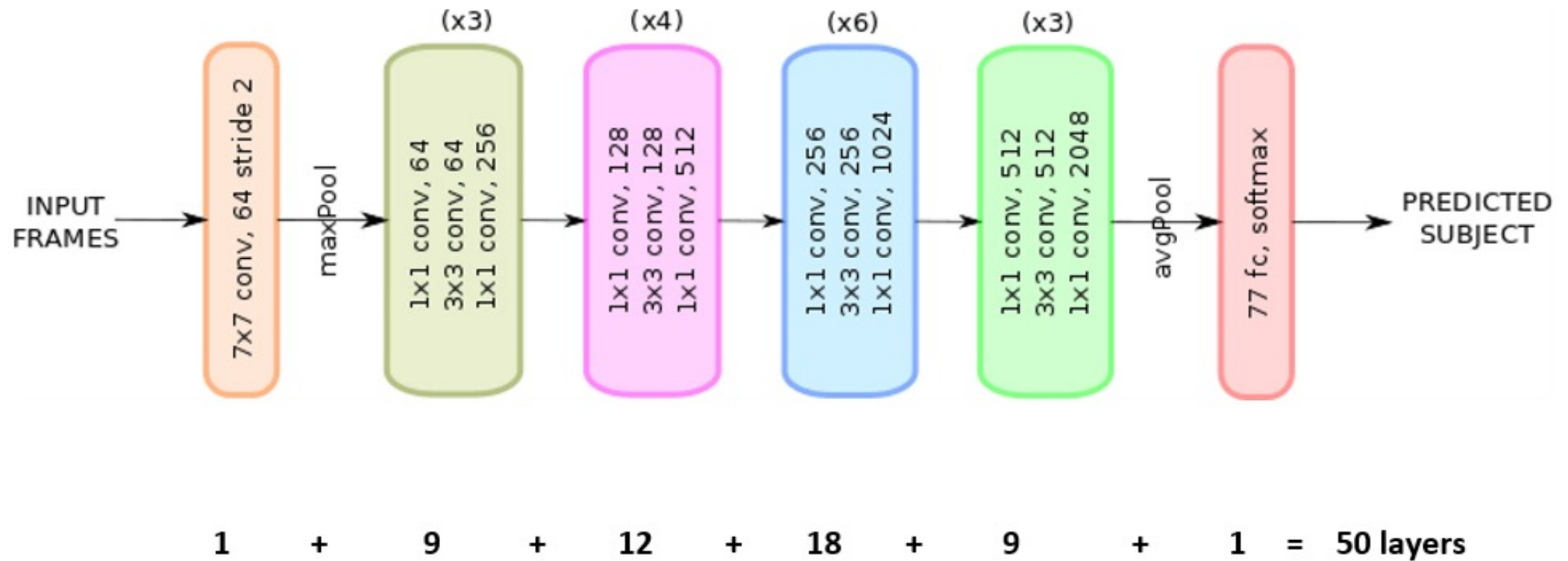


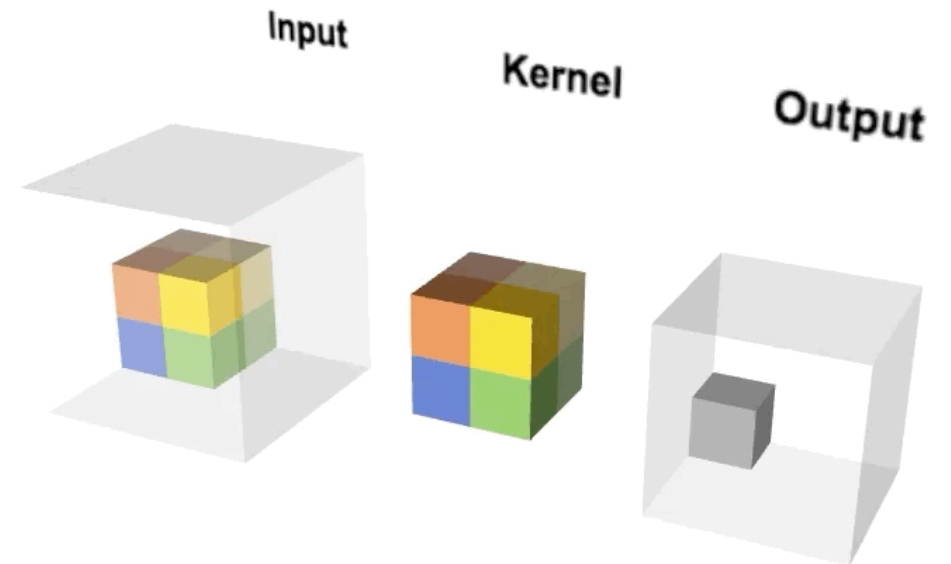
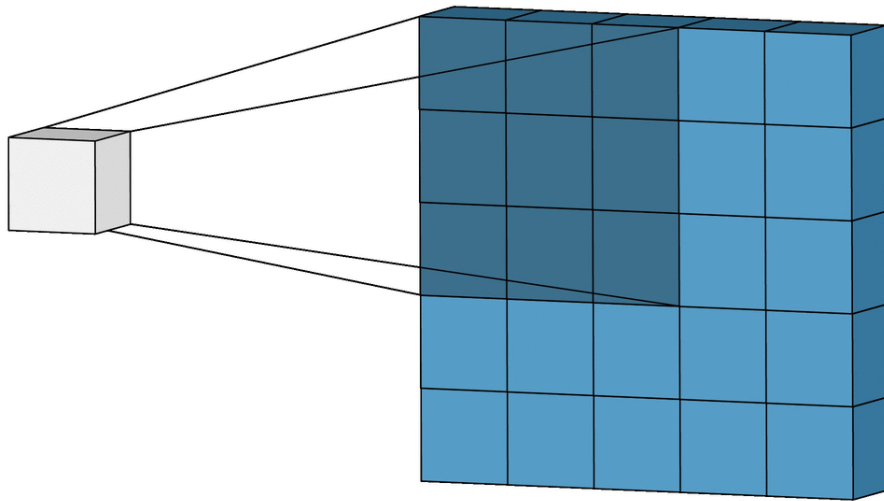
Part 2

Video Models

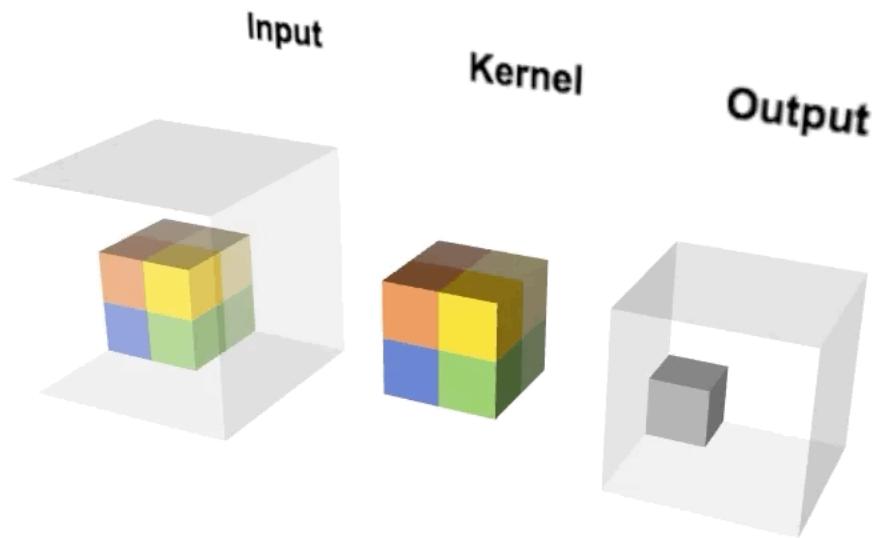
CNNs (ResNet50)



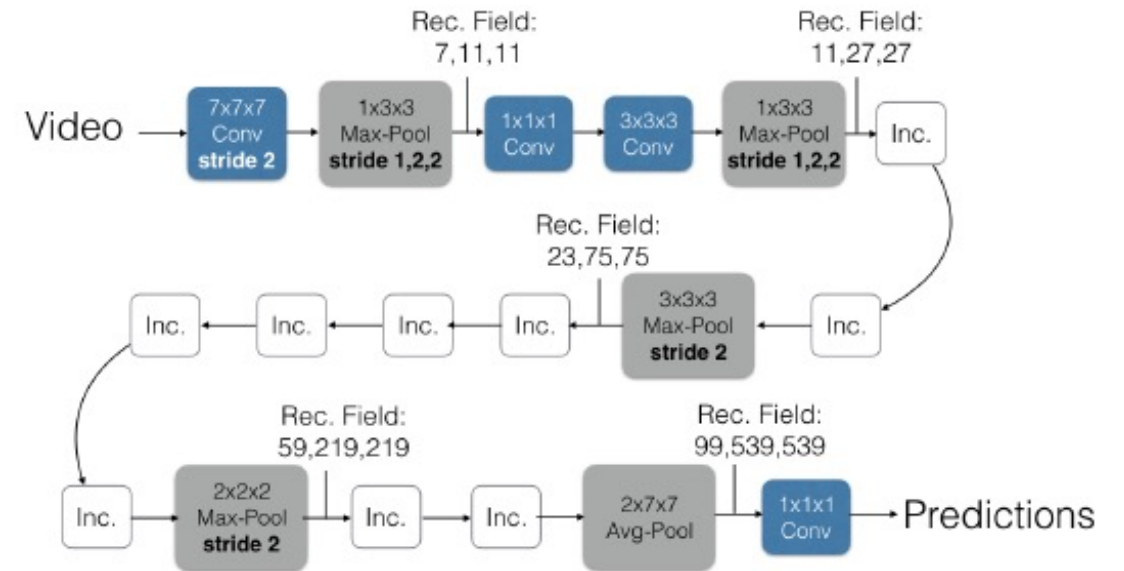
2D to 3D Convolutions



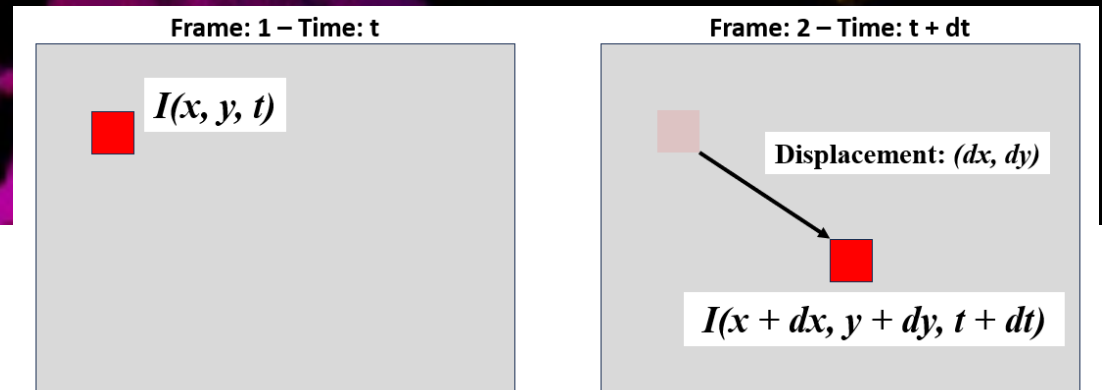
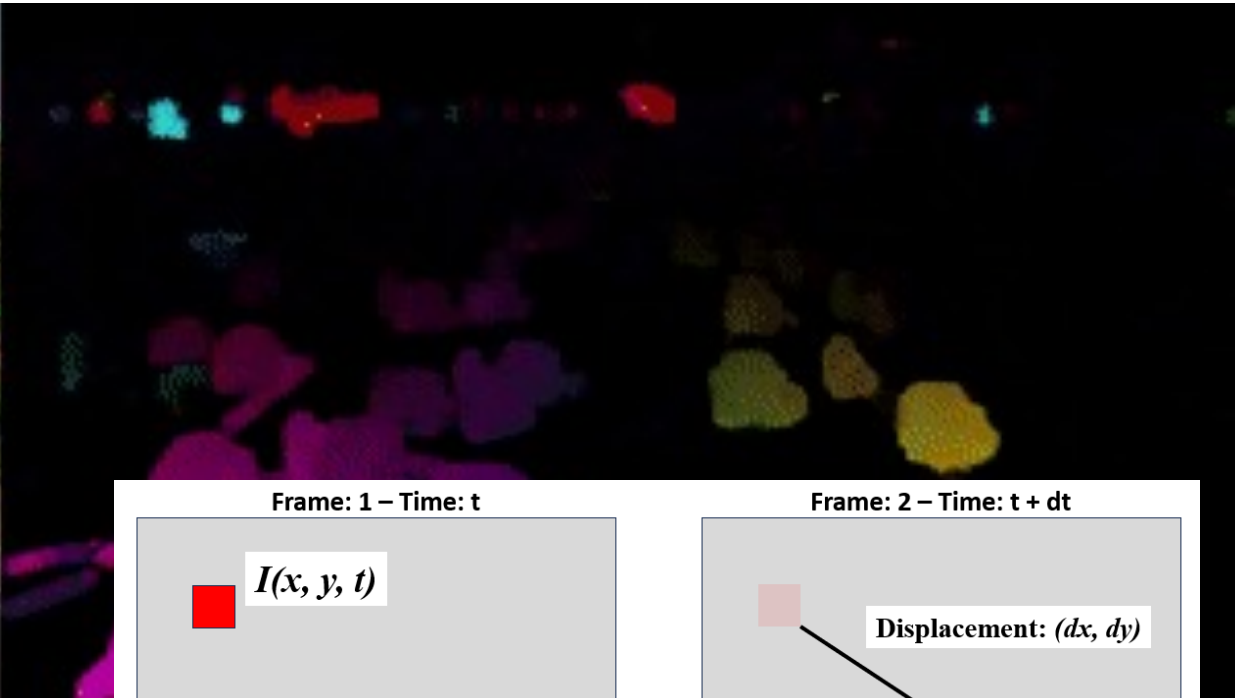
I3D Network



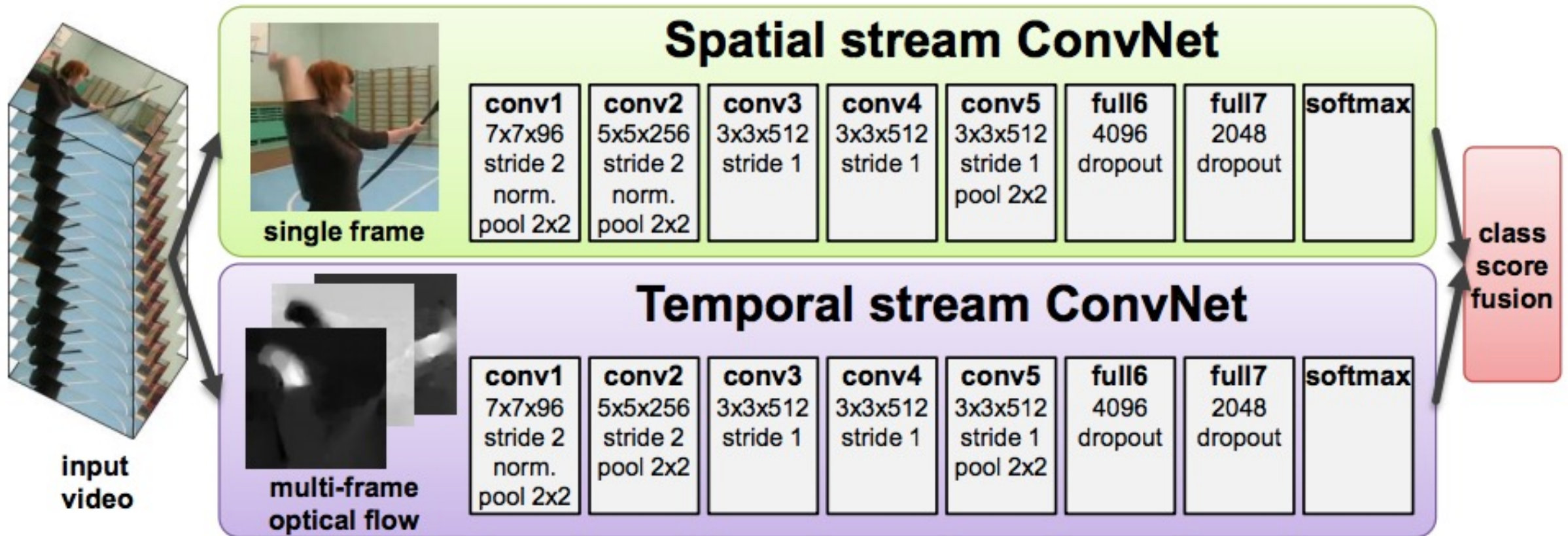
Inflated Inception-V1



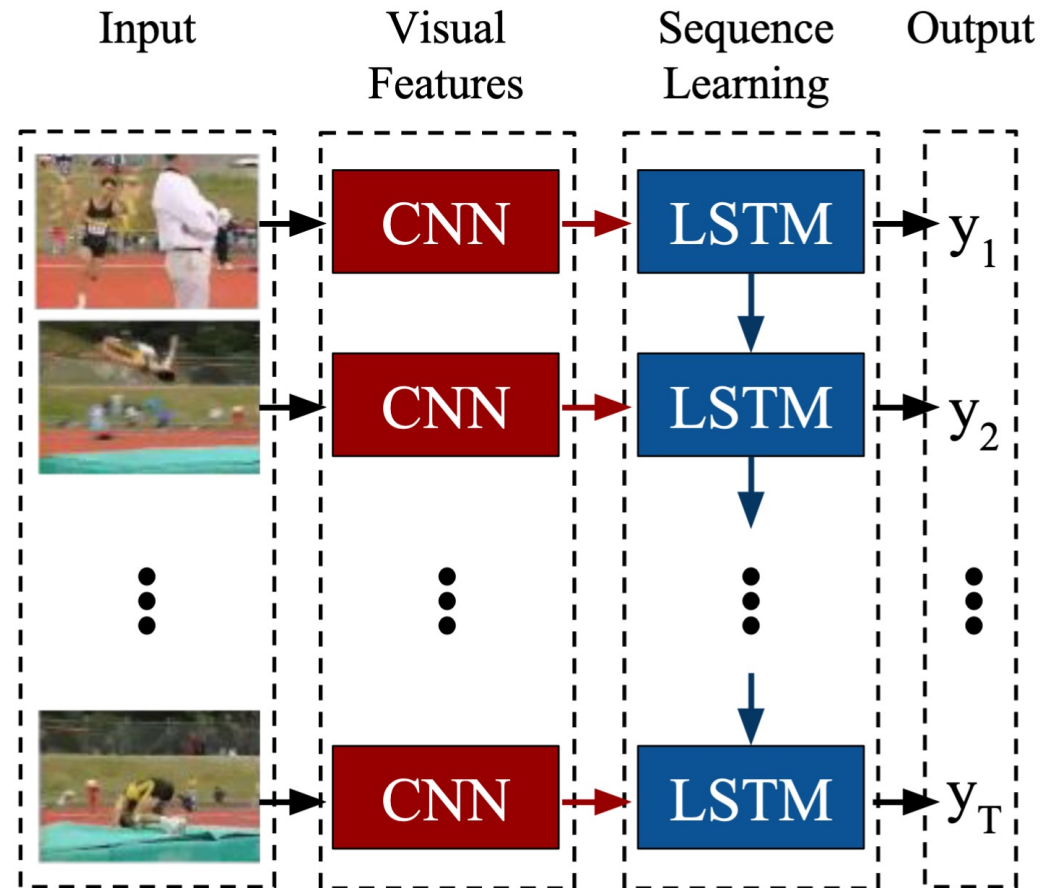
Optical Flow (review)



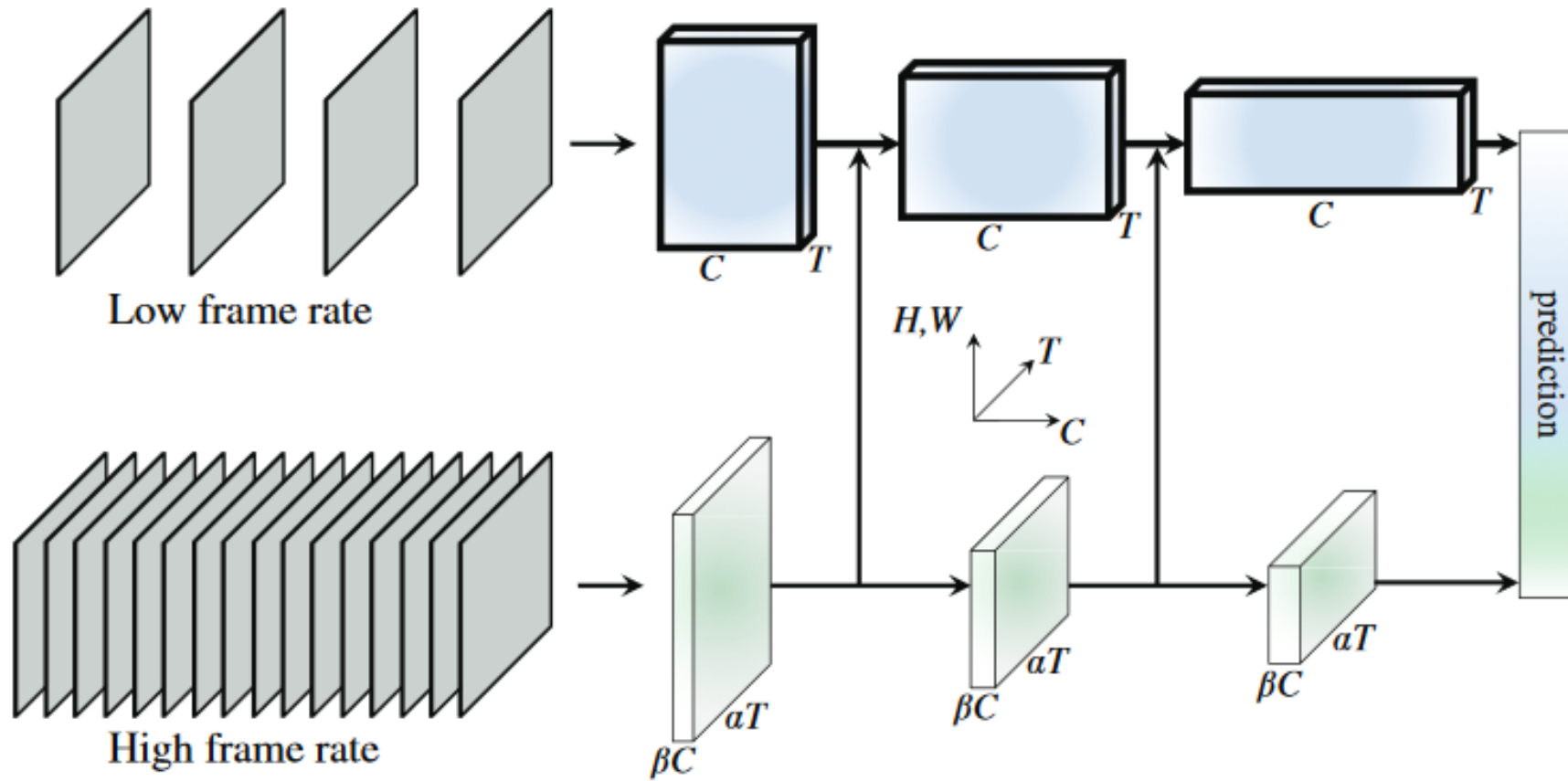
Two-stream networks



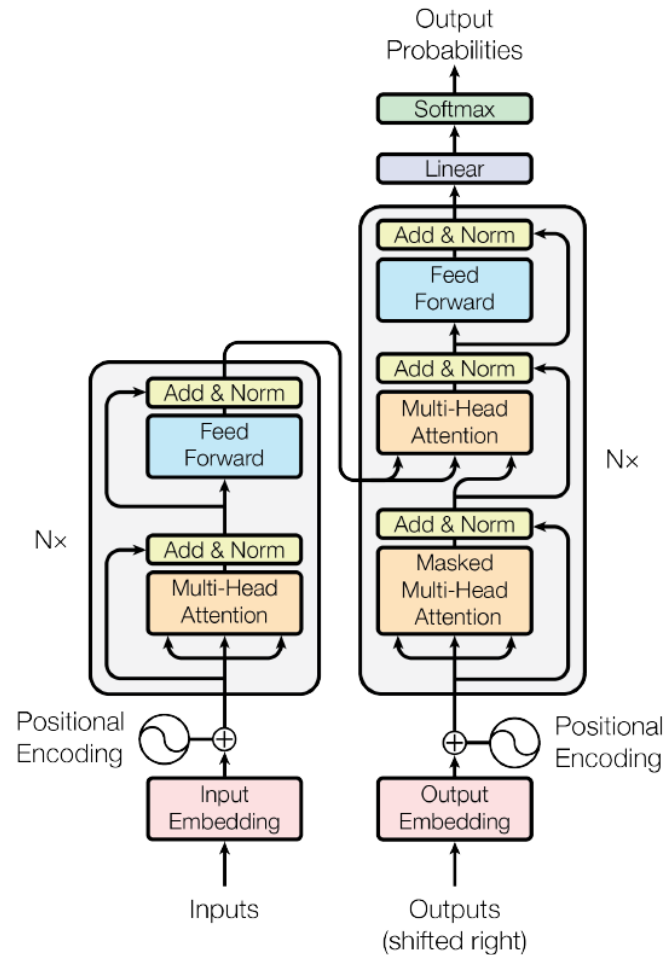
Convolutional + Recurrent Neural Networks



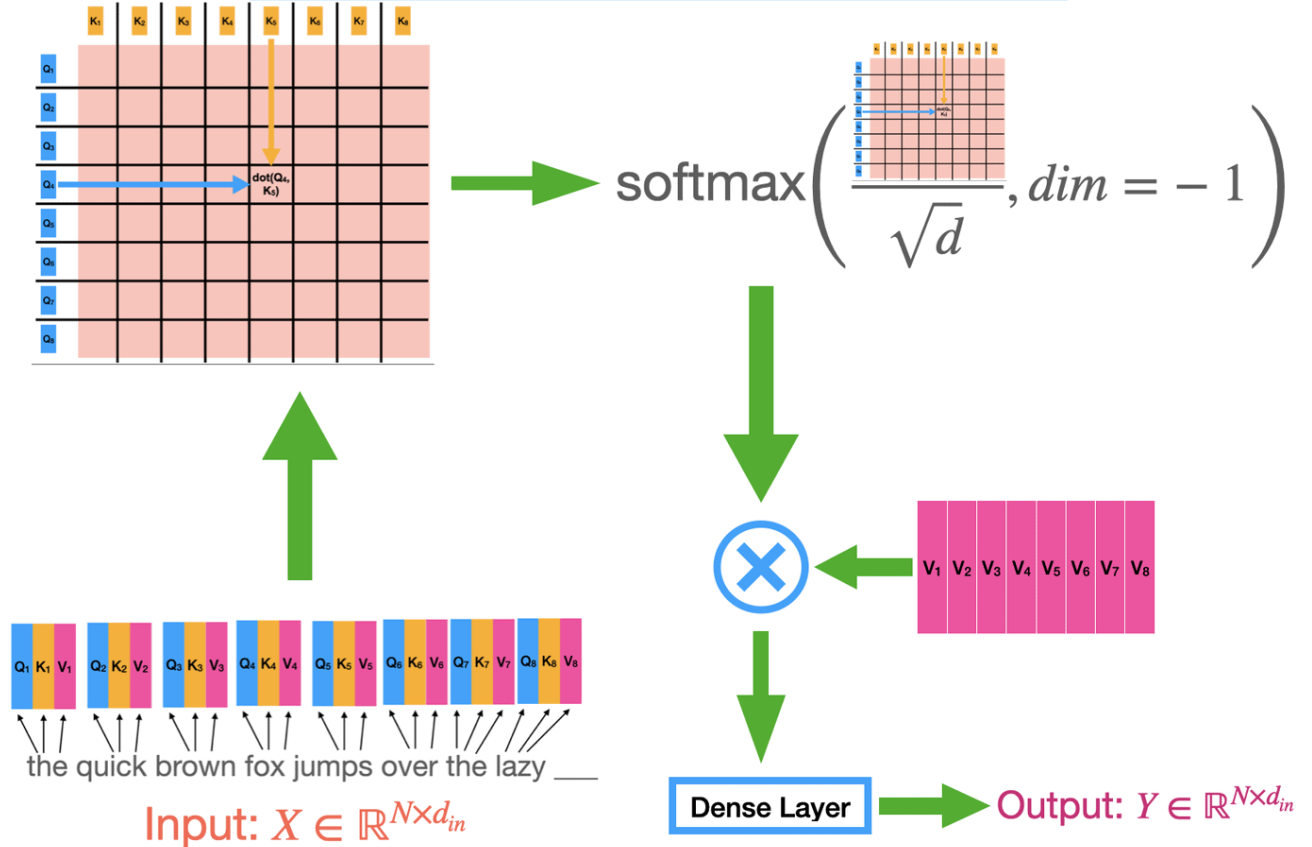
SlowFast Networks



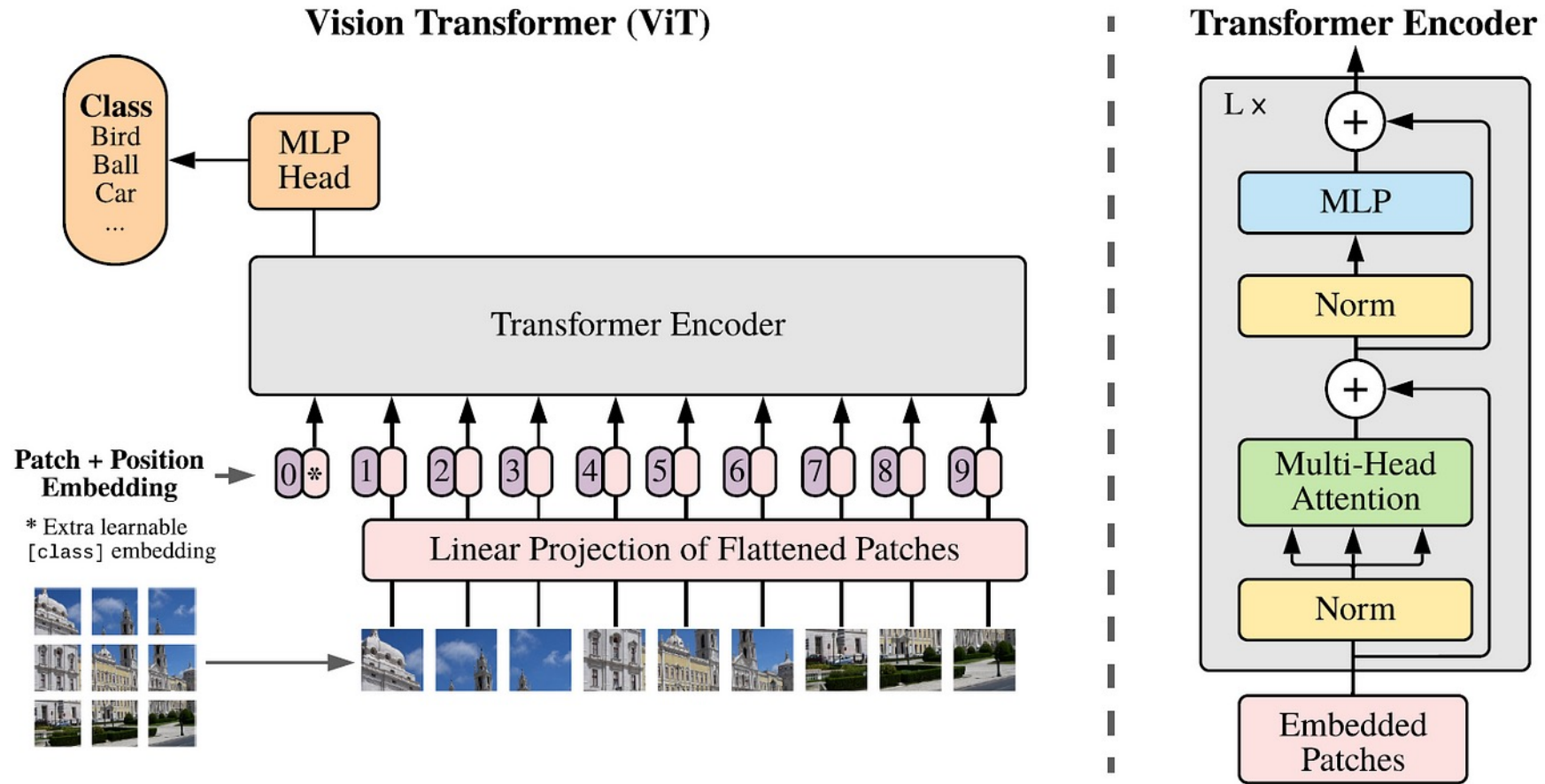
Transformer (review)



Neural Self Attention Mechanism

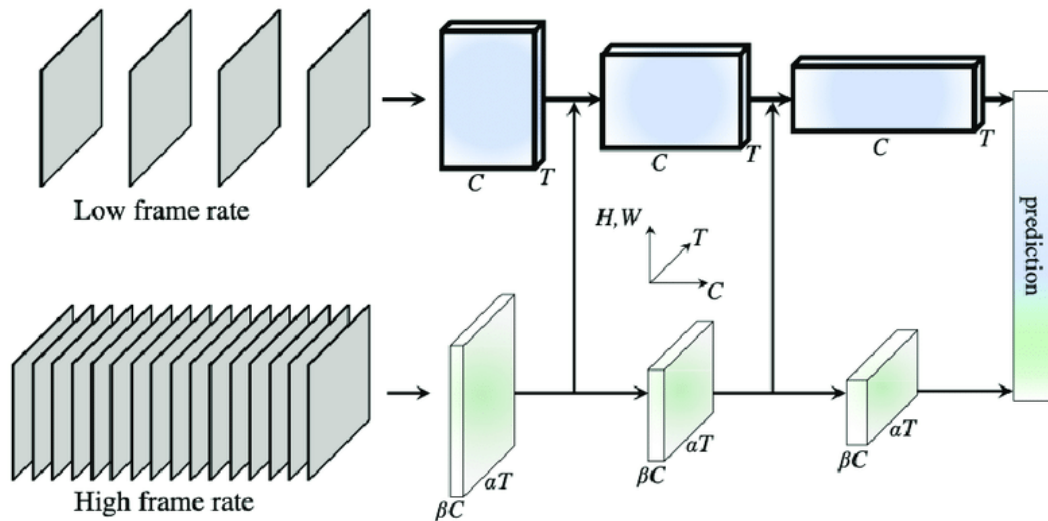


Vision Transformer, ViT (review)

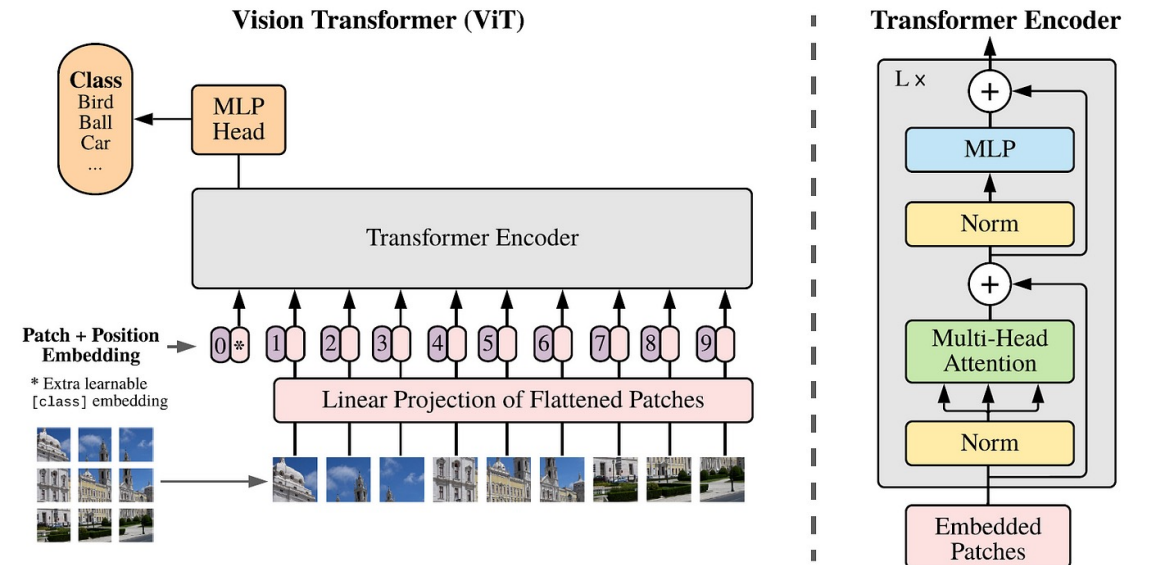


CNNs vs Transformer models

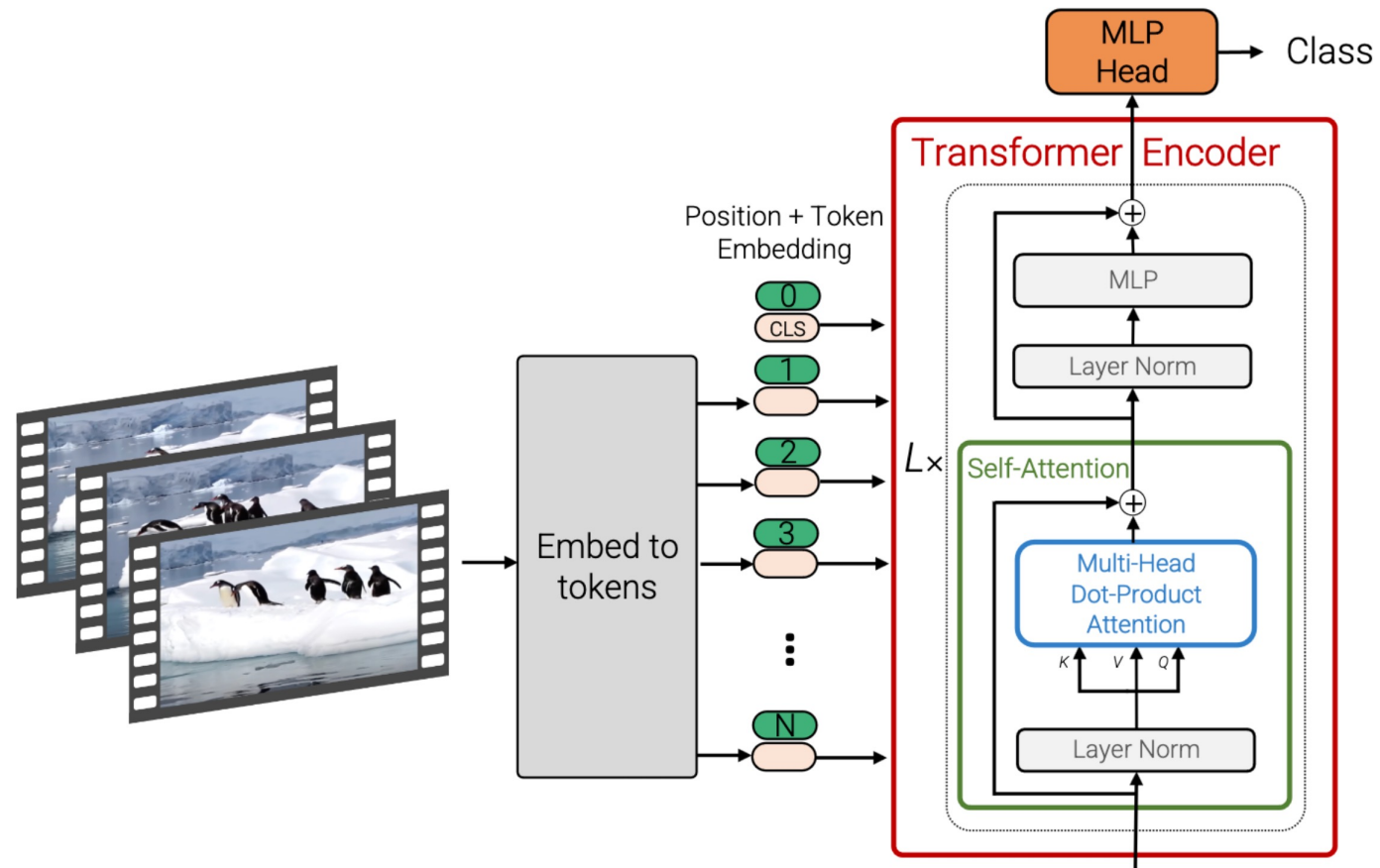
- Convolutions
- Pooling layers
- Local receptive fields



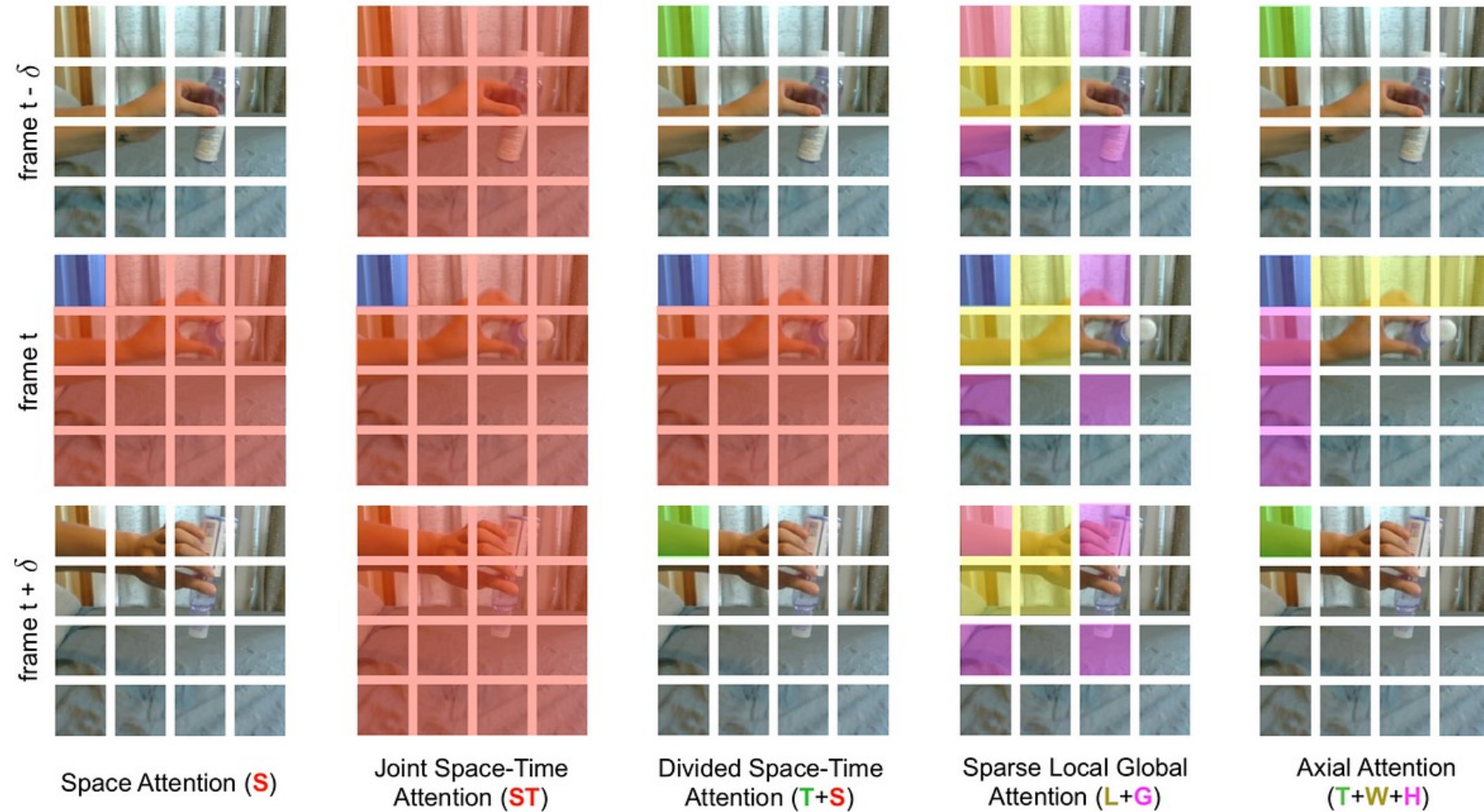
- Patches as tokens
- Self-attention across tokens
- Modeling across the image



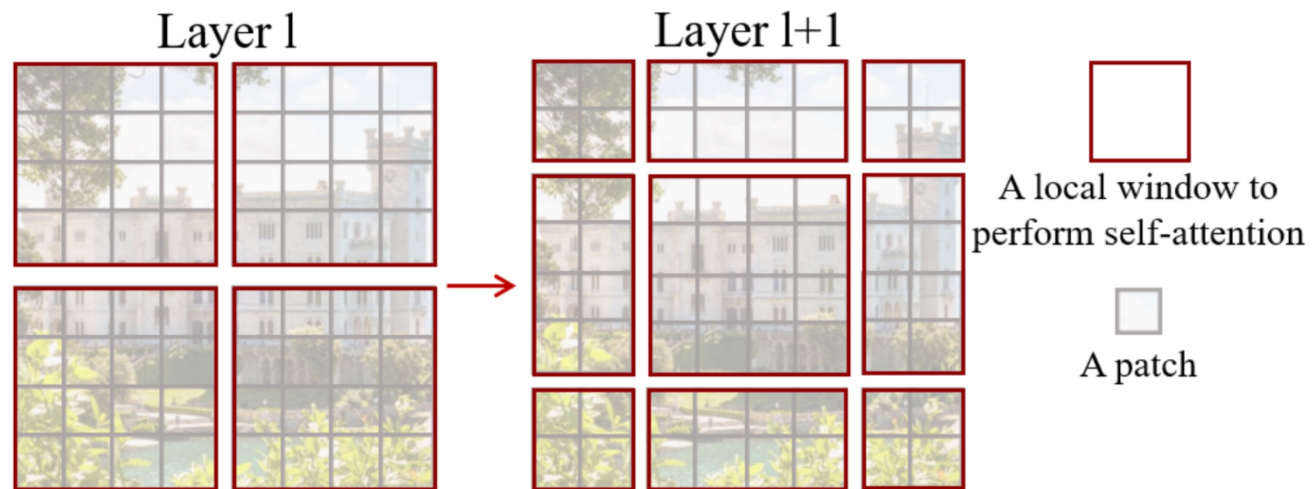
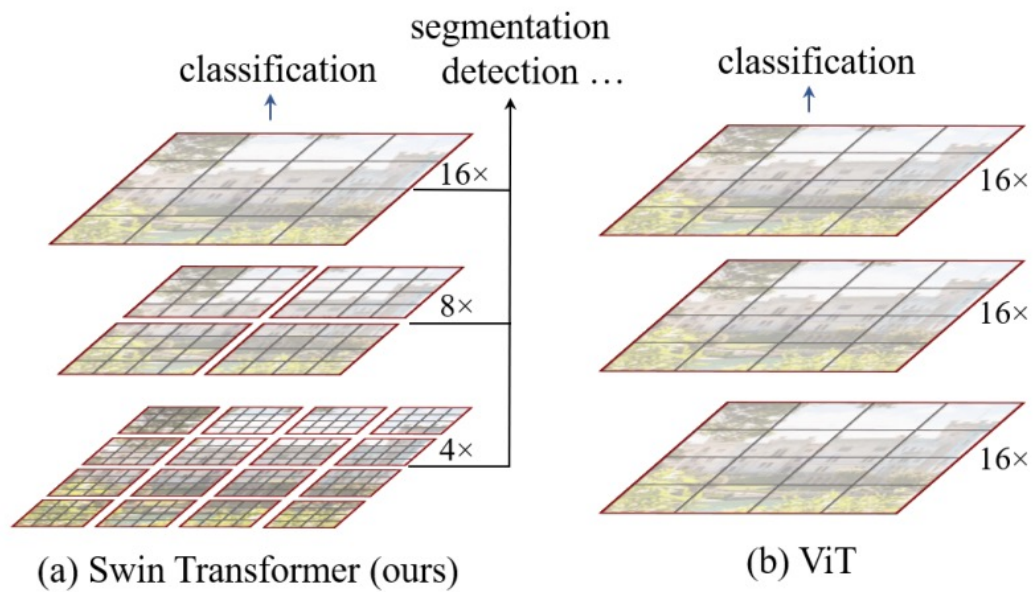
Video Vision Transformer



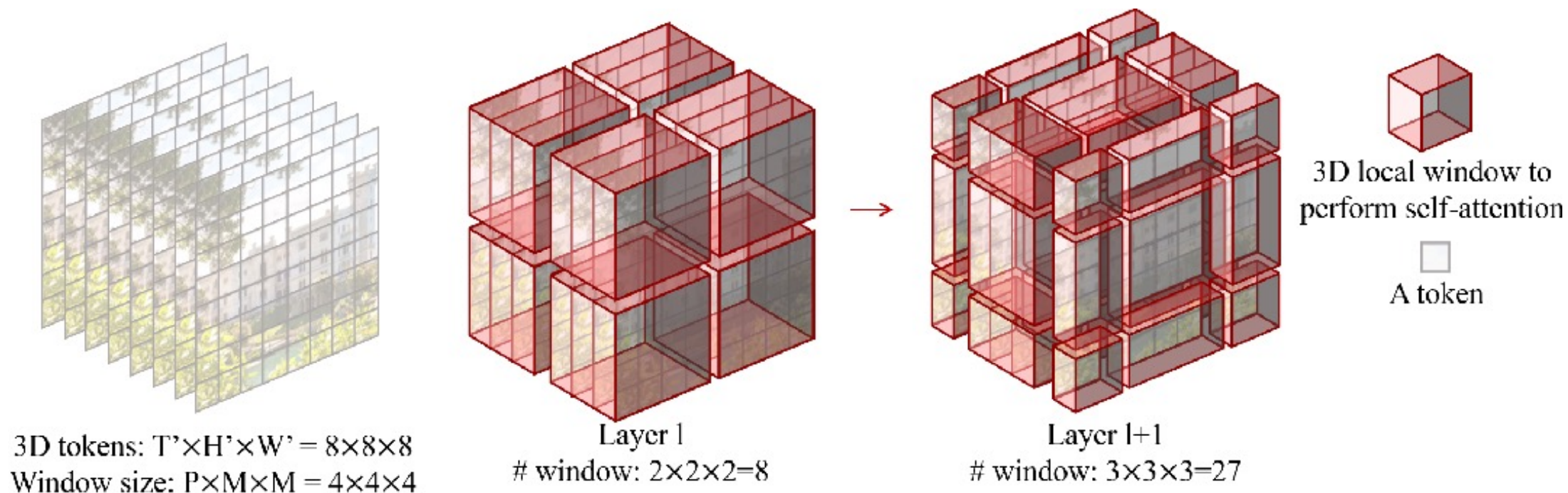
TimeSFormer: Where to attend?



Swin Transformer



Video Swin Transformer



Questions?