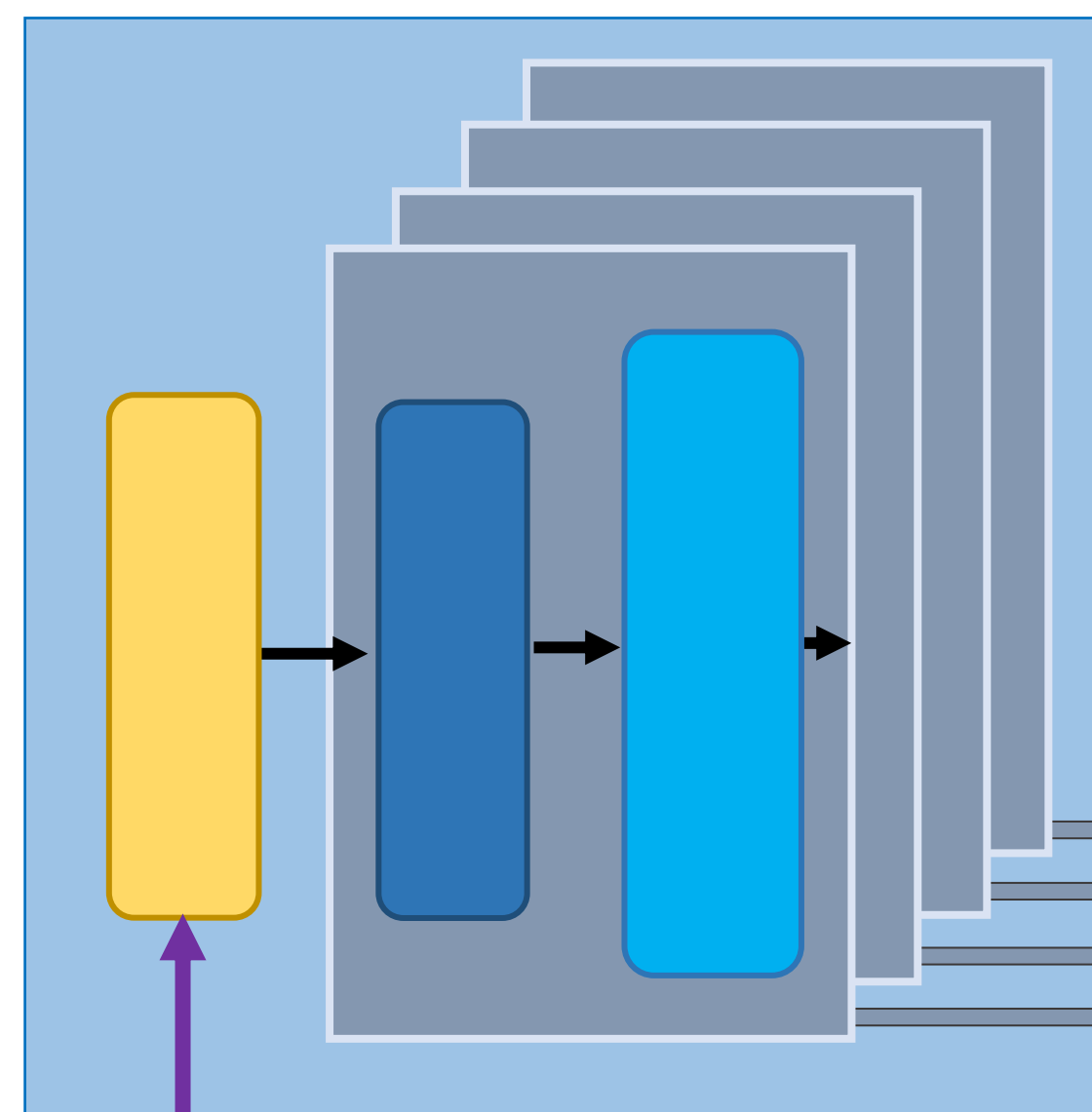
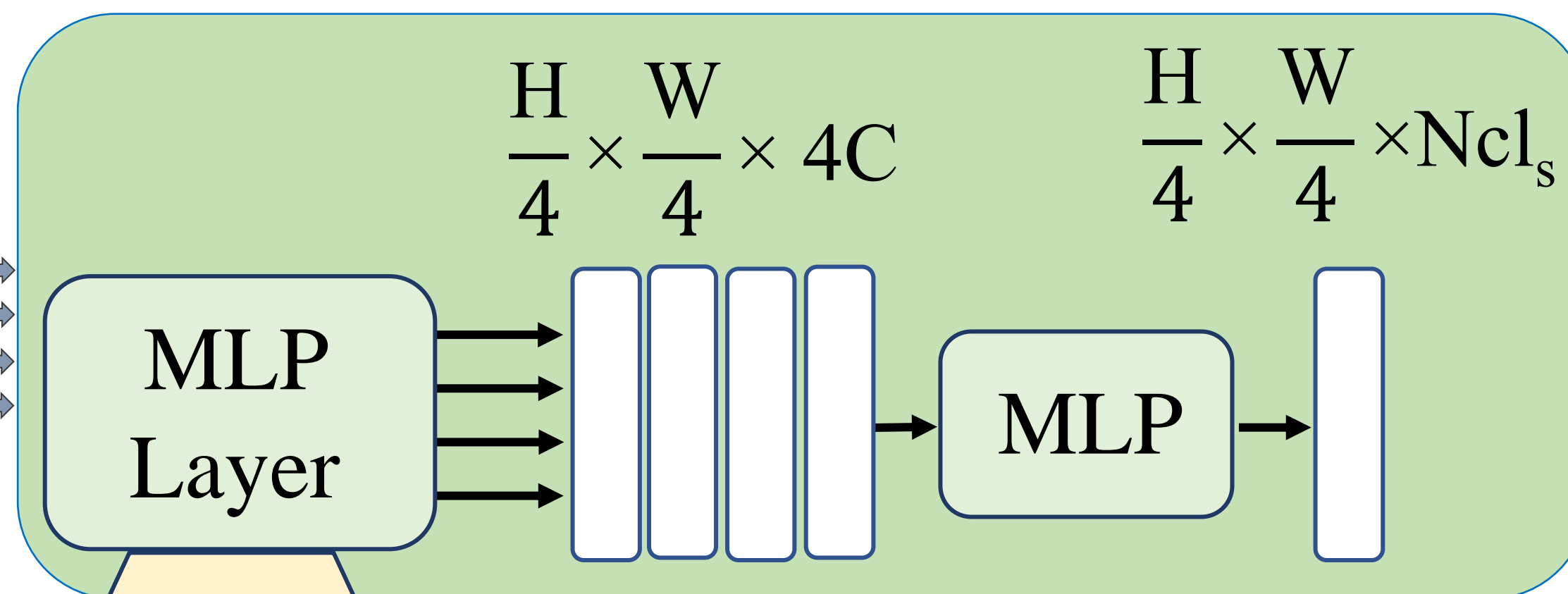


# Encoder



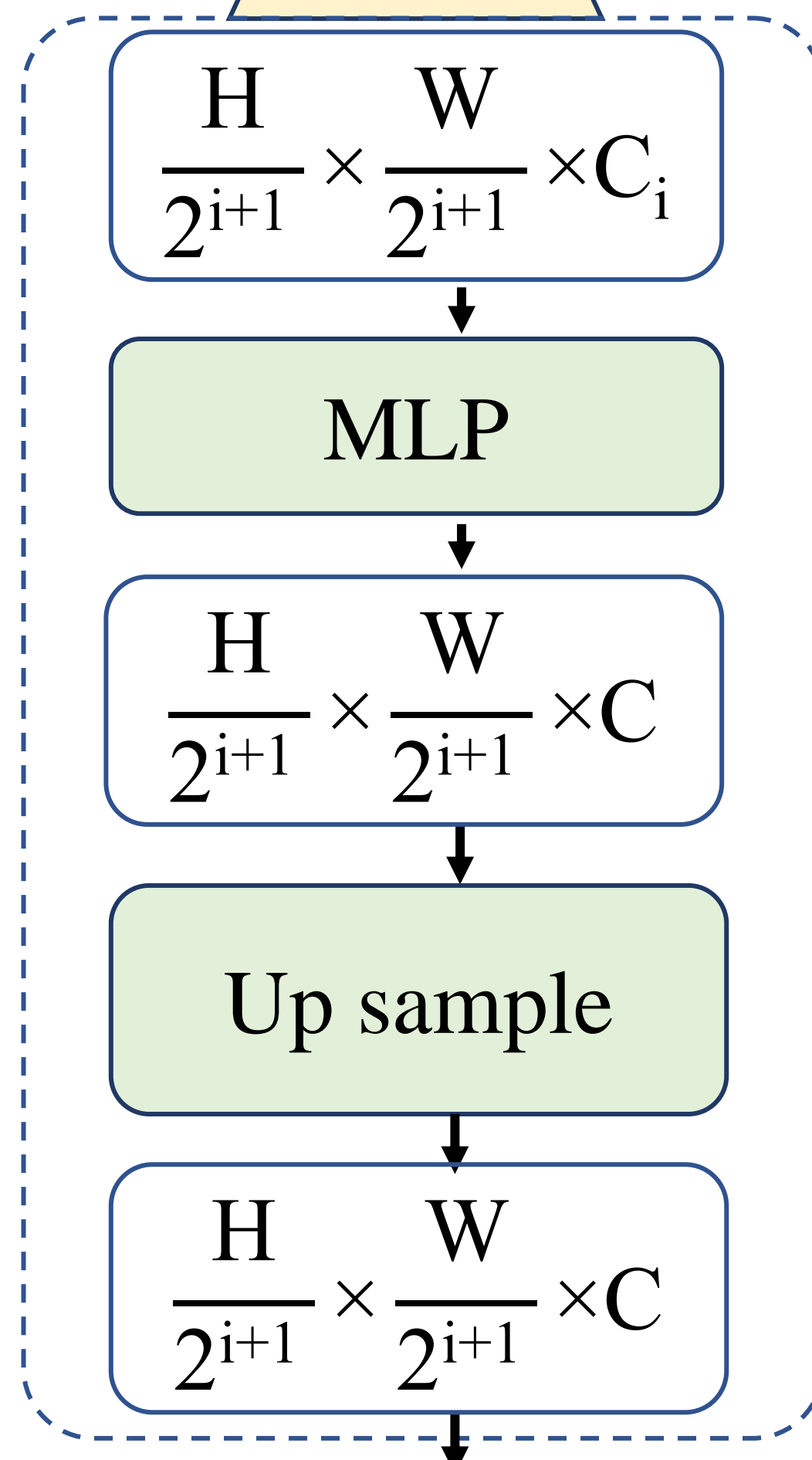
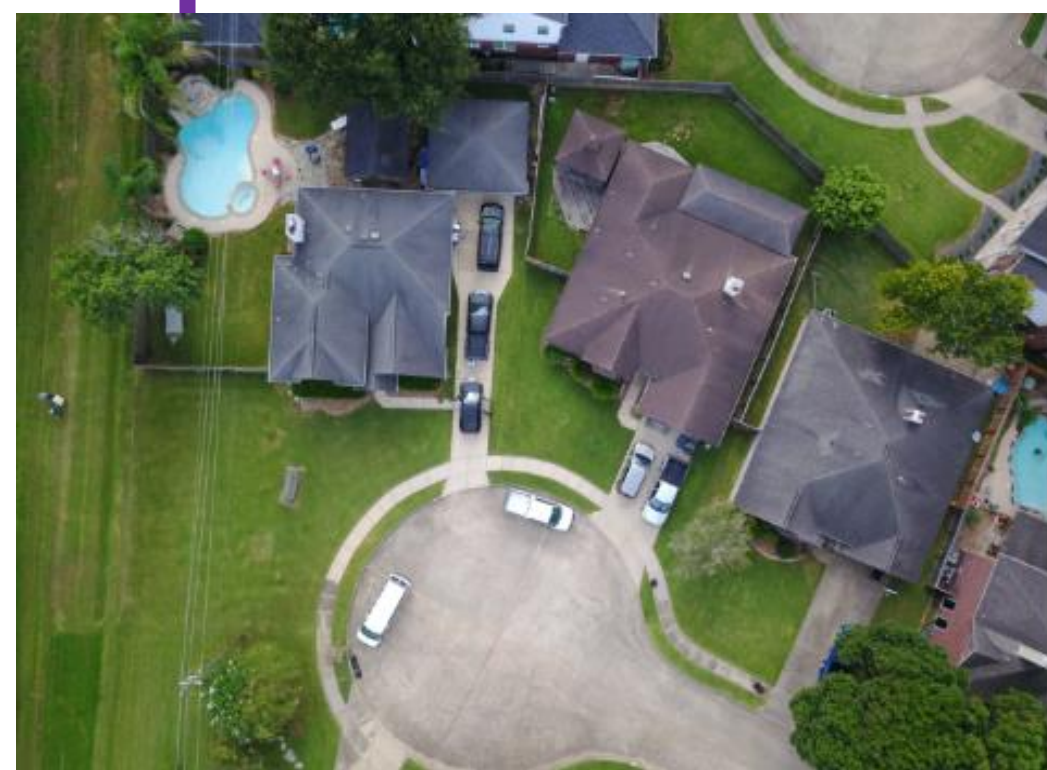
# Decoder



# Predicted mask



# Input image



Patch partitioning



Linear embedding



SwinT block



Stages:  $S_i$  where  $i \in \{1, 2, 3, 4\}$ ,

Feature size at  $S_i$ :  $\frac{H}{2^{i+1}} \times \frac{W}{2^{i+1}} \times C_i$