

## Appendix A. The Mathematical Derivation of Structural Reparameterization Process

To clearly demonstrate the equivalence between the training-time dual-path architecture and its inference-time reparameterized structure, we provide a step-by-step mathematical derivation of structural reparameterization process.

For notational clarity, we define the following terms:

- Input feature map:  $f \in \mathbb{R}^{C \times H \times W}$ .
- Convolution kernel:  $K_{c_o, c_i}^{(s)}$ , where  $s$  denotes the kernel size, and  $C_i, C_o$  denote input and output channels, respectively.
- Convolution operator:  $*$  denotes valid convolution; zero-padding, if needed, is explicitly written in the kernel.
- Element-wise multiplication:  $\odot$ .
- Sigmoid activation:  $\sigma(\cdot)$ .

### A.1. Dual-Path Formula in Training Phase

The Transformer branch employs a window-based multi-head self-attention mechanism, which can be formally expressed as:

$$f_{\text{att}} = \text{W-MSA}(f) \in \mathbb{R}^{C \times H \times W} \quad (5)$$

In parallel, the convolutional branch applies a standard  $3 \times 3$  convolution operation to the input feature map:

$$g = \text{Conv}_{3 \times 3}(f) \in \mathbb{R}^{C \times H \times W} \quad (6)$$

The outputs of the two branches are fused through element-wise multiplication, followed by a  $1 \times 1$  convolution:

$$f_{\text{out}} = \text{Conv}_{1 \times 1}^{(1)} \left( \text{Conv}_A(f_{\text{att}}) \odot \sigma(\text{Conv}_B(g)) \right) \quad (7)$$

Both  $\text{Conv}_A$  and  $\text{Conv}_B$  consist solely of  $1 \times 1$  convolutions.

### A.2. Converting W-MSA's Linear Projection into Convolution

The linear projections used for  $(Q, K, V)$  in W-MSA can be viewed as  $1 \times 1$  convolutions. This equivalence allows the W-MSA operation to be expressed as:

$$\text{W-MSA}(f) = K_{\text{msa}} * f, \quad K_{\text{msa}} \in \mathbb{R}^{C \times C \times 1 \times 1}. \quad (8)$$

Therefore, the output can be written as:

$$f_{\text{att}} = K_{\text{msa}} * f \quad (9)$$

### A.3. Merging Two $1 \times 1$ Convolutions into Single Convolution

For convenience, we define:

$$K_A = \text{kernel}(\text{Conv}_A), \quad K_B = \text{kernel}(\text{Conv}_B), \quad K_{3 \times 3} = \text{kernel}(\text{Conv}_{3 \times 3}) \quad (10)$$

These kernels have the following dimensions:

$$K_A, K_B \in \mathbb{R}^{C \times C \times 1 \times 1}, \quad K_{3 \times 3} \in \mathbb{R}^{C \times C \times 3 \times 3} \quad (11)$$

Substituting Eq. 9 into Eq. 7, we obtain:

$$\begin{aligned} f_{\text{out}} &= \text{Conv}_{1 \times 1}^{(1)} \left( (K_A * K_{\text{msa}} * f) \odot \sigma(K_B * (K_{3 \times 3} * f)) \right) \\ &= \text{Conv}_{1 \times 1}^{(1)} \left( [(K_A * K_{\text{msa}}) * f] \odot \sigma[(K_B * K_{3 \times 3}) * f] \right) \end{aligned}$$

### A.4. Transforming Element-wise Multiplication into $1 \times 1$ Convolution

Element-wise multiplication  $\odot$  is equivalent to a depthwise convolution with  $1 \times 1$  kernel and group size of 1. The corresponding weight tensor is constructed as follows:

$$K_{\text{mult}}^{(c,c)} = [(K_A * K_{\text{msa}})]^{(c)} \cdot \sigma[(K_B * K_{3 \times 3})]^{(c)} \quad (12)$$

Consequently, the overall transformation in Eq. 12 can be expressed as a single convolution operation:

$$f_{\text{out}} = K_{\text{mult}} * f \quad (13)$$

where the fused kernel  $K_{\text{mult}}$  is obtained via:

$$K_{\text{mult}} = \text{Conv}_{1 \times 1}^{(1)} \left( (K_A * K_{\text{msa}}) \cdot \sigma(K_B * K_{3 \times 3}) \right) \in \mathbb{R}^{C \times C \times 1 \times 1} \quad (14)$$

### A.5. Decomposing the $1 \times 1$ Convolution into Cascaded Convolutions

Although the kernel  $K_{\text{mult}}$  in Eq. 13 consists solely of  $1 \times 1$  convolutions, hardware accelerators during inference often favor  $3 \times 3$  convolutions to better leverage algorithmic optimizations such as Winograd or Fast Fourier Transform (FFT). By exploiting the additivity and homogeneity properties of convolution operations, the  $1 \times 1$  convolution kernel can be equivalently decomposed into a sequence of nested convolutions:

$$K_{\text{mult}} = K_1^{(1)} * K_3^{(2)} * K_3^{(3)} * K_1^{(4)} \quad (15)$$

where

- $K_1^{(1)}$  is  $1 \times 1$  convolution that reduces input channels for dimensionality compression,
- $K_3^{(2)}$  and  $K_3^{(3)}$  are  $3 \times 3$  standard convolutions.
- $K_1^{(4)}$  is  $1 \times 1$  convolution that expands channels back to dimension  $C$ .

These kernels are obtained via numerical decomposition (e.g., Singular Value Decomposition) after training and are loaded directly during inference without further decomposition. The inference process can thus be expressed as:

$$f_{\text{out}} = \text{Conv}_{1 \times 1}^{(4)} \left( \text{Conv}_{3 \times 3}^{(3)} \left( \text{Conv}_{3 \times 3}^{(2)} \left( \text{Conv}_{1 \times 1}^{(1)}(f) \right) \right) \right) \quad (16)$$

The transformations from Eq. 12 to Eq. 13 and finally to Eq. 16 are strictly equivalent at each step. This demonstrates the functional equivalence between the original dual-branch structure and the reparameterized convolutional sequence. Furthermore, the use of nested convolutions reduces computational complexity from  $O(N^2)$  to  $O(N)$ .

## Appendix B. Additional Ablation Study Results and Model Complexity Breakdown

Moreover, we provide BD-MS-SSIM and BD-Rate comparison on both datasets in Table 3. Superior compression performance is indicated by positive BD-MS-SSIM and negative BD-Rate, reflecting higher quality preservation and improved coding efficiency, respectively. To evaluate the lightweight nature of the proposed framework and demonstrate the necessity of structural reparameterization, a component-level complexity breakdown is presented in Table 4.

Table 3: BD-Rate and BD-MS-SSIM comparison of architectural variants on InStereo2K and Cityscapes dataset.

Architectural Variants	InStereo2K		Cityscapes	
	BD-Rate	BD-MS-SSIM	BD-Rate	BD-MS-SSIM
Backbone	0	0	0	0
Backbone + RSB (Reparam)	-4.60	0.13 (+0.26%)	-9.40	0.08 (+0.32%)
Backbone + RSB (Dual-Path)	-9.41	0.27 (+0.85%)	-14.35	0.64 (+1.16%)
Backbone + CFEMs	-12.69	0.33 (+1.73%)	-26.96	0.87 (+1.45%)
RSTSIC	<b>-14.55</b>	<b>0.56 (+2.58%)</b>	<b>-34.48</b>	<b>1.40 (+1.60%)</b>

Table 4: Component-level complexity breakdown of the proposed framework.

Component	Parameters (M)	FLOPs (G)
RSB	0.20	25.15
CFEM	0.49	2.11
Joint Decoder	1.60	33.81
Independent Encoder	0.42	1.11
RSTSIC	2.73	40.79