

Learning 3D representations, disparity estimation, and structure from motion

Thomas Brox University of Freiburg, Germany

Research funded by the ERC Starting Grant VideoLearn, the German Research Foundation, and the Deutsche Telekom Stiftung







3D shape and texture from a single image



FlowNet: end-to-end optical flow



DispNet: end-to-end disparities



DeMoN: end-to-end structure from motion

Single-view to multi-view

Single-V

Maxim Tatarchenko Alexey Dosovitskiy ECCV 2016

Up-convolutional part



Single-view to multi-view



Synthetic images



Real images



Multi-view looks like 3D





Reconstructing explicit 3D models

Input images











Multiview morphing



Other interesting work

Yang et al. NIPS 2015 Recurrent network, incrementally rotates the object





Ours for comparison

Kar et al. CVPR 2015 Choy et al. 2016







3D shape and texture from a single image



FlowNet: end-to-end optical flow



DispNet: end-to-end disparities



DeMoN: end-to-end structure from motion

FlowNet: estimating optical flow with a ConvNet



- Can networks learn to find correspondences?
- New learning task!

(very different from classification, etc.)

Dosovitskiy et al. ICCV 2015



Thomas Brox



\rightarrow Help the network with an explicit correlation layer



Dosovitskiy et al. ICCV 2015

Enough data to train such a network?

- UNI FREIBURG Getting ground truth optical flow for realistic videos is hard
 - Existing datasets are small: •

	Frames with ground truth		
Middlebury	8		
ΚΙΤΤΙ	194		
Sintel	1041		
Needed	>10000		

Realism is overrated: the "flying chairs" dataset



Image pair

Optical flow



Synthetic 3D datasets

Mayer et al. CVPR 2016



Driving, Monkaa, FlyingThings3D datasets publicly available



Generalization: it works!



FlowNetSimple

FlowNetCorr

Although the network has only seen flying chairs for training, it predicts good optical flow on other data

Optical flow estimation in 18ms





Major changes:

- Improved data and training schedules
- Stacking of networks with motion compensation
- Special small displacements and fusion network



FlowNet vs. FlowNet 2.0





	Sintel	KITTI	runtime
DeepFlow (Weinzaepfel et al. 2013)	7.21	5.8	51940 ms
FlowFields (Bailer et al. 2015)	5.81	3.5	22810 ms
PCA Flow (Wulff & Black 2015)	8.65	6.2	140 ms
FlowNet (Dosovitskiy et al. 2015)	7.52	-	18 ms
FlowNet 2.0	5.74	1.8	123 ms

DispNet: disparity estimation



Mayer et al. CVPR 2016





DispNet: disparity estimation









3D shape and texture from a single image



FlowNet: end-to-end optical flow



DispNet: end-to-end disparities



DeMoN: end-to-end structure from motion



DeMoN: Structure from motion with a network



Egomotion estimation and depth estimation are mutually dependent







Estimates optical flow

Estimates depth and egomotion

Iterative refinement



Input images



Ground truth Optical Flow



Estimated optical flow





Ground truth Depth



Estimated depth



Outperforms two-frame SfM baselines

Motion & Pointcloud Comparison

MVS South-Building



DeMoN

Base-Oracle



Pointcloud Comparison

Sculpture



DeMoN

Eigen and Fergus ICCV 2015



Two images generalize better than one image

Pointcloud Comparison

NYU Test 578



DeMoN

Eigen and Fergus ICCV 2015



Structure from motion at 7fps

Bike Ride



Т

T+10

predicted depth





Example from RGB-D SLAM dataset (Sturm et al.) Red: DeMoN. Black: Ground truth.

Deep learning for 3D Vision is promising



3D shape and texture from a single image



FlowNet: end-to-end optical flow



DispNet: end-to-end disparities



DeMoN: end-to-end structure from motion