Video Autoencoder: self-supervised disentanglement of static 3D structure and motion



3D Trajectory

Zihang Lai CMU

Sifei Liu **NVIDIA**

Alexei A. Efros UC Berkeley

Xiaolong Wang UC San Diego



Disentangle the visual world





Depth





Structure and viewpoint



Prior work From the very beginnings of computer vision, ...



FIGURE 3 A SET OF INTRINSIC IMAGES DERIVED FROM A SINGLE MONOCHROME INTEN-SITY IMAGE

The images are depicted as line drawings, but, in fact, would contain values at every point. The solid lines in the intrinsic images represent discontinuities in the scene characteric; the dashed lines represent discontinuities in its derivative. In the input image, intensities correspond to the reflected light flux received from the visible points in the scene. The distance image gives the range along the line of sight from the center of projection to each visible point in the scene. The orientation image gives a vector representing the direction of the surface normal at each point. The reflectance image gives the albedo (the ratio of total reflected to total incident illumination) at each point.



(c) REFLECTANCE

(a) ORIGINAL SCENE

FIGURE 3 A SET OF INTRINSIC IMAGES DERIVED FROM A SINGLE MONOCHROME INTEN-SITY IMAGE

The images are depicted as line drawings, but, in fact, would contain values at every point. The solid lines in the intrinsic images represent discontinuities in the scene characteric; the dashed lines represent discontinuities in its derivative. In the input image, intensities correspond to the reflected light flux received from the visible points in the scene. The distance image gives the range along the line of sight from the center of projection to each visible point in the scene. The orientation image gives a vector representing the direction of the surface normal at each point. The reflectance image gives the albedo (the ratio of total reflected to total incident illumination) at each point.







Barrow and Tenenbaum, Comput. Vis. Syst., 1978





Prior work Learning disentangled visual representation from auto-encoders



Prior work Learning 3D from 2D supervisions





Objective In this work, we learn to separate 3D structure from Camera Motion without any human annotations



Input raw video



3D Structure



Camera trajectory



Test time The features obtained (3D structure, Camera Motion) can be used for several downstream tasks:



*Actual results



Learning from temporal continuity of videos

Spatio-temporal continuity



No spatio-temporal continuity





Learning from temporal continuity of videos Assume that a local snippet of video is capturing a static scene





Learning from temporal continuity of videos Assume that a local snippet of video is capturing a static scene





Learning from temporal continuity of videos Assume that a local snippet of video is capturing a static scene



11









3D Encoder



Input image







3D Encoder



Stacked image pair











Training loss

- Reconstruction loss:
 - $L_r(I_t, \hat{I}_t) = ||I_t \hat{I}_t||_1 + L_p(I_t, \hat{I}_t)$
- GAN loss:

•
$$L_g(\hat{I}_t) = -F_D(\hat{I})$$

- Consistency loss between deep voxels:
 - $L_c(V_{t1}, V_{t2}) = ||V_{t1} R(V_{t2}, P_{t2 \to t1})||_1$





Video Autoencoder

Results

17

Datasets



RealEstate10K

Matterport3D

Replica



Results **Novel view synthesis**

Single input image





Output video



RealEstate10K dataset





Results Novel view synthesis (Out-of-domain results)

Single input image



A Japanese living room (out-of-distribution)

Output video







Results Novel view synthesis (Out-of-domain results)

Single input image





Output video

"Spirited Away"



Results Novel view synthesis (Out-of-domain results)

Single input image



Bedroom in Arles, Vincent van Gogh

Output video





Results Comparison with previous methods



Input image RealEstate10K dataset



Mustikovela et al.

More artefacts



Wiles et al.



Tung et al.



Yu et al.



Ours



Ground truth 23



Results Comparison with previous methods



Input image RealEstate10K dataset



Mustikovela et al.



Wiles et al.



Tung et al.



Yu et al.



Ours



Ground truth²⁴



Results **Novel view synthesis (RealEstate10K)** Novel view synthesis task with RealEstate10K





camera intrinsics

Methods trained with full camera poses



Results Novel view synthesis (Matterport3D & Replica)

		Matterport 3D		MP3D → Replica	
Method	Pose	PSNR 1	SSIM 1	PSNR 1	SSIM 1
Methods without any camera supervision					
Ours	×	20.58	0.64	21.72	0.77
Methods with full camera supervision					
Dosovisky et al.	\checkmark	14.79	0.57	14.36	0.68
Appearance Flow	\checkmark	15.87	0.53	17.42	0.66
SynSin (w/ voxel)	\checkmark	20.62	0.70	19.77	0.75
SynSin (w/ point cloud)	√	20.91	0.72	21.94	0.81



Results Camera pose estimation

Input video



Trajectory Prediction





Results **Camera pose estimation**



Comparisons on 30-frame videos



Results Video following

Input image

Followed Video





Result

Camera Shaking





Results Video following

Input image

Followed Video



Result

Rotating Right



Please see paper and website for details

https://zlai0.github.io/VideoAutoencoder

Thanks!

