

---

# EXPLAINABILITY-DRIVEN DIMENSIONALITY REDUCTION FOR HYPERSPPECTRAL IMAGING

---

Salma Haidar \*    José Oramas 

Department of Computer Science  
University of Antwerp, imec - IDLab  
Sint-Pietersvliet 7,2000, Belgium  
{salma.haidar,jose.oramass}@uantwerpen.be

## ABSTRACT

Hyperspectral imaging (HSI) provides rich spectral information for precise material classification and analysis; however, its high dimensionality introduces a computational burden and redundancy, making dimensionality reduction essential. We present an exploratory study into the application of post-hoc explainability methods in a model-driven framework for band selection, which reduces the spectral dimension while preserving predictive performance. A trained classifier is probed with explanations to quantify each band’s contribution to its decisions. We then perform deletion–insertion evaluations, recording confidence changes as ranked bands are removed or reintroduced, and aggregate these signals into influence scores. Selecting the highest-influence bands yields compact spectral subsets that maintain accuracy and improve efficiency. Experiments on two public benchmarks (Pavia University and Salinas) demonstrate that classifiers trained on as few as 30 selected bands match or exceed full-spectrum baselines while reducing computational requirements. The resulting subsets align with physically meaningful, highly discriminative wavelength regions, indicating that model-aligned, explanation-guided band selection is a principled route to effective dimensionality reduction for HSI.

**Keywords** Hyperspectral Image Analysis, Band Selection, Explainable Artificial Intelligence, LRP, SHAP, RISE.

## 1 Introduction

Hyperspectral imaging (HSI) is an advanced imaging technology that employs sensors to capture data across hundreds of narrow, contiguous spectral bands spanning a wide range of electromagnetic wavelengths well beyond the visible spectrum [1]. Every pixel in a hyperspectral image records a complete light spectrum, yielding fine-grained spectral data that enables precise material identification and characterisation based on unique chemical and physical properties—capabilities that exceed those of conventional imaging systems. However, high-dimensionality presents significant challenges for analysis and modelling. It increases computational complexity, introduces redundancy, and exacerbates the risk of overfitting. Moreover, data quality is often compromised by noise from sensor imperfections, environmental factors, and external interferences, which can obscure meaningful spectral patterns.

To address these issues, researchers have widely adopted dimensionality reduction techniques [2, 3] aimed at eliminating redundancy, suppressing noise and retaining the the most informative spectral features. These techniques are broadly categorised into feature extraction and feature (band) selection, with different trade-offs in interpretability, efficiency, and downstream performance.

Feature extraction [4, 5] transforms the original high-dimensional data into a lower-dimensional space where new representative feature vectors are selected but can obscure the physical and chemical significance of spectral bands—a critical limitation for applications that require domain-specific validation.

---

\*Corresponding author: salma.haidar@uantwerpen.be

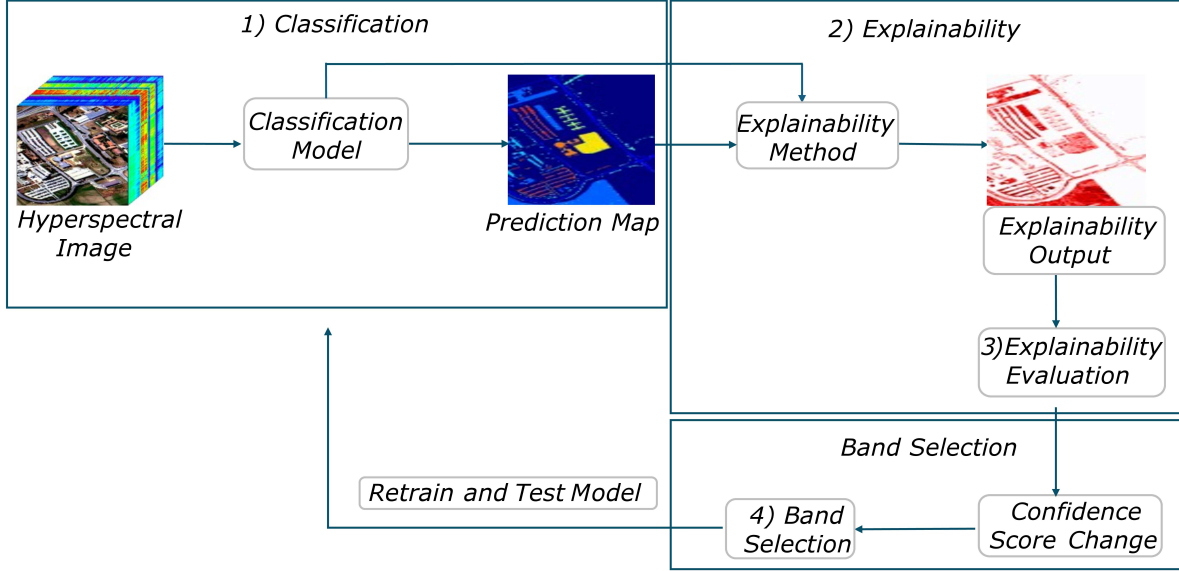


Figure 1: Workflow using explainability-based selection of key spectral bands and model retraining with reduced dimensionality.

In contrast, band selection [6, 7], identifies and retains a subset of the original bands without altering their structure, thereby preserving their physical interpretability. Nonetheless, existing band-selection methods may mistakenly prioritise noisy bands as informative or capture specific spectral clusters, overlooking the broader spectral variability inherent in hyperspectral data.

Meanwhile, there has been a growing emphasis on understanding models decision-making processes [8–10], aiming to make complex models more transparent and interpretable. Explaining why a model makes a specific decision has driven research into model explainability [11–13]. While explainability methods have been widely adopted to justify predictions, their integration into HSI analysis remains limited [14, 15]. Beyond explaining model outputs, these methods have the potential to improve feature selection by identifying bands that significantly influence model predictions. This explainability-driven band selection enables dimensionality reduction while maintaining robust performance. From an industrial perspective, isolating the most discriminative bands supports real-time applications by enabling faster image acquisition and enhancing overall efficiency.

This paper is *exploratory*: we examine whether post-hoc explanation methods (LRP, SHAP, and RISE) can guide hyperspectral band selection. For that we propose a four-stage methodology where we integrate explainability methods with Convolutional Neural Networks (CNNs) classifiers and prioritises influential spectral bands in line with model behaviour, avoiding statistical or heuristic criteria. This framework offers a promising direction towards advancing hyperspectral image analysis and presents the following key contributions are:

- (1) **Exploratory evaluation of explainability for band selection:** We conduct one of the first comparative studies of three established post-hoc explainability methods (LRP, SHAP, RISE) for band selection in HSI, assessed on two CNNs (TRI-CNN and HSI-CNN) and two benchmarks (PaviaU and Salinas), addressing a methodological gap in the literature.
- (2) **Physically faithful band selection:** Band relevance is computed in the native wavelength domain, preserving the original band–wavelength correspondence and requiring no feature–space transformations or architectural changes.
- (3) **Accuracy-preserving dimensionality reduction:** Classifiers using as few as 30 explainability-selected bands match or exceed full-spectrum baselines.

We organise the remainder of the paper as follows: Section 2 reviews related work and positions our research with respect to existing efforts. Section 3 describes the proposed method. In Section 4, we evaluate the considered explanation methods and empirically validate the explainability-based bands on model performance. Finally, Section 5 offers concluding remarks and outlines future work prospects.

## 2 Related Work

We position our work along two main dimensions:

**Band selection methods.** Dimensionality reduction constitutes a key preprocessing step in hyperspectral image analysis, with band selection often preferred due to its ability to retain original spectral information while reducing redundancy and computational complexity. Broadly, band-selection methods fall into three categories, search-based, learning-driven, and hybrid, each with distinct mechanisms, benefits and trade-offs.

*Search-based methods.* These include (1) Ranking, wherein each spectral band is assigned a score based on a pre-defined criterion, such as mutual information [16], and the top scoring bands are selected. Although straightforward, ranking can be computationally expensive and may overlook inter-band dependencies, resulting in the selection of redundant bands. (2) Clustering groups spectrally similar bands and selects representative bands from each cluster [17]. While effective in reducing redundancy, performance of these methods relies heavily on the choice of clustering algorithm and its parameters. (3) Heuristic search, such as genetic algorithms, iteratively refines subsets of bands to maximise performance [18, 19], albeit with significant computational demands.

*Learning-driven methods.* These approaches integrate band selection directly into model training, allowing the model to highlight the most relevant bands [20]. (1) Embedded regularisation applies sparsity constraints to the model’s input weights, ablating uninformative bands and flagging the remainder as important [21, 22]. However, these selections may not generalise across models. (2) Deep learning, such as autoencoders, compress and reconstruct the full spectrum; analysing reconstruction errors or bottleneck activations reveals the most informative bands [23]. Alternatively, attention mechanisms in CNNs or transformers learn per-band weights during training, selecting those with consistently high scores [24, 25]. Both mechanisms effectively learn compact and performant spectral representations, yet they introduce computational overhead.

*Hybrid methods.* These methods combine elements from search-based and learning-driven strategies aiming to balance computational costs and selection accuracy [26, 27].

In contrast to conventional approaches, we harness explainability methods to identify spectral bands that strongly drive predictions. This allows us to avoid the drawbacks of high computational overhead and inter-band correlations typically associated with conventional band selection.

**Model Explainability.** We focus on the post-hoc explainability methods that justify the decisions of a pre-trained model. In image analysis, these methods assign importance to pixels or regions, highlighting features that drive model predictions. In the hyperspectral imagery, the spectral dimension maps each input channel to a specific wavelength, enabling band-level attributions. For instance, [14] applies the model agnostic SHAP algorithm to a CNN-based hyperspectral remote sensing image classifier to explain model outputs. A preprocessing step using Principal Component Analysis (PCA) [28] is applied to reduce dimensionality by projecting the original pixel-wise band values onto a new set of orthogonal components. This step removes the direct band–wavelength correspondence, thus preventing us from attributing the model decisions back to the individual bands and undermining spectral interpretability. Similarly, [15] employs back-propagation-based methods such as GradCAM [29], GradCAM++ [30] and Guided Back-propagation [31] to highlight critical spectral regions, but reports mixed results regarding band-specific importance. Their results indicate that while Guided Back-propagation identifies critical spectral bands for certain classes, the overall band importance remained relatively uniform, underscoring the need for sharper wavelength-level attributions. While [32] integrates multiple explainability methods for band selection, it primarily evaluates combinations of CNN architectures and normalisation strategies. Although it reports performance with selected bands, it lacks a direct comparison of the relative effectiveness of the explainability methods used, their contributions, quality or suitability for guiding the selection.

Our research leverages explainability to guide band selection in HSI. By identifying and validating the bands that have the highest influence on the model’s decisions, our approach simultaneously reduces dimensionality, maintains interpretability, and enhances computational efficiency.

## 3 Model design and description

In this section, we outline our methodology for band selection using explainability techniques. Figure 1 provides an overview of the four stages involved in the process.

### Stage 1: Classification

We begin by densely extracting overlapping patches from the hyperspectral image  $X \in \mathbb{R}^{h \times w \times b}$  where  $h$  and  $w$  are the spatial dimensions and  $b$  the number of spectral bands. Each patch  $X_p \in \mathbb{R}^{h' \times w' \times b}$  corresponds to a small spatial region of  $X$  of size  $h'$  and  $w'$ , retains all  $b$  bands, and is centred on pixel  $(i, j)$  from which it inherits its ground-truth label. We define a classifier as the mapping

$$F : \mathbb{R}^{h' \times w' \times b} \longrightarrow \mathbb{R}^c, \quad (1)$$

where  $c$  is the number of classes. For each input patch  $X_p$ , the classifier produces a confidence (logit or probability) vector

$$F(X_p) = (F(X_p)^{(1)}, \dots, F(X_p)^{(c)}) \in \mathbb{R}^c \quad (2)$$

We obtain the predicted class by the maximum confidence rule:

$$Y_p = \arg \max_{k \in \{1, \dots, c\}} F(X_p)^{(k)}. \quad (3)$$

### Stage 2: Explainability

In the second stage, an explainability method  $E$  is applied to analyse the classifier and its decision in order to assign importance values to each spectral band. We formalise it as

$$E : \mathcal{F} \times \mathbb{R}^{h' \times w' \times b} \times \mathcal{R} \longrightarrow \mathbb{R}^b \quad (4)$$

where  $\mathcal{F}$  denotes the space of classifiers,  $\mathbb{R}^{h' \times w' \times b}$  is the input patch space, and  $\mathcal{R}$  represents the internal model representations (e.g. feature maps or latent activations). The output  $O = E(F, X_p, R) \in \mathbb{R}^b$  is a band-relevance vector that assigns an importance score to each spectral band for the model's decision on patch  $X_p$ .

### Stage 3: Evaluation

We assess each explanation method by systematically measuring how the model's confidence changes as bands are progressively perturbed according to their relevance scores. Specifically, bands are either *deleted* (replaced by the training-set mean) or *inserted* (restoring the original values). At each step, we modify 20% of the bands, in descending order of relevance, and record the resulting change in confidence. Formally, we define the evaluation operator

$$V : \mathbb{R}_{\geq 0}^b \times \mathbb{R}^{h' \times w' \times b} \times \mathcal{F} \longrightarrow \mathbb{R}, \quad (5)$$

which, when applied to a relevance vector  $O \in \mathbb{R}_{\geq 0}^b$ , a patch  $X_p \in \mathbb{R}^{h' \times w' \times b}$ , and a classifier  $F \in \mathcal{F}$ , yields

$$Q = V(O, X_p, F) \in \mathbb{R}, \quad (6)$$

where  $Q$  is the evaluation score that quantifies the confidence drop (for deletion) or confidence gain (for insertion). This procedure reveals how faithfully the relevance vector  $O$  captures the model's reliance on individual bands.

### Stage 4: Band Selection and Retraining

Based on the evaluation results, bands whose removal elicit the largest confidence drop or whose insertion yields the largest gain, are deemed most relevant. We aggregate and normalise these per-band confidence changes into a single influence score and select the top  $b'$  bands. Formally, let  $\mathcal{B}' \subset \{1, \dots, b\}$ ,  $|\mathcal{B}'| = b'$  denote the selected subset of bands, where  $b' < b$ . For each original patch  $X_p$ , we construct  $X'_p \in \mathbb{R}^{h' \times w' \times b'}$ , by restricting  $X_p$  to the bands in  $\mathcal{B}'$ . We then retrain an adjusted classifier  $F' : \mathbb{R}^{h' \times w' \times b'} \longrightarrow \mathbb{R}^c$ , so that for each reduced patch  $X'_p$  the prediction is  $Y'_p = F'(X'_p)$ .

## 4 Experiments

### 4.1 Data

We validate our methodology on two publicly available hyperspectral benchmarks: Pavia University (PaviaU) and Salinas [33]. PaviaU, acquired by the ROSIS sensor, comprises 103 spectral bands spanning  $0.43 \mu\text{m}$  to  $0.86 \mu\text{m}$  at  $1.3 \text{ m}$  per-pixel spatial resolution, with image dimensions of  $610 \times 340$  pixels and 9 ground-truth classes. Salinas, captured by the AVIRIS sensor, originally contains 224 spectral bands spanning  $0.4 \mu\text{m}$  to  $2.5 \mu\text{m}$ ; after removing

water-absorption bands, 204 bands remain. Its spatial resolution is 3.7 m per pixel, and has  $512 \times 217$  pixels. It encompasses 16 ground-truth classes.

For both datasets, we apply the same dense, overlapping patch extraction described in section 3, Stage 1: Classification. This yields approximately 42,318 labelled patches for PaviaU and 54,129 for Salinas, each patch inheriting the class of its centre pixel.

## 4.2 Neural Network Architecture

To validate our approach, we consider two CNN architectures for hyperspectral remote sensing classification, using their predictions as baseline for evaluation. These architectures are not intended for benchmarking or state-of-the-art performance; rather, they serve as test beds to validate the utility of explainability methods in guiding band selection in our experiments.

**Tri-CNN** [34]. A multi-scale 3D-CNN with three parallel branches (spectral, spatial, and spectral-spatial). Feature maps of each branch are flattened, concatenated, and fed to fully connected layers and a softmax layer. We use  $13 \times 13 \times b$  (PaviaU) and  $11 \times 11 \times b$  (Salinas) patches retaining the full spectral dimension  $b$ , unlike the original work which reduces the spectral dimension to 15 and 35 components, respectively, using PCA. For PaviaU, we follow the original split of 1% train, 99% test. For Salinas, due to its higher spectral dimensionality, we use 75% train, 25% test, to manage the computational load.

**HSI-CNN** [35]. A two-dimensional (2D) CNN that converts one-dimensional spectral vectors into 2D matrices and processes them by standard 2D convolutions to extract spectral-spatial features. We use  $3 \times 3 \times b$  patches on both datasets and the original split of 30% train/ 70% test.

Under both methods, we split each model’s designated test set into two equal subsets. One is used for assessing the classification accuracy, and the other is reserved exclusively for generating and evaluating the explanation maps.

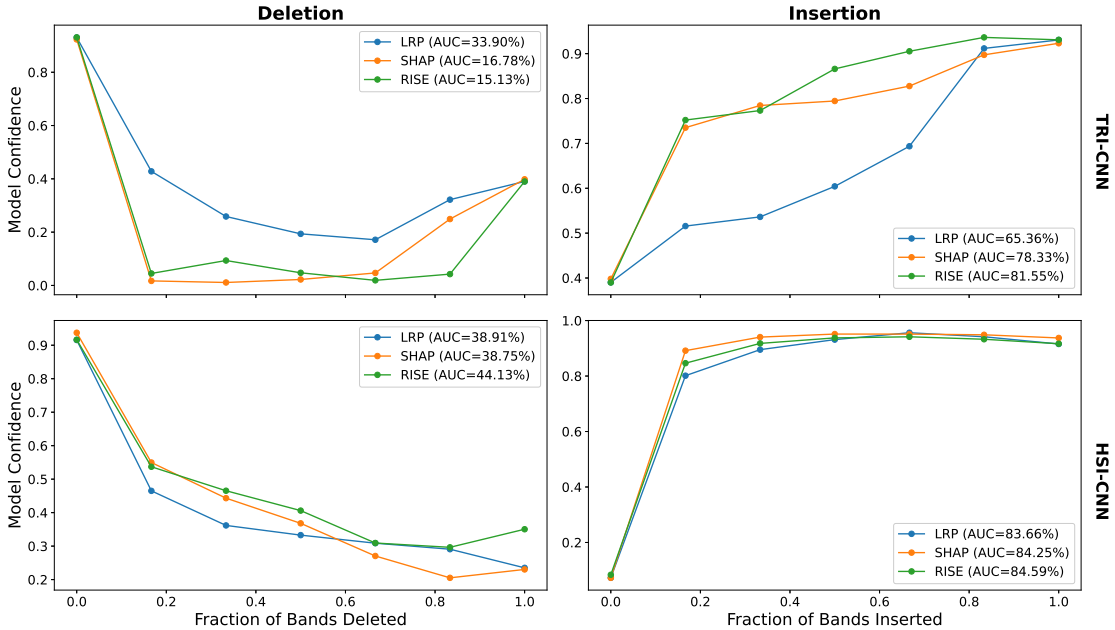


Figure 2: PaviaU: Deletion / Insertion AUC curves. TRI-CNN (top), HSI-CNN (bottom).

For Stage 2, we implement three explainability methods:

**Layer-wise Relevance propagation (LRP)** [36] attributes the neural network’s decision by redistributing its output backwards through each layer according to predefined rules described in [37, 38]. We first perform a standard forward pass to store the activation maps. Relevance is then initialised at the networks’ output layer using the logit of the true class and propagated backward. We apply different LRP rules at various layers: (1) the **LRP-0** rule for fully connected layers, which distributes relevance proportionally to neuron activations; (2) the  **$\epsilon$ -rule** for intermediate layers, introduces a small stabiliser  $\epsilon$  to suppress noise and avoid division by zero; (3) the  **$\gamma$ -rule** for early layers,

which scales positive weights by a factor  $\gamma > 0$  to emphasise positive contributions (set to 0.25 in our experiments). The output is a relevance map  $O_{LRP} \in \mathbb{R}^{h' \times w' \times b}$  where each value reflects the contribution of a spectral band and spatial location to the model’s decision.

**SHAP (SHapley Additive exPlanations)** [39] assigns a Shapley value to each input feature, quantifying its contribution to the model’s output. The values are grounded in cooperative game theory and are computed by evaluating every possible feature subset. However, this exhaustive computation is infeasible for high-dimensional data. To balance accuracy and efficiency, we sample 30 random subsets of bands. For each subset, we compute two predictions: one including the band,  $p_{incl} = F(X_p^{S \cup \{\text{band}\}})$ , and one excluding it,  $p_{excl} = F(X_p^S)$ . The band’s marginal contribution for that subset is computed as  $(p_{incl} - p_{excl})$ . Averaging these differences over the 30 subsets yields an approximate Shapley value for each band. Repeating this for all  $X_p$  patches produces a band-level importance matrix  $O_{SHAP} \in \mathbb{R}^{X_p \times b}$ .

**Randomized Input Sampling for Explanation (RISE)** [40] explains the model output using perturbations to generate saliency maps. Unlike gradient-based methods, RISE is architecture-agnostic and requires no access to the internal parameters of the model. In our implementation, we sample (5000) one-dimensional spectral masks that randomly occlude different bands while preserving the patch spatial layout. We apply a **density** ratio of 0.5, indicating that 50% of the input remains visible in each mask. Because masks are constant over the spatial extent, the procedure yields *band-level* (spatially uniform) relevance scores for each instance. Hence RISE returns a *band-level* relevance vector  $O_{RISE}(x) \in \mathbb{R}^b$ , i.e., one score per band. Stacking across  $X_p$  patches  $\{x_i\}_{i=1}^{X_p}$  gives  $O_{RISE} \in \mathbb{R}^{X_p \times b}$ .

### 4.3 Explainability Assessment Framework

We evaluate the faithfulness of the explanation-maps using two perturbation-based metrics [40], which quantify how changes to the most relevant features impact the prediction.

**Deletion – Insertion** [41]: measures how the model’s confidence changes as we gradually remove or add the most relevant bands.

*Ranking bands.* For a given input, we sum each band’s positive relevance scores over the spatial dimension, this gives a single relevance value per band. We then sort bands from the most to the least relevant.

*Deletion.* Starting with the intact input, we replace bands in 20% increments with their training-set mean values in the ranked order, continuing until all bands are replaced. After each increment, we record the confidence in the true class; a steeper decline indicates that the explanation method has correctly identified bands critical to the model’s decision.

*Insertion.* Starting from an empty input (all bands are set to zero), we reintroduce band values in descending order of relevance, progressively in 20% increments, tracking the increase in the confidence gain. A steeper rise in confidence shows that early-added bands drive prediction, indicating a more faithful explanation.

We summarise each process by the area under its confidence–fraction curve (AUC): lower AUC is better for deletion (downward trend), and higher AUC is better for insertion (upward trend). Together, these AUC values quantify how effectively an explanation method identifies the most important bands.

**Average–Drop.** This metric measures how much the model’s confidence in the true class falls when low-relevance bands are masked out. Specifically, it measures whether an explanation map effectively identifies critical features. Concretely, we perform an element-wise multiplication between the original input and its relevance map, preserving the most influential bands while suppressing the weakly relevant ones. We then feed the result through the model and measure the decrease in the model’s confidence for the true class. The percentage drop in class confidence reflects the overall quality of the explanation: smaller drops imply that the preserved bands capture the core discriminative information and thus a higher-quality explanations.

### 4.4 Explainability Evaluation Results

Figures 2, 3 display the AUC curves obtained from the *Deletion–Insertion* tests. In each case, the confidence decreases under deletion and increases under insertion, exhibiting the expected trends. The magnitude and rate of these changes differ between the models, which we attribute to differences in their architectural complexity. Replacing bands with their mean values rather than entirely removing them, prevents the confidence from reaching zero. Moreover, using mean values aligns more closely with the distribution learned by the model, causing a partial rebound in confidence after the replacement of the final fraction of bands.

Table 1 reports the *AverageDrop* (%) in model confidence across the three explainability methods. It is critical to highlight that band-level explanations may lack the granularity provided by pixel-level analyses, potentially leading to the omission or under-representation of truly important features. This loss of detail can lead to a larger drop in model

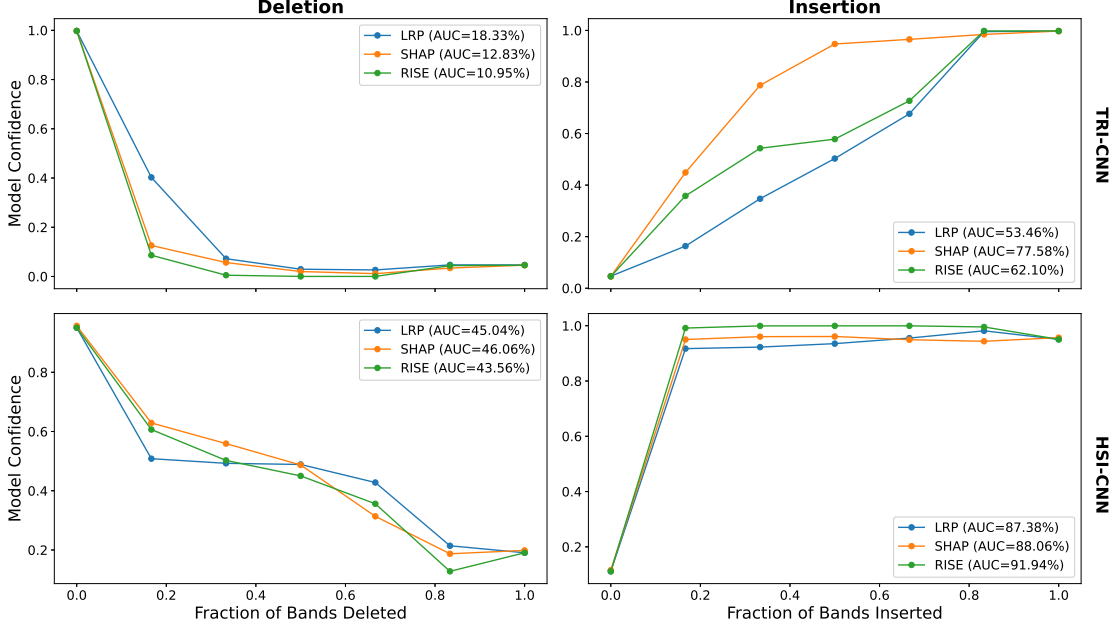


Figure 3: Salinas: Deletion / Insertion AUC curves. TRI-CNN (top), HSI-CNN (bottom).

Table 1: Average Drop (%) comparison across explainability methods for Pavia and Salinas.

Model	Dataset	LRP (%)	SHAP (%)	RISE (%)
TRI-CNN	Pavia	2.006	3.286	12.381
	Salinas	26.568	7.082	15.022
HSI-CNN	Pavia	0.218	7.527	1.317
	Salinas	1.795	4.981	0.250

confidence. Nonetheless, this band-wise evolution remains effective in comparing the relative ability of the different methods to identify the most influential spectral information.

#### 4.5 Explainability-Driven Top-30 Band Selection

Table 2: Average test accuracy (%) using the full bands and a subset of 30 bands selected based on explainability methods for PaviaU and Salinas datasets.

Model	Dataset	Full Bands	LRP-30	SHAP-30	RISE-30
TRI-CNN	Pavia	91.94 ± 1.67	91.68 ± 0.56	91.47 ± 0.46	91.21 ± 0.65
	Salinas	99.85 ± 0.02	99.82 ± 0.09	99.86 ± 0.03	99.86 ± 0.03
HSI-CNN	Pavia	95.25 ± 0.16	92.92 ± 0.27	92.98 ± 0.26	92.94 ± 0.14
	Salinas	97.30 ± 0.25	94.89 ± 0.13	94.75 ± 0.38	95.94 ± 0.25

Building on the procedure in Section 3—Stage 4, we employed the *Deletion – Insertion* evaluation method to identify the 30 bands with the highest influence on model confidence. Each network was then retrained and evaluated over five independent runs using the same hyperparameters as in the original training on the full spectral bands, adjusting only the kernel sizes and strides to accommodate the reduced spectral dimensionality.

Table 2 compares the test accuracies of TRI-CNN and HSI-CNN trained on the full set of spectral bands versus the 30 bands selected using the explainability methods. Despite the reduced spectral dimensionality, TRI-CNN maintains nearly identical performance across all band-selection strategies on both datasets, suggesting it can effectively adapt to fewer bands without a significant loss in accuracy. By contrast, the HSI-CNN, which relies more heavily on the

full spectral range, exhibits a noticeable decline in performance, albeit to varying extents. On PaviaU, accuracy drops by approximately 2.27%–2.33%, and on Salinas 1.36%–2.55%, indicating a higher sensitivity to spectral reduction. Given the intricate correlations among spectral bands, one might expect an even larger accuracy drop. However, our

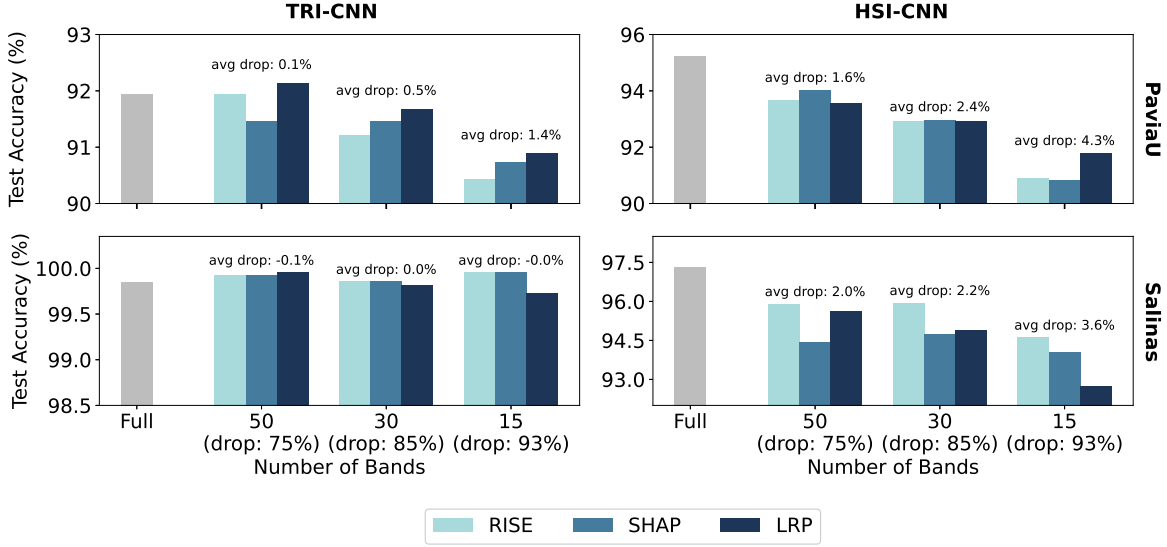


Figure 4: Test Accuracy Performance (%) of TRI-CNN (left) and HSI-CNN (right) with Different Subsets of Selected Bands (Full, 50, 30, 15). PaviaU (top), Salinas (bottom).

explainability-driven selection mitigates this effect by preserving the most informative bands, thereby limiting the loss of critical information. Overall, the results demonstrate that our explainability-based approach achieves effective dimensionality reduction while preserving robust performance.

#### 4.6 Effect of the Number of Selected Bands

To evaluate the suitability of the 30-band subset used previously, we compare model performance with 15- and 50-band subsets for both datasets (Figure 4). As expected, both models exhibited accuracy losses compared to using the full spectral range. However, the explainability-driven band-selection successfully isolated the most informative bands, keeping high accuracy even with only 15 bands. On PaviaU, both models maintained over 90% accuracy at 15 bands, while on Salinas the TRI-CNN remained near 100% and the HSI-CNN above 94%.

When we increased the bands to 50, both models improved on both datasets, yet with diminishing returns beyond 30 bands. For instance, on PaviaU, TRI-CNN’s accuracy rose approximately 0.7%–1.5% across RISE, SHAP, and LRP when going from 15 to 30 bands, yet only yielded another 0.7% when bands were increased to 50. A similar trend appeared on Salinas, where the TRI-CNN’s accuracy change after 30 bands was negligible ( $\leq 0.1\%$ ), highlighting its robustness to spectral-band reduction. The HSI-CNN showed larger overall gains, 1.8%–3.2% for PaviaU and 2.0%–3.6% for Salinas from 15 to 50 bands, but again we saw only marginal gains (0.65%–1.0%) beyond 30 bands.

Thus, while adding bands can boost accuracy, benefits plateau after 30 bands, suggesting that a 30-band subset strikes an optimal balance between computational efficiency and classification performance for both datasets.

#### 4.7 Distribution of Selected Bands’ Wavelengths

To investigate whether the different explainability methods converge on similar spectral regions (i.e. whether they follow similar wavelength distributions), we performed a kernel density estimation (KDE) over the wavelengths of the top 30 bands selected by each method (Figure 5). We computed KDEs using `gaussian_kde` with Scott’s rule  $h = \sigma n^{-\frac{1}{d+4}}$  which automatically determines a data-dependent bandwidth based on the spread of values within each subset. As a result, each curve may have a slightly different bandwidth, reflecting variations in the selected wavelength distributions across methods.



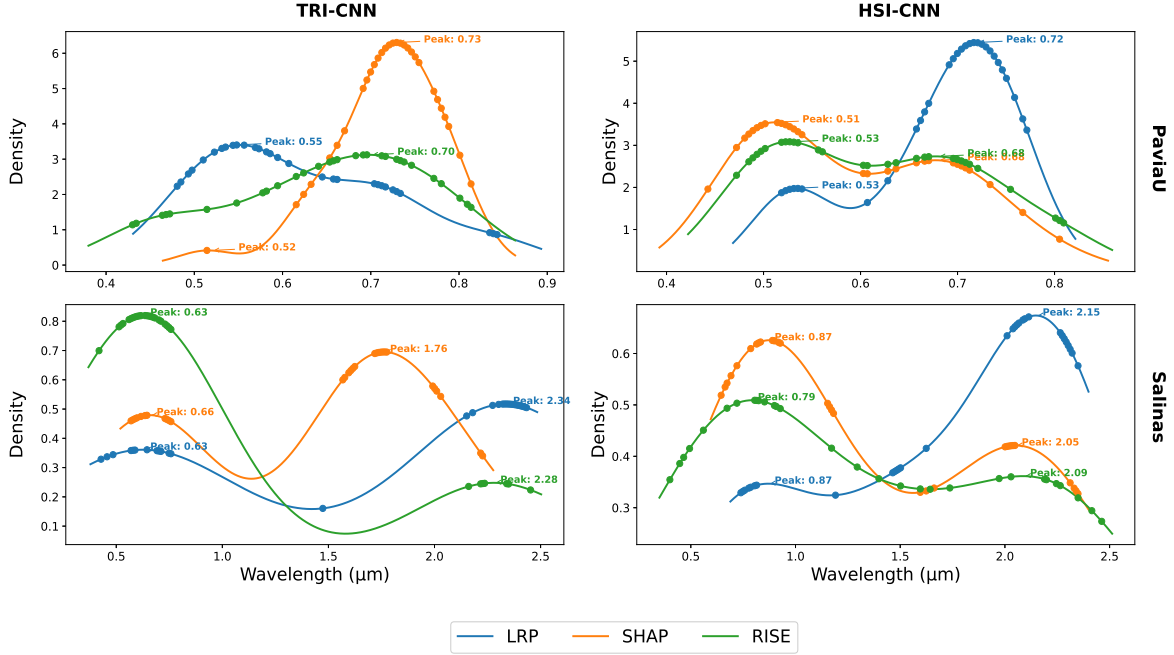


Figure 5: Distribution of selected bands wavelength-ranges (in  $\mu\text{m}$ ) by each explanation method for TRI-CNN (left), HSI-CNN (right), PaviaU (top), Salinas (bottom).

**Pavia University (Urban scene).** Although each explanation method selects slightly different spectral regions, all three methods consistently concentrate their selections within the visible–red-edge interval ( $0.43\text{--}0.84\mu\text{m}$ ). According to [42], this region coincides with the chlorophyll absorption peak ( $0.67\mu\text{m}$ ) and the red-edge ( $0.68\text{--}0.75\mu\text{m}$ ), which underpin vegetation indicators and help differentiate vegetated (natural) from man-made (built) surfaces. Bands in this interval are highly discriminative for urban land-cover analysis (e.g. asphalt, concrete, roofing, and vegetation).

For TRI-CNN, the KDE peaks occur at  $0.55\mu\text{m}$  (LRP),  $0.73\mu\text{m}$  (SHAP), and  $0.70\mu\text{m}$  (RISE), with all three distributions concentrated in the mid-visible region. Although LRP’s selections extend numerically to  $0.84\mu\text{m}$ , its primary focus lies between  $0.48\text{--}0.76\mu\text{m}$ .

For HSI-CNN, dual peaks appear near  $0.53\mu\text{m}$  and  $0.72\mu\text{m}$  for LRP,  $0.51\mu\text{m}$  and  $0.68\mu\text{m}$  for SHAP, and  $0.53\mu\text{m}$  and  $0.68\mu\text{m}$  for RISE. These broader selections ( $0.44\text{--}0.81\mu\text{m}$ ) still emphasise the mid-visible region, despite their wider span of the spectrum.

**Salinas dataset (Agricultural Scene).** All three methods show prominent selections in the visible/red edge range around  $0.5\text{--}0.75\mu\text{m}$  and the short-wave infrared (SWIR) region  $1\text{--}2.15\mu\text{m}$ . The visible–red-edge bands capture key pigment and structural variations in vegetation canopies, while the SWIR bands probe a range critical for characterising vegetation health, moisture content levels, and soil properties. According to [42, 43], SWIR wavelengths penetrate deeper into canopy layers and are sensitive to dry matter and leaf internal structure, making them particularly informative for agricultural modelling. TRI-CNN’s KDE nearly covers the entire visible–NIR range, while HSI-CNN shows two distinct peaks, one in the high visible–low-NIR region,  $0.79\text{--}0.87\mu\text{m}$ , and another in the SWIR ( $2.00\text{--}2.15\mu\text{m}$ ).

Although each explanation method selects slightly different spectral regions, they consistently highlight overlapping bands. For PaviaU, all methods converge on mid-visible wavelength ( $0.43\text{--}0.84\mu\text{m}$ ), while for Salinas, the selections span the visible into NIR, underscoring the importance of both spectral regions. These patterns indicate a dual alignment: the selected bands are *model-aligned*, reflecting the classifiers’ decision process, and *domain-meaningful*, coinciding with established spectral features (e.g. chlorophyll absorption, SWIR sensitivities); supporting explainability-driven band selection as a principled route to compact, interpretable spectral subsets.

## 5 Conclusion

We demonstrate that post-hoc explainability methods can effectively drive dimensionality reduction in hyperspectral imaging by identifying the most informative spectral bands and wavelength ranges. Despite their differing complexities, both CNN architectures maintained robust classification performance when retrained on the reduced spectral subset. Furthermore, all three methods consistently converged on similar wavelength regions, highlighting the robustness and physical relevance of the selected bands. Unlike conventional band selection techniques, which risk discarding critical signals, overlooking data diversity, or erroneously prioritising noise, our approach preserves the chemical and spectral integrity of the data while reducing redundancy. In future work, we will validate these findings on additional datasets and evaluate other families of explainability methods to further generalise our results.

## Acknowledgments

The research presented in this article is part of the project "Learning-based representations for the automation of hyperspectral microscopic imaging and predictive maintenance" funded by the Flanders Innovation & Entrepreneurship-VLAIO, under grant number HBC.2020.2266.

## References

- [1] A. Bhargava, A. Sachdeva, K. Sharma, M. H. Alsharif, P. Uthansakul, and M. Uthansakul. Hyperspectral imaging and its applications: A review. *Heliyon*, 10, 2024.
- [2] D. Al-Alimi, M. A. Al-qaness, Z. Cai, and E. A. Alawamy. Ida: Improving distribution analysis for reducing data complexity and dimensionality in hyperspectral images. *Pattern Recognition*, 134:109096, 2023.
- [3] H. Li, J. Cui, X. Zhang, Y. Han, and L. Cao. Dimensionality reduction and classification of hyperspectral remote sensing image feature extraction. *Remote Sensing*, 14(18), 2022.
- [4] Behnood Rasti, D. Hong, R. Hang, P. Ghamisi, X. Kang, J. Chanussot, and J. A. Benediktsson. Feature extraction for hyperspectral imagery: The evolution from shallow to deep: Overview and toolbox. *IEEE Geoscience and Remote Sensing Magazine*, 8(4):60–88, 2020.
- [5] A. Fejjari, K. Saheb Ettabaa, and O. Korbaa. Feature extraction techniques for hyperspectral images classification. In *Soft Computing Applications*, pages 174–188. Springer International Publishing, 2021.
- [6] O. S. C. Goud. Band selection methods for hyperspectral imagery analysis – a critical comparison. *International Journal of Intelligent Systems and Applications in Engineering*, 12(3):3093–3109, Mar 2024.
- [7] S. S. Sawant, P. Manoharan, and A. Loganathan. Band selection strategies for hyperspectral image classification based on machine learning and artificial intelligent techniques – survey. *Arabian Journal of Geosciences*, 14:646, 2021.
- [8] F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.
- [9] S. Ali, T. Abuhmed, S. El-Sappagh, K. Muhammad, J. M. Alonso-Moral, R. Confalonieri, R. Guidotti, J. Del Ser, N. Díaz-Rodríguez, and F. Herrera. Explainable artificial intelligence (xai): What we know and what is left to attain trustworthy artificial intelligence. *Information Fusion*, 99:101805, 2023.
- [10] R. Marcinkevičs and J. E. Vogt. Interpretable and explainable machine learning: A methods-centric overview with concrete examples. *WIREs Data Mining and Knowledge Discovery*, 13(3):e1493, 2023.
- [11] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.
- [12] L. Longo, M. Brcic, and F. Cabitza. Explainable artificial intelligence (xai) 2.0: A manifesto of open challenges and interdisciplinary research directions. *Information Fusion*, 106:102301, 2024.
- [13] A. Barredo Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020.
- [14] A. H. Oveis, E. Giusti, G. Meucci, S. Ghio, and M. Martorella. Explainability in hyperspectral image classification: A study of xai through the shap algorithm. In *Workshop on Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing*, pages 1–5, 2023.

- [15] D. E. Turan, E. Aptoula, A. Ertürk, and G. Taskin. Interpreting hyperspectral remote sensing image classification methods via explainable artificial intelligence. In *IEEE International Geoscience and Remote Sensing Symposium*, pages 5950–5953, 2023.
- [16] S. S. Sawant and M. Prabukumar. A survey of band selection techniques for hyperspectral image classification. *Journal of Spectral Imaging*, 2020.
- [17] R. Yang, L. Su, X. Zhao, H. Wan, and J. Sun. Representative band selection for hyperspectral image classification. *Journal of Visual Communication and Image Representation*, 48:396–403, 2017.
- [18] M. Esmaeili, D. Abbasi-Moghadam, A. Sharifi, A. Tariq, and Q. Li. Hyperspectral image band selection based on cnn embedded ga (cnnga). *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 16:1927–1950, 2023.
- [19] K. Deep and M. Thakur. Hyperspectral band selection using a decomposition based multiobjective wrapper approach. *Infrared Physics & Technology*, 136:105053, 2024.
- [20] C. O. Ayna, R. Mdrafi, Q. Du, and A. C. Gurbuz. Learning-based optimization of hyperspectral band selection for classification. *Remote Sensing*, 15(18), 2023.
- [21] W. Sun, G. Yang, J. Peng, and Q. Du. Hyperspectral band selection using weighted kernel regularization. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(9):3665–3676, 2019.
- [22] M. Ma, S. Mei, F. Li, Y. Ge, and Q. Du. Spectral correlation-based diverse band selection for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–13, 2023.
- [23] Z. Dou, K. Gao, X. Zhang, H. Wang, and L. Han. Band selection of hyperspectral images using attention-based autoencoders. *IEEE Geoscience and Remote Sensing Letters*, 18(1):147–151, 2021.
- [24] J. Wang, J. Zhou, and W. Huang. Attend in bands: Hyperspectral band weighting and selection for image classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(12):4712–4727, 2019.
- [25] S. Kumar, H. Aetesam, A. Saha, and S. K. Maji. Attention-based deep autoencoder for hyperspectral image denoising. In B. Raman, S. Murala, A. Chowdhury, A. Dhall, and P. Goyal, editors, *Computer Vision and Image Processing (CVIP 2021). Communications in Computer and Information Science*, vol. 1568, pages 171–183. Springer, Cham, 2022.
- [26] H. Fu, A. Zhang, G. Sun, J. Ren, X. Jia, Z. Pan, and H. Ma. A novel band selection and spatial noise reduction method for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–13, 2022.
- [27] Y. Zimmer and O. Glickman. Embedded hyperspectral band selection with adaptive optimization for image semantic segmentation. Preprint, 2024.
- [28] M. A. M. P. Uddin and M. A. Hossain. Pca-based feature reduction for hyperspectral remote sensing image classification. *IETE Technical Review*, 38(4):377–396, 2021.
- [29] R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128:336–359, 2020.
- [30] A. Chattopadhyay, A. Sarkar, P. Howlader, and V. N. Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847, 2018.
- [31] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net, 2015.
- [32] G. Zhang and W. Abdulla. Explainable ai-driven wavelength selection for hyperspectral imaging of honey products. *Food Chemistry Advances*, 3:100491, 2023.
- [33] Hyperspectral remote sensing scenes, salinas, Pavia university. <https://www.ehu.eus/ccwintco/index.php/Hyperspectral> 1998. Accessed: 10-02-2025.
- [34] M. Q. Alkhatib, M. Al-Saad, N. Aburaed, S. Almansoori, J. Zabalza, S. Marshall, and H. Al-Ahmad. Tri-cnn: A three branch model for hyperspectral image classification. *Remote Sensing*, 15(2), 2023.
- [35] Y. Luo, J. Zou, C. Yao, X. Zhao, T. Li, and G. Bai. Hsi-cnn: A novel convolution neural network for hyperspectral image. In *International Conference on Audio, Language and Image Processing*, pages 464–469, 2018.
- [36] A. Höhl, I. Obadic, M.-A. Fernández-Torres, H. Najjar, D. A. B. Oliveira, Z. Akata, A. Dengel, and X. X. Zhu. Opening the black box: A systematic review on explainable artificial intelligence in remote sensing. *IEEE Geoscience and Remote Sensing Magazine*, 12(4):261–304, 2024.

- [37] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE*, 10(7):1–46, 2015.
- [38] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, and K.-R. Müller. Layer-wise relevance propagation: An overview. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 193–209. Springer International Publishing, 2019.
- [39] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30. Curran Associates, Inc., 2017.
- [40] V. Petsiuk, A. Das, and K. Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.
- [41] Y. Wang and X. Wang. Benchmarking deletion metrics with the principled explanations. In *International Conference on Machine Learning (ICML)*, 2024.
- [42] Sadegh Ranjbar, Danielle Losos, Sophie Hoffman, Shiva Arabi, Ankur Desai, and Paul C. Stoy. Near real-time mapping of all-sky land surface temperature from goes-r using machine learning. *Journal of Geophysical Research: Machine Learning and Computation*, 2(2):e2024JH000464, 2025.
- [43] Freek van der Meer. Analysis of spectral absorption features in hyperspectral imagery. *International Journal of Applied Earth Observation and Geoinformation*, 5(1):55–68, 2004.