

### Re-param for Plain Inference-time Model

In this subsection, we describe how to convert a trained block into a single  $3 \times 3$  conv layer for inference. Note that we use BN in each branch before the addition (Fig. 4). Formally, we use  $W^{(3)} \in \mathbb{R}^{C_2 \times C_1 \times 3 \times 3}$  to denote the kernel of a  $3 \times 3$  conv layer with  $C_1$  input channels and  $C_2$  output channels, and  $W^{(1)} \in \mathbb{R}^{C_2 \times C_1}$  for the kernel of  $1 \times 1$  branch. We use  $\mu^{(3)}, \sigma^{(3)}, \gamma^{(3)}, \beta^{(3)}$  as the accumulated mean, standard deviation and learned scaling factor and bias of the BN layer following  $3 \times 3$  conv,  $\mu^{(1)}, \sigma^{(1)}, \gamma^{(1)}, \beta^{(1)}$  for the BN following  $1 \times 1$  conv, and  $\mu^{(0)}, \sigma^{(0)}, \gamma^{(0)}, \beta^{(0)}$  for the identity branch. Let  $M^{(1)} \in \mathbb{R}^{N \times C_1 \times H_1 \times W_1}$ ,  $M^{(2)} \in \mathbb{R}^{N \times C_2 \times H_2 \times W_2}$  be the input and output, respectively, and  $*$  be the convolution operator. If  $C_1 = C_2, H_1 = H_2, W_1 = W_2$ , we have

$$\begin{aligned} M^{(2)} = & \text{bn}(M^{(1)} * W^{(3)}, \mu^{(3)}, \sigma^{(3)}, \gamma^{(3)}, \beta^{(3)}) \\ & + \text{bn}(M^{(1)} * W^{(1)}, \mu^{(1)}, \sigma^{(1)}, \gamma^{(1)}, \beta^{(1)}) \\ & + \text{bn}(M^{(1)}, \mu^{(0)}, \sigma^{(0)}, \gamma^{(0)}, \beta^{(0)}). \end{aligned} \quad (1)$$

Otherwise, we simply use no identity branch, hence the above equation only has the first two terms. Here bn is the inference-time BN function, formally,  $\forall 1 \leq i \leq C_2$ ,

$$\text{bn}(M, \mu, \sigma, \gamma, \beta)_{:,i,:,:) = (M_{:,i,:,:) - \mu_i) \frac{\gamma_i}{\sigma_i} + \beta_i. \quad (2)$$

We first convert every BN and its preceding conv layer into a conv with a bias vector. Let  $\{W', b'\}$  be the kernel and bias converted from  $\{W, \mu, \sigma, \gamma, \beta\}$ , we have

$$W'_{i,:,:) = \frac{\gamma_i}{\sigma_i} W_{i,:,:), \quad b'_i = -\frac{\mu_i \gamma_i}{\sigma_i} + \beta_i. \quad (3)$$

Then it is easy to verify that  $\forall 1 \leq i \leq C_2$ ,

$$\text{bn}(M * W, \mu, \sigma, \gamma, \beta)_{:,i,:,:) = (M * W')_{:,i,:,:) + b'_i. \quad (4)$$

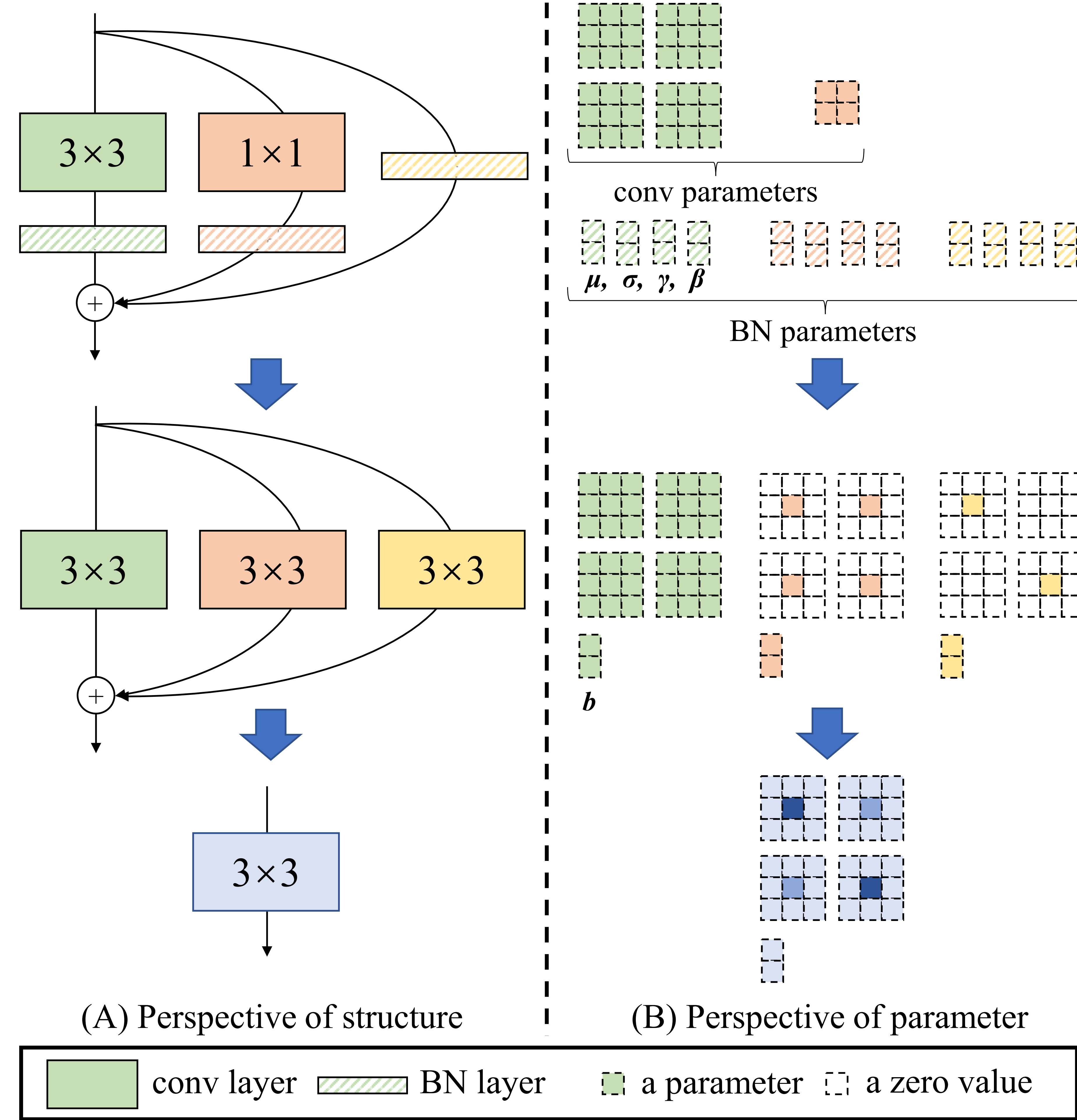


Figure 4: Structural re-parameterization of a RepVGG block. For the ease of visualization, we assume  $C_2 = C_1 = 2$ , thus the  $3 \times 3$  layer has four  $3 \times 3$  matrices and the kernel of  $1 \times 1$  layer is a  $2 \times 2$  matrix.

The above transformation also applies to the identity branch because an identity mapping can be viewed as a  $1 \times 1$  conv with an identity matrix as the kernel.

After such transformations, we will have one  $3 \times 3$  kernel, two  $1 \times 1$  kernels, and three bias vectors.

Then we obtain the final bias by adding up the three bias vectors, and the final  $3 \times 3$  kernel by adding the  $1 \times 1$  kernels onto the central point of  $3 \times 3$  kernel, which can be easily implemented by first zero-padding the two  $1 \times 1$  kernels to  $3 \times 3$  and adding the three kernels up, as shown in Fig. 4.

Note that the equivalence of such transformations requires the  $3 \times 3$  and  $1 \times 1$  layer to have the same stride, and the padding configuration of the latter shall be one pixel less than the former. *E.g.*, for a  $3 \times 3$  layer that pads the input by one pixel, the most common case, the  $1 \times 1$  layer should have padding = 0.