

GFR-CAM: Gram-Schmidt Feature Reduction for Hierarchical Class Activation Maps

Kaveh Safavigerdini¹Bahram Yaghooti²Amir Erfan Zareei Shams Abadi¹Kannappan Palaniappan¹¹University of Missouri-Columbia, USA ²Washington University in St. Louis, USA

{ksggh2, azzgg, pal}@missouri.edu {byaghooti}@wustl.edu

Abstract

Deep learning models have achieved remarkable success in computer vision tasks, yet their decision-making processes remain largely opaque, limiting their adoption especially in safety-critical applications. While Class Activation Maps (CAMs) have emerged as a prominent solution for visual explanation, existing methods suffer from a fundamental limitation: they produce single, consolidated explanations leading to “explanatory tunnel vision.” Current CAM methods fail to capture the rich, multi-faceted reasoning that underlies model predictions, particularly in complex scenes with multiple objects or intricate visual relationships. We introduce the Gram-Schmidt Feature Reduction Class Activation Map (GFR-CAM), a novel gradient-free framework that overcomes this limitation through hierarchical feature decomposition that provides a more holistic view of the architecture’s explanatory power. Unlike existing feature reduction methods that rely on Principal Component Analysis (PCA) and generate a single dominant explanation, GFR-CAM leverages Gram-Schmidt orthogonalization to systematically extract a sequence of orthogonal, information rich components from model feature maps. The subsequent orthogonal components are shown to be meaningful explanations — not mere noise, that decomposes single objects into semantic parts and systematically disentangles multi-object scenes to identify co-existing entities. We show that GFR-CAM on ResNet-50 and Swin Transformer architectures across ImageNet and PASCAL VOC datasets achieves competitive performance with state-of-the-art methods.

1. Introduction

Deep architectures such as Convolutional Neural Networks (CNNs) [15] and vision transformers [5] continue to revolutionize computer vision by achieving remarkable performance on tasks from visual object detection to video object tracking. Yet, despite their successes, CNNs and

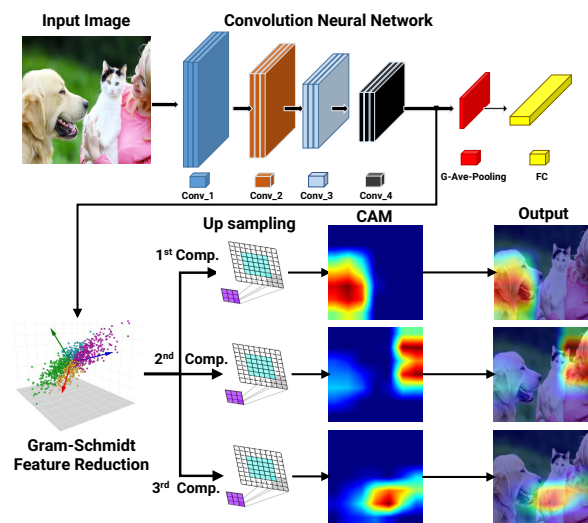


Figure 1. GFR-CAM integrated within a CNN. GFR-CAM uses Gram-Schmidt orthogonalization to generate a hierarchical explanation, sequentially identifying the primary object (Component 1) and distinct secondary objects (Components 2 and 3) to provide a disentangled hierarchical view of the model’s explainability.

deep learning models often remain opaque “black boxes,” providing little insight into the decision-making processes behind their predictions [17, 19, 45]. In safety-critical domains such as medical diagnostics [26, 37], autonomous and control systems [2, 43], and manufacturing [34, 46], this lack of interpretability remains a significant barrier, fueling a growing need for methods that can reliably elucidate the inner workings of these complex models. We were motivated by our own applications that could benefit from better deep architecture explainability, including malaria red blood cell segmentation [38], aerial vehicle tracking [24, 32], vessel segmentation [40], and carbon nanotube property analysis [30, 31, 36].

Class Activation Maps (CAM) [45] have emerged as a leading approach to this challenge, visually highlighting image regions that influence a model’s output [11, 32]. While early methods relied on gradients (e.g., Grad-CAM[33]), a new class of more robust, gradient-free techniques has gained prominence. Methods like Eigen-CAM [22] and Kernel-PCA CAM (KPCA-CAM) [14] leverage powerful decomposition techniques such as Principal Component Analysis (PCA) [23] and its nonlinear extension, Kernel PCA [20]. By directly analyzing convolutional feature maps, these methods produce clean explanations without gradient-related noise and vanishing gradient issues.

However, despite their advancements, these state-of-the-art decomposition methods [1, 8, 21, 35] share a fundamental limitation: they are designed to produce a *single, consolidated explanation*. By projecting all feature map information onto a single principal component, they generate a monolithic heatmap that captures only the most dominant evidence for a prediction. This “winner-takes-all” approach creates a form of explanatory tunnel vision. In multi-object scenes, it highlights one object while ignoring others that may be relevant to the model’s overall understanding. For single-object images, it obscures the rich “visual subtext”—crucial secondary features like texture, context, or distinct object parts—that are essential for a complete and robust interpretation of the model’s reasoning.

This limitation is a direct consequence of the feature reduction techniques employed. PCA and its variants are designed to find the single direction of maximum variance; subsequent components, by definition, capture progressively less information and often represent noise rather than meaningful, distinct concepts. The core challenge, therefore, is not to simply find the best single explanation, but to *decompose the feature space into a hierarchy of interpretable parts*. To achieve this, we move away from PCA and turn to a different, theoretically-grounded decomposition method: *Gram-Schmidt orthogonalization* [41, 42, 44]. Unlike PCA, the Gram-Schmidt process is explicitly designed to construct a set of linearly independent basis vectors. This inherent property makes it an ideal candidate for systematically extracting a sequence of distinct, information-rich, and non-redundant components from the feature space, paving the way for a multi-faceted explanation.

Building on this foundation, we introduce the *Gram-Schmidt Feature Reduction Class Activation Map (GFR-CAM)*, a novel framework that leverages Gram-Schmidt orthogonalization to deconstruct a model’s reasoning into a hierarchy of visual evidence. As illustrated in Figure 1, GFR-CAM analyzes the final convolutional feature maps to generate a sequence of activation maps. The first component reliably identifies the primary evidence, achieving per-

formance competitive with state-of-the-art methods. However, the true innovation lies in the subsequent components. Because the Gram-Schmidt process isolates linearly independent and information-rich features, the second, third, and further activations are not noise; they are meaningful explanations that solve the tunnel vision problem. In single-object images, they isolate supporting evidence like texture and distinct parts. In multi-object scenes, they systematically disentangle the visual field, identifying co-existing objects one by one to reveal a comprehensive scene understanding that monolithic CAMs cannot provide.

Our contributions are threefold:

- We introduce GFR-CAM, a novel, gradient-free, and computationally efficient CAM framework built on Gram-Schmidt orthogonalization, which is theoretically suited for hierarchical feature decomposition.
- We demonstrate that GFR-CAM overcomes the single-explanation limitation of existing methods by generating a hierarchy of interpretable activation maps. We show that its subsequent components are not noise but instead reveal secondary features, contextual clues, and distinct objects in complex scenes.
- Through extensive quantitative and qualitative evaluation across multiple models and rigorous benchmarks, we prove that GFR-CAM provides a more complete and disentangled understanding of a model’s decision-making, particularly in complex, multi-object scenes, thereby setting a new standard for comprehensive AI explanation.

2. Methodology

In this section, we present the Gram-Schmidt Class Activation Maps (GFR-CAM) framework. Our method enhances class activation maps for interpretability using Gram-Schmidt orthogonalization and feature reduction [41, 44]. This approach identifies discriminative image regions with higher precision and reduced redundancy, revealing the model’s hierarchical visual evidence.

The overall GFR-CAM workflow is summarized as follows (see Figure 1). The key innovation is in Step 3, where hierarchical, orthogonal components are generated to produce a sequence of interpretable activation maps:

1. **Feature Extraction:** Extract the feature maps $\{A^k\}_{k=1}^K$ from the final convolutional layer of a pre-trained CNN or Vision Transformer (ViT) for a given input image.
2. **Data Preparation:** Reshape the feature maps into a matrix $F \in \mathbb{R}^{K \times L}$, where K is the number of channels and $L = H \times W$ is the product of the spatial dimensions.
3. **Hierarchical Component Generation:** Apply the Gram-Schmidt Functional Reduction (GFR) procedure (detailed in Section 3.2) to the feature matrix F . This iteratively computes a set of orthogonal direction vectors $\{\nu_1, \nu_2, \dots, \nu_m\}$ and their corresponding activation maps $\{M_1, M_2, \dots, M_m\}$.

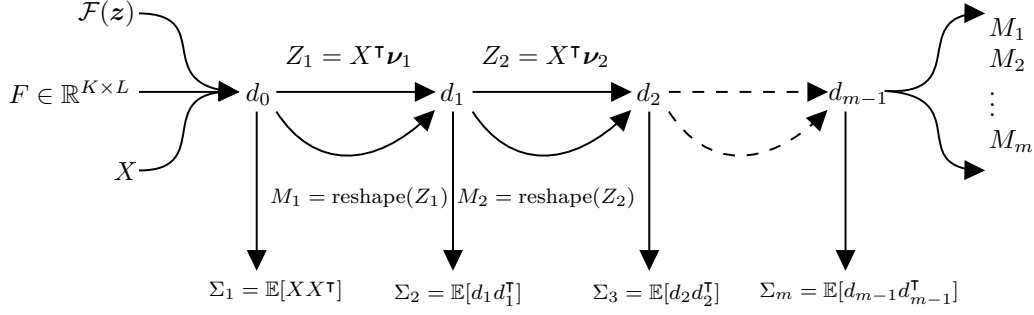


Figure 2. Schematic overview of the GFR-CAM framework. It begins with a family of functions and independent CNN feature samples, then iteratively substitutes fixed parameters with linear projections of the data, orthogonalizes the resulting functions, and subtracts their contributions from the original features [41, 44].

4. **Post-processing:** Normalize each M_j to $[0,1]$, then up-sample to original image size via bilinear interpolation.
5. **Visualization:** Overlay each processed map onto the original image to generate sequential explanations.

In what follows, we first describe the foundational orthogonalization algorithm before detailing how it is employed within the GFR procedure to generate our hierarchical GFR-CAMs.

2.1. Gram-Schmidt Orthogonalization: Core Mechanism

The foundation of our method is the Gram-Schmidt process, which we use to systematically decompose the CNN/Vit’s feature space. This process is formalized in Algorithm 1. At each iteration j , the algorithm takes a new feature component Z_j and makes it orthogonal to all previously computed components. This is achieved by subtracting its projection onto the already-orthogonalized function set \mathcal{T} (Line 4) and then normalizing the result (Line 5).

This process can be concisely represented by the classical Gram-Schmidt projection formula. For a new component Z_j and a set of existing orthonormal components $\{\hat{Z}_i\}_{i=1}^{j-1}$, the orthogonalization is:

$$\tilde{Z}_j = Z_j - \sum_{i=1}^{j-1} \langle Z_j, \hat{Z}_i \rangle \hat{Z}_i, \quad \text{and} \quad \hat{Z}_j = \frac{\tilde{Z}_j}{\|\tilde{Z}_j\|} \quad (1)$$

where \tilde{z}_j is the orthogonalized component and \hat{Z}_j is its normalized version. A critical output of Algorithm 1 is the residual covariance matrix Σ_{j+1} (Line 9). This matrix captures the variance structure of the data *after* the information from the first j components has been removed, which is essential for finding the next most important feature direction.

2.2. Hierarchical Activation Map Generation with GFR

We operationalize the orthogonalization process using Gram-Schmidt Functional Reduction (GFR), as detailed in

Algorithm 1: Orthogonalize($\hat{\mathcal{F}}(Z_1, \dots, Z_{j-1}), \mathcal{F}(z_{[j]}) \setminus \mathcal{F}(z_{[j-1]}), \{Z_i\}_{i=1}^j$)

Input:

- $\hat{\mathcal{F}}(Z_1, \dots, Z_{j-1})$: an orthogonalized function family in random variables Z_1, \dots, Z_{j-1} .
- $\mathcal{F}(z_{[j]}) \setminus \mathcal{F}(z_{[j-1]})$: all the functions in $\mathcal{F}(z)$ which depend on $z_{[j]}$ but not only on $z_{[j-1]}$.
- $\{Z_i\}_{i=1}^j$: the random variables Z_1, \dots, Z_{j-1} .
- Z_j : a new random variable.

Output:

- Orthogonalized function family $\hat{\mathcal{F}}(Z_1, \dots, Z_j)$
- A covariance matrix Σ_{j+1} .

- 1 **Initialize:** $\mathcal{T} \leftarrow \hat{\mathcal{F}}(Z_1, \dots, Z_{j-1})$. Denote $\mathcal{F}(z_{[j]}) \setminus \mathcal{F}(z_{[j-1]}) \triangleq \{f_1(z), \dots, f_\ell(z)\}$, where $f_1(z) = z_j$.
 - 2 **for** $k \leftarrow 1$ **to** ℓ **do**
 - 3 $g(Z_1, \dots, Z_j) \leftarrow f_k(Z_1, \dots, Z_j)$.
 - 4 $g(Z_1, \dots, Z_j) \leftarrow g(Z_1, \dots, Z_j) - \sum_{f \in \mathcal{T}} \mathbb{E}[g f] f$;
 - 5 $g(Z_1, \dots, Z_j) \leftarrow \frac{g(Z_1, \dots, Z_j)}{\sqrt{\mathbb{E}[g(Z_1, \dots, Z_j)^2]}}$.
 - 6 $\mathcal{T} \leftarrow \mathcal{T} \cup \{g(Z_1, \dots, Z_j)\}$.
 - 7 Define $\hat{\mathcal{F}}(Z_1, \dots, Z_j) = \mathcal{T}$.
 - 8 Define $d_j(X) = X - \sum_{f \in \mathcal{T}} \mathbb{E}[X f] f$.
 - 9 Define $\Sigma_{j+1} = \mathbb{E}[d_j(X) d_j^T(X)]$.
 - 10 **Return** $\Sigma_{j+1}, \mathcal{T}$.
-

Algorithm 2, to produce our sequence of activation maps.

The GFR algorithm iteratively identifies orthogonal components that maximize variance. At each step j , it finds the principal eigenvector ν_j of the current (residual) covariance matrix Σ_j . This eigenvector represents the direction of maximum variance in the information that has not yet been explained by previous components.

Algorithm 2: Gram-Schmidt Functional Reduction (GFR) for GFR-CAM

Input: Random variable X taken from feature maps $F \in \mathbb{R}^{K \times (H,W)}$, function family $\mathcal{F}(z)$, and threshold $\epsilon > 0$.

Output: Orthogonal activation maps $\{M_j\}_{j=1}^m$.

```
1 Initialize:  $\Sigma_1 \leftarrow \mathbb{E}[XX^\top]$ .
2 for  $j \leftarrow 1$  to  $K$  do
3    $\nu_j \leftarrow \arg \max_{\|\nu\|=1} \nu^\top \Sigma_j \nu$ ;
4   if  $\nu_j^\top \Sigma_j \nu_j \leq \epsilon^2$  then
5     break.
6    $Z_j \leftarrow X^\top \nu_j$ .
7    $M_j \leftarrow \text{reshape}(Z_j, H, W)$ .
8    $\Sigma_{j+1}, \hat{\mathcal{F}}(Z_{[j]}) \leftarrow$ 
     Orthogonalize( $\Sigma_j, \hat{\mathcal{F}}(Z_{[j-1]}), \mathcal{F}(z_{[j]}) \setminus$ 
      $\mathcal{F}(z_{[j-1]}), \{Z_i\}_{i=1}^j$ ).
9 Return  $\{M_j\}_{j=1}^m$ 
```

Primary Activation Map (M_1): The first activation map, analogous to the output of Eigen-CAM, is generated by projecting the feature maps F onto the principal eigenvector ν_1 of their covariance matrix $\Sigma_1 = \mathbb{E}[XX^\top]$. This map, M_1 , highlights the most dominant visual evidence for the model’s prediction:

$$M_1 = F^\top \nu_1 \quad (2)$$

Subsequent Activation Maps (M_j): Beyond the primary map, subsequent maps are generated by iteratively applying the Gram Schmidt algorithm. For each new map M_j , the algorithm first projects out the information captured by previous components (ν_1, \dots, ν_{j-1}). It then finds the principal eigenvector ν_j of the remaining data. This vector is guaranteed to be orthogonal to its predecessors and thus highlights a new, independent facet of the model’s reasoning, visualized as:

$$M_j = F^\top \nu_j \quad (3)$$

This iterative decomposition continues until the residual variance is negligible (below ϵ), providing a richer, multi-faceted explanation that goes beyond the single focus of existing methods.

3. Experiment Results

Our comprehensive evaluation is twofold. First, we establish GFR-CAM’s primary component as a competitive baseline by benchmarking it against state-of-the-art methods across diverse architectures and metrics, including ROAD. Second, and more critically, we demonstrate its unique explanatory power by analyzing its subsequent compo-

nents. We show qualitatively and quantitatively how GFR-CAM surpasses other decompositional methods in providing richer, more disentangled explanations, particularly for complex scenes. The results confirm that GFR-CAM is not only a competitive general-purpose CAM but also offers a superior method for generating multi-faceted, interpretable insights into a model’s reasoning.

3.1. Experimental Setup

Models and Datasets. We conduct a comprehensive evaluation on two representative architectures: ResNet-50 [10], a standard CNN, and the Swin Transformer [16]. For ResNet-50, feature maps are extracted from the final convolutional block (layer4), while for Swin-T, we use the output of the last normalization layer in the final stage. Our evaluation is performed on two standard benchmarks: a 5,000-image subset of the ImageNet (ILSVRC2012) validation set [29] (1,000 classes) and a 4,000-image subset from PASCAL VOC 2012 [7] (20 classes). For preprocessing, all images are resized and center-cropped to 224×224 pixels and normalized using the standard ImageNet mean and standard deviation.

Evaluation Metrics. To assess the quality of the generated CAMs, we employ a standard suite of metrics that measure an explanation’s faithfulness and interpretability [4, 13, 25]. Faithfulness is measured by quantifying the change in model confidence as regions identified by the CAM are manipulated. Specifically, we use the following six established metrics.

Average Drop (AD) [4] measures the average percentage drop in confidence when only the region highlighted by the CAM is provided as input. A lower AD is better, indicating the map preserves class-discriminative features:

$$\text{AD} = \frac{1}{N} \sum_{i=1}^N \frac{\max(0, y_i^c - o_i^c)}{y_i^c} \times 100 \quad (4)$$

Here, y_i^c and o_i^c are the model’s confidence scores for class c on the original image and the explanation-preserved image, respectively.

Coherency (Coh) [25] measures the stability of an explanation by calculating the Pearson correlation between the original heatmap and a heatmap generated from the explanation itself. Higher values indicate more consistent explanations:

$$\text{Coh}(x) = \frac{1}{2} \frac{\text{Cov}(H_c(x \otimes H_c(x)), H_c(x))}{\sigma_{H_c(x \otimes H_c(x))} \sigma_{H_c(x)}} + \frac{1}{2} \quad (5)$$

where $H_c(\cdot)$ is the CAM generation function, \otimes is element-wise multiplication, and $\text{Cov}(\cdot, \cdot)$ denotes the covariance between two random variables.

Complexity (Com) [25] measures the sparsity of an explanation via its L_1 norm. Lower complexity signifies a more focused and simpler explanation:

$$\text{Com}(x) = \|H_c(x)\|_1$$

Average Drop, Coherency, and Complexity (ADCC) [25] is a unified metric combining the above via their harmonic mean. A higher ADCC score indicates a better overall balance of explanation qualities:

$$\text{ADCC} = 3 \left(\frac{1}{1 - \text{AD}} + \frac{1}{\text{Coh}} + \frac{1}{1 - \text{Com}} \right)^{-1} \quad (6)$$

Increase in Confidence (IC) [4] measures the percentage of images for which the model’s confidence increases when shown only the CAM-highlighted region:

$$\text{IC} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i^c < o_i^c) \times 100 \quad (7)$$

Here, $\mathbb{I}(\cdot)$ is the indicator function, and y_i^c and o_i^c are the confidences for the original and explanation-preserving images, respectively.

Average Drop in Deletion (ADD) [13] measures the drop in model confidence when the CAM-highlighted region is removed from the image. A higher ADD is better, as it indicates the removed region was critical to the prediction:

$$\text{ADD} = \frac{1}{N} \sum_{i=1}^N \frac{\max(0, y_i^c - d_i^c)}{y_i^c} \times 100 \quad (8)$$

Here, y_i^c is the original confidence and d_i^c is the confidence after deleting the highlighted pixels.

3.2. Evaluation of the Primary GFR-CAM Component against General CAM Baselines

In this section, we validate the effectiveness of our GFR-CAM by evaluating its *primary component*, which captures the most salient features identified by the model. We benchmark its performance against a comprehensive set of state-of-the-art CAM methods. These include gradient-based approaches like Grad-CAM [33], Grad-CAM++ [4], and XGradCAM [9]; and other notable techniques including LayerCAM [12], HiResCAM [6], KPCA CAM [14], and ShapleyCAM [3]. The comparison is conducted on two distinct architectures: ResNet-50 on a subset of the ImageNet validation set and the Swin Transformer on a subset of PASCAL VOC. The quantitative results demonstrate that GFR-CAM’s *primary component* achieves performance that is highly competitive with these leading baselines.

Qualitative Evaluation of Visual Explanations: In our qualitative analysis, we employ ResNet-50 as the backbone network and generate visual explanations using its final convolutional layer. As depicted in Figure 3, GFR-CAM produces the most comprehensive activation maps in comparison to other methods. Grad-CAM struggles to identify the main target objects, often highlighting irrelevant regions, while both KPCA-CAM and Grad-CAM++ correctly localize the objects. However, our method exceeds these approaches in accurately pinpointing the relevant regions.

Performance on ResNet-50: Table 1 presents the results on the ResNet-50 architecture with ImageNet data. Our GFR-CAM demonstrates excellent faithfulness to the model’s decision-making, achieving the best scores in Average Drop (**11.17**), Increase in Confidence (**43.31**), and ADD (**45.57**). This indicates that the primary GFR-CAM component is highly effective at identifying the most critical regions for the model’s prediction. While Shapley-based methods excel in Coherency and ADCC, our method provides a competitive performance on the core faithfulness metrics, establishing it as a highly competitive CAM.

Table 1. Quantitative comparison of CAM methods on ResNet-50: leveraging the final convolutional layer

Method	AD ↓	Coh ↑	Com ↑	ADCC ↑	IC ↑	ADD ↑
GradCAM	12.73	88.12	40.13	67.22	41.97	44.19
GradCAM++[4]	13.79	96.49	48.22	71.23	42.85	37.02
HiResCAM	14.35	61.74	20.49	5.07	40.52	43.71
ShapleyCAM-E	12.28	98.43	41.27	82.59	39.33	36.23
KPCA CAM	13.98	90.21	35.67	16.23	40.01	41.11
ShapleyCAM	11.75	89.36	39.58	83.74	38.29	43.57
GFR-CAM (Ours)	11.17	91.74	30.91	20.27	43.31	45.57

Performance on Swin Transformer. In Table 2, the evaluation on the Swin Transformer with PASCAL VOC data further confirms the robustness of GFR-CAM on non-CNN architectures. Our method achieves top performance in Complexity (**32.81**) and ADD (**27.43**), producing maps that are both informative and sparse. It also remains highly competitive in Coherency (62.17), closely following the leading method, ShapleyCAM. GFR-CAM demonstrates robust performance, achieving state-of-the-art results in Complexity and ADD while remaining competitive in other key metrics like Coherency. This solidifies its efficacy as a general-purpose explanation method, applicable to diverse model architectures.

RemOve And Debias (ROAD) Benchmark [28]: The ROAD framework provides an efficient evaluation strategy for attribution methods by measuring how well saliency maps identify relevant features. The core metric is defined as:

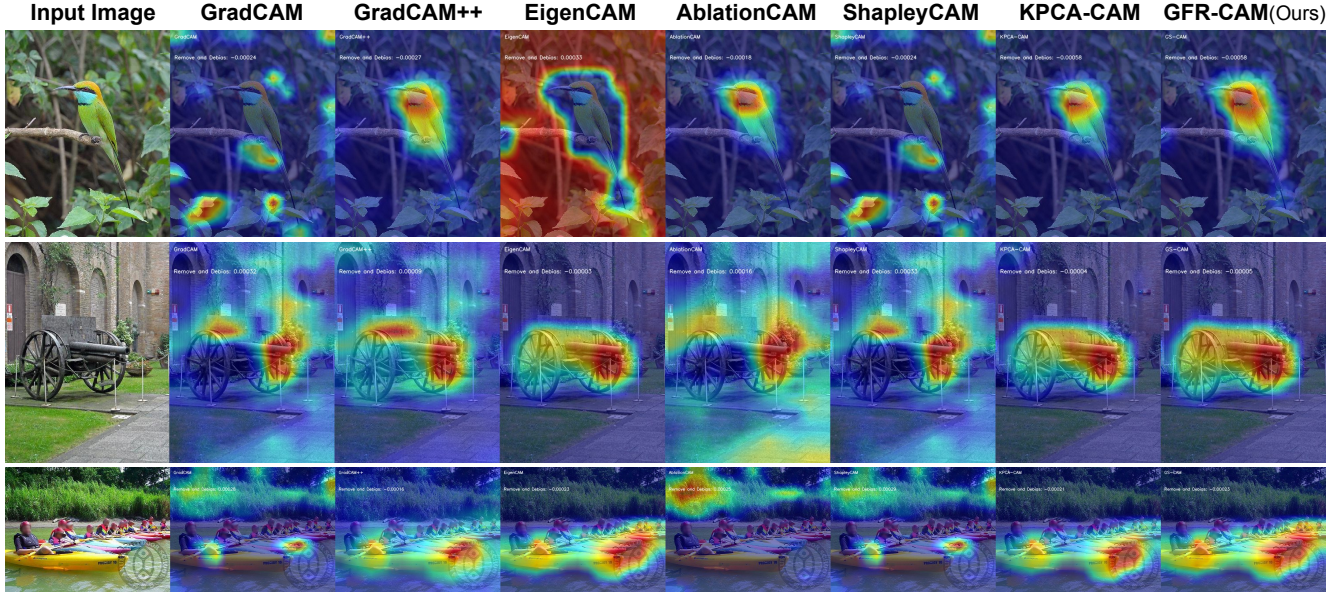


Figure 3. Qualitative comparison of CAM methods on a ResNet-50 model. The figure demonstrates that our proposed GFR-CAM generates more precise and comprehensive activation maps, accurately localizing target objects compared to other leading methods.

Table 2. Quantitative comparison of CAM methods on the Swin Transformer: targeting the initial normalization layer in the final transformer block

Method	AD ↓	Coh ↑	Com ↑	ADCC ↑	IC ↑	ADD ↑
GradCAM	57.76	41.38	23.77	6.16	5.45	23.17
GradCAM++	68.15	49.67	25.03	3.65	12.42	24.85
XGradCAM[9]	83.73	50.48	31.63	4.06	21.71	19.58
LayerCAM[12]	71.28	58.94	28.26	5.17	—	—
HiResCAM	82.54	51.48	17.25	6.19	25.12	21.45
ShapleyCAM	67.22	64.93	30.20	7.49	17.42	27.18
GFR-CAM (Ours)	72.39	62.17	32.81	6.53	13.82	27.43

$$\text{ROAD} = \frac{1}{2} (\text{LeRF} - \text{MoRF}), \quad (9)$$

where LeRF and MoRF represent model confidence scores after perturbing least and most relevant features respectively. ROAD employs a perturbation scheme that removes features in order of importance (MoRF: Most Relevant First; LeRF: Least Relevant First) and uses linear imputation with neighbor-weighted interpolation to fill perturbed regions. A more effective attribution map will identify features that cause a rapid degradation in performance, resulting in a lower Area Under the Curve (AUC) of the performance curve. Therefore, for this benchmark, a lower score is better, signifying a more accurate identification of crucial features.

Our quantitative evaluation using the ROAD benchmark is presented in Table 3. The results highlight the exceptional performance of GFR-CAM, particularly on CNN architec-

tures. Our method achieves state-of-the-art (lowest) scores on VGG-16 (59.38), ResNet-50 (-58.7), and DenseNet-161 (13.11), demonstrating its superior ability to pinpoint the most influential features in convolutional models. While AblationCAM [27] is specialized for and performs best on the transformer-based DeiT-Base, GFR-CAM’s strong and consistent results across multiple backbones confirm its robustness. The qualitative visualization in Figure 4 illustrates why GFR-CAM excels: it identifies a compact and semantically critical region, whose removal causes significant information loss and thus explains the rapid performance drop measured by the ROAD metric.

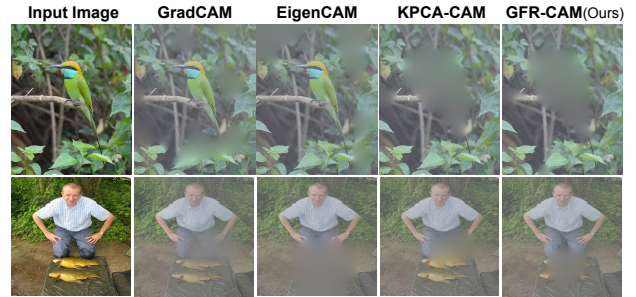


Figure 4. Qualitative results for the ROAD benchmark. GFR-CAM identifies a more compact and critical region, leading to a greater performance drop when blurred and thus a superior quantitative score.

Table 3. ROAD benchmark scores (lower is better, \downarrow) across various architectures. GFR-CAM achieves state-of-the-art performance on CNN-based models. All values are $\times 10^2$.

Model	VGG-16	ResNet50	DenseNet-161	DeiT-Base
GradCAM	80.23	-24.3	67.69	5.85
GradCAM++	82.34	-27.8	82.63	9.84
HiResCAM	65.73	3.14	85.40	16.63
ScoreCAM [39]	—	12.1	15.23	10.10
AblationCAM[27]	68.51	-18.3	32.45	3.89
EigenCAM	75.23	-16.2	23.39	11.76
ShapleyCAM	87.57	-24.2	86.42	—
KPCA-CAM	82.03	-58.3	85.32	12.51
GFR-CAM (Ours)	59.38	-58.7	13.11	10.17

3.3. Beyond the First Component: The Explanatory Power of GFR-CAM

A key limitation of existing decompositional CAMs like EigenCAM, KPCA-CAM, and UMAP-CAM [18] is their reliance on dimensionality reduction techniques that prioritize variance preservation, where subsequent components capture progressively less variance and often degrade into noise. Our GFR-CAM, however, uses Gram-Schmidt orthogonalization, which we hypothesize does not merely rank features by importance but instead isolates distinct, semantically meaningful concepts. This section validates this core advantage, demonstrating that GFR-CAM’s subsequent components provide richer, more structured explanations than their dimensionality reduction-based counterparts, effectively decomposing a model’s reasoning into a set of interpretable parts.

The qualitative results visually confirm this hypothesis. As shown in Figure 5a, when applied to a single object, GFR-CAM’s components successfully partition the object into its constituent parts, such as a lizard’s head and limbs, while competing methods produce diffuse or redundant maps. This ability extends to more complex scenarios, as seen in Figure 5b, where GFR-CAM effectively disentangles a multi-object scene by assigning each of its primary components to a distinct object instance, a task where other methods struggle and often merge the explanations.

This visual superiority is substantiated by the quantitative analysis in Table 4. For both the 2nd and 3rd components, GFR-CAM consistently outperforms EigenCAM and KPCA-CAM, leading in the majority of metrics. Notably, it achieves the best faithfulness scores (lowest Average Drop) while maintaining high coherency. While the performance of PCA-based methods degrades sharply for the 3rd component, GFR-CAM remains robust, proving its unique capability to generate a stable and high-fidelity set of orthogonal explanations beyond the primary one.

Table 4. Quantitative comparison of CAM methods using the Swin-T model for the 2nd and 3rd extracted components. We evaluate our proposed GFR-CAM against EigenCAM and KPCA-CAM across six different metrics.

2 nd Component						
Method	AD \downarrow	Coh \uparrow	Com \uparrow	ADCC \uparrow	IC \uparrow	ADD \uparrow
EigenCAM	91.31	44.61	21.61	4.74	10.95	21.73
KPCA-CAM	84.69	59.03	22.45	5.92	11.85	24.93
GFR-CAM (Ours)	80.43	58.46	22.11	6.14	12.57	26.03
3 rd Component						
Method	AD \downarrow	Coh \uparrow	Com \uparrow	ADCC \uparrow	IC \uparrow	ADD \uparrow
EigenCAM	94.23	39.27	12.45	3.56	7.16	10.30
KPCA-CAM	89.76	43.79	15.70	4.29	10.44	17.44
GFR-CAM (Ours)	87.21	46.51	16.40	4.41	9.91	18.63

3.4. Ablation Study

To validate the design choices of GFR-CAM, we conducted an ablation study on two key hyperparameters: the function family used for orthogonalization and the number of components to generate.

Choice of Orthogonalization Function. The computational cost of GFR-CAM is influenced by the complexity of the function family, $f(z)$, used in the Gram-Schmidt process. Following the methodology in [41], we evaluated linear functions (e.g., $f_1(z) = z_i$ for $i \in [d]$), polynomial functions up to degree 2 ($f_2(z) = z_i z_j$, for $i, j \in [d]$), and polynomial functions up to degree 3 ($f_3(z) = z_i z_j z_k$, for $i, j, k \in [d]$), where d is the number of features in the feature map. polynomials. As summarized in Table 5, the linear option achieves the best trade-off between runtime and attribution quality, so we adopt $f_1(z)$.

Table 5. Computational cost comparison (seconds) for different orthogonalization functions in GFR-CAM.

Family Function	$f_1(z) = z_i$	$f_2(z) = z_i z_j$	$f_3(z) = z_i z_j z_k$
GFR-CAM($f(z)$)	0.45	0.91	2.21

Number of Components. We also analyzed the trade-off between the number of generated components (m) and their explanatory value. We found that the 2nd and 3rd components consistently provided valuable, distinct visual evidence, as detailed in Section 3.3. However, components beyond the third ($m \geq 4$) introduced significant computational latency while failing to provide a commensurate improvement in interpretable information; these higher-order maps often highlighted redundant features. Based on this analysis, we determined that focusing on the first three components offers the most effective compromise between detailed multi-faceted explanations and practical efficiency.

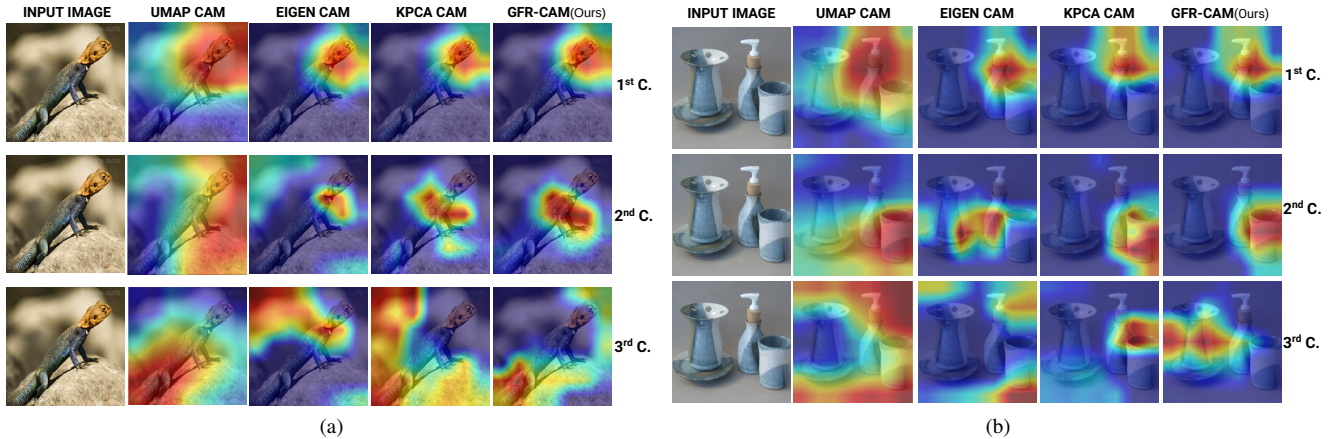


Figure 5. GFR-CAM component analysis. (a) Single object decomposition: Our method’s 2nd and 3rd components isolate distinct semantic parts (e.g., lizard’s head vs. hand/tail), providing finer-grained explanations than UMAP CAM, EigenCAM and KPCA-CAM. (b) Multi-object scene separation: GFR-CAM components disentangle complex scenes by highlighting distinct object instances, unlike competing methods that merge or fail to separate objects clearly.

4. Conclusion and Future Work

We introduced GFR-CAM, a novel gradient-free framework that addresses the fundamental limitation of existing Class Activation Maps: their inability to provide comprehensive, multi-faceted explanations of CNN decision-making. By leveraging Gram-Schmidt orthogonalization instead of PCA-based decomposition, GFR-CAM generates hierarchical activation maps where each component captures distinct, meaningful visual evidence rather than progressively noisier approximations.

Our extensive evaluation demonstrates that GFR-CAM’s primary component achieves state-of-the-art performance across multiple architectures and benchmarks, while its subsequent components provide unprecedented explanatory power. We showed that these additional maps are not artifacts but genuine insights that decompose single objects into semantic parts and systematically disentangle complex multi-object scenes. This capability represents a significant advancement over existing methods that suffer from “explanatory tunnel vision.”

The theoretical foundation of Gram-Schmidt orthogonalization ensures that each component is linearly independent and information-rich, making GFR-CAM particularly valuable for safety-critical applications where comprehensive understanding of model reasoning is essential. By providing a more complete picture of CNN decision-making, GFR-CAM sets a new standard for interpretable AI in computer vision.

Several promising directions emerge from this work. First, extending GFR-CAM to other vision architectures, including newer transformer variants and hybrid CNN-transformer models, could further validate its generalizability. Second, investigating adaptive component selec-

tion mechanisms could automatically determine the optimal number of meaningful components for different image complexities, potentially improving efficiency.

The integration of GFR-CAM with interactive visualization tools presents opportunities for enhanced human-AI collaboration, particularly in medical imaging and autonomous systems where detailed explanatory feedback is crucial. Additionally, exploring temporal extensions for video understanding and applying the hierarchical decomposition principle to other modalities (e.g., natural language processing) could broaden the impact of our orthogonalization-based approach.

Finally, investigating the theoretical connections between Gram-Schmidt decomposition and other interpretability methods, as well as developing quantitative metrics specifically designed to evaluate multi-component explanations, would strengthen the foundation for next-generation explainable AI systems that move beyond single-focus explanations toward comprehensive model understanding.

References

- [1] Leila Arras, Bruno Puri, Patrick Kahardiprja, Sebastian Lapuschkin, and Wojciech Samek. A close look at decomposition-based XAI-methods for transformer language models. *arXiv preprint arXiv:2502.15886*, 2025. 2
- [2] Shahin Atakishiyev, Mohammad Salameh, Hengshuai Yao, and Randy Goebel. Explainable artificial intelligence for autonomous driving: A comprehensive overview and field guide for future research directions. *IEEE Access*, 2024. 1
- [3] Huaiguang Cai. Cams as shapley value-based explainers. *arXiv preprint arXiv:2501.06261*, 2025. 5
- [4] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-CAM++: Gener-

- alized gradient-based visual explanations for deep convolutional networks. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847, 2018. 4, 5
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 1
- [6] Rachel Lea Draelos and Lawrence Carin. Use hirescam instead of grad-cam for faithful explanations of convolutional neural networks. *arXiv preprint arXiv:2011.08891*, 2020. 5
- [7] Mark Everingham, Luc Van Gool, Christopher K.I. Williams, John Winn, and Andrew Zisserman. The PASCAL visual object classes (VOC) challenge. *International Journal of Computer Vision (IJCV)*, 88:303–338, 2010. 4
- [8] Qizhang Feng, Ninghao Liu, Fan Yang, Ruixiang Tang, Mengnan Du, and Xia Hu. DEGREE: Decomposition based explanation for graph neural networks. In *10th Int. Conf. Learning Representations (ICLR)*, 2022. 2
- [9] Ruigang Fu, Qingyong Hu, Xiaohu Dong, Yulan Guo, Yinghui Gao, and Biao Li. Axiom-based Grad-CAM: Towards accurate visualization and explanation of CNNs. In *British Machine Vision Conference (BMVC)*, 2020. 5, 6
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 4
- [11] Habib Irani and Vangelis Metsis. Enhancing time-series prediction with temporal context modeling: A bayesian and deep learning synergy. In *The International FLAIRS Conference Proceedings*, 2024. 2
- [12] Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. LayerCAM: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing*, 30:5875–5888, 2021. 5, 6
- [13] Hyungsik Jung and Youngrock Oh. Towards better explanations of class activation mapping. In *IEEE CVF International Conference on Computer Vision*, pages 1336–1344, 2021. 4, 5
- [14] Sachin Karmani, Thanushon Sivakaran, Gaurav Prasad, Mehmet Ali, Wenbo Yang, and Sheyang Tang. KPCA-CAM: Visual explainability of deep computer vision models using Kernel PCA. In *IEEE International Workshop on Multimedia Signal Processing (MMSP)*, pages 1–5, 2024. 2, 5
- [15] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. 1
- [16] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10012–10022, 2021. 4
- [17] Luca Longo, Mario Brcic, Federico Cabitza, Jaesik Choi, Roberto Confalonieri, Javier Del Ser, Riccardo Guidotti, Yoichi Hayashi, Francisco Herrera, Andreas Holzinger, et al. Explainable artificial intelligence (xai) 2.0: A manifesto of open challenges and interdisciplinary research directions. *Information Fusion*, 106:102301, 2024. 1
- [18] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. Umap: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29):861, 2018. 7
- [19] Melkamu Mersha, Khang Lam, Joseph Wood, Ali K Alshami, and Jugal Kalita. Explainable artificial intelligence: A survey of needs, techniques, applications, and future direction. *Neurocomputing*, 599:128111, 2024. 1
- [20] Sebastian Mika, Bernhard Schölkopf, Alex Smola, Klaus-Robert Müller, Matthias Scholz, and Gunnar Rätsch. Kernel pca and de-noising in feature spaces. *Advances in Neural Information Processing Systems (NeurIPS)*, 11, 1998. 2
- [21] Juan C Moreno, VB Prasath, Dmitry Vorotnikov, Hugo Proença, and Kannappan Palaniappan. Adaptive diffusion constrained total variation scheme with application to cartoon+ texture+ edge image decomposition. *arXiv preprint arXiv:1505.00866*, 2015. 2
- [22] Mohammed Bany Muhammad and Mohammed Yeasin. Eigen-cam: Class activation map using principal components. In *IEEE International Joint Conference on Neural Networks (IJCNN)*, pages 1–7, 2020. 2
- [23] Karl Pearson. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901. 2
- [24] Mahdiah Poostchi, Kannappan Palaniappan, and Guna Seetharaman. Spatial pyramid context-aware moving vehicle detection and tracking in urban aerial imagery. In *IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6, 2017. 1
- [25] Samuele Poppi, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Revisiting the evaluation of class activation mapping for explainability: A novel metric and experimental analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2299–2304, 2021. 4, 5
- [26] Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, et al. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225*, 2017. 1
- [27] Harish Guruprasad Ramaswamy et al. Ablation-CAM: Visual explanations for deep convolutional network via gradient-free localization. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 983–991, 2020. 6, 7
- [28] Yao Rong, Tobias Leemann, Vadim Borisov, Gjergji Kasneci, and Enkelejd Kasneci. A consistent and efficient evaluation strategy for attribution methods. *arXiv preprint arXiv:2202.00449*, 2022. 5
- [29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 115:211–252, 2015. 4
- [30] Kaveh Safavigerdini, Koundinya Nouduri, Ramakrishna Surya, Andrew Reinhard, Zach Quinlan, Filiz Bunyak, Matthew R Maschmann, and Kannappan Palaniappan. Pre-

- dicting mechanical properties of carbon nanotube (CNT) images using multi-layer synthetic finite element model simulations. In *IEEE International Conference on Image Processing (ICIP)*, pages 3264–3268, 2023. 1
- [31] Kaveh Safavigerdini, Ramakrishna Surya, Andrew Reinhard, Zach Quinlan, Filiz Bunyak, Matthew R Maschmann, and Kannappan Palaniappan. Creating semi-quanta multi-layer synthetic CNT images using CycleGAN. In *IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*, pages 1–6, 2023. 1
- [32] K Safavigerdini, J Collins, R Huynh, J Fraser, and K Palaniappan. EpiX 2.0: an enhanced 3D measurement tool with integrated triangulation for cultural heritage and aerial photogrammetry applications. In *Geospatial Informatics XIV*, page PC1303706, 2024. 1, 2
- [33] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017. 2, 5
- [34] Saleh Valizadeh Sotubadi, Shyam Sasi Pallissery, and Vinh Nguyen. Multi-modal explainable artificial intelligence for neural network-based tool wear detection in machining. *Engineering Applications of Artificial Intelligence*, 144: 110141, 2025. 1
- [35] Mingzhai Sun, Jiaqing Huang, Filiz Bunyak, Kristyn Gumpfer, Gejing De, Matthew Sermersheim, George Liu, Pei-Hui Lin, Kannappan Palaniappan, and Jianjie Ma. Superresolution microscope image reconstruction by spatiotemporal object decomposition and association: Application in resolving t-tubule structure in skeletal muscle. *Optics Express*, 22(10): 12160–12176, 2014. 2
- [36] Ramakrishna Surya, Gordon L Koerner, Taher Hajilounezhad, Kaveh Safavigerdini, Martin Spies, Prasad Calyam, Filiz Bunyak, Kannappan Palaniappan, and Matthew R Maschmann. Cnt forest self-assembly insights from in-situ esem synthesis. *Carbon*, 229:119439, 2024. 1
- [37] Dan Tang, Jinjing Chen, Lijuan Ren, Xie Wang, Daiwei Li, and Haiqing Zhang. Reviewing cam-based deep explainable methods in healthcare. *Applied Sciences*, 14(10):4124, 2024. 1
- [38] Deniz Kavzak Ufuktepe, Feng Yang, Yasmin M Kassim, Hang Yu, Richard J Maude, Kannappan Palaniappan, and Stefan Jaeger. Deep learning-based cell detection and extraction in thin blood smears for malaria diagnosis. In *2021 IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*, pages 1–6. IEEE, 2021. 1
- [39] Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu. Score-CAM: Score-weighted visual explanations for convolutional neural networks, 2020. 7
- [40] Yang Yang Wang, O. V. Glinskii, Filiz Bunyak, and Kannappan Palaniappan. Ensemble of deep learning cascades for segmentation of blood vessels in confocal microscopy images. In *IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*, pages 1–7, 2021. 1
- [41] Bahram Yaghooti, Netanel Raviv, and Bruno Sinopoli. Beyond PCA: A Gram-Schmidt approach to feature extraction. In *Allerton Conference on Communication, Control, and Computing*, pages 1–8, 2024. 2, 3, 7
- [42] Bahram Yaghooti, Netanel Raviv, and Bruno Sinopoli. Gram-Schmidt methods for unsupervised feature selection. In *IEEE Conference on Decision and Control (CDC)*, pages 7700–7707, 2024. 2
- [43] Bahram Yaghooti, Kaveh Safavigerdini, Reza Hajiloo, and Hassan Salarieh. Stabilizing unstable periodic orbit of unknown fractional-order systems via adaptive delayed feedback control. *Proceedings of the Institution of Mechanical Engineers, Part I: Journal of Systems and Control Engineering*, 238(4):693–703, 2024. 1
- [44] Bahram Yaghooti, Netanel Raviv, and Bruno Sinopoli. Gram-Schmidt methods for unsupervised feature extraction and selection. *IEEE Transactions on Information Theory*, 2025. 2, 3
- [45] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2921–2929, 2016. 1, 2
- [46] Erfan Ziad, Zhuo Yang, Yan Lu, and Feng Ju. Knowledge constrained deep clustering for melt pool anomaly detection in laser powder bed fusion. In *International Conference on Automation Science and Engineering (CASE)*, pages 670–675. IEEE, 2024. 1