

Transformer : ViT with patch embedding -----> : Feature flow across two blocks : 2D feature maps