**Low-Resolution Stage**

Reference U-Net

K,V concat

Denoising U-Net

$\tilde{x}_g$

$\tilde{x}_p$ $\tilde{\epsilon}$

$\tilde{x}_r$

Res Block    Self-Attention Block

**High-Resolution Stage**

Reference U-Net

K,V concat

Denoising U-Net

$x_g$

**Resize ratio:** $\sigma$

$\times \beta$ $\times \alpha$

$\tilde{x}_r$ $\epsilon$ $x_p$

$x_r$