input audio stylized speech input style reference animation style basis a style conditional encoder latent (VAE) diffusion b