

Method 2

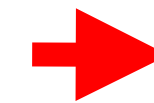
Compressing the Convolution

$$\mathcal{Q}(f, h, w, r_4) = \sum_{s=1}^S \mathbf{U}^{(4)}(s, r_4) \mathcal{X}(f, h, w, s)$$



$1 \times 1 \times 1$ convolution with S input channels and R_4 output channels

$$\mathcal{Q}'(f', h', w', r_5) = \sum_{i=1}^{D_F} \sum_{j=1}^{D_H} \sum_{l=1}^{D_W} \sum_{r_4=1}^{R_4} \mathcal{C}(i, j, l, r_4, r_5) \mathcal{Q}(f_i, h_j, w_l, r_4)$$

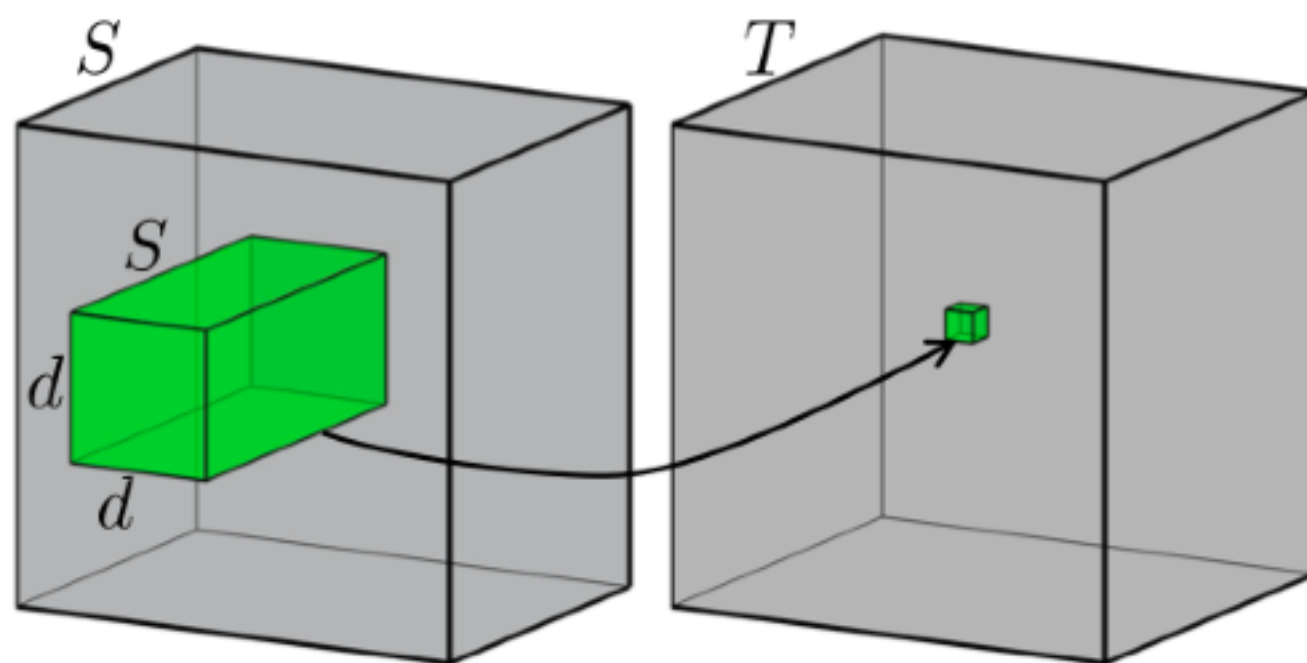


$D_F \times D_H \times D_W$ convolution with R_4 input channels and R_5 output channels

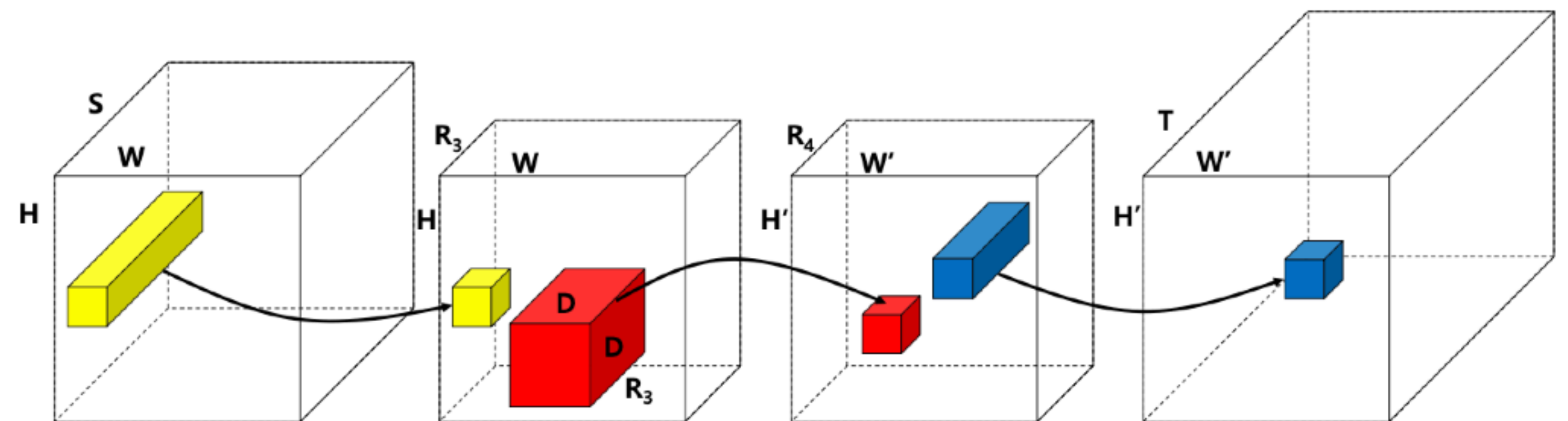
$$\mathcal{Y}(f', h', w', t) = \sum_{r_5=1}^{R_5} \mathbf{U}^{(5)}(t, r_5) \mathcal{Q}'(f', h', w', r_5)$$



$1 \times 1 \times 1$ convolution with R_5 input channels and T output channels



Original convolution



Sequence of convolutions of the compressed version

Method 2

Compressing the Linear Layer

The original layer

$$\mathbf{y}^{N_{out}} = \mathbf{W}^{N_{out} \times N_{in}} \mathbf{x}^{N_{in}} + \mathbf{b}^{N_{out}}$$

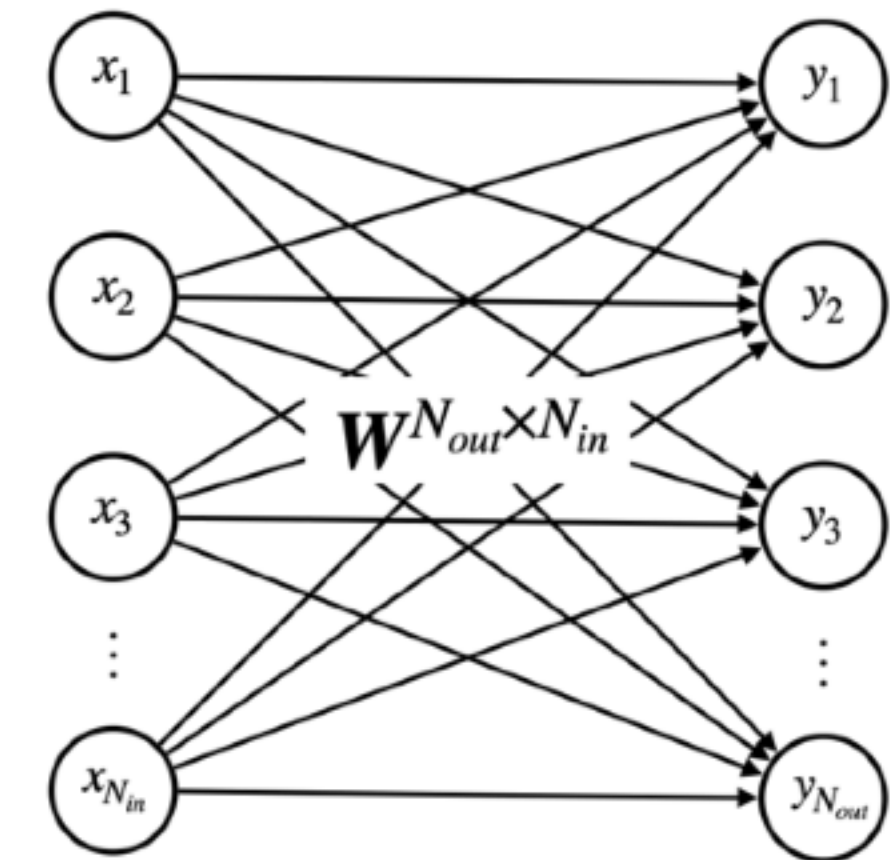
Tucker-2 approximation
of weight matrix

$$\mathbf{W} \approx \mathbf{G} \times_1 \mathbf{A} \times_2 \mathbf{B} = \mathbf{A} \mathbf{G} \mathbf{B}^\top$$

The compressed layer

$$\mathbf{y}^{N_{out}} \approx \underbrace{\mathbf{A}^{N_{out} \times R_A}}_{\mathbf{q}_2} \underbrace{\mathbf{G}^{R_A \times N_{in}} \underbrace{\mathbf{B}^\top R_B \times N_{in}}_{\mathbf{q}_1} \mathbf{x}^{N_{in}}}_{\mathbf{q}_1} + \mathbf{b}^{N_{out}}$$

The original layer



The compressed layer

