$ , the preprocessing function  $p(\cdot)$  performs:

$$x' = p(x) = \text{Resize}(\text{CenterCrop}(x, s), s) \quad (1)$$

where  $s$  is the model-specific input size. For instance, an image of size  $3264 \times 2448$  is first center cropped to  $2448 \times 2448$  (i.e., the smaller dimension to create a square image) before being resized to  $336 \times 336$  for the ViT-L/14 model. This two-step process helps preserve the proportions of the central region while standardizing the input size.

This standardized input is then divided into non-overlapping patches for Vision Transformer (ViT) models. Given patch size  $P$  (e.g.,  $32 \times 32$  for ViT-B/32 or  $14 \times 14$  for ViT-L/14), the image is segmented into a sequence of  $N = (s/P)^2$  patches, each flattened into a vector:

$$\{x'_1, x'_2, \dots, x'_N\} \quad \text{where} \quad x'_i \in \mathbb{R}^{3P^2}. \quad (2)$$

The image encoder  $f_\theta$  and text encoder  $g_\phi$  transform their inputs into a shared embedding space  $\mathbb{R}^d$ . These embeddings are learned representations that capture semantic meaning in dense vectors. CLIP’s key innovation lies in jointly training these encoders so that semantically similar concepts end up close together in the embedding space, regardless of their modality (image or text).

For comparing embeddings, CLIP uses cosine similarity because it (a) normalizes for vector magnitude, focusing on the directional similarity; (b) bounds similarity scores between -1 and 1, providing numerical stability; and (c) remains differentiable, enabling gradient-based optimization. For a given image embedding  $f_\theta(x)$  and text embedding  $g_\phi(y)$ , their cosine similarity is:

$$s(x, y) = \text{cosSim}(f_\theta(x), g_\phi(y)) = \frac{f_\theta(x) \cdot g_\phi(y)}{\|f_\theta(x)\| \|g_\phi(y)\|}. \quad (3)$$

During training, CLIP optimizes these similarities directly through a contrastive loss function. While the original CLIP implementation includes a temperature parameter  $\tau$  for scaling similarities during training, the core functionality relies on the raw similarities themselves.

CLIP enables both zero-shot and few-shot classification approaches. In zero-shot classification, the model directly compares an image embedding with text embeddings of possible classes:

$$c^* = \arg \max_{c \in \mathcal{C}} \text{cosSim}(f_\theta(x), g_\phi(t_c)) \quad (4)$$

where  $t_c$  represents the text template for class  $c$ , and  $\mathcal{C}$  is the set of possible classes. For few-shot classification, given a set of example images  $\{x_k^c\}_{k=1}^K$  for each class  $c$ , we compute the maximum similarity between a new image  $x$  and all examples within each class:

$$s_c(x) = \max_{k \in \{1, \dots, K\}} \text{cosSim}(f_\theta(x), f_\theta(x_k^c)). \quad (5)$$

New images are then classified based on the class with the highest maximum similarity:

$$c^* = \arg \max_{c \in \mathcal{C}} s_c(x). \quad (6)$$

One can then convert these similarities into a probability using the classical `softmax` transformation (3).

### 2.3 An Illustrative Example

To illustrate how CLIP operates in practice, consider the image in Figure 2 showing a metallic component with a simulated defect. When this image (originally  $3264 \times 2448$  pixels) is processed by the ViT-L/14 model, it undergoes several transformations:

- (1) The image is first preprocessed to the model’s required input size ( $336 \times 336$  pixels).
- (2) The preprocessed image is embedded into a 768-dimensional feature vector as follows:  $\mathbf{v} \in \mathbb{R}^{768} = [0.3713, 0.8574, -0.0626, 0.5874, \dots, -0.0016, -0.6489, -0.1530, -0.1975]$ .
- (3) For zero-shot classification, we provide five text descriptions: (a) “A defective metal component”, (b) “A nominal metal component”, (c) “An industrial part”, (d) “A piece of sheet metal”, and (e) “An artistic photograph”. Each description is first broken down into individual tokens (words and subwords) and padded to a fixed length of 77 tokens. This tokenization step transforms text into a format the model can process.
- (4) Each tokenized description is then encoded into a 768-dimensional embedding. For example, the text “A defective metal component” is encoded into:  $\mathbf{t} \in \mathbb{R}^{768} = [-0.0135, -0.0178, -0.0316, 0.0183, \dots, -0.0169, 0.0006, 0.0008, 0.0638]$ .

- (5) The model computes cosine similarities between the normalized image embedding and each text embedding, which are then converted to probabilities using the softmax function.



Figure 2: Example defective image used to illustrate CLIP’s embedding process and zero-shot classification limitations.

For this example, despite the visible defect (a black marker line simulating a crack), the following probabilities (extracted from the cosine similarities computed by the CLIP model) are presented to each of the pre-assigned five possible text descriptions/categories:

- (a)  $Pr(\text{A defective metal component}) = 0.0616$
- (b)  $Pr(\text{A nominal metal component}) = 0.8506$
- (c)  $Pr(\text{An industrial part}) = 0.0495$
- (d)  $Pr(\text{A piece of sheet metal}) = 0.0380$
- (e)  $Pr(\text{An artistic photograph}) = 0.0001$

Hence, CLIP would select “A nominal metal component” as the label since it had the highest similarity score and probability.

This misclassification, where the model assigns the highest probability (0.8506) to “A nominal metal component” rather than “A defective metal component” (0.0616), illustrates a key limitation: while CLIP excels at general visual-textual understanding through its pretraining on 400 million image-text pairs, it may struggle with domain-specific tasks where the visual features differ significantly from its training distribution. Manufacturing defects, being relatively rare in general web images, represent exactly such a domain gap.

This limitation motivates our few-shot learning approach. Instead of relying on text descriptions alone, we provide CLIP with example images of both nominal and defective components. By computing embeddings for these examples, we create more reliable reference points in the embedding space. This allows CLIP to better distinguish between nominal and defective components based on visual similarity rather than text descriptions.

### 3 Methods

Using CLIP’s visual encoder capabilities, our methodological framework introduces a deliberately simplified approach to few-shot learning in manufacturing quality control. While CLIP offers both visual and textual encoding pathways, we specifically use only its visual encoder for few-shot classification. We deliberately avoid text descriptions in the classification process. This design choice, though not leveraging CLIP’s full potential, establishes a minimal yet effective benchmark for the quality engineering community. More sophisticated approaches, such as Tip-Adapter (34) or Prompt-Generate-Cache (33) methods, could potentially enhance performance by incorporating both visual and textual features. However, the collection of text captions by industrial machine vision systems is uncommon (which is why there is a need for image classification). Therefore, we prioritize establishing this more straightforward baseline for this expository paper. As illustrated in Figure 3, our implementation consists of three primary components: learning set creation, CLIP embedding generation, and test image classification.

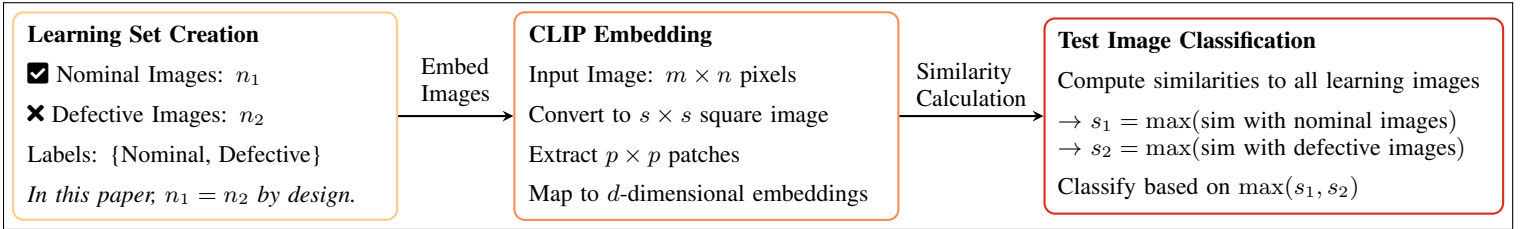


Figure 3: Workflow for learning and classification using CLIP embeddings for image quality control.

#### 3.1 Implementation Framework

The learning set creation phase involves curating balanced learning and testing sets with equal numbers of nominal and defective examples ( $n_1 = n_2$ ). This balanced design choice minimizes the impact of class imbalance (20), allowing us to focus on CLIP’s inherent classification capabilities. Each image undergoes standardized preprocessing to meet CLIP’s input requirements through center cropping (preserving the central region of interest) and resizing to model-specific dimensions.

The CLIP embedding process transforms each preprocessed image into a high-dimensional feature vector. For instance, the ViT-L/14 model maps each  $336 \times 336$  pixel image into a 768-dimensional embedding space, while ViT-B/32 produces 512-dimensional embeddings from  $224 \times 224$  pixel inputs. This transformation occurs through CLIP’s patching mechanism, where images are divided into  $p \times p$  patches ( $p = 14$  for ViT-L/14,  $p = 32$  for ViT-B/32) before being processed by the transformer architecture. Note that the patching idea is similar to the region of interest (ROI) idea in image monitoring applications; however, unlike typical implementations of ROIs (12, 21, 24), the patches are all equal in size and are non-overlapping.

The classification phase uses a maximum similarity approach that compares each test image against the learning set examples. Rather than relying on text descriptions, our few-shot implementation computes similarities between the test image’s embedding and those of the learning set. The classification decision stems from identifying whether the



highest similarity score corresponds to a nominal or defective example. This approach’s computational complexity scales linearly with the learning set size during both the learning and testing phases.

### 3.2 Rationale for the Selected Cases and Computational Experiment Overview

We curated five case studies to systematically examine CLIP’s capabilities and limitations in manufacturing quality control.

- (1) The first case study demonstrates CLIP’s effectiveness in applications where few-shot learning can work well with minimal modification, using metallic pan surface inspection. Here, we compare our approach against the results from Megahed and Camelio (18), showing that CLIP can outperform the original approach with easy-to-use code using the same experimental data.
- (2) The second case study explores CLIP’s performance on lower-resolution extrusion profile images. We investigate how increasing the learning set size improves classification accuracy. The study examines the model’s adaptability to images smaller than its standard  $336 \times 336$  input size. This provides insights into practical applicability in resource-constrained environments.
- (3) The third case study focuses on stochastic textured surfaces. We investigate how model choice impacts performance when defects are subtle or visually challenging to detect. The comparison between ViT-L/14 and ViT-B/32 variants demonstrates that model selection becomes crucial for nuanced defects requiring fine-grained feature detection.
- (4) The fourth case study examines automotive assembly images. It illustrates the limitations of our simplified implementation when dealing with complex images containing multiple components. The study highlights scenarios requiring more sophisticated approaches. It provides guidance for practitioners considering alternative methods.
- (5) The fifth case study examines CLIP’s application for automated classification of microstructural properties. We introduce this example to highlight how binary classification results can be seamlessly extended to multi-class outputs by leveraging the model’s existing capabilities, without retraining or rerunning the model. This scalability demonstrates the practicality of using CLIP for numerous image-based quality engineering applications.

Throughout these studies, we maintain consistent evaluation metrics, including accuracy, sensitivity/recall, specificity, precision, F1-score (harmonic mean of precision and recall), and AUC (area under the receiver operating characteristic curve), while adapting our experimental approach to address each case’s unique challenges and objectives. We assume the reader is familiar with those commonly used classification metrics (otherwise, see Lever et al. (16) for a quick introduction). Our implementation emphasizes reproducibility through Python-based code that interfaces directly with CLIP’s package and OpenAI’s API, requiring minimal dependencies beyond core requirements.

## 4 Case Study 1: Metallic Pan Surface

### 4.1 Dataset Description

Our first case study examines the metallic pan surface inspection dataset from Megahed and Camelio (18). Figure 4 depicts representative nominal and defective images from the dataset, highlighting the nature of the simulated defects.

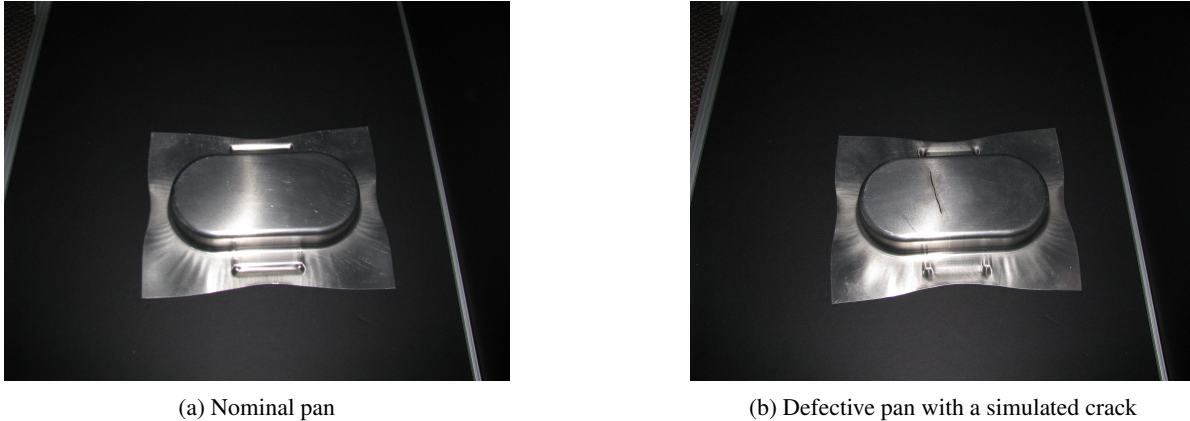


Figure 4: Example images from the metallic pan surface inspection data: (a) shows a nominal pan surface without defects, while (b) displays a pan surface with a simulated crack marked in black.

The original study captured images using a CANON SX 100 IS PowerShot 8.0 Mega-Pixels camera mounted on an aluminum frame (3ft  $\times$  4ft  $\times$  3ft). While conducted in a lab setting, the experimental design deliberately introduced variability to simulate shop floor conditions through:

(1) Part location variations:

- Rotations around the central axis up to  $\pm 15^\circ$
- Translations up to  $\pm 0.6$  cm in both vertical and horizontal directions
- Variable lighting conditions across three different schemes

(2) Defect characteristics:

- Three fault categories: splits, cracks, and inclusions
- Crack and split widths ranging from 1.2 to 1.8 mm
- Defect lengths ranging from 1.8 to 4 cm

We used their dataset structure, which includes a learning set (10 nominal and 10 defective images) and a testing set (50 nominal and 50 defective images). This setup allows us to directly compare our CLIP-based approach with their original results while examining both zero-shot and few-shot learning capabilities. It should be noted that these images were provided directly by the first author of Megahed and Camelio (18) and were not available online before our analysis. This ensures that the images were not part of CLIP’s pre-training dataset, and that this dataset can be used for evaluating CLIP’s generalization capabilities on truly unseen manufacturing data.

## 4.2 Zero-Shot Classification Performance

Building upon our illustrative example in Section 2, we expanded our zero-shot classification experiment in two ways. First, we simplified the classification to a binary problem (nominal vs. defective) rather than the five categories used in the illustration. Second, instead of analyzing a single image, we evaluated CLIP’s performance on the 100 test images (50 nominal, 50 defective). The use of multiple test images enables us to compute comprehensive classification metrics rather than just individual prediction probabilities (as we did in our illustrative example).

We provided CLIP with only textual descriptions: “A metallic pan free of black scuff marks” for the nominal class and “A metallic pan with a simulated scuff mark drawn by a black marker” for the defective class. Table 2 presents the classification metrics, revealing poor prediction performance.

Table 2: Zero-shot classification performance metrics for the metallic pan inspection case study.

Accuracy	Sensitivity (Recall)	Specificity	Precision	F1-Score	AUC
0.630	0.260	1.000	1.000	0.413	0.619

The perfect specificity indicates that CLIP has correctly identified all nominal cases. Moreover, the perfect precision (1) suggests that it was always correct when CLIP identified a defect. However, the low sensitivity (0.260) indicates that CLIP missed many defective cases (i.e., classifying them as nominal) when relying solely on textual descriptions. This performance aligns with our earlier discussion about CLIP’s limitations in zero-shot manufacturing applications due to the domain gap between its training data and industrial applications.

## 4.3 Few-Shot Classification Performance

We then implemented our few-shot learning approach using the same 10 nominal and 10 defective images for learning as those used in Megahed and Camelio (18). Table 3 displays these results alongside the original benchmark results from Megahed and Camelio (18). Our approach shows substantial improvements across most metrics. Moreover, there is no metric where the original approach outperformed the proposed method in predictive performance.

Table 3: Few-shot classification performance metrics for the metallic pan inspection case study.

Metric	Original Results	CLIP Few-Shot	Improvement
Accuracy	0.880	0.940	+0.060
Sensitivity (Recall)	0.760	0.880	+0.120
Specificity	1.000	1.000	0.000
Precision	1.000	1.000	0.000
F1-Score	0.857	0.936	+0.079
AUC	--	0.999	--

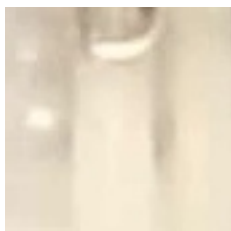
The AUC values were not reported in Megahed and Camelio (18).

These results highlight several key findings. First, our CLIP-based few-shot approach demonstrated significant improvement over the original benchmark across all metrics. Our results show gains in sensitivity (+0.120) while maintaining perfect specificity and precision. Second, the significant improvement in classification accuracy from zero-shot (0.630) to few-shot (0.940) validates our hypothesis that few-shot learning can effectively bridge the domain gap between CLIP’s pre-training and manufacturing applications. Finally, our simple methodology achieves better results without the need for the custom image processing pipelines and feature engineering required by the original approach. Instead, we rely on the standard image processing done by the CLIP model.

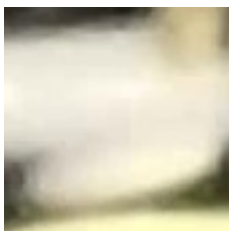
## 5 Case Study 2: 3D Printing Extrusion Profile Inspection

### 5.1 Dataset Description

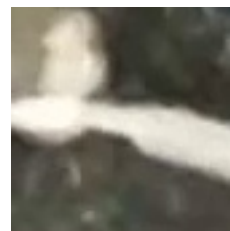
Our second case study examines the direct ink writing (DIW) 3D printing extrusion profile consistency inspection. Through the DIW printing process, complex devices made of composite materials can be obtained to realize different properties, such as high stretchability, heat insulation, and electrical conductivity. Guaranteeing the consistency of the process profile is critical to DIW. Figure 5 depicts representative nominal extrusion profile, over-extruded profile, and under-extruded profile images from the extrusion profile dataset.



(a) Nominal extrusion profile



(b) Over-extruded profile



(c) Under-extruded profile

Figure 5: Examples from the extrusion profile dataset showing (a) a nominal profile with proper extrusion, (b) an over-extruded profile with excessive material flow, and (c) an under-extruded profile with insufficient material flow. Note: Images are shown at their original  $85 \times 85$  pixel resolution.

The printer (Ender-3, Shenzhen Creality 3D Technology Co., Ltd.) is equipped with a custom fixture and an endoscope camera during the experiment. The syringe loaded with the ink is mounted on the fixture. In particular, the silica composite is used as the ink material. First, 3 grams of polyvinyl acetate (PVA) (360627, Sigma-Aldrich) is dissolved into 50 grams of deionized water, and the solution is heated under  $90^\circ\text{C}$  for 30 minutes with magnetic stirring. Then, 12.5 grams of silica nanoparticles (SkySpring Inc.) are added to the solution and mixed for 12 hours. Next, 0.35 grams of a viscosifier (Propylene carbonate, P52652, Sigma-Aldrich) are added into the silica solution and mixed for another two hours. Lastly,  $500\ \mu\text{L}$  of 1-Octanol (A15977, Thermo Fisher Scientific Inc.) are added to remove the bubbles.

The motor pushes out the material and moves along the pre-defined route simultaneously. The endoscope camera monitors the printing condition in real-time, as shown in Figure 5. We annotated the videos during printing and got the extrusion profile dataset. The data were not publicly released before and thus were not part of CLIP’s pre-training dataset.

## 5.2 Zero-Shot Classification Performance

We evaluated CLIP’s zero-shot classification performance using the following text prompts: “An image of a normal extrusion profile from an additive manufacturing process” for the nominal class and “An image of an over-extruded or under-extruded profile from an additive manufacturing process” for the defective class. Table 4 presents the zero-shot classification results.

Table 4: Zero-shot classification performance metrics for the extrusion profile case study.

Accuracy	Sensitivity	Specificity	Precision	F1-Score	AUC
0.620	0.880	0.360	0.579	0.698	0.736

Unlike the metallic pan case study, where zero-shot classification showed high specificity but low sensitivity, here we observe the opposite pattern. The model achieved relatively high sensitivity (0.880) but poor specificity (0.360), suggesting a tendency to over-classify samples as defective. While zero-shot CLIP resulted in a moderate accuracy of 0.620, it cannot be used for this application in practice since it would produce too many false alarms.

## 5.3 Few-Shot Learning and the Impact of Learning Set Size

We first evaluated CLIP’s few-shot learning performance using a baseline set of 20 examples (10 nominal, 10 defective). This setup was selected since it is similar to what we did in the previous case study. Table 5 presents these initial results, where we can see that the predictive performance is poor. We suspect that this is mainly due to the low resolution images available in this case study; the extrusion profile images were only  $85 \times 85$  pixels, significantly smaller than CLIP’s expected input dimensions ( $336 \times 336$  for the ViT-L/14 encoder model). Note that throughout this case study, the number of nominal images was equal to the number of defective images (which was equally split between over and under extruded defects). Hereafter, we use the term “per class” to denote the nominal and overarching defective classes.

Table 5: Few-shot classification performance metrics for the extrusion profile case study using 10 examples per class.

Accuracy	Sensitivity	Specificity	Precision	F1-Score	AUC
0.570	0.400	0.740	0.600	0.480	0.620

To investigate whether additional examples could improve these baseline results, we conducted a comprehensive analysis varying the learning set size from 10 to 350 nominal examples (i.e., total learning sample size is double these numbers) while maintaining a fixed test set for consistent comparison. Our results, depicted in Figure 6, revealed three key findings. First, even with just 10 examples per class, the model achieved an accuracy of 0.57 with 40% of the defects (sensitivity of 0.4) and 74% of the nominal cases (specificity of 0.74) correctly identified. Second, we observed improved performance between 10 and 100 per class examples for some classification metrics. For example,

accuracy improved from 0.57 to 0.71, and AUC increased from 0.62 to 0.75; however, specificity remained relatively stable. Finally, performance improvements plateaued beyond 200 per class examples. The maximum values for our predictive metrics were achieved at 350 nominal learning examples, with an accuracy of 0.80, specificity of 0.81, and sensitivity of 0.78.

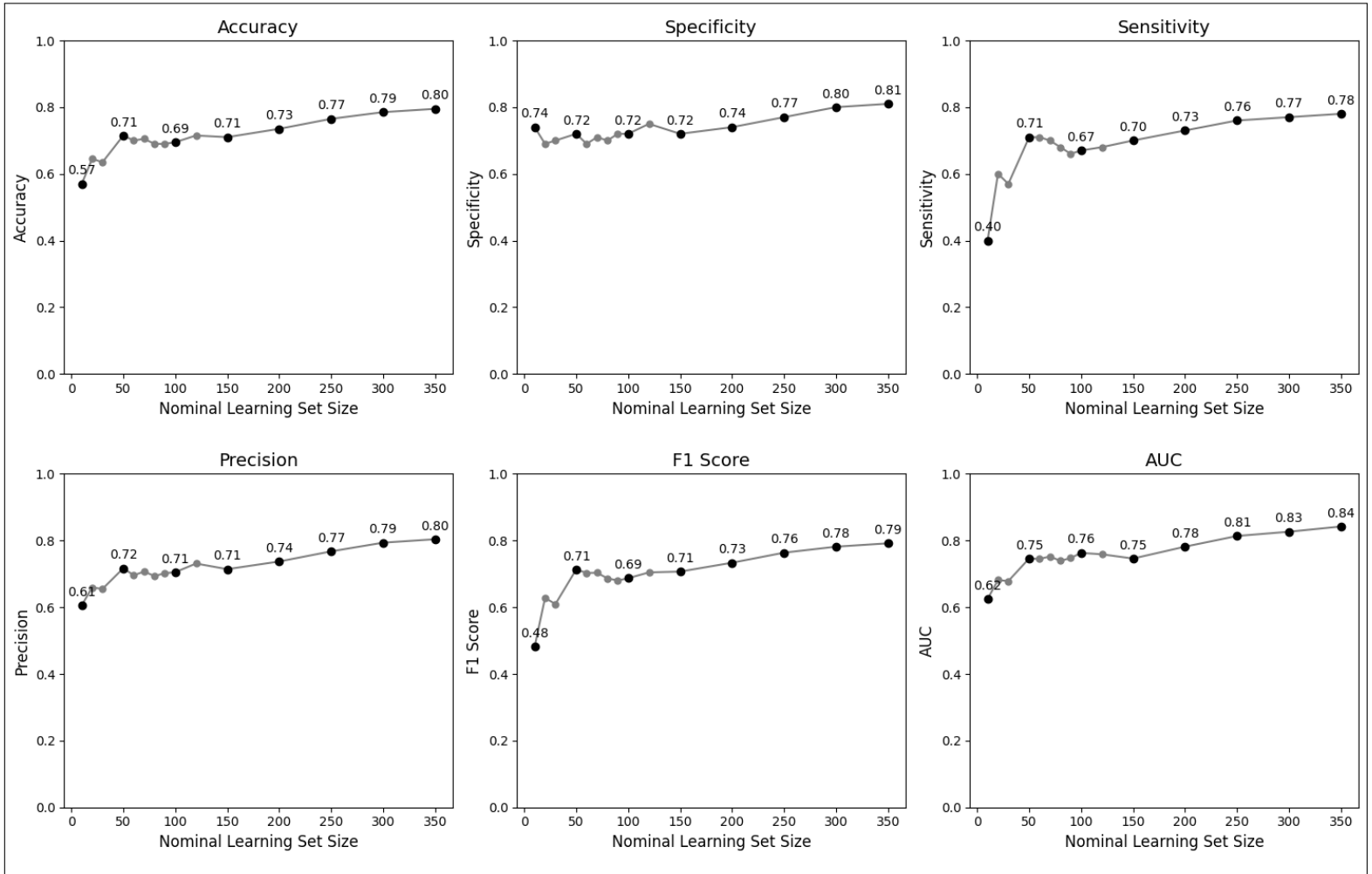


Figure 6: Impact of learning set size on classification metrics for the extrusion profile case study. Black markers highlight specific measurements at key per class sample sizes (10, 50, 100, 150, 200, 250, 300, and 350 nominal examples).

These findings are notable given that the  $85 \times 85$  pixel resolution falls well below CLIP’s designed input specifications. Despite requiring significant upscaling, the model achieved performance gains with modest data requirements. The improvement with just 100 per class examples contrasts with traditional deep learning approaches, which typically require thousands of training samples.


The results suggest that while CLIP can effectively adapt to lower-resolution inputs through few-shot learning, optimal performance likely requires images closer to the model’s intended dimensions. This limitation warrants consideration when applying this approach to manufacturing applications with low-resolution imaging systems.

## 6 Case Study 3: Stochastic Textured Surfaces

### 6.1 Dataset Description

Our third case study examines stochastic textured surfaces (STS), a unique manufacturing quality control data class that presents distinct challenges from traditional profile monitoring. STS data include woven textiles, surface metrology data from machined, cast, or formed metal components, and microscopy images of material microstructures (4). As noted by Bui and Apley (4), STS data present several challenges: (a) unlike our first case study of metallic surfaces where the gold standard is simply a non-textured surface of constant intensity or cases with deterministically repeated patterns where computer-aided design models can serve as references, STS has no definitive “gold standard” image; (b) under normal process conditions, there exist infinitely many valid images that exhibit identical statistical properties while being completely different at the pixel level; and (c) these images cannot be easily aligned, transformed, or warped into a standard reference image due to their inherent stochasticity. These challenges render both traditional profile monitoring and image inspection/monitoring approaches inapplicable (4, 22). This fundamental limitation of traditional approaches motivates our investigation of CLIP’s potential for STS image inspection, i.e., can vision transformers (with few-shot examples) help alleviate the aforementioned limitation of traditional approaches?

To systematically evaluate CLIP’s performance on STS image classification, we utilized the `spc4sts` package (5) to generate a comprehensive dataset of simulated textile images. This approach allowed for the controlled generation of a large-scale data set, including nominal and defective samples with known characteristics and precise manipulation of defect types (local vs. global) and magnitudes. Our experimental dataset included 1,000 nominal images representing standard weave patterns, 500 images with local defects characterized by localized disruptions in the weave pattern, and 500 images with global defects involving systematic shifts in weave parameters. Each image was generated at  $250 \times 250$  pixels following Bui and Apley (5)’s recommended parameters. Nominal images were generated using spatial autoregressive parameters  $\phi_1 = 0.6$  and  $\phi_2 = 0.35$ . Local defects were imposed using the package’s defect generation function, and global defects were simulated by reducing both parameters by 5%. Figure 7 shows representative examples of the nominal and two defect categories. Note the subtle visual differences between the cases.

We used simulated images because Bui and Apley (4) only provided six defect images, which would make it challenging to split the data into few-shot learning and testing sets effectively. We provide the  code (with a fixed seed) and the generated images on our GitHub repository to facilitate the reproduction of our image generation process and/or analysis.

### 6.2 Zero-Shot Classification Performance

We used the following two captions to evaluate CLIP’s zero-shot classification performance: “An image of a textile material with consistent weave patterns, showing no visible defects or irregularities” for the nominal class; and “A woven textile, featuring a visible tear, defect, or a disruption in the material’s weave structure” for the defective class. Table 6 presents these results, revealing significant limitations in CLIP’s zero-shot capabilities for STS inspection. The perfect specificity but zero sensitivity indicates that CLIP classified all images as nominal, i.e., failing to detect any

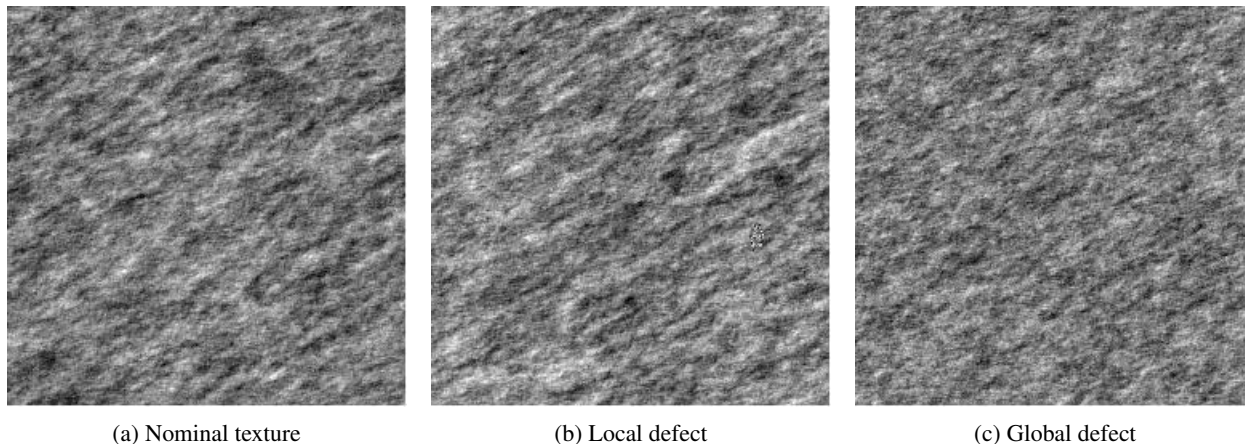


Figure 7: Examples of simulated stochastic textured surfaces (STS): (a) nominal texture with spatial autoregressive parameters  $\phi_1 = 0.6$  and  $\phi_2 = 0.35$ , (b) texture with a localized defect disrupting the pattern, and (c) texture with a global defect simulated by reducing both parameters by 5%.

defects. This performance is particularly poor compared to our previous case studies, likely due to the subtle nature of STS defects and their significant deviation from CLIP’s training data distribution.

Table 6: Zero-shot classification performance metrics for the STS case study.

Accuracy	Sensitivity	Specificity	Precision	F1-Score	AUC
0.500	0.000	1.000	0.000	0.000	0.171

### 6.3 Model Selection Impact on Few-Shot Learning

We suspect that one contributor to the poor performance of the zero-shot approach is the complexity of the STS images. To investigate whether a more detailed vectorization of the images might improve performance, we compare two image encoders: ViT-L/14 and ViT-B/32. The ViT-L/14 encoder uses a  $14 \times 14$  pixel patch size, giving a more fine-tuned vectorization of the images than the ViT-B/32 which uses a  $32 \times 32$  pixel patch. Figure 8 shows how classification metrics evolved with increasing learning set size for both models. Note that throughout this case study, the number of nominal images was equal to the number of defective images (which was equally split between local and global defects).

The results reveal several key findings. The ViT-L/14 model demonstrated superior performance across all metrics, achieving an impressive accuracy of 0.97 with just 50 nominal and 50 defective learning examples. Its performance remained consistently high, exceeding 0.95 for all metrics once the learning set size surpassed 100 nominal and 100 defective examples. It also achieved a near-perfect AUC of 0.99 across most of the learning range. In contrast, the ViT-B/32 model exhibited significantly weaker performance. It had a maximum accuracy of only 0.76, even when provided with 350 nominal and 350 defective examples. Its improvement with increased learning set size was inconsistent. The ViT-B/32 achieved markedly lower sensitivity (0.66) and precision (0.78). The performance gap



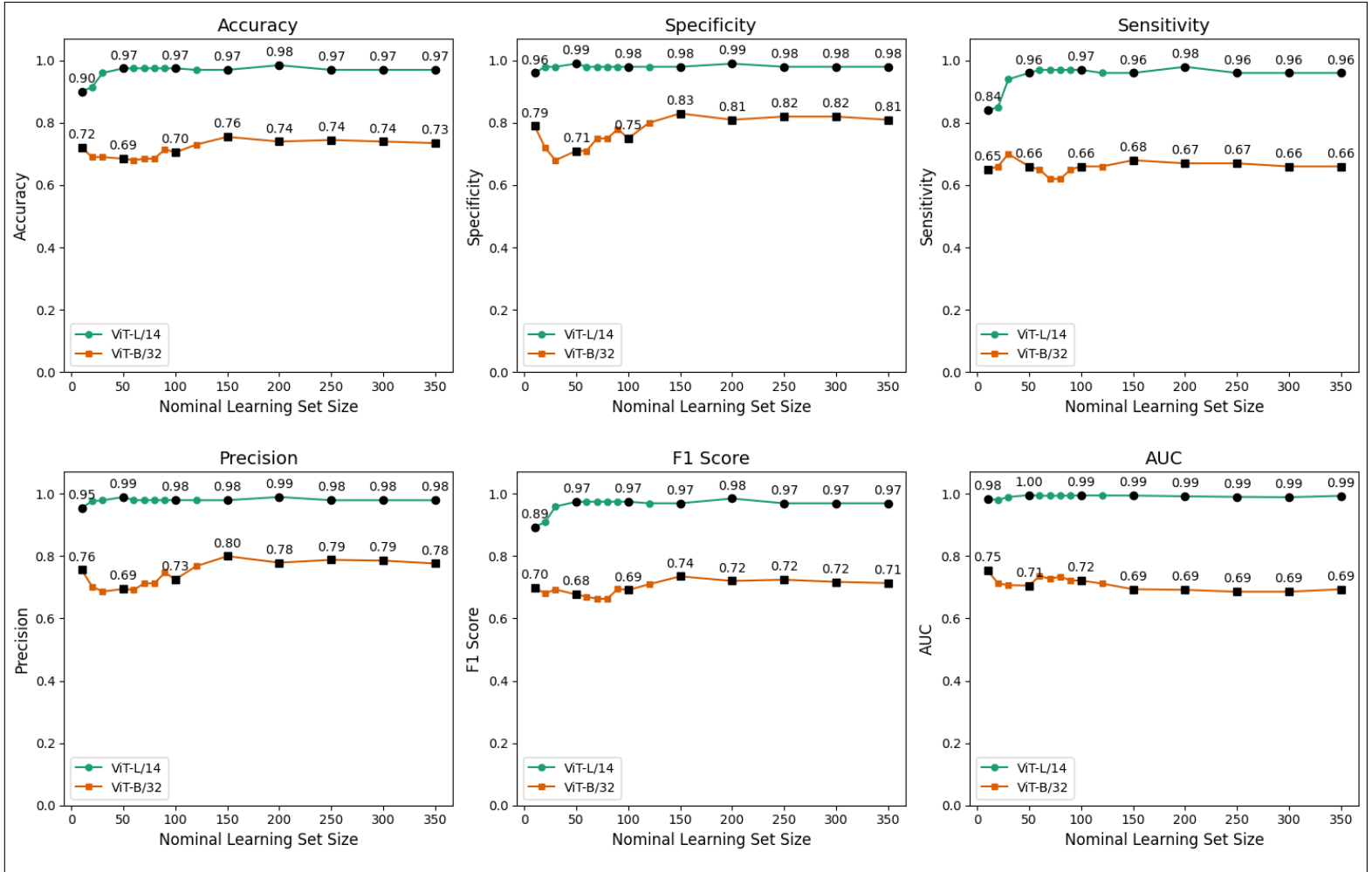


Figure 8: Impact of encoder model and learning set size on classification metrics for the STS case study. Black markers highlight specific measurements at key sample sizes (10, 50, 100, 150, 200, 250, 300, and 350 nominal examples).

between the models persisted even at the maximum learning set size, with ViT-L/14 outperforming ViT-B/32 by 24 percentage points in accuracy (0.97 vs. 0.73), 30 percentage points in sensitivity (0.96 vs. 0.66), and 17 percentage points in specificity (0.98 vs. 0.81).

These findings highlight the importance of encoder model selection for STS inspection tasks. The superior performance of ViT-L/14 suggests that its larger model capacity and finer-grained patch size ( $14 \times 14$  vs  $32 \times 32$ ) allow for capturing the subtle patterns that characterize STS defects. This aligns with theoretical expectations, as smaller patch sizes should better preserve the local spatial correlations that define STS patterns. While the input images ( $250 \times 250$  pixels) were smaller than ViT-L/14’s designed input size ( $336 \times 336$  pixels), the model still performed well. This finding suggests that the model’s architecture is robust to moderate downscaling when the resolution remains sufficient to capture relevant texture patterns.

The results also demonstrate that the ViT-L/14 CLIP-based few-shot learning model can detect local and global STS defects without requiring explicit statistical modeling of the underlying stochastic process. This finding is particularly

significant as it suggests that ViT-L/14 can maintain high classification performance even with smaller-than-designed input image sizes, provided they retain sufficient resolution to capture the textural characteristics. This capability, combined with the elimination of the need for robust parameter estimation and monitoring scheme design, represents a significant practical advantage over traditional approaches to STS inspection.

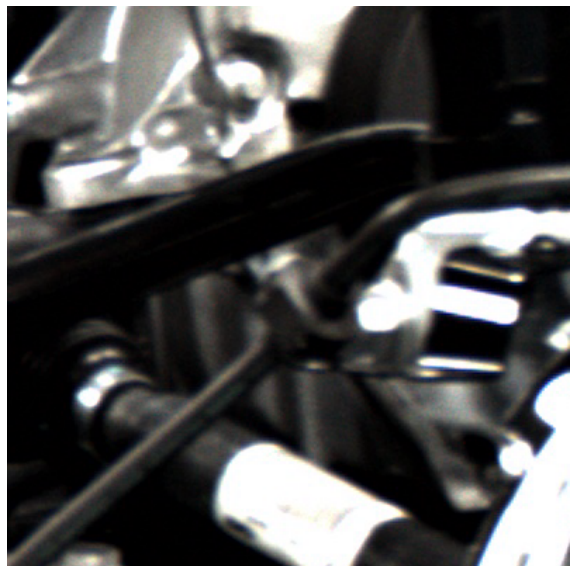
## 7 Case Study 4: Renault’s Pipe Staple

### 7.1 Dataset Description

Our fourth case study examines the recently published Renault pipe staple dataset (6), which presents a real-world automotive assembly inspection challenge. The dataset consists of images capturing the assembly of flexible cables secured by metallic rectangular clamps in automotive manufacturing. Figure 9 shows representative examples of nominal and defective cases from the dataset.



(a) Nominal assembly



(b) Assembly with missing clamp

Figure 9: Examples from the pipe staple dataset: (a) shows a properly secured cable with the metallic clamp in place, while (b) displays a defective assembly where the securing clamp is missing.

This dataset is particularly interesting for our study for several reasons. First, published in July 2023 (approximately two years after CLIP’s development), these images were definitively not part of CLIP’s training data. Second, the dataset was originally constructed for unsupervised classification research, with their test dataset containing both nominal and defective examples. Given that their datasheet indicates the data were collected in chronological order, we leveraged this temporal structure in our study. Specifically, we used the first 50 nominal and 50 defective images from their test dataset for learning, reserving the final 50 images of each class for testing. This chronological splitting approach simulates a realistic implementation scenario where earlier examples inform the classification of later cases.

The inspection task in this dataset presents several unique challenges (6). The images capture varying camera angles and perspectives of the assembled components, contain complex backgrounds with multiple automotive parts visible, and exhibit real manufacturing environment lighting variations and shadows. These characteristics make the visual detection of missing clamps particularly challenging, as the differences between properly secured and unsecured cables can be quite subtle depending on the viewing angle and lighting conditions.

## 7.2 Zero-Shot Classification Performance

For zero-shot classification, we evaluated CLIP using the following two prompts: “A close-up of an automotive assembly process showing a flexible cable secured by a metallic rectangular clamp with rounded edges” for the nominal class, and “A close-up of an automotive assembly process focusing on a flexible cable, where the clamp that typically secures the pipe during assembly is notably absent” for the defective class. Table 7 presents the zero-shot classification results.

Table 7: Zero-shot classification performance metrics for the pipe staple case study.

Accuracy	Sensitivity	Specificity	Precision	F1-Score	AUC
0.500	1.000	0.000	0.500	0.667	0.710

The zero-shot classification results revealed significant limitations in CLIP’s ability to handle this complex inspection task without example images. The perfect sensitivity but zero specificity indicates that CLIP classified all images as defective, suggesting an inability to distinguish the visual cues that differentiate properly secured cables from those missing clamps.

## 7.3 Few-Shot Classification Performance

Given the dataset’s size limitations and our chronological splitting approach, we implemented few-shot learning using a fixed few-shot example set of 50 nominal and 50 defective examples from the earlier portion of the dataset. Table 8 presents these results.

Table 8: Few-shot classification performance metrics for the pipe staple case study using 50 examples per class.

Accuracy	Sensitivity	Specificity	Precision	F1-Score	AUC
0.580	0.560	0.600	0.583	0.571	0.608

While these results show modest improvement over zero-shot classification, they fall short of the performance levels achieved in our previous case studies and state-of-the-art approaches. For example, Carvalho et al. (6) demonstrated that specialized computer vision methods can achieve mean AUC values as high as 98.9% ( $\pm 0.6\%$ ) on this dataset. This substantial performance gap suggests that our simplified CLIP-based approach, while effective for some inspec-

tion tasks, may be insufficient for complex assembly inspection scenarios involving multiple components, varying viewpoints, and subtle defect characteristics. The results reinforce that some manufacturing inspection tasks require more sophisticated approaches, particularly when dealing with real-world assembly operations where defect detection involves complex spatial relationships between multiple components.

## 8 Case Study 5: Microstructure Images

### 8.1 Dataset Description

Manufacturing applications with stringent functional requirements demand consistent, repeatable microstructural properties across parts. Recent advancements in metal additive manufacturing have demonstrated that microstructure can be easily tuned by altering the scanning strategy, offering significant flexibility in microstructure customization (17). Current inspection methods rely on destructive testing and Electron Backscattered Diffraction (EBSD) imaging (2) to assess material properties at the microstructural level.

Manual evaluation by human experts introduces significant limitations in industrial practice (15). The growing volume of material characterization data makes manual labeling resource-intensive. Expert analysis inherently contains subjective biases, limiting the detection of subtle variations within natural data variance. Therefore, the materials science and manufacturing communities have increasingly adopted automated classification methods for microstructure analysis (9, 26). Our investigation of this case is motivated by the suitability of our proposed few-shot learning approach with CLIP to applications with limited training data. If successful, our proposed approach can automate the inspection of microstructure images to reduce the cost associated with expert-driven microstructural measurements.

We simulated EBSD microstructure images by combining random Voronoi tessellation with crystal orientation distribution functions (ODFs) (1). Voronoi edges represent grain size and morphology. The ODF assigns crystal orientations to grains based on predefined probability functions. Crystal orientations map to colors using standard EBSD measurement schemes. Figure 10 illustrates the four example simulated microstructures investigated in our study: (a) uniform (nominal) with random crystal orientations from a uniform distribution, (b) clustered orientation showing a band of grains with predominant crystal orientation surrounded by uniform ODF grains, (c) bimodal with crystal orientations from “bimodal” distribution, and (d) single crystal-like where crystals align toward one orientation. These classes reflect EBSD patterns in metal additive manufacturing. The uniform ODF represents nominal conditions. The other defective classes indicate defects from local or global heating and cooling dynamics variations during manufacturing. Our ODF simulations used weighted schemes to generate three severity levels per class (“low”, “medium”, “high”). These levels represent different degrees of deviation from the nominal uniform ODF. The simulated microstructural images are 400 x 400 pixels. For technical implementation details, we refer the interested reader to (31).

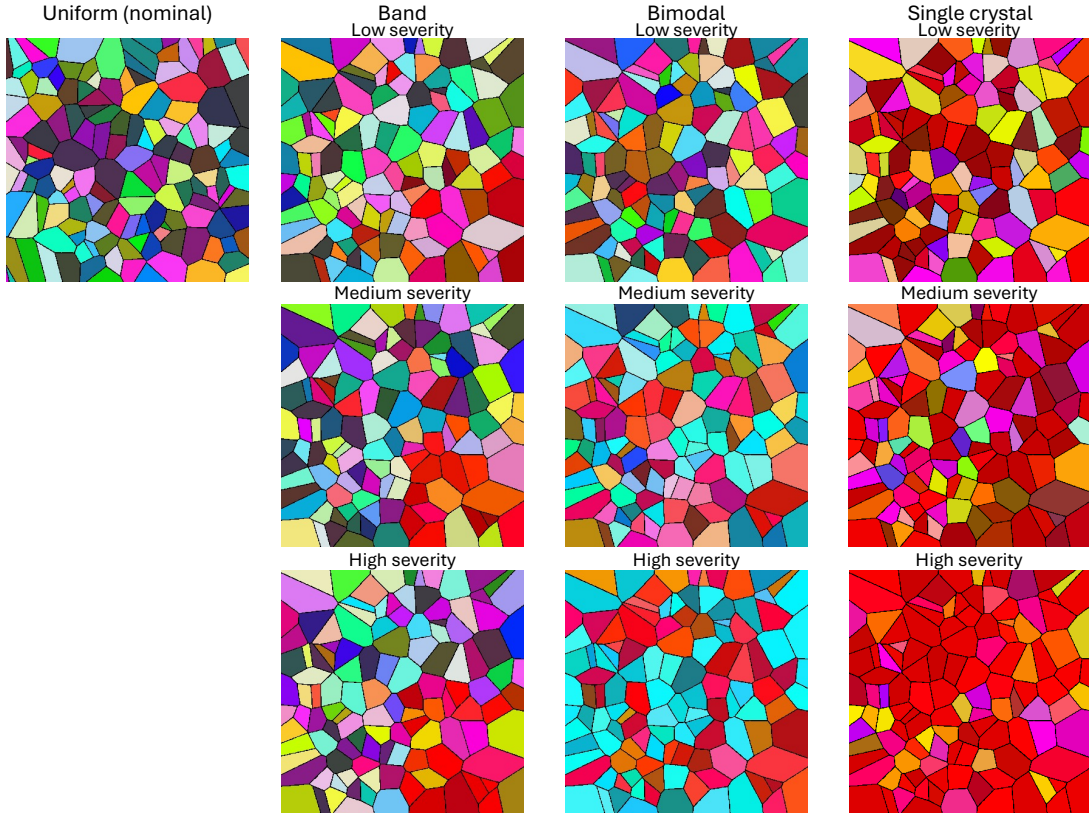


Figure 10: Examples of simulated microstructures with different crystal orientation distributions

## 8.2 Few-Shot Classification Performance

Our study used 72 simulated nominal images for few-shot learning, with eight images for each defective class–severity combination (3 classes [band, bimodal, single-crystal]  $\times$  3 severity levels). The test data had 225 images for the nominal class and 25 for each defect–severity level combination. The training and testing datasets were balanced based on the overarching nominal and defective classes, i.e., the total number of nominal images was equal to the total number of defective images (irrespective of their category and severity). Furthermore, the number of images within each primary defective class (“band”, “bimodal”, and “single-crystal”) was equal. Similarly, the number of images within the sub-labels of primary defective class and severity level combinations was equal. For the sake of conciseness, we only report the few-shot classification performance to highlight the potential of using CLIP with few-shot learning in multi-class image classification applications.

Similar to case studies 1-4, Table 9 presents the performance of the few-shot learning for the binary image classification experiment (nominal vs. defective). Here, we define an image as defective if it belongs to any of the three primary defective classes (“band”, “bimodal”, and “single-crystal”). This few-shot learning binary classification results show that the model excelled at detecting defects with a sensitivity of 0.987, i.e., 222 out of the 225 defective images were correctly classified. Furthermore, nominal images were correctly classified as nominal 80.9% of the time, with 182 out

of 225 nominal images classified correctly. Due to the balanced nature of the dataset, the overall accuracy represents the average of those two numbers and is approximately 90%.

Table 9: Few-shot classification performance metric for the microstructure classification case study.

Test mode	Accuracy	Sensitivity	Specificity	Precision	F1-Score	AUC
Two-class	0.898	0.987	0.809	0.838	0.906	0.952

Given the relatively large number of potential defect classes (up to nine if we consider the class–severity combination), it would be interesting to investigate the model’s multi-class classification performance. Note that we can extract its multiclass classification performance without the need to “retrain” or “rerun” the model. Recall that our approach takes every test image and computes its cosine similarity with every image within the “nominal” and “defective” training dataset. From this comparison, we identify the most relevant image within each nominal and defective class, their match probability, and training-based captions (detailed image description) for its most similar training “nominal” and “defective” images. Consider our test set image `microstructure_039.png` to explain this concept further. The output from our CLIP-based implementation is shown in Figure 11a. The Figure shows that the overall classification result is defective since the probability of being defective (0.505) is slightly larger than non-defective (0.495). We can also see that its most similar defective image is `microstructure_001.png`, which had a description/caption that ends with “low single crystal structure”. We extracted this description using simple text extraction strategies to obtain our multi-class results. We utilized two extraction strategies: (a) focusing only on whether it is “band”, “bimodal”, or “single-crystal”, and (b) assigning it to one of the nine defective sub-labels shown in Figure 10 by extracting severity (“low”, “medium”, and “high”) and the primary defect class (“band”, “bimodal”, and “single-crystal”) indicators.

Figure 12 captures the performance of our CLIP model for both extraction approaches. The first row captures the confusion matrix for the primary labels (nominal, “band”, “bimodal”, and “single-crystal”) alongside its corresponding classification metrics table. The confusion matrix highlights the distribution of actual and predicted labels, with correct classifications appearing on the diagonal and off-diagonal values representing specific classification errors. Unsurprisingly, Figure 12a shows that 182 out of 225 nominal cases were correctly classified as nominal (precisely equal to the specificity values reported in Table 9). This is expected since we did not retrain or rerun the model. However, the extraction of labels allows us to see how the 43 misclassified nominal images were labeled, with 25 labeled as “band” and 18 labeled as “bimodal”. Furthermore, we can see that all 75 (100% of) “single-crystal” defective images were classified correctly. We can also see that 68 out of 75 ( $\simeq 91\%$ ) “band” defects and 60 out of 75 (80%) “bimodal” defects were correctly classified. The accompanying metrics table (Figure 12b) shows an overall classification accuracy of 85.6%, and reports the values for several macro-averaged metrics.

To help explain the idea of “macro-average” computation, let us consider our reported macro-averaged sensitivity of 87.9%. This is computed by performing a one-vs-all comparison for each class. For each class, sensitivity is calculated as the ratio of true positives (diagonal values) to the sum of true positives and false negatives (row total minus the diagonal). The sensitivity for the nominal class is  $\frac{182}{182+25+18} = 0.809$ , and we have shown that sensitivities for “band”, “bimodal”, and “single-crystal” are 0.907, 0.8, and 1, respectively. Hence, the macro-averaged sensitivity

```

Reformatted CSV-Line Output for Image microstructure_039.png

datetime_of_operation: 2025-01-15T17:04:03.471179,
num_few_shot_nominal_imgs: 72,
image_name: microstructure_039.png,
classification_result: Defective,
non_defect_prob: 0.495,
defect_prob: 0.505,

nominal_description: Image microstructure_234.png:
  An microstructure image with uniformly distributed grain sizes and colors.,
defective_description: Image microstructure_001.png:
  An microstructure image with a low single crystal structure.

```

(a) Our implementation’s reformatted output for test image “microstructure\_039.png”.

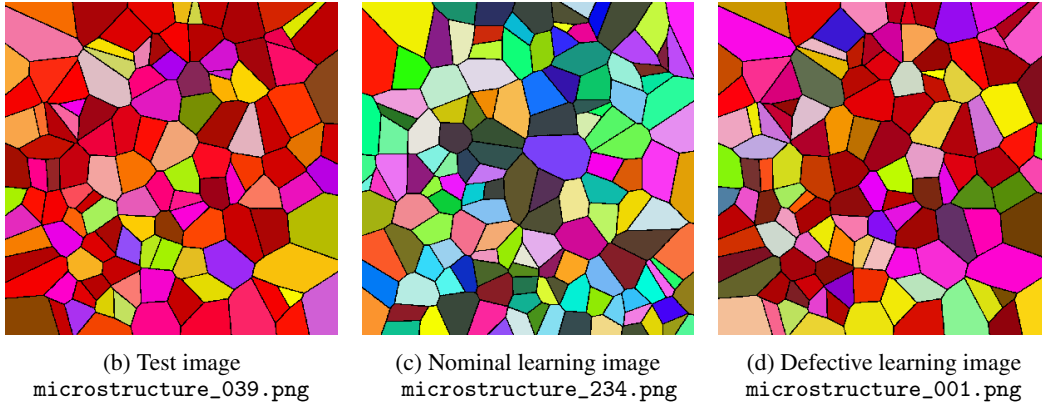
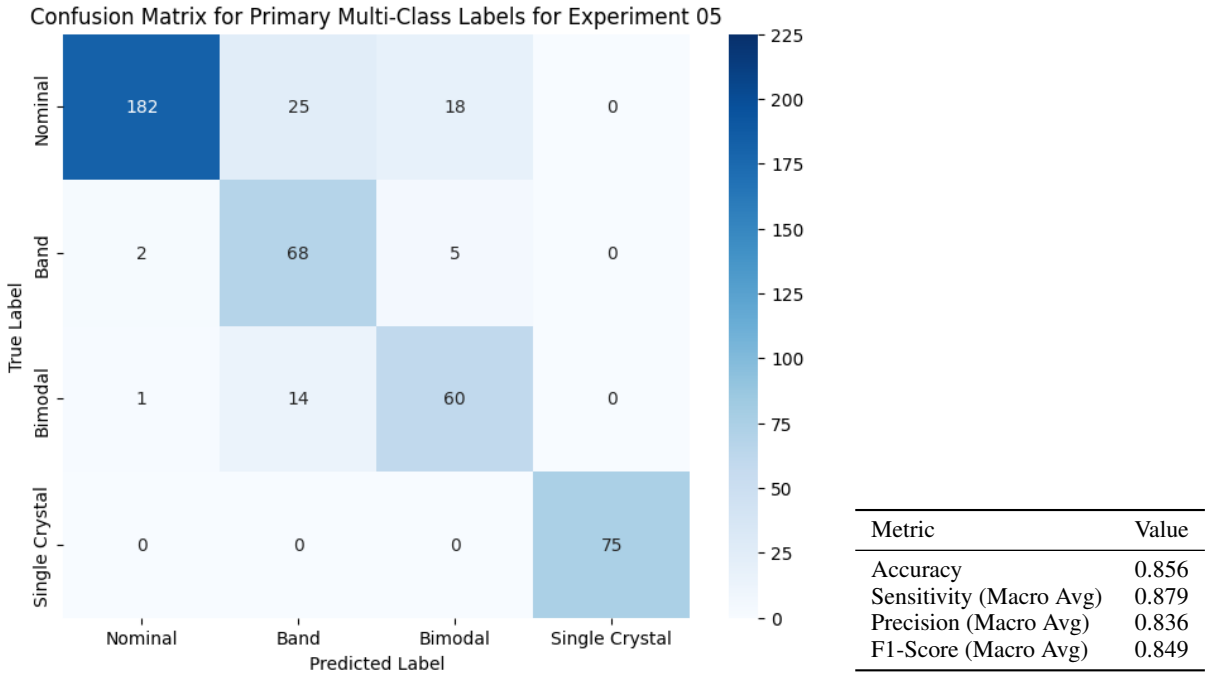


Figure 11: (a) shows the reformatted output for the test image `microstructure_039.png`, including its classification as defective. It is more similar to `microstructure_001.png`, whose corresponding description includes “low single-crystal structure.” (b) presents the test image `microstructure_039.png`. (c) displays the most similar nominal image, `microstructure_234.png`, identified by our model. (d) shows the most similar defective image, `microstructure_001.png`, identified by our model. The binary classification result of defective is correct. In our four-class classification problem (nominal, “band”, “bimodal”, and “single-crystal”), the model correctly classifies our test image as “single-crystal.” However, in our 10-class classification (nominal + nine defect types/severity combinations) problem, the classification of “low single-crystal” is incorrect since the true label of `microstructure_234.png` is “medium single-crystal.” Note that we do not need to rerun the model for the 4-class and 10-class classification problems.

is computed by averaging these values across all classes:  $\frac{0.809+0.907+0.8+1.0}{4} = 0.879$  (10). This computational approach was applied to the remaining metrics shown in Figure 12b.

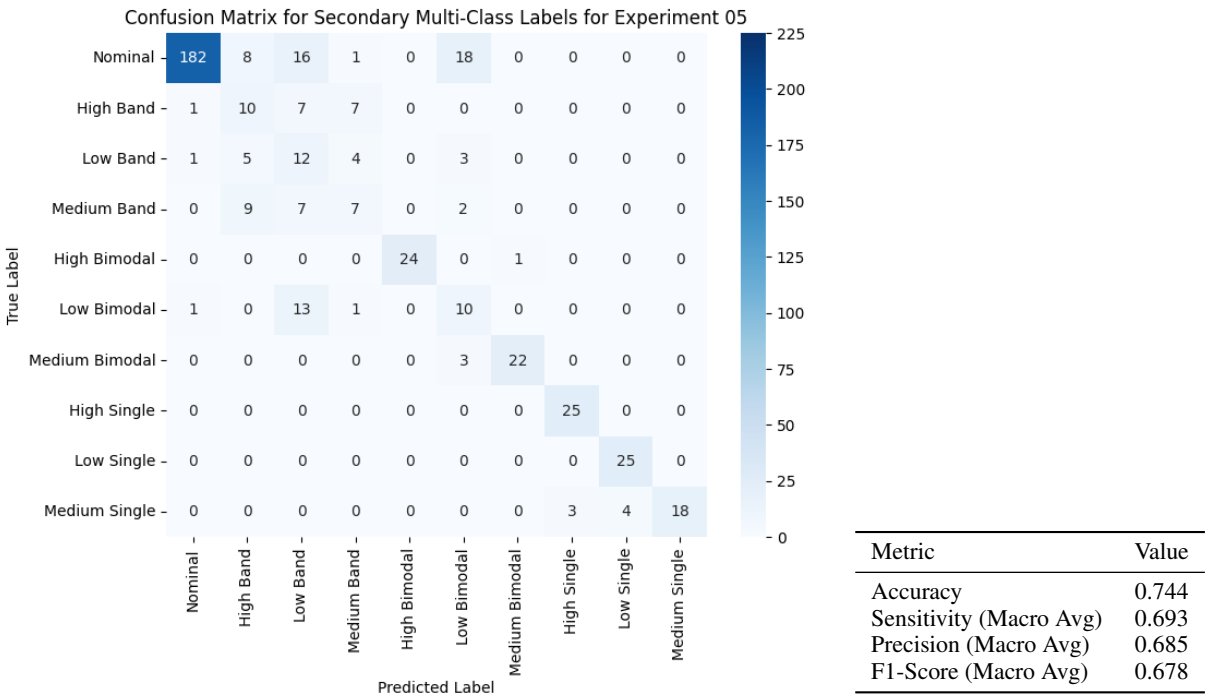
The second row in Figure 12 captures the confusion matrix for the nominal class and the nine defective sublabels alongside its corresponding classification metrics table. As expected, with the increase in the classification degree of difficulty (the model must get both the primary defect class and its severity correct), the overall accuracy drops to 74.4%. However, compared to the primary classification example, the decrease in performance can only be attributed



(a) Primary Confusion Matrix.

Metric	Value
Accuracy	0.856
Sensitivity (Macro Avg)	0.879
Precision (Macro Avg)	0.836
F1-Score (Macro Avg)	0.849

(b) Primary Metrics.



(c) Secondary Label (Sublabel) Confusion Matrix.

Metric	Value
Accuracy	0.744
Sensitivity (Macro Avg)	0.693
Precision (Macro Avg)	0.685
F1-Score (Macro Avg)	0.678

(d) Sublabel Metrics.

Figure 12: Confusion matrices and the corresponding classification performance metrics for primary and secondary multi-class labels for Case Study 5.



to misclassification errors among different severities of the same type of defect (crystal orientation) rather than errors between one crystal orientation and another.

A closer investigation of the results, reported in Figure 12, also reveals that a higher misclassification error was obtained for the “bimodal” ODF compared to the other classes. It was misclassified mainly as “band”. This makes practical sense since the “band” and “bimodal” distributions involve a “bimodal” orientation. The difference in labeling arises from the spatial distribution of orientations: in the “band” class, orientations are spatially clustered, whereas in the “bimodal” class, they are randomly distributed.

These results highlight the proposed approach’s capability to classify complex microstructure data through few-shot classification. The results also show the simplicity of extending our binary classification approach to multiclass classification without rerunning the model. This allows practitioners to perform classifications sequentially: first distinguishing between nominal and defective data and then identifying specific types of defects. Note that our results for this case study will likely be improved if we increase the number of learning examples, as shown in case studies 2 and 3. However, we do not explore this here for three reasons: (a) the obtained results demonstrate the approach’s potential and were satisfactory; (b) our primary rationale for presenting this case study was to show the ease of translating the binary classification results into a multi-class classification output, without the need to rerun or “retrain” the model; and (c) conciseness.

## 9 Discussion

### 9.1 Insights from the Experiments

Our case studies demonstrate that CLIP, when adapted for few-shot learning, can serve as a powerful yet simple baseline for image-based quality control. The model consistently performed well with relatively small learning sets across most applications. For instance, in the metallic pan case study, we achieved 94% accuracy with just 10 examples per class, while the stochastic textured surfaces case demonstrated 97% accuracy with 50 examples per class. Even in the more challenging extrusion profile case, the model reached 80% accuracy with 350 examples per class, despite working with very low-resolution ( $85 \times 85$  pixel) images. This dataset is still much smaller than the 10,000+ examples used in deep learning applications (e.g., see 13).

CLIP successfully processed images ranging from  $250 \times 250$  pixels (STS case) to full-resolution  $3264 \times 2448$  images (metallic pan case). While the ViT-L/14 model is designed for  $336 \times 336$  pixel inputs, it can effectively handle smaller and larger images through appropriate preprocessing. The impact of resolution on performance appears application-dependent, with optimal results generally achieved when image dimensions are larger or relatively close to the model’s design specifications. The results were acceptable in the extrusion case ( $85 \times 85$  pixels) with 350 few-shot examples per case.

Performance varied significantly based on the complexity of the image. The model excelled with single-component images, as demonstrated in the metallic pan. Particularly noteworthy was its success with stochastic textured surfaces,

a traditionally challenging domain for automated inspection. This suggests that transformer-based architectures may inherently capture the statistical properties that define such surfaces, offering a potential new paradigm for texture analysis in manufacturing. Performance in the extrusion case was acceptable with 350 nominal examples. However, the model struggled with multi-component scenarios, as evidenced by the pipe staple case, indicating a limitation in handling complex scenes with multiple interacting elements.

The relationship between learning set size and performance was nonlinear across all cases. For some performance metrics, initial improvements were often large, with substantial gains typically observed in the first 50-100 per class examples. However, diminishing returns appeared after 150-200 per class examples, though the specific inflection point varied by application. This pattern suggests that practitioners can often achieve satisfactory performance with relatively modest data collection efforts, especially when working with single-component or texture-based applications. Not surprisingly, the performance gains across learning set size were nonmonotonic. The variation across learning set size was not likely due to the model’s classification variation; i.e., the same image was classified consistently for a given learning set size. The dips and peaks in performance are likely due to the addition of learning set images, with some additional learning cases making it easier (or harder) to classify the test images.

## 9.2 Advice to Practitioners

Based on our findings and implementation experience, we propose a workflow for practitioners interested in deploying CLIP-based quality control systems (Figure 13). The workflow begins with data collection and preparation, emphasizing consistent imaging conditions across samples. We recommend standardizing camera positions and lighting conditions while maintaining the highest possible image resolution for the application. For the learning set, collect at least 50 examples per class through stratified random sampling across different production runs, shifts, and expected operating conditions to capture natural process variation. When dealing with multiple defect types (as in our extrusion profile case), maintain equal proportions across defect categories. Additionally, prepare a similarly sized independent test set using the same sampling strategy to validate the model’s performance.

The initial evaluation should leverage the provided Python notebook (available in the online materials) to assess the feasibility of using CLIP-based inspection for the intended application. Practitioners only need to update the file paths to their training and testing datasets to quickly determine whether CLIP-based few-shot learning meets their application requirements or aligns with the performance of existing approaches. If the initial results are promising but fall short of expectations, performance can be enhanced by increasing the number of few-shot examples (as illustrated in our extrusion and STS cases) or by transitioning to the ViT-L/14 encoder model. This evaluation process requires minimal setup and provides clear guidance on whether to refine and research few-shot learning with CLIP or proceed to full implementation.

## 9.3 Open Research Questions

Below, we highlight three directions for future research:

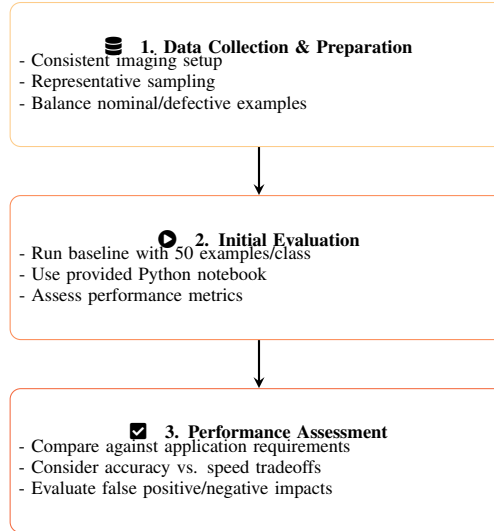


Figure 13: Implementation workflow for CLIP-based quality control systems.

- (1) Exploration of more advanced vision-language models: Beyond our comparison of ViT-L/14 and ViT-B/32, investigating other encoder models could unveil more effective architectures for manufacturing quality control applications. These comparisons should evaluate both performance metrics and computational efficiency to provide actionable insights for industrial implementations.
- (2) Enhancing performance in complex scenarios: Addressing challenging cases, like the pipe staple example, represents another critical research direction. One promising strategy is combining visual and textual embeddings through frameworks like Tip-Adapter (34), enabling more nuanced feature extraction. Alternatively, integrating advanced segmentation models, such as the Segment Anything Model 2 (28), could isolate relevant components before applying the few-shot learning approach with CLIP for classification. While our current method’s simplicity offers practical advantages, these hybrid approaches might be needed for tackling intricate manufacturing environments.
- (3) Monitoring CLIP probability outputs with control charts: Employing simple univariate control charting techniques to monitor the nominal class probabilities derived from our CLIP-based approach over time presents a third research opportunity. This strategy could be particularly valuable in manufacturing settings where out-of-control behavior manifests as a gradual drift rather than an abrupt step change.

## 10 Concluding Remarks


Our expository study demonstrates that our CLIP-based few-shot learning approach can serve as an effective baseline for image-based quality control in manufacturing. Through five case studies, we showed that CLIP can achieve high classification accuracy with relatively small learning sets, often requiring only 50-100 examples per class compared to the thousands typically needed for deep learning approaches. The microstructure case study further highlighted the model’s ability to adapt to multi-class classification tasks without retraining, demonstrating promising perfor-

mance across multiple levels of classification granularity (from nominal vs. defective images to identifying specific defect types and their severities). Overall, the CLIP-based few-shot learning approach excelled at single-component inspection, texture analysis, and microstructure classification. However, its performance was poor in complex multi-component scenes, where more advanced computer vision and/or statistical monitoring approaches were needed. Our findings establish CLIP-based few-shot learning as a practical first baseline for quality engineers and researchers. It offers a simple yet powerful solution that should be evaluated before pursuing more complex implementations.

## Data and Code Availability

To facilitate reproducibility and encourage further research in this area, we provide comprehensive materials in our GitHub repository. The repository contains three main components:

- (1) a data folder with separate subfolders for each case study, organized to clearly distinguish between training and testing datasets, with images categorized by their condition (nominal vs. defective) and defect type where applicable;
- (2) a notebook folder containing our Python implementation using the CLIP model, which we will share as a GitHub Gist (once the paper is accepted) for easy integration with platforms like Google Colab or DeepNote; and
- (3) a results folder storing all experimental outputs, including confusion matrices, animated visualizations of classification performance across different learning set sizes, and both raw and aggregated classification metrics in CSV format.

Our repository includes the  code used to generate the stochastic textured surfaces dataset. These materials are hosted in the GitHub repository ([https://github.com/fmegahed/qe\\_genai](https://github.com/fmegahed/qe_genai)).

To upload the Python notebook from GitHub to Google Colab, start by opening Google Colab in a web browser. Once the Colab interface is open, click the *File* menu located in the top left corner of the screen. From the dropdown options, select *Upload notebook*, which will bring up a dialog box for choosing the source of the notebook. In this dialog, switch to the *GitHub* tab and paste the URL of the GitHub repository: <https://github.com/fmegahed/>. After pasting the URL, press Enter or click the Search button to let Colab locate the notebooks in the repository. Once the repository list appears, locate and select *fmegahed/qe\_genai*. Navigate to the "Path" section and find the notebook titled *notebook/image\_inspection\_with\_clip.ipynb*. Click on the notebook file to load it into Colab. After the notebook is loaded, we can start working on it immediately. Colab allows us to edit, run, and interact with the notebook as needed.

## Acknowledgement

The work of Hongyue Sun is partially supported by the NSF Grant No. FM-2134409 and CMMI-2412678.

## References

- [1] Barrett, T. J., A. Eghtesad, R. J. McCabe, B. Clausen, D. W. Brown, S. C. Vogel, and M. Knezevic (2019). A generalized spherical harmonics-based procedure for the interpolation of partial datasets of orientation distributions to enable crystal mechanics-based simulations. Materialia 6, 100328.
- [2] Bostanabad, R., Y. Zhang, X. Li, T. Kearney, L. C. Brinson, D. W. Apley, W. K. Liu, and W. Chen (2018). Computational microstructure characterization and reconstruction: Review of the state-of-the-art techniques. Progress in Materials Science 95, 1–41.
- [3] Bridle, J. S. (1990). Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition. In F. F. Soulié and J. Héroult (Eds.), Neurocomputing, Berlin, Heidelberg, pp. 227–236. Springer Berlin Heidelberg.
- [4] Bui, A. T. and D. W. Apley (2018). A monitoring and diagnostic approach for stochastic textured surfaces. Technometrics 60(1), 1–13.
- [5] Bui, A. T. and D. W. Apley (2021). spc4sts: Statistical process control for stochastic textured surfaces in r. Journal of Quality Technology 53(3), 219–242.
- [6] Carvalho, P., M. Lafou, A. Durupt, A. Leblanc, and Y. Grandvalet (2024). Detecting visual anomalies in an industrial environment: Unsupervised methods put to the test on the AutoVI dataset. Computers in Industry 163, 104151.
- [7] Colosimo, B. M. (2018). Modeling and monitoring methods for spatial and image data. Quality Engineering 30(1), 94–111.
- [8] Colosimo, B. M., Q. Huang, T. Dasgupta, and F. Tsung (2018). Opportunities and challenges of quality engineering for additive manufacturing. Journal of Quality Technology 50(3), 233–252.
- [9] Gola, J., D. Britz, T. Staudt, M. Winter, A. S. Schneider, M. Ludovici, and F. Mücklich (2018). Advanced microstructure classification by data mining methods. Computational Materials Science 148, 324–335.
- [10] Grandini, M., E. Bagli, and G. Visani (2020). Metrics for multi-class classification: an overview. arXiv preprint arXiv:2008.05756.
- [11] He, K., X. Zhang, S. Ren, and J. Sun (2016). Deep residual learning for image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778.
- [12] He, Z., L. Zuo, M. Zhang, and F. M. Megahed (2016). An image-based multivariate generalized likelihood ratio control chart for detecting and diagnosing multiple faults in manufactured products. International Journal of Production Research 54(6), 1771–1784.
- [13] Kang, Y., Y. Jiao, X. Geng, and M. Nagarajan (2024). Deep vision in smart manufacturing: MODERN framework for intelligent quality monitoring and diagnosis. Research Paper 4856345, University of Miami Business School. Available at <https://dx.doi.org/10.2139/ssrn.4856345>.

- [14] Krizhevsky, A., I. Sutskever, and G. E. Hinton (2012). Imagenet classification with deep convolutional neural networks. In F. Pereira, C. Burges, L. Bottou, and K. Weinberger (Eds.), Advances in Neural Information Processing Systems, Volume 25. Curran Associates, Inc.
- [15] Larmuseau, M., M. Sluydts, K. Theuwissen, L. Duprez, T. Dhaene, and S. Cottenier (2021). Race against the machine: can deep learning recognize microstructures as well as the trained human eye? Scripta Materialia 193, 33–37.
- [16] Lever, J., M. Krzywinski, and N. Altman (2016). Classification evaluation. Nature Methods 13(8), 603–604.
- [17] Liu, Z., D. Zhao, P. Wang, M. Yan, C. Yang, Z. Chen, J. Lu, and Z. Lu (2022). Additive manufacturing of metals: Microstructure evolution and multistage control. Journal of Materials Science & Technology 100, 224–236.
- [18] Megahed, F. M. and J. A. Camelio (2012). Real-time fault detection in manufacturing environments using face recognition techniques. Journal of Intelligent Manufacturing 23, 393–408.
- [19] Megahed, F. M., Y.-J. Chen, L. A. Jones-Farmer, S. E. Rigdon, M. Krzywinski, and N. Altman (2024). Comparing classifier performance with baselines. Nature Methods 21(4), 546—548.
- [20] Megahed, F. M., Y.-J. Chen, A. Megahed, Y. Ong, N. Altman, and M. Krzywinski (2021). The class imbalance problem. Nature Methods 18(11), 1270–1272.
- [21] Megahed, F. M., L. J. Wells, J. A. Camelio, and W. H. Woodall (2012). A spatiotemporal method for the monitoring of image data. Quality and Reliability Engineering International 28(8), 967–980.
- [22] Megahed, F. M., W. H. Woodall, and J. A. Camelio (2011). A review and perspective on control charting with image data. Journal of Quality Technology 43(2), 83–98.
- [23] Menafoglio, A., M. Grasso, P. Secchi, and B. M. Colosimo (2018). Profile monitoring of probability density functions via simplicial functional pca with application to image data. Technometrics 60(4), 497–510.
- [24] Okhrin, Y., W. Schmid, and I. Semeniuk (2024). A control chart for monitoring image processes based on convolutional neural networks. Statistica Neerlandica Early View, 1–26.
- [25] OpenAI (2022). CLIP model card. Available at <https://github.com/openai/CLIP/blob/main/model-card.md>. Accessed: 2025-01-06.
- [26] Pratap, A. and N. Sardana (2022). Machine learning-based image processing in materials science and engineering: A review. Materials Today: Proceedings 62, 7341–7347.
- [27] Radford, A., J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. (2021). Learning transferable visual models from natural language supervision. In International Conference on Machine Learning, pp. 8748–8763. PMLR.
- [28] Ravi, N., V. Gabeur, Y.-T. Hu, R. Hu, C. Ryalı, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, et al. (2024). SAM 2: Segment anything in images and videos. arXiv preprint arXiv:2408.00714v2.

- [29] Simonyan, K. and A. Zisserman (2015). Very deep convolutional networks for large-scale image recognition. In 3rd International Conference on Learning Representations (ICLR 2015), pp. 1–14. Computational and Biological Learning Society.
- [30] Tan, M. and Q. Le (2019). Efficientnet: Rethinking model scaling for convolutional neural networks. In International Conference on Machine Learning, pp. 6105–6114. PMLR.
- [31] Wei, Y., M. Grasso, M. N. Bisheh, K. Paynabar, and B. M. Colosimo (2025). A novel low-dimensional learning approach for automated classification of microstructure data with application to additive manufacturing. Under review.
- [32] Yan, H., K. Paynabar, and J. Shi (2018). Real-time monitoring of high-dimensional functional data streams via spatio-temporal smooth sparse decomposition. Technometrics *60*(2), 181–197.
- [33] Zhang, R., X. Hu, B. Li, S. Huang, H. Deng, Y. Qiao, P. Gao, and H. Li (2023). Prompt, generate, then cache: Cascade of foundation models makes strong few-shot learners. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15211–15222.
- [34] Zhang, R., W. Zhang, R. Fang, P. Gao, K. Li, J. Dai, Y. Qiao, and H. Li (2022). Tip-adapter: Training-free adaption of clip for few-shot classification. In European Conference on Computer Vision, pp. 493–510. Springer.