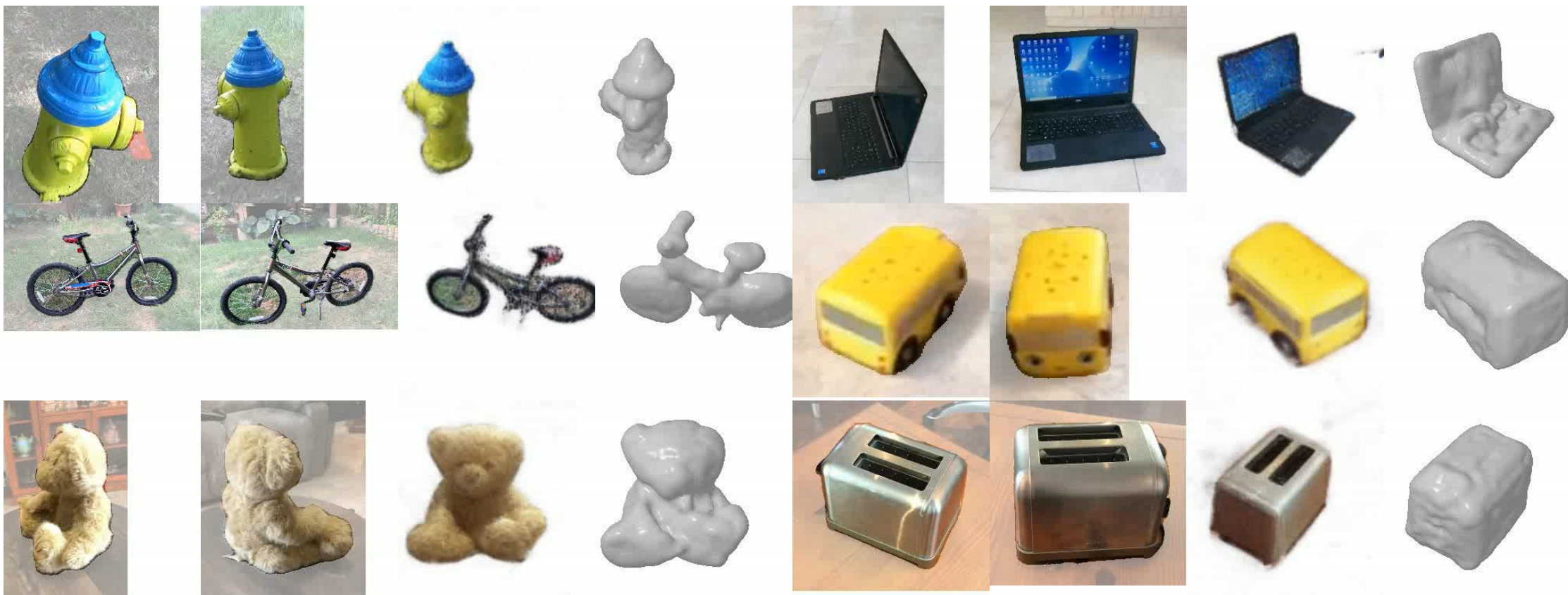
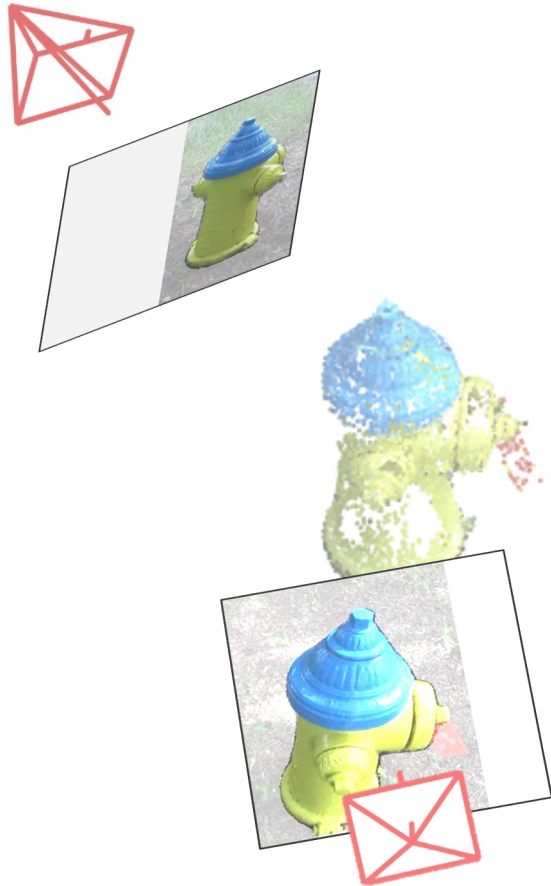


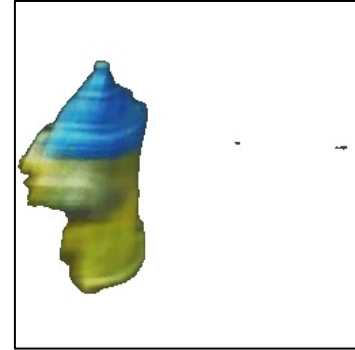
SparseFusion: Distilling View-conditioned Diffusion for 3D Reconstruction



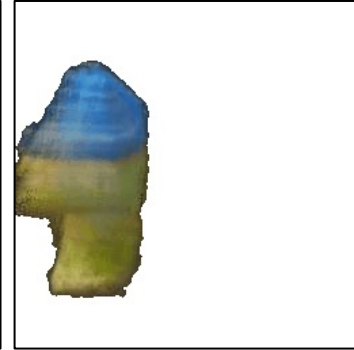
Task: Sparse-view Reconstruction



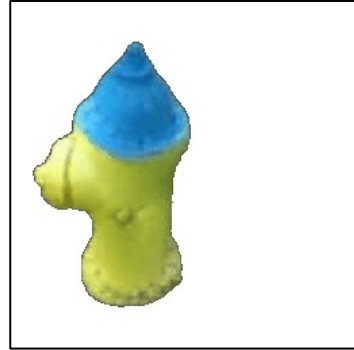
PixelNeRF



NerFormer



ViewFormer

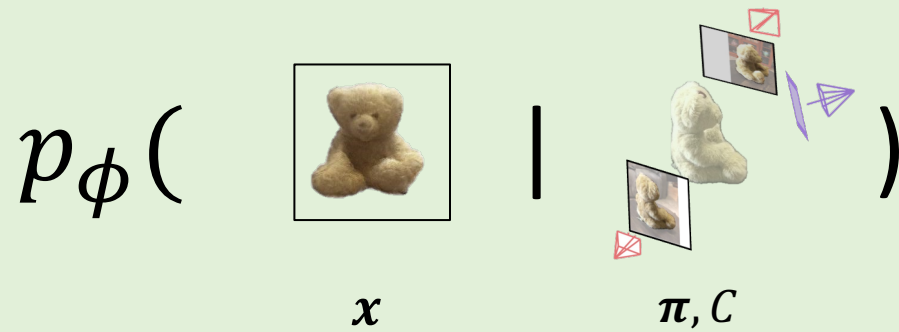


SparseFusion



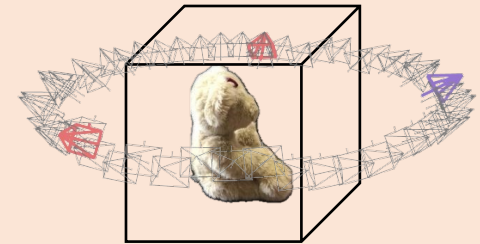
SparseFusion Overview

View-conditioned Latent Diffusion



Diffusion Distillation

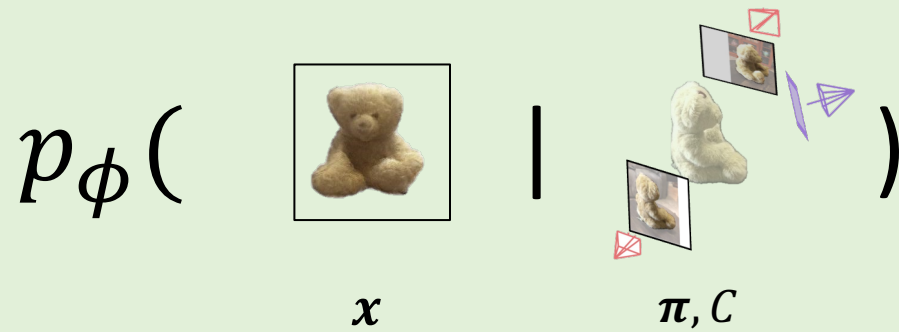
$$(x, y, z) \rightarrow f_\theta \rightarrow (rgb, \sigma)$$



$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_\phi(f_\theta(\pi) | \pi, C)]$$

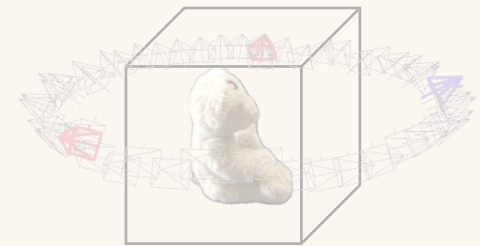
SparseFusion Overview

View-conditioned Latent Diffusion



Diffusion Distillation

$$(x, y, z) \rightarrow f_\theta \rightarrow (rgb, \sigma)$$



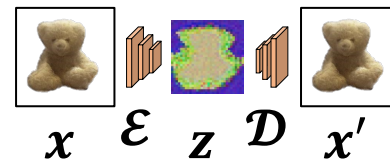
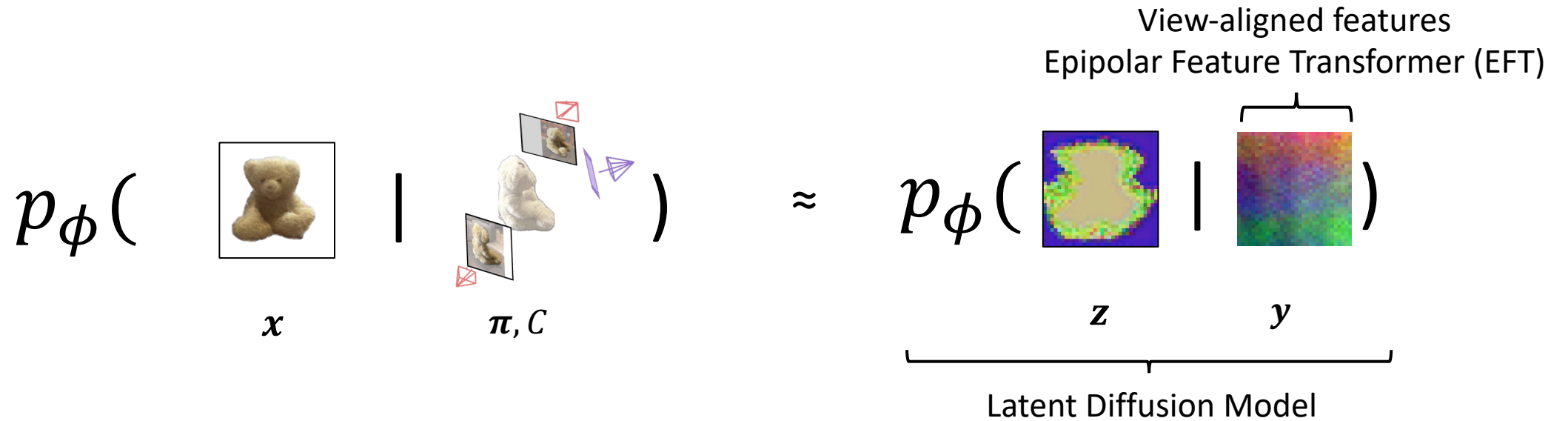
$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_{\phi}(f_{\theta}(\pi) | \pi, C)]$$

View-conditioned Latent Diffusion (VLDM)

$$p_{\phi} \left(\begin{array}{c} \boxed{\text{Image of teddy bear}} \\ x \end{array} \mid \begin{array}{c} \text{3D model of teddy bear} \\ \text{with camera frustums} \\ \pi, C \end{array} \right) \approx p_{\phi} \left(\begin{array}{c} \boxed{\text{Image of teddy bear}} \\ x \end{array} \mid \begin{array}{c} \boxed{\text{Latent space visualization}} \\ y \end{array} \right)$$

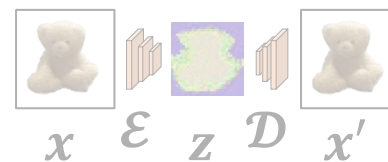
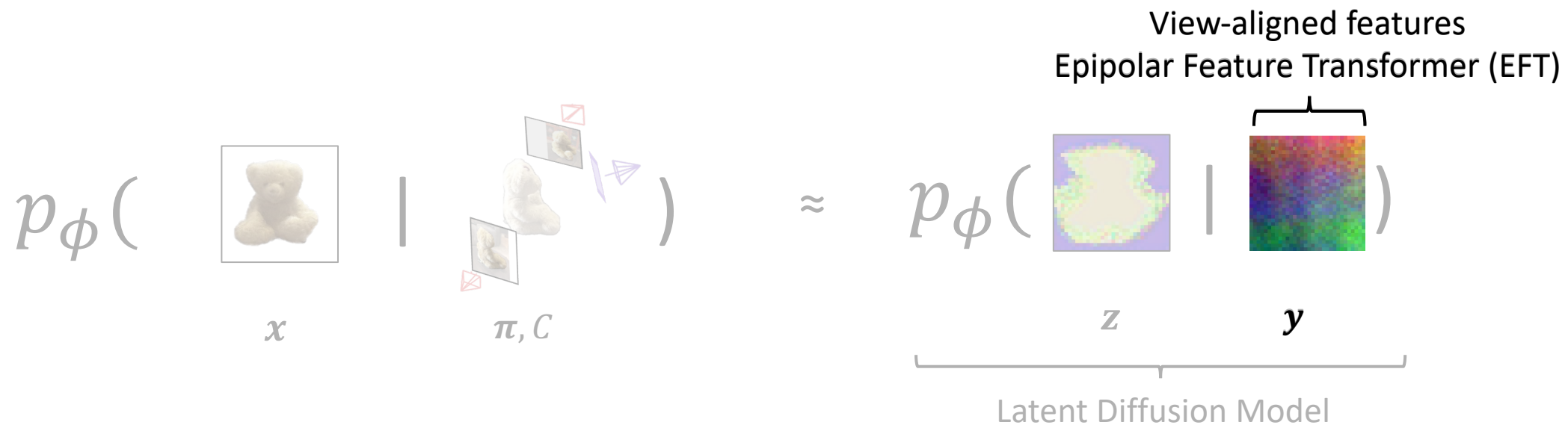
The diagram illustrates the VLDM process. On the left, a probability distribution p_{ϕ} is shown conditioned on an input image x (a teddy bear) and a set of camera parameters π, C (represented by a 3D model and camera frustums). This is approximately equal to a probability distribution p_{ϕ} conditioned on the same input image x and a latent space visualization y (a colorful noise pattern).

View-conditioned Latent Diffusion (VLDM)



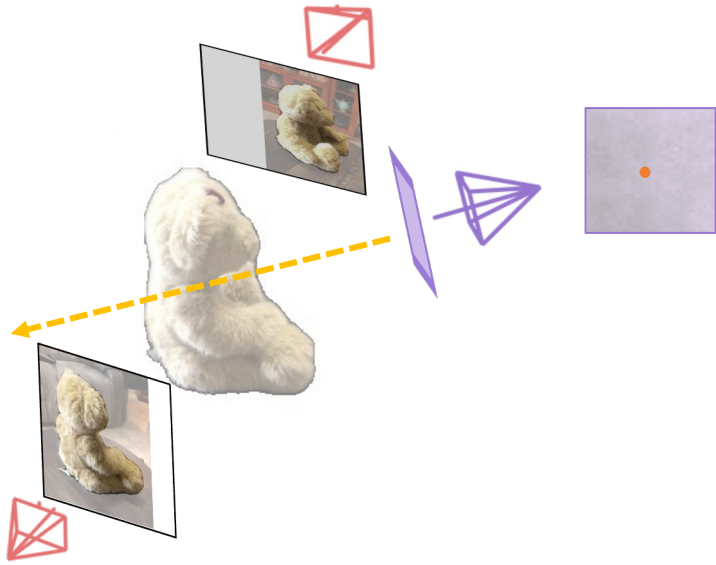
VAE from High-Resolution Image Synthesis
with Latent Diffusion Model
Rombach et al. CVPR 2022.

View-conditioned Latent Diffusion (VLDM)

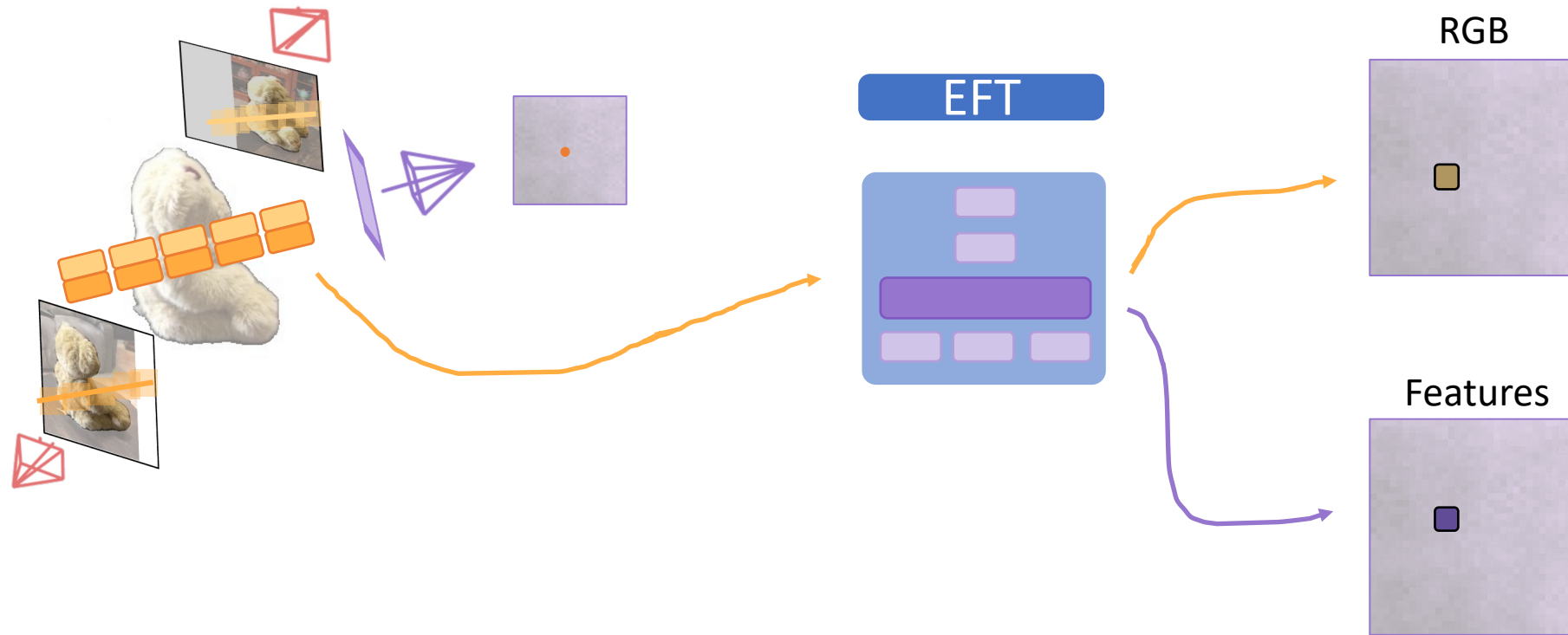


VAE from High-Resolution Image Synthesis
with Latent Diffusion Model
Rombach et al. CVPR 2022.

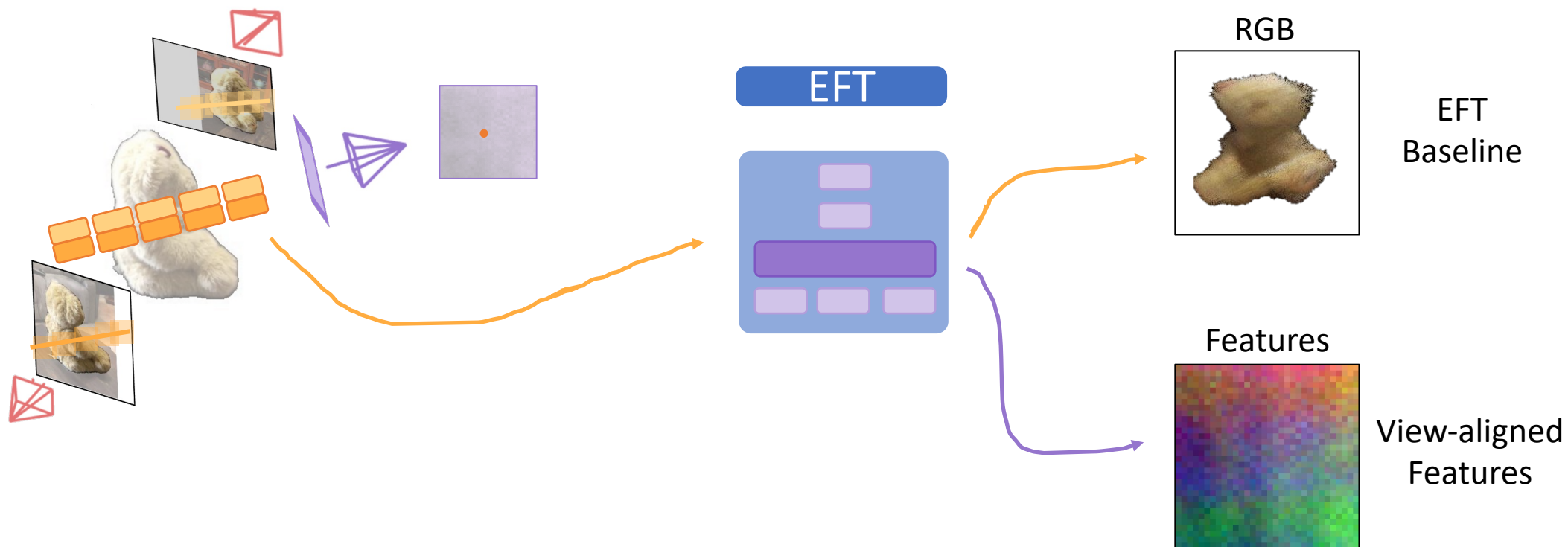
Epipolar Feature Transformer (EFT)



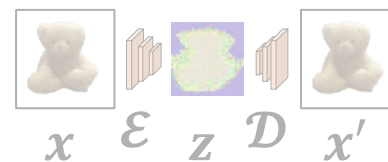
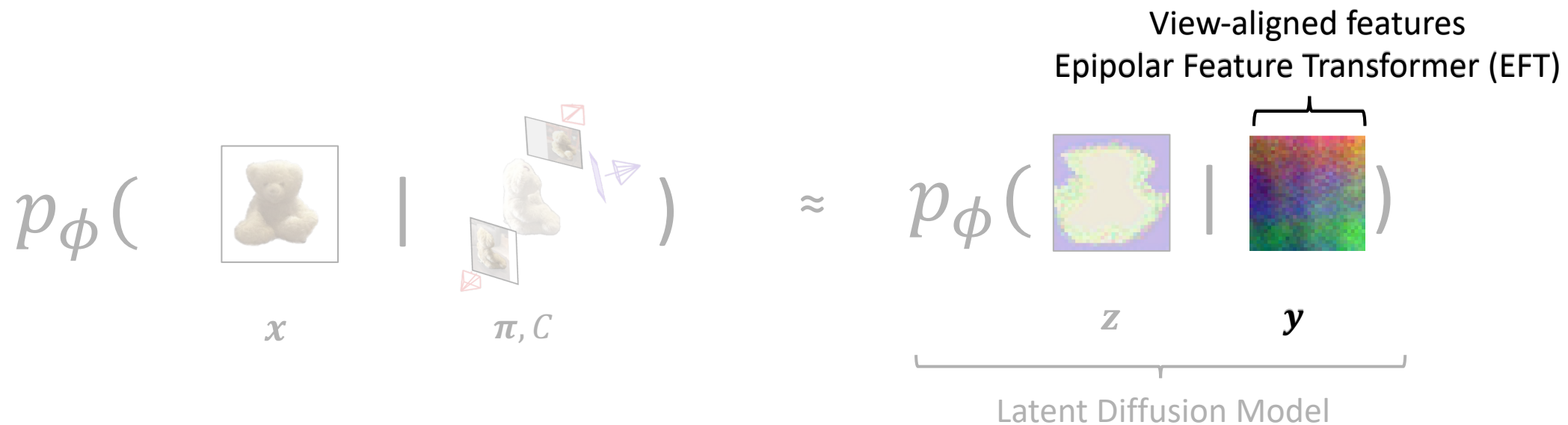
Epipolar Feature Transformer (EFT)



Epipolar Feature Transformer (EFT)

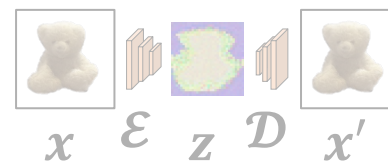
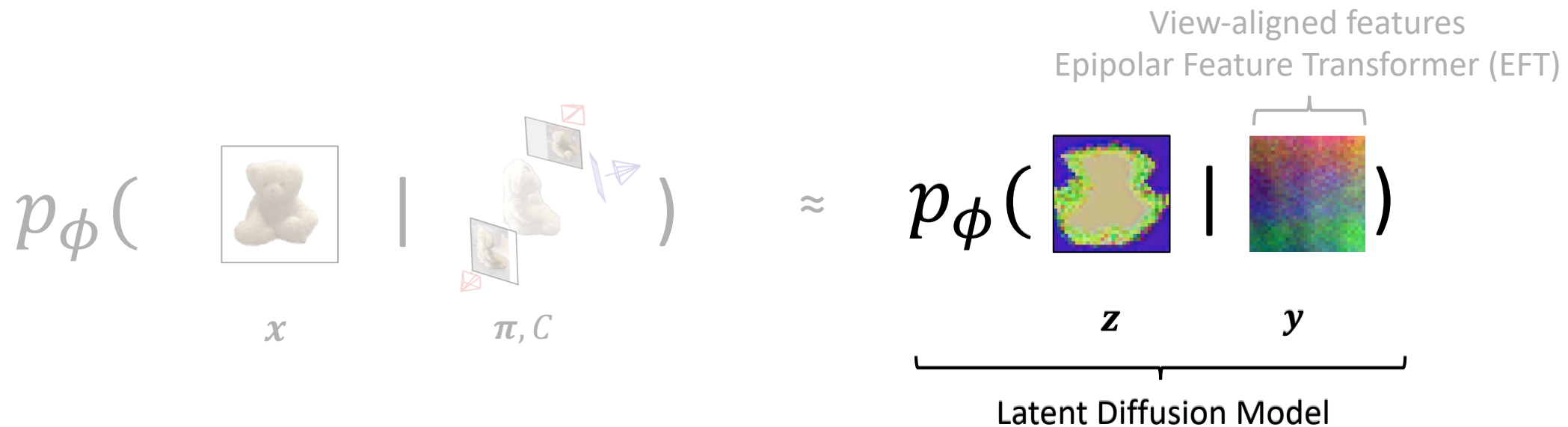


View-conditioned Latent Diffusion (VLDM)



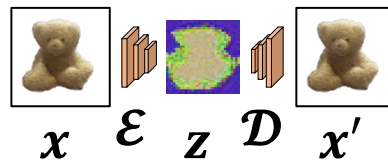
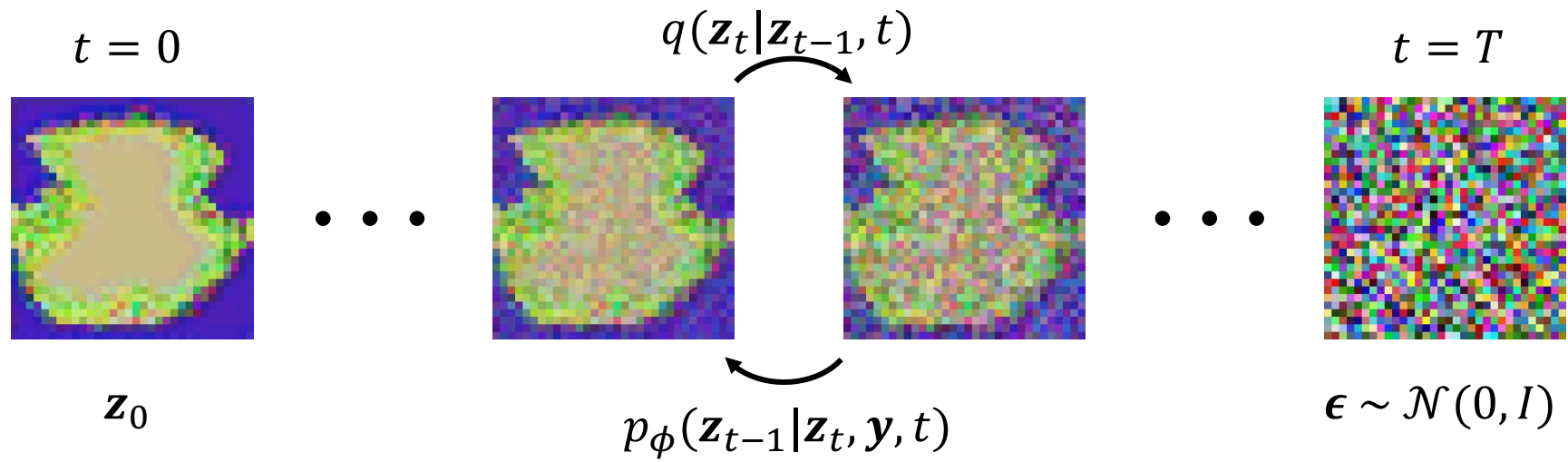
VAE from High-Resolution Image Synthesis
with Latent Diffusion Model
Rombach et al. CVPR 2022.

View-conditioned Latent Diffusion (VLDM)



VAE from High-Resolution Image Synthesis
with Latent Diffusion Model
Rombach et al. CVPR 2022.

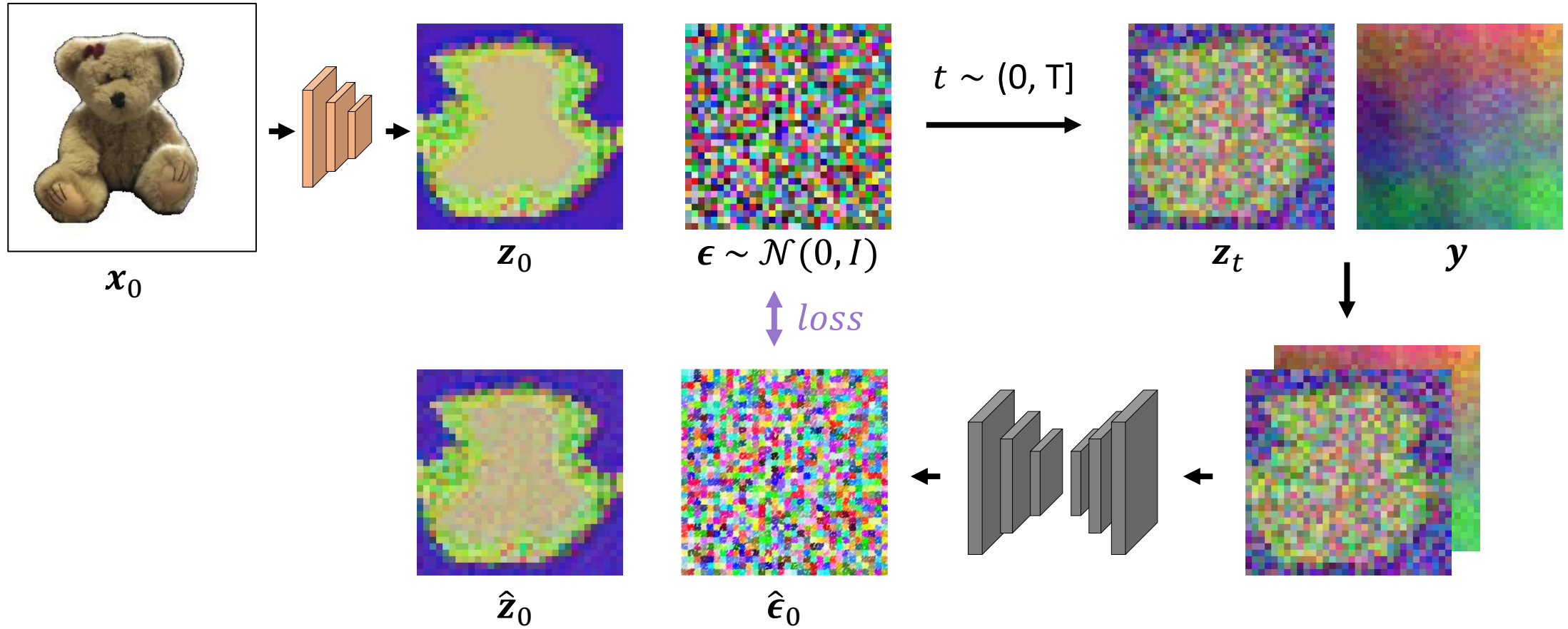
Latent Diffusion Model



Denoising Diffusion Probabilistic Models
Ho et al. NeurIPS 2020.

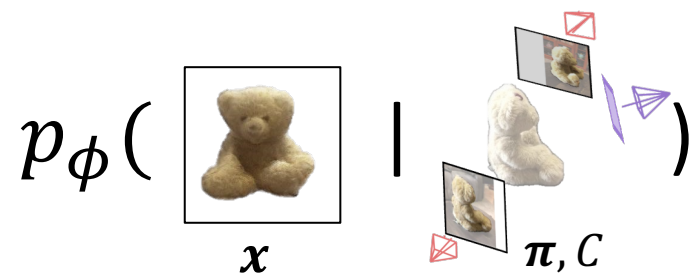
VAE from High-Resolution Image Synthesis with Latent Diffusion Model
Rombach et al. CVPR 2022.

View-conditioned Latent Diffusion (VLDM)



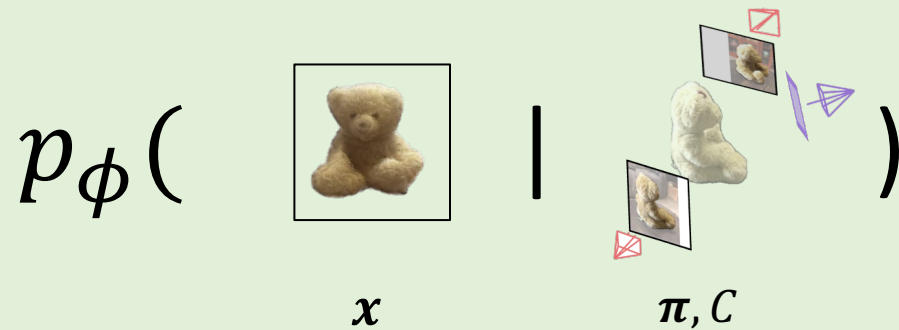
$$p_{\phi} \left(\begin{array}{c} \text{[Latent Image]} \\ z \end{array} \mid \begin{array}{c} \text{[Reconstructed Image]} \\ y \end{array} \right)$$

View-conditioned Latent Diffusion (VLDM)



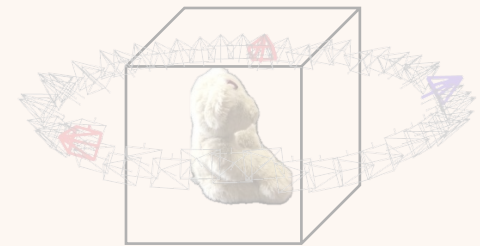
SparseFusion Overview

View-conditioned Latent Diffusion



Diffusion Distillation

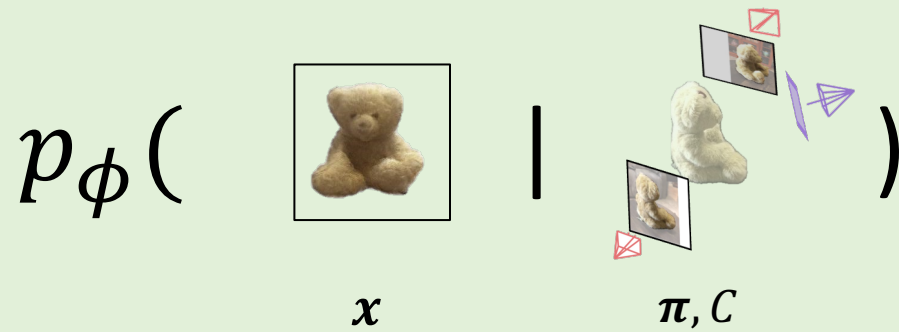
$$(x, y, z) \rightarrow f_\theta \rightarrow (rgb, \sigma)$$



$$\min_{\theta} \mathbb{E}_{\boldsymbol{\pi}} [-\log p_\phi(f_\theta(\boldsymbol{\pi}) \mid \boldsymbol{\pi}, \mathcal{C})]$$

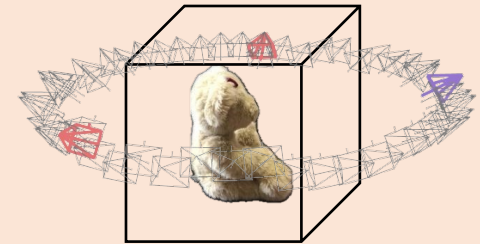
SparseFusion Overview

View-conditioned Latent Diffusion



Diffusion Distillation

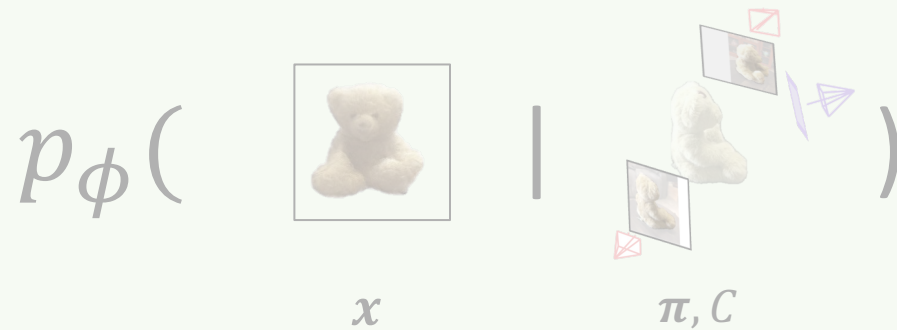
$$(x, y, z) \rightarrow f_\theta \rightarrow (rgb, \sigma)$$



$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_\phi(f_\theta(\pi) | \pi, C)]$$

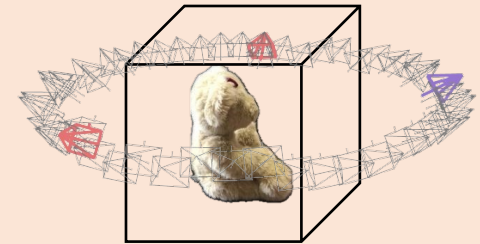
SparseFusion Overview

View-conditioned Latent Diffusion



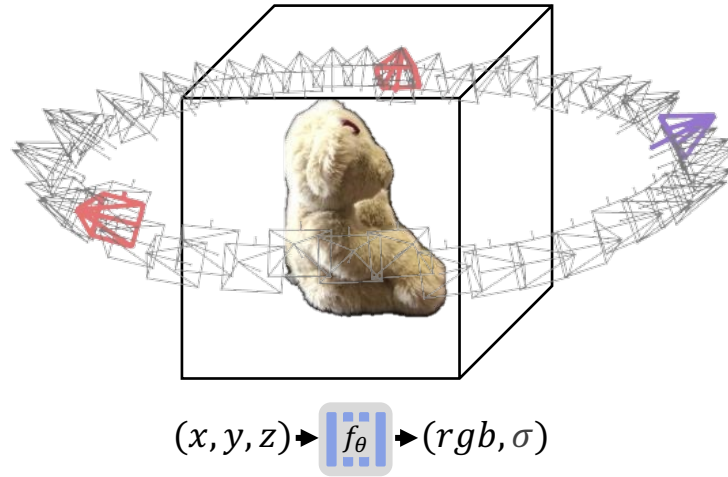
Diffusion Distillation

$$(x, y, z) \rightarrow f_\theta \rightarrow (rgb, \sigma)$$



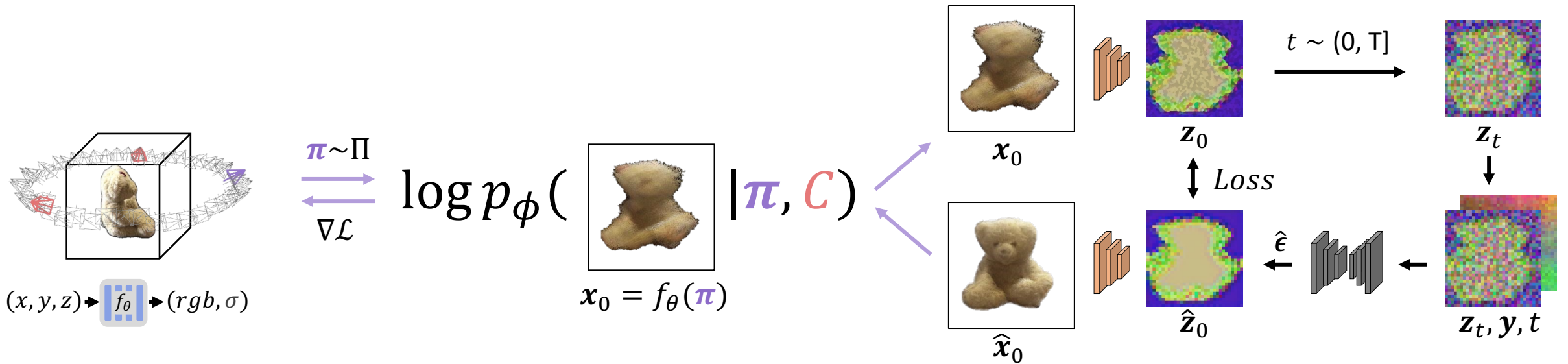
$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_\phi(f_\theta(\pi) | \pi, C)]$$

Diffusion Distillation



$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_{\phi}(f_{\theta}(\pi) | \pi, \mathcal{C})]$$

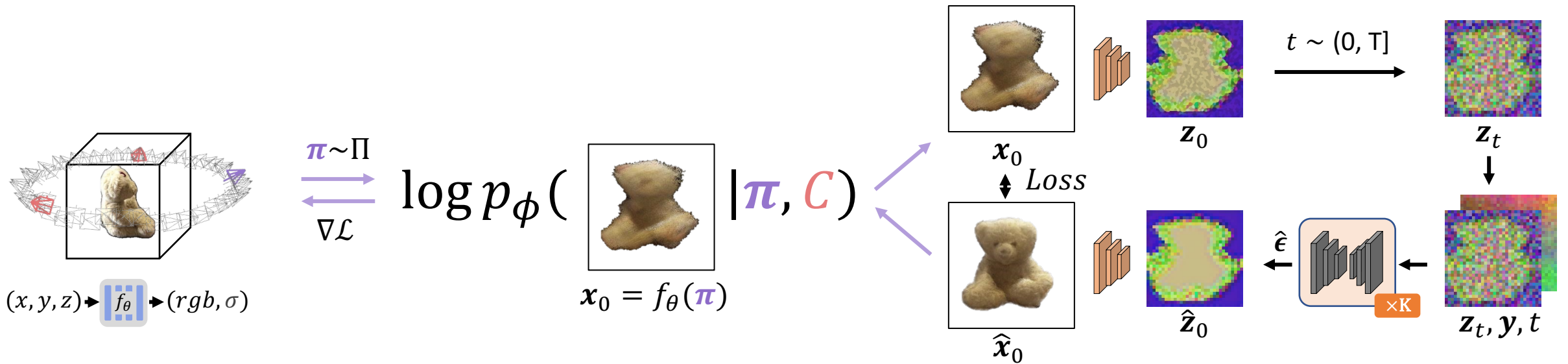
Diffusion Distillation



$$\log p_\phi(x_0 | \pi, C) \approx \mathbb{E}_t \|z_0 - \hat{z}_0\|^2$$

$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_\phi(f_\theta(\pi) | \pi, C)]$$

Diffusion Distillation

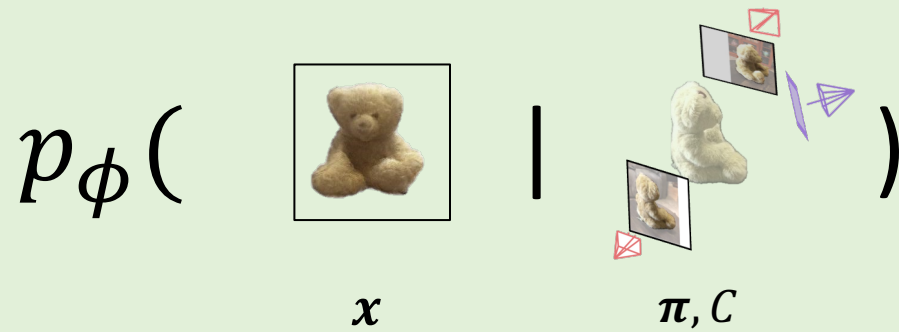


$$\log p_\phi(x_0 | \pi, C) \approx \mathbb{E}_t \|x_0 - \hat{x}_0\|^2$$

$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_\phi(f_\theta(\pi) | \pi, C)]$$

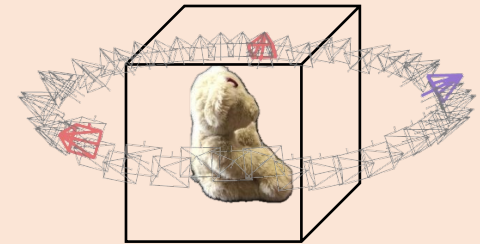
SparseFusion Overview

View-conditioned Latent Diffusion



Diffusion Distillation

$$(x, y, z) \rightarrow f_\theta \rightarrow (rgb, \sigma)$$



$$\min_{\theta} \mathbb{E}_{\pi} [-\log p_\phi(f_\theta(\pi) | \pi, C)]$$

Results



Does not hallucinate unseen regions



Does not respect geometry



Hallucinates unseen region + respects geometry

Results



Does not hallucinate unseen regions



Does not respect geometry



Hallucinates unseen region + respects geometry

Results

CO3Dv2 Dataset



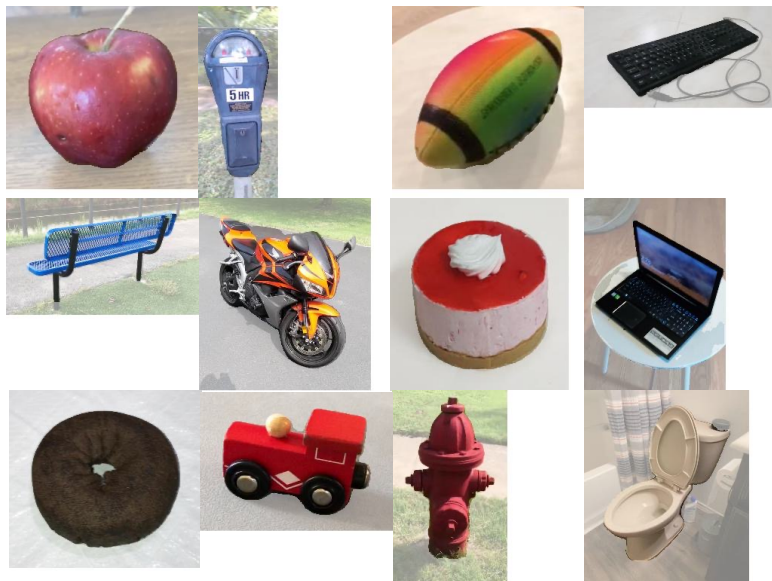
Reizenstein et al. ICCV 2021.

2-view Novel View Synthesis on 10 Categories

	PSNR \uparrow	LPIPS \downarrow
PixelNeRF	19.52	0.327
NerFormer	17.88	0.382
ViewFormer	18.37	0.282

Results

CO3Dv2 Dataset

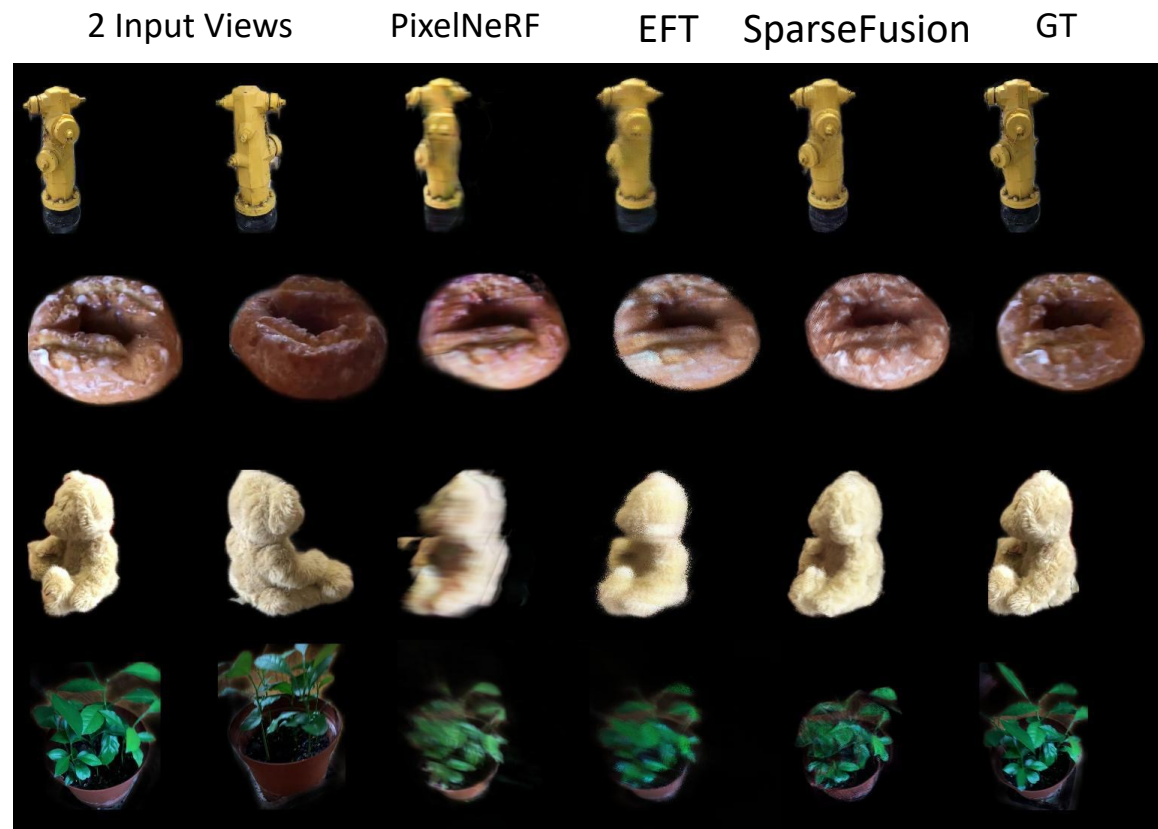
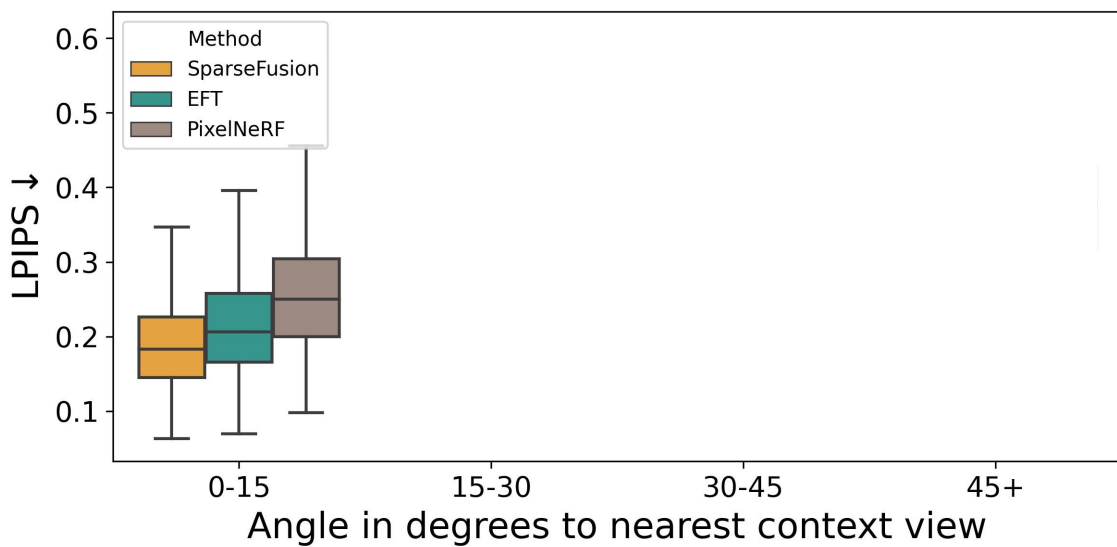
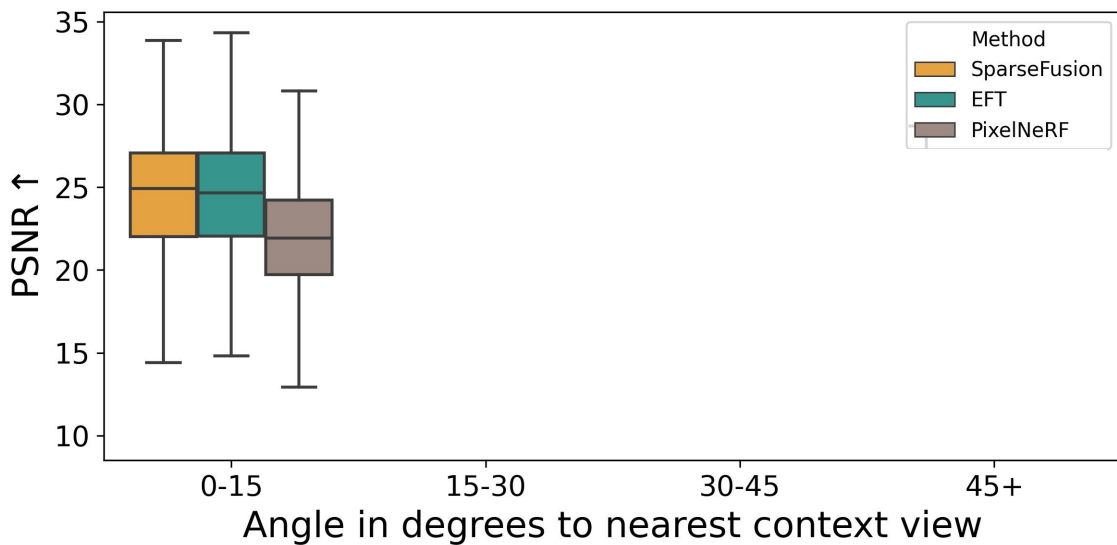


Reizenstein et al. ICCV 2021.

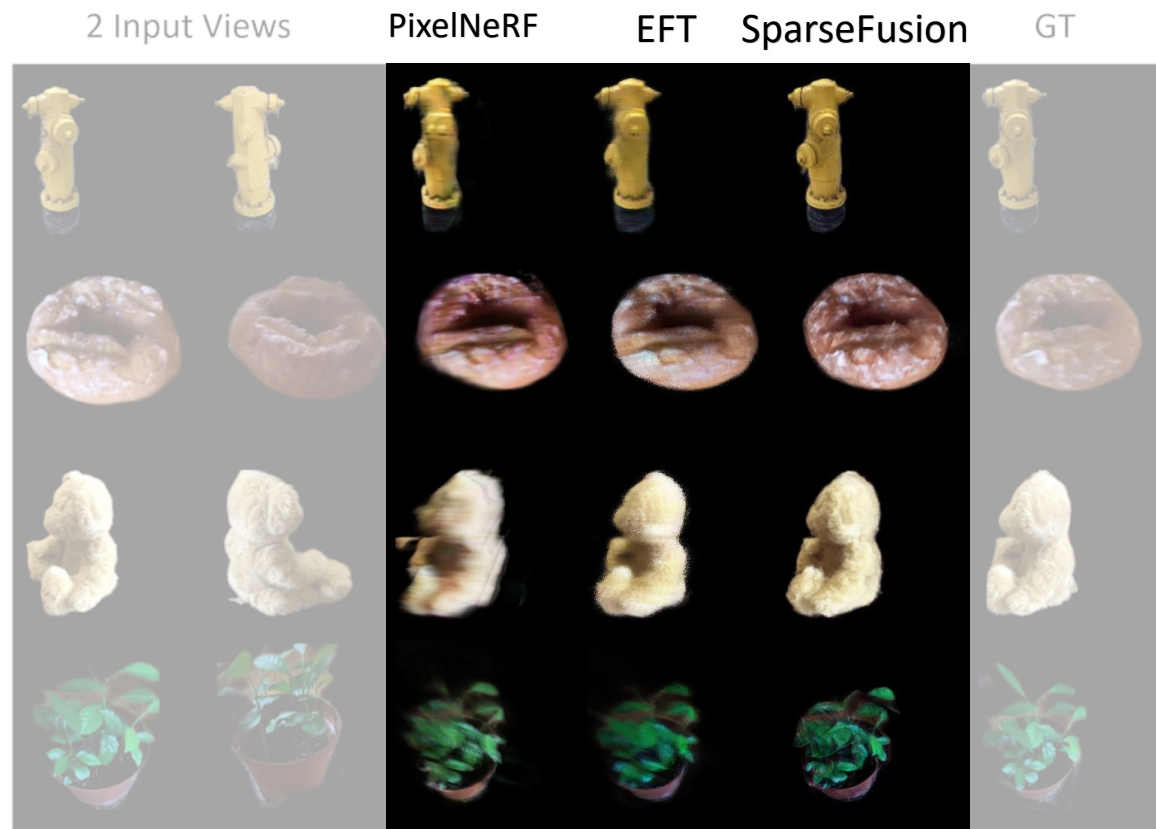
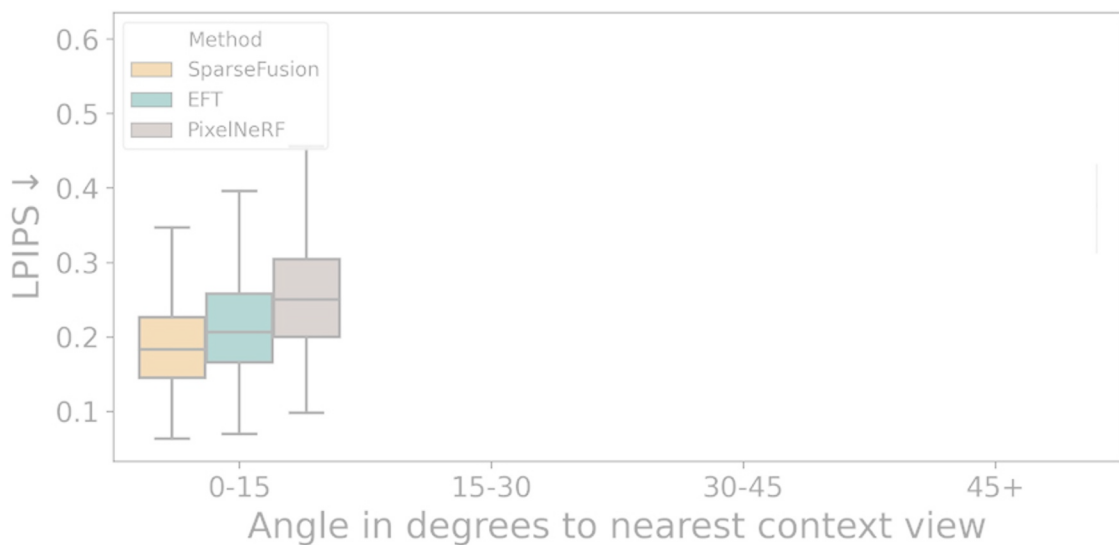
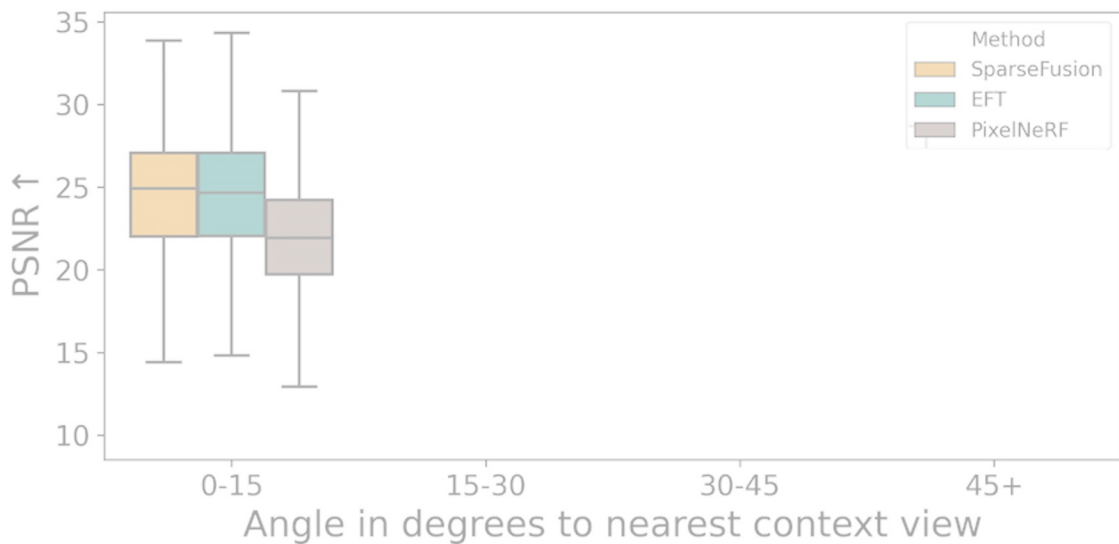
2-view Novel View Synthesis on 10 Categories

	PSNR \uparrow	LPIPS \downarrow
PixelNeRF	19.52	0.327
NerFormer	17.88	0.382
ViewFormer	18.37	0.282
EFT	20.85	0.289
VLDM	19.55	0.247
SparseFusion	21.34	0.225

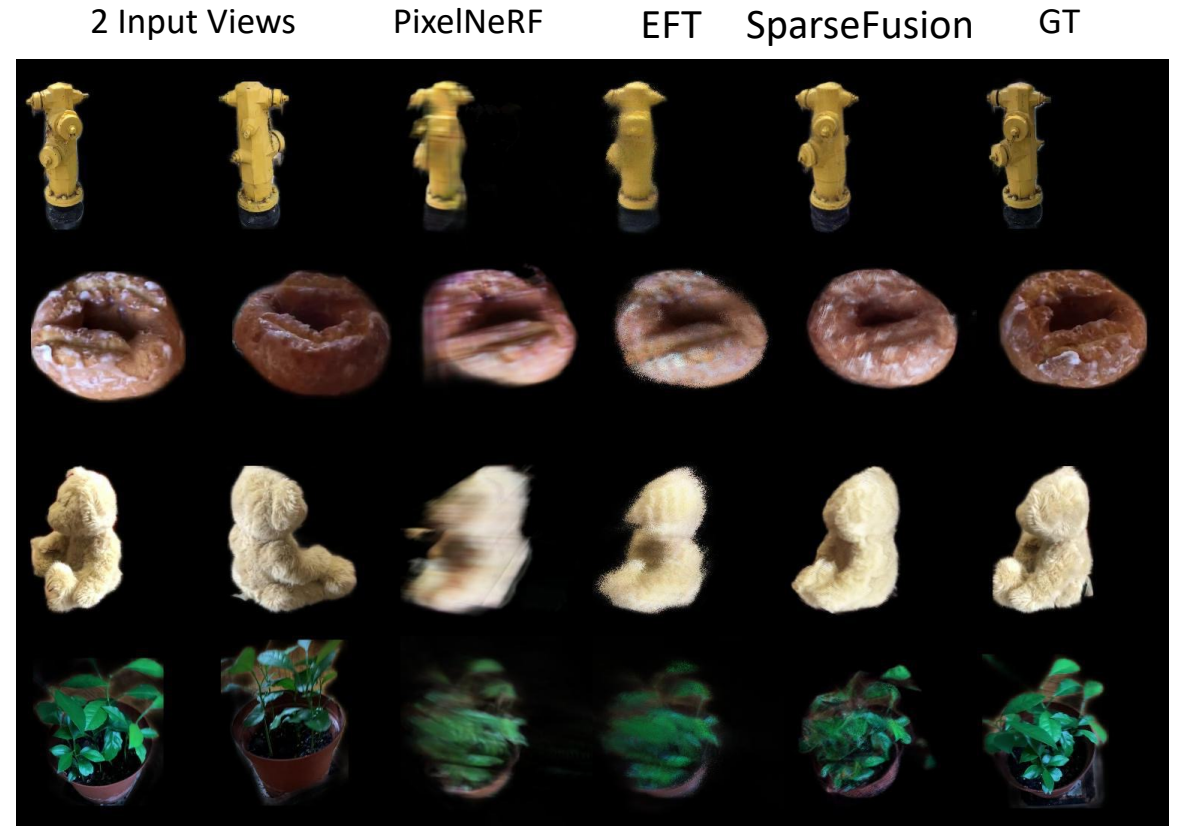
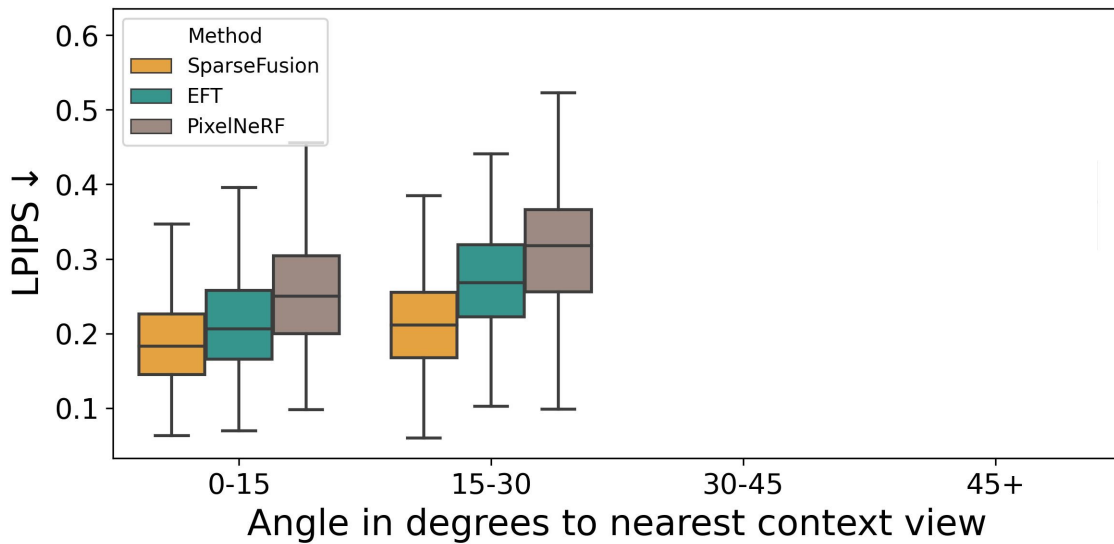
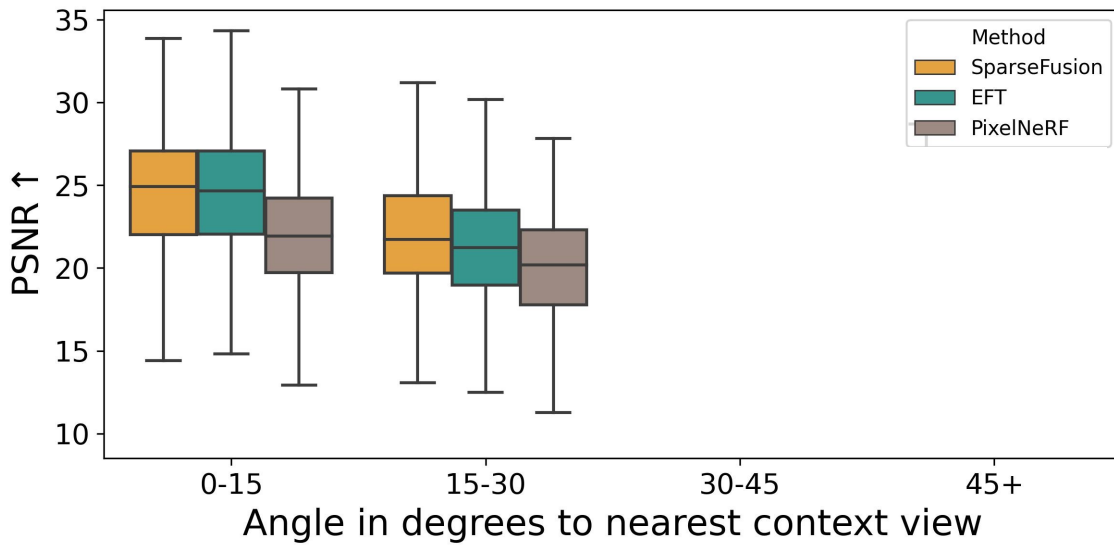
Results



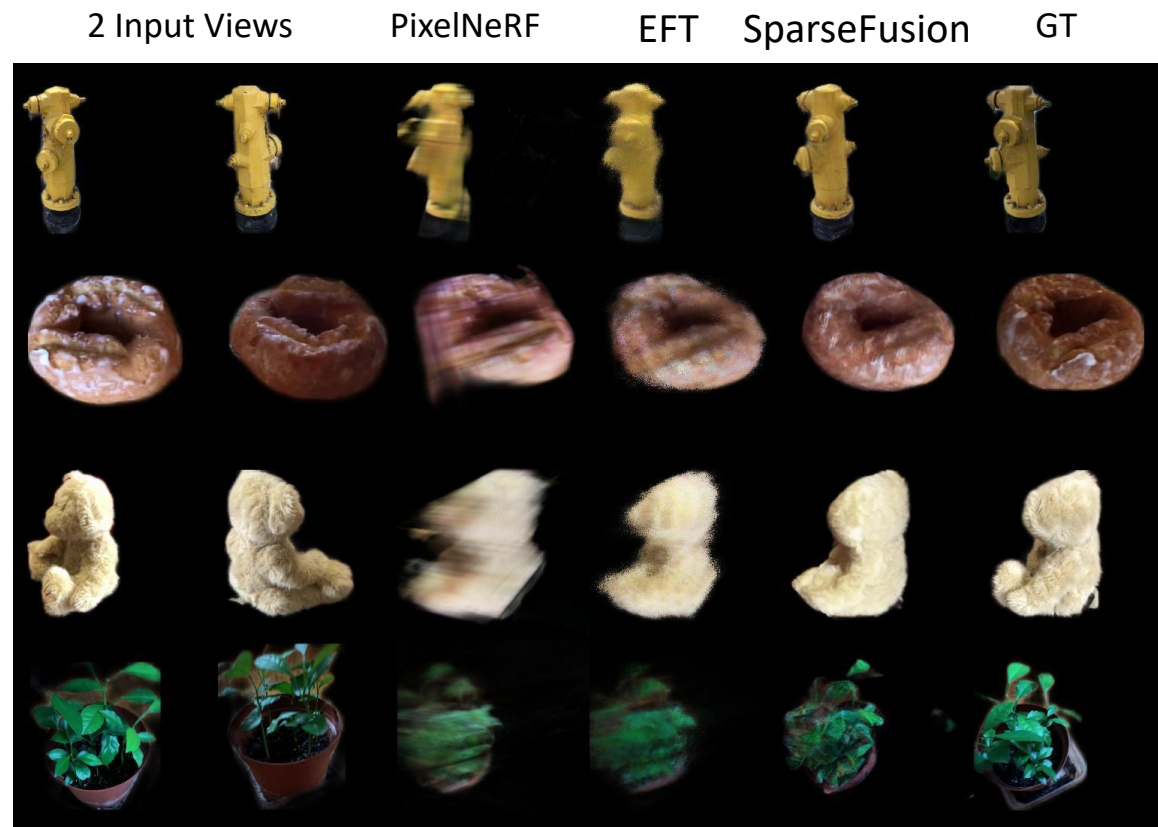
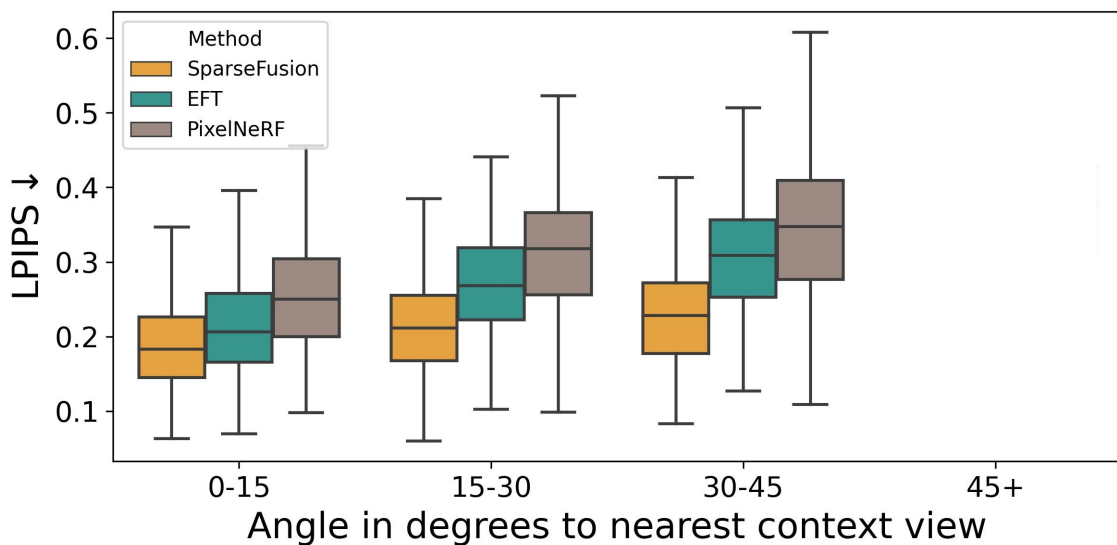
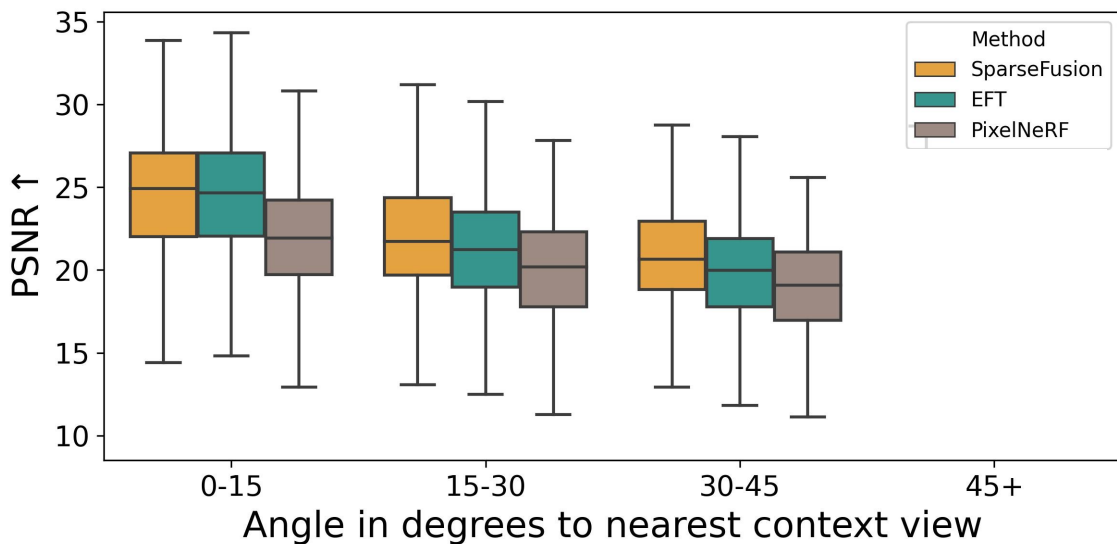
Results



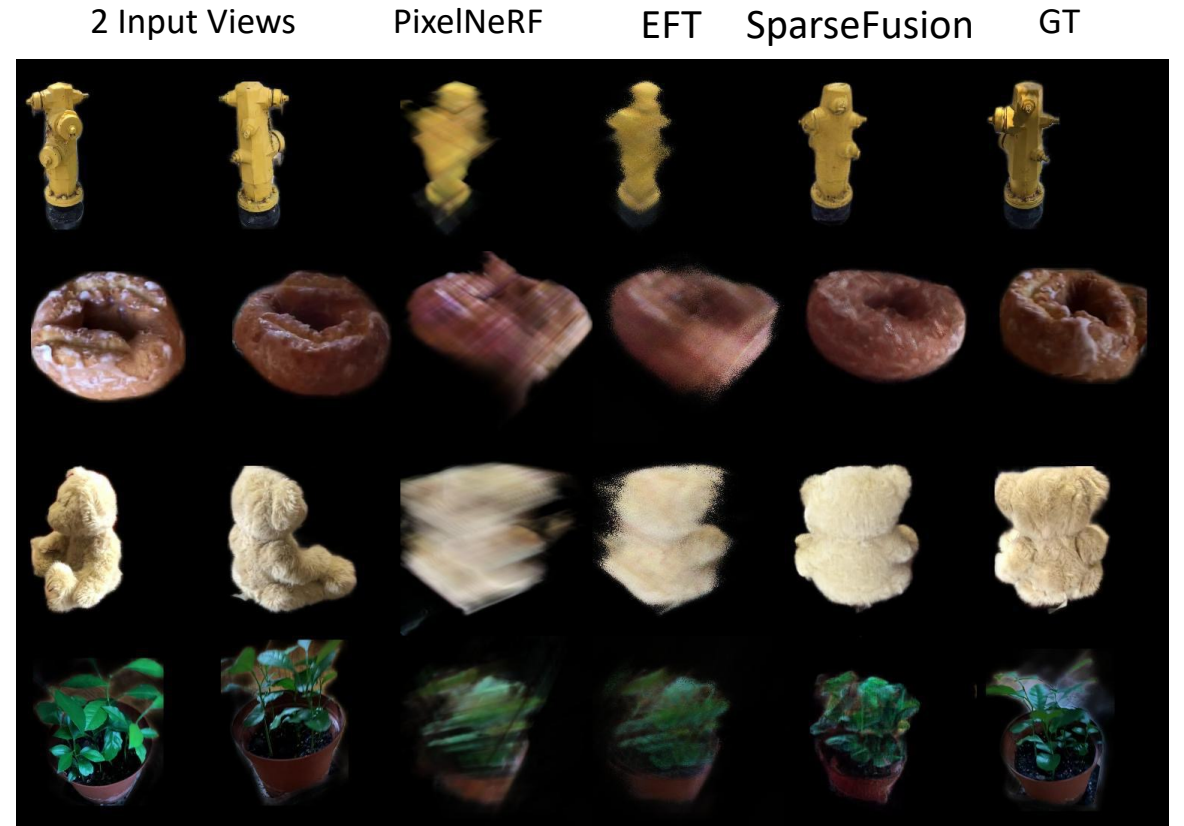
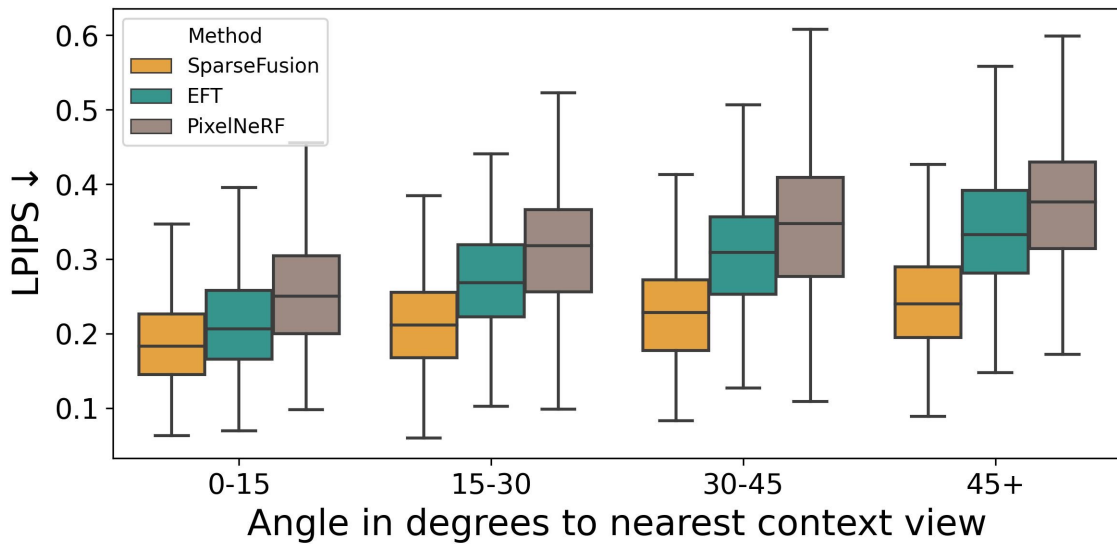
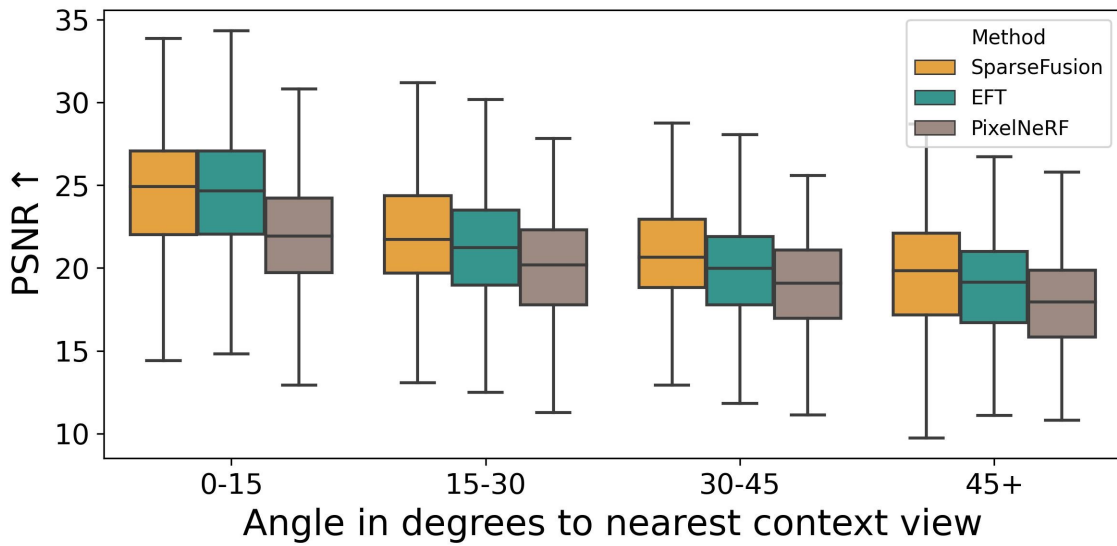
Results



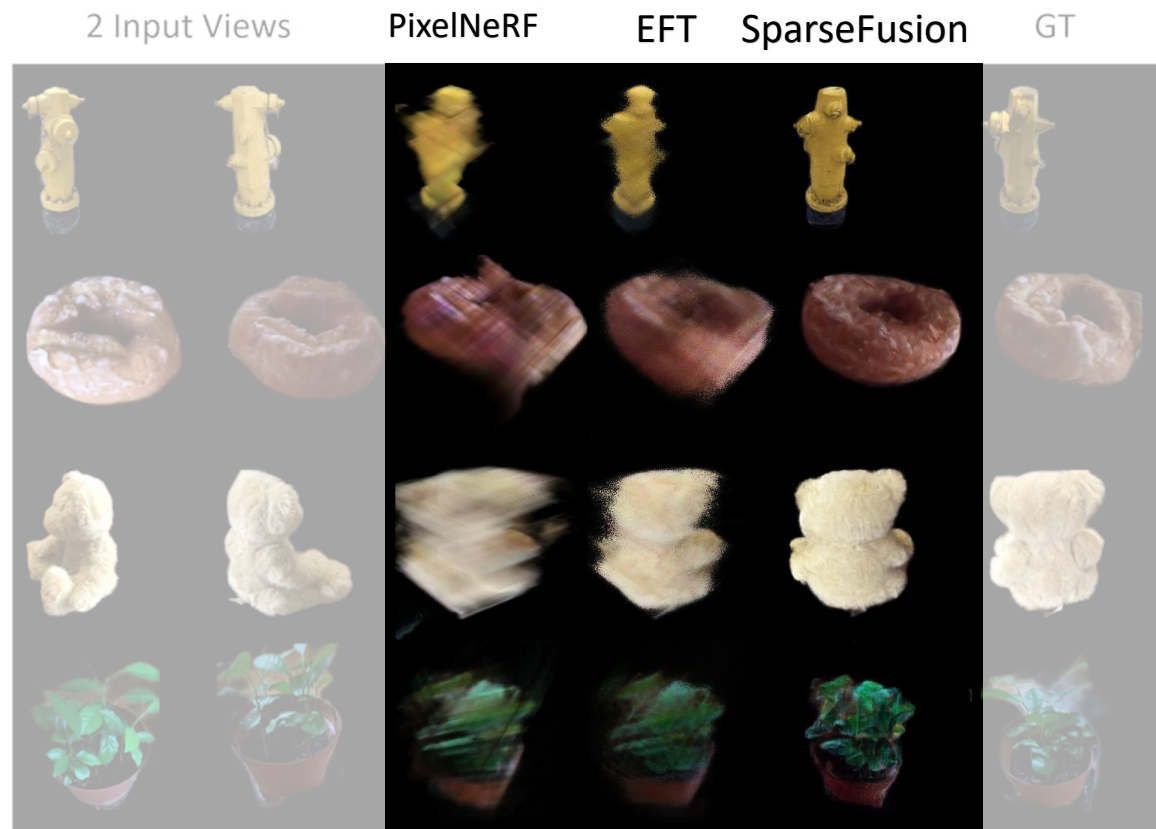
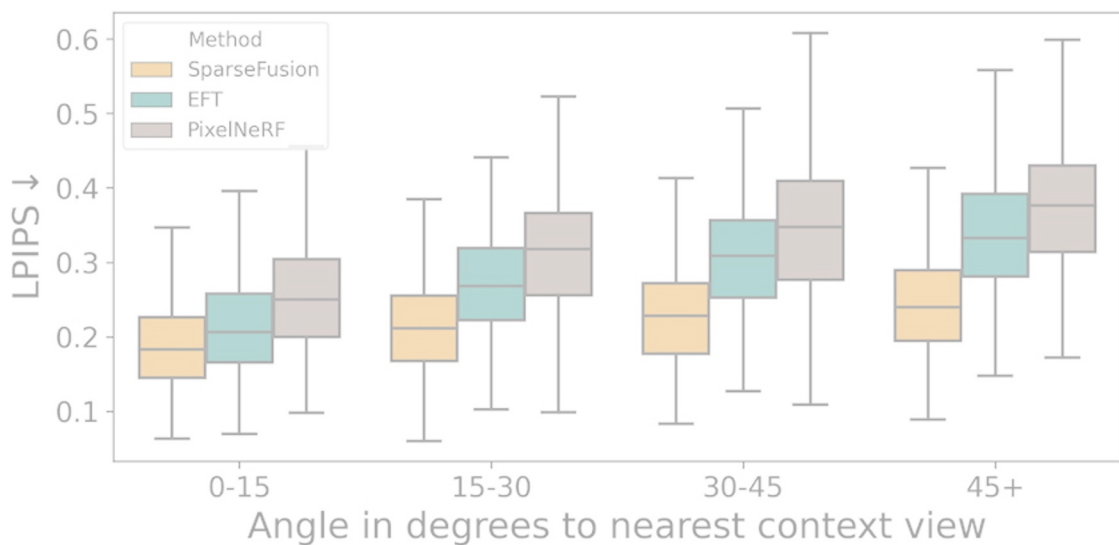
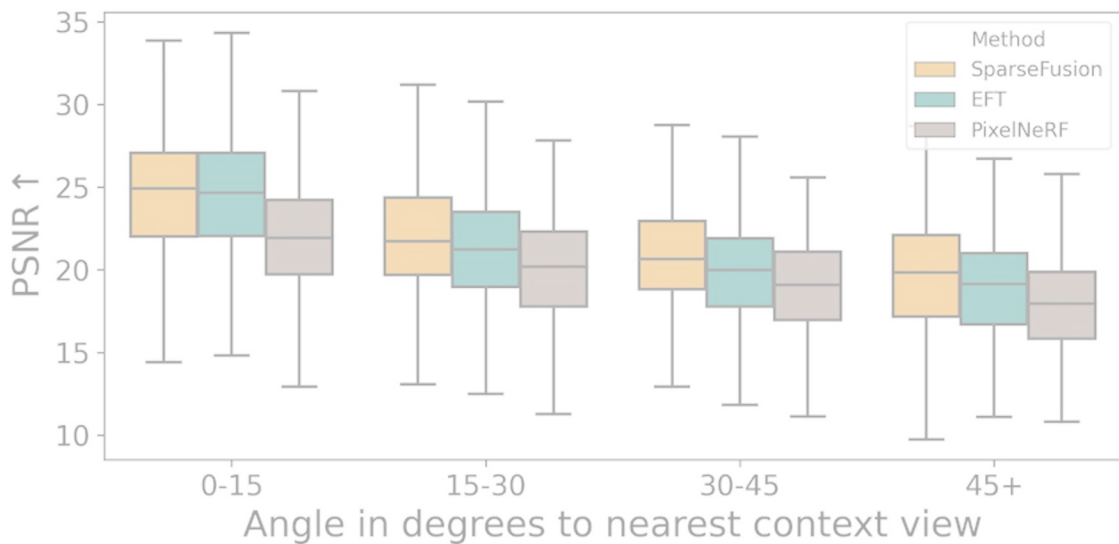
Results



Results



Results



sparsefusion.github.io



Zhizhuo (Z) Zhou, Shubham Tulsiani