

# Symbolic Complexity of the Colorization Model

(Akshith / 2022102048)

## Assumptions and Notation

- Input is grayscale of spatial size  $H \times W$ , with channels NIC (typically NIC = 1).
- Convolution / transpose-convolution kernels are  $k \times k$  with  $k^2 = k \cdot k$ , padding  $p = k/2$ .
- Strides:
  - Encoder: three stride-2 convs:  $32 \rightarrow 16 \rightarrow 8 \rightarrow 4$  (or  $64 \rightarrow 32 \rightarrow 16 \rightarrow 8$  for doubled input).
  - Decoder: three stride-2 transpose-convs:  $4 \rightarrow 8 \rightarrow 16 \rightarrow 32$  (or  $8 \rightarrow 16 \rightarrow 32 \rightarrow 64$ ).
- Channels per layer:
$$\begin{aligned} L1: \text{NIC} &\rightarrow \text{NF}, & L2: \text{NF} &\rightarrow 2\text{NF}, & L3: 2\text{NF} &\rightarrow 4\text{NF}, \\ D1: 4\text{NF} &\rightarrow 2\text{NF}, & D2: 2\text{NF} &\rightarrow \text{NF}, & D3: \text{NF} &\rightarrow \text{NF}, & \text{CLS } (1 \times 1): \text{NF} &\rightarrow \text{NC}. \end{aligned}$$
- Ignore biases and BatchNorm parameters.

**Per-layer formulae.** For any (transpose-)convolution with  $C_{\text{in}} \rightarrow C_{\text{out}}$ , kernel  $k \times k$ , output map  $H_\ell \times W_\ell$ :

$$\begin{aligned} \text{Weights: } &C_{\text{out}} C_{\text{in}} k^2, \\ \text{Outputs (activation elements): } &C_{\text{out}} H_\ell W_\ell, \\ \text{Connections: } &(C_{\text{out}} C_{\text{in}} k^2) \cdot (H_\ell W_\ell). \end{aligned}$$

For the  $1 \times 1$  classifier: weights = NC · NF; outputs = NC ·  $H_{\text{out}} W_{\text{out}}$ ; connections =  $(\text{NC} \cdot \text{NF}) \cdot (H_{\text{out}} W_{\text{out}})$ .

## Case A: Input $32 \times 32$

Output spatial sizes per layer:

$$\begin{aligned} L1: 16 \times 16 & (= 256), & L2: 8 \times 8 & (= 64), & L3: 4 \times 4 & (= 16), \\ D1: 8 \times 8 & (= 64), & D2: 16 \times 16 & (= 256), & D3: 32 \times 32 & (= 1024), \\ & & & & \text{CLS: } 32 \times 32 & (= 1024). \end{aligned}$$

### 1) Total Number of Weights

Layer-wise weights:

$$\begin{aligned} L1: \text{NF} \cdot \text{NIC} \cdot k^2, & L2: 2\text{NF}^2 k^2, & L3: 8\text{NF}^2 k^2, \\ D1: 8\text{NF}^2 k^2, & D2: 2\text{NF}^2 k^2, & D3: \text{NF}^2 k^2, & \text{CLS: } \text{NF} \cdot \text{NC}. \end{aligned}$$

Sum:

$$W_{\text{total}, 32} = k^2 (\text{NF} \cdot \text{NIC} + 21 \text{NF}^2) + \text{NF} \cdot \text{NC}.$$

## 2) Total Number of Outputs (Activation Elements)

Layer-wise outputs:

$$\begin{aligned} L1: NF \cdot 256, \quad L2: 2NF \cdot 64 = NF \cdot 128, \quad L3: 4NF \cdot 16 = NF \cdot 64, \\ D1: 2NF \cdot 64 = NF \cdot 128, \quad D2: NF \cdot 256, \quad D3: NF \cdot 1024, \\ \text{CLS: } NC \cdot 1024. \end{aligned}$$

Sum:

$$O_{\text{total}, 32} = NF \cdot 1856 + NC \cdot 1024.$$

## 3) Total Number of Connections

Layer-wise connections (weights  $\times$  output positions):

$$\begin{aligned} L1: (NF \cdot NIC \cdot k^2) \cdot 256, \\ L2: (2NF^2 k^2) \cdot 64 = 128 NF^2 k^2, \\ L3: (8NF^2 k^2) \cdot 16 = 128 NF^2 k^2, \\ D1: (8NF^2 k^2) \cdot 64 = 512 NF^2 k^2, \\ D2: (2NF^2 k^2) \cdot 256 = 512 NF^2 k^2, \\ D3: (NF^2 k^2) \cdot 1024 = 1024 NF^2 k^2, \\ \text{CLS: } (NF \cdot NC) \cdot 1024. \end{aligned}$$

Sum:

$$C_{\text{total}, 32} = k^2 (256 NF \cdot NIC + 2304 NF^2) + 1024 NF \cdot NC.$$

## Case B: Input $64 \times 64$ (spatial doubled)

With the same stride pattern (three  $\times 2$  downs/ups), every feature-map area  $H_\ell W_\ell$  is  $4 \times$  larger than in the  $32 \times 32$  case; *weights do not depend on spatial size*.

Therefore:

$$W_{\text{total}, 64} = W_{\text{total}, 32} = k^2 (NF \cdot NIC + 21 NF^2) + NF \cdot NC.$$

$$O_{\text{total}, 64} = 4 O_{\text{total}, 32} = NF \cdot 7424 + NC \cdot 4096.$$

$$C_{\text{total}, 64} = 4 C_{\text{total}, 32} = k^2 (1024 NF \cdot NIC + 9216 NF^2) + 4096 NF \cdot NC.$$

## Notes

- If the actual channel schedule or number of stages differs, re-apply the per-layer formulas above with your  $C_{\text{in}}$ ,  $C_{\text{out}}$ ,  $H_\ell$ ,  $W_\ell$  and sum.
- Transpose-conv layers use the same weight formula  $C_{\text{out}} C_{\text{in}} k^2$ ; their output sizes follow the stride-2 upsampling with padding  $p = k//2$ .