

RobustFormer: Noise-Robust Pre-training for Images and Videos

Ashish Bastola^{1,*}, Nishant Luitel^{2,*}, Hao Wang¹, Danda Pani Paudel², Roshni Poudel², Abolfazl Razi¹
¹Clemson University ²NAAMII

{abastol, hao9, arazi}@clemson.edu, {nishant.luitel, danda.paudel, roshani.poudel}@naamii.org.np

Abstract

While deep learning-based models like transformers, have revolutionized time-series and vision tasks, they remain highly susceptible to noise and often overfit on noisy patterns rather than robust features. This issue is exacerbated in vision transformers, which rely on pixel-level details that can easily be corrupt. To address this, we leverage the discrete wavelet transform (DWT) for its ability to decompose into multi-resolution layers, isolating noise primarily in the high frequency domain while preserving essential low-frequency information for resilient feature learning. Conventional DWT-based methods, however, struggle with computational inefficiencies due to the requirement for a subsequent inverse discrete wavelet transform (IDWT) step. In this work, we introduce RobustFormer, a novel framework that enables noise-robust masked autoencoder (MAE) pre-training for both images and videos by using DWT for efficient downsampling, eliminating the need for expensive IDWT reconstruction and simplifying the attention mechanism to focus on noise-resilient multi-scale representations. To our knowledge, RobustFormer is the first DWT-based method fully compatible with video inputs and MAE-style pre-training. Extensive experiments on noisy image and video datasets demonstrate that our approach achieves up to 8% increase in Top-1 classification accuracy under severe noise conditions in Imagenet-C and up to 2.7% in Imagenet-P standard benchmarks compared to the baseline and up to 13% higher Top-1 accuracy on UCF-101 under severe custom noise perturbations while maintaining similar accuracy scores for clean datasets. We also observe the reduction of computation complexity by up to 4.4% through IDWT removal compared to VideoMAE baseline without any performance drop.

1. Introduction

Addressing robustness in deep learning models for images and video data is crucial for practical applications, in-

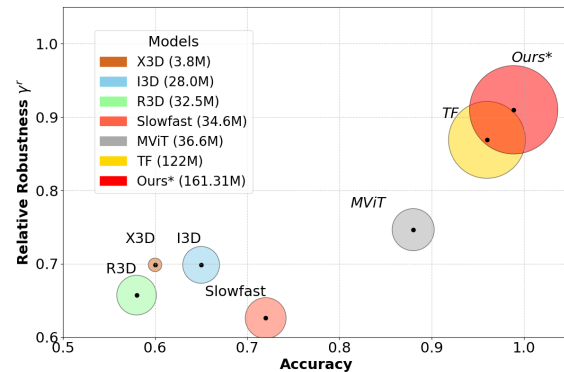


Figure 1. Accuracy vs. relative robustness (performance on corrupted vs. clean data) of action recognition models on UCF-101.

cluding surveillance, autonomous driving, and multimedia content analysis [34]. Images can suffer from distortions such as blurring, noise, and artifacts, while video data experience similar issues, compounded by the additional variability introduced over time [38]. Deep learning models, especially those trained on clean, high-quality data, are highly sensitive to these disruptions because they rely heavily on precise patterns within the data to make accurate predictions [49]. Even minor noise or inconsistencies in an image can cause the model to misinterpret important features, as these models are not inherently robust to unexpected distortions [43]. With video data, the challenge becomes even greater due to the temporal nature of the information, where each frame builds upon previous ones to create a coherent sequence [55]. For instance, temporal inconsistencies such as random noise across frames, motion blur, flickering shadows, and heat-induced turbulence can lead to abrupt changes between frames, confusing the model as it attempts to track and interpret motion or recognize actions [15, 55]. In tasks like autonomous driving, these inconsistencies are highly prevalent and pose significant risks, as the model may fail to accurately detect objects, assess distances, or predict movement patterns [41].

A primary reason for these issues lies in the pixel-based design of most deep-learning models. Although RGB pixel

*Equal contribution

representations are the standard in vision applications because they capture detailed color and spatial information, they are computationally demanding and highly sensitive to noise and domain shifts. This sensitivity arises since RGB data represent exact pixel-level details, so even slight distortions can disrupt deep learning models that depend on these patterns for accuracy. Additionally, RGB representations process both essential and irrelevant features equally, limiting their efficiency. Inspired by image compression techniques like JPEG, frequency domain transformations such as Discrete Cosine Transform (DCT) provide a solution by isolating low-frequency components (representing significant information) from high-frequency ones (less relevant details) [32, 39]. This separation reduces memory demands and increases noise resistance without compromising performance [11, 16, 32].

While pure frequency domain transformations offer these benefits, they also lead to a loss of spatial or temporal locality due to the change in basis, a consequence of the Heisenberg uncertainty principle. In contrast, Discrete Wavelet Transform (DWT) provides a more balanced approach, decomposing the signal into both high- and low-frequency components across multiple scales. This multi-resolution representation allows DWT to retain spatial and temporal information, preserving important features and contextual details while still isolating high-frequency information. Image-based methods that leverage DWT have shown promise in improving robustness against spatial noise by allowing models to focus on frequency-based features that are less affected by disturbances [30]. Applying these principles to images and videos can address both spatial and temporal (for videos) corruptions, as DWT can help models focus on low-frequency information, improving robustness to noise and other distortions across frames.

In this paper, we employ a masked autoencoding approach [50, 51], additionally, incorporating DWT for pre-training and fine-tuning on the uncorrupted datasets, followed by evaluation on their corrupted counterparts. We experiment on widely used datasets for both images and videos, where the generated corruptions include both spatial and temporal types reflecting real-world scenarios [57].

In summary, our contributions are as follows:

- We present a novel Discrete Wavelet Transform (DWT) based masked autoencoder architecture that is robust to spatial and temporal corruptions in both video and image data. To our knowledge, this is the first work to implement DWT in a masked autoencoder setting.
- We perform a comprehensive evaluation of several real-world noise types with varying severity levels in large-scale benchmark datasets.
- We demonstrate that our method performs equally well and in many cases better than the commonly used IDWT variant, which requires comparatively more compute

resource and makes architectures overly complicated.

2. Related Works

2.1. Masked Video Training

BERT introduced the concept of masked language modeling (MLM), a novel pre-training objective for natural language understanding, which led to significant advances in NLP [26] and further inspiring extensions including RoBERTa [35], ALBERT [28], and ELECTRA [6], extended this masked training approach. Also inspired by BERT's masked token prediction, the concept of masked image modeling (MIM) emerged for computer vision tasks. iGPT [5] first adapted the Transformer model for image data by treating images as pixel sequences and using a similar masked prediction task. BEiT [1] extended this concept by proposing a discrete variational autoencoder (dVAE) to tokenize image patches, predicting masked tokens in a BERT-style pre-training. These approaches demonstrated that MIM could effectively learn transferable representations for various downstream tasks. Masked Autoencoders (MAE) [19] refined MIM by reconstructing masked regions in images, achieving state-of-the-art results on multiple benchmarks using Vision Transformers (ViTs) and large-scale unlabeled data.

Building on masked training in language and image domains, Masked Video Modeling (MVM) has emerged to capture both temporal and spatial representations from video data. VideoMAE [50] extends MAE to videos by applying random masking across both spatial and temporal dimensions, enabling the model to learn robust spatio-temporal representations, as shown in Figure 2 (a). VideoMAE achieves state-of-the-art performance on numerous video classification tasks, highlighting the effectiveness of masked training for video data. Similarly, BEVT (BEiT for Video) [52] extends BEiT's tokenization and masked prediction strategy to video patches, achieving strong results in action recognition tasks.

2.2. Discrete Wavelet transform(DWT)

Wavelets are essential in time-frequency analysis and signal processing tasks, such as anti-aliasing and detail restoration, through Discrete Wavelet Transform (DWT) and its inverse (IDWT) [36]. Early studies combined wavelets with shallow neural networks for function approximation and classification, optimizing wavelet parameters within the network [48]. Recently, deeper networks have adopted wavelets for image classification, although these integrations can be computationally intensive [7]. The Multilevel Wavelet CNN (MWCNN) [33] applies Wavelet Packet Transform (WPT) for image restoration, handling both low- and high-frequency components. Similarly, the Convolutional-Wavelet Neural Network (CWNN) [10] uti-

lizes dual-tree complex wavelet transform to reduce noise in SAR images while preserving crucial features, though within a simplified two-layer structure. Wavelet Pooling [54] employs a two-level DWT for pooling and combines DWT/IDWT for backpropagation, deviating from conventional gradient methods.

A key advantage of wavelets is their ability to provide sparse representations, which enables more efficient processing and faster model training [31, 54, 55]. DWT-based training reduces redundancy, making it well-suited for real-time video classification tasks where small shifts in data can significantly affect model output. Unlike other approaches, our method avoids reconstruction after the wavelet transform, thus saving considerable computational resources, especially when handling images and videos. This streamlined process enhances efficiency while maintaining robust performance in noisy environments.

2.3. Noise Robustness

Robustness in video classification has been an ongoing research focus, initially centered on adapting image robustness techniques to video data. For instance, [22] introduced benchmarks for evaluating image classifiers’ robustness against corruptions, a framework later adapted to videos by [24] to assess model performance under various perturbations. Addressing the challenge of temporal consistency in video data, models such as the two-stream network [45] and 3D convolutions [3] were developed, with data augmentation methods like temporal jittering and frame dropping [44] further enhancing robustness.

Since DWT has been widely applied in noisy data reconstruction, achieving notable gains in both noise removal and complex tasks like deblurring [31, 54], our approach builds on the methodologies of WaveCNets [30] by incorporating DWT for noise-robust representation. Figure 2(b) shows IDWT configuration which computes attention using only low-frequency components of the query and key for noise filtering (similar to [55]) during attention computation. Similar to WaveCNets, we discard high-frequency components during the initial downsampling layer but extend this technique to 3D-DWT to handle temporal decomposition. Our framework leverages 3D-DWT’s ability to decompose both spatially and temporally, enabling effective noise handling across frames and enhancing robustness over the entire video sequence. As shown in Figure 1, our approach improves both accuracy and relative robustness compared to several action recognition models on the UCF-101 dataset.

3. Methods

Compared to the classical VideoMAE [49], our proposed method enhances robustness by employing 3D-DWT for spatio-temporal analysis; in contrast to other IDWT-based methods, our approach eliminates the computationally ex-

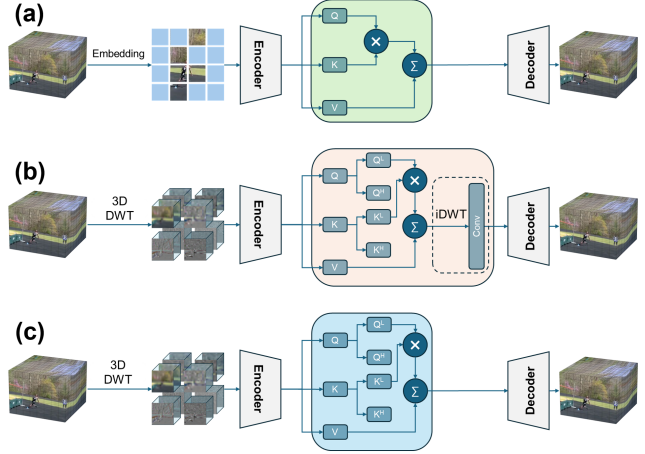


Figure 2. Comparison between different architectures designed for video tasks. (a) is the regular masked autoencoder [49], (b) is the DWT-based architecture with IDWT module, and (c) is our proposed method.

pensive IDWT step, significantly improving efficiency, as illustrated in Figure 2(c).

Our implementation incorporates a two-step wavelet transform. First, we utilize DWT’s efficient downsampling capability to generate embeddings for video patches. By leveraging 3D-DWT, we can process video data more effectively than stacking coefficients from 2D-DWT. This approach enables simultaneous handling of both temporal and spatial noise. The robust feature extraction achieved is thus the result of noise-adaptive compression in the initial DWT phase, noise filtering during attention computation well as masked pretraining which allows to recover back to RGB without requiring additional IDWT step. The forward and backward propagation processes for this step are detailed below. For experiments with images we use 2D-DWT instead on 3-dimensional version while every other aspects remaining same.

3.1. Wavelet Embedding

Let $\mathbf{X}^{t \times h \times w \times c}$ be the original downsampled and clipped video sequence of sequence length t , height h , width w and channels c . We start with the generation of transformation matrices based on the selected wavelet filters. This step is critical as it defines the low-pass and high-pass filters that are applied in the transformation process. We define the wavelet by low-pass (L) and high-pass (H) filter coefficients as $\mathbf{L} = [l_1, l_2, \dots, l_n]$, and $\mathbf{H} = [h_1, h_2, \dots, h_n]$. In this work, we employ Haar wavelets due to their computational efficiency and effectiveness in noise reduction.

Now we construct the transformation matrices for each dimension (depth, height, and width). For simplicity, only the depth (time) dimension’s matrix construction is shown here. We can replicate this for the height and width dimen-

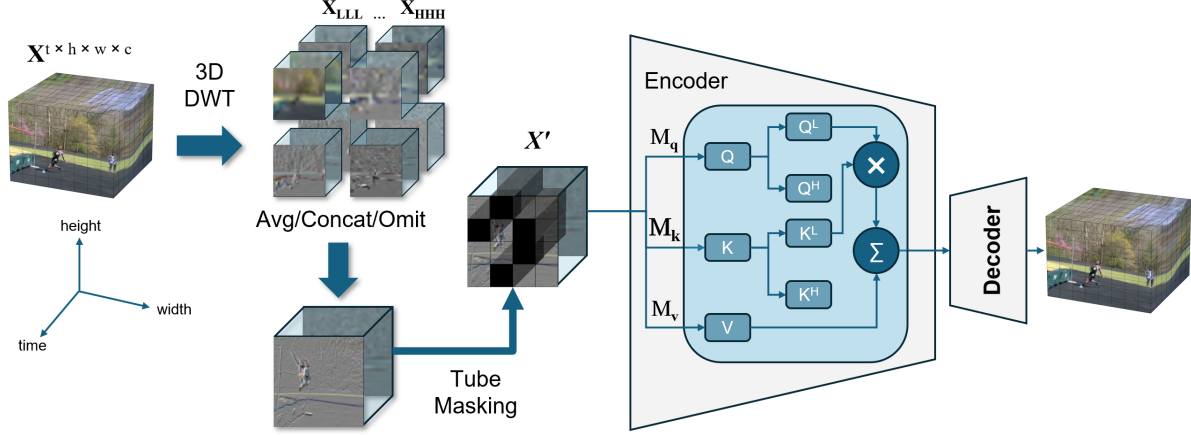


Figure 3. The framework of RobustFormer. Our approach integrates spatio-temporal tube masking as well as multi-resolution feature transformation using DWT to handle real-world noise types.

sions:

$$\begin{aligned} \mathcal{L}_k &= \text{ConstructMatrix}(\mathbf{L}, k) \\ \mathcal{H}_k &= \text{ConstructMatrix}(\mathbf{H}, k) \end{aligned} \quad (1)$$

where, $k \in \{H, W, D\}$ represent each of the individual dimension. The forward pass thus involves applying these matrices to decompose the input data into its respective frequency sub-bands. The input tensor \mathbf{X} is thus decomposed into eight sub-bands using the transformation matrices for depth ($\mathcal{L}_D, \mathcal{H}_D$), height ($\mathcal{L}_H, \mathcal{H}_H$), and width ($\mathcal{L}_W, \mathcal{H}_W$), as shown in Figure 3.

$$\begin{aligned} \mathbf{X}_{LLL} &= \mathcal{L}_D^T \mathcal{L}_H^T \mathcal{L}_W^T \mathbf{X}, \\ \mathbf{X}_{LLH} &= \mathcal{L}_D^T \mathcal{L}_H^T \mathcal{H}_W^T \mathbf{X}, \\ &\dots, \\ \mathbf{X}_{HHH} &= \mathcal{H}_D^T \mathcal{H}_H^T \mathcal{H}_W^T \mathbf{X} \end{aligned} \quad (2)$$

where, the matrix construction leads to generating matrix $\mathcal{M}_k \in \{\mathcal{L}_k, \mathcal{H}_k\}$ for dimension corresponding to k . The resulting matrix \mathcal{M}_k will have dimensions $(\lceil \frac{k}{2} \rceil, k)$, where $\lceil \cdot \rceil$ is the ceiling function that handles cases where k is odd. for filter $\mathbf{F} \in \{\mathbf{L}, \mathbf{H}\}$ and $\mathbf{F} = [f_1, f_2, \dots, f_n]$. The ConstructMatrix function works as follows:

$$\mathcal{M}_k^{ij} = \begin{cases} f_{t+1} & \text{if } j = 2i + t, t < n \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where, where t indexes the filter coefficients $[f_1, f_2, \dots, f_n]$.

3.2. Attention Computation

In the wavelet embedding step, we average, concatenate, or omit high-frequency components to reduce noise and smooth the input representation. However, the subsequent latent transformation (3D convolution with tube masking) may amplify residual noise during attention calculation. To minimize redundant details, we aim to reduce noise as much

Algorithm 1: RobustFormer Pseudocode

Function PatchEmbedDWT(x):

```

for  $i \leftarrow 1$  to  $\text{dims}$  do
   $dwt[i] \leftarrow \text{DWT\_2D}(x[i])$ 
 $\text{embeddings} \leftarrow \text{Masking}(\text{Conv}(\text{wavelet\_ops}(dwt)))$ 
return  $\text{embeddings}$ 

```

Main:

```

 $\text{embeddings} \leftarrow \text{PatchEmbedDWT}(x)$ 
 $Z \leftarrow \text{Encoder}(\text{embeddings})$ 

```

Pre-training:

```

 $\hat{x} \leftarrow \text{Decoder}(Z)$ 
 $\mathcal{L}_{\text{recon}} \leftarrow \text{MSE}(x, \hat{x})$ 

```

Fine-tuning:

```

 $\hat{y} \leftarrow \text{RobustFormer}(Z)$ 
 $\mathcal{L}_{\text{task}} \leftarrow \text{CrossEntropy}(\hat{y}, y)$ 

```

Update model parameters via backpropagation

as possible during training. Prior work [55] has shown that noise before the attention mechanism distorts correlation scores by increasing values for unrelated pairs (due to random alignment) and decreasing values for related pairs (as noise weakens true alignment), ultimately impairing classification accuracy. To mitigate this, we follow [55] and fully omit high-frequency components, ensuring a smoother latent representation before attention calculation.

Suppose $\mathbf{X}' = \{x'_1, x'_2, \dots, x'_T\}$ denotes the output of the tube-masked 3D convolution, representing deep encoded feature representations. Each x'_i is a tensor in $\mathbb{R}^{\frac{t}{2} \times e \times \frac{h}{p} \times \frac{w}{p}}$, where t' denotes the tubelet size, e the embedding dimension, and p the patch size. The attention mechanism elements \mathbf{Q} , \mathbf{K} , and \mathbf{V} are computed by passing \mathbf{X}' through three distinct linear layers.

We then perform 1D-DWT for each \mathbf{Q} , \mathbf{K} and \mathbf{V} by obtaining the low and high pass filters corresponding to a specific wavelet(Haar) similar to [55] as follows,

$$\begin{aligned} \mathbf{Q}_L &= \mathcal{L}_Q^T \mathbf{Q}, \mathbf{Q}_H = \mathcal{H}_Q^T \mathbf{Q} \\ \mathbf{K}_L &= \mathcal{L}_K^T \mathbf{K}, \mathbf{K}_H = \mathcal{H}_K^T \mathbf{K} \end{aligned} \quad (4)$$

We calculate these segregated components similar to eq (2) and omit \mathbf{Q}_H and \mathbf{K}_H . Note that we also omit DWT computation for \mathbf{V} , so that the attention computation can attend to some useful high level components from \mathbf{V} (that are omitted in \mathbf{Q} and \mathbf{K}) in the input, as shown in Figure 3. With this we can avoid performing additional IDWT step as in [55] which makes the architecture complicated and computationally intensive.

We finally compute our attention score as follows:

$$\text{Attention} = \text{softmax}\left(\frac{\mathbf{Q}_L \mathbf{K}_L^T}{\sqrt{d_k}}\right) \mathbf{V} \quad (5)$$

The backward function is defined for each DWT computation similar to 3D DWT computation.

With this we create various configurations of RobustFormer that involves various operations over DWT blocks. **RF-A** indicates RobustFormer with DWT averaging, **RF-AA** refers to averaging as well as DWT attention to every attention layer as mentioned in 5. Similarly we define **RF-O**, **RF-OA**, **RF-C** and **RF-CA** for omit, omit using DWT attention, concat and concat using DWT attention.

4. Experiments

4.1. Datasets

We evaluate our method across five datasets: **Kinetics-400** [25], **UCF-101** [47] and **HMDB-51** for video tasks, and **ImageNet-1K** [8] and **ImageNet-Tiny-200** [29], for image tasks. The **Kinetics-400** dataset includes 240k training videos and 20k validation videos, each clip lasting 10 seconds. The dataset focuses on human-centered actions, covering a wide range of interactions. The **UCF-101** dataset, while smaller with 13.3k action videos across 101 categories, offers high diversity in actions and significant variations in camera motion, object appearance, pose, background, and lighting conditions. For image-based experiments, we use the standard **ImageNet-1K** dataset.

For noisy evaluation on video datasets, we assess our models under real-world perturbations following [43], along with additional noise types and intensity configurations. Video noise types include four categories: **Noise** (Gaussian, shot, impulse, speckle), **Blur** (zoom, motion, defocus), **Digital** (JPEG and MPEG compression artifacts), **Temporal** (jumble, box jumble), and **Camera Motion** (static rotation, dynamic rotation, translation). Each noise type is tested across five severity levels, resulting in a total of **70** unique noise configurations. In the image experiments, we evaluate model robustness using the ImageNet-P

Table 1. Comparison of robustness of various Models using mean corruption error (mCE), normalized by AlexNet values, on ImageNet-C. ‘*’ represents models with wavelet based strategy.

| Models | IN-C (mCE) ↓ |
|------------------------------|--------------|
| CNNs | |
| ResNet-50 [21] | 76.7 |
| ResNeXt50-32x4d [56] | 64.7 |
| VGG-11 [46] | 93.5 |
| VGG-19 [46] | 88.9 |
| VGG-19 + BN [46] | 81.6 |
| ANT [42] | 63.0 |
| EWS [17] | 63.0 |
| WRResNet-18* [31] | 80.8 |
| WRResNet-101* [31] | 65.8 |
| ViTs | |
| PVT-Large [53] | 59.8 |
| BiT-mr101x3 [27] | 58.3 |
| MAE-ViT-B [50] | 58.8 |
| Robust Formers (Ours) | |
| RobustFormer-A | 55.3 |
| RobustFormer-AA | 55.7 |
| RobustFormer-O | 55.3 |
| RobustFormer-OA | 55.9 |
| RobustFormer-C | 57.7 |
| RobustFormer-CA | 58.3 |

Table 2. Comparison of RobustFormer models with varying configurations using absolute mean Flip Probability(mFP) metric on Imagenet-P. ‘↓’ indicates lower is better.

| Model | Params | Flops | IN-P(mFP) ↓ |
|-----------------|--------|-------|-------------|
| MAE-ViT-B [50] | 111.7M | 36.7G | 13.4 |
| RobustFormer-A | 111.2M | 36.6G | 12.7 |
| RobustFormer-AA | 111.2M | 36.6G | 12.7 |
| RobustFormer-O | 111.6M | 37.5G | 12.2 |
| RobustFormer-OA | 111.6M | 37.5G | 12.4 |
| RobustFormer-C | 111.2M | 38.0G | 12.9 |
| RobustFormer-CA | 111.2M | 38.0G | 12.9 |

and ImageNet-C benchmarks [23]. ImageNet-C includes 15 types of corruptions with five severity levels, totaling **75** distinct corruptions. ImageNet-P adds 10 perturbation types with 30 intensity variations, leading to **300** different noise scenarios. All pre-training, fine-tuning, and validation are conducted on clean datasets, without noise augmentation.

4.2. Training

We use ViT-Base (ViT-B/16, input size 224) [9] as the backbone for all training and ablation studies, covering both image and video datasets. ViT-B, while over three times larger than ResNet-50 [21], is lightweight compared to other transformer architectures, making it an efficient yet powerful choice for our experiments. For image tasks, pre-

Table 3. Comparison of robustness of various networks under every type of distortion in imagenet-C using the corruption error metric(CE) [22] using AlexNet as the baseline. “*” represents models with Wavelet based strategy. The best models are marked as **BOLD** for each corruption category. Lower values are better.

| Network | Noise | | | Blur | | | Weather | | | | Digital | | | | |
|---------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | Gauss. | Shot | Impulse | Defocus | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
| AlexNet | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| SqueezeNet | 107 | 106 | 105 | 100 | 103 | 101 | 100 | 101 | 103 | 97 | 97 | 98 | 106 | 109 | 134 |
| VGG-11 | 97 | 97 | 100 | 92 | 99 | 93 | 91 | 92 | 91 | 84 | 75 | 86 | 97 | 107 | 100 |
| VGG-19 | 89 | 91 | 95 | 89 | 98 | 90 | 90 | 89 | 86 | 75 | 68 | 80 | 97 | 102 | 94 |
| VGG-19+BN | 82 | 83 | 88 | 82 | 94 | 84 | 86 | 80 | 78 | 69 | 61 | 74 | 94 | 85 | 83 |
| ResNet-18 | 87 | 88 | 91 | 84 | 91 | 87 | 89 | 86 | 84 | 78 | 69 | 78 | 90 | 80 | 85 |
| ResNet-50 | 80 | 82 | 83 | 75 | 89 | 78 | 80 | 78 | 75 | 66 | 57 | 71 | 85 | 77 | 77 |
| WRResNet-18* | 80.2 | 80.5 | 80.5 | 79.7 | 89.8 | 83.6 | 84.8 | 84.9 | 80.8 | 73.9 | 66.3 | 75.1 | 88.2 | 75.1 | 88.6 |
| WRResNet-101* | 64.3 | 65.9 | 65.0 | 63.6 | 80.4 | 68.2 | 71.8 | 71.6 | 67.5 | 60.2 | 50.1 | 62.9 | 74.6 | 52.3 | 67.9 |
| AugMix | 67 | 66 | 68 | 64 | 79 | 59 | 64 | 69 | 68 | 65 | 54 | 57 | 74 | 60 | 66 |
| MAE-ViT | 56.7 | 57.3 | 56.1 | 61.9 | 75.4 | 58.1 | 68.6 | 57.7 | 63.4 | 53.7 | 45.4 | 45.0 | 68.9 | 55.1 | 58.7 |
| RF-O | 51.9 | 52.3 | 52.0 | 58.0 | 71.7 | 56.0 | 65.2 | 57.1 | 58.9 | 54.4 | 44.2 | 45.4 | 63.4 | 45.6 | 56.5 |
| RF-OA | 51.6 | 51.8 | 51.8 | 59.2 | 72.2 | 55.9 | 64.8 | 56.9 | 60.2 | 55.8 | 44.4 | 45.3 | 64.2 | 47.9 | 56.4 |
| RF-A | 52.9 | 53.5 | 52.4 | 58.7 | 68.7 | 55.9 | 64.4 | 56.6 | 58.6 | 55.3 | 44.5 | 43.7 | 61.6 | 47.8 | 55.2 |
| RF-AA | 52.8 | 53.9 | 52.8 | 58.9 | 69.1 | 57.4 | 64.1 | 56.8 | 59.5 | 55.5 | 44.9 | 44.8 | 61.9 | 48.0 | 55.9 |
| RF-C | 53.6 | 54.6 | 52.1 | 62.9 | 76.6 | 57.4 | 69.2 | 55.9 | 61.9 | 52.8 | 44.2 | 45.4 | 69.1 | 52.2 | 57.3 |
| RF-CA | 54.6 | 55.6 | 53.4 | 62.8 | 75.9 | 57.2 | 69.2 | 57.4 | 62.6 | 53.3 | 44.4 | 47.1 | 69.4 | 53.5 | 57.5 |

Table 4. Comparison of robustness of Video Classification models on 5 different corruption categories: Noise, Blur, Temporal, Digital and Camera [43] evaluated on Kinetic-400P benchmark dataset. γ^a and γ^r are absolute and relative robustness scores of the models averaged across all the severity level and noise types in a particular category. The best models are marked as **BOLD** for each corruption category. For both, higher is better.

| Network | Params | Flops | Noise | | Blur | | Temporal | | Digital | | Camera | | Mean | |
|---------------|--------|--------|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|
| | | | γ^a | γ^r | γ^a | γ^r | γ^a | γ^r | γ^a | γ^r | γ^a | γ^r | γ^a | γ^r |
| R3D [18] | 32.5M | 55.1G | .71 | .61 | .78 | .70 | .98 | .97 | .91 | .88 | .89 | .85 | .85 | .80 |
| I3D [4] | 28.0M | 75.1G | .72 | .61 | .80 | .72 | .97 | .96 | .91 | .87 | .89 | .85 | .86 | .80 |
| SF [14] | 34.6M | 66.6G | .73 | .64 | .81 | .73 | .96 | .94 | .91 | .87 | .88 | .84 | .85 | .80 |
| X3D [13] | 3.8M | 5.15G | .71 | .62 | .81 | .75 | .96 | .94 | .90 | .86 | .85 | .84 | .85 | .80 |
| TF [2] | 122M | 196G | .87 | .84 | .91 | .87 | .98 | .96 | .94 | .94 | .95 | .93 | .91 | .88 |
| MViT [12] | 36.6M | 70.7G | .93 | .91 | .86 | .82 | .96 | .95 | .94 | .93 | .92 | .92 | .93 | .91 |
| VideoMAE [50] | 94.8M | 167.7G | .95 | .92 | .89 | .79 | .97 | .95 | .97 | .95 | .84 | .69 | .92 | .86 |
| RF-A | 93.2M | 162.8G | .94 | .90 | .90 | .82 | .97 | .95 | .97 | .94 | .85 | .69 | .93 | .86 |
| RF-AA | 93.2M | 162.8G | .94 | .91 | .90 | .80 | .96 | .95 | .97 | .95 | .85 | .70 | .92 | .86 |
| RF-O | 93.2M | 160.4G | .94 | .91 | .90 | .81 | .97 | .97 | .96 | .96 | .96 | .68 | .95 | .87 |
| RF-OA | 93.2M | 160.4G | .95 | .91 | .90 | .80 | .97 | .96 | .97 | .94 | .86 | .70 | .93 | .86 |

training is conducted on ImageNet using a masking ratio of 0.75 and norm-pix-loss as the reconstruction objective, following the MAE framework [20]. Training spans 400 epochs with a batch size of 64, distributed across 20 nodes equipped with P100 GPUs (2 GPUs and 28 cores per node). Fine-tuning is performed for 100 additional epochs in full float-32 precision without autocast.

For video datasets, we apply the same foundational framework with adaptations tailored to temporal data. Pre-training on UCF-101 and HMDB-51 involves 800 epochs using 16-frame sequences, a batch size of 8, a tube masking

ratio of 0.9, a sampling rate of 4, and a decoder depth of 4. The same node configuration is used as in the image tasks. Fine-tuning for these datasets is conducted over 100 epochs. For larger-scale datasets like Kinetics-400P, training is distributed across 12 nodes with V100 GPUs (2 GPUs and 40 cores per node) for 400 epochs, maintaining a tube masking ratio of 0.9, sampling rate of 4, and decoder depth of 4 but with an increased batch size of 16. We use DeepSpeed [40] during video fine-tuning for memory-saving optimizations. Our evaluation focuses on the model’s robustness to real-world noise, including shot noise, rain noise, Gaussian

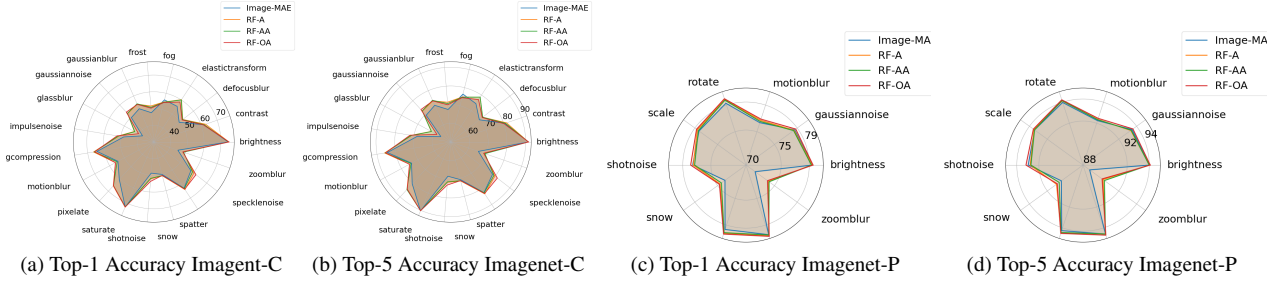


Figure 4. Comparison of Top-1 and Top-5 Accuracy for Imagenet-C and Imagenet-P. Accuracy for Imagenet-C are averaged across severity levels

Table 5. Relative robustness for individual corruptions (γ_p^r) averaged across all severity levels for Kinetics-400 dataset. The best models are marked as **BOLD** for each perturbation.

| Perturbation | R3D | I3D | SF | X3D | TF | MViT | VideoMAE | RF-A | RF-AA | RF-OA |
|---------------|------------|-----|-----|-----|------------|------|------------|------------|------------|------------|
| Defocus Blur | .77 | .75 | .80 | .85 | .80 | .83 | .87 | .88 | .87 | .88 |
| Motion Blur | .63 | .60 | .64 | .66 | .75 | .82 | .85 | .89 | .89 | .88 |
| Zoom Blur | .71 | .74 | .80 | .85 | .89 | .90 | .65 | .65 | .65 | .65 |
| Gaussian | .47 | .46 | .36 | .49 | .75 | .87 | .88 | .85 | .85 | .87 |
| Shot | .78 | .79 | .82 | .79 | .94 | .95 | .96 | .96 | .98 | .96 |
| Impulse | .44 | .42 | .34 | .46 | .76 | .87 | .87 | .83 | .84 | .86 |
| Speckle | .75 | .75 | .70 | .75 | .91 | .95 | .95 | .95 | .95 | .96 |
| Compression | .93 | .90 | .92 | .89 | .94 | .94 | .95 | .95 | .95 | .96 |
| Static Rotate | .67 | .65 | .70 | .71 | .82 | .87 | .85 | .86 | .86 | .87 |
| Rotate | .92 | .93 | .83 | .87 | .97 | .90 | .77 | .80 | .83 | .82 |
| Translate | .97 | .96 | .89 | .93 | .99 | .95 | .43 | .41 | .41 | .42 |
| Jumbling | .97 | .96 | .89 | .91 | .95 | .91 | .89 | .89 | .92 | .90 |
| Box Jumbling | .99 | .97 | .96 | .96 | .97 | .94 | .99 | .99 | .99 | .99 |

Table 6. Parameters and FLOPs for Video models.

| Model | Flops |
|---------------|--------|
| VideoMAE [49] | 167.7G |
| RF-A | 162.8G |
| RF-AA | 162.8G |
| RF-O | 160.4G |
| RF-OA | 160.4G |
| RF-I | 166G |
| RF-IA | 167.5G |

noise, packet loss, speckle noise, and tampering noise, as identified in prior benchmarks [37, 57]. These perturbations mirror real-world conditions, where noise is often diverse and unstructured. For instance, practical noise types like rain, critical for autonomous driving, and packet loss, which may affecting video-based applications. Additionally, testing on image datasets shows that our DWT-based approach is not only robust but also achieves a 2.5% accuracy improvement on clean datasets.

4.3. Evaluation

For evaluating noise robustness on Imagenet-C we compute the mean corruption error(mCE) as follows:

$$CE_c^f = \left(\sum_{s=1}^5 E_{s,c}^f \right) / \left(\sum_{s=1}^5 E_{s,c}^{\text{Baseline}} \right) \quad (6)$$

where, $E_{s,c}^f$ is the top-1 error for corruption c and severity s for the model f and $E_{s,c}^{\text{Baseline}}$ is the same for the baseline model for which we use AlexNet. The mean of the corruption error(mCE) across corruptions of RobutFormer and other comparable models is shown in Table 1, whereas, CE for the individual corruptions is shown in Table 3. Variants of RobustFormer methods always beats MAE [20] and

CNN based approaches. Similarly, for Imagenet-P we compute the flip probability of the model ‘ f ’ on ‘ m ’ perturbation sequences S as

$$FP_p^f = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j=2}^n \mathbb{1} \left(f(x_j^{(i)}) \neq f(x_1^{(i)}) \right) \quad (7)$$

where $x_1^{(i)}$ is the clean image and $x_j^{(i)}$ with ($j > 1$) are the perturbed images of $x_1^{(i)}$. We then average across all perturbations to get mean flip error(mFP). The comparison of RobustFormer variants with Image-MAE [20] is shown in Table 2 where all the RF perform better. Figure 4, shows both the average top-1 and top-5 accuracies on Imagenet-P and Imagenet-C noise benchmarks which demonstrates that our method is robust to almost all evaluated noise types for both benchmarks.

To measure video robustness we use two metrics for relative and absolute accuracy drops. After training model f , we first compute the accuracy A_c^f on the clean set and accuracy $A_{p,s}^f$ for perturbation p and severity s . The absolute robustness and relative robustness are then computed as $\gamma_{p,s}^a = 1 - (A_c^f - A_{p,s}^f) / 100$ and $\gamma_{p,s}^r = 1 - (A_c^f - A_{p,s}^f) / A_c^f$. The aggregated performance for all models can be thus achieved by averaging across severity levels to get γ_p^a and

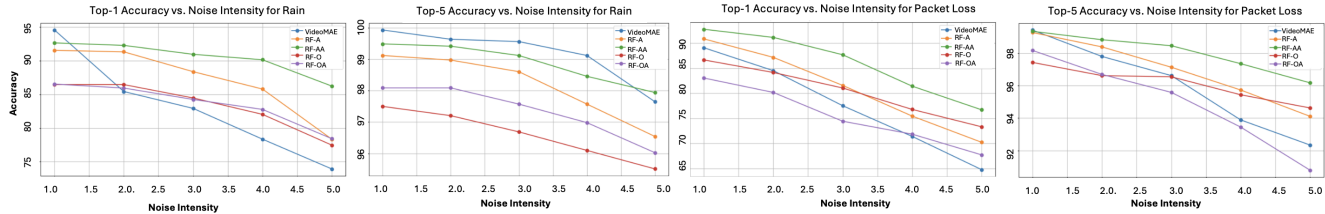


Figure 5. Comparison of Top-1 and Top-5 Accuracy for different severity levels across rain and packet loss noise for UCF-101 dataset

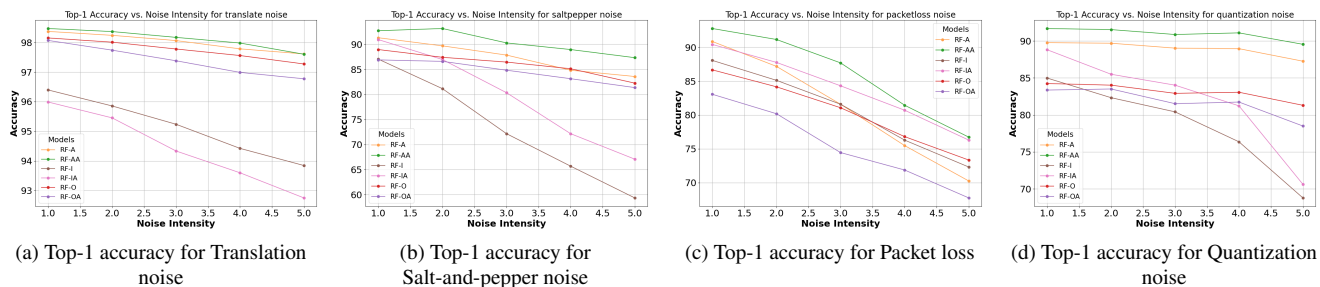


Figure 6. Ablation experiments: Comparison of Top-1 accuracy across various noise types for variants with average, IDWT, and DWT attention.

γ_p^r . Table 4 shows γ_p^a and γ_p^r , which are averaged γ_p^a and γ_p^r across the perturbation categories p for Kinetics-400 dataset. Our RF variants even without augmented training demonstrated significant robustness compared to some of the methods that used augmentation during training. Similarly, Table 5 shows the relative robustness scores γ_p^r for individual perturbations averaged across 5 severity levels. RF models here demonstrated better relative robustness to recent state-of-the-art methods. Moreover, top-1 and top-5 accuracy performance on additional noise types like rain, packetloss for UCF-101 dataset are shown in Figure 5, where RF methods are significantly superior especially in noise with high severity levels. Additional results (including for ImageNet-Tiny-200 dataset evaluated on our custom noise) are kept in supplementary.

4.4. Ablation Study

We ablate the effects of the IDWT step and the proposed DWT-based attention module across different RobustFormer variants. Our preferred model, **RF-AA**, achieves competitive and often superior top-1 and top-5 accuracy compared to all other configurations. The variants **RF-I** and **RF-IA** correspond to adding the IDWT reconstruction step and adding both IDWT and DWT-attention, respectively. While the IDWT step is commonly used in traditional DWT pipelines, it offers only limited gains when used alone.

Our experiments show that the **DWT-attention module is the primary source of robustness improvements**. Averaging alone struggles with high-impulse corruptions such as salt-and-pepper noise, whereas adding DWT attention

yields substantial accuracy gains. Moreover, DWT attention continues to improve performance even within IDWT-based variants (Fig. 6b–d). Removing the IDWT step also reduces computational cost by **7.1 GFLOPs** when comparing RF-OA with RF-IA (Table 6). Overall, these ablations highlight the central role of DWT-based attention in achieving strong robustness under diverse noise conditions.

5. Discussion and Conclusion

Noise-robust models are essential for many real-world applications, where noise patterns can significantly impact performance. While some deep learning models rely on heavy augmentation and others on rule-based filtering, these approaches are limited in scalability due to the inherent randomness and diversity of noise. Our approach enables efficient decomposition of spatio-temporal information, selectively focusing on low-frequency components to enhance robustness across diverse noise types without the need for exhaustive augmentation or complex filtering rules, thus making few assumptions about the evaluation noises. We validated our method on both image and video benchmark datasets with benchmark noise types for both and show our model’s significance compared to prior works.

References

- [1] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers, 2022. 2
- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, volume 2, page 4, 2021. 6

- [3] João Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4733, 2017. 3
- [4] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. 6
- [5] Mark Chen, Alec Radford, Jeff Wu, Heewoo Jun, Prafulla Dhariwal, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International Conference on Machine Learning*, 2020. 2
- [6] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: Pre-training text encoders as discriminators rather than generators. In *ICLR*, 2020. 2
- [7] DDN De Silva, HWMK Vithanage, KSD Fernando, and I Thilini S Piyatilake. Multi-path learnable wavelet neural network for image classification. In *Twelfth International Conference on Machine Vision (ICMV 2019)*, volume 11433, pages 459–467. SPIE, 2020. 2
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 5
- [10] Yiping Duan, Fang Liu, Licheng Jiao, Peng Zhao, and Lu Zhang. Sar image segmentation based on convolutional-wavelet neural network and markov random field. *Pattern Recognition*, 64:255–267, 2017. 2
- [11] Ahmet M Eskicioglu and Paul S Fisher. Image quality measures and their performance. *IEEE Transactions on communications*, 43(12):2959–2965, 1995. 2
- [12] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6824–6835, 2021. 6
- [13] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 203–213, 2020. 6
- [14] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. 6
- [15] Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal learners. *Advances in neural information processing systems*, 35:35946–35958, 2022. 1
- [16] Ge Gao, Pei You, Rong Pan, Shunyuan Han, Yuanyuan Zhang, Yuchao Dai, and Hojae Lee. Neural image compression via attentional multi-scale back projection and frequency decomposition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14677–14686, 2021. 2
- [17] Yong Guo, David Stutz, and Bernt Schiele. Improving robustness by enhancing weak subnets. In *European Conference on Computer Vision*, pages 320–338. Springer, 2022. 5
- [18] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 3154–3160, 2017. 6
- [19] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16000–16009, June 2022. 2
- [20] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 6, 7
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [22] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2019. 3, 6
- [23] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. 5
- [24] Tomu Hirata, Yusuke Mukuta, and Tatsuya Harada. Making video recognition models robust to common corruptions with supervised contrastive learning. In *Proceedings of the 3rd ACM International Conference on Multimedia in Asia*, pages 1–6, 2021. 3
- [25] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 5
- [26] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, volume 1, page 2, 2019. 2
- [27] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 491–507. Springer, 2020. 5
- [28] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *ArXiv*, abs/1909.11942, 2019. 2

- [29] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 5
- [30] Qiufu Li, Linlin Shen, Sheng Guo, and Zhihui Lai. Wavelet integrated cnns for noise-robust image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 3
- [31] Qiufu Li, Linlin Shen, Sheng Guo, and Zhihui Lai. Wavelet integrated cnns for noise-robust image classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7245–7254, 2020. 3, 5
- [32] Xinyu Li, Yanyi Zhang, Jianbo Yuan, Hanlin Lu, and Yibo Zhu. Discrete cosin transformer: Image modeling from frequency domain. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5468–5478, 2023. 2
- [33] Pengju Liu, Hongzhi Zhang, Kai Zhang, Liang Lin, and Wangmeng Zuo. Multi-level wavelet-cnn for image restoration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018. 2
- [34] Xinyu Liu, Wuyang Li, Qiushi Yang, Baopu Li, and Yixuan Yuan. Towards robust adaptive object detection under noisy annotations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14207–14216, 2022. 1
- [35] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *ArXiv*, abs/1907.11692, 2019. 2
- [36] Stephane G Mallat. A theory for multiresolution signal decomposition: the wavelet representation. *IEEE transactions on pattern analysis and machine intelligence*, 11(7):674–693, 1989. 2
- [37] Mohammad Momeny, Ali Mohammad Latif, Mehdi Agha Sarram, Razieh Sheikhpour, and Yu Dong Zhang. A noise robust convolutional neural network for image classification. *Results in Engineering*, 10:100225, 2021. 7
- [38] Yi Pan, Jun-Jie Huang, Zihan Chen, Wentao Zhao, and Ziyue Wang. Svastin: Sparse video adversarial attack via spatio-temporal invertible neural networks. *arXiv preprint arXiv:2406.01894*, 2024. 1
- [39] Jeongsoo Park and Justin Johnson. Rgb no more: Minimally-decoded jpeg vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22334–22346, 2023. 2
- [40] Jeff Rasley, Samyam Rajbhandari, Olatunji Ranjan, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020. 6
- [41] Abolfazl Razi, Xiwen Chen, Huayu Li, Hao Wang, Brendan Russo, Yan Chen, and Hongbin Yu. Deep learning serves traffic safety analysis: A forward-looking review. *IET Intelligent Transport Systems*, 17(1):22–71, 2023. 1
- [42] Evgenia Rusak, Lukas Schott, Roland S Zimmermann, Julian Bitterwolf, Oliver Bringmann, Matthias Bethge, and Wieland Brendel. A simple way to make neural networks robust against diverse image corruptions. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 53–69. Springer, 2020. 5
- [43] Madeline Chantry Schiappa, Naman Biyani, Prudvi Kamtam, Shruti Vyas, Hamid Palangi, Vibhav Vineet, and Yogesh S Rawat. A large-scale robustness analysis of video action recognition models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14698–14708, 2023. 1, 5, 6
- [44] Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6:1–48, 2019. 3
- [45] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. 3
- [46] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5
- [47] K Soomro. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 5
- [48] Harold H Szu, Brian A Telfer, and Shubha L Kadambe. Neural network adaptive wavelets for signal representation and classification. *Optical Engineering*, 31(9):1907–1916, 1992. 2
- [49] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Advances in neural information processing systems*, 35:10078–10093, 2022. 1, 3, 7
- [50] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 10078–10093. Curran Associates, Inc., 2022. 2, 5, 6
- [51] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinnan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14549–14560, June 2023. 2
- [52] Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Yu-Gang Jiang, Luwei Zhou, and Lu Yuan. Bevt: Bert pretraining of video transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14733–14743, June 2022. 2
- [53] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense

- prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 568–578, 2021. [5](#)
- [54] Travis Williams and Robert Li. Wavelet pooling for convolutional neural networks. In *International conference on learning representations*, 2018. [3](#)
- [55] Zhiliang Wu, Changchang Sun, Hanyu Xuan, Gaowen Liu, and Yan Yan. Waveformer: Wavelet transformer for noise-robust video inpainting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 6180–6188, 2024. [1](#), [3](#), [4](#), [5](#)
- [56] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. [5](#)
- [57] Chenyu Yi, Siyuan Yang, Haoliang Li, Yap-peng Tan, and Alex Kot. Benchmarking the robustness of spatial-temporal models against corruptions. *arXiv preprint arXiv:2110.06513*, 2021. [2](#), [7](#)