

Awesome Conditional 3D Generation

Tianhang Cheng

Update: 2023.11.29

Category

<https://github.com/yyeboah/Awesome-Text-to-3D>

A. Input

Text, Image, Mesh, DMTeT

<https://github.com/topics/text-to-3d>

B. Output

3D model, PBR Material, styled output, segmentation

C. Representation

Point, 3D Gaussian Splattings, Mesh, NeRF

D. Optimization

Intermediate representation: multiview image/normal

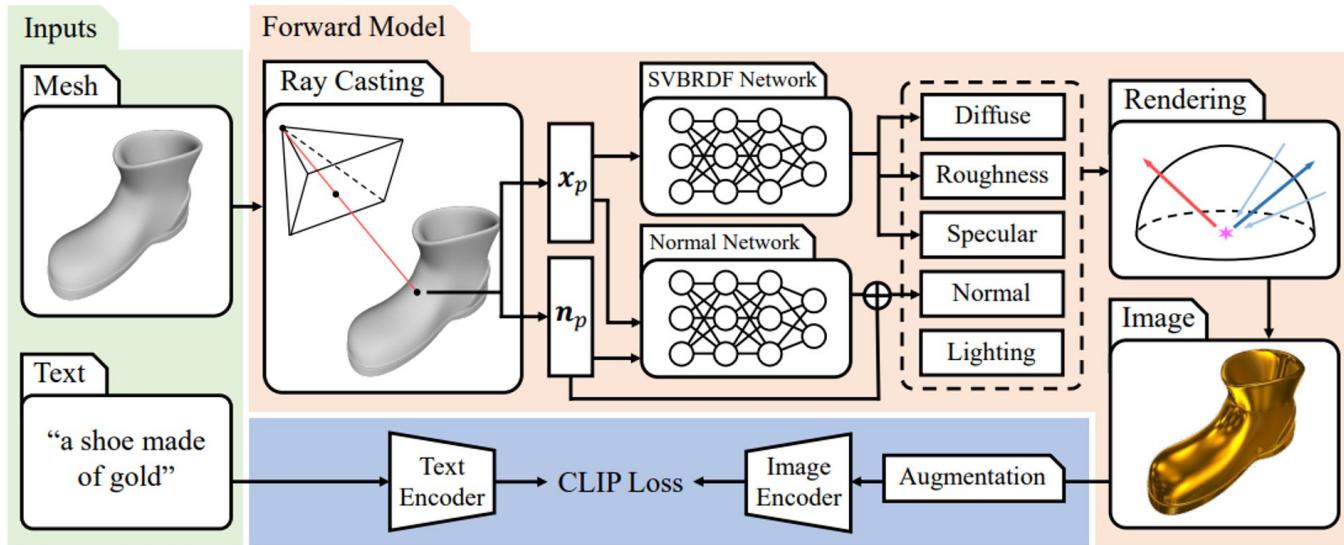
Diffusion: SDS, VSD, SJC, NFSD, ISM, etc.

GAN: Descriminator

Forward: Triplane, cost-volume

Other: GTP+Blender, etc.

A. Input: Mesh



B.Output: PBR Material

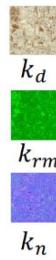
*"a highly detailed stone bust
of Theodoros Kolokotronis"*



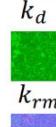
geometry



appearance



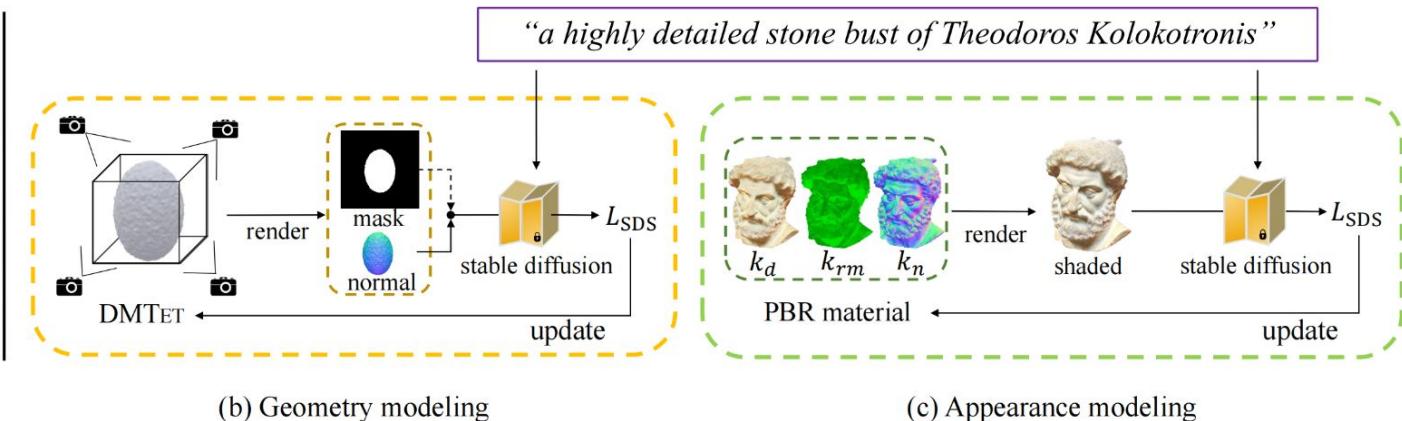
k_d



k_{rm}

k_n

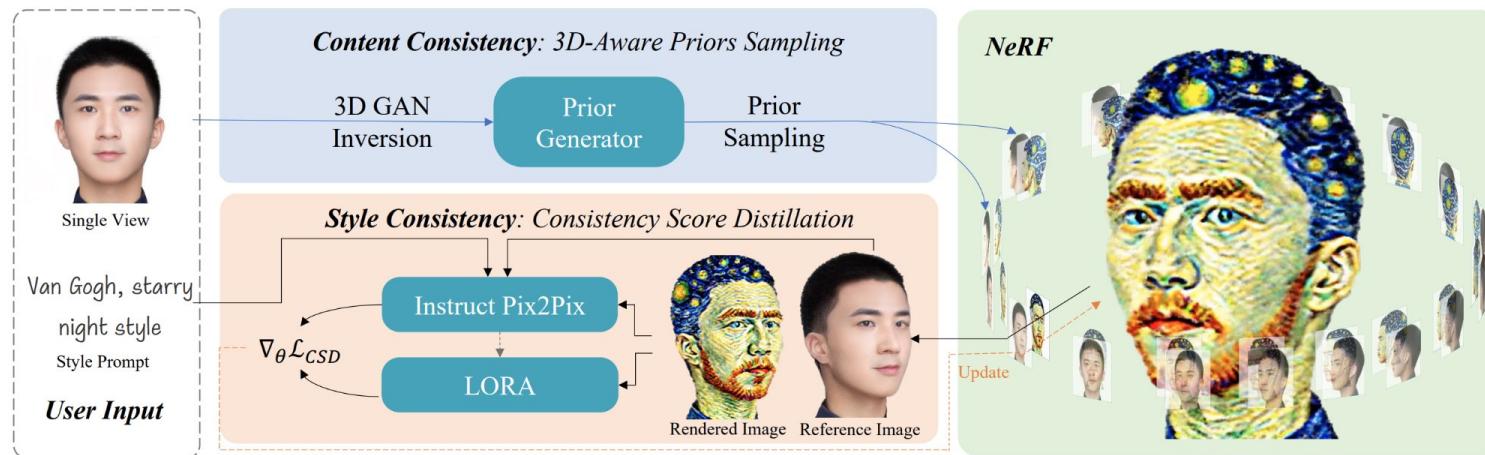
(a) Disentangled representation



(b) Geometry modeling

(c) Appearance modeling

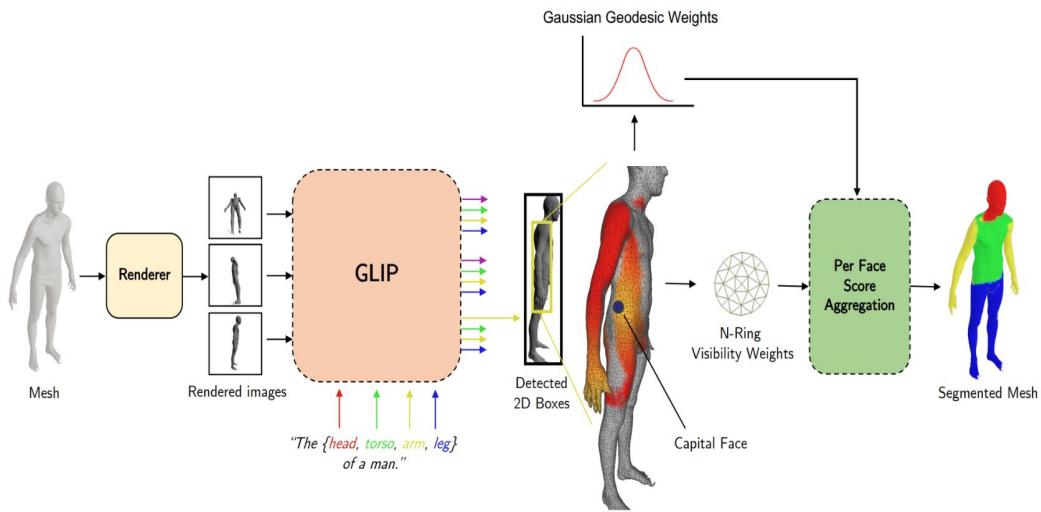
B.Output: styled output



**STYLEDREAMER: MAKE YOUR 3D STYLE AVATAR FROM A SINGLE VIEW WITH
CONSISTENCY SCORE DISTILLATION**

B.Output: segmentation

zero-shot segmentation from mesh



C. Representation: Overview

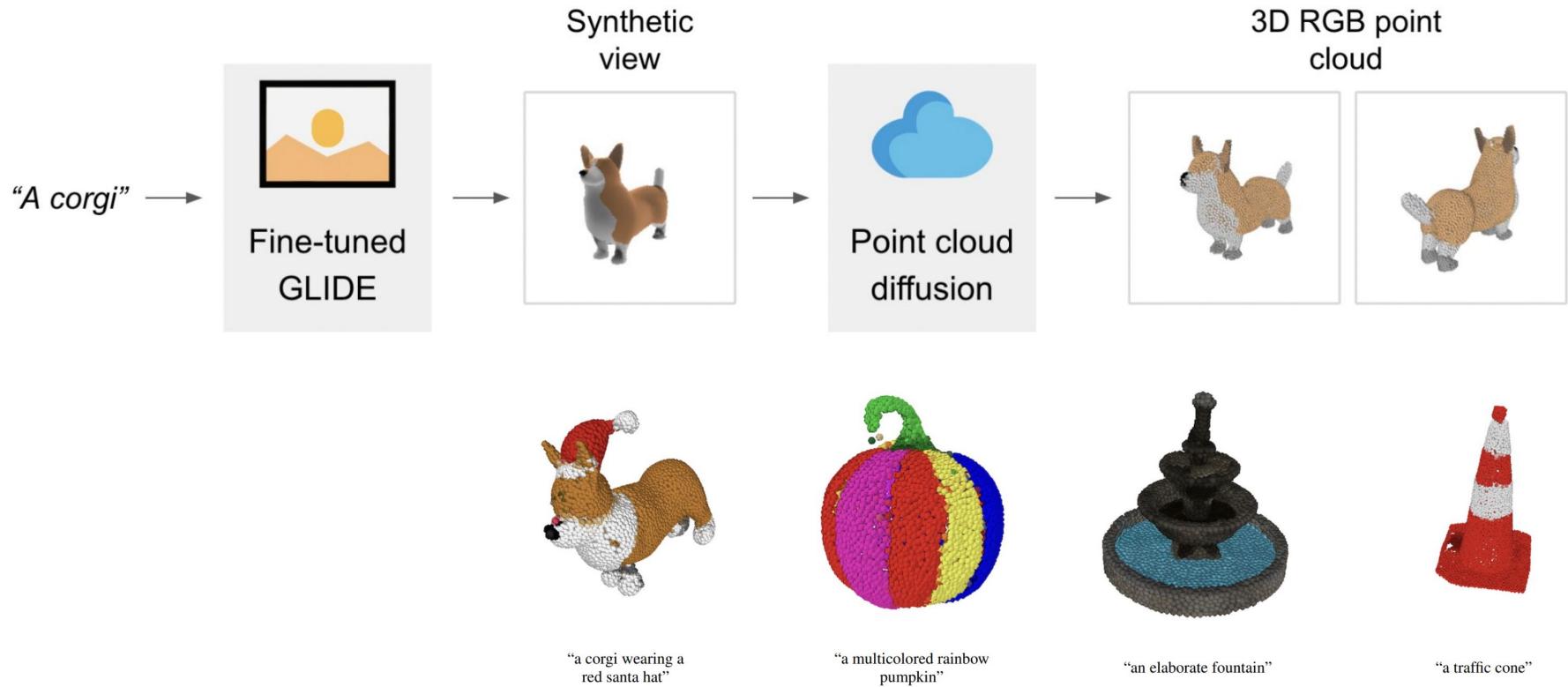
Point, 3D Gaussian Splattings, Mesh, NeRF, etc.

Table 1: Comparison of text-conditioned 3D generation methods on efficiency and applications.

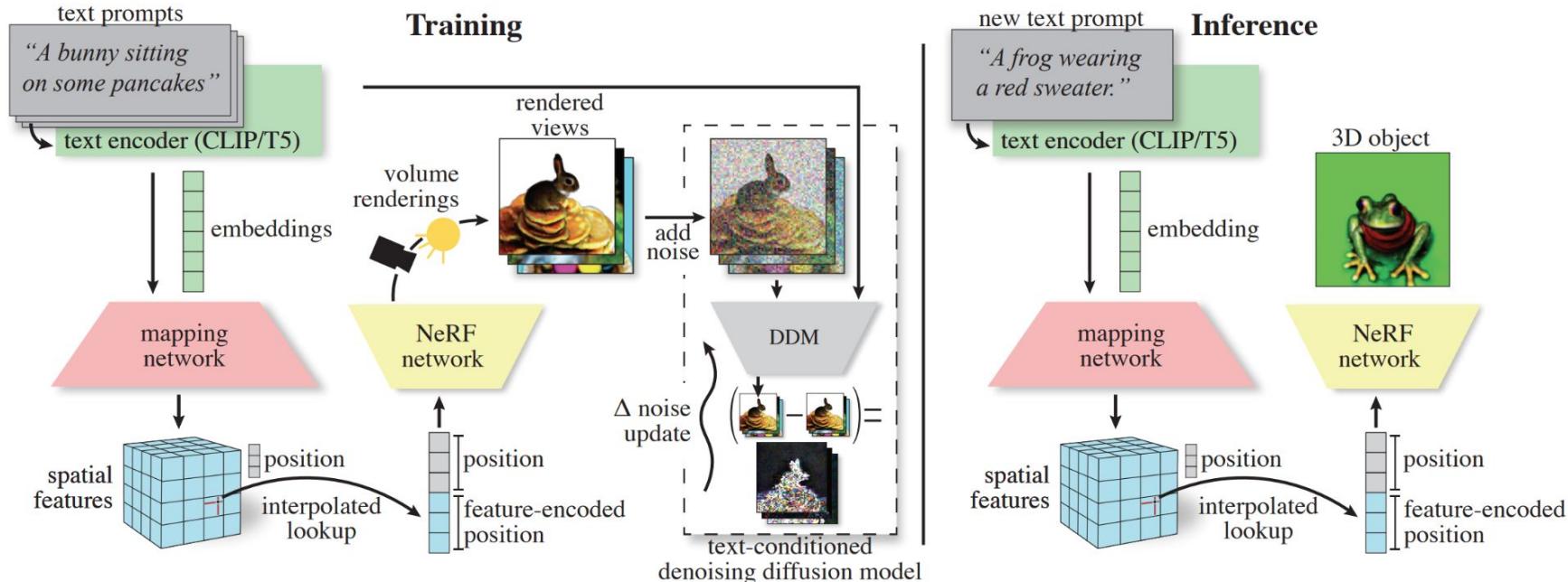
Method	Latency	Device	Generation Category	Downstream Task
DreamFields [37]	1.2h	8xTPUv4	Multi-Category	Generation
DreamFusion [65]	1.5h	4xTPUv4	Multi-Category	Generation
CLIP-Mesh [58]	30min	V100	Multi-Category	Generation
CLIP-Sculptor [77]	0.9s	V100	Multi-Category	Generation
Point-E (40M) [59]	25s	V100	Multi-Category	Generation
ShapeCrafter [24]	-	-	Single-Category	Generation, Editing
LION [99]	27s	V100	Single-Category	Generation, Completion
SDFusion [9]	-	-	Single-Category	Generation, Completion

figure from VPP : *Efficient Conditional 3D Generation via Voxel-Point Progressive Representation*

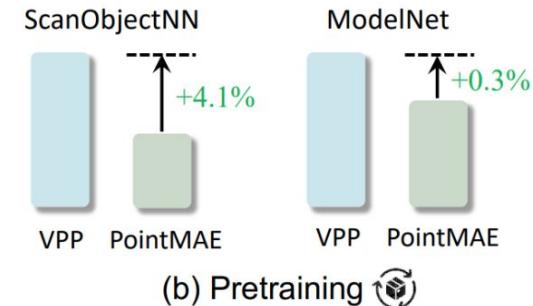
C. Representation: Point



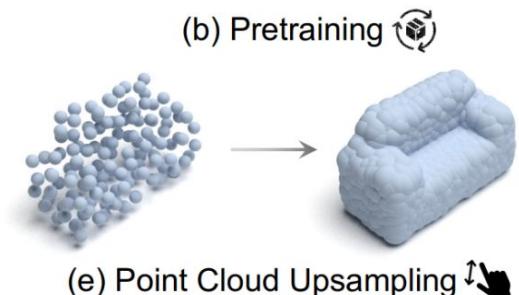
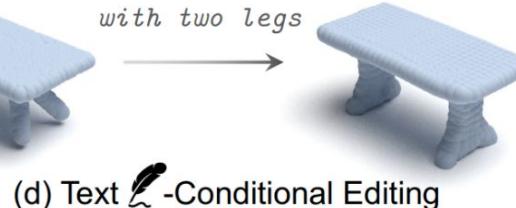
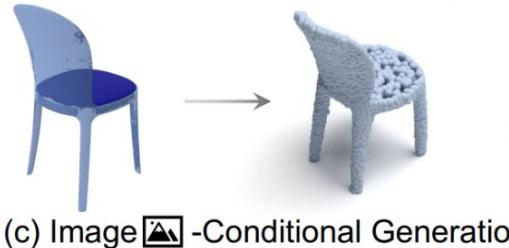
C. Representation: Voxel NeRF



C. Representation: Voxel & Point

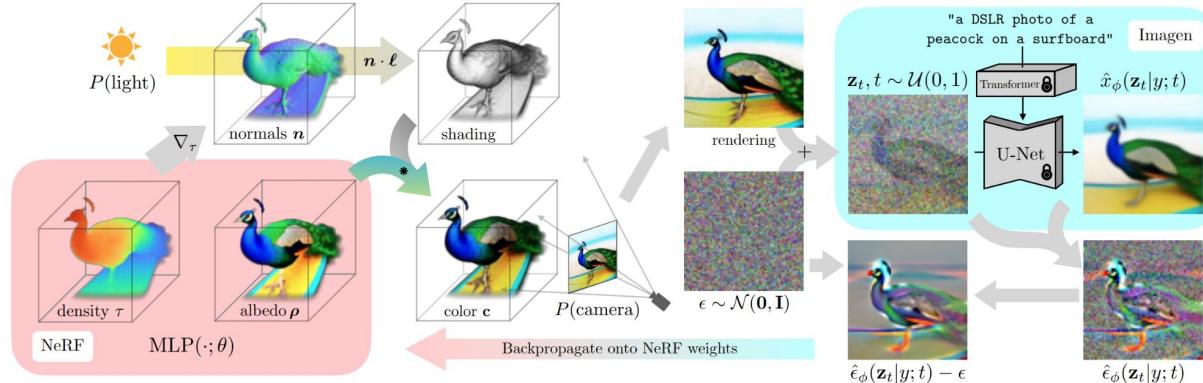


(a) Text -Conditional Generation

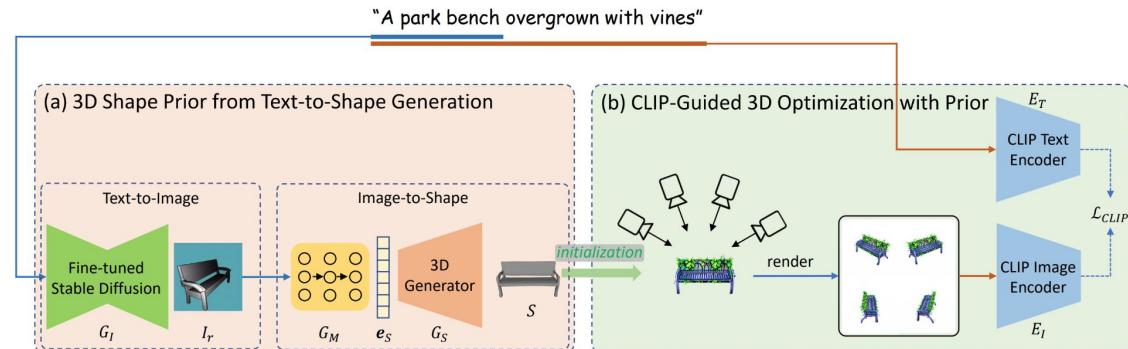


C. Representation: NeRF

NeRF (volume render)



Dream3D: Zero-Shot Text-to-3D Synthesis Using 3D Shape Prior and Text-to-Image Diffusion Models



C.Representation: NeRF

NeRF (volume render) + instant nfp

Magic3D: High-Resolution Text-to-3D Content Creation

Chen-Hsuan Lin* Jun Gao* Luming Tang* Towaki Takikawa* Xiaohui Zeng* Xun Huang Karsten Kreis Sanja Fidler† Ming-Yu Liu† Tsung-Yi Lin

*†: equal contributions

NVIDIA
JUNE 18-22, 2023
CVPR
VANCOUVER, CANADA

Overview

Magic3D is a new text-to-3D content creation tool that generates high-quality 3D meshes.

We provide users with new ways to control 3D synthesis, opening up new avenues to various creative applications.

Major benefits

- Fast (40 minutes, 2x faster than DreamFusion^[1])
- Uses high-resolution (512x512) diffusion priors

Potential applications

- Movie/game production
- Metaverse
- Robot/self-driving simulation
- 3D designs

Our goal

- Turbocharge expert 3D artists
- Facilitate 3D content creation for novices

Approach

We optimize the 3D content over pretrained text-to-image diffusion models with SDS (Score Distillation Sampling) loss in a coarse-to-fine fashion, with the following procedure:

1. Use a **low-resolution** image diffusion prior and optimize a neural field representation to obtain the coarse model.
2. Extract a textured 3D mesh from the density/color neural fields.
3. Fine-tune the mesh with a **high-resolution** latent diffusion model.

Results

[1] Poole et al. "DreamFusion: Text-to-3D using 2D Diffusion." ICLR 2023.

Magic3D: High-Resolution Text-to-3D Content Creation

C. Representation: DMTeT

*"a highly detailed stone bust
of Theodoros Kolokotronis"*



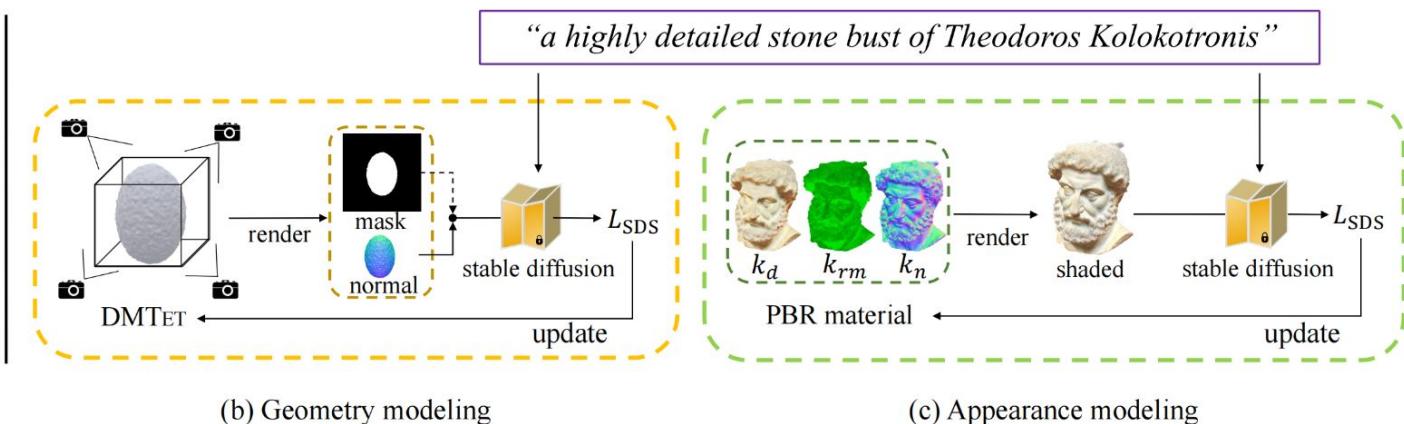
geometry



appearance



(a) Disentangled representation

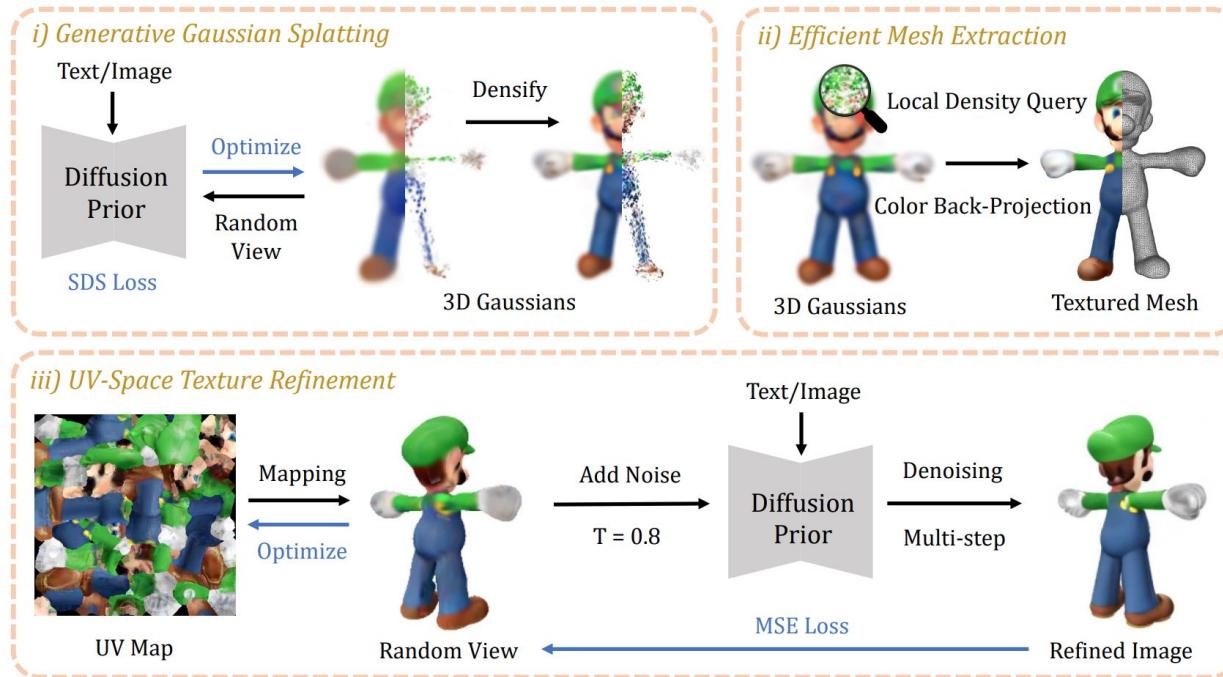


(b) Geometry modeling

(c) Appearance modeling

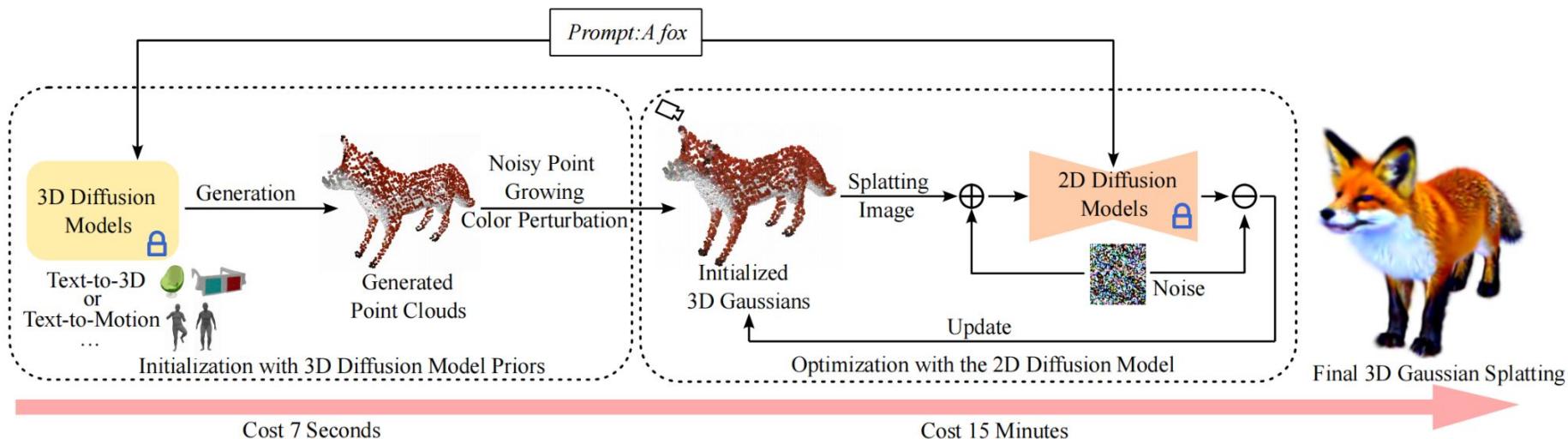
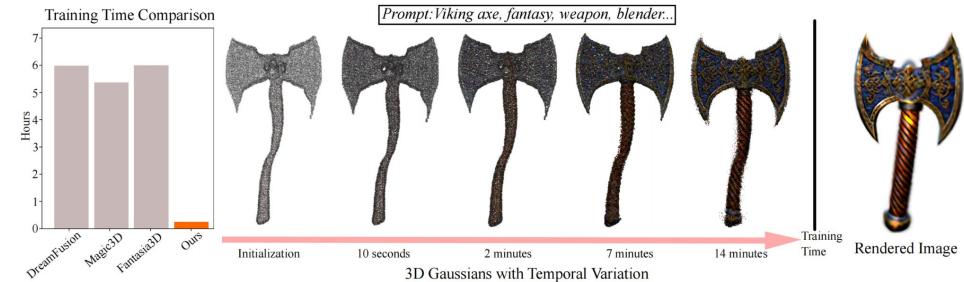
C.Representation: Mesh

3dgs as mesh initilization



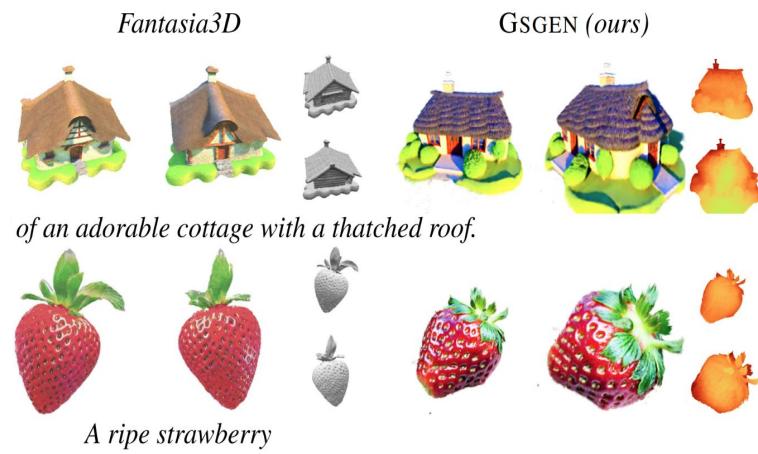
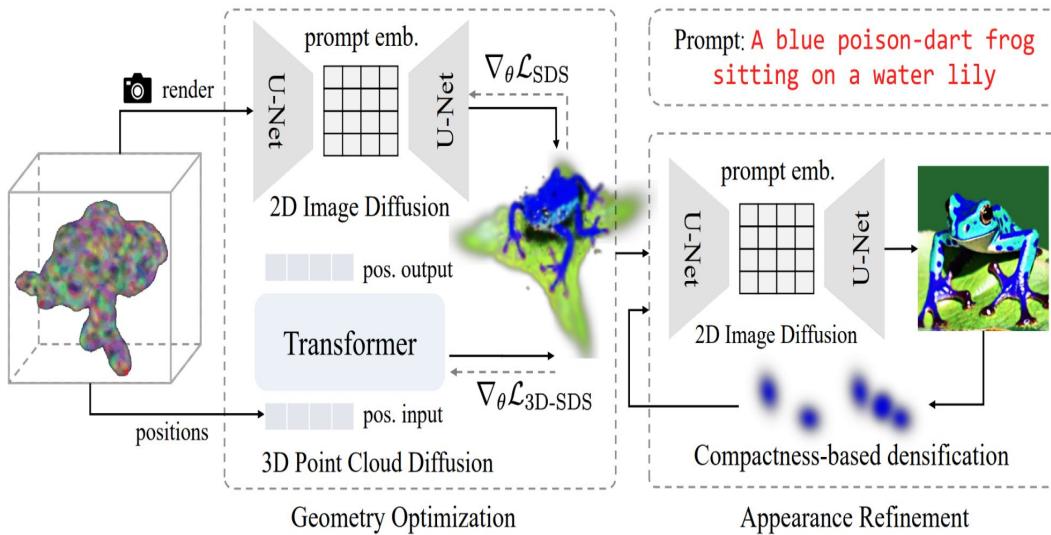
C. Representation: 3DGS

prompt -> point cloud, optimized by 2D diffusion



C. Representation: 3DGS

3DGS + 3D SDS Loss



D.Optimization: Overview

- Multiview image/surface normal as intermediate representation

Use (and finetune) pretrained generative model to generate multiview information from input; Then transfer multiview to 3d model (usually with post-processing if images are not perfect, like 3d conv)

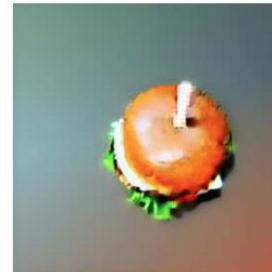
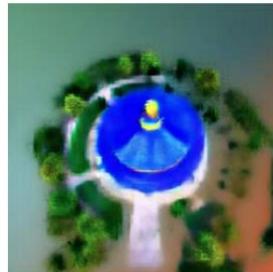
- SDS, VSD, SJC, NFSD, ISM, GAN, etc.

Use pretrained generative model to implicitly modify the probability / adversarial training

- Forward model (reference is fast)

Generate triplane, cost-volume, 3D U-Net

D.Optimization: Diffusion: SJC



A high quality photo of a Victorian style wooden chair with velvet upholstery

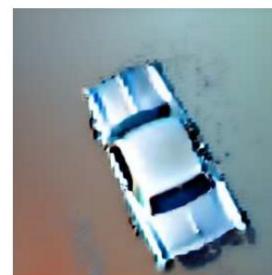
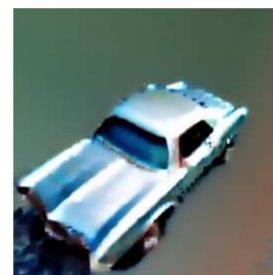
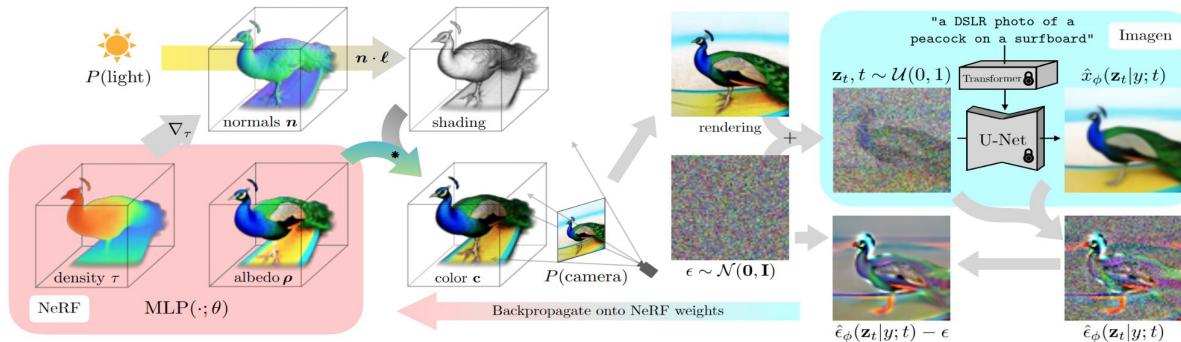
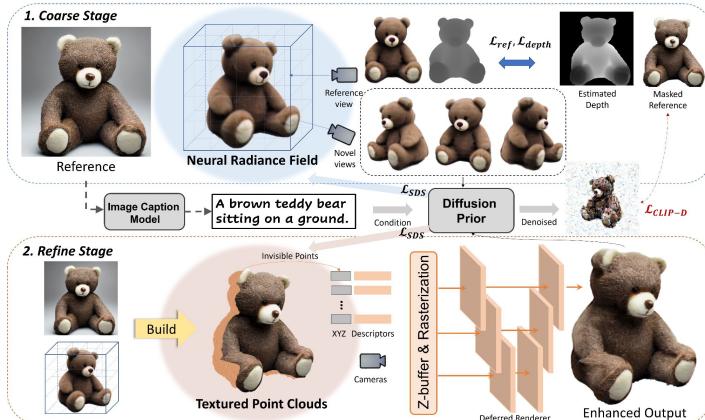


Figure 1. Results for text-driven 3D generation using Score Jacobian Chaining with Stable Diffusion as the pretrained model.

D.Optimization: Diffusion: SDS



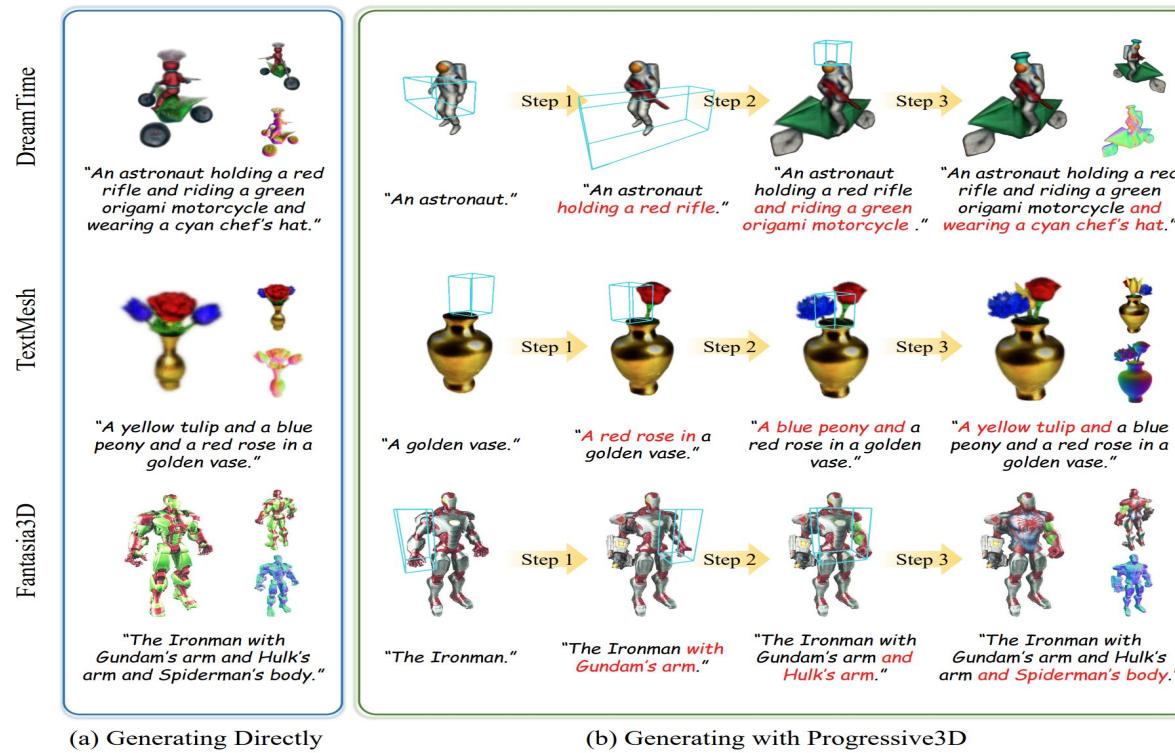
**DREAMFUSION: TEXT-TO-3D
USING 2D DIFFUSION**



depth, clip, sds as supervision.
Not perform very well on back face

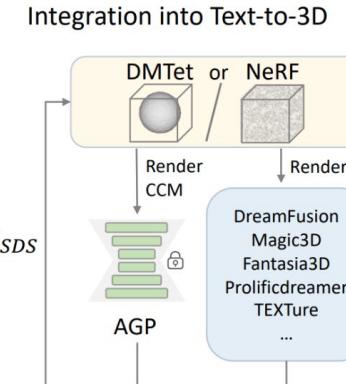
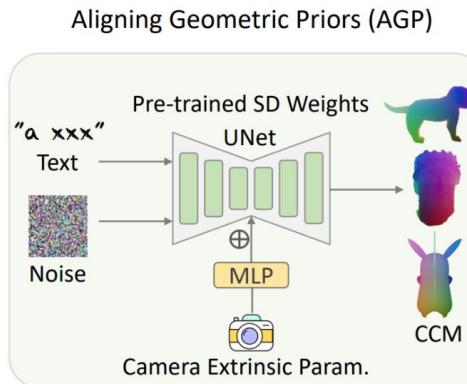
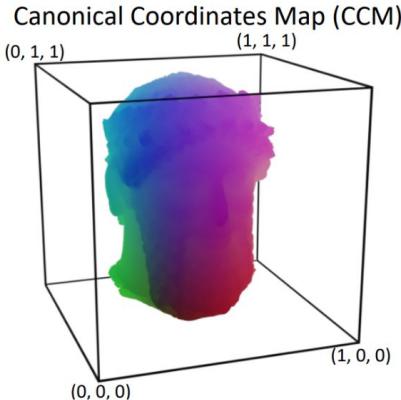
Make-It-3D: High-Fidelity 3D Creation from A Single Image with Diffusion Prior

D.Optimization: Diffusion: (progressive) SDS



PROGRESSIVE3D: PROGRESSIVELY LOCAL EDITING FOR TEXT-TO-3D CONTENT CREATION WITH COMPLEX SEMANTIC PROMPTS

D.Optimization: Diffusion: (aligned) SDS



align 2D diffusion prior by output CCM

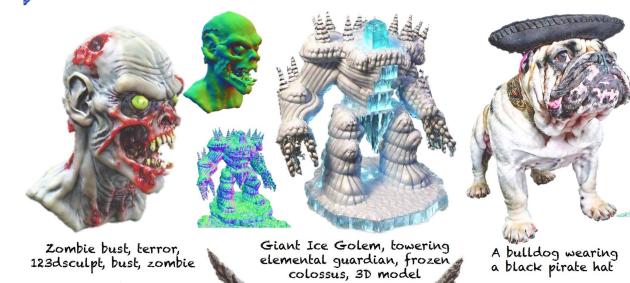
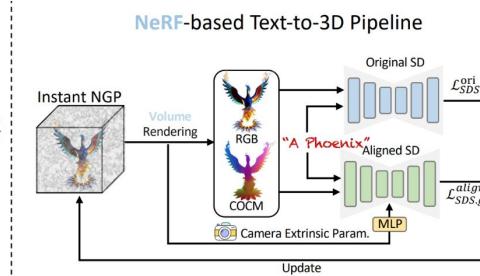
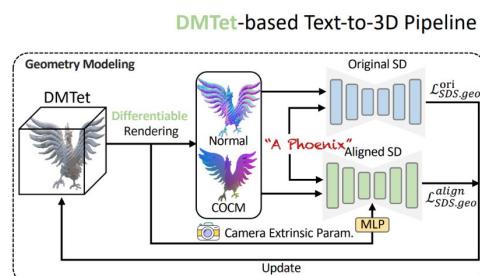
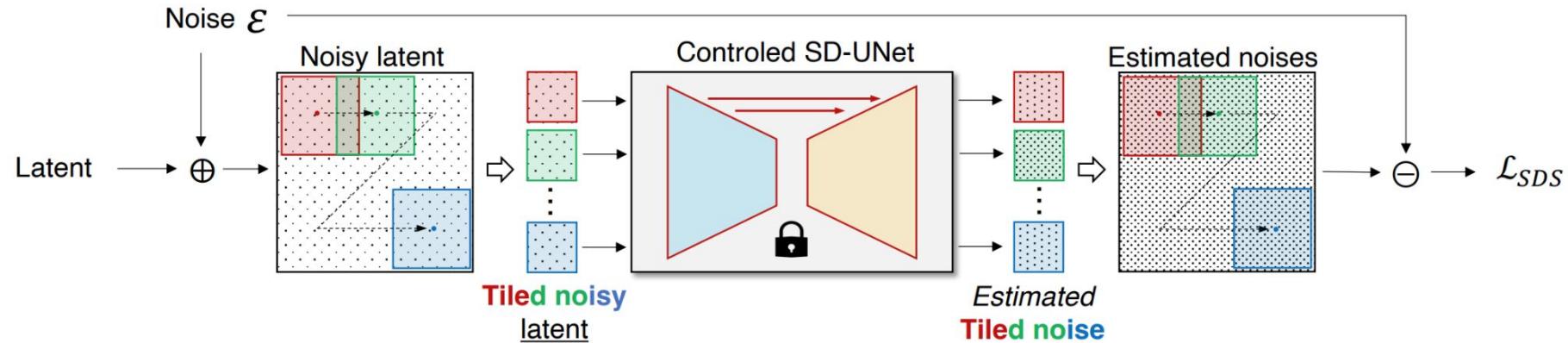


Figure 3: Seamless integration of our AGP in various text-to-3D pipelines.

D.Optimization: Diffusion: (multi-patch) SDS



multi-patch reduce unaligned problem of 2d prior

like the relation between *Triplane* vs *Sparse Tri-Vector Radiance Fields*



D.Optimization: Diffusion: (constrain) SDS

constrain the semantic, geometry and saturation inconsistency in SDS

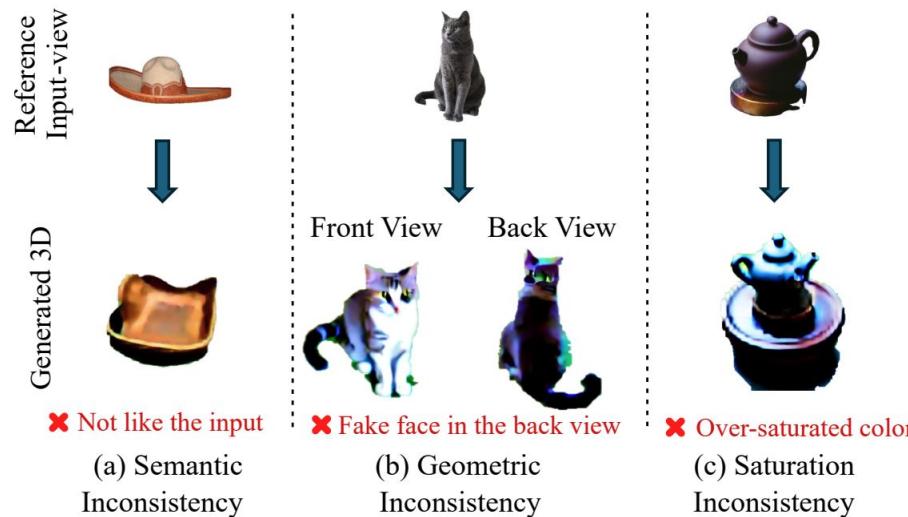
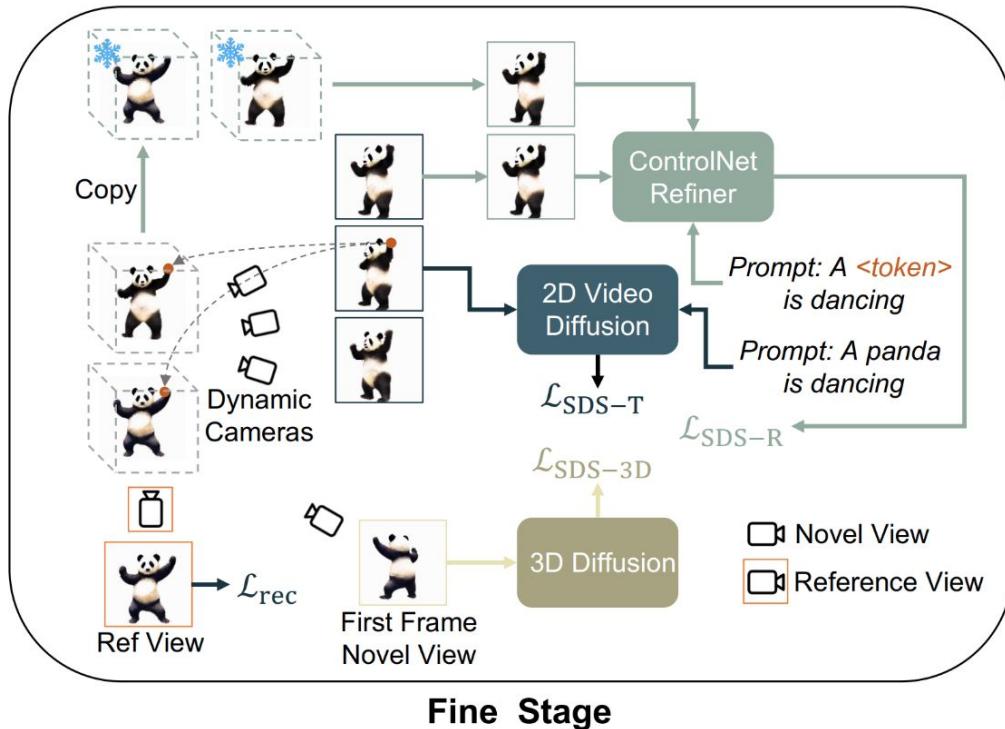
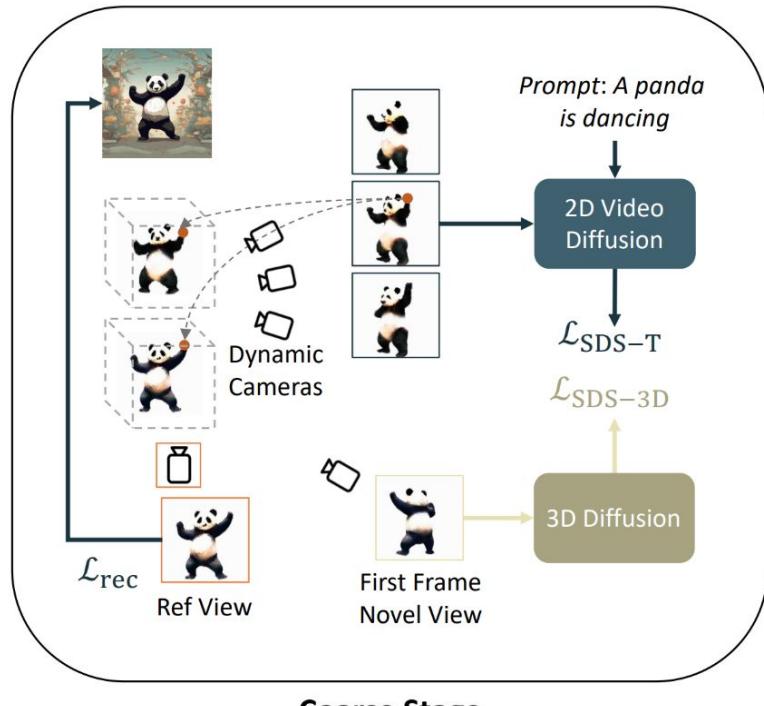


Figure 1: **Inconsistency issues.** (a) The semantic inconsis-

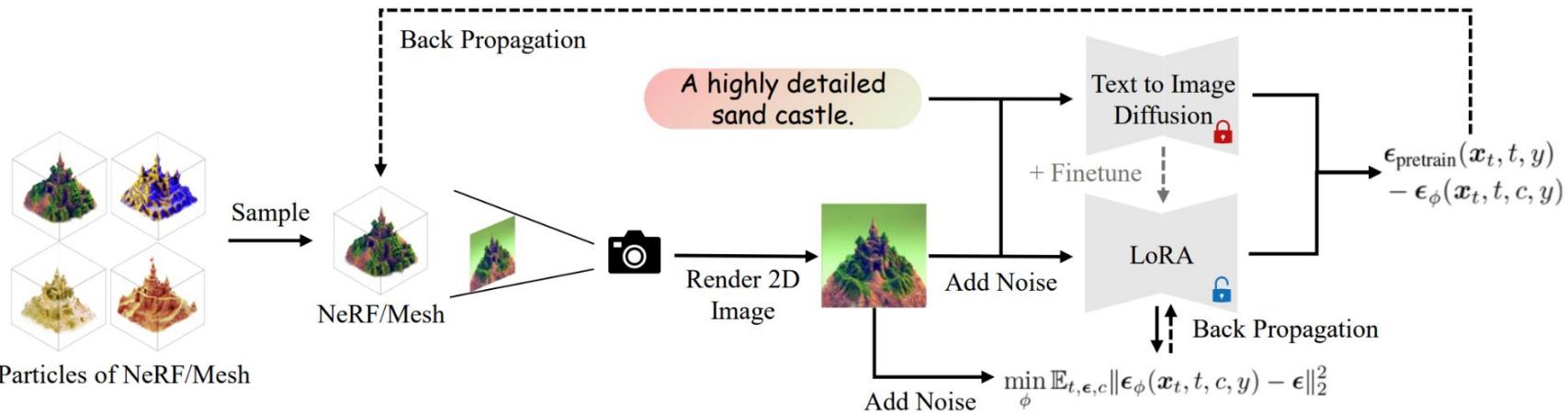
D.Optimization: Diffusion: multiple diffusion

4D



D.Optimization: Diffusion: VDS

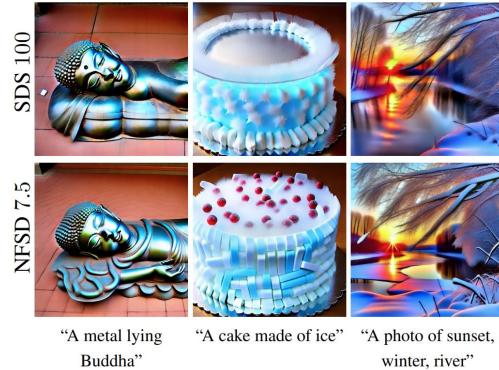
model variant probability by training several nerf “particle” together



D.Optimization: Diffusion: NFSD



Figure 6: NeRFs optimized with NFSD.



“A metal lying Buddha”

“A cake made of ice”

“A photo of sunset, winter, river”

: 2D image generation with SDS and NFSD. We



D.Optimization: Diffusion: ISM

Interval Score Matching



"A portrait of Hatsune Miku, robot"

"A beautiful cyborg with brown hair"

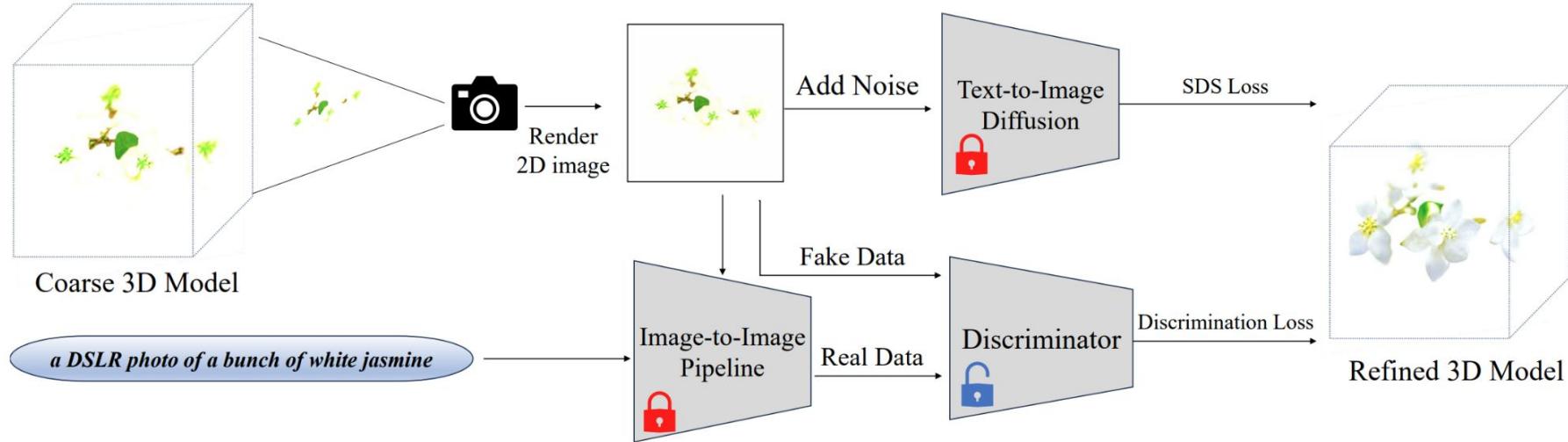
"An armored green-skin orc warrior riding a vicious hog"

"A warrior with red cape riding a horse"

"A forbidden castle high up in the mountains"

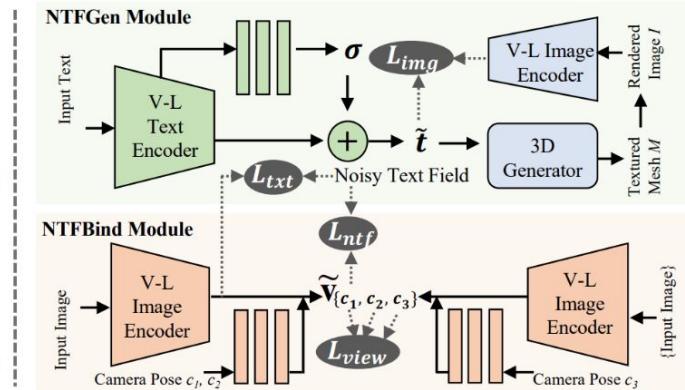
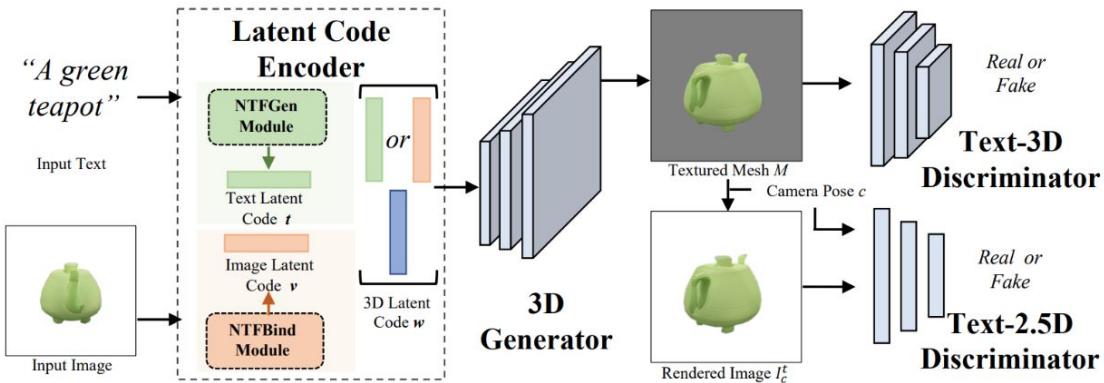
"A highly-detailed sandcastle"

D.Optimization: GAN: Descriminator

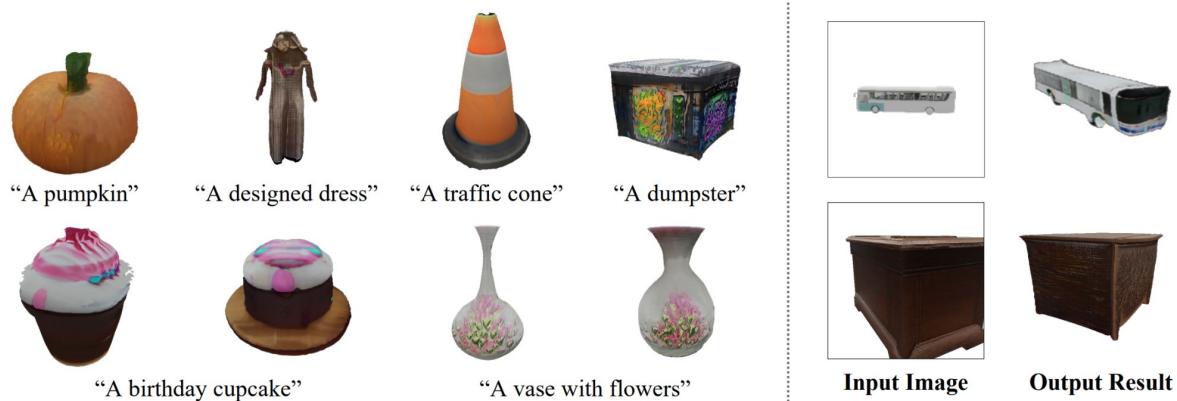


D.Optimization: GAN: (multi-modal) Descriminator

text3D discriminator and a text-2.5D discriminator



*IT3D: Improved
Text-to-3D Generation
with Explicit View
Synthesis*



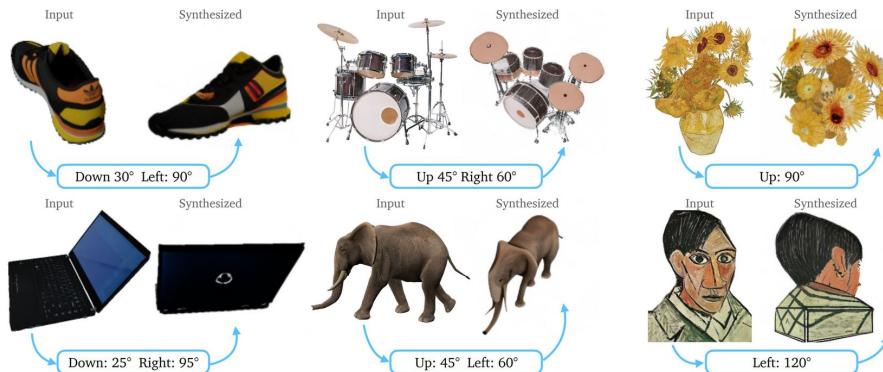
(a) Text-to-3D

(b) Image-to-3D

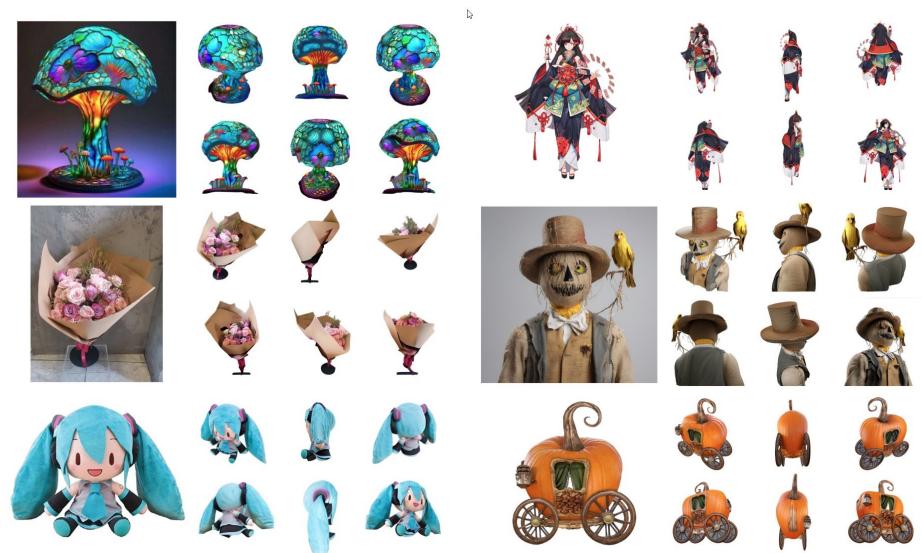
D.Optimization: Intermediate representation: multiview image

finetune 2D diffusion model to get a pose-dependent model

finetune 2D diffusion model to generate more consistent multiview image



Zero-1-to-3: Zero-shot One Image to 3D Object



Zero123++: a Single Image to Consistent Multi-view Diffusion Base Model

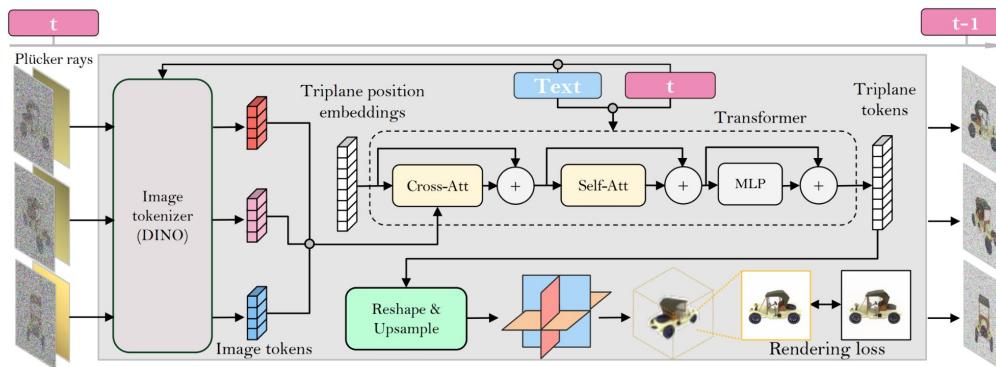
D.Optimization: Intermediate representation: multiview image



finetune 2D diffusion model to get a multiview RGB images

<https://mv-dream.github.io/index.html>

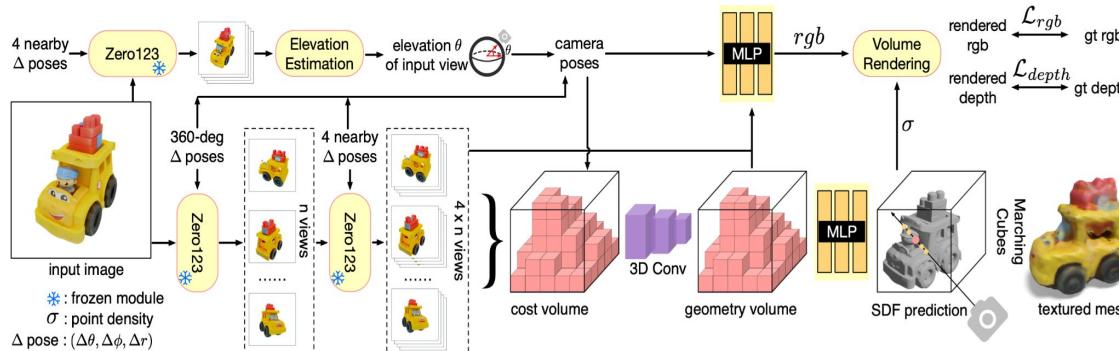
MVDream: Multi-view Diffusion for 3D Generation



DMV3D: DENOISING MULTI-VIEW DIFFUSION USING 3D LARGE RECONSTRUCTION MODEL

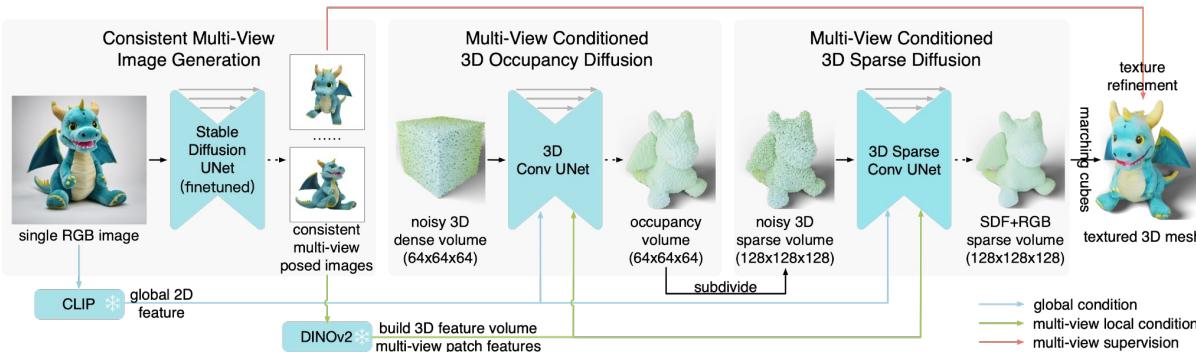
D.Optimization: Intermediate representation: multiview image

add cost-volume based on zero123 to mitigate the in-consistency



One-2-3-45: Any Single Image to 3D Mesh in 45 Seconds without Per-Shape Optimization

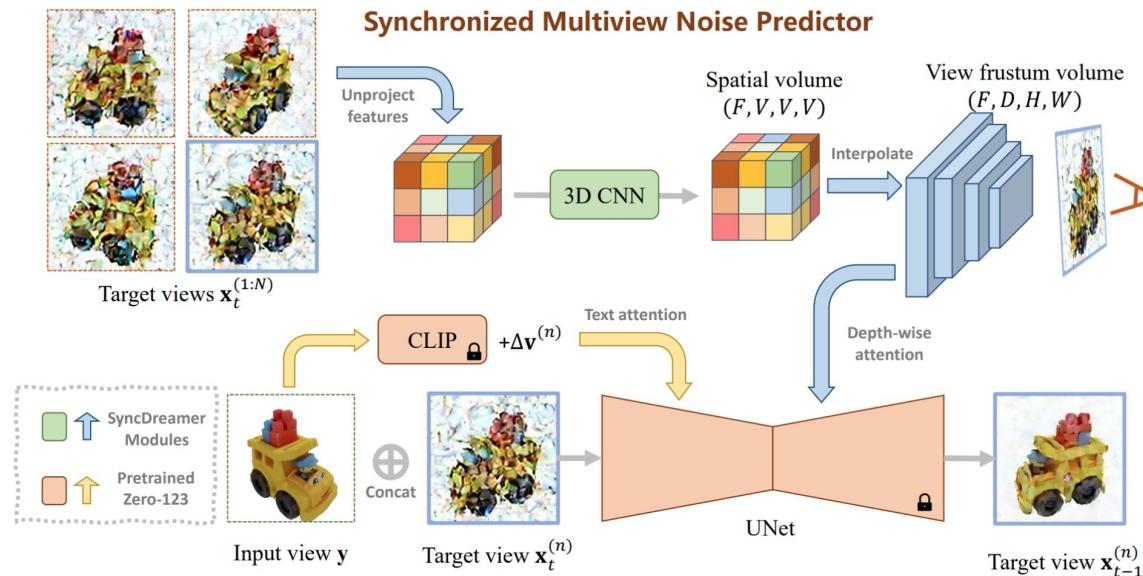
use better diffusion model and feature volume



One-2-3-45++: Fast Single Image to 3D Objects with Consistent Multi-View Generation and 3D Diffusion

D.Optimization: Intermediate representation: multiview image

constrain multi-view feature by unprojected to a volume

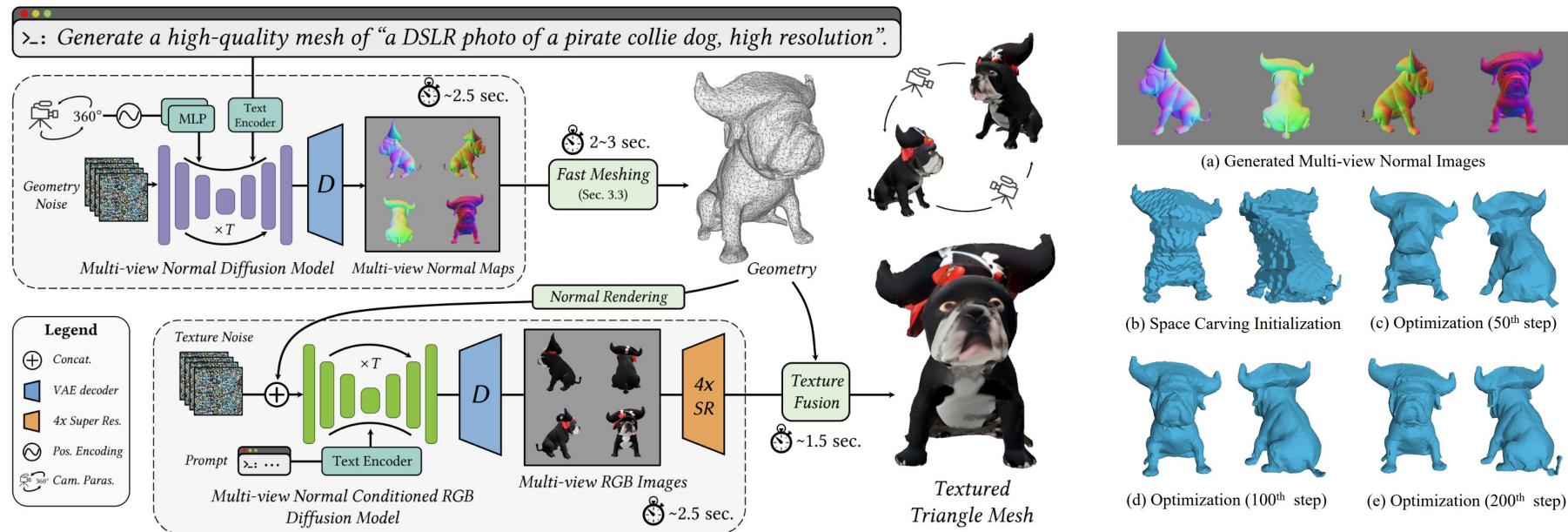


SyncDreamer: Generating Multiview-consistent Images from a Single-view Image

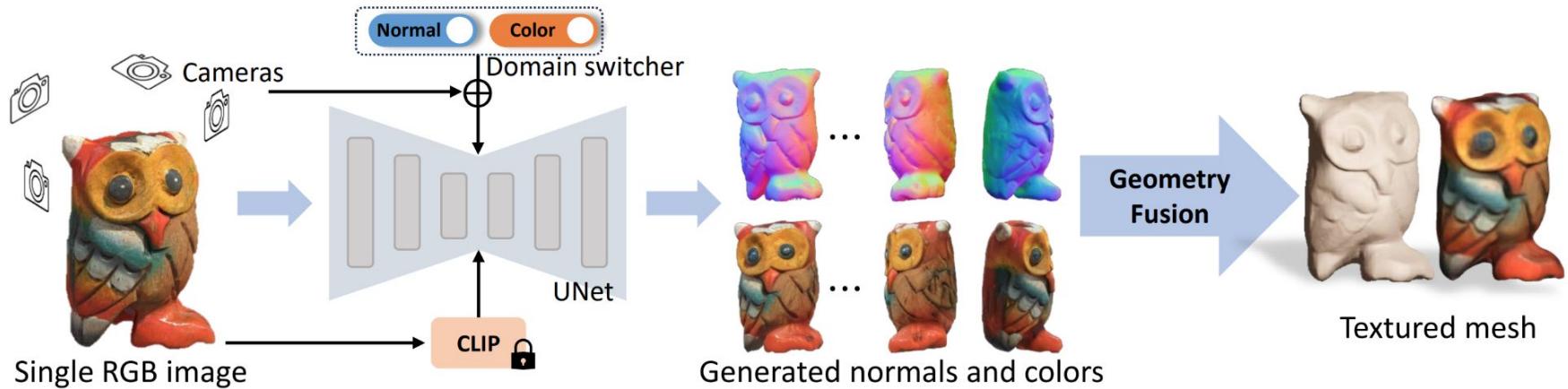
D.Optimization: Intermediate representation: multiview normal

finetune 2D diffusion model to get multiview normal images

TL;DR: Diverse, Janus-free, and high-fidelity 3D content generation in only 10 seconds.

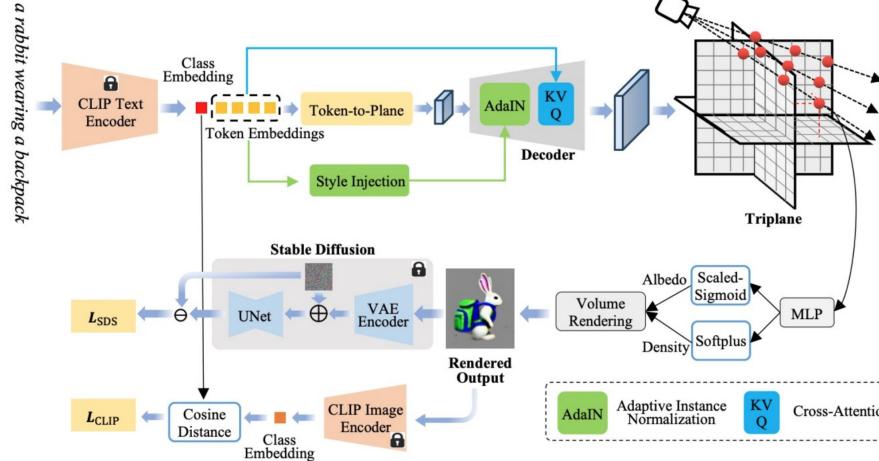


D.Optimization: Intermediate representation: multiview image & normal



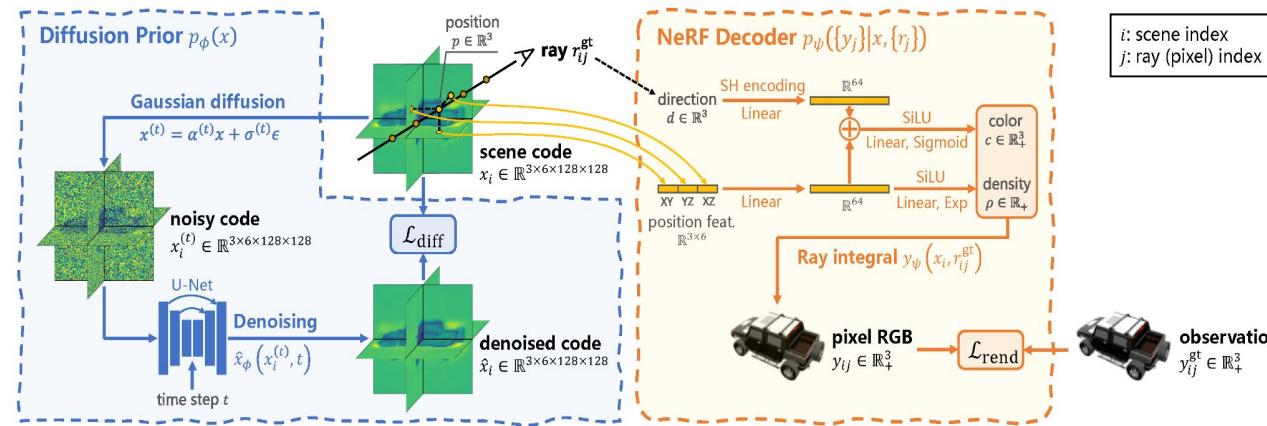
Wonder3D: Single Image to 3D using Cross-Domain Diffusion

D.Optimization: Forward model: triplane



Text-to-3D generation without per-prompt training,
taking only **under a second**.
But the quality is not very realistic

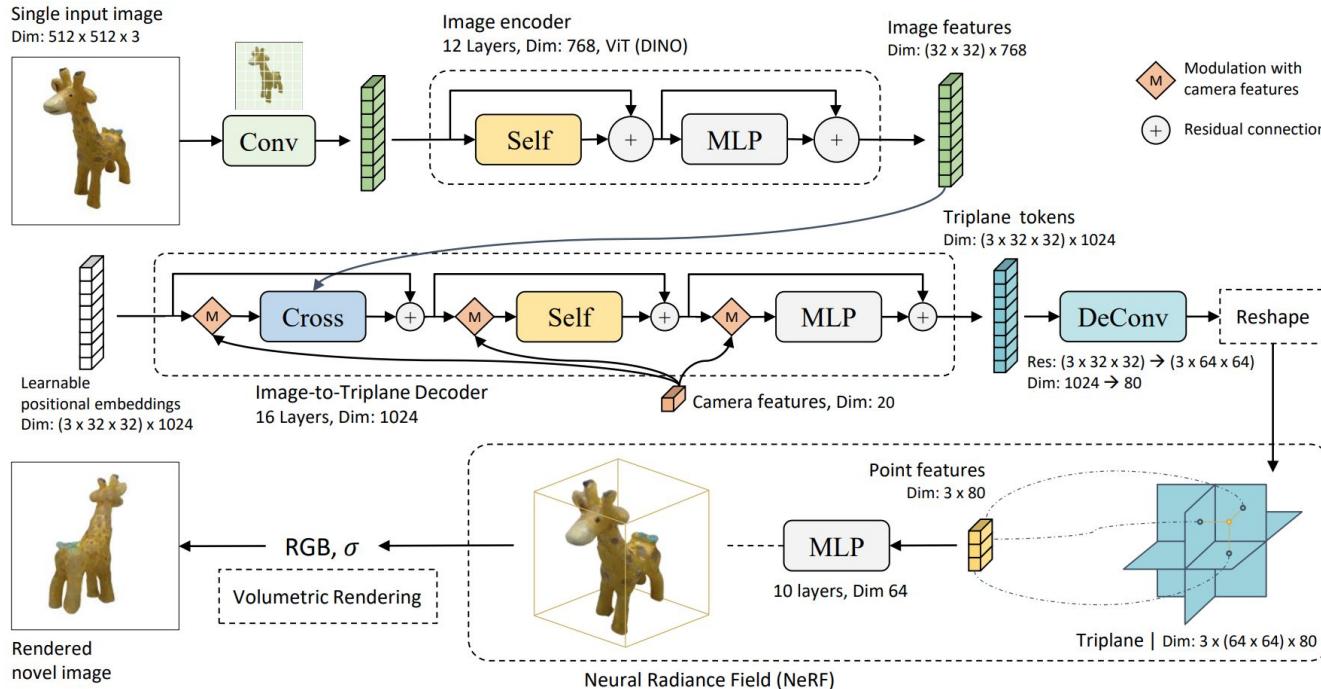
Instant3D: Instant Text-to-3D Generation



Under the unconditional generation setting (50 DDIM steps), sampling a batch of 8 scenes takes **4.63 sec** on a single RTX 3090 GPU.

*Single-Stage Diffusion
NeRF: A Unified Approach to
3D Generation and
Reconstruction*

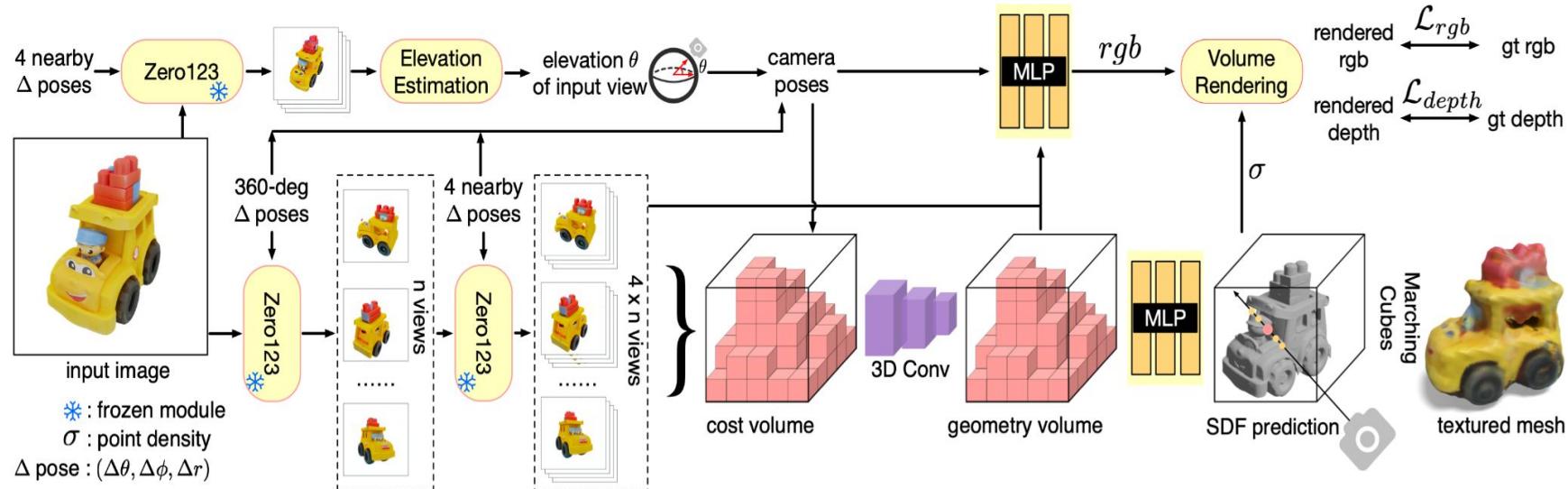
D.Optimization: Forward model: triplane



LRM: LARGE RECONSTRUCTION MODEL FOR SINGLE IMAGE TO 3D

D.Optimization: Forward model: cost-volume

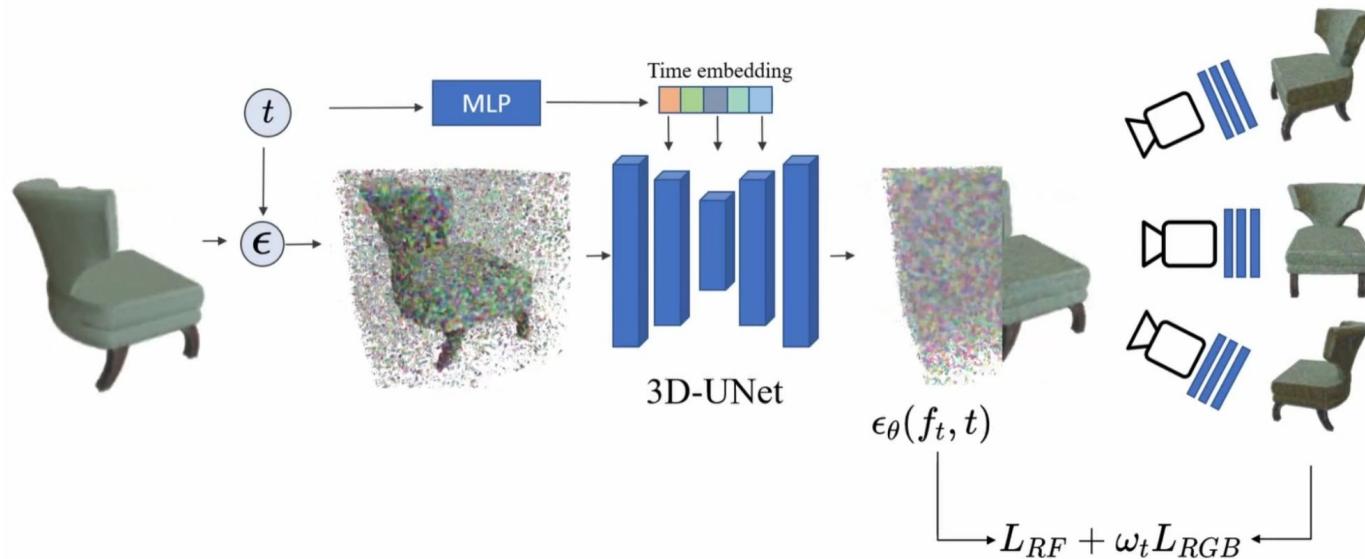
e.g. zero123



One-2-3-45: Any Single Image to 3D Mesh in 45 Seconds without Per-Shape Optimization

D.Optimization: Forward model: train 3D U-Net

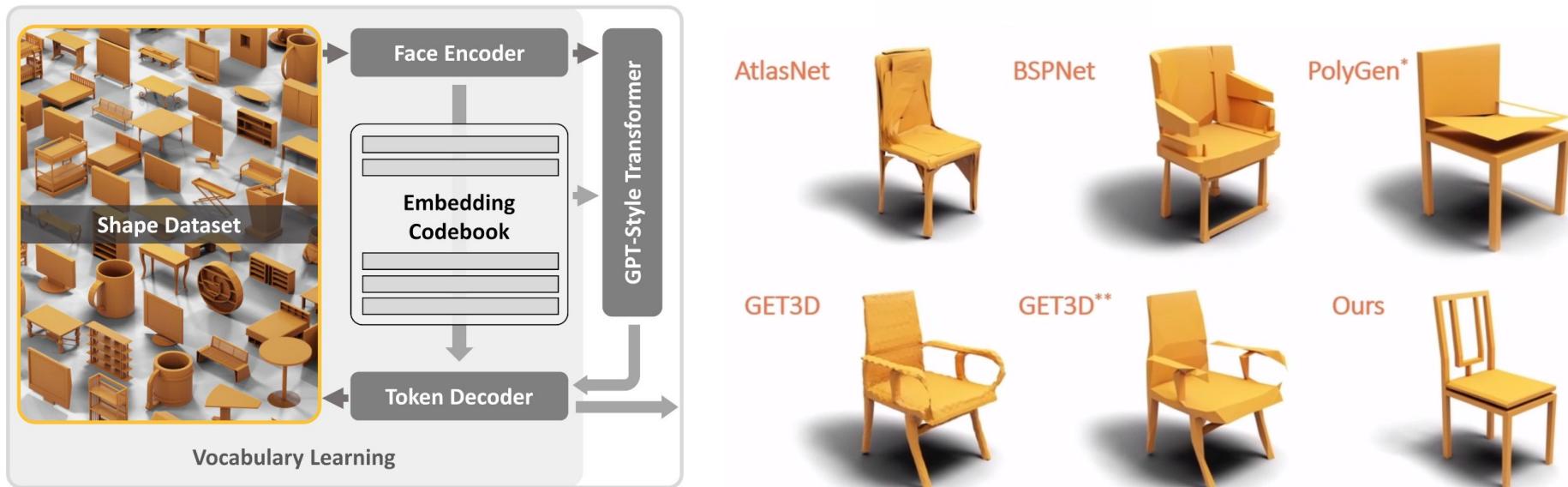
directly generate nerf from 3D U-Net



DiffRF is a denoising diffusion probabilistic model directly operating on 3D radiance fields and trained with an additional volumetric rendering loss. This enables learning strong radiance priors with high rendering quality and accurate geometry.

D.Optimization: Forward model: train a decoder

train a decoder transformer based on large 3d mesh dataset



D.Optimization: GPT-control

(a). "The sun shines brightly in a clear blue sky.
Not a cloud in sight."



(b). "Thick gray clouds cover the entire sky,
diffusing the sunlight."



(c). "The sky appears hazy, with a soft,
diffuse sunlight filtering through."



(d). "The sky is painted in shades of orange
and pink as the sun sets behind the island."



(e). "The sky is a mixture of soft purples
and pinks as the sun rises."



(f). "A canopy of stars fills the dark night
sky, with the island bathed in moonlight."



(a). **Solution:** From the description of "clear blue sky" and "not a cloud in sight", we can infer that the sun_intensity should be high, the sun_elevation should be around 90 (daytime), the sun_rotation can be set to 0. The air_density and dust_density should be low, so we can set them to 1 and 0 respectively. Since there are no clouds, the cloud_density should be 0. To make the sky appear bluer, we can increase the ozone value, so let's set it to 2. We can model the sky by calling the following function:

```
```python
sky_texture_node(sun_intensity='high', sun_elevation=90, sun_rotation=0, air_density=1, dust_density=0, ozone=2, cloud_density=0)
```
```