

## Computer Vision 2 WS 2018/19

### Part 17 – CNNs for Video Analysis II 22.01.2019

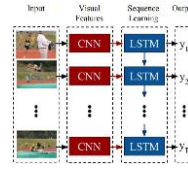
Prof. Dr. Bastian Leibe

RWTH Aachen University, Computer Vision Group  
<http://www.vision.rwth-aachen.de>




### Course Outline

- Single-Object Tracking
- Bayesian Filtering
- Multi-Object Tracking
- Visual Odometry
- Visual SLAM & 3D Reconstruction
  - Online SLAM methods
  - Full SLAM methods
- Deep Learning for Video Analysis
  - CNNs for video analysis
  - CNNs for motion estimation
  - Video object segmentation




2 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis



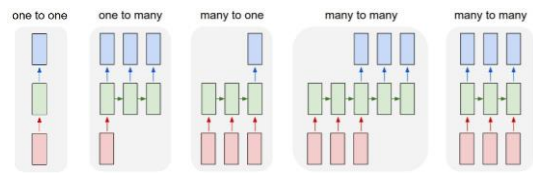
### Topics of This Lecture

- Recap: CNNs for Video Analysis
- Matching and correspondence estimation
  - Metric learning
  - Spatial Transformer Networks
  - Correspondence networks
- Optical Flow Estimation
  - FlowNet
  - FlowNet2

3 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis




### Recap: Recurrent Networks

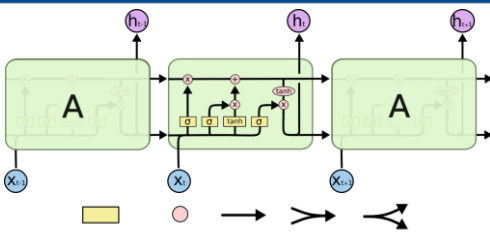


- Feed-forward networks
  - Simple neural network structure: 1-to-1 mapping of inputs to outputs
- Recurrent Neural Networks
  - Generalize this to arbitrary mappings

4 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis




### Recap: Long Short-Term Memory (LSTM)

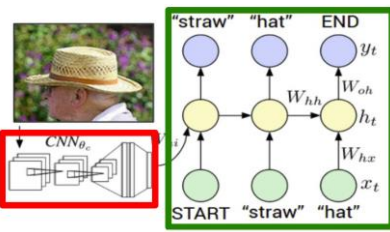


- LSTMs
  - Inspired by the design of memory cells
  - Each module has 4 layers, interacting in a special way.
  - Effect: LSTMs can learn longer dependencies (~100 steps) than RNNs

5 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis




### Recap: Image Tagging



- Simple combination of CNN and RNN
  - Use CNN to define initial state  $h_0$  of an RNN.
  - Use RNN to produce text description of the image.

6 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis



### Recap: Video to Text Description

Our LSTM network is connected to a CNN for RGB frames or a CNN for optical flow images.

Source: Subhrajit Mukherjee, ICCV14

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
RWTH AACHEN UNIVERSITY

### Topics of This Lecture

- Recap: CNNs for Video Analysis
- Matching and correspondence estimation
  - Metric learning
  - Spatial Transformer Networks
  - Correspondence networks
- Optical Flow Estimation
  - FlowNet
  - FlowNet2

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
RWTH AACHEN UNIVERSITY

### Learning Similarity Functions

- Siamese Network
  - Present the two stimuli to two identical copies of a network (with shared parameters)
  - Train them to output similar values if the inputs are (semantically) similar.
- Used for many matching tasks
  - Face identification
  - Stereo estimation
  - Optical flow
  - ...

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
RWTH AACHEN UNIVERSITY

### Metric Learning: Contrastive Loss

- Mapping an image to a metric embedding space
  - Metric space: distance relationship = class membership

$$\|f(x) - f(x_+)\| \rightarrow 0$$

$$\|f(x) - f(x_-)\| \geq m$$

Yi et al., LIFT: Learned Invariant Feature Transform, ECCV 16

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Slide credit: Christopher Chou  
RWTH AACHEN UNIVERSITY

### Metric Learning: Triplet Loss

- Learning a discriminative embedding
  - Present the network with triplets of examples
- Apply triplet loss to learn an embedding  $f(\cdot)$  that groups the positive example closer to the anchor than the negative one.

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
RWTH AACHEN UNIVERSITY

### Patch Normalization with Spatial Transformer Nets

- Patch Normalization
  - Key component of local feature matching
  - Finding the scale and rotation
  - Invariant to perspective transformation
- Spatial Transformer Network
  - Adaptively apply transformation

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Slide credit: Christopher Chou  
ledesma et al., Spatial Transformer Network, NIPS 2015  
RWTH AACHEN UNIVERSITY

### Universal Correspondence Network

- Computing a patch descriptor

Fully Convolutional NN    Convolutional Spatial Transformer    L2-Normalization

13 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Slide credit: Christopher Choy

### Universal Correspondence Network

- Siamese architecture for matching patches

14 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Slide credit: Christopher Choy

### Universal Correspondence Network

- UCN Training

$f(x_+)$   
 $f(x'_+)$   
 $f(x_-)$   
 $f(x'_-)$

$\|f(x_+) - f(x'_+)\| \rightarrow 0$   
 $\|f(x_-) - f(x'_-)\| > m$

- Contrastive loss

15 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Slide credit: Christopher Choy

### Semantic Correspondences with UCN

Ground truth      UCN      VGG Conv4

16 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Slide credit: Christopher Choy

### Exact Correspondences with UCN (Disparity Estimation)

C. Choy, J.Y. Gwak, S. Savarese, M. Chandraker, *Universal Correspondence Network*, NIPS'16

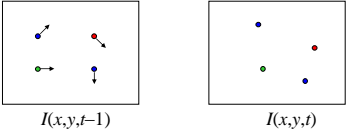
17 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Slide credit: Christopher Choy

### Topics of This Lecture

- Recap: CNNs for Video Analysis
- Matching and correspondence estimation
  - Metric learning
  - Spatial Transformer Networks
  - Correspondence networks
- Optical Flow Estimation
  - FlowNet
  - FlowNet2

18 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis

### Recap: Estimating Optical Flow

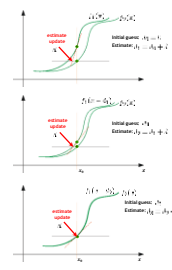


- **Optical Flow**
  - Given two subsequent frames, estimate the apparent motion field  $u(x,y)$  and  $v(x,y)$  between them.
- **Key assumptions**
  - **Brightness constancy:** projection of the same point looks the same in every frame.
  - **Small motion:** points do not move very far.
  - **Spatial coherence:** points move like their neighbors.

19 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Slide credit: Sverdlow, Lapshin

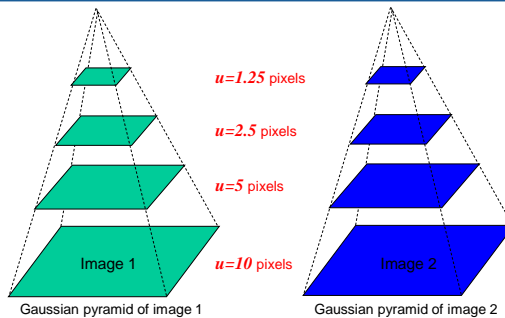
### Recap: Iterative LK Refinement

- Estimate velocity at each pixel using one iteration of LK estimation.
- Warp one image toward the other using the estimated flow field.
- Refine estimate by repeating the process.
- **Iterative procedure**
  - Results in subpixel accurate localization.
  - Converges for small displacements.



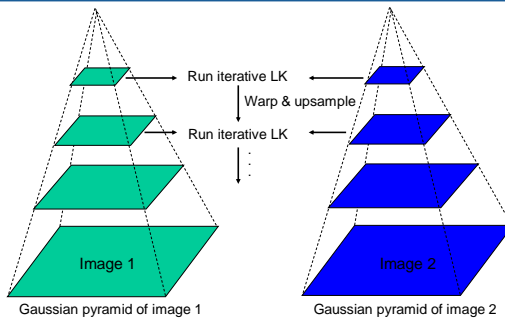
20 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Slide adapted from Steve Saitz

### Recap: Coarse-to-fine Optical Flow Estimation



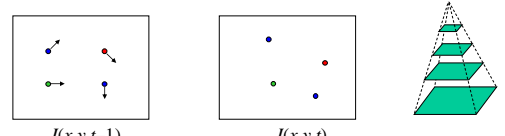
21 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Slide credit: Steve Saitz

### Recap: Coarse-to-fine Optical Flow Estimation



22 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Slide credit: Steve Saitz

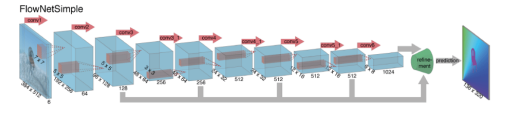
### CNNs for Optical Flow Estimation



- **How can we achieve this with Deep Networks?**
  - Intuition: need to match local image patches
  - CNNs can capture local context, so spatial smoothing should not be necessary
  - But iterative and coarse-to-fine estimation may be necessary.

23 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis

### FlowNet: FlowNetSimple Design



- **Simple initial design**
  - Simply stack two sequential images together and feed them through the network
  - In order to compute flow, the network has to compare image patches
  - But it has to figure out on its own how to do that...

24 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Image source: Fischer et al., ICCV15

### FlowNet: FlowNetCorr Design

- Correlation network
  - Central idea: compute a correlation score between two feature maps

$$c(x_1, x_2) = \sum_{o \in [-k, k] \times [-k, k]} \langle f_1(x_1 + o), f_2(x_2 + o) \rangle$$

- Then refine the correlation scores and turn them into flow predictions

25 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Image source: Fischer et al., ICCV'15

### FlowNet

- Flow refinement stage (both network designs)
  - After series of conv and pooling layers, the resolution has been reduced
  - Refine the coarse pooled representation by upconvolution layers (unpooling + upconvolution)
  - Skip connections to preserve high-res information from early layers

26 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Image source: Fischer et al., ICCV'15

### FlowNet: Training

- Training on FlyingChairs dataset
  - Synthetic dataset with known ground-truth
- Example prediction
  - Both networks can capture fine details

27 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Image source: Fischer et al., ICCV'15

### FlowNet: Comparing the two designs

| Method        | Sintel Clean |      | Sintel Final |      | KITTI |      | Middlebury train |       | Middlebury test |      | Chairs test | Time (sec) |      |
|---------------|--------------|------|--------------|------|-------|------|------------------|-------|-----------------|------|-------------|------------|------|
|               | train        | test | train        | test | train | test | AEE              | AAE   | AEE             | AAE  |             | CPU        | GPU  |
| EpicFlow [30] | 2.40         | 4.12 | 3.70         | 6.29 | 3.47  | 3.8  | 0.31             | 3.21  | 0.39            | 3.55 | 2.94        | 16         | -    |
| DeepFlow [35] | 3.31         | 5.38 | 4.56         | 7.21 | 4.58  | 5.8  | 0.21             | 3.04  | 0.42            | 4.22 | 3.53        | 17         | -    |
| EPPM [3]      | -            | 6.49 | -            | 8.38 | -     | 9.2  | -                | -     | 0.33            | 3.36 | -           | -          | 0.2  |
| LDOF [6]      | 4.29         | 7.56 | 6.42         | 9.12 | 13.73 | 12.4 | 0.45             | 4.97  | 0.56            | 4.55 | 3.47        | 65         | 2.5  |
| FlowNetS      | 4.50         | 7.42 | 5.45         | 8.43 | 8.26  | -    | 1.09             | 13.28 | -               | -    | 2.71        | -          | 0.08 |
| FlowNetS+v    | 3.66         | 6.45 | 4.76         | 7.07 | 6.50  | -    | 0.33             | 3.87  | -               | -    | 2.86        | -          | 1.05 |
| FlowNetS+fl   | (3.66)       | 6.96 | (4.44)       | 7.76 | 7.52  | 9.1  | 0.98             | 15.20 | -               | -    | 3.04        | -          | 0.08 |
| FlowNetS+fl+v | (2.97)       | 6.16 | (4.07)       | 7.22 | 6.07  | 7.6  | 0.32             | 3.84  | 0.47            | 4.58 | 3.03        | -          | 1.05 |
| FlowNetC      | 4.31         | 7.28 | 5.87         | 8.81 | 9.35  | -    | 1.15             | 15.64 | -               | -    | 2.19        | -          | 0.15 |
| FlowNetC+v    | 3.57         | 6.27 | 5.25         | 8.01 | 7.45  | -    | 0.34             | 3.92  | -               | -    | 2.61        | -          | 1.12 |
| FlowNetC+fl   | (3.78)       | 6.85 | (5.28)       | 8.51 | 8.79  | -    | 0.93             | 12.33 | -               | -    | 2.27        | -          | 0.15 |
| FlowNetC+fl+v | (3.20)       | 6.08 | (4.83)       | 7.88 | 7.31  | -    | 0.33             | 3.81  | 0.50            | 4.52 | 2.67        | -          | 1.12 |

- Comparison (avg endpoint errors)
  - Both FlowNetS and FlowNetC can effectively learn to estimate flow
  - FlowNetC overfits to the training data slightly more
  - Finetuning (+fl) and variational refinement (+v) improve results further
  - Performance close to pre-CNN methods, but much faster to compute

28 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Image source: Fischer et al., ICCV'15

### FlowNet: Results

P. Fischer, A. Dosovitskiy, E. Ilg, P. Häusser, C. Hazirbas, V. Golkov  
P. v.d. Smagt, D. Cremers, T. Brox

# FlowNet: Learning Optical Flow with Convolutional Networks

P. Fischer et al., FlowNet: Learning Optical Flow with Convolutional Networks, ICCV 2015.

29 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Image source: Fischer et al., CVPR'12

### FlowNet 2.0: Improved KITTI

- Stacked architecture
  - Several instances of FlowNetC and FlowNetS stacked together to estimate large-displacement flow
  - Sub-network specialized on small motions
  - Fusion layer

30 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 - CNNs for Video Analysis  
Image source: Ilg et al., CVPR'12

### FlowNet 2.0: Detailed View

- Stacked FlowNets
  - Estimates large motion in a coarse-to-fine approach
  - Second image is warped at each level with the intermediate optical flow
  - Intermediate flow and (warped brightness) error are concatenated
  - Difficulty of the learning task is reduced at each level

31 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Image source: Ili et al., CVPR17

### FlowNet 2.0: Detailed View

- Small Displacement Module and Fusion
  - For small displacements, FlowNet2-CSS is not accurate
  - Separate FlowNet2-SD module replaces 5x5 and 7x7 by multiple 3x3 kernels and assumes a stride 1 instead of stride 2 at the first layer
  - Small and simple network to fuse the outputs

32 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Image source: Ili et al., CVPR17

### FlowNet 2.0: Comparison

| Method             | Sintel clean |                   | Sintel final |                   | KITTI 2012 |        | KITTI 2015 |         | Middlebury |      | Runtime      |           |       |
|--------------------|--------------|-------------------|--------------|-------------------|------------|--------|------------|---------|------------|------|--------------|-----------|-------|
|                    | AEE          | Flow              | AEE          | Flow              | AEE        | F-full | AEE        | F-full  | AEE        | Flow | ms per frame | CPU / GPU |       |
| EpFlow [22]        | 2.27         | 4.12              | 3.56         | 6.29              | 3.09       | 3.8    | 9.27       | 27.18%  | 0.31       | 0.30 | 32,600       | -         |       |
| DeepFlow [23]      | 2.60         | 5.38              | 3.57         | 7.21              | 4.48       | 5.8    | 10.03      | 26.32%  | 29.15%     | 0.25 | 0.42         | 51,940    | -     |
| FlowFields [2]     | 1.86         | 3.75              | 3.06         | 5.81              | 3.33       | 3.5    | 8.33       | 24.43%  | -          | 0.27 | 0.33         | 22,810    | -     |
| LDOF (CPU) [7]     | 4.64         | 7.56              | 5.96         | 9.12              | 10.94      | 12.4   | 18.19      | 38.11%  | -          | 0.44 | 0.56         | 64,900    | -     |
| LDOF (GPU) [27]    | 4.76         | -                 | 6.22         | -                 | 10.43      | -      | 18.20      | 38.05%  | -          | 0.36 | -            | -         | 6,270 |
| PCA-Layers [33]    | 3.22         | 5.73              | 4.52         | 7.89              | 5.99       | 5.2    | 12.74      | 27.26%  | -          | 0.66 | -            | 3,300     | -     |
| EPPM [1]           | 6.49         | -                 | 8.38         | -                 | 8.88       | 9.2    | -          | -       | -          | 0.70 | 0.33         | -         | 200   |
| PCA-Flow [13]      | 4.84         | 6.83              | 5.18         | 8.65              | 5.48       | 6.2    | 14.01      | 39.59%  | -          | 0.70 | -            | 140       | -     |
| DIS-Fast [16]      | 5.61         | 9.35              | 6.31         | 10.13             | 11.01      | 11.4   | 21.20      | 53.73%  | -          | 0.92 | -            | 70        | -     |
| FlowNets [11]      | 4.50         | 6.96 <sup>†</sup> | 5.45         | 7.92 <sup>†</sup> | 8.26       | -      | -          | -       | -          | 1.09 | -            | -         | 18    |
| FlowNet2 [11]      | 4.31         | 6.80 <sup>†</sup> | 5.87         | 8.01 <sup>†</sup> | 9.35       | -      | -          | -       | -          | 1.15 | -            | -         | 32    |
| FlowNet2-s         | 4.55         | -                 | 5.21         | -                 | 8.89       | -      | 16.42      | 56.81%  | -          | 1.27 | -            | -         | 7     |
| FlowNet2-ss        | 3.22         | -                 | 3.85         | -                 | 5.45       | -      | 12.84      | 41.09%  | -          | 0.68 | -            | -         | 14    |
| FlowNet2-ss        | 2.51         | -                 | 3.54         | -                 | 4.09       | -      | 11.01      | 35.19%  | -          | 0.54 | -            | -         | 31    |
| FlowNet2-ss-ft-nd  | 2.50         | -                 | 3.50         | -                 | 4.71       | -      | 11.18      | 34.10%  | -          | 0.43 | -            | -         | 31    |
| FlowNet2-CSS       | 2.10         | -                 | 3.23         | -                 | 3.55       | -      | 8.94       | 29.77%  | -          | 0.44 | -            | -         | 69    |
| FlowNet2-CSS-ft-nd | 2.08         | -                 | 3.17         | -                 | 4.05       | -      | 10.07      | 30.73%  | -          | 0.38 | -            | -         | 69    |
| FlowNet2           | 2.02         | 3.96              | 3.14         | 6.02              | 4.09       | -      | 10.06      | 30.37%  | -          | 0.35 | 0.52         | -         | 123   |
| FlowNet2-ft-sintel | (1.45)       | 4.16              | (2.01)       | 5.74              | 3.61       | -      | 9.81       | 28.20%  | -          | 0.35 | -            | -         | 123   |
| FlowNet2-ft-kitti  | 3.43         | -                 | 4.66         | -                 | (1.29)     | 1.8    | (2.30)     | (6.01%) | 11.48%     | 0.56 | -            | -         | 123   |

- Comparison (avg endpoint errors)
  - Similar accuracy as best pre-CNN methods (but much faster)

34 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Image source: Ili et al., CVPR17

### References and Further Reading

- RNNs
  - R. Pascanu, T. Mikolov, Y. Bengio, [On the difficulty of training recurrent neural networks](#), JMLR, Vol. 28, 2013.
  - A. Karpathy, [The Unreasonable Effectiveness of Recurrent Neural Networks](#), blog post, May 2015.
- LSTM
  - S. Hochreiter, J. Schmidhuber, [Long short-term memory](#), Neural Computation, Vol. 9(8): 1735–1780, 1997.
  - A. Graves, [Generating Sequences With Recurrent Neural Networks](#), ArXiv 1308.0850v5, 2014.
  - C. Olah, [Understanding LSTM Networks](#), blog post, August 2015.

35 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Image source: Ili et al., CVPR17

### References and Further Reading

- Optical Flow
  - P. Fischer, A. Dosovitskiy, E. Ilg, P. Haeusser, C. Hazirbas, V. Golkov, P. v.d. Smagd, D. Cremers, T. Brox, [FlowNet: Learning Optical Flow with Convolutional Networks](#), ICCV 2015.
  - E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, T. Brox, [FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks](#), CVPR 2017.
  - A. Ranjan, M.J. Black, [Optical Flow Estimation using a Spatial Pyramid Network](#), CVPR 2017.

36 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 17 – CNNs for Video Analysis  
Image source: Ili et al., CVPR17