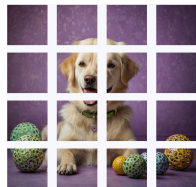


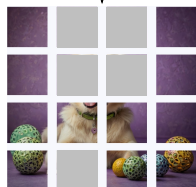
Original image



Image Patches



Blockwise Masking



Flatten

[S]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

[M]

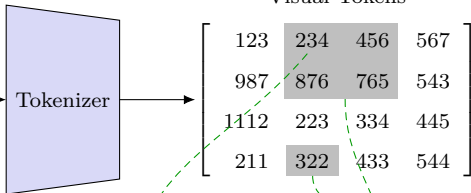
[M]

[M]

[M]

[M]

[M]

Position  
EmbeddingPatch  
EmbeddingUnused During  
Pre-TrainingReconstructed  
image

234 456

876 765

322

Masked Image Modeling Head

 $h_2^L$  $h_3^L$  $h_6^L$  $h_7^L$  $h_{14}^L$ 

BEiT Encoder