**Computer Vision 2
WS 2018/19**

**Part 17 – CNNs for Video Analysis I**
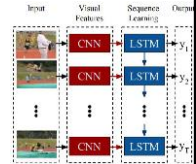**15.01.2019**

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
  – Optical flow
  – Video object segmentation



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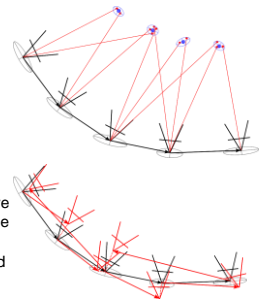
---

## Topics of This Lecture

• Recap: Full SLAM methods

• CNNs for Video Analysis
  – Motivation
  – Example: Video classification

• CNN + RNN
  – RNN, LSTM
  – Example: Video captioning

• Matching and correspondence estimation
  – Metric learning
  – Correspondence networks

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

---

## Recap: Full SLAM Approaches

• SLAM graph optimization:
  – Joint optimization for poses and map elements from image observations of map elements and control inputs

• Pose graph optimization:
  – Optimization of poses from relative pose constraints deduced from the image observations
  – Map recovered from the optimized poses

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

---

## Pose Graph Optimization

• Optimization of poses
  – From relative pose constraints deduced from the image observations
  – Map recovered from the optimized poses
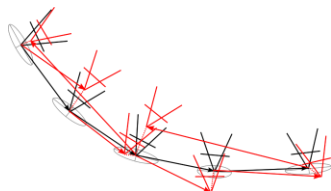
• Deduce relative constraints between poses from image observations, e.g.,
  – 8-point algorithm
  – Direct image alignment

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

---

## Pose Graph Optimization Example

**Dense Visual SLAM
for RGB-D Cameras**

Christian Kerl, Jürgen Sturm,
Daniel Cremers

Computer Vision and Pattern Recognition Group
Department of Computer Science
Technical University of Munich

Kerl et al., Dense Visual SLAM for RGB-D Cameras, IROS 2013

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

## Topics of This Lecture

- Recap: Full SLAM methods

- CNNs for Video Analysis
  – Motivation
  – Example: Video classification

- CNN + RNN
  – RNN, LSTM
  – Example: Video captioning

- Matching and correspondence estimation
  – Metric learning
  – Correspondence networks

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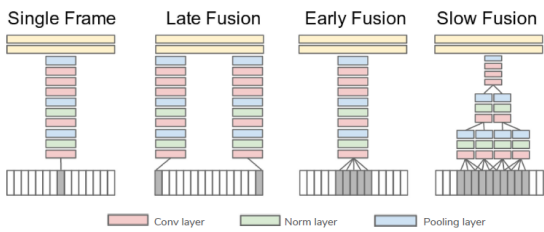Part 17 – CNNs for Video Analysis

---

## Video Analysis with CNNs



- Modeling perspective
  – What architecture to use to best capture temporal patterns?

- Computational perspective
  – Video processing is expensive!
  – How to reduce computation cost without sacrificing accuracy

8 | Visual Computing Institute | Prof. Dr . Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Slide credit: Fei-Fei Li

---

## Large-Scale Video Classification with CNNs

- Architecture
  – Different ways to fuse features from multiple frames

Single Frame    Late Fusion    Early Fusion    Slow Fusion



Conv layer    Norm layer    Pooling layer

9 | Visual Computing Institute | Prof. Dr . Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Slide credit: Fei-Fei Li
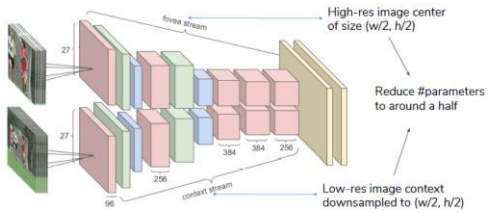
---

## Large-Scale Video Classification with CNNs

- Computational cost
  – Reduce spatial dimension to reduce model complexity
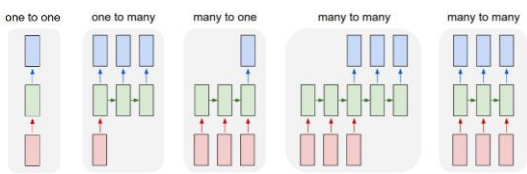  – Multi-resolution: low-res context + high-res foveate



High-res image center of size (w/2, h/2)

Reduce #parameters to around a half

Low-res image context downsampled to (w/2, h/2)

10 | Visual Computing Institute | Prof. Dr . Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Image source: Andrej Karpathy

---

## Topics of This Lecture

- Recap: Full SLAM methods

- CNNs for Video Analysis
  – Motivation
  – Example: Video classification

- CNN + RNN
  – RNN, LSTM
  – Example: Video captioning

- Matching and correspondence estimation
  – Metric learning
  – Correspondence networks

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

---

## Recap: Recurrent Networks

one to one    one to many    many to one    many to many    many to many



- Feed-forward networks
  – Simple neural network structure: 1-to-1 mapping of inputs to outputs

- Recurrent Neural Networks
  – Generalize this to arbitrary mappings

12 | Visual Computing Institute | Prof. Dr . Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Image source: Andrej Karpathy

## Recap: RNNs

- RNNs are regular NNs whose hidden units have additional forward connections over time.
  - You can unroll them to create a network that extends over time.
  - When you do this, keep in mind that the weights for the hidden units are shared between temporal layers.



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

Image source: Andrei Karpathy

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



| | | | | |
|---|---|---|---|---|
| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

- 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

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

Image source: Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/

## Recap: RNNs for Text Generation

- RNN for text generation



10,001D class scores (Softmax over 10k words and a special <END> token)

$$\mathbf{y}_4 = \mathbf{W}_{hy}\mathbf{h}_4$$

Hidden layer (e.g., 500D vectors)

$$\mathbf{h}_4 = \max\{0, \mathbf{W}_{xh}\mathbf{x}_4 + \mathbf{W}_{hh}\mathbf{h}_3\}$$

**15** Computer Vision 2
Part 17 – CNNs for Video Analysis

Slide credit: Andrei Karpathy, Fei-Fei Li
Image source: Andrei Karpathy

## Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$p(next\ word\ |$
$\quad previous\ words)$



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

Slide credit: Andrei Karpathy, Fei-Fei Li

## Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$p(next\ word\ |$
$\quad previous\ words)$



sample!

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

Slide credit: Andrei Karpathy, Fei-Fei Li

## Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$p(next\ word\ |$
$\quad previous\ words)$



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

Slide credit: Andrei Karpathy, Fei-Fei Li

Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$$p(next\ word\ |\ previous\ words)$$

sample!

---

Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$$p(next\ word\ |\ previous\ words)$$

---

Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$$p(next\ word\ |\ previous\ words)$$

sample

---

Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$$p(next\ word\ |\ previous\ words)$$

---

Recap: RNNs for Text Generation

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$$p(next\ word\ |\ previous\ words)$$

sample!

---

Recap: RNNs for Text Generation

samples <END>? Done!

- Training this on a lot of sentences would give us a language model.

- I.e., a way to predict

$$p(next\ word\ |\ previous\ words)$$

## Applications: Image Tagging



- Simple combination of CNN and RNN
  – Use CNN to define initial state $\mathbf{h}_0$ of an RNN.
  – Use RNN to produce text description of the image.

25 | Visual Computing Institute | Prof. Dr. Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Slide adapted from Andrej Karpathy

---

## Applications: Image Tagging

- Setup
  – Train on corpus of images with textual descriptions
  – E.g. Microsoft CoCo
    - 120k images
    - 5 sentences each



26 | Visual Computing Institute | Prof. Dr. Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Slide adapted from Andrej Karpathy

---

## Results: Image Tagging



a group of people standing around a room with remotes
logprob: -9.17

a young boy is holding a baseball bat
logprob: -7.61

a cow is standing in the middle of a street
logprob: -8.84

*Spectacular results!*

27 | Visual Computing Institute | Prof. Dr. Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Slide adapted from Andrej Karpathy

---

## Results: Image Tagging



a baby laying on a bed with a stuffed bear
logprob: -8.66

a young boy is holding a baseball bat
logprob: -7.65

a cat is sitting on a couch with a remote control
logprob: -12.45

- Wrong, but one can still see why those results were selected...

28 | Visual Computing Institute | Prof. Dr. Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Slide adapted from Andrej Karpathy

---

## Application: Video to Text Description



29 | Visual Computing Institute | Prof. Dr. Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Source: Subhashini Venugopalan, ICCV15

---

## Video-to-Text Results



30 | Visual Computing Institute | Prof. Dr. Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
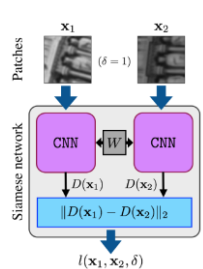Source: Subhashini Venugopalan, ICCV15

## Topics of This Lecture

- Recap: Full SLAM methods
- CNNs for Video Analysis
  - Motivation
  - Example: Video classification
- CNN + RNN
  - RNN, LSTM
  - Example: Video captioning
- Matching and correspondence estimation
  - Metric learning
  - Correspondence networks

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

---

## 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
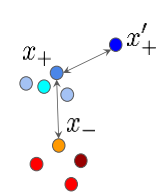  - …



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

---

## Metric Learning: Contrastive Loss

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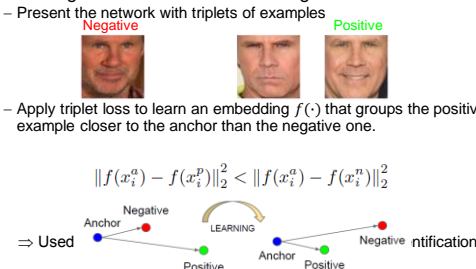  - Metric space: distance relationship = class membership

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

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



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

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

---

## 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.

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$$
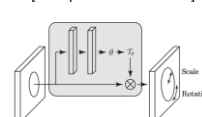
$\Rightarrow$ Used ... ntification

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

---

## Patch Normalization with Spatial Transformer Nets

- Patch Normalization
  - Key component of local feature matching
  - Finding the scale and rotation
  - Invariant to perspective transformation

[SIFT patch normalization]

- Spatial Transformer Network
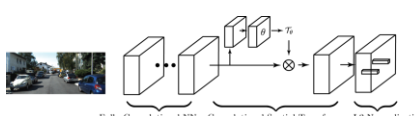  - Adaptively apply transfomation

[Spatial Transformer Network]

**35** Visual Computing Institute | Prof. Dr . Bastian Leibe
Computer Vision 2
Part 17 – CNNs for Video Analysis
Slide credit: Christopher Choy    Jaderberg et al., Spatial Transformer Network, NIPS 2015

---

## Universal Correspondence Network

- Computing a patch descriptor



Fully Convolutional NN    Convolutional Spatial Transformer    L2-Normalization

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

16.01.2019

---

## Universal Correspondence Network

- Siamese architecture for matching patches



**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



- Contrastive loss
$$\|f(x_+) - f(x'_+)\| \to 0$$
$$\|f(x_-) - f(x'_-)\| > m$$

**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

**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

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

---

## 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.

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