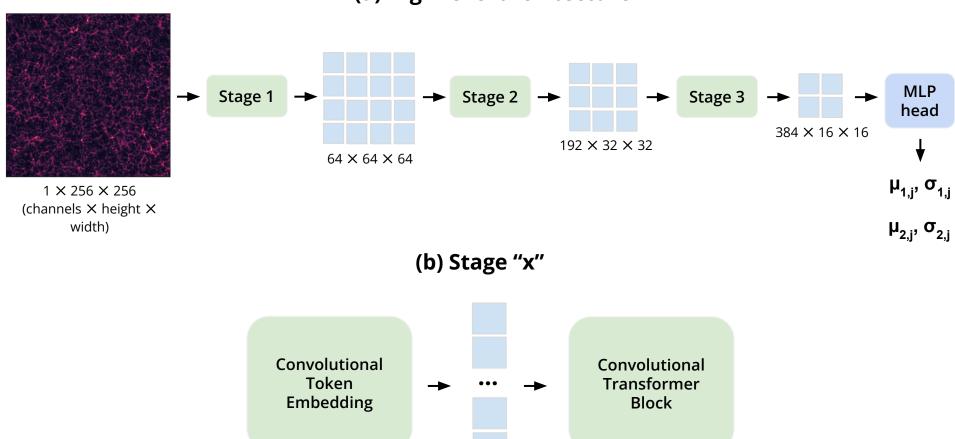
## (a) High-level architecture



tokens

2D convolution operation

 $\times$  N