Render for CNN:

Viewpoint Estimation in Images Using CNNs Trained with Rendered 3D Model Views



Hao Su^{*} Charles R. Qi^{*} Yangyan Li Leonidas J. Guibas



ILSVRC Image Classification Top-5 Error (%)



2010 2011 2012 2013 2014 2015

Go beyond 2D Image Classification



- 3D bounding box
- 3D alignment
- 3D model retrieval

Go beyond 2D Image Classification



3D Viewpoint Estimation



3D Viewpoint Estimation in the Wild







Images in the Wild Models unknown







3D Perception in the Wild



Wild Models unknown

What's the camera viewpoint angles to the SUV in the image?



PASCAL₃D+ dataset [Xiang et al.]











High-cost Label Acquisition

High-capacity Model



30K images with viewpoint labels in PASCAL₃D+ dataset [Xiang et al.]

6oM parameters. AlexNet [Krizhevsky et al.]

How to get MORE images with ACCURATE viewpoint labels?



Good News: ShapeNet



Key Idea: Render for CNN



Synthetic Images

Key Idea: Render for CNN



Real Images













- 8oK rendered chair images
- Metric: 16-view classification accuracy **tested on real images**

At beginning..

- Lighting: 4 **fixed** point light sources on the sphere
- Background: clean



95% on synthetic val set47% on real test set 8

ConvNet: Ah ha, I know! Viewpoint is just the brightness pattern!



95% on synthetic val set47% on real test set 8





ConvNet: hmm.. viewpoint is not the brightness pattern. Maybe it's the contour?



ConvNet: hmm.. viewpoint is not the brightness pattern. Maybe it's the contour?



ConvNet: It becomes really hard! Let me look more into the picture.

bbox crop texture

86% -> 93%





Key Lesson: Don't give CNN a chance to "cheat" - it's very good at it. When there is no way to cheat, true learning starts.

Render for CNN Image Synthesis Pipeline



Render for CNN Image Synthesis Pipeline



Rendering

Lighting params Randomly sampled

- Number of light sources
- Light distances
- Light energies
- Light positions
- Light types

Camera params KDE from PASCAL3D+ train set



Render for CNN Image Synthesis Pipeline



Background Composition

- Simple but effective!
- Backgrounds randomly sampled from SUN397 dataset [Xiao et al.]
- Alpha blending composition for natural boundaries



Render for CNN Image Synthesis Pipeline



Image Cropping

Cropping patterns KDE from PASCAL3D+ train set



Image Cropping

Cropping patterns KDE from PASCAL₃D+ train set



2.4M Synthesized Images for 12 Categories

- High scalability
- High quality
 - Overfit-resistant
 - Accurate labels





Results

Metric: median angle error (lower the better)

Real test images from PASCAL₃D+ dataset

	aero	bike	boat	bottle	bus	car	chair	table	mbike	sofa	train	tv	mean
$Acc_{\frac{\pi}{6}}$ (Tulsiani, Malik)	0.78	0.74	0.49	0.93	0.94	0.90	0.65	0.67	0.83	0.67	0.79	0.76	0.76
$Acc_{\frac{\pi}{6}}$ (Ours-Render)	0.74	0.83	0.52	0.91	0.91	0.88	0.86	0.73	0.78	0.90	0.86	0.92	0.82
MedErr (Tulsiani, Malik)	14.7	18.6	31.2	13.5	6.3	8.8	17.7	17.4	17.6	15.1	8.9	17.8	15.6
MedErr (Ours-Render)	15.4	14.8	25.6	9.3	3.6	6.0	9.7	10.8	16.7	9.5	6.1	12.6	11.7

Our model **trained on rendered images** outperforms state-of-the-art model **trained on real images** in PASCAL₃D+.



How many 3D models are necessary?



3D Viewpoint Estimation

























Azimuth Viewpoint Estimation



Azimuth Viewpoint Estimation



Failure Cases

sofa occluded by people



0 360	0	90	b
(0	90	D

car occluded by motorbike



ambiguous car viewpoint



ambiguous chair viewpoint



multiple cars



multiple chairs

90

0



180

270

360

Limitations of Current Synthesis Pipeline

- Modeling Occlusions?
- Modeling Background Context?
- Shape database augmentation by interpolation?

Render for CNN – Beyond Viewpoint

- 3D model retrieval
 - Joint Embedding [Li et al sigasia15]
- Object detection
- Segmentation
- Intrinsic image decomposition and an experimentation of the second second
- Controlled experiments for DL
- Vision algorithm verification



Conclusion

Images rendered from 3D models can be effectively used to train CNNs, especially for 3D tasks. State-of-the-art result has been achieved.

Keys to success

- Quantity: Large scale 3D model collection (ShapeNet)
- Quality: Overfit-resistant, scalable image synthesis pipeline



THE END

THANKYOU!