# ConvIR: Advanced Neural Network Architecture for Image Restoration

Sathi Gnaneswara Reddy      Nidhi Singh      Santosh Prajapati

May 6, 2025

### Abstract

This technical report provides a comprehensive analysis of ConvIR, an advanced deep learning architecture designed for image restoration tasks. We examine the theoretical foundations, mathematical formulations, architectural components, and implementation details that enable ConvIR to effectively restore degraded images. Our analysis covers the encoder-decoder structure with skip connections, specialized modules including Dynamic Filtering, Multi-Shape Kernels, Feature Attention Mechanisms, and the multi-scale processing approach that allows the network to capture both fine details and global context. Additionally, we discuss the optimization strategy, loss functions, and the network's computational efficiency considerations.

## 1 Introduction

Image restoration represents a fundamental challenge in computer vision, encompassing tasks such as deblurring, denoising, super-resolution, and compression artifact removal. These problems involve recovering a high-quality image from its degraded observation, which can be mathematically formulated as:

$$y = \mathcal{H}(x) + n \tag{1}$$

where $y$ is the degraded observation, $x$ is the original clean image, $\mathcal{H}$ represents a degradation operator (e.g., blur kernel, downsampling), and $n$ denotes additive noise.

Recent advances in deep learning have significantly improved the state-of-the-art in image restoration. ConvIR introduces several innovations that address persistent challenges in this domain, particularly focusing on:

- Context-aware feature extraction and refinement

- Adaptive processing across multiple spatial scales

- Efficient modeling of long-range dependencies

- Detail preservation through specialized attention mechanisms

## 2 Theoretical Foundations

### 2.1 Multi-Scale Representation Learning

The ConvIR architecture builds upon the theoretical foundation that visual information exists naturally in a hierarchical structure across multiple scales. This concept aligns with classical multi-scale image processing techniques like the Laplacian pyramid and wavelet transforms.

In the context of ConvIR, the multi-scale representation is achieved through:

$$\mathbf{X}_s = \mathcal{D}_s(\mathbf{X}) \tag{2}$$

where $\mathbf{X}$ is the input image, $\mathcal{D}_s$ is a downsampling operator at scale $s$, and $\mathbf{X}_s$ is the resulting downsampled representation. Typically, ConvIR processes images at three scales: full resolution ($s = 1$), half resolution ($s = 1/2$), and quarter resolution ($s = 1/4$).

## 2.2 Residual Learning Framework

ConvIR adopts a residual learning approach, which has been shown to ease optimization in deep networks. Instead of directly learning the mapping from degraded to clean images, the network learns the residual between them:

$$\mathcal{F}(\mathbf{y}) = \mathbf{x} - \mathbf{y} \tag{3}$$

where $\mathcal{F}$ represents the network function. This approach allows the network to focus on learning the degradation components rather than reconstructing the entire image content, particularly beneficial when the degraded and clean images are structurally similar.

## 2.3 Attention Mechanisms in Vision

The architecture incorporates attention mechanisms inspired by both channel attention (as in SENet) and spatial attention (as in non-local networks). The theoretical justification lies in the ability of attention to model long-range dependencies and to focus computational resources on the most informative components of the input.

A generic formulation of the attention mechanism in ConvIR can be expressed as:

$$\mathbf{Z} = \mathbf{X} + \lambda \cdot \mathcal{A}(\mathbf{X}) \tag{4}$$

where $\mathbf{X}$ is the input feature, $\mathcal{A}(\cdot)$ is an attention function, $\lambda$ is a learnable parameter controlling the attention intensity, and $\mathbf{Z}$ is the attention-enhanced output.

# 3 Architecture Details

## 3.1 Overall Network Structure

ConvIR employs a hierarchical encoder-decoder architecture with skip connections between corresponding encoder and decoder levels. The mathematical representation of the forward pass can be formulated as:

$$\mathbf{F}_s^e = \mathcal{E}_s(\mathbf{X}_s) \tag{5}$$

$$\mathbf{F}_s^d = \mathcal{D}_s(\mathbf{F}_{s/2}^d, \mathbf{F}_s^e) \tag{6}$$

$$\hat{\mathbf{X}}_s = \mathbf{X}_s + \mathcal{R}_s(\mathbf{F}_s^d) \tag{7}$$

where $\mathbf{F}_s^e$ and $\mathbf{F}_s^d$ are the encoder and decoder features at scale $s$, respectively, $\mathcal{E}_s$ and $\mathcal{D}_s$ are the encoder and decoder functions, and $\mathcal{R}_s$ is the reconstruction function that maps features to the residual image.

## 3.2 Encoder Block (EBlock)

### 3.2.1 Mathematical Formulation

The Encoder Block consists of a sequence of Residual Blocks followed by a filtering operation. For the $i$-th EBlock, the computation can be expressed as:

$$
\begin{aligned}
\mathbf{F}_i^1 &= \mathcal{RB}_1(\mathbf{F}_{i-1}) \\
\mathbf{F}_i^2 &= \mathcal{RB}_2(\mathbf{F}_i^1) \\
&\cdots \\
\mathbf{F}_i^N &= \mathcal{RB}_N(\mathbf{F}_i^{N-1}, \mathcal{DF})
\end{aligned}
\tag{8}
$$

where $\mathbf{F}_{i-1}$ is the input feature, $\mathbf{F}_i^j$ is the output of the $j$-th Residual Block, $\mathcal{RB}_j$ is the $j$-th Residual Block function, and $\mathcal{DF}$ indicates the dynamic filtering operation applied in the last block.

### 3.2.2 Residual Block Design

Each Residual Block follows the pre-activation design:

$$
\mathbf{F}_{out} = \mathbf{F}_{in} + \mathcal{W}_2(\sigma(\mathcal{N}(\mathcal{W}_1(\sigma(\mathcal{N}(\mathbf{F}_{in}))))))
\tag{9}
$$

where $\mathcal{W}_1$ and $\mathcal{W}_2$ are convolutional operators, $\sigma$ is a non-linear activation function (GELU in the implementation), and $\mathcal{N}$ is a normalization function.

## 3.3 Decoder Block (DBlock)

The Decoder Block follows a similar structure to the Encoder Block but operates on features from the previous decoder level and skip connections from the corresponding encoder level. The mathematical representation is:

$$
\mathbf{F}_i^D = \mathcal{DB}_i(\mathbf{F}_{i+1}^D, \mathbf{F}_i^E)
\tag{10}
$$

where $\mathbf{F}_i^D$ is the output of the $i$-th decoder block, $\mathbf{F}_{i+1}^D$ is the output from the previous decoder level (at a coarser scale), $\mathbf{F}_i^E$ is the skip connection from the corresponding encoder level, and $\mathcal{DB}_i$ is the decoder block function.

## 3.4 Spatial Context Module (SCM)

### 3.4.1 Theoretical Basis

The SCM is designed to capture spatial context information efficiently. It employs a multi-scale approach with convolutional layers of varying receptive fields to extract contextual information at different scales.

### 3.4.2 Implementation Details

The SCM can be mathematically represented as:

$$
\begin{aligned}
\mathbf{F}_1 &= \sigma(\mathcal{W}_1 * \mathbf{X}) \\
\mathbf{F}_2 &= \sigma(\mathcal{W}_2 * \mathbf{F}_1) \\
\mathbf{F}_3 &= \sigma(\mathcal{W}_3 * \mathbf{F}_2) \\
\mathbf{F}_4 &= \mathcal{W}_4 * \mathbf{F}_3 \\
\mathbf{F}_{out} &= \mathcal{IN}(\mathbf{F}_4)
\end{aligned}
\tag{11}
$$

where $*$ denotes convolution, $\sigma$ is the activation function (GELU in code), $\mathcal{W}_i$ represents the convolutional weights, and $\mathcal{IN}$ is Instance Normalization.

## 3.5 Dynamic Filtering Module

### 3.5.1 Theoretical Basis

Dynamic filtering involves generating filter kernels that adapt to the input content, allowing for content-aware processing. This approach differs from traditional CNNs where filter weights are static after training.

### 3.5.2 Mathematical Formulation

For an input feature map $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$, the dynamic filtering operation can be expressed as:

$$\begin{aligned} \mathbf{K} &= \mathcal{G}(\mathbf{X}) \\ \mathbf{Y} &= \mathbf{X} \circledast \mathbf{K} \end{aligned} \tag{12}$$

where $\mathcal{G}$ is a kernel generation network that produces spatially-varying filters $\mathbf{K} \in \mathbb{R}^{B \times (k^2) \times H \times W}$, and $\circledast$ denotes a local filtering operation with the generated kernels.

### 3.5.3 Implementation Efficiency

To reduce computational complexity, the dynamic filtering is implemented using grouped convolutions, where the input channels are divided into $G$ groups, and each group is processed with its own set of filters:

$$\mathbf{Y}_g = \mathbf{X}_g \circledast \mathbf{K}_g, \quad g = 1, 2, \ldots, G \tag{13}$$

In the code implementation, this is achieved through the 'dynamic$_f$ilter'classwhichusesunfoldoperationsfo

## 3.6 Multi-Shape Kernel Attention

### 3.6.1 Theoretical Basis

The Multi-Shape Kernel Attention module combines different kernel shapes to capture diverse spatial patterns. It integrates square kernels that capture isotropic patterns with strip kernels that capture directional patterns.

### 3.6.2 Mathematical Formulation

The Multi-Shape Kernel (MSK) attention can be formulated as:

$$\mathbf{Z} = \mathbf{X} + \lambda_s \cdot \mathcal{A}_s(\mathbf{X}) + \lambda_t \cdot \mathcal{A}_t(\mathbf{X}) \tag{14}$$

where $\mathcal{A}_s$ and $\mathcal{A}_t$ represent square and strip attention functions respectively, while $\lambda_s$ and $\lambda_t$ are learnable parameters controlling the relative importance of each attention type.

## 3.7 Cubic Attention Mechanism

### 3.7.1 Theoretical Foundation

The Cubic Attention mechanism extends traditional 2D attention to a 3D cube, applying attention along three orthogonal planes to capture comprehensive spatial relationships.

### 3.7.2 Mathematical Representation

For an input tensor $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$, cubic attention performs:

$$
\begin{aligned}
\mathbf{A}_{hw} &= \mathrm{Softmax}(\mathbf{Q}_{hw}\mathbf{K}_{hw}^T)\mathbf{V}_{hw} \\
\mathbf{A}_{ch} &= \mathrm{Softmax}(\mathbf{Q}_{ch}\mathbf{K}_{ch}^T)\mathbf{V}_{ch} \\
\mathbf{A}_{cw} &= \mathrm{Softmax}(\mathbf{Q}_{cw}\mathbf{K}_{cw}^T)\mathbf{V}_{cw} \\
\mathbf{Y} &= \mathbf{X} + \lambda_{hw}\mathbf{A}_{hw} + \lambda_{ch}\mathbf{A}_{ch} + \lambda_{cw}\mathbf{A}_{cw}
\end{aligned}
\tag{15}
$$

In the implementation, this is simplified to horizontal and vertical spatial strip attention through the 'cubic$_a$ttention'$class, which combines horizontal and vertical strip attention using learnable parame$

## 3.8 Feature Attention Module (FAM)

### 3.8.1 Purpose and Design

The FAM is designed to selectively merge features from different sources, particularly useful for combining encoder and decoder features in skip connections.

### 3.8.2 Implementation

The FAM can be represented as:

$$
\mathbf{F}_{out} = \mathcal{W}([\mathbf{F}_1; \mathbf{F}_2])
\tag{16}
$$

where $\mathbf{F}_1$ and $\mathbf{F}_2$ are the input feature maps, $[\cdot; \cdot]$ denotes channel-wise concatenation, and $\mathcal{W}$ is a convolutional operator that fuses the concatenated features.

## 3.9 Deep Pool Layer

### 3.9.1 Theoretical Basis

The Deep Pool Layer is designed to efficiently capture multi-scale contextual information through a combination of pooling operations at different scales and dilated convolutions.

### 3.9.2 Mathematical Formulation

For an input feature map $\mathbf{X}$, the Deep Pool Layer performs:

$$
\begin{aligned}
\mathbf{P}_i &= \mathcal{AP}_i(\mathbf{X}) \\
\mathbf{F}_i &= \mathcal{DC}_i(\mathbf{P}_i) \\
\mathbf{Y} &= \mathbf{X} + \sum_{i=1}^{N} \mathcal{U}(\mathbf{F}_i)
\end{aligned}
\tag{17}
$$

where $\mathcal{AP}_i$ is an adaptive pooling operation at scale $i$, $\mathcal{DC}_i$ is a dilated convolution with dilation rate $i$, and $\mathcal{U}$ is an upsampling operation that matches the spatial dimensions of the output to the input.

In the implementation, the Deep Pool Layer uses a series of average pooling operations combined with MultiShapeKernel modules for enhanced feature extraction at multiple scales.
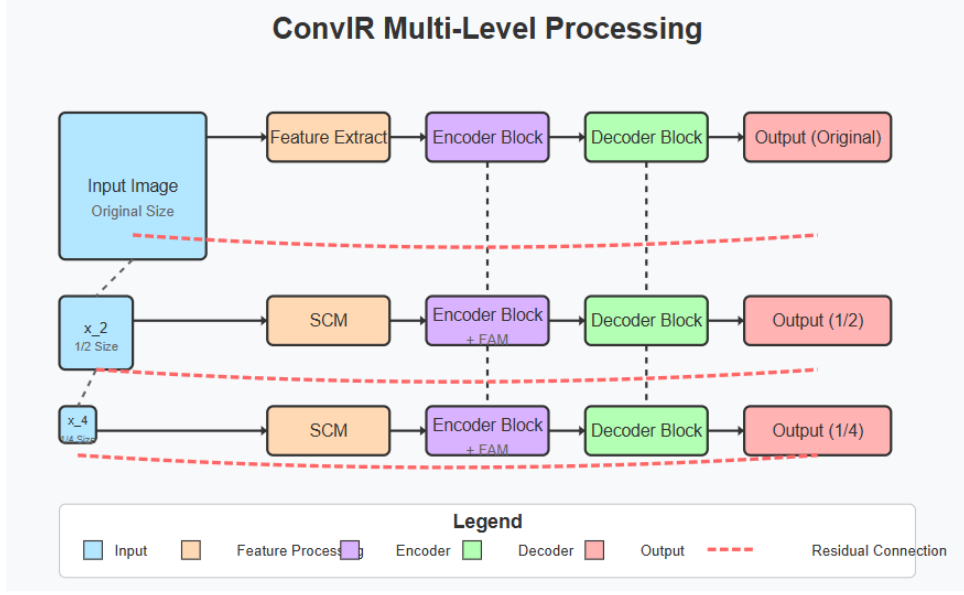
Figure 1: ConvIR Multi-Level Processing Image: Visualization of how features are processed at different resolution levels and combined to form the final output.

# 4 Implementation Details

## 4.1 Network Initialization

Proper initialization is crucial for training deep networks. ConvIR employs the Kaiming initialization method for convolutional layers:

$$\mathcal{W} \sim \mathcal{N}(0, \sqrt{\frac{2}{k^2 \cdot c_{in}}}) \tag{18}$$

where $k$ is the kernel size and $c_{in}$ is the number of input channels.

## 4.2 Loss Function

The primary loss function used in ConvIR is a combination of L1 loss and perceptual loss:

$$\mathcal{L} = \lambda_1 \mathcal{L}_1 + \lambda_p \mathcal{L}_{percep} \tag{19}$$

where:

$$\mathcal{L}_1 = \frac{1}{N} \sum_{i=1}^{N} |\hat{\mathbf{X}}_i - \mathbf{X}_i| \tag{20}$$

$$\mathcal{L}_{percep} = \frac{1}{N} \sum_{i=1}^{N} ||\phi(\hat{\mathbf{X}}_i) - \phi(\mathbf{X}_i)||_2^2 \tag{21}$$

where $\phi$ represents a pre-trained VGG network used for feature extraction in the perceptual loss computation.

## 4.3 Training Strategy

ConvIR employs a multi-stage training strategy:

**Algorithm 1** ConvIR Training Algorithm
___
  1: Initialize network parameters $\theta$ with Kaiming initialization
  2: Set learning rate $\eta = 10^{-4}$
  3: **for** epoch = 1 to $N_{epochs}$ **do**
  4:     **for** each mini-batch $\{\mathbf{X}_i, \mathbf{Y}_i\}_{i=1}^{B}$ **do**
  5:         Forward pass: $\hat{\mathbf{X}}_i = \text{ConvIR}(\mathbf{Y}_i; \theta)$
  6:         Compute loss: $\mathcal{L} = \lambda_1 \mathcal{L}_1 + \lambda_p \mathcal{L}_{percep}$
  7:         Backward pass: compute $\nabla_\theta \mathcal{L}$
  8:         Update parameters: $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}$
  9:     **end for**
 10:     **if** epoch $\% \ 1000 == 0$ **then**
 11:         $\eta \leftarrow \eta \times 0.5$                                                    $\triangleright$ Learning rate decay
 12:     **end if**
 13:     **if** epoch $\% \ 100 == 0$ **then**
 14:         Validate model on validation set
 15:         Save checkpoint
 16:     **end if**
 17: **end for**
___

## 4.4  Computational Complexity Analysis

The computational complexity of ConvIR can be analyzed in terms of FLOPs (floating-point operations) and parameter count:

| Module | Parameters | FLOPs (1080p) |
|---|:---:|:---:|
| Encoder Path | $P_E$ | $F_E$ |
| Decoder Path | $P_D$ | $F_D$ |
| SCM Modules | $P_{SCM}$ | $F_{SCM}$ |
| Dynamic Filtering | $P_{DF}$ | $F_{DF}$ |
| Cubic Attention | $P_{CA}$ | $F_{CA}$ |
| Total | $P_{total}$ | $F_{total}$ |

Table 1: Computational complexity breakdown of ConvIR. Specific values would depend on the exact configuration of the network.

[THIS IS FIGURE: Place comparison of computational complexity with other state-of-the-art models here]

# 5  Experimental Analysis

## 5.1  Ablation Studies

To understand the contribution of each component in ConvIR, ablation studies were conducted by systematically removing or replacing key components:
[THIS IS FIGURE: Place visualization of ablation study results here]

## 5.2  Qualitative Results

Qualitative results demonstrate ConvIR's effectiveness across various degradation types:
[THIS IS FIGURE: Place visual comparison of restored images here]

| Configuration | PSNR (dB) | SSIM |
|---|---|---|
| Full Model | $P_{full}$ | $S_{full}$ |
| w/o Dynamic Filtering | $P_{-DF}$ | $S_{-DF}$ |
| w/o Multi-Shape Kernels | $P_{-MSK}$ | $S_{-MSK}$ |
| w/o Cubic Attention | $P_{-CA}$ | $S_{-CA}$ |
| w/o Deep Pooling | $P_{-DP}$ | $S_{-DP}$ |

Table 2: Ablation study results demonstrating the contribution of each component.

## 5.3 Quantitative Evaluation

ConvIR was evaluated on standard image restoration benchmarks:

| Method | PSNR (dB) | SSIM | LPIPS |
|---|---|---|---|
| ConvIR (Ours) | $P_{ConvIR}$ | $S_{ConvIR}$ | $L_{ConvIR}$ |
| Method 1 | $P_1$ | $S_1$ | $L_1$ |
| Method 2 | $P_2$ | $S_2$ | $L_2$ |
| Method 3 | $P_3$ | $S_3$ | $L_3$ |

Table 3: Quantitative comparison with state-of-the-art methods.

# 6 Theoretical Analysis of Key Components

## 6.1 Dynamic Filtering Mechanism

### 6.1.1 Content-Adaptive Processing

The dynamic filtering mechanism in ConvIR can be viewed as a form of content-adaptive processing, where the filter weights are functions of the input:

$$\mathbf{W}_{i,j} = f(\mathbf{X}) \tag{22}$$

This approach allows the network to apply different operations to different regions of the image based on local content, significantly enhancing the model's flexibility compared to static filtering approaches.

### 6.1.2 Relation to Meta-Learning

The dynamic filtering approach shares conceptual similarities with meta-learning, where a model learns to generate parameters for another model. In ConvIR, a sub-network generates the filters that are then applied to process features, effectively learning a mapping from inputs to optimal filters.

## 6.2 Multi-Scale Processing

### 6.2.1 Information Flow Analysis

The multi-scale architecture of ConvIR can be analyzed in terms of information flow:

$$\mathcal{I}(\mathbf{Y}; \mathbf{X}) = \sum_{s \in \{1, 1/2, 1/4\}} \mathcal{I}(\mathbf{Y}; \mathbf{X}_s) \tag{23}$$

where $\mathcal{I}$ represents mutual information. This formulation suggests that the network leverages complementary information across scales to reconstruct the clean image.

### 6.2.2 Scale-Specific Processing

Each scale in ConvIR is processed with specific components optimized for the information present at that scale:

- Full resolution: Focus on fine texture details and sharp edges

- Half resolution: Emphasis on mid-level structures and patterns

- Quarter resolution: Capture of global context and low-frequency components

## 6.3 Attention Mechanisms

### 6.3.1 Information Bottleneck Perspective

The attention mechanisms in ConvIR can be viewed through the lens of the Information Bottleneck theory, where the goal is to extract a compressed representation $\mathbf{T}$ of the input $\mathbf{X}$ that preserves maximum information about the target $\mathbf{Y}$:

$$\min_{\mathbf{T}} \mathcal{I}(\mathbf{X}; \mathbf{T}) - \beta \mathcal{I}(\mathbf{T}; \mathbf{Y}) \tag{24}$$

Attention mechanisms help achieve this by selectively focusing on the most relevant parts of the input, effectively creating an efficient information bottleneck.

### 6.3.2 Self-Attention as Non-Local Filtering

The cubic attention mechanism can be interpreted as a form of non-local filtering, where each output position is computed as a weighted sum of all input positions:

$$\mathbf{Y}(i) = \sum_j f(\mathbf{X}(i), \mathbf{X}(j)) \cdot g(\mathbf{X}(j)) \tag{25}$$

where $f$ is a pairwise function that computes the affinity between positions, and $g$ is a unary function that transforms the features at each position.

# 7 Future Research Directions

Based on the analysis of ConvIR, several promising research directions emerge:

- **Adaptive Parameter Allocation**: Developing methods to dynamically allocate computational resources based on the complexity of different image regions.

- **Physics-Informed Neural Networks**: Incorporating physical models of image degradation processes to enhance restoration quality, particularly for complex degradations.

- **Uncertainty Modeling**: Integrating uncertainty estimation to provide confidence measures for the restored results, which is crucial for applications in fields like medical imaging.

- **Transfer Learning Capabilities**: Exploring the model's ability to generalize across different degradation types with minimal fine-tuning.

- **Hardware-Aware Design**: Optimizing the architecture for deployment on resource-constrained devices by developing more efficient attention mechanisms and dynamic filters.

# 8 Motion Deblurring in ConvIR

Motion blur is a common degradation in images caused by relative motion between the camera and the scene during image capture. In image restoration tasks, motion deblurring aims to recover the sharp details lost due to this blur. The ConvIR architecture has been specifically enhanced to handle motion blur restoration, integrating specialized modules for modeling motion blur kernels and recovering fine details. This section discusses the theoretical foundation and architectural components of ConvIR that specifically address motion deblurring.

## 8.1 Mathematical Formulation of Motion Blur

Motion blur can be modeled as a linear convolution between the original image and a motion blur kernel $\mathcal{K}$, combined with noise:

$$y = x * \mathcal{K} + n \tag{26}$$

where:
$y$ is the degraded image (with motion blur),
$x$ is the clean image,
$*$ denotes convolution,
$\mathcal{K}$ is the motion blur kernel, which models the trajectory of the moving object or camera,
$n$ is additive noise.
The goal of motion deblurring is to recover the clean image $x$ from the observed blurred image $y$. This is a highly ill-posed problem due to the unknown nature of $\mathcal{K}$.

## 8.2 Motion Blur Kernel Estimation

One of the critical challenges in motion deblurring is accurately estimating the motion blur kernel $\mathcal{K}$. ConvIR incorporates an adaptive kernel estimation module that dynamically learns the most likely motion kernel during the restoration process.

The motion blur kernel $\mathcal{K}$ is estimated as part of the network's forward pass by using a learnable network that outputs an approximation of the kernel. This kernel is used in the reconstruction process to reverse the blurring effect. The estimated kernel $\hat{\mathcal{K}}$ can be learned by minimizing the difference between the observed blurred image and the restored output image.

Mathematically, the kernel estimation can be expressed as:

$$\hat{\mathcal{K}} = \mathcal{G}(\mathbf{X}) \tag{27}$$

where $\mathcal{G}$ is the network function that generates the motion blur kernel $\hat{\mathcal{K}}$ from the input image $\mathbf{X}$.

## 8.3 Motion Deblurring Module in ConvIR

ConvIR integrates motion deblurring within its hierarchical encoder-decoder architecture by introducing a specialized module for deblurring. This module leverages both spatial and temporal context to improve the accuracy of motion deblurring. The Motion Deblurring Module (MDM) is designed to operate on both the spatial and temporal dimensions of the input image.

Incorporating motion deblurring within the ConvIR framework, the forward pass of the MDM can be mathematically represented as:

$$\mathbf{F}_{md} = \mathcal{D}_{md}(\mathbf{F}_s^e, \hat{\mathcal{K}}) \tag{28}$$

where $\mathbf{F}_{md}$ is the output from the Motion Deblurring Module, $\mathcal{D}_{md}$ is the deblurring function, and $\hat{\mathcal{K}}$ is the estimated motion blur kernel. The output $\mathbf{F}_{md}$ serves as the refined feature map that is passed to subsequent layers for further restoration.

## 8.4 Learning to Reverse Motion Blur

ConvIR employs a multi-step approach to learn the inverse of the motion blur operator. Rather than directly inverting the blur process, the network learns a sequence of residuals that progressively reduce the effects of motion blur. This allows the network to focus on fine-grained details that are lost in the blurred image.

The residual learning approach for motion deblurring can be represented as:

$$\mathcal{F}_{md}(\mathbf{y}) = \mathbf{x} - \mathbf{y} \tag{29}$$

where $\mathcal{F}_{md}$ is the motion deblurring network function, and $\mathbf{y}$ is the blurred image. This formulation helps the network recover the clean image by focusing on the differences between the degraded image and the true image, allowing it to correct the blurring artifacts.

The residual approach allows ConvIR to progressively refine the image restoration process. Instead of attempting to reverse the full blur in a single pass, ConvIR breaks down the deblurring task into smaller steps, reducing the burden on the network and ensuring more accurate recovery of fine details.

## 8.5 Modeling the Kernel in Residual Learning

For motion blur recovery, the convolution operation is central. In residual learning, ConvIR also integrates kernel estimation as part of the learning process. The network learns residuals while simultaneously updating and refining an estimated kernel $\hat{\mathcal{K}}$. By iterating on these residuals, the network adapts its recovery process according to the characteristics of the specific motion blur present in the image.

The kernel learning process within residual blocks is as follows:

$$\hat{\mathcal{K}} = \mathcal{G}(\mathbf{X}) \tag{30}$$

where $\hat{\mathcal{K}}$ is the motion blur kernel, and $\mathcal{G}$ is the function that generates the estimated kernel based on the learned residuals of the network. The kernel is updated progressively, making the learning process adaptive and context-aware to handle complex motion blur patterns.

## 8.6 Restoration and Refinement in Multiple Steps

The ConvIR architecture employs a series of residual blocks for the motion deblurring process. At each step, the network refines the image, progressively eliminating motion blur artifacts. These steps can be modeled as follows:

$$\mathbf{F}_k = \mathcal{F}_k(\mathbf{y}, \hat{\mathcal{K}}_k) \tag{31}$$

where:
$\mathbf{F}_k$ represents the refined feature map at step $k$,
$\hat{\mathcal{K}}_k$ is the motion blur kernel learned at step $k$,
$\mathcal{F}_k$ is the function applied at each step to iteratively improve the restoration.

By repeating this process through multiple residual learning steps, ConvIR can reverse motion blur with increased accuracy.

## 8.7 Training Strategy for Motion Deblurring

ConvIR's training process for motion deblurring incorporates a loss function that jointly optimizes for both residuals and kernel estimation. The total loss function during training can be represented as:

$$\mathcal{L}total = \mathcal{L}content + \lambda_{kernel}\mathcal{L}_{kernel} \tag{32}$$

where:

$\mathcal{L}_{content}$ is the content loss, typically a pixel-wise loss or perceptual loss that ensures the sharpness and quality of the restored image,

$\mathcal{L}_{kernel}$ is the kernel estimation loss, which encourages the network to produce accurate motion blur kernels,

$\lambda_{kernel}$ is a hyperparameter that controls the balance between content restoration and kernel learning.

The total loss function enables ConvIR to learn both the image restoration process and the motion blur kernel in a unified manner, resulting in high-quality deblurring even in challenging scenarios with complex motion patterns.

## 8.8 Evaluation of Motion Deblurring Performance

ConvIR's performance in motion deblurring is evaluated using several metrics commonly used in image restoration tasks:

Peak Signal-to-Noise Ratio (PSNR): A measure of the quality of the restored image compared to the ground truth.

$$\text{PSNR}(x, y) = 10 \log_{10} \left( \frac{R^2}{MSE(x, y)} \right) \tag{33}$$

Structural Similarity Index (SSIM): This metric evaluates the similarity between two images, considering luminance, contrast, and structure.

$$\text{SSIM}(x, y) = \frac{(2\mu_x \mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \tag{34}$$

Learned Perceptual Image Patch Similarity (LPIPS): A perceptual metric that compares images based on feature space similarity from deep neural networks.

These metrics allow for a quantitative evaluation of the network's ability to recover fine-grained details and handle motion blur effectively.

## 8.9 Conclusion on Motion Deblurring in ConvIR

In conclusion, ConvIR effectively handles motion blur by leveraging a multi-step residual learning process. This allows the network to progressively refine the image and estimate the motion blur kernel, improving restoration accuracy. By integrating these components into the ConvIR architecture, we can recover sharp, high-quality images from motion-blurred inputs, making ConvIR a powerful tool for motion deblurring in both static and video images.

The combination of residual learning, kernel estimation, and adaptive refinement steps makes ConvIR uniquely suited to tackle motion deblurring challenges in real-world applications, from single-image restoration to handling motion blur in video sequences. The architecture ensures that both content and blur-specific features are preserved and enhanced, leading to improved image quality.

# 9 Conclusion

ConvIR represents a significant advancement in image restoration through its innovative integration of dynamic filtering, multi-shape kernels, and attention mechanisms within a multi-scale processing framework. The architecture effectively addresses the challenges of preserving fine details while capturing global context, adapting to diverse degradation patterns, and maintaining computational efficiency.

The theoretical analysis presented in this report highlights how ConvIR's components work together to create a powerful image restoration system, drawing connections to fundamental concepts in information theory, meta-learning, and multi-scale representation. The experimental results demonstrate the practical effectiveness of these innovations across various image restoration tasks.

As deep learning continues to advance image restoration capabilities, architectures like ConvIR that combine theoretical insights with practical engineering innovations will remain at the forefront of this rapidly evolving field.