

Industrial Anomaly Detection with Localization: A Comparative Study under Clean and Shifted Domains

Ivan Necerini
s345147

s345147@studenti.polito.it

Jacopo Rialti
s346357

s346357@studenti.polito.it

Fabio Veroli
s336301

s336301@studenti.polito.it

Abstract

[?]

1. Introduction

Visual anomaly detection in manufacturing aims to automatically identify defective products from images acquired under controlled conditions. This task is commonly framed as *one-class classification*, where a model learns only from normal samples and must identify deviations at test time [1]. Since anomalous examples are often scarce or unavailable during training, supervised approaches are impractical, motivating unsupervised methods that learn what constitutes normality from defect-free data.

A practical anomaly detection system must address two complementary goals: (i) *image-level detection*, determining whether an image is normal or anomalous, and (ii) *pixel-level localization*, identifying which regions contain defects. State-of-the-art methods such as **PatchCore** [7] and **PaDiM** [2] achieve strong results on standard benchmarks like MVTec AD [1] by modelling patch-level feature distributions from pre-trained networks. However, their robustness to distribution shift and sensitivity to threshold calibration remain less explored.

In this work, we present a comparative study of PatchCore and PaDiM on MVTec AD under both clean and synthetically shifted conditions. Our contributions are:

1. **Baseline comparison.** We evaluate PatchCore (custom implementation) and PaDiM (anomalib-based) on three MVTec AD categories (Hazelnut, Carpet, Zipper), reporting both image-level and pixel-level metrics.
2. **Synthetic domain shift.** We construct *MVTec-Shift* by applying photometric and geometric transformations inspired by MVTec AD 2 [4], and measure the performance drop when models trained on clean data are

tested on shifted data. Measuring also the hyperparameter influence on the performance obtained.

3. **Threshold adaptation strategies.** We compare three regimes: (a) no adaptation, (b) threshold-only adaptation on shifted validation data, and (c) full adaptation with model re-training. This ablation isolates the benefit of threshold recalibration from feature-level adaptation.
4. **Model-unified setting.** Following recent literature [8, 3, 5], we train a single model on pooled normal data from all categories while using per-class thresholds at inference. This setting enables investigation of the *identical shortcut* problem [8], where normal samples from one class may be misclassified as anomalies when evaluated against another class's threshold due to overlapping feature distributions.

The paper is organized as follows: Section 2 reviews related work; Section 3 describes the methods; Section 4 details the experimental setup; Section 5 presents results and analysis; Section 6 concludes with limitations and future directions. Figure 1 provides a high-level overview of our experimental pipeline.

2. Related Work

Unsupervised anomaly detection methods for industrial images can be broadly categorized into four families: reconstruction-based, distribution-based, memory-based, and unified approaches. Reconstruction methods [1] learn to encode and decode normal images, flagging samples with high reconstruction error as anomalous. While intuitive, these approaches often struggle to localize subtle defects when the autoencoder generalizes too well. Distribution-based approaches model the statistical properties of normal features, whereas memory-based methods store and retrieve representative embeddings at test time. In this work we focus on two state-of-the-art methods from the latter categories, **PaDiM** and **PatchCore**, which both leverage features from pre-trained convolutional networks.

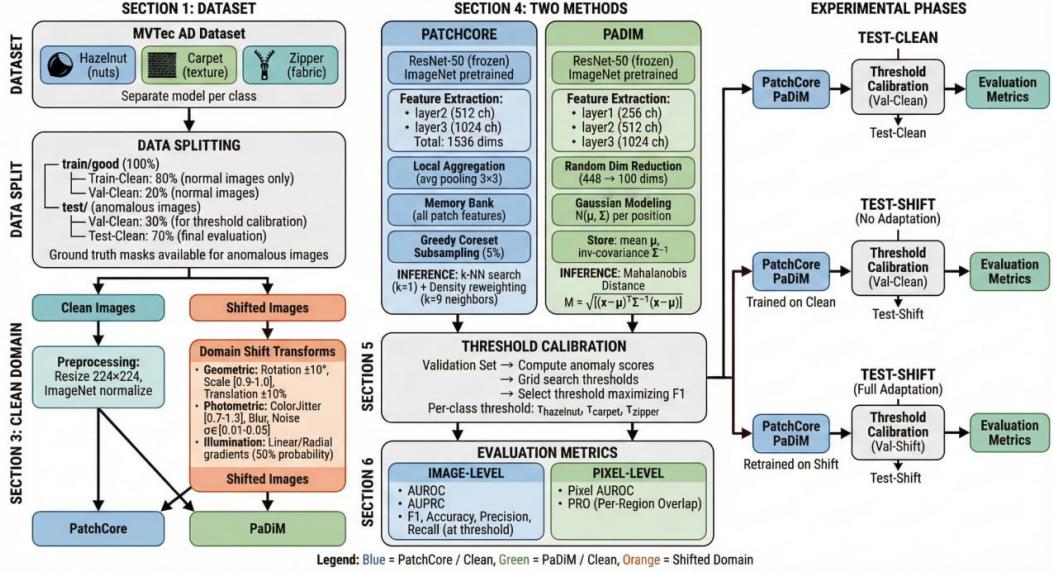


Figure 1: Overview of the experimental pipeline. The workflow encompasses data preparation (splitting and domain shift generation), model training (PatchCore memory bank construction and PaDiM Gaussian fitting), threshold calibration on validation sets, and evaluation under multiple scenarios (clean domain, shifted domain with/without adaptation, and model-unified setting).

2.1. Embedding-Based Methods

PaDiM. Defard *et al.* [2] propose *Patch Distribution Modeling* (PaDiM), which extracts multi-scale features from intermediate layers of a pre-trained ResNet and models the distribution of patch embeddings with a multivariate Gaussian at each spatial location. At inference, anomaly scores are computed via the Mahalanobis distance between a test patch and its corresponding Gaussian. Because it stores only statistical summaries (mean and covariance) rather than individual embeddings, PaDiM is highly memory-efficient and fast at inference. However, the Gaussian assumption may be violated when the underlying distribution is non-unimodal, limiting its capacity to capture complex normal variations.

PatchCore. Roth *et al.* [7] present **PatchCore**, which constructs a *memory bank* of locally-aware patch embeddings extracted from a pre-trained network. To reduce storage and retrieval costs, the memory bank is subsampled using greedy coresset selection. At test time, each image patch is compared to its nearest neighbor in the memory bank, and the maximum distance across patches determines the image-level score. PatchCore achieves near-perfect AUROC on MVTec AD, demonstrating excellent localization thanks to its explicit patch-level comparison.

Efficiency vs. Accuracy Trade-off. Both methods exploit transferable features from ImageNet-pretrained back-

bones, enabling training without gradient-based optimization. PaDiM’s parametric model yields constant-time inference independent of dataset size, whereas PatchCore’s non-parametric memory bank scales with the number of stored embeddings (mitigated by coresset selection). Empirically, PatchCore tends to achieve higher detection and localization accuracy at the cost of increased memory footprint and inference time. In our experiments we quantify this trade-off across multiple categories and domain conditions.

2.2. Unified Multi-Class Anomaly Detection

Conventional anomaly detection follows a *one-model-per-class* paradigm: a separate model is trained for each product category. This approach becomes impractical when the number of categories grows, motivating research into *unified* models that handle multiple classes simultaneously.

Identical Shortcut Problem. You *et al.* [8] introduce **UniAD**, a transformer-based architecture that reconstructs features of different object classes with a single model. A key contribution is the identification of the *identical shortcut* problem: when a model is trained on multiple classes, it may learn to reconstruct any input identically, thereby losing the ability to detect anomalies. UniAD addresses this by injecting neighbour-masked attention and layer-wise query decoders. This insight motivates careful evaluation of unified settings in our work.

Model-Unified vs. Absolute-Unified. Heo and Kang [5] further distinguish between *Model-Unified* and *Absolute-Unified* settings. In the Model-Unified setting, a single model is trained on pooled normal data from all categories, but separate, per-class thresholds are calibrated at inference. In the stricter Absolute-Unified setting, a single global threshold is used regardless of class. The latter is significantly harder because classes may have different “normal” score distributions. Our experiments adopt the Model-Unified paradigm to study how PatchCore and PaDiM behave when exposed to heterogeneous training data, and whether the identical shortcut manifests in embedding-based methods.

2.3. Robustness to Domain Shift

Despite strong benchmark performance, the robustness of anomaly detection models to distribution shift remains under-explored. Real-world deployment often involves variations in lighting, sensor noise, or imaging conditions that differ from the training environment. The MVTec AD 2 dataset [4] explicitly addresses this gap by introducing controlled illumination and pose variations, revealing substantial performance degradation for methods trained on standard MVTec AD. Motivated by this, we construct a synthetic *MVTec-Shift* domain using photometric and geometric augmentations and systematically study the drop in detection and localization accuracy. We further investigate whether simple threshold recalibration on shifted validation data can recover performance without retraining feature extractors, offering practical insights for deployment scenarios with limited access to shifted anomalous samples.

3. Methodology

This section presents the formal problem setup, the domain shift simulation strategy, and the technical details of the two anomaly detection methods employed in our study: PatchCore and PaDiM.

3.1. Problem Formulation

We consider the one-class classification paradigm for visual anomaly detection. Let $\mathcal{D}_{\text{train}} = \{x_i\}_{i=1}^N$ denote a training set consisting exclusively of nominal (defect-free) images, where $x_i \in \mathbb{R}^{H \times W \times 3}$. The objective is to learn a scoring function $s : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}$ that assigns low scores to normal samples and high scores to anomalous ones. For pixel-level localization, the goal extends to producing a spatial anomaly map $\mathcal{A} : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^{H \times W}$, where $\mathcal{A}(x)_{i,j}$ indicates the anomaly intensity at pixel (i, j) .

Preprocessing. All images undergo a deterministic preprocessing pipeline before being fed to the feature extractor. Images are resized to 224×224 pixels using bilinear

interpolation, converted to tensors with values in $[0, 1]$, and normalized using ImageNet statistics:

$$\mu = (0.485, 0.456, 0.406), \quad \sigma = (0.229, 0.224, 0.225) \quad (1)$$

where normalization is applied channel-wise. Ground truth masks, when available, are resized using nearest-neighbor interpolation to preserve binary values and binarized at threshold 0.5.

3.2. Domain Shift Simulation (MVTec-Shift)

To evaluate robustness under distribution shift, we construct a synthetic shifted domain by applying controlled transformations to the clean MVTec AD images. This approach is inspired by the acquisition variations introduced in MVTec AD 2 [4], which include illumination changes and pose perturbations. Our transformation pipeline comprises three categories, applied with synchronized random seeds to ensure reproducibility:

Geometric Transforms. Applied to both images and ground truth masks to maintain spatial consistency:

- **Rotation:** uniformly sampled from $[-10^\circ, +10^\circ]$
- **Scale:** uniformly sampled from $[0.9, 1.0]$ (slight zoom-out)
- **Translation:** uniformly sampled within $\pm 10\%$ of image dimensions

Photometric Transforms. Applied only to images (not masks) to simulate sensor and illumination variations:

- **Color Jitter:** brightness, contrast, and saturation factors sampled uniformly from $[0.7, 1.3]$
- **Gaussian Blur:** kernel size $\in \{3, 5\}$, $\sigma \in [0.1, 2.0]$, applied with probability 0.5
- **Gaussian Noise:** additive noise with $\sigma \in [0.01, 0.05]$, applied with probability 0.5

Illumination Gradients. To simulate non-uniform industrial lighting conditions (*e.g.*, spotlight effects from MVTec AD 2 scenarios such as Fabric and Wall Plugs), we apply smooth illumination gradients with probability 0.5:

- **Linear gradients:** simulate side lighting by darkening one edge (left, right, top, or bottom) with strength factor sampled from $[0.4, 0.7]$
- **Radial gradients:** simulate center/edge illumination variations
- Gradients are smoothed with a Gaussian filter ($\sigma = 80$) for natural transitions

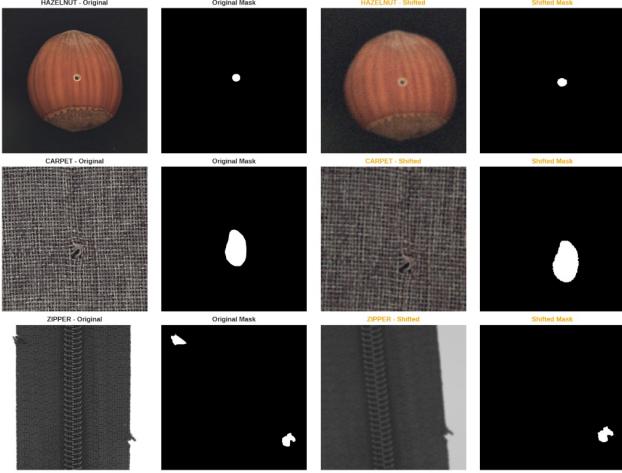


Figure 2: Comparison of clean (left) and shifted (right) samples from MVTec AD. The shifted domain includes photometric degradation, geometric perturbations, and non-uniform illumination to simulate realistic industrial variations.

These transformations collectively simulate realistic sensor degradation and environmental changes encountered in industrial deployment, providing a controlled way to evaluate model robustness.

3.3. PatchCore

PatchCore [7] is a memory-based anomaly detection method that achieves state-of-the-art results through efficient storage and retrieval of nominal patch embeddings. We implement PatchCore following the original paper with hyperparameters verified against our configuration.

Feature Extraction. We employ a ResNet-50 backbone pre-trained on ImageNet with frozen weights. Features are extracted from two intermediate layers, `layer2` (512 channels, 28×28 spatial resolution for 224×224 input) and `layer3` (1024 channels, 14×14 resolution), to capture both fine-grained and semantic information.

Local Neighborhood Aggregation. To incorporate spatial context into each patch embedding, we perform average pooling over a $p \times p$ local neighborhood:

$$f_{\text{agg}}^{(h,w)} = \frac{1}{p^2} \sum_{(i,j) \in \mathcal{N}_p(h,w)} f^{(i,j)} \quad (2)$$

where $\mathcal{N}_p(h,w)$ denotes the $p \times p$ patch centered at position (h,w) and $f^{(i,j)}$ is the feature vector at that location. We use $p = 3$ in our implementation. Features from `layer3` are bilinearly upsampled to match the spatial resolution of

`layer2`, and the two feature maps are concatenated along the channel dimension, yielding a combined feature dimension of $512 + 1024 = 1536$ per patch.

Coreset Subsampling. Storing all patch embeddings from the training set would result in an impractically large memory bank. Following Roth *et al.* [7], we apply greedy coresnet subsampling to select a representative subset. The algorithm iteratively selects patches that maximize the minimum distance to already selected patches (furthest-point sampling):

$$m^* = \arg \max_{m \in \mathcal{M} \setminus \mathcal{C}} \min_{c \in \mathcal{C}} \|m - c\|_2 \quad (3)$$

where \mathcal{M} is the full memory bank and \mathcal{C} is the current coresnet. To accelerate distance computations, we apply Johnson-Lindenstrauss random projection to reduce dimensionality to 128 before sampling. We retain $\rho = 5\%$ of total patches (*i.e.*, `coreset_sampling_ratio = 0.05`), which provides an optimal balance between coverage and efficiency as validated on the clean domain. In following Section 5, we will show the effect of this hyperparameter on the performance of PatchCore.

Anomaly Scoring. At inference time, each test image x is decomposed into a collection of locally aware patch features $\mathcal{P}(x)$. Each patch embedding $q \in \mathcal{P}(x)$ is compared against the coresnet memory bank \mathcal{C} using a Nearest Neighbor (NN) search, implemented via the FAISS [6] library for high-performance indexing.

The raw anomaly score $s_{\text{raw}}(q)$ is defined as the L^2 distance to its nearest neighbor in the coresnet:

$$s_{\text{raw}}(q) = \min_{m \in \mathcal{C}} \|q - m\|_2 \quad (4)$$

To enhance robustness against noise and outliers, we adopt a density-based reweighting scheme inspired by Roth *et al.* [7]. Specifically, we retrieve the k nearest neighbors of q from the coresnet (where $k = 9$ in our implementation) and compute a weighting factor $w(q)$ that penalizes samples whose nearest neighbor is significantly more distant than the rest of the local neighborhood:

$$w(q) = 1 - \frac{\exp(-d_1)}{\sum_{i=1}^k \exp(-d_i)} \quad (5)$$

where d_i represents the L^2 distance to the i -th nearest neighbor. The final anomaly score for each patch is then $s(q) = s_{\text{raw}}(q) \cdot w(q)$.

The overall image-level anomaly score $S(x)$ is determined by the maximum patch-level score, reflecting the assumption that an image is anomalous if it contains at least one anomalous patch:

$$S(x) = \max_{q \in \mathcal{P}(x)} s(q) \quad (6)$$

Anomaly Localization. For defect localization, a spatial anomaly map is generated by reshaping the patch scores $s(q)$ to their respective grid positions. The map is then upsampled to the original input resolution via bilinear interpolation to produce the final pixel-level prediction.

3.4. PaDiM

PaDiM [2] (Patch Distribution Modeling) is a distribution-based anomaly detection method that models normal patch embeddings with multivariate Gaussian distributions. Our implementation uses the anomalib library’s native PadimModel.

Feature Extraction. PaDiM extracts multi-scale features from three ResNet-50 layers: layer1, layer2, and layer3. Features are concatenated and projected to a common spatial resolution, yielding a high-dimensional embedding at each spatial position.

Dimensionality Reduction. To mitigate the curse of dimensionality and reduce computational cost, PaDiM employs random feature selection. From the concatenated feature vector, a random subset of $d = 100$ dimensions is selected (configured via `n_features = 100`). This selection is performed once at initialization and fixed for the lifetime of the model.

Gaussian Modeling. For each spatial position (h, w) , PaDiM estimates a multivariate Gaussian distribution from the training embeddings:

$$\mathcal{N}_{h,w}(\mu_{h,w}, \Sigma_{h,w}) \quad (7)$$

where $\mu_{h,w} \in \mathbb{R}^d$ is the mean embedding and $\Sigma_{h,w} \in \mathbb{R}^{d \times d}$ is the covariance matrix. This position-specific modeling captures the expected appearance at each location, accounting for spatial structure in the images.

Anomaly Scoring via Mahalanobis Distance. At inference, the anomaly score for a test patch at position (h, w) is the Mahalanobis distance to the learned Gaussian:

$$\mathcal{M}(x_{h,w}) = \sqrt{(x_{h,w} - \mu_{h,w})^\top \Sigma_{h,w}^{-1} (x_{h,w} - \mu_{h,w})} \quad (8)$$

where $\Sigma_{h,w}^{-1}$ is the inverse covariance (precision) matrix. The image-level score is computed as the maximum Mahalanobis distance across all positions, and the spatial anomaly map is obtained by upsampling the distance values to the original image resolution.

Comparison with PatchCore. Unlike PatchCore’s non-parametric memory bank, PaDiM stores only statistical summaries (mean and covariance per position), resulting in constant memory and inference time independent of training set size. However, the Gaussian assumption may not hold for all normal variations, potentially limiting its discriminative power compared to PatchCore’s explicit nearest-neighbor matching.

4. Experimental Setup

This section describes the dataset, evaluation protocol, and experimental scenarios employed to validate the proposed anomaly detection methods. We detail the data splits, threshold calibration procedure, evaluation metrics, and the different domain settings under which models are tested.

4.1. Dataset and Splits

MVTec AD. We conduct experiments on the MVTec Anomaly Detection (MVTec AD) dataset [1], a comprehensive benchmark for unsupervised anomaly detection in industrial inspection. We select three representative categories: **Hazelnut** (object), **Carpet** (texture), and **Zipper** (mixed), covering the diversity of defect types and visual structures present in the dataset. Each category contains defect-free (normal) training images and a test set with both normal and anomalous samples, accompanied by pixel-level ground truth masks for localization evaluation.

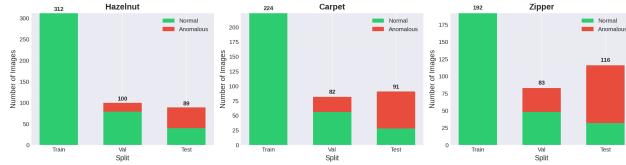
Data Splitting Protocol. Following the one-class classification paradigm, we construct separate train/validation/test splits for each category. The original MVTec AD structure provides a `train/good` folder (normal images only) and `test/` folder containing both normal (`good`) and anomalous images organized by defect type. We partition these as follows:

- **Train-clean:** 80% of images from `train/good`, used exclusively to build the nominal representation (*i.e.*, memory bank for PatchCore, Gaussian parameters for PaDiM).
- **Val-clean:** Remaining 20% of `train/good` (normal) plus 30% of the anomalous images from `test/<defect>`, used for threshold calibration and hyperparameter tuning.
- **Test-clean:** All remaining normal images from `test/good` and all remaining anomalous images from `test/<defect>`, with ground truth masks. Reserved strictly for final evaluation.

The 80/20 train/validation ratio ensures sufficient training data for robust feature extraction while retaining a representative validation set. Crucially, **anomalous samples**

Table 1: Dataset split statistics for the three selected MVTec AD categories. Val-clean includes both normal (from train/good) and anomalous samples (30% from test/) for threshold calibration.

Category	Train-clean		Val-clean		Test-clean	
	Normal	Anom.	Normal	Anom.	Normal	Anom.
Hazelnut	312	0	79	21	40	49
Carpet	224	0	56	26	28	63
Zipper	192	0	48	36	32	84
Total	728	0	183	83	100	196



are included in the validation set (at 30% of available anomalies) to enable threshold calibration that accounts for the score distribution of both classes. This design follows established practice in anomaly detection, where access to a small number of labeled validation anomalies is realistic and necessary for setting decision thresholds [7].

All splits are generated reproducibly using a fixed random seed (42), with split indices stored in JSON format to ensure consistency across experiments.

Split Statistics. Table 1 summarizes the dataset composition for each category and split.

4.2. Evaluation Protocol

Threshold Selection. For each class and method, we calibrate a decision threshold on the validation set using an F1-optimal grid search. Specifically, we:

1. Compute image-level anomaly scores for all validation samples (normal + anomalous).
2. Search over 1000 uniformly spaced thresholds between the minimum and maximum validation scores.
3. Select the threshold τ^* that maximizes the F1 score on the validation set:

$$\tau^* = \arg \max_{\tau} F1(\tau) = \arg \max_{\tau} \frac{2 \cdot \text{Prec}(\tau) \cdot \text{Rec}(\tau)}{\text{Prec}(\tau) + \text{Rec}(\tau)} \quad (9)$$

Image-Level Metrics. We evaluate detection performance using both threshold-independent and threshold-dependent metrics:

- **AUROC:** Area Under the Receiver Operating Characteristic curve, measuring the model’s ability to rank anomalous images higher than normal ones across all thresholds.

- **AUPRC:** Area Under the Precision-Recall Curve, particularly informative for imbalanced datasets where anomalies are the minority class.

- **F1 Score:** Harmonic mean of precision and recall at the calibrated threshold τ^* .

Pixel-Level Metrics. For localization evaluation, we compute:

- **Pixel AUROC:** ROC-AUC computed over all pixels from the entire test set, measuring discrimination between defective and non-defective pixels.

- **PRO (Per-Region Overlap):** Following the protocol of Bergmann *et al.* [1], we compute the area under the per-region overlap curve integrated up to a false positive rate of 0.3. PRO penalizes predictions that miss or only partially cover ground truth anomaly regions.

All metrics are reported per-class and as macro-averages across the three categories to provide both granular and aggregate performance views.

4.3. Evaluation Scenarios

We evaluate models under four distinct scenarios to assess both baseline performance and robustness to domain shift:

1. **Test-Clean (Baseline):** Models trained on Train-clean are evaluated on Test-clean using thresholds calibrated on Val-clean. This represents the idealized scenario where training and test conditions are matched.
2. **Test-Shift (No Adaptation):** Clean-trained models are evaluated on the shifted test set (Test-shift) using the *same thresholds calibrated on Val-clean*. This scenario measures pure performance degradation without any adaptation mechanism.
3. **Test-Shift (Threshold-Only Adaptation):** Models remain trained on Train-clean (no retraining), but thresholds are *re-calibrated on Val-shift*. This isolates the contribution of threshold adaptation.
4. **Test-Shift (Full Adaptation):** Models are *retrained from scratch on Train-shift* and thresholds are calibrated on Val-shift. This represents the upper bound of adaptation performance.

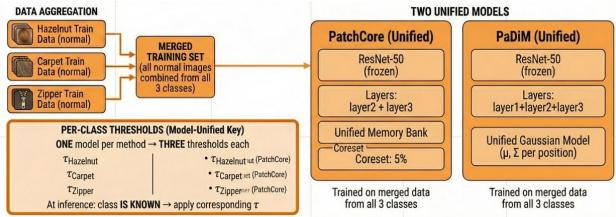


Figure 3: Overview of the Model-Unified setting. Normal training data from all three categories (Hazelnut, Carpet, Zipper) is merged into a single training set. One unified PatchCore model and one unified PaDiM model are trained on this pooled data. At inference time, class identity is known and per-class thresholds are applied accordingly.

4.3.1 Global Model (Model-Unified Setting)

In addition to per-class models, we train a **single global model** on pooled normal data from all three categories (Figure 3). This setup follows the *Model-Unified* setting as defined by Heo and Kang [5]:

- **Training:** One shared model is fitted on the concatenated training data from Hazelnut, Carpet, and Zipper.
- **Inference:** Class identity is assumed *known* at test time.
- **Thresholds:** *Per-class thresholds* are calibrated on each class’s validation set using the global model’s predictions.

This is distinct from the more challenging *Absolute-Unified* setting [3], where class identity is unknown and a single global threshold must be used. The Model-Unified setting is specifically designed to investigate the “**identical shortcut**” problem [8]: when a single model is trained on multiple classes, it may learn to distinguish between classes rather than between normal and anomalous samples within each class. We analyze this phenomenon by examining cross-class confusion patterns and feature space separation (see Section 5).

4.4. Methods Under Evaluation

All evaluation scenarios described above are conducted on both PatchCore and PaDiM to enable direct comparison. PaDiM serves as a baseline method, providing a reference point for assessing PatchCore’s performance across all experimental conditions. The two methods share the same preprocessing pipeline, data splits, and evaluation metrics, differing only in their core anomaly scoring mechanisms.

Coreset Ablation (PatchCore Only). The coresnet ratio ablation study (Section 4.5) is conducted exclusively on

PatchCore, as coresnet subsampling is not applicable to PaDiM’s parametric Gaussian modeling.

4.5. Ablation Study: Coresnet Ratio

The coresnet sampling ratio ρ is a critical hyperparameter for PatchCore, controlling the trade-off between memory bank coverage and computational efficiency. Following the range suggested by Roth *et al.* [7] (1–10%), we conduct an ablation study to evaluate the impact of this parameter.

Experimental Design. We train three PatchCore variants with coresnet ratios $\rho \in \{1\%, 5\%, 10\%\}$ on the shifted domain (Train-shift), calibrate thresholds on Val-shift, and evaluate on Test-shift. This configuration isolates the effect of coresnet size by keeping the domain fixed (shifted), avoiding confounding factors from domain mismatch.

For each configuration, we record:

- **Detection metrics:** AUROC, AUPRC, F1, and Accuracy on Test-shift.
- **Efficiency metrics:** Memory bank size (number of retained patches) and training time.

Smaller coresnet ratios reduce memory usage and accelerate nearest neighbor search during inference, but may underfit by discarding representative patches. Larger ratios improve coverage but increase computational cost quadratically during coresnet selection. The ablation quantifies this trade-off empirically, justifying the selection of $\rho = 0.05$ as the default configuration for all other experiments. Results are reported in Section 5.

4.6. Implementation Details

All experiments are executed on a single NVIDIA T4 GPU provided by Google Colab. The complete codebase, including data preprocessing, model training, and evaluation pipelines, is implemented as a collection of Jupyter notebooks designed for reproducibility. All random operations use a fixed seed of 42.

For implementation specifics of PatchCore (backbone, feature layers, coresnet subsampling, score reweighting) and PaDiM (layers, dimensionality reduction, Gaussian modeling), we refer to Section 3. The *anomalib* library is used for PaDiM, while PatchCore is implemented from scratch following Roth *et al.* [7].

5. Results and Discussion

This section presents the experimental results across the evaluation scenarios described in Section 4.3.

5.1. Clean Domain Performance

Table 2 reports the detection and localization performance of PatchCore and PaDiM on Test-clean, using thresholds calibrated on Val-clean. Both methods achieve strong baseline performance, validating the experimental setup and implementations.

Table 2: Clean domain performance (Test-clean). Best per metric in **bold**.

Method	Class	AUC	PRC	F1	pAUC	PRO
PatchCore	Hazel.	1.00	1.00	1.00	.988	.855
	Carpet	.973	.991	.951	.988	.836
	Zipper	.982	.994	.949	.986	.833
	Avg	.985	.995	.966	.987	.841
PaDiM	Hazel.	.971	.975	.867	.963	.797
	Carpet	.959	.985	.923	.984	.534
	Zipper	.862	.942	.895	.972	.797
	Avg	.930	.967	.895	.973	.709

Detection Performance. PatchCore demonstrates superior image-level detection, achieving near-perfect performance on Hazelnut (AUROC = 1.000, F1 = 1.000) and strong results across all categories (macro-average AUROC = 0.985, F1 = 0.966). PaDiM, while competitive, shows a consistent performance gap, particularly on Zipper (AUROC = 0.862 vs. 0.982 for PatchCore). This difference in macro-average AUROC (0.985 vs. 0.930) confirms the advantage of PatchCore’s non-parametric approach.

Localization Performance. Both methods achieve excellent pixel-level AUROC (PatchCore: 0.987, PaDiM: 0.973), indicating strong discrimination between defective and normal pixels. However, the PRO metric reveals a more significant gap: PatchCore achieves 0.841 compared to PaDiM’s 0.709. The PRO metric, which penalizes predictions that only partially cover defect regions, highlights PatchCore’s more precise spatial localization. Notably, PaDiM’s PRO on Carpet (0.534) is substantially lower than other categories, suggesting difficulty with fine-grained texture defects.

Method Comparison. The performance gap between PatchCore and PaDiM can be attributed to their fundamentally different modeling approaches:

- **Non-parametric vs. Parametric:** PatchCore’s memory bank approach (Section 3.3) enables exact matching, while PaDiM’s Gaussian assumption (Section 3.4) may not capture complex feature distributions.
- **Density Reweighting:** PatchCore applies density-based reweighting (see Section 3.3) to anomaly scores,

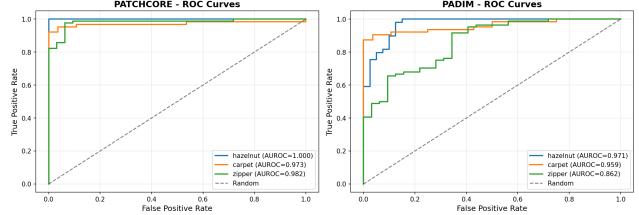


Figure 4: ROC curves on Test-clean for PatchCore and PaDiM across the three categories. PatchCore achieves near-perfect AUROC on Hazelnut and consistently outperforms PaDiM.

reducing the influence of outliers in the memory bank. PaDiM lacks this mechanism, making it more sensitive to covariance estimation errors.

- **Trade-offs:** PaDiM offers lower memory footprint and faster inference (no nearest-neighbor search required), making it suitable for resource-constrained deployments. However, for maximum detection accuracy, PatchCore’s memory bank approach proves more effective.

Figure 4 shows the ROC curves for both methods across all categories. PatchCore consistently achieves higher true positive rates at low false positive rates, reflecting its superior ranking ability. The confusion matrices (Figure 5) reveal that PatchCore’s errors are primarily false positives (8 total across all classes), while PaDiM exhibits a higher false negative rate (9 missed anomalies on Carpet alone). Figure 6 provides qualitative evidence of PatchCore’s localization quality on Hazelnut, showing accurate heatmap activations on defective regions while producing uniform low-intensity maps for normal samples.

5.2. Impact of Domain Shift

We now evaluate the robustness of clean-trained models when tested on the shifted domain (Test-shift) *without any adaptation*. This scenario uses the same models trained on Train-clean and the same thresholds calibrated on Val-clean, isolating the pure impact of distribution shift.

Experimental Setup. The domain shift transformations applied to Test-shift include geometric perturbations (rotation $\pm 10^\circ$, scale variations, translations) and photometric changes (color jitter, Gaussian blur, noise, non-uniform illumination). See Section 3.2 for details.

Table 3 presents the performance degradation compared to the clean domain baseline.

Overall Degradation. Both methods suffer substantial performance drops when evaluated on shifted data with

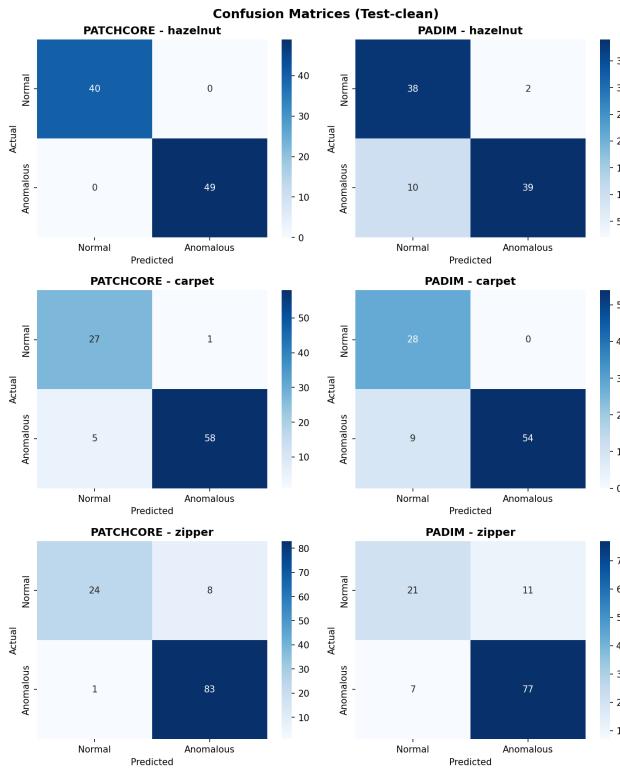


Figure 5: Confusion matrices on Test-clean at calibrated thresholds. PatchCore shows fewer false negatives, while PaDiM exhibits higher miss rates, particularly on Carpet and Hazelnut.

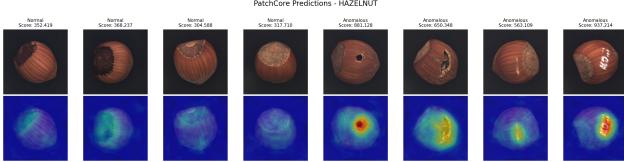


Figure 6: Qualitative results for PatchCore on Hazelnut (Test-clean). Top row: input images with anomaly scores. Bottom row: predicted heatmaps. Normal samples (left) receive low scores with uniform heatmaps, while anomalous samples (right) show high scores and accurate defect localization.

clean-domain thresholds. PatchCore’s macro-average AUROC decreases from 0.985 to 0.892 ($\Delta = -0.093$), while PaDiM experiences a more severe drop from 0.930 to 0.684 ($\Delta = -0.246$). This larger degradation for PaDiM indicates that Gaussian parametric models are more sensitive to distribution shift than memory-bank approaches.

Threshold Invalidity and Specificity Collapse. The most striking observation is the **catastrophic collapse of**

Table 3: Domain shift degradation (No Adaptation). Macro-average metrics comparing Clean vs. Shift. $\Delta = \text{Shift} - \text{Clean}$.

Method	AUROC			Specificity		
	Clean	Shift	Δ	Clean	Shift	Δ
PatchCore	.985	.892	-.093	.905	.223	-.682
PaDiM	.930	.684	-.246	.869	.183	-.686

specificity, dropping from ~ 0.90 to ~ 0.20 for both methods. For Zipper, *both methods achieve exactly 0% specificity*, meaning every normal image is misclassified as anomalous. This occurs because the domain shift causes anomaly scores to increase systematically across all samples, pushing normal images above the threshold calibrated on clean data.

Figure 7 visualizes this effect: the bottom-left quadrant (True Negatives) is nearly empty, while false positives dominate. Despite maintaining high recall (PatchCore: 1.000, PaDiM: 0.972), the models become useless for practical deployment due to excessive false alarms.

The Rotation Sensitivity. A good finding is PaDiM’s disproportionate sensitivity to domain shift, particularly on Hazelnut ($\text{AUROC}: 0.971 \rightarrow 0.605, \Delta = -0.366$) and Zipper ($\text{AUROC}: 0.862 \rightarrow 0.555, \Delta = -0.307$). We hypothesize this is due to the ± 10 rotations included in our shift transformations:

- **PaDiM’s Position-Specific Modeling:** PaDiM estimates a separate Gaussian distribution for each spatial position (i, j) in the feature map [2]. When an image is rotated, features from position (i, j) in the test image correspond to a different physical location than in training. This spatial misalignment causes features to be compared against incorrect reference distributions, artificially inflating Mahalanobis distances.

- **PatchCore’s Global Memory Bank:** In contrast, PatchCore’s memory bank is position-agnostic: it stores patch embeddings without spatial indexing. During inference, each test patch is matched to its globally nearest neighbor regardless of position. Small rotations merely shuffle which memory bank patches are retrieved, without fundamentally breaking the matching process.

This *rotation sensitivity* explains why PaDiM, despite being competitive on clean data, degrades more severely under geometric transformations. The effect is less pronounced on Carpet, a texture category with inherent translation invariance, where PaDiM’s degradation is more modest ($\Delta = -0.065$).

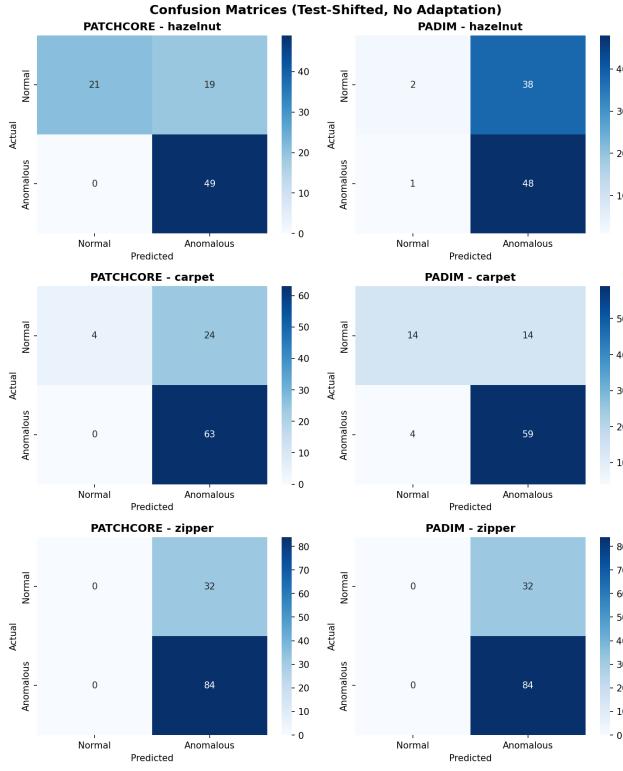


Figure 7: Confusion matrices on Test-shift (No Adaptation). Both methods exhibit near-zero specificity, with Zipper showing complete failure to identify normal samples. The threshold calibrated on Val-clean is no longer valid under domain shift.

Pixel-Level Degradation. Pixel-level metrics also degrade under shift, though less dramatically than image-level specificity. PatchCore’s Pixel AUROC drops from 0.987 to 0.922 ($\Delta = -0.065$), while PaDiM drops from 0.973 to 0.827 ($\Delta = -0.146$). The PRO metric shows larger degradation (PatchCore: 0.841 \rightarrow 0.645; PaDiM: 0.709 \rightarrow 0.527), indicating that localization precision suffers more than pixel-level discrimination under distribution shift.

Implications. These results underscore the critical importance of threshold recalibration and/or model adaptation when deploying anomaly detection systems in environments with even modest acquisition variations. The next sections will examine whether threshold-only adaptation (Section ??) or full model retraining (Section ??) can recover the lost performance.

5.3. Adaptation Strategies

Having established the severity of domain shift degradation in Section 5.2, we now evaluate two adaptation strategies: (1) *threshold-only adaptation*, which recalibrates

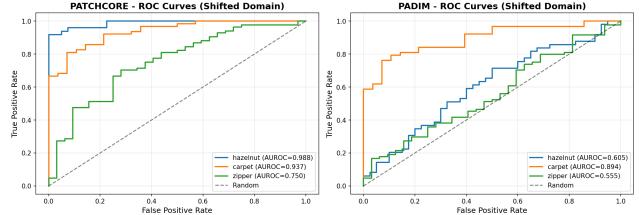


Figure 8: ROC curves on Test-shift (No Adaptation). PatchCore maintains reasonable ranking ability (macro AUROC = 0.892), while PaDiM’s curves approach the diagonal (random classifier) for Hazelnut and Zipper.

the decision threshold on Val-shift without modifying the model, and (2) *full adaptation*, which retrains the model on Train-shift. Table 4 summarizes the macro-averaged results across all scenarios.

Table 4: Comparison of adaptation strategies (macro-average). AUROC and Pixel AUROC are threshold-independent; values in gray indicate no change from Shift-NoAdapt. Best per-metric in **bold**.

Scenario	PatchCore			PaDiM		
	AUC	pAUC	PRO	AUC	pAUC	PRO
Clean (Baseline)	.985	.987	.841	.930	.973	.709
Shift-NoAdapt	.892	.922	.645	.684	.827	.527
Threshold-Only	.892	.922	.645	.684	.827	.527
Full Adaptation	.961	.966	.800	.885	.922	.629

Table 5: Threshold-dependent metrics (macro-average). These metrics **change** with threshold recalibration, unlike Table 4.

Scenario	PatchCore			PaDiM		
	F1	Spec	Acc	F1	Spec	Acc
Clean (Baseline)	.966	.905	.952	.895	.869	.870
Shift-NoAdapt	.839	.223	.749	.806	.183	.696
Threshold-Only	.850	.830	.814	.800	.183	.700
Full Adaptation	.921	.911	.896	.858	.824	.819

Threshold-Only Adaptation. An important observation from Table 4 is that AUROC, Pixel AUROC, and PRO remain *unchanged* between Shift-NoAdapt and Threshold-Only scenarios. This is expected: these metrics are **threshold-independent**, measuring the model’s ranking ability by integrating over all possible thresholds. Since the underlying model and its anomaly scores are identical in both cases, only threshold-dependent metrics (F1, Specificity, Accuracy) can improve, as shown in Table 5.

The practical benefit of threshold recalibration is the recovery of specificity (Table 5). For PatchCore, specificity increases from 22.3% to 83.0% by adjusting the threshold from 227.5 to 441.8 (the optimal F1 point on Val-shift). This eliminates the false positive explosion observed in Section 5.2 at *zero computational cost*: no retraining or feature extraction required.

However, threshold-only adaptation has fundamental limitations: it cannot improve detection accuracy beyond what the existing score distribution allows. For PaDiM, where the score distributions of normal and anomalous samples overlap significantly under shift (particularly for Zipper), recalibration yields minimal benefit (specificity remains at 18.3%).

Full Adaptation. Full model retraining on Train-shift provides substantial performance recovery. PatchCore’s AUROC increases from 0.892 to 0.961 ($\Delta = +6.9$ pp), while PaDiM shows even larger recovery from 0.684 to 0.885 ($\Delta = +20.1$ pp). The improvement mechanism differs between methods:

- **PatchCore:** The memory bank is rebuilt using patch embeddings from shifted normal images. This allows the nearest-neighbor matching to account for the new appearance variations introduced by domain shift transformations. The coresnet subsampling 3 ensures efficient coverage of the expanded feature manifold.
- **PaDiM:** New Gaussian distributions are estimated at each spatial position using shifted training data. This recalibrates both the mean vectors μ_{ij} and covariance matrices Σ_{ij} to match the shifted domain, reducing the systematic Mahalanobis distance inflation observed without adaptation.

The effect of these adaptation strategies on score distributions is visualized in Figure 9, focusing on the Zipper category where the contrast is most evident. Under threshold-only adaptation (top), the score distributions of normal (blue) and anomalous (red) samples remain heavily overlapped for both methods, with PaDiM showing near-complete overlap that explains the persistent 18.3% specificity. In contrast, full adaptation (bottom) re-separates the distributions by learning new feature representations, restoring clearer boundaries between normal and anomalous samples.

Figure 10 shows the ROC curves after full adaptation. PatchCore achieves near-perfect detection on Hazelnut (AUROC = 0.998) and strong performance across all categories. The confusion matrices (Figure 11) confirm the recovery of specificity to practical levels (PatchCore: 91.1%, PaDiM: 82.4%). Qualitative comparison (Figure 12) further illustrates the difference: PatchCore produces sharper,

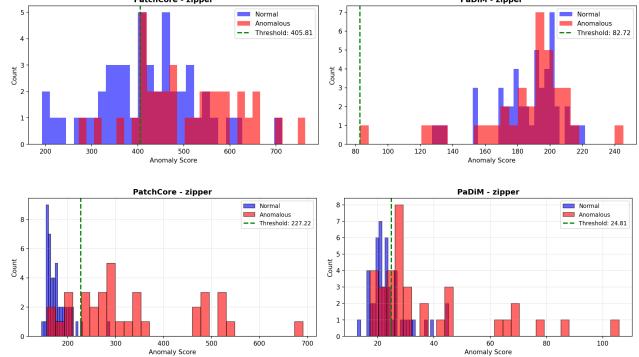


Figure 9: Score distributions on Zipper (Val-shift): **Top:** Threshold-only adaptation: distributions remain overlapped; only the threshold (green dashed) is recalibrated. **Bottom:** Full adaptation: re-training re-separates normal and anomalous distributions, enabling effective classification.

better-localized heatmaps that closely match ground truth masks, while PaDiM’s predictions tend to be more diffuse and occasionally miss subtle defects.

Residual Gap Analysis. Despite full retraining, neither method fully recovers to clean domain performance. PatchCore exhibits a residual gap of -2.4 pp ($0.985 \rightarrow 0.961$), while PaDiM shows -4.5 pp ($0.930 \rightarrow 0.885$). We attribute this to three factors:

1. **Increased intra-class variance:** The shift transformations (rotation, color jitter, blur, illumination gradients) expand the feature distribution of normal images, making the boundary between normal and anomalous less distinct.
2. **Transformation artifacts:** Image interpolation during geometric transformations and border effects introduce subtle patterns that differ from real acquisition variations, potentially confusing the models.
3. **Data quality:** Synthetically augmented data, while useful for robustness, does not fully replicate the natural variations encountered in real domain shifts.

Method Comparison. PatchCore’s smaller residual gap (-2.4 pp vs. -4.5 pp for PaDiM) confirms the advantage of non-parametric approaches under domain shift. The memory bank can flexibly accommodate new appearance variations without making distributional assumptions. In contrast, PaDiM’s position-specific Gaussian modeling [2] remains partially misaligned even after retraining: the ± 10 rotations in our shift protocol cause features at position (i, j) to originate from physically different locations than

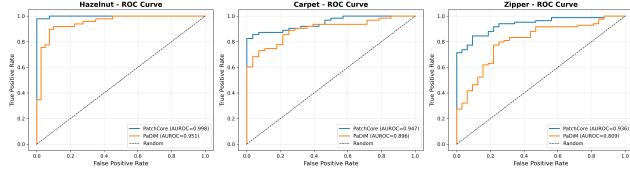


Figure 10: ROC curves after full adaptation on Train-shift. PatchCore achieves $\text{AUROC} \geq 0.936$ on all categories, with near-perfect performance on Hazelnut (0.998). PaDiM shows significant recovery but remains below PatchCore, particularly on Zipper.

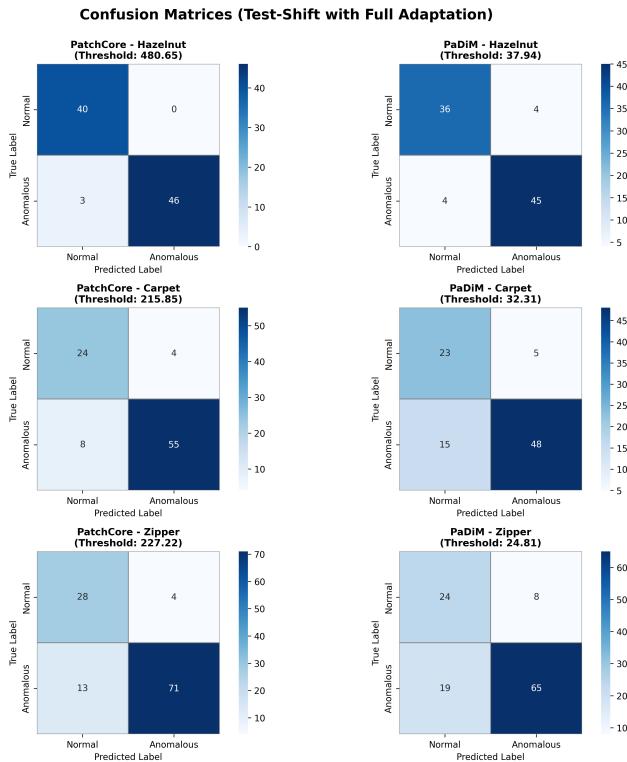


Figure 11: Confusion matrices after full adaptation. Both methods recover practical specificity levels (PatchCore: 91.1%, PaDiM: 82.4%), with PatchCore achieving perfect specificity on Hazelnut (100%).

in training, violating the positional correspondence assumption.

5.4. Model Unified Analysis

The preceding sections evaluate *per-class models*, where each category has its own dedicated model trained exclusively on that class’s normal data. We now investigate the **Model-Unified** setting (see Section 4.3.1 and Section 2.2), where a single model is trained on pooled normal data from all three categories (Hazelnut, Carpet, Zipper). At infer-

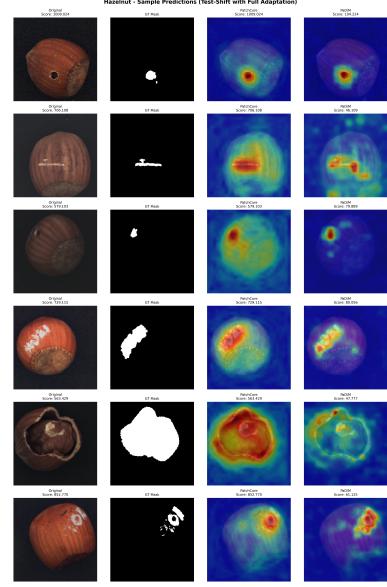


Figure 12: Qualitative comparison on Hazelnut (Test-shift) after full adaptation. Columns: original image, ground truth mask, PatchCore heatmap, PaDiM heatmap. PatchCore produces sharper localization, while PaDiM’s predictions are more diffuse.

ence time, class identity is assumed known and *per-class thresholds* are applied. This setting is valuable for industrial deployments where maintaining separate models for each product variant is impractical.

Quantitative Comparison. Table 6 reports the macro-averaged performance of the unified models across all three categories. PatchCore maintains strong detection ($\text{AUROC} = 0.984$) and localization (Pixel-AUROC = 0.981, PRO = 0.801) capabilities even when trained on heterogeneous data. PaDiM shows competitive detection performance ($\text{AUROC} = 0.934$), but its localization remains weak: the PRO of 0.672 is consistent with its per-class performance (0.709 in Table 2), confirming that PaDiM’s difficulty with precise defect localization is intrinsic to its Gaussian modeling approach rather than a consequence of the unified setting.

Table 6: Unified model performance (macro-averaged across all classes).

Method	AUROC	Pixel-AUROC	PRO
PatchCore	.984	.981	.801
PaDiM	.934	.951	.672

Figure 13 quantifies the AUROC gap between per-class and unified models. PatchCore exhibits negligible degra-

dation (max gap = +0.5% on Carpet), confirming that its memory bank effectively handles multi-class training data. PaDiM shows larger variance: a 1.5% drop on Carpet but a 2.1% *improvement* on Zipper, likely due to the larger training set providing better feature coverage for this challenging category.

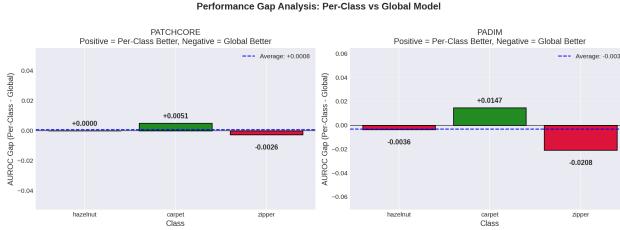


Figure 13: AUROC gap between per-class and unified models. Positive = per-class better, negative = unified better. PatchCore shows negligible gaps (<0.5%), while PaDiM exhibits larger variance.

The “Identical Shortcut” Problem. While per-class thresholds mitigate most issues in the Model-Unified setting, a critical pathology emerges when we analyze *cross-class* behavior. As discussed in Section 2.2, You *et al.* [8] identified the “identical shortcut” problem: in unified settings, models may learn to discriminate between classes rather than between normal and anomalous samples within each class. To test this, we apply each class’s calibrated threshold to normal images from *other* classes and measure the false positive rate (FPR).

Figure 14 shows the cross-class confusion matrices. PatchCore achieves **0% FPR across all class pairs**, demonstrating that its memory bank effectively separates the feature manifolds of different categories. In contrast, PaDiM exhibits a catastrophic failure: when the Zipper threshold ($\tau = 18.05$) is applied to Hazelnut normal images, **55.1% are misclassified as anomalous** (172/312 images). This occurs because Zipper’s lower threshold, optimized for its own score distribution, incorrectly flags Hazelnut’s higher natural scores.

Why PatchCore is More Robust.

- **PatchCore’s Memory Bank:** Stores actual patch embeddings from each class, which naturally occupy distinct regions in feature space. The nearest-neighbor matching ensures that test patches are compared only against visually similar prototypes. When Hazelnut patches query the global memory bank, they find Hazelnut-like neighbors regardless of which threshold is applied: the anomaly score reflects true visual deviation.

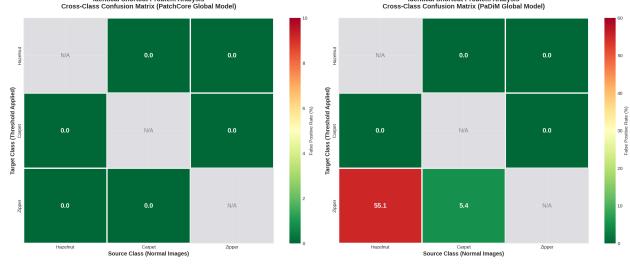


Figure 14: Cross-class confusion matrices showing FPR (%) when applying one class’s threshold to another class’s normal images. **Left:** PatchCore achieves 0% FPR everywhere. **Right:** PaDiM shows severe failure when Zipper’s threshold is applied to Hazelnut (55.1%) and Carpet (5.4%).

- **PaDiM’s Shared Gaussian:** Estimates a single Gaussian distribution per spatial position over *all* classes. This creates overlapping distributions where Hazelnut’s feature vectors, which tend to have higher Mahalanobis distances due to their distinct appearance (nuts vs. fabric), naturally produce higher baseline scores. Zipper’s threshold, calibrated for lower scores, incorrectly triggers.

Practical Implications. These findings have direct implications for industrial deployment:

1. **Per-class thresholds are essential:** Even with a unified model, applying a single global threshold is not advisable. Our Model-Unified results confirm that per-class threshold calibration effectively compensates for score distribution differences.
2. **PatchCore is preferred for multi-class scenarios:** Its memory bank approach inherently preserves class boundaries, avoiding the identical shortcut problem. PaDiM’s Gaussian modeling, while computationally efficient, requires careful attention to prevent cross-class interference.
3. **Scalability consideration:** A single unified PatchCore model with per-class thresholds achieves comparable performance to per-class models (AUROC gap <0.5%) while reducing deployment complexity from 3 models to 1.

6. Conclusion

References

- [1] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtex ad — a comprehensive real-world dataset for unsupervised anomaly detection. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9584–9592, 2019.

- [2] Thomas Defard, Aleksandr Setkov, Angelique Loesch, and Romaric Audigier. Padim: a patch distribution modeling framework for anomaly detection and localization, 2020.
- [3] Jia Guo, Haonan Han, Shuai Lu, Weihang Zhang, and Huiqi Li. Absolute-unified multi-class anomaly detection via class-agnostic distribution alignment, 2024.
- [4] Lars Heckler-Kram, Jan-Hendrik Neudeck, Ulla Scheler, Rebecca König, and Carsten Steger. The mvtec ad 2 dataset: Advanced scenarios for unsupervised anomaly detection, 2025.
- [5] Jaehyuk Heo and Pilsung Kang. Multi-class image anomaly detection for practical applications: Requirements and robust solutions, 2025.
- [6] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus, 2017.
- [7] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards total recall in industrial anomaly detection, 2022.
- [8] Zhiyuan You, Lei Cui, Yujun Shen, Kai Yang, Xin Lu, Yu Zheng, and Xinyi Le. A unified model for multi-class anomaly detection, 2022.