# DeMoN: Depth and Motion Network for Learning Monocular Stereo

Benjamin Ummenhofer,  Huizhong Zhou,  Jonas Uhrig,  Nikolaus Mayer,  Eddy Ilg, Alexey Dosovitskiy,  Thomas Brox

{ummenhof, zhouh, uhrigj, mayern, ilge, dosovits, brox}@cs.uni-freiburg.de

## Contribution

**Two view geometry problem**
Retrieving the camera motion and scene structure from two images is a fundamental problem in Structure from Motion (SfM).

DeMoN is a ConvNet architecture solving this problem.



**Contributions**
- **A computer algorithm for reconstructing a scene from two projections**
- A network architecture exploiting motion parallax for depth prediction
- An iterative network part for refinement
- A scale invariant gradient loss for improved depth predictions
- Artificial datasets complementing shortcomings of real data

## Depth & Motion Parameterization

**Inverse depth**
Depth uncertainty grows with increasing distance. Thus, we directly estimate the inverse depth (reciprocal of the depth values) to account for this.

- inverse depth  $\xi = \frac{1}{z}$
- can represent points at infinity
- close objects are more important

**Motion**
We present the camera motion from the first to the second frame as:
- 3D translation vector  $\mathbf{t}$
- 3D angle axis vector  $\mathbf{r} = \theta \mathbf{v}$

**Angle axis representation**
- Minimum parameterization
→ Network cannot generate invalid values

**Scale ambiguity**
Scene scale cannot be obtained from images in the **general case**. We resolve the ambiguity by normalizing translations such that

$$\|\mathbf{t}\| = 1$$

Estimated depth values need to correspond to the normalized translation. To facilitate adjusting the depth values we predict a scale factor along with the motion estimate and obtain $s\xi$.

## Project Page

- Paper
- Videos
- Code (Tensorflow)

https://goo.gl/cXf4ct

## Network Architecture

DeMoN consists of three subnets:
- **bootstrap net:**
  computes the initial depth and motion estimates
- **iterative net:**
  successively refines the previous estimates
- **Base-Net:**
  increases the resolution of the final depth map



The bootstrap and iterative net use an encoder-decoder pair:
- **1st encoder-decoder:**
  estimates optical flow and its confidence
- **2nd encoder-decoder:**
  predicts depth and surface normals
- **a fully connected network appended to the 2nd encoder:**
  computes camera motion and a depth scale factor, which relates the scale of the depth values to the camera motion

**Two images are better than one?!**
A single encoder-decoder network does not make use of the second image and prefers to directly infer depth from a single image.

| Method | L1-inv | sc-inv | L1-rel |
|---|---|---|---|
| Single image | 0.080 | 0.159 | 0.696 |
| Naive image pair | 0.079 | 0.165 | 0.722 |
| DeMoN | 0.012 | 0.131 | 0.097 |

**A naive architecture does not use the 2nd image**

DeMoN explicitly solves the more difficult correspondence problem by computing optical flow in the first encoder-decoder.

## Iterative Refinement

The *iterative net* can improve and correct estimates from the *bootstrap net* or from previous iterations.

**Wrong scale**

first image  GT depth  iter 0  iter 1  iter 2  iter 3

**Wrong depth**

first image  GT depth  iter 0  iter 1  iter 2  iter 3

**Iterative refinement on SUN3D**
While we use 4 iterations during training, we find that 3 iterations on average gives the best results with respect to depth *and* motion.

Performance slightly decays with many more iterations (>10) but remains stable.

| Iteration | L1-inv | sc-inv | L1-rel | Depth δ<1.25 | δ<1.25² | δ<1.25³ | Motion rot | tran | Flow EPE |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.029 | 0.145 | 0.244 | 0.587 | 0.844 | 0.940 | 2.18 | 20.27 | 0.030 |
| 1 | 0.024 | 0.130 | 0.207 | 0.679 | 0.891 | 0.961 | 1.94 | 17.25 | 0.020 |
| 2 | 0.022 | 0.131 | 0.187 | 0.688 | 0.900 | 0.982 | 1.87 | 18.31 | 0.019 |
| 3 | 0.021 | 0.132 | 0.179 | 0.698 | 0.912 | 0.981 | 1.83 | 18.81 | 0.019 |
| 4 | 0.021 | 0.133 | 0.184 | 0.690 | 0.908 | 0.975 | 1.79 | 18.94 | 0.019 |
| 5 | 0.021 | 0.133 | 0.185 | 0.692 | 0.910 | 0.970 | 1.79 | 19.65 | 0.019 |

## Scale Invariant Gradient

**Operator Definition**
We define a finite differences operator invariant to scale changes:

$$\mathbf{g}_h[f](i,j) = \left( \frac{f(i+h,j)-f(i,j)}{|f(i+h,j)|+|f(i,j)|}, \frac{f(i,j+h)-f(i,j)}{|f(i,j+h)|+|f(i,j)|} \right)^\top$$

ground truth

$\xi$

prediction

depth  scale invariant gradient images
$\xi$

**Loss Design is Crucial!**
In addition to L1 loss on the depth we compute a loss on the scale invariant gradient images (sig).

The loss on the sig images
- emphasizes importance of depth discontinuities
- stimulates spatial comparisons


L1 loss on depth     + L1 loss on scale invariant gradient images

$$\mathcal{L}_{depth} = \sum_{i,j} |\xi(i,j) - \hat{\xi}(i,j)| \qquad \mathcal{L}_{sig\,\xi} = \sum_{h\in\{1,2,4,8,16\}} \sum_{i,j} \left\| \mathbf{g}_h[\xi](i,j) - \mathbf{g}_h[\hat{\xi}](i,j) \right\|_2$$

## Point Cloud Comparison

Our method produces fewer depth artifacts, which can be seen if we visualize the depth as point clouds.


GT     Eigen[3]     DeMoN

## Training

The network training is based on the caffe[17] framework. It is trained from scratch on all the datasets jointly with Adam[18] using a momentum of 0.9 and a weight decay of 0.0004. We train sequentially the three subnets for 3200k iterations in total. A multistep learning rate policy is applied during training.



## Flow Confidence

We train the flow confidence in a supervised manner.  The ground truth confidence for the x component is given by the flow ground truth $w_x(i,j)$ and the flow prediction $\hat{w}_x(i,j)$:

$$\hat{c}_x(i,j) = e^{-|w_x(i,j) - \hat{w}_x(i,j)|}$$

Flow confidence helps the motion estimation since egomotion only requires sparse but high-quality correspondences.



| | Motion rot | tran | Flow EPE |
|---|---|---|---|
| Confidence | | | |
| no | 2.830 | 25.262 | 0.027 |
| yes | 2.479 | 24.372 | 0.027 |

## Depth Comparison

The depth maps produced by DeMoN are more detailed and more regular than the ones produced by other methods.


Image1  Image2  GT  Base-O  Eigen[3]  Liu[7]  DeMoN
NYU / Sun3D / RGBD / MVS / Scenes11

**Generalization to new data**
DeMoN exploits the geometric relations between a pair of images and therefore generalizes better to unknown scenes for example close-ups of people and objects, images rotated by 90 degrees .


Image  GT  Eigen[3]  Liu[7]  DeMoN

## Datasets

Networks can easily overfit to training data, i.e. training on one dataset is not enough for a method trying to be as general as possible. We train on synthetic and real datasets with complementary properties to improve generalization.


SUN3D [19]  RGBD [14]  Scenes11  MVS  Blendswap

| Dataset | Perfect GT | Photorealistic | Outdoor scenes | Rot. avg | Rot. stddev | Tri. angle avg | Tri. angle stddev |
|---|---|---|---|---|---|---|---|
| SUN3D | no | no | no | 16.6 | 7.5 | 5.2 | 4.6 |
| RGBD | no | no | no | 10.4 | 5.3 | 6.8 | 4.5 |
| Scenes11 | yes | no | no | 3.3 | 2.1 | 13.4 | 6.3 |
| MVS | yes | (yes) | (yes) | 34.3 | 24.7 | 28.9 | 17.5 |
| Blendswap | yes | (yes) | (yes) | 23.1 | 17.1 | 20.1 | 13.6 |

**SUN3D & RGBD**
- Depth from structured light sensor
- Camera pose from SfM (SUN3D) or external tracking (RGBD)

**Scenes11**
- Randomly generated scenes and objects from ShapeNet[2]

**MVS**
- Collection of Multi View Stereo datasets
- Depth and camera poses from SfM pipelines[4,10,11,16]

**Blendswap**
- About 150 distinct scenes from blendswap.com
- Annotated to enable automatic generation of image pairs


scene with annotations     generated images

## Motion

Our method can estimate the camera motion in scenarios difficult for traditional approaches like low texture or small motions.

**Failure cases for Base-FF**



homogeneous region
DeMoN: tran 4.804,  rot 1.237
Base-FF: tran 56.948, rot 2.087

homogeneous region
DeMoN: tran 11.725,  rot 1.628
Base-FF: tran 110.516, rot 15.197

small camera motion
DeMoN: tran 24.096, rot 0.878
Base-FF: tran 71.871, rot 2.564

**Concatenated pairwise motions**
The local pairwise camera poses are consistent with the ground truth.


DeMoN
GT

## Quantitative Comparison

We compare against several traditional methods as well as CNN based image methods.
- DeMoN outperforms all baseline methods on most datasets.
- Besides visual quality, we quantitatively perform as good or better than the single image methods.

| | Correspondences | E-Matrix Method | Depth Method |
|---|---|---|---|
| Base-Oracle | | 8-point algorithm[5] | SGM NCC[6] |
| Base-FF | SIFT[8] | 8-point algorithm[5] | SGM NCC[6] |
| Base-SIFT | Flowfields[1] | 8-point algorithm[5] | SGM NCC[6] |
| Base-Matlab | KLT[12,15] | 5-point algorithm[9] | - |
| Base-Mat-F | DeMoN | 5-point algorithm[9] | - |

| | | Depth | | | Motion | | | Depth |
|---|---|---|---|---|---|---|---|---|
| | Method | L1-inv | sc-inv | L1-rel | rot | trans | Method | sc-inv |
| MVS | Base-Oracle | 0.019 | 0.197 | 0.105 | 0 | 0 | | |
| | Base-SIFT | 0.056 | 0.309 | 0.361 | 21.180 | 60.516 | Liu indoor | 0.260 |
| | Base-FF | 0.055 | 0.308 | 0.322 | 4.834 | 17.252 | Liu outdoor | 0.341 |
| | Base-Mat-F | | | | 10.843 | 32.736 | Eigen VGG | 0.225 |
| | DeMoN | 0.047 | 0.202 | 0.305 | 5.156 | 14.447 | DeMoN | 0.203 |
| Scenes11 | Base-Oracle | 0.023 | 0.618 | 0.349 | 0 | 0 | | |
| | Base-SIFT | 0.051 | 0.900 | 1.027 | 6.179 | 56.650 | Liu indoor | 0.816 |
| | Base-FF | 0.038 | 0.793 | 0.776 | 1.309 | 19.425 | Liu outdoor | 0.814 |
| | Base-Mat-F | | | | 0.917 | 14.639 | Eigen VGG | 0.763 |
| | DeMoN | 0.019 | 0.315 | 0.248 | 0.809 | 8.918 | DeMoN | 0.303 |
| RGB-D | Base-Oracle | 0.026 | 0.398 | 0.336 | 0 | 0 | | |
| | Base-SIFT | 0.050 | 0.577 | 0.703 | 12.010 | 56.021 | | |
| | Base-FF | 0.045 | 0.548 | 0.613 | 4.709 | 46.058 | Liu indoor | 0.338 |
| | Base-Mat-F | | | | 12.831 | 49.612 | Liu outdoor | 0.428 |
| | | | | | 2.917 | 22.523 | Eigen VGG | 0.272 |
| | DeMoN | 0.028 | 0.130 | 0.212 | 2.641 | 20.585 | DeMoN | 0.134 |
| Sun3D | Base-Oracle | 0.020 | 0.241 | 0.220 | 0 | 0 | | |
| | Base-SIFT | 0.029 | 0.290 | 0.286 | 7.702 | 41.825 | | |
| | Base-FF | 0.029 | 0.284 | 0.297 | 3.681 | 33.301 | Liu indoor | 0.214 |
| | Base-Mat-F | | | | 5.920 | 32.298 | Liu outdoor | 0.401 |
| | | | | | 2.230 | 26.338 | Eigen VGG | 0.175 |
| | DeMoN | 0.019 | 0.114 | 0.172 | 1.801 | 18.811 | DeMoN | 0.126 |
| NYU2 | Base-oracle | - | - | - | - | - | | |
| | Base-SIFT | - | - | - | - | - | Liu indoor | 0.210 |
| | Base-FF | - | - | - | - | - | Liu outdoor | 0.421 |
| | Base-Matlab | - | - | - | - | - | Eigen VGG | 0.148 |
| | Base-Mat-F | - | - | - | - | - | DeMoN | 0.180 |
| | DeMoN | - | - | - | - | - | | |

classic methods          single image methods

## References

[References list — small print, not legible]

## Acknowledgements