

modifier

θ



latent space
Z

**Encoder
Network**

Input

**images,
sentence,**

**Decoder
Network**

output

**images,
sentence,**

