# 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

**Visual Computing Institute**

**RWTH**AACHEN
UNIVERSITY

# 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

# 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

# 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

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

- ## 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

Slide credit: Fei-Fei Li

# Large-Scale Video Classification with CNNs

- Architecture
  - Different ways to fuse features from multiple frames

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

Image source: Andrej Karpathy

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

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

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

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: Andrej Karpathy

# Extension: 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

**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/

- 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)

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

Slide credit: Andrej Karpathy, Fei-Fei Li

Image source: Andrej 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

$$p(\textcolor{red}{next\ word}\ |\ \textcolor{blue}{previous\ words})$$

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

Slide credit: Andrej 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!

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

$$p(\textcolor{red}{next\ word}\ |\ \textcolor{blue}{previous\ words})$$

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

Slide credit: Andrej 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!

**19**

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

Slide credit: Andrej Karpathy, Fei-Fei Li

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

- I.e., a way to predict

$$p(\textcolor{red}{next\ word}\ |\ \textcolor{blue}{previous\ words})$$

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

Slide credit: Andrej 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

$$p(\textcolor{red}{next\ word}\ |\ \textcolor{blue}{previous\ words})$$



sample!

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

Slide credit: Andrej 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

$p(\textcolor{red}{next\ word}\ |$
$\textcolor{blue}{previous\ words})$

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

$$p(\textcolor{red}{next\ word}\ |\ \textcolor{blue}{previous\ words})$$



sample!

y0  y1  y2  y3

h0 → h1 → h2 → h3

x0 <START>  x1 "cat"  x2 "sat"  x3 "on"  x4 "mat"

Slide credit: Andrej Karpathy, Fei-Fei Li

# 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

$$p(\textcolor{red}{next\ word}\ |\ \textcolor{blue}{previous\ words})$$

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

# 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.

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



a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.

Slide adapted from Andrej Karpathy

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!*

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

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...

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

Source: Subhashini Venugopalan, ICCV'15

**Correct descriptions.**

S2VT: A man is doing stunts on his bike.

2VT: A herd of zebras are walking in a field.

S2VT: A young woman is doing her hair.

S2VT: A man is shooting a gun at a target.

**Relevant but incorrect descriptions.**

S2VT: A small bus is running into a building.

S2VT: A man is cutting a piece of a pair of a paper.

S2VT: A cat is trying to get a small board.

S2VT: A man is spreading butter on a tortilla.

**Irrelevant descriptions.**

S2VT: A man is pouring liquid in a pan.

S2VT: A polar bear is walking on a hill.

S2VT: A man is doing a pencil.

S2VT: A black clip to walking through a path.

**30**

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

**Visual Computing Institute**

**RWTH AACHEN UNIVERSITY**

Source: Subhashini Venugopalan, ICCV'15

# 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

# 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
  - …

# 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

# Metric Learning: Triplet Loss

- Learning a discriminative embedding
  - Present the network with triplets of examples

    Negative                                    Positive

  - Apply triplet loss to learn an embedding $f(\cdot)$ that groups the positive example closer to the anchor than the negative one.

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$$

Anchor     Negative          LEARNING          Negative

Positive                     Anchor    Positive

$\Rightarrow$ Used                                     ntification

# 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]

Jaderberg et al., Spatial Transformer Network, NIPS 2015

- Computing a patch descriptor



Fully Convolutional NN    Convolutional Spatial Transformer    L2-Normalization

Slide credit: Christopher Choy

# Universal Correspondence Network

- Siamese architecture for matching patches

Slide credit: Christopher Choy

# Universal Correspondence Network

- ## UCN Training



- ## Contrastive loss

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

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

# Semantic Correspondences with UCN



Ground truth                    UCN                    VGG Conv4

Slide credit: Christopher Choy

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