text-to-3D in Colab! <u>ajayj.com/dreamfields</u>



Text-to-image generation now supports incredible quality and diversity.

G

conditioning informatio

(latent, pose)



an armchair in the shape of an avocado." (DALL-E)



chromeplated cat sculpture placed on a ersian rug." (Imagen)



Neural Radiance Fields are photorealistic, but need many multiview photo captures. How to design novel content?



Can 3D reconstruction systems be trained with **no photos**? We perform text-to-3D synthesis by training Neural Radiance Fields with *zero input photos, only a caption.*



Problems / Contributions

1) Lack of 3D data

Problem: There aren't many 3D models available, making it difficult to train 3D generation.

Solution: Captioned photographs are abundant and diverse. We repurpose scalable pre-trained image-text representations like CLIP.

2) Underconstrained scene representation

Problem: Photorealistic 3D scene representations like NeRF are too flexible without photos or learned priors.

Solution: Regularize scene to be opaque but compact. Improve MLP for easier optimization.

Zero-Shot Text-Guided Object Generation with Dream Fields

Ben Mildenhall Pieter Abbeel Ben Poole Jon Barron Ajay Jain ajayj@berkeley.edu, {bmild, barron}@google.com, pabbeel@berkeley.edu, pooleb@google.com

Training Neural Radiance Fields with CLIP

DietNeRF regularized NeRF with a feature space loss, supervising Rendering novel views during training and transferring 2D knowledge.



Intuition: All views should have consistent semantics. "A bulldozer is a bulldozer from any perspective."

Our insight: Train NeRF entirely in feature space. An aligned image-text representation like CLIP allows us to **control** synthesis with a caption.



Regularizing Dream Fields

Dream Fields lacks a 3D prior, and NeRF can fit highly degenerate geometry (floaters, occlusions, translucency).



To fix geometry, we remove view dependence to model the shape without lighting and regularize mip-NeRF to be sparse and compact.



Augment backgrounds + random crop





+ Maximize transmittance

 $\sigma_{ heta}(\mathbf{r}(s)) ds$ $T(\mathbf{r}, \theta, t) = \exp\left(-\right)$ $\mathcal{L}_T = -\min(\tau, \operatorname{mean}(T(\theta, \mathbf{p})))$

"an illustration of a pumpkin on the vine."



Differentiable Image Parameterizations, Mordvintsev, et al., Distill, 2018.