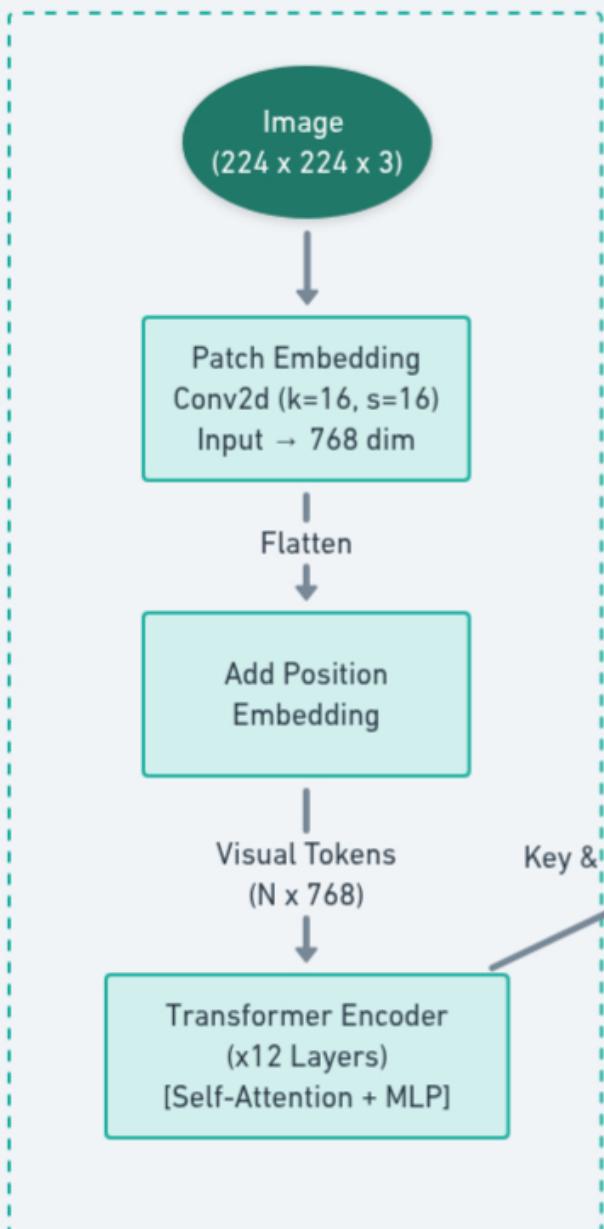


BERT Decoder (Generative)

ViT Encoder (Vision Transformer)



Key & Values (Visual Features)

