

Neural fields for 3D Vision (a MAP perspective)

Andrea Tagliasacchi ( @taiyasaki)

Associate Professor – Simon Fraser University
Staff Research Scientist – Google Deepmind

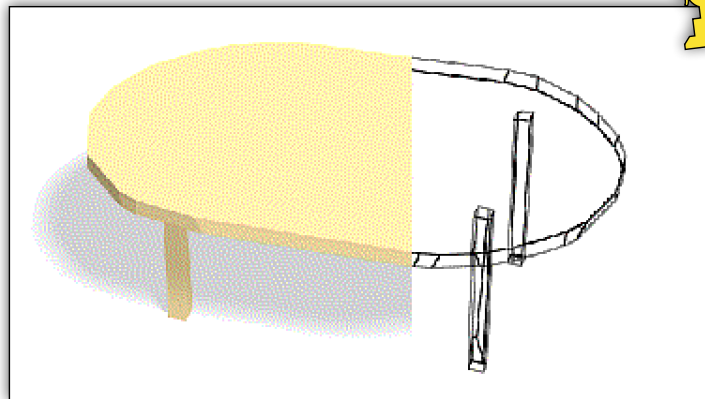


Google AI

SFU

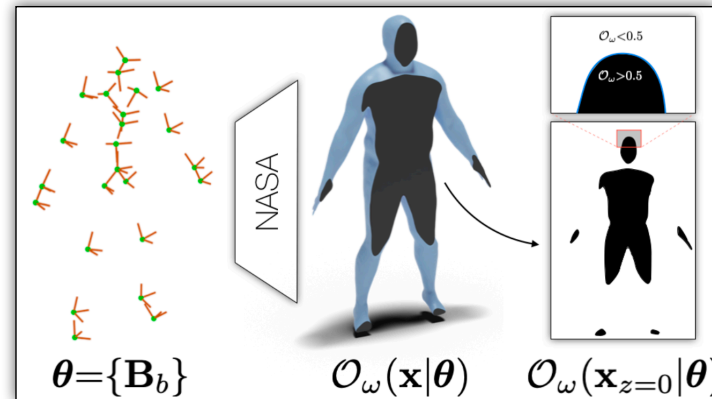
Recent: bridge graphics to vision

BSPNet @ CVPR'20



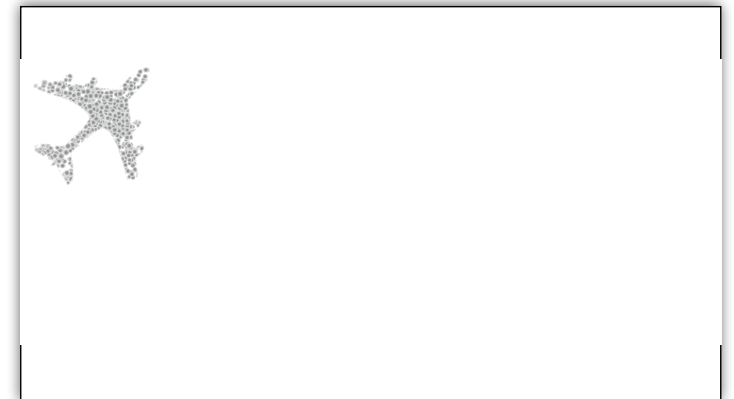
...represent meshes as fields

NASA @ ECCV'20



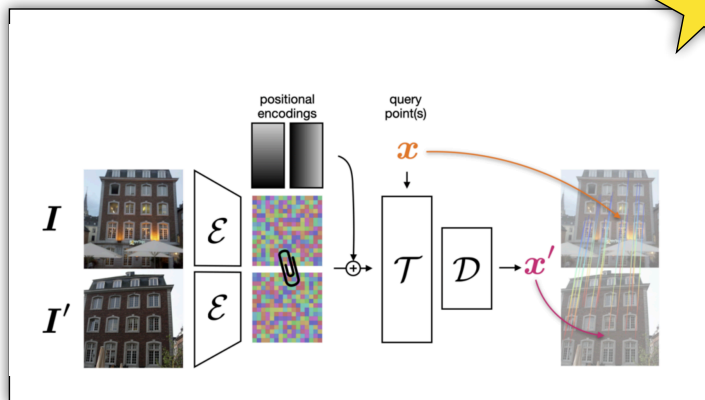
...represent humans as fields

Capsules @ NIPS'21



...canonicalize data

COTR @ ICCV'21



...image correspondences as fields

LoLNeRF @ CVPR'22



...multi-identity face field

Urban Fields @ CVPR'22



...leverage multi-sensor training data

CVPR 23: fields, fields, fields

MobileNeRF



what if I don't have a 4090?

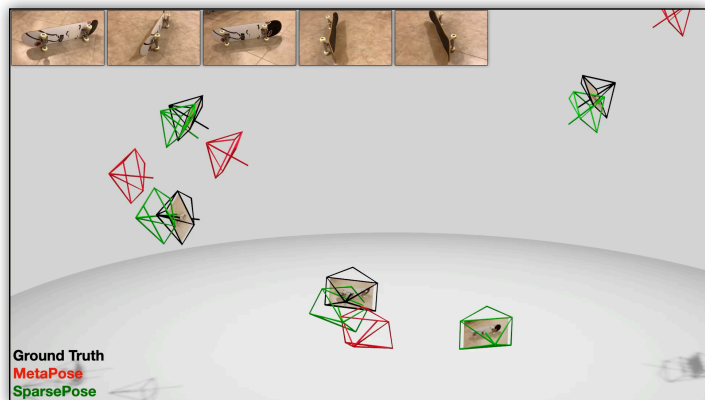
RobustNeRF



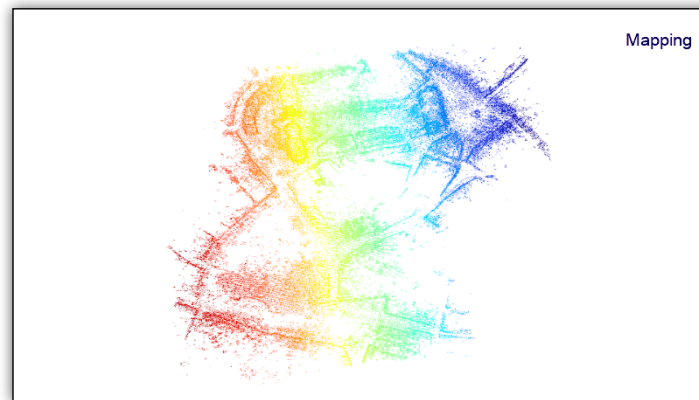
BlendFields



SparsePose



NeuMap



Continuous Upsampling



Neural Radiance Fields

- “NeRF freed us from the shackles of synthetic data”

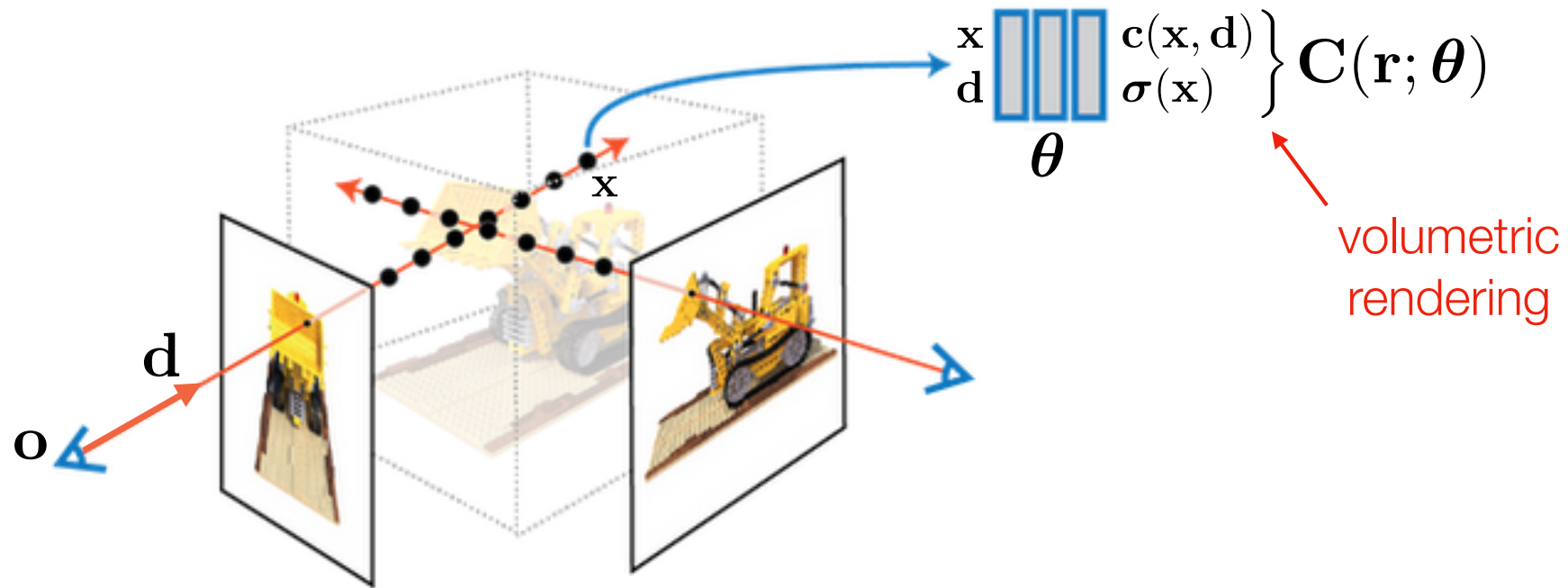


set of posed images



novel view synthesis

Basics of NeRF



$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{C} \sim \{\mathbf{C}_i\}} \mathbb{E}_{\mathbf{r} \sim \mathbf{C}} \|\mathbf{C}(\mathbf{r}; \boldsymbol{\theta}) - \mathbf{C}_i(\mathbf{r})\|_2^2$$

select image
select pixel/ray
predicted color
ground truth

NeRF as MAP (maximum a posteriori)

$$\begin{aligned}
 \arg \max_{\mathcal{M}} p(\mathcal{D}, \mathcal{M}) &\equiv \arg \max_{\mathcal{M}} p(\mathcal{D}|\mathcal{M}) \cdot p(\mathcal{M}) && \text{Bayes' rule} \\
 &\equiv \arg \max_{\mathcal{M}} \sum_i \underbrace{p(\mathbf{d}_i|\mathcal{M})}_{\text{likelihood}} \cdot \underbrace{p(\mathcal{M})}_{\text{prior}} && \text{i.i.d. pixels}
 \end{aligned}$$

↓
↓

photometric loss

$$\|\mathbf{C}(\mathbf{r}; \boldsymbol{\theta}) - \mathbf{C}_i(\mathbf{r})\|_2^2$$

spectral bias of MLPs

$$\left. \begin{array}{l} \mathbf{x} \\ \mathbf{d} \end{array} \right\} \left. \begin{array}{l} \mathbf{c}(\mathbf{x}, \mathbf{d}) \\ \sigma(\mathbf{x}) \end{array} \right\} \mathbf{C}(\mathbf{r}; \boldsymbol{\theta})$$

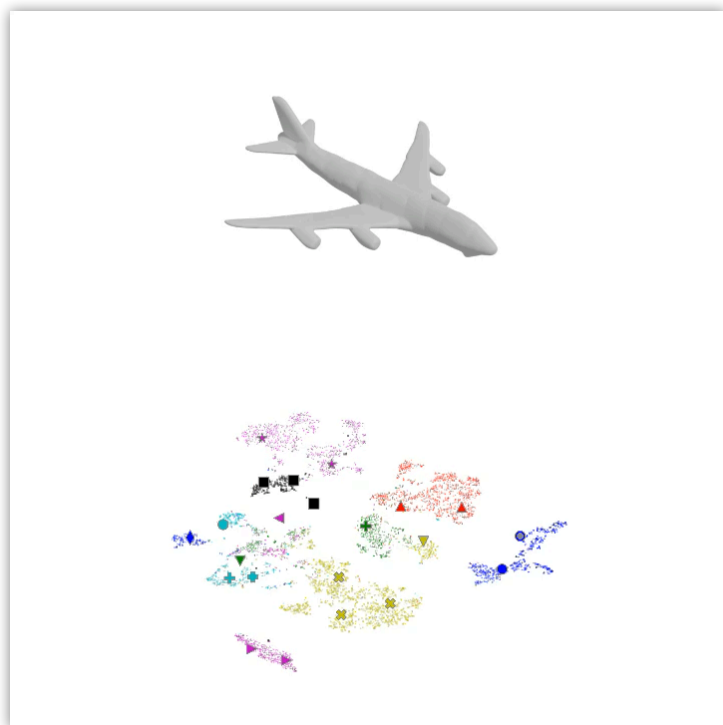
$\boldsymbol{\theta}$

A historical perspective



$$\arg \max_{\mathcal{M}} p(\mathcal{D}|\mathcal{M}) \cdot p(\mathcal{M})$$

encoder-decoder



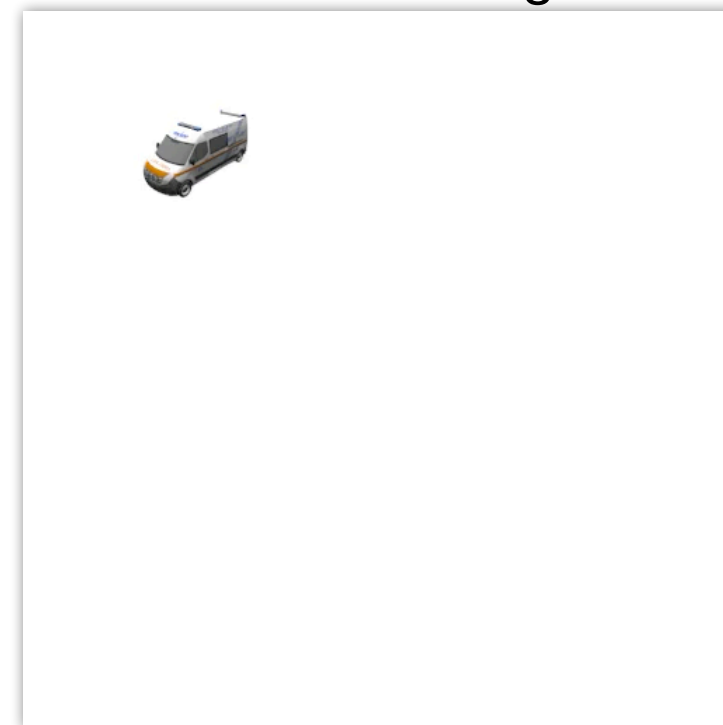
CvxNet @ CVPR'20

auto-decoder



LoLNeRF @ CVPR'22

meta-learning



3DiM @ ICLR'23

A historical perspective

$$\arg \max_{\mathcal{M}} p(\mathcal{D}|\mathcal{M}) \cdot p(\mathcal{M})$$

encoder-decoder



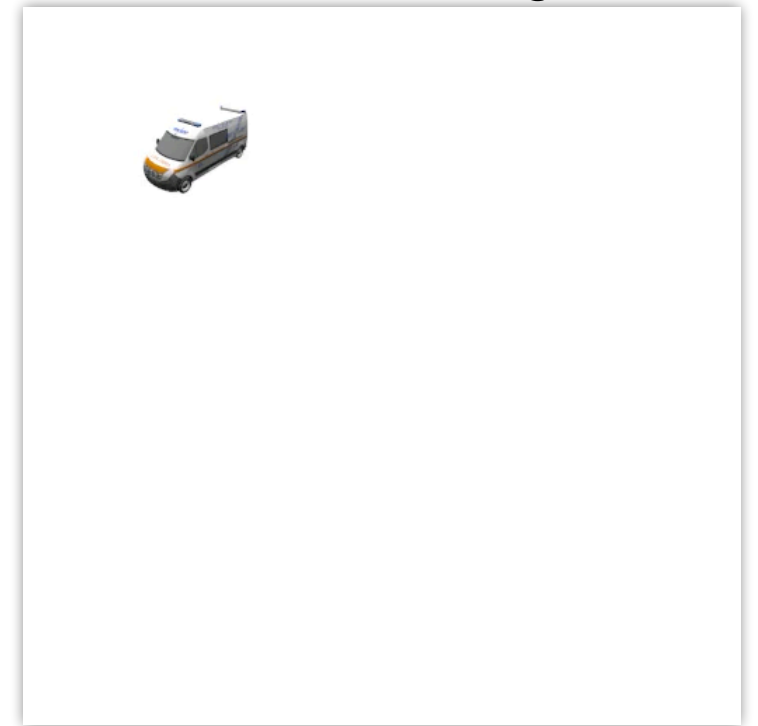
CvxNet @ CVPR'20

auto-decoder



LoLNeRF @ CVPR'22

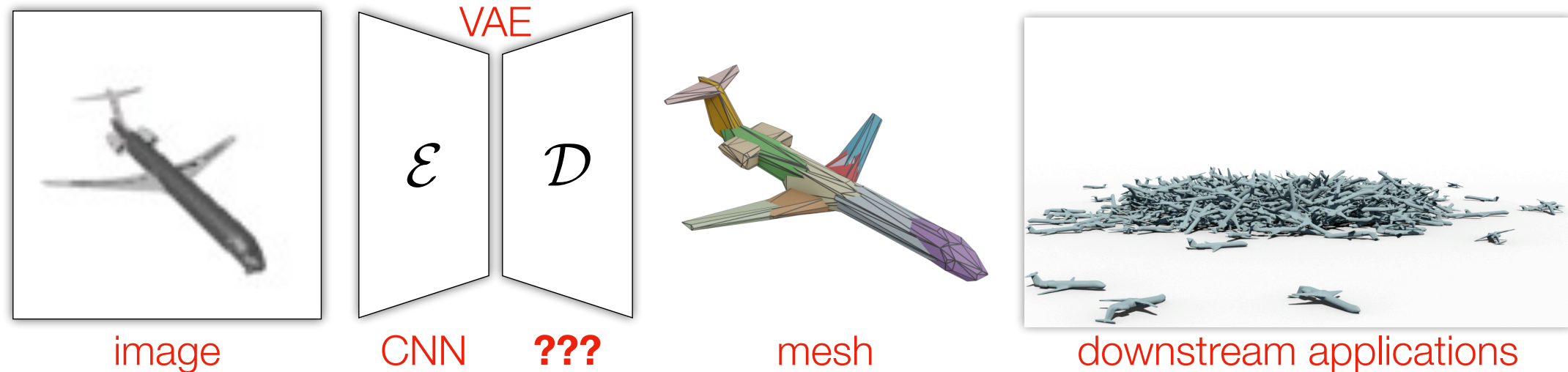
meta-learning



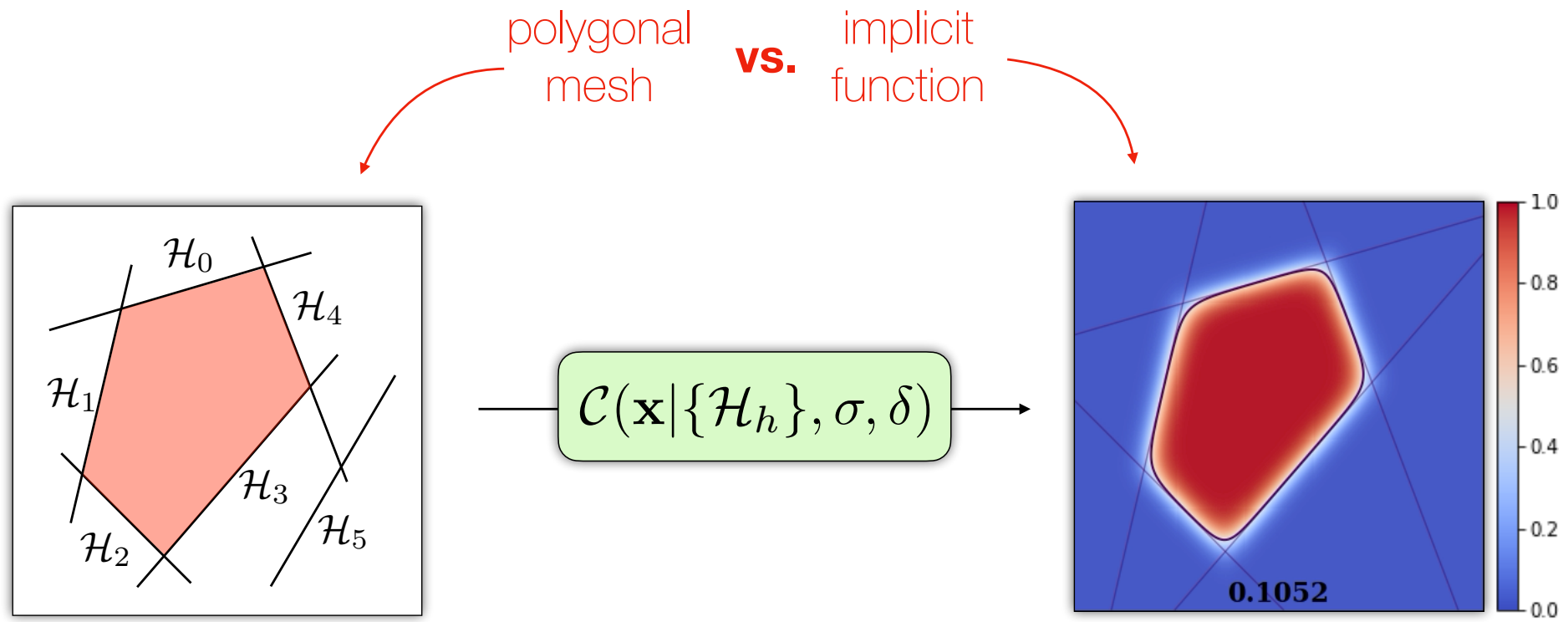
3DiM @ ICLR'23

Problem Statement

- What type of decoder should we use? ...a 3D CNN?

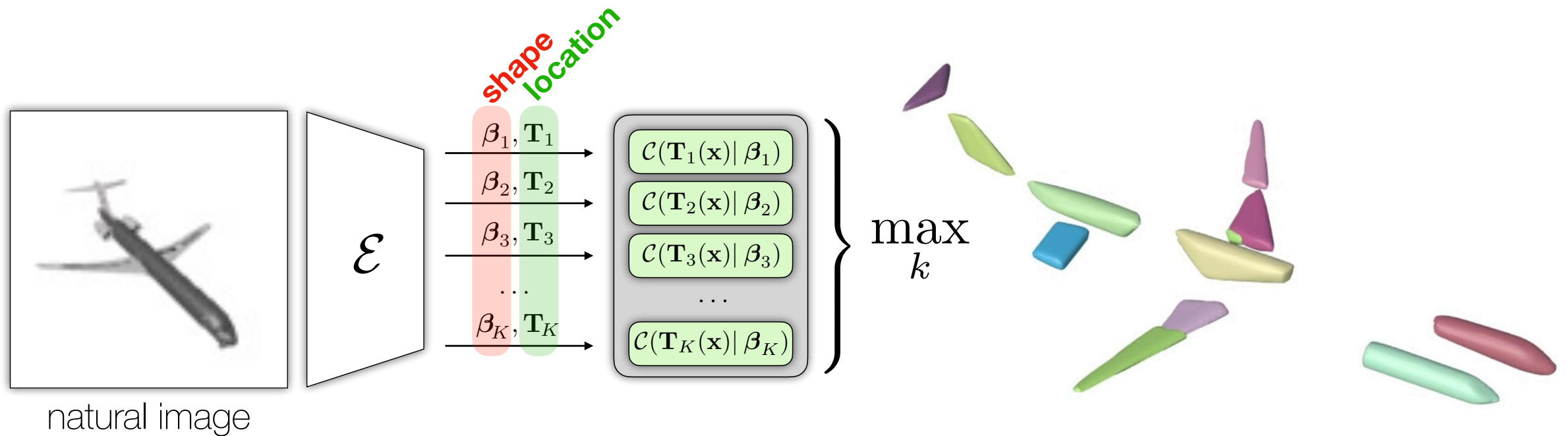


Convexes as differentiable fields

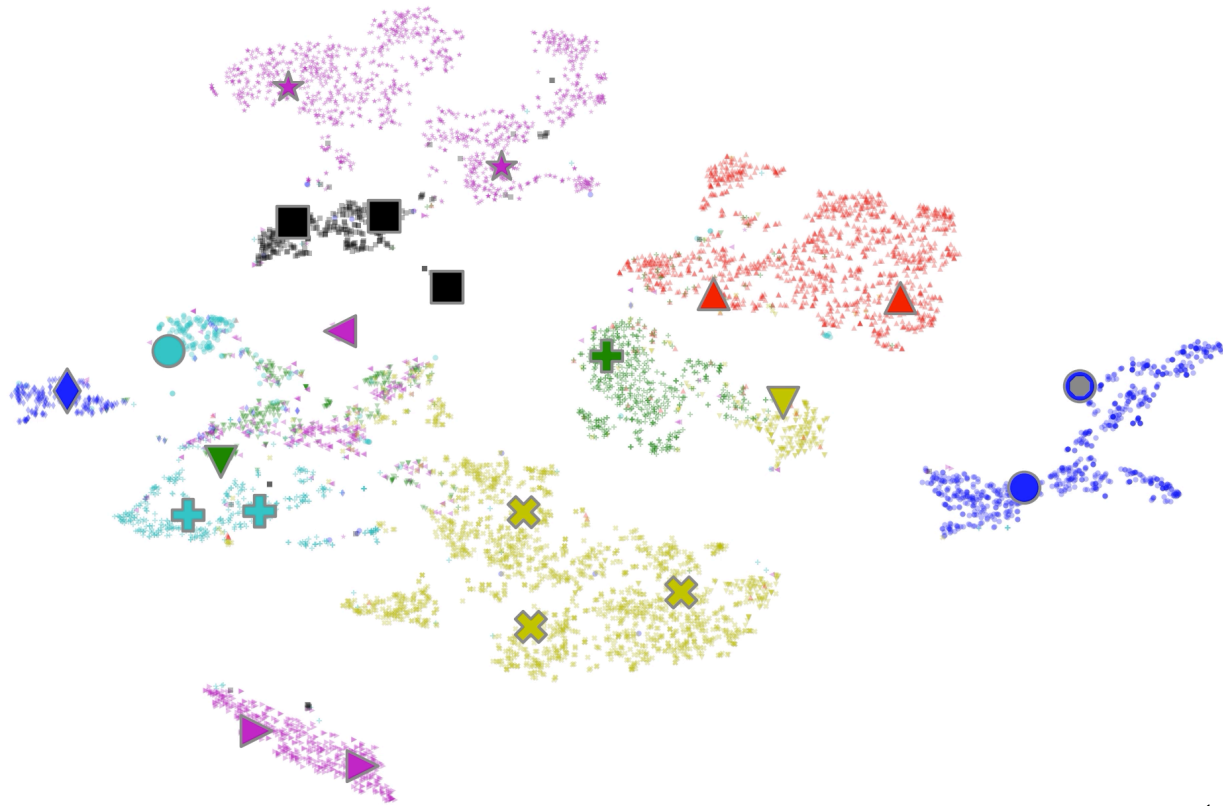


$$\text{Sigmoid}(-\sigma \text{LogSumExp}\{\delta(\mathbf{n}_h \cdot \mathbf{x} + d_h)\})$$

Convexes for generative modeling



Navigating the latent space



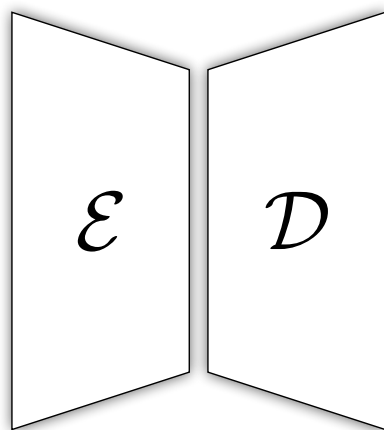
$$\arg \max_{\mathcal{M}} p(\mathcal{D}|\mathcal{M}) \cdot p(\mathcal{M})$$

So why this path “dried up”?

- Required 3D supervision and paired image/3D
- Shapes lack crisp details
 - Latent code causes representation bottleneck
 - NeRF positional encoding... does not help

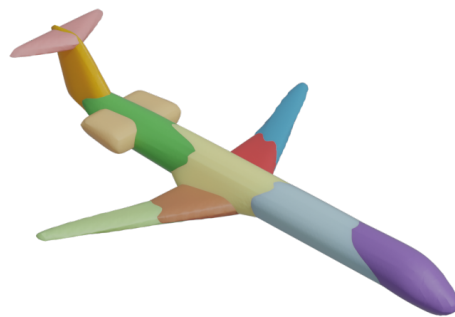


image



CNN

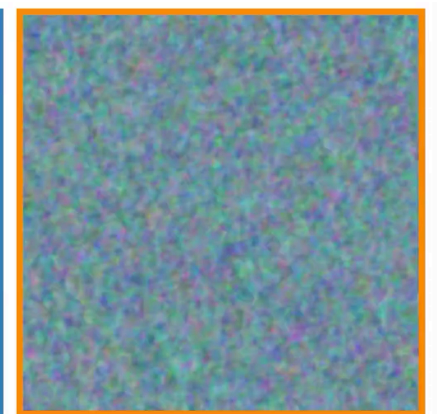
field



(smooth) mesh



w/o posenc

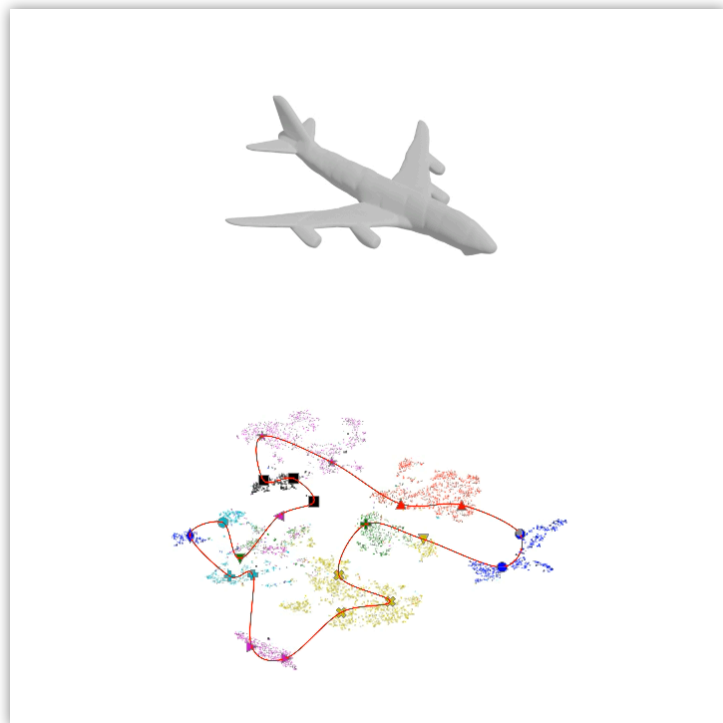


w/ posenc

A historical perspective

$$\arg \max_{\mathcal{M}} p(\mathcal{D}|\mathcal{M}) \cdot p(\mathcal{M})$$

encoder-decoder



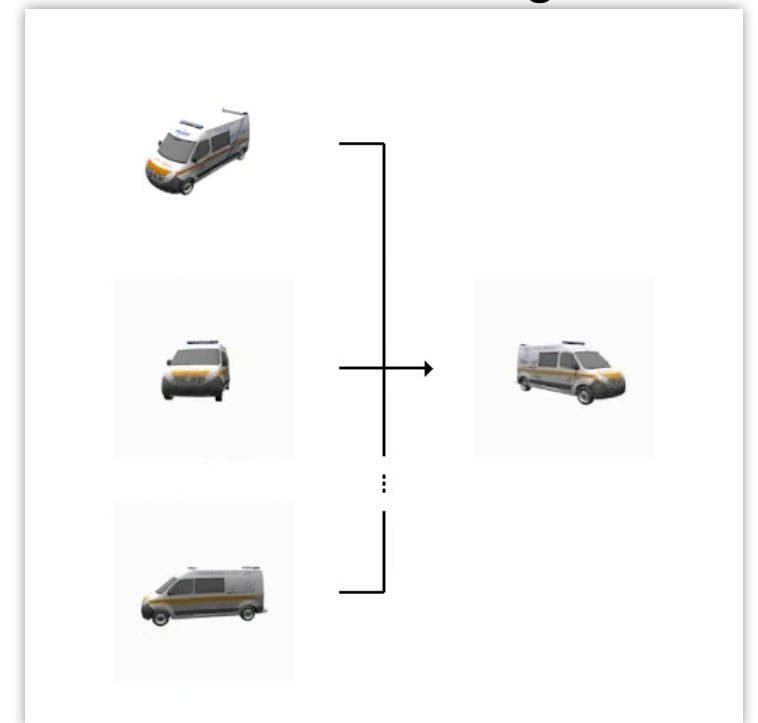
CvxNet @ CVPR'20

auto-decoder



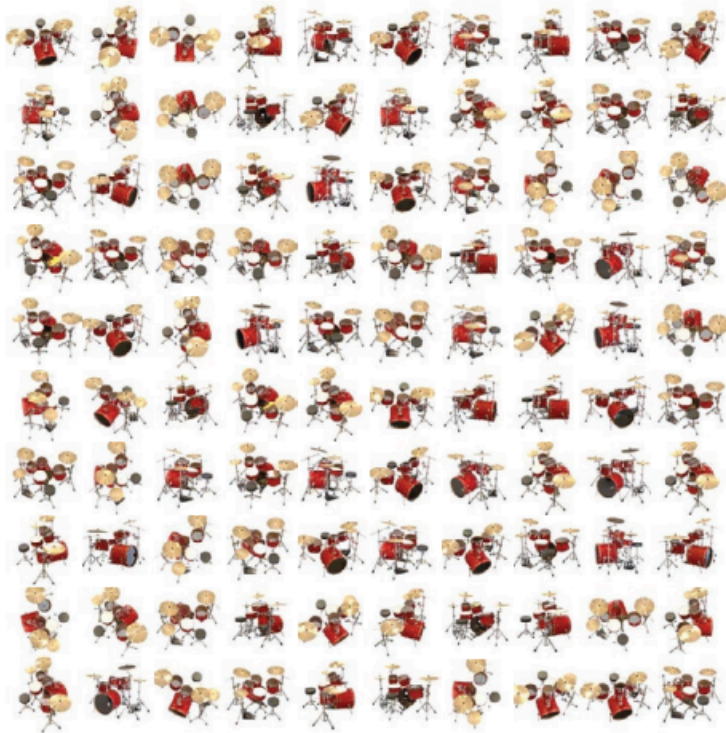
LoLNeRF @ CVPR'22

meta-learning



3DiM @ ICLR'23

Problem Statement



many images / scene
one scene

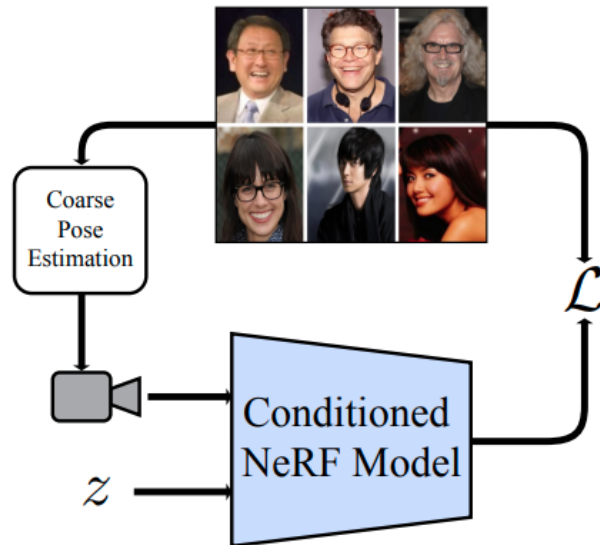
vs.



single image / scene
many scenes

Conditional Fields

- Lift 3D out 2D – via a large collection of 2D images (e.g. CelebA-HQ)
- Avoid relying on GANs to enforce 3D consistency (e.g. pi-GAN)



$$\mathcal{L}_{\text{rgb}}^{\mathbf{r},i}(\mathbf{z}_i, \boldsymbol{\theta}) = \|\mathbf{C}(\mathbf{r}; \mathbf{z}_i, \boldsymbol{\theta}) - \mathbf{C}_i(\mathbf{r})\|_2^2$$

$$\mathcal{L}_{\text{rgb}}(\boldsymbol{\theta}, \{\mathbf{z}_i\}) = \sum_i \mathbb{E}_{\mathbf{r} \sim \mathbf{C}_i} \left[\mathcal{L}_{\text{rgb}}^{\mathbf{r},i}(\mathbf{z}_i, \boldsymbol{\theta}) \right]$$

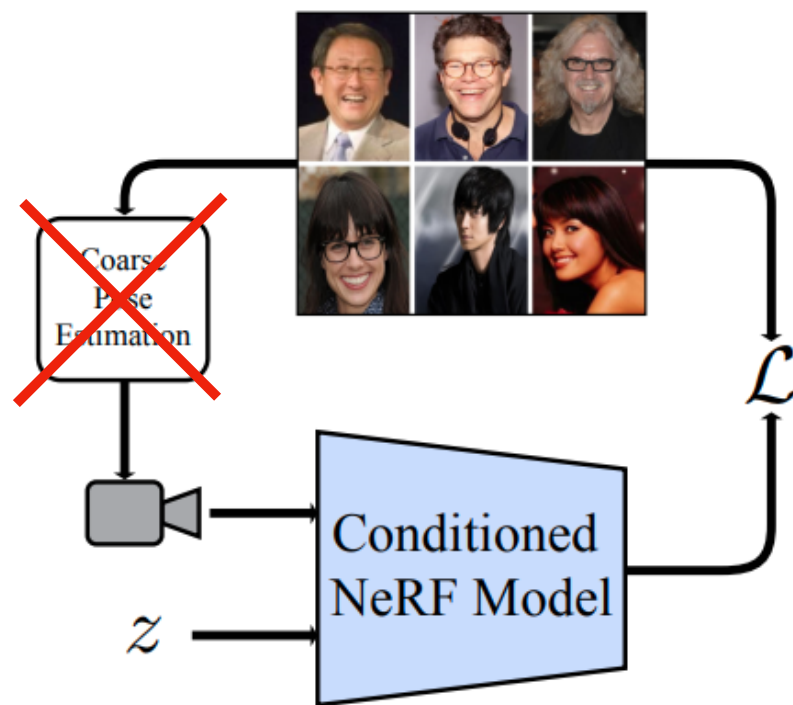
$$\arg \min_{\{\mathbf{z}_i\}, \boldsymbol{\theta}} \mathcal{L}_{\text{rgb}}(\{\mathbf{z}_i\}, \boldsymbol{\theta})$$

conditioning latent

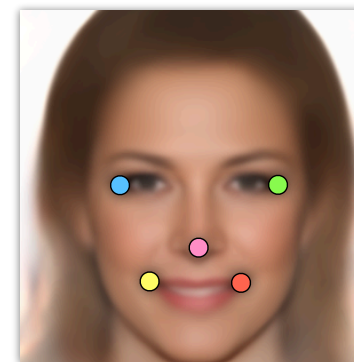
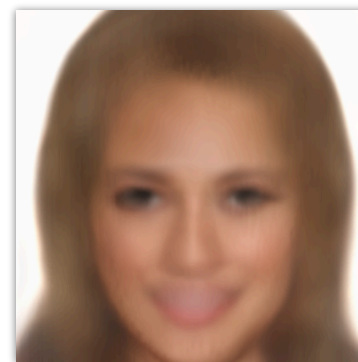
auto-decoder

Flatland mode collapse

- What happens if **no camera pose** is given?

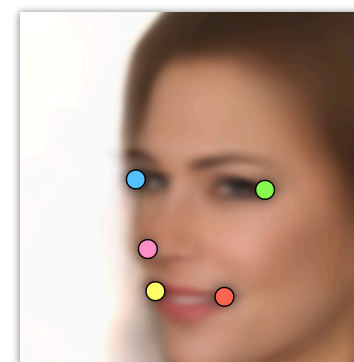
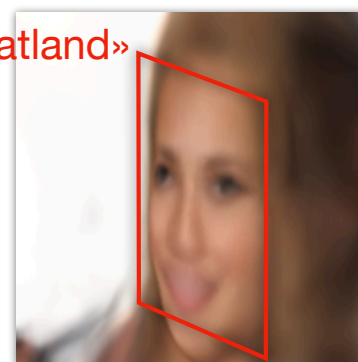


training
view



«flatland»

test
view



basic mode

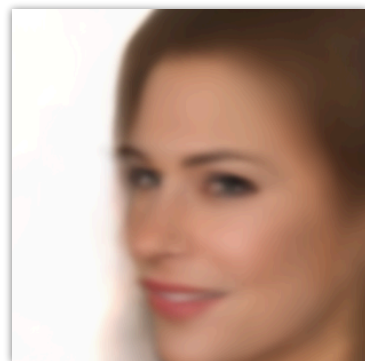
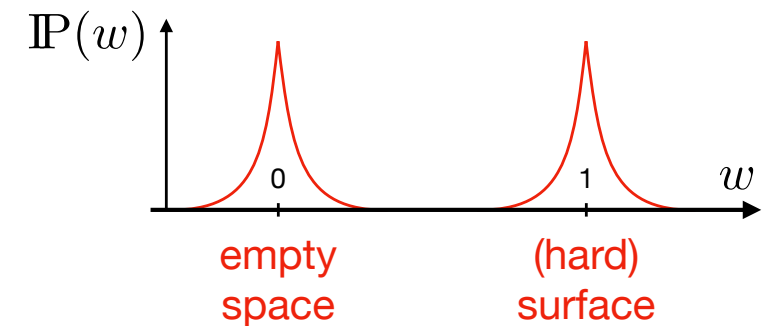
+camera pose

{empty | solid} prior model

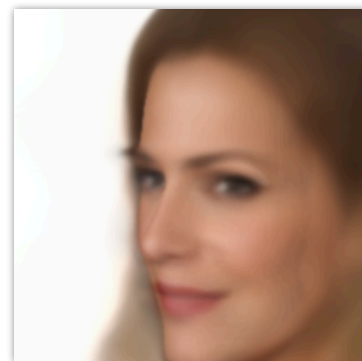
- What happens if most cameras are frontal?
- How to regularize this behavior?
 - skin ~ hard / solid media

$$\mathbf{C}(\mathbf{r}) = \int_{t_n}^{t_f} w(\mathbf{x}) \cdot \mathbf{c}(\mathbf{x}, \mathbf{d}) dt$$

$$\mathcal{L}_{\text{hard}} = -\log(\underbrace{e^{-|w|} + e^{-|1-w|}}_{\mathbb{P}(w)})$$



skin ≠ smoke

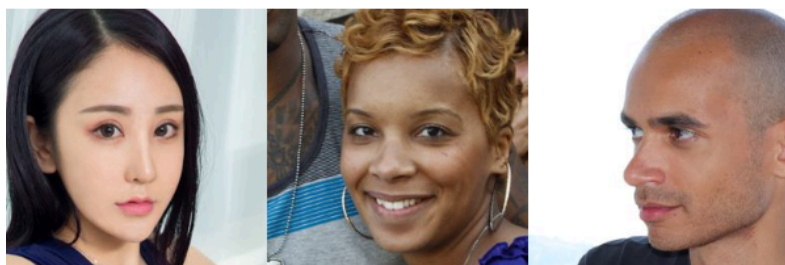


regularizes output geometry

Fitting *new* identities

- Seek the *latent-code* that minimizes the losses (i.e. frozen weights)

Input Image
(FFHQ)



piGAN fit
(CelebA@128²)



LoLNeRF fit
(CelebA@256²)



testing dataset

CelebA



FFHQ



Method	PSNR↑	SSIM↑	LPIPS↓	Res.
π -GAN [11] (CelebA)	21.8	0.796	0.412	256 ²
Ours (CelebA-HQ)	26.2	0.856	0.363	
π -GAN [11] (CelebA)	20.9	0.795	0.522	512 ²
Ours (CelebA-HQ)	25.1	0.831	0.501	
Ours (FFHQ)	25.3	0.836	0.491	

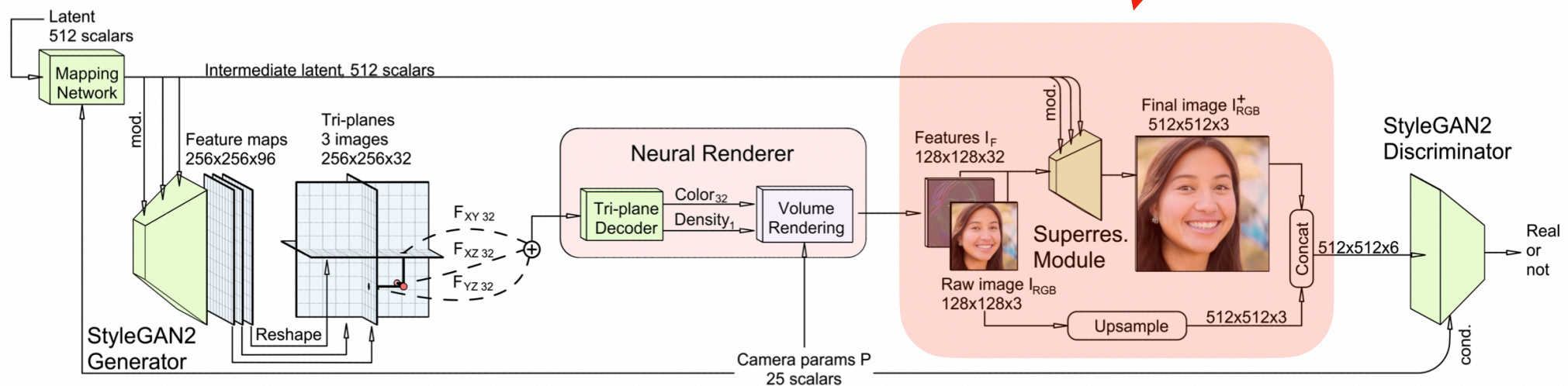
testing dataset

Can we move away from 2D?

- Hybrid of NeRF / GAN models
- supervised by adversarial 2D losses 🥲🥲

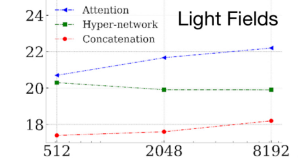


sweep: after/before upsampling

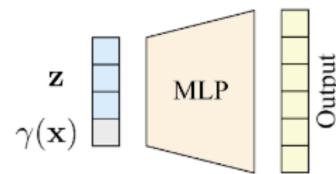


More powerful conditioning?

- Are **MLPs architectures** for **conditional fields** the way to go?

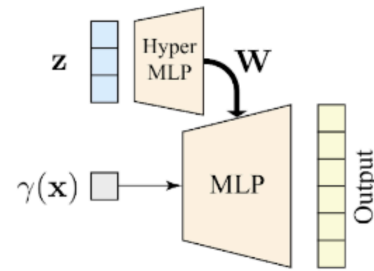


concatenate



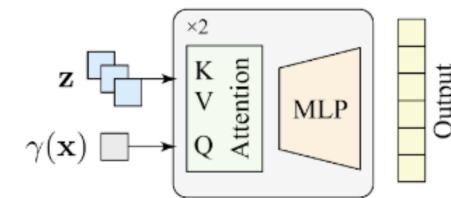
+2dB

hyper-network



+2dB

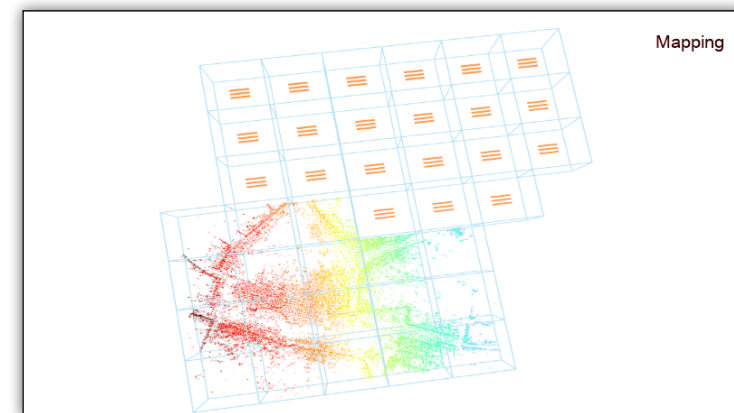
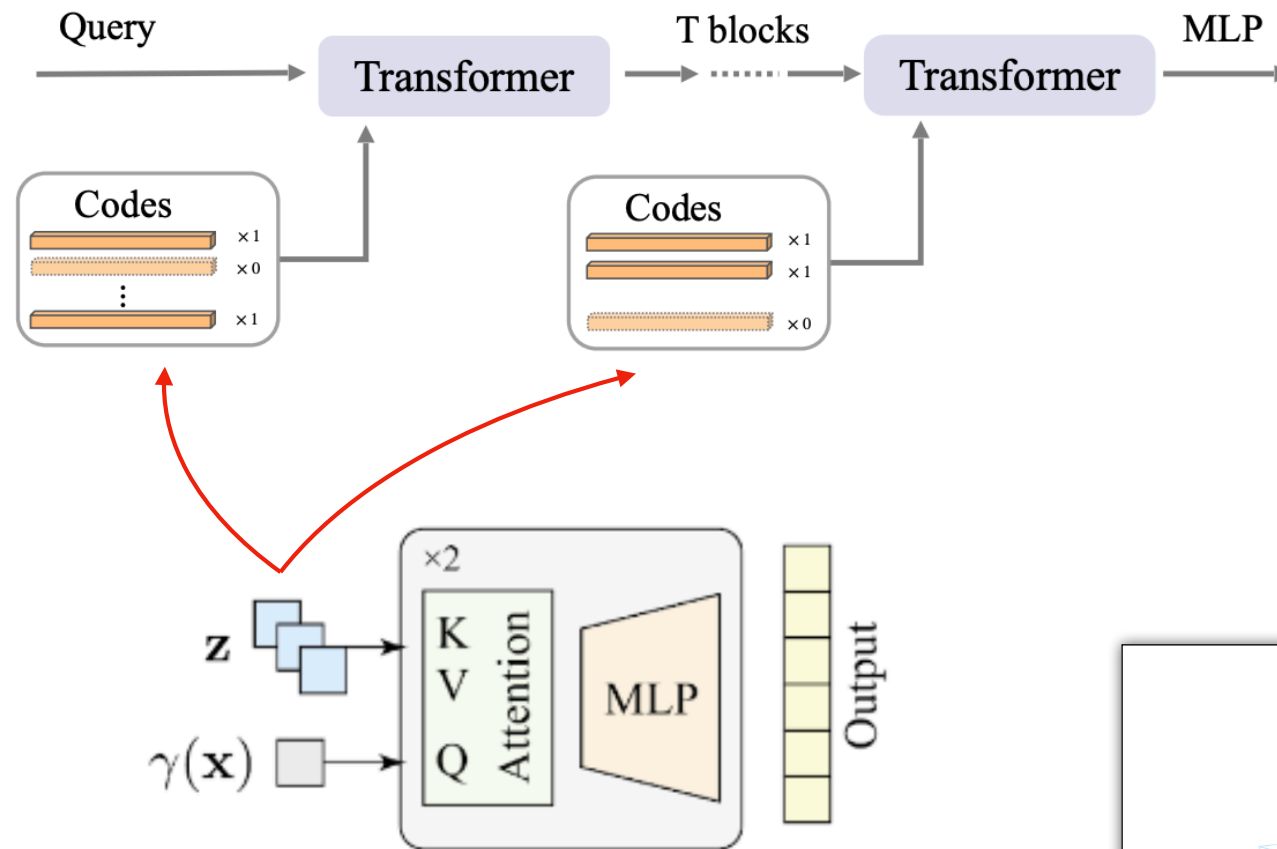
attention



Hierarchical Neural Field

hierarchical
transformer
field decoder

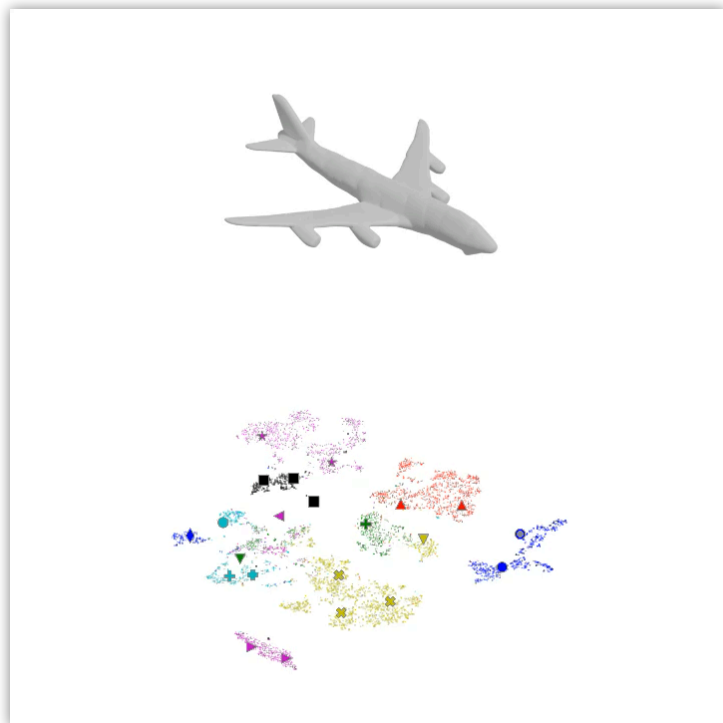
transformer
field decoder



A historical perspective

$$\arg \max_{\mathcal{M}} p(\mathcal{D}|\mathcal{M}) \cdot p(\mathcal{M})$$

encoder-decoder



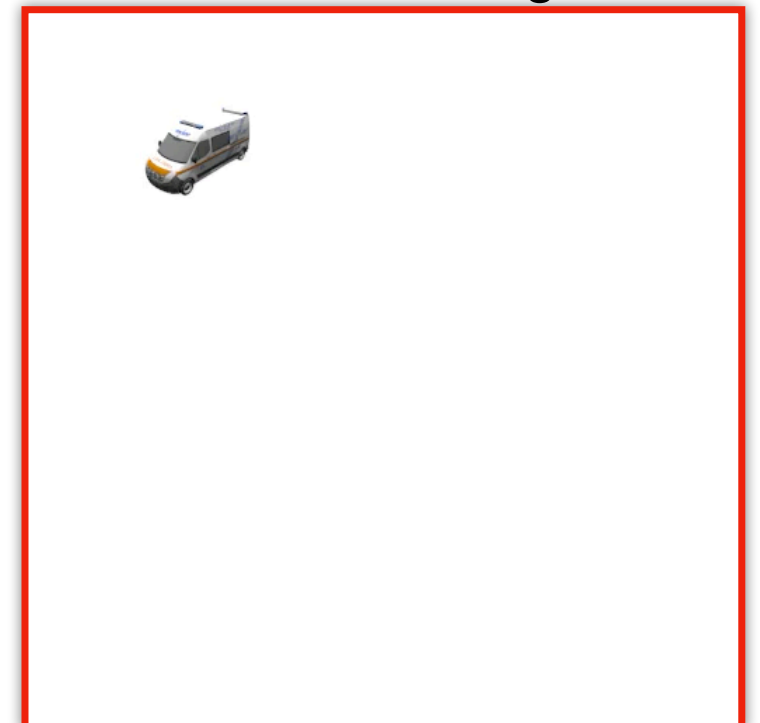
CvxNet @ CVPR'20

auto-decoder



LoLNeRF @ CVPR'22

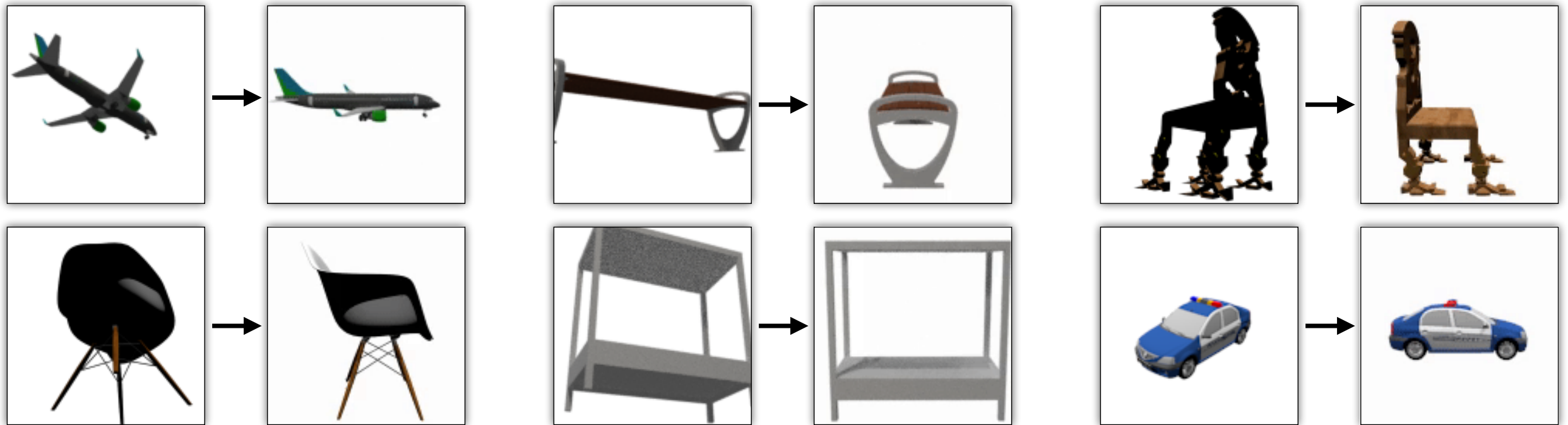
meta-learning



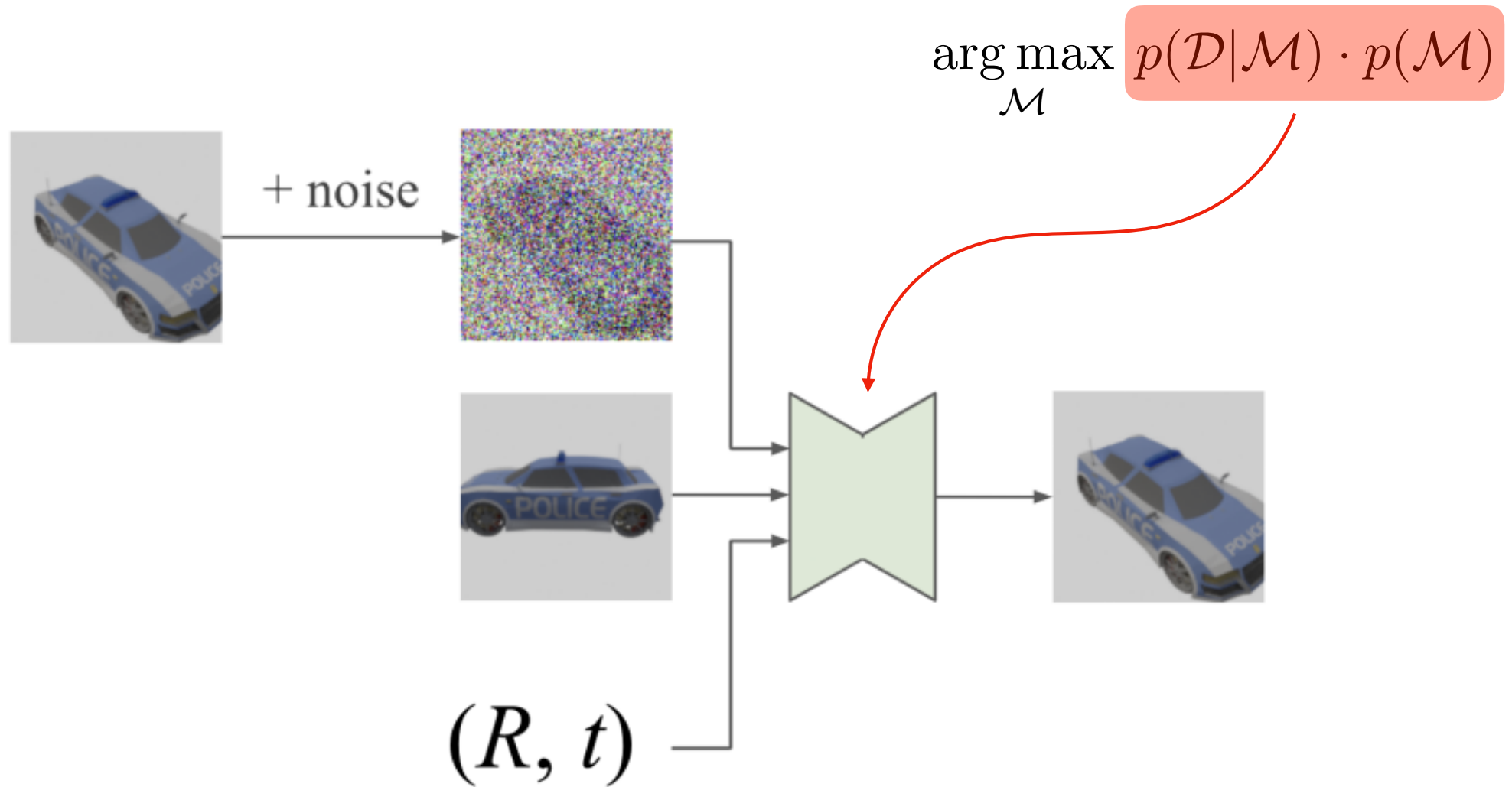
3DiM @ ICLR'23

Problem Statement

- Input: single image
- Output: novel view given relative camera
- Objective: can diffusion models learn 3D? ...without NeRF!

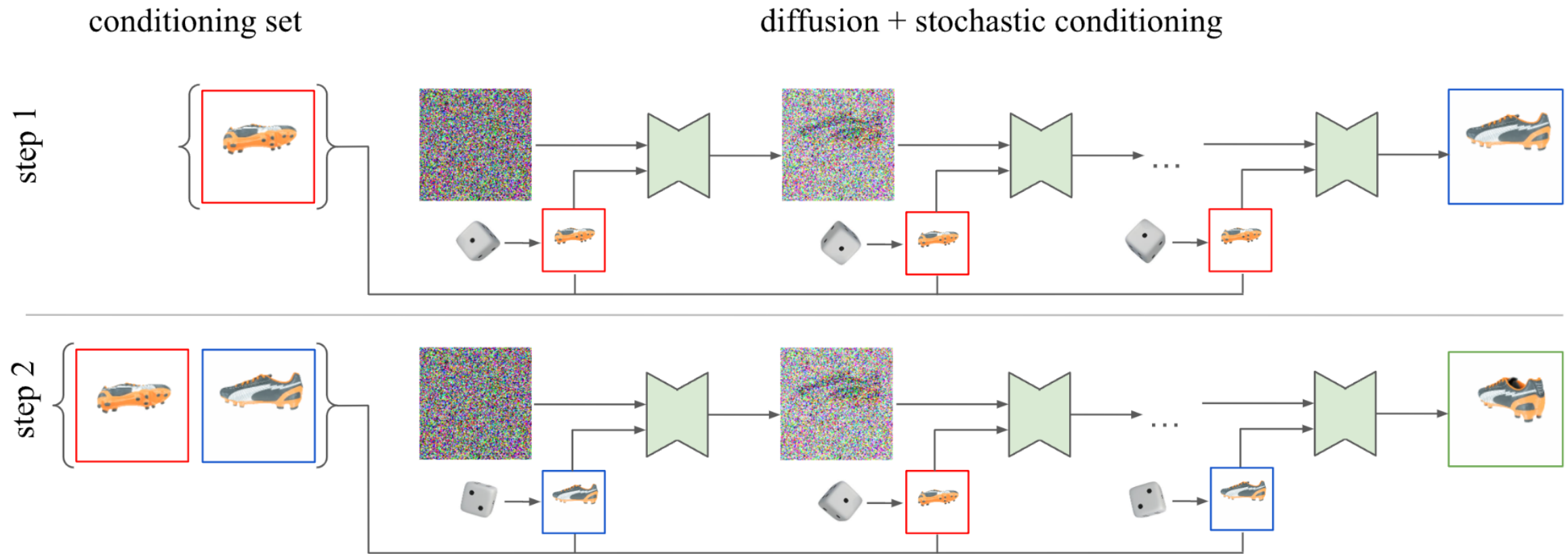


Novel View Synthesis by Diffusion



3DiM – Stochastic Conditioning

- Question: how to make sure $(n+1)$ -th image consistent with $\{0, 1, \dots, n\}$?



In-the-wild generalization

- Trained on ShapeNet renderings (Kubric), tested on in-the-wild images
- Outcome: diffusion model “learnt 3D” from raw data

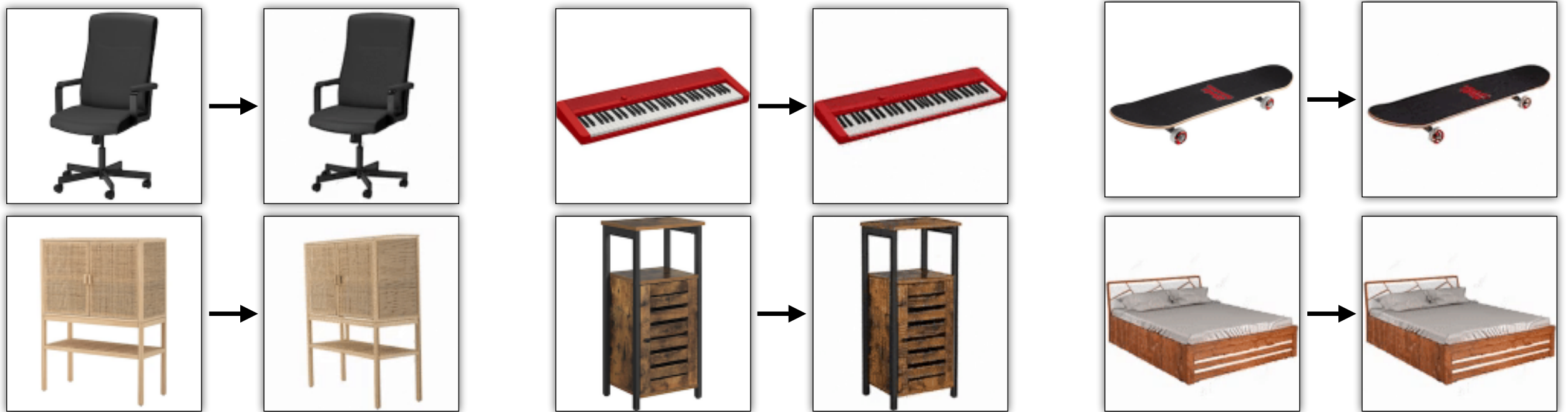
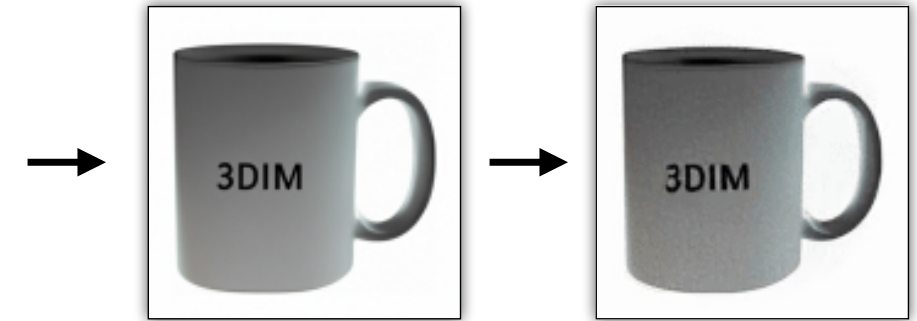


Imagen to 3D

“A mug with '3DiM' written on it, white background, no shadow.”



“An upright piano, at an angle, white background, no shadow.”



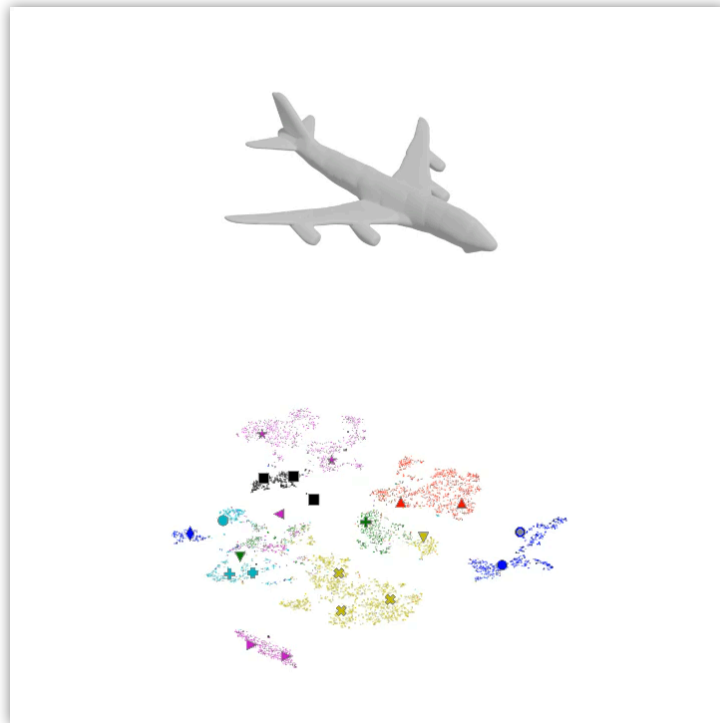
“A bed made of stone, white background, no shadows, at an angle.”



Foundation Models

$$\arg \max_{\mathcal{M}} \overset{\text{NeRF}}{p(\mathcal{D}|\mathcal{M})} \cdot \overset{\text{???}}{p(\mathcal{M})}$$

encoder-decoder



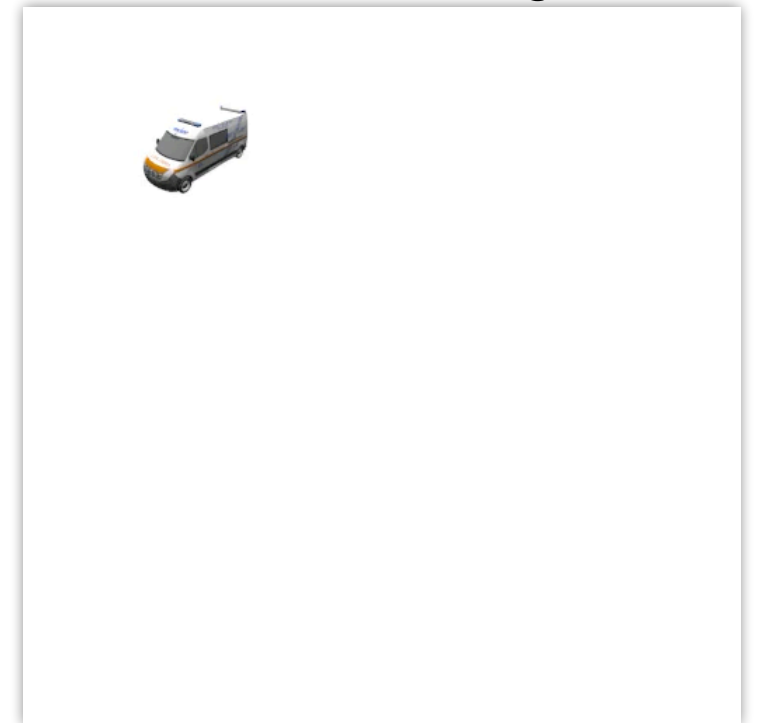
CvxNet @ CVPR'20

auto-decoder



LoLNeRF @ CVPR'22

meta-learning



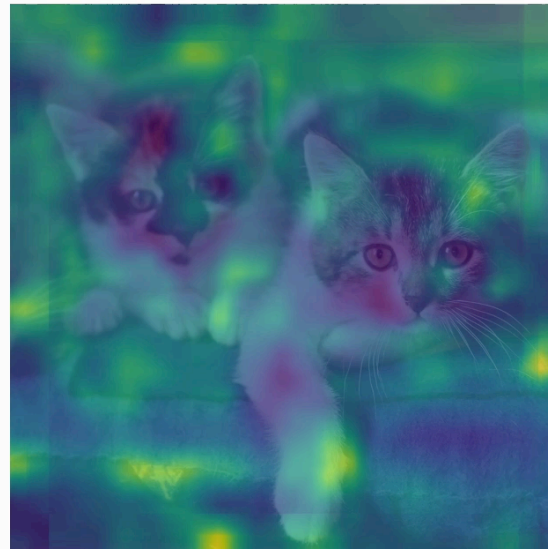
3DiM @ ICLR'23

Exciting new ideas

- Distill **structural understanding** from pre-trained image models
 - can these empower unsupervised 3D object discovery?



Source Image



Target Image



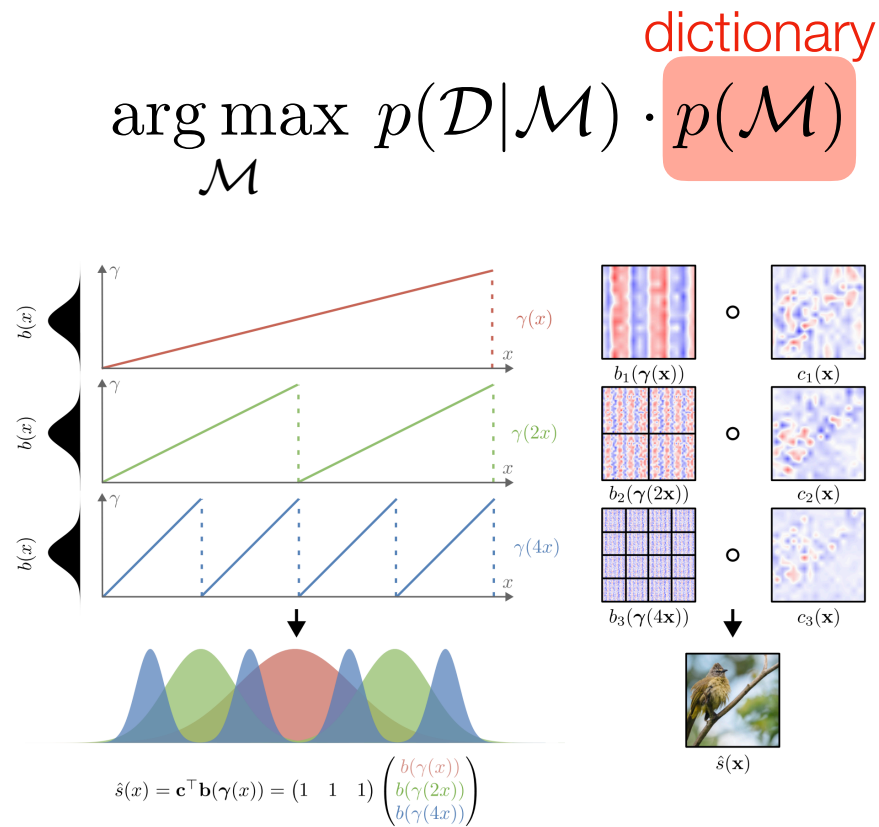
Target Image



Target Image

Exciting new ideas

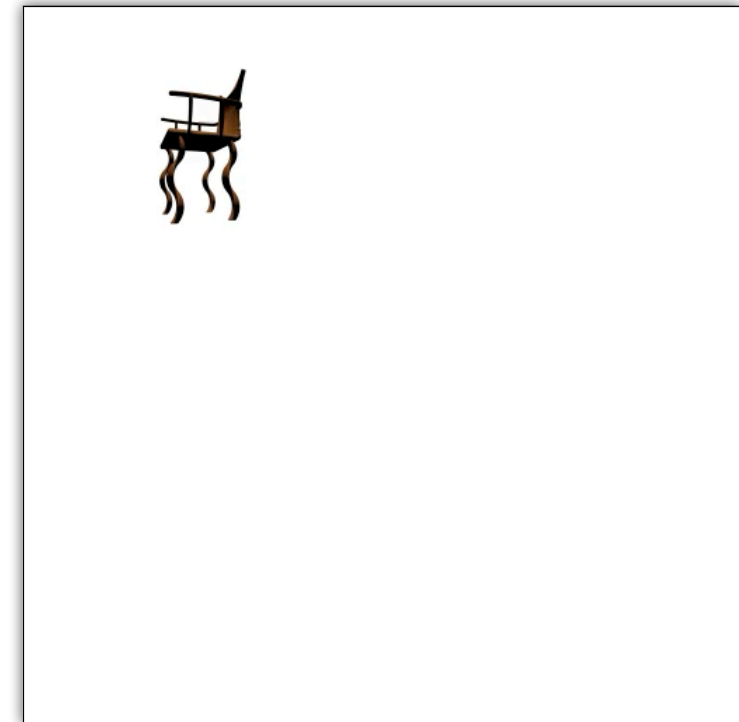
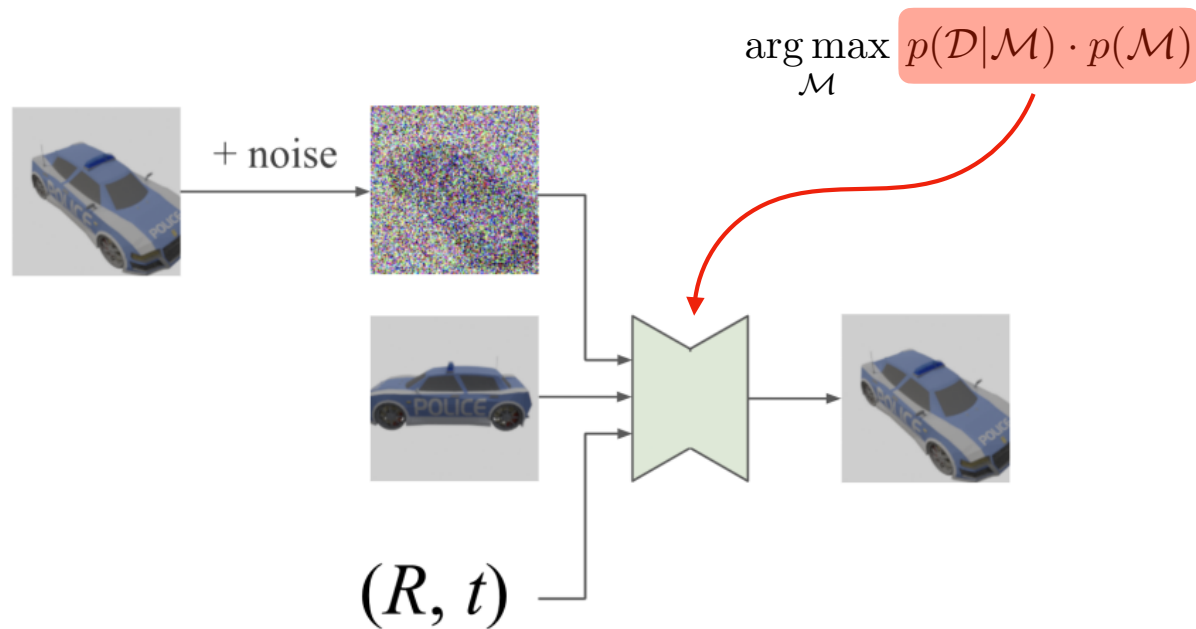
- Write (unobserved) fields as weighted combinations of (local) functions



few shot generalization

Exciting new ideas

- ...or perhaps learning from video at scale will “solve” 3D vision
 - is all we need large video datasets with calibrated cameras?



Neural fields for 3D Vision

Andrea Tagliasacchi ( @taiyasaki)

Associate Professor – Simon Fraser University
Staff Research Scientist – Google Deepmind



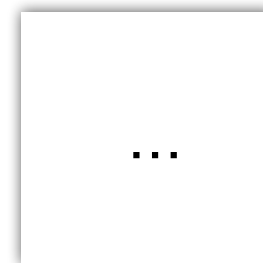
Boyang
Deng



Daniel
Rebain



Daniel
Watson



many
others