

DDPM (V-Net, noise( $\varepsilon$ ), etc.)

↪ DDIM (non-stoch)  $\rightarrow$  improved DDPM better noise timesteps, etc.

Latent (LDM): pixel space  $\rightarrow$  latent space (w/ percept. compulsion autoencoder)  
intra, cross-affiliation ( $A_{t,t} = \text{softmax}\left(\frac{Q_{t,t}}{T}\right) \cdot V$ )  
make w/ input reverse trained weights  
add back (res-10.4)  
encode noisy img  $\rightarrow$  time  $\rightarrow$  transformer block  $\rightarrow$  unet  $\rightarrow$   $\varepsilon$  latent  $\rightarrow$  image  
text input  $\rightarrow$  (tokenized, encoded)  $\rightarrow$  (tr. w/ on low-res steps)

stable Diffusion: using LDM, CLIP

↪ SDXL; larger model, additional conditioning parameters (h/w)  $\cup$  { $t_{-B}$ }  
(cropping)  
better training data (more sizes), better text encoder

transformer: add position encoding to img patches, sent through transformer blocks,  
condition on aCaLN, 1-patch  
uses Latent.

video: 16 frames / training data, similar mechanics as DDPM V-Net,