

Lab CudaVision

Learning Vision Systems on Graphics Cards (MA-INF 4308)

# Introduction Session

---

17.10.2023

PROF. SVEN BEHNKE, ANGEL VILLAR-CORRALES

Contact: [villar@ais.uni-bonn.de](mailto:villar@ais.uni-bonn.de)

# About Me

---

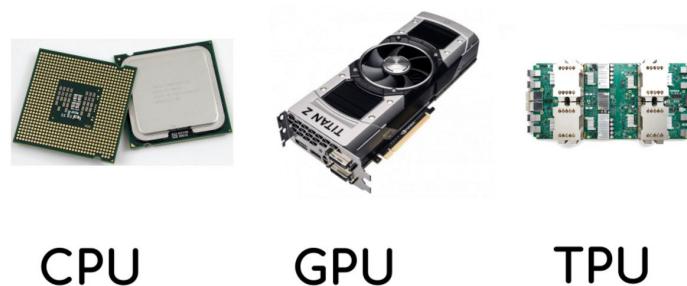
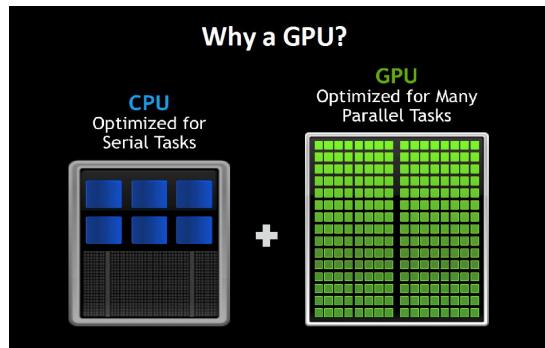
- Member of research staff at AIS since 02.2021
- Before:
  - M. Sc. at FAU Erlangen-Nürnberg
  - B. Sc. at University of Vigo (Spain)
- Research interests:
  - Object-Centric Learning
  - Video Prediction
  - Computer Vision and Deep Learning

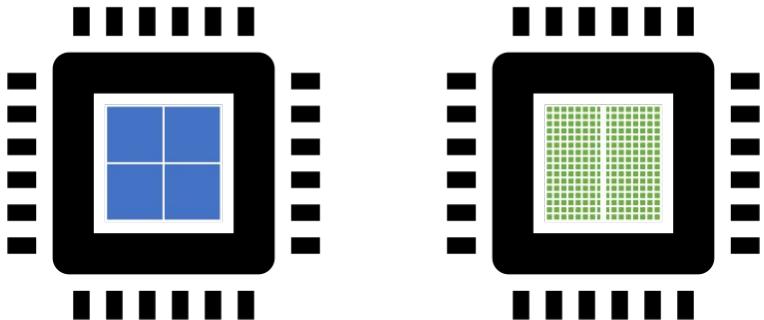


# Motivation

# Why Image Processing on GPUs?

- Image processing and analysis algorithms are inherently parallel:
  - Convolutions
  - Filtering
- Advancements on parallel computing devices: GPUs or TPUs
- Availability of programming interfaces: CUDA, OpenCL, ...



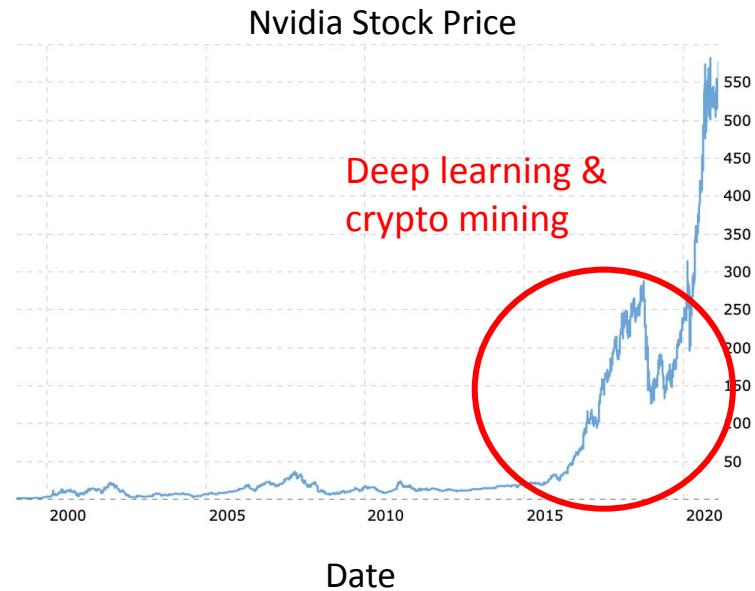
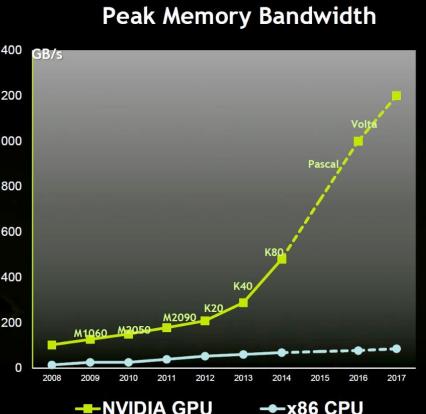
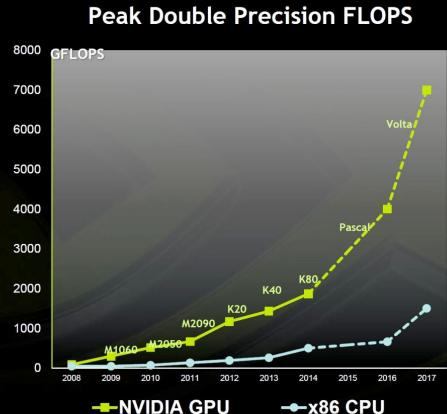


CPU	GPU
Central Processing Unit	Graphics Processing Unit
4-8 Cores	100s or 1000s of Cores
Low Latency	High Throughput
Good for Serial Processing	Good for Parallel Processing
Quickly Process Tasks That Require Interactivity	Breaks Jobs Into Separate Tasks To Process Simultaneously
Traditional Programming Are Written For CPU Sequential Execution	Requires Additional Software To Convert CPU Functions to GPU Functions for Parallel Execution

<https://towardsdatascience.com/parallel-computing-upgrade-your-data-science-with-a-gpu-bba1cc007c24>

# CPU vs. GPU Performance

## GPU Motivation (I): Performance Trends



7

# ImageNet Challenge

---

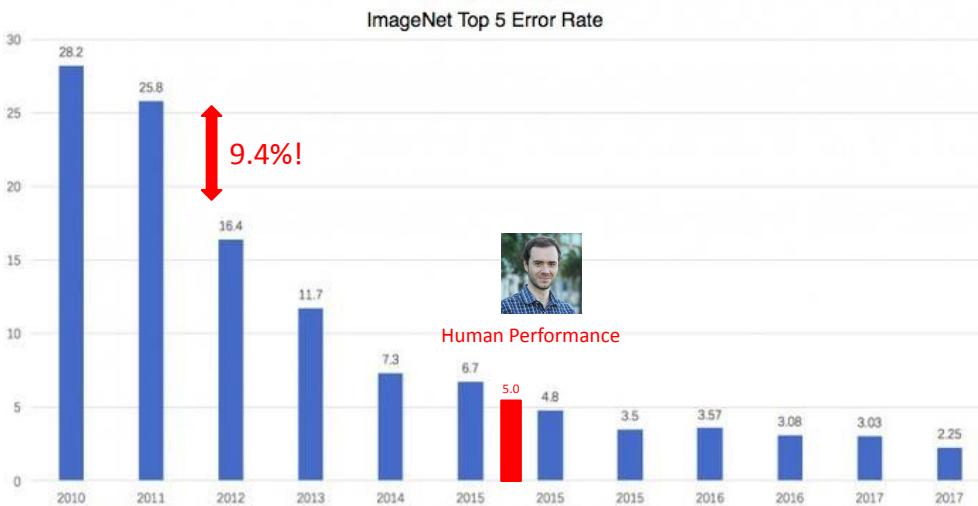
- $\approx$ 14 million natural images labelled into  $\approx$ 1000 classes
- **2012:** Deep learning breakthrough by Krizhevsky et al.

**IMAGENET**

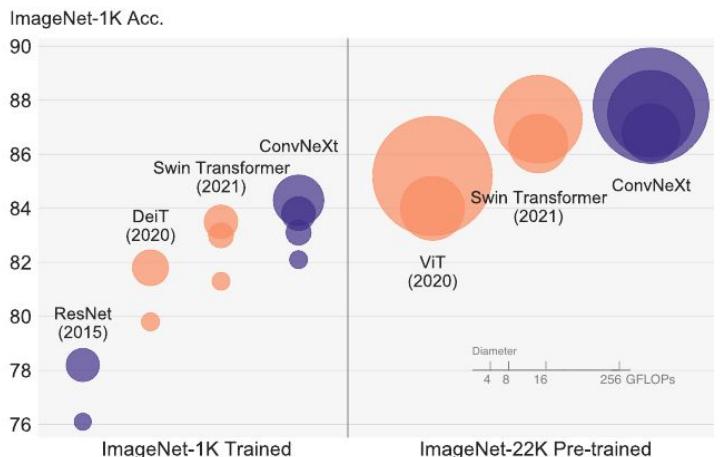


[Krizhevsky et al. NeurIPS 20212]

# ImageNet SOTA

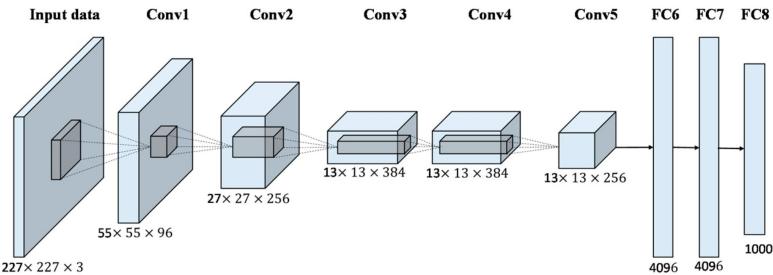


First CNN

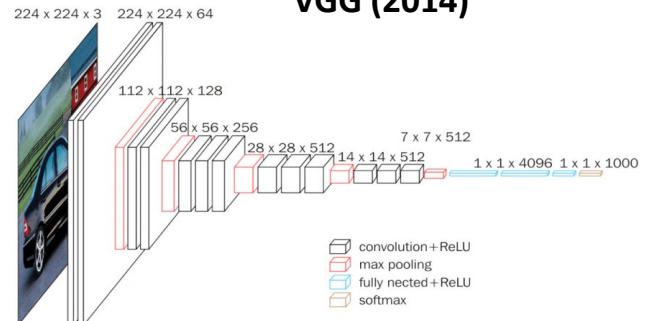


# Architectures

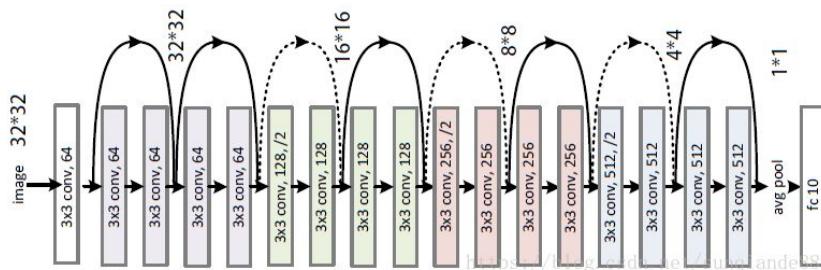
**AlexNet (2012)**



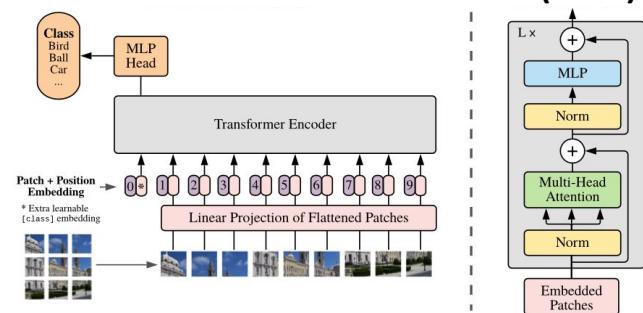
**VGG (2014)**



**ResNet (2015)**

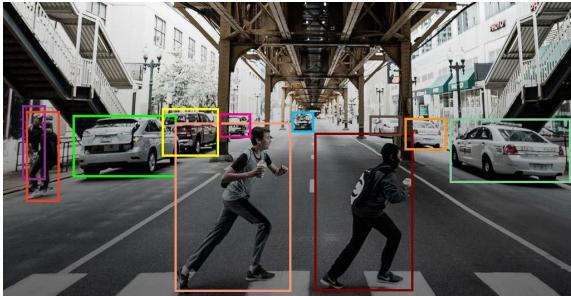


**Vision Transformers (2020)**

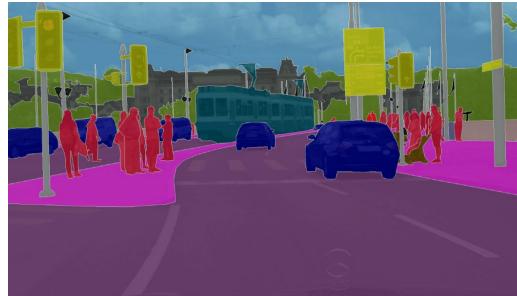


# Applications in Vision

Object Detection



Semantic Segmentation



Human Pose Estimation

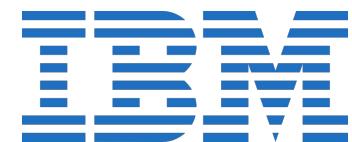


Image Synthesis & Manipulation



# Deep Learning Users

---



SIEMENS



# In this Lab...

# In this Lab...

---

- Implementing deep learning algorithms for visual pattern recognition
  - Python programming language
  - PyTorch framework
  - Deep learning basics
- 8 Lab sessions (30%) in 13 weeks
- Final project (70%) in lecture free period
  - Code/Results
  - Technical Report (8-10 pages)

# Organization

---

- Meeting time: (Bi)Weekly 2 hours meeting, in-person
  - Discuss solution to previous assignment
  - Review some theory
  - Run sample code
  - Provide next assignment
  - Questions
- Accessible GPUs (Informatik ID):
  - Room 0.057 accessible during working hours
  - 8 GPU computers (cuda3 - cuda12) with RTX 2080/3090/Titan
  - Free online resources (Google Colab/ Kaggle kernels)



# Assignments

---

- Each session covers one topic
- Take-home assignment
  - Similar to what we do during the session
  - Due shortly before follow up session
- Assignments & project can be done in pairs
  - Highly recommended!
- Send me a mail with the name of your partner or tell me if you need one

# Topics Covered

---

1. PyTorch basics, Autograd and Fully Connected Neural Networks
2. Optimization & Convolutional Neural Networks (CNNs)
3. Popular Architectures and Transfer Learning
4. Recurrent Neural Networks (RNNs & LSTM)
5. Autoencoders (AEs), Denoising and Variational AEs
6. Generative Adversarial Networks (GANs)
7. Deep Metric and Similarity Learning
8. Semantic Segmentation and Introduction to Final Project



**May be updated to fit  
final project**

# Registration (more on this later)

---

- 16 slots assigned at random
- Please fill this form before 18.10 (select the time slots that work for you)  
<https://forms.gle/AZ7CiXx3SeMV8wAu8>
- Registration (contact me first):
  - **Uni Bonn:** in BASIS
  - **Bonn-Rhein-Sieg & Others:** Contact me

	Monday	Tuesday	Wednesday	Thursday	Friday
<b>9-11 or 10-12</b>					
<b>11-13 or 12-14</b>					
<b>13-15 or 14-16</b>					
<b>15-17</b>					

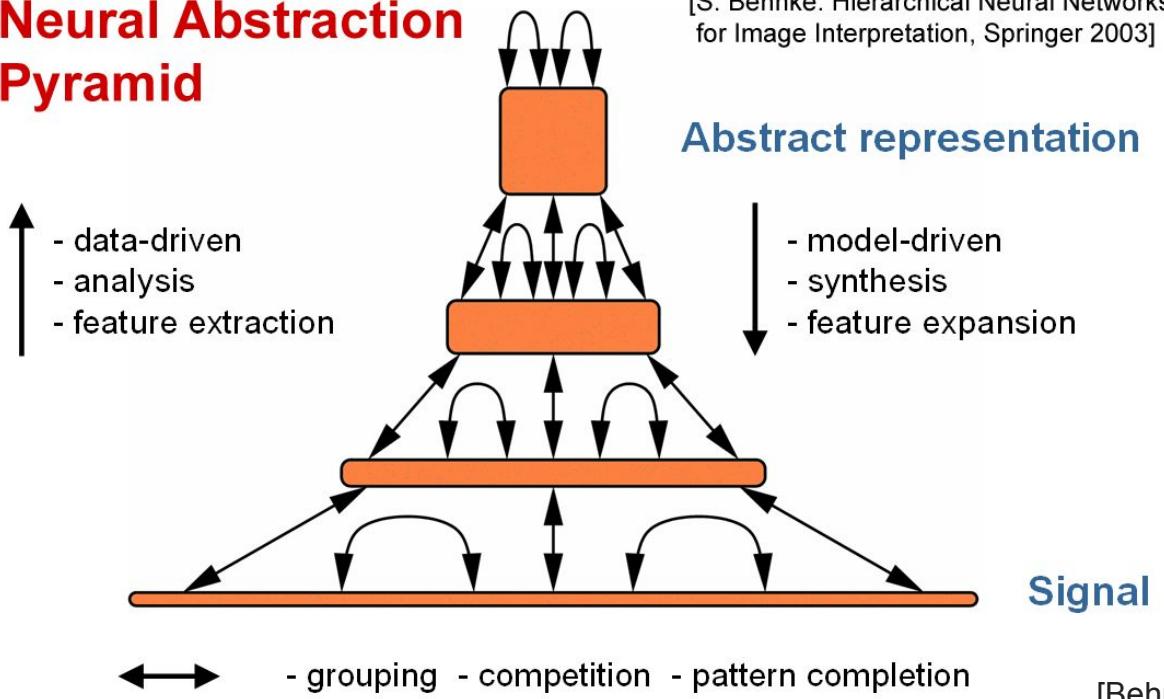
# Some of our Research...

# Neural Abstraction Pyramid

## Neural Abstraction Pyramid

[S. Behnke: Hierarchical Neural Networks for Image Interpretation, Springer 2003]

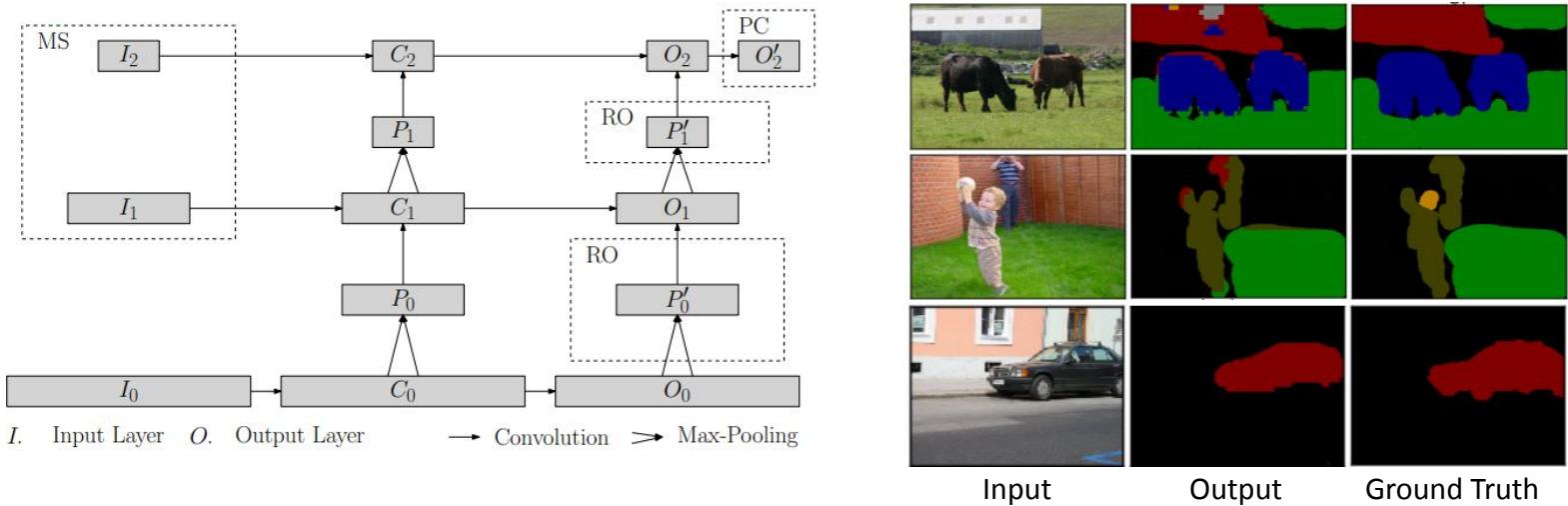
Abstract representation



[Behnke. IJCNN 1998]

# Object-class Segmentation

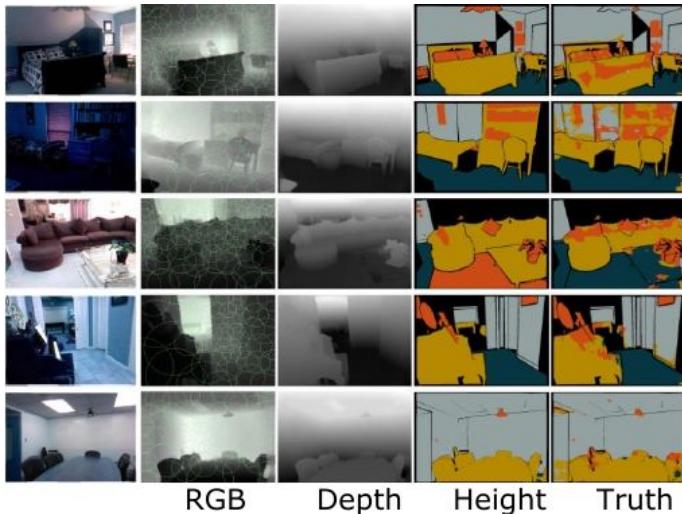
- Multi-scale CNN for RGB-pixel segmentation



[Schulz and Behnke. ESANN 2012]

# RGB-D Object Segmentation

- Use of kinect-like sensors to obtain depth values



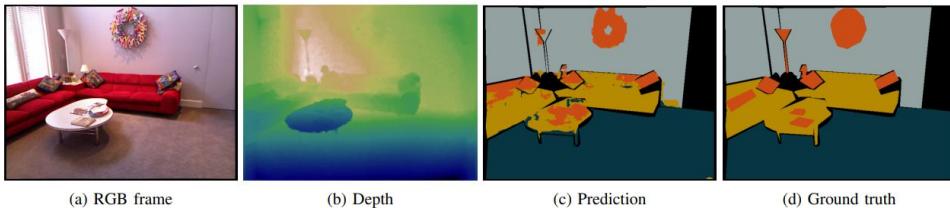
Method	floor	struct	furnit	prop	Class Avg.	Pixel Acc.
CW	84.6	70.3	58.7	52.9	66.6	65.4
CW+DN	87.7	70.8	57.0	53.6	67.3	65.5
CW+H	78.4	74.5	55.6	62.7	67.8	66.5
CW+DN+H	93.7	72.5	61.7	55.5	70.9	70.5
CW+DN+H+SP	91.8	74.1	59.4	63.4	72.2	71.9
CW+DN+H+CRF	93.5	80.2	66.4	54.9	<b>73.7</b>	<b>73.4</b>
Müller et al.[8]	94.9	78.9	71.1	42.7	71.9	72.3
Random Forest [8]	90.8	81.6	67.9	19.9	65.1	68.3
Couprise et al.[9]	87.3	86.1	45.3	35.5	63.6	64.5
Höft et al.[10]	77.9	65.4	55.9	49.9	62.3	62.0
Silberman [12]	68	59	70	42	59.7	58.6

CW is covering windows, H is height above ground, DN is depth normalized patch sizes. SP averaged within superpixels and SVM-reweighted. CRF is a conditional random field over superpixels [8]. Structure class numbers are optimized for class accuracy.

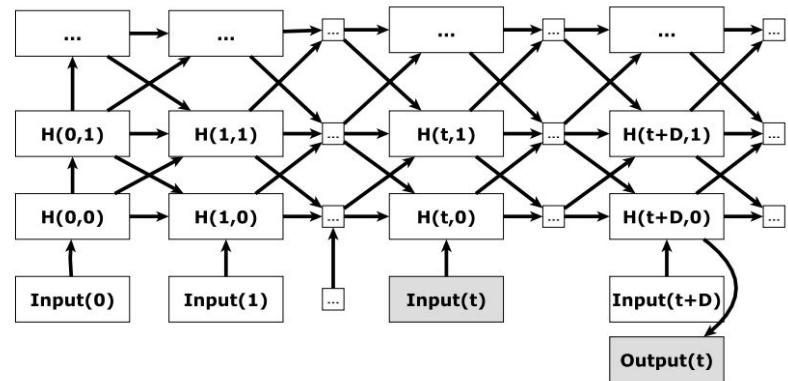
[Schulz, Höft and Behnke. ESANN 2015]

# Object Segmentation from RGB-D Video

- Video processing with multi-scale Convolutional RNNs
- Iterative refinement through different time steps



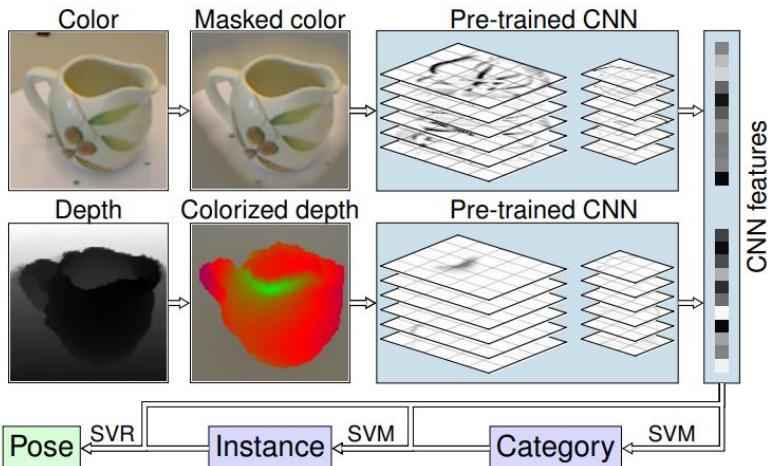
Method	Class Accuracies (%)				Average (%)	
	ground	struct	furnit	prop	Class	Pixel
Unidirectional + SW	90.0	76.3	52.1	<b>61.2</b>	69.9	67.5
Schulz <i>et al.</i> [20]	93.6	80.2	66.4	54.9	<b>73.7</b>	<b>73.4</b>
Müller and Behnke [22]	<b>94.9</b>	78.9	<b>79.7</b>	55.1	71.9	72.3
Stückler <i>et al.</i> [21]	90.8	81.6	67.9	19.9	65.0	68.3
Couprie <i>et al.</i> [23]	87.3	<b>86.1</b>	45.3	35.5	63.5	64.5
Höft <i>et al.</i> [19]	77.9	65.4	55.9	49.9	61.1	62.0
Silberman <i>et al.</i> [17]	68	59	70	42	59.6	58.6



[Pavel, Schulz, and Behnke. IJCNN 2015, Neural Networks 2017]

# Computer Vision with Pretrained Features

- Object recognition and pose estimation
- Pretrained features from ImageNet
- Improved classification and estimation performance



Evaluation on the Washington RGB-D Objects dataset

Method	Category Accuracy (%)		Instance Accuracy (%)	
	RGB	RGB-D	RGB	RGB-D
Lai <i>et al.</i> [12]	$74.3 \pm 3.3$	$81.9 \pm 2.8$	59.3	73.9
Bo <i>et al.</i> [14]	$82.4 \pm 3.1$	$87.5 \pm 2.9$	<b>92.1</b>	92.8
PHOW[18]	$80.2 \pm 1.8$	—	62.8	—
<b>Ours</b>	<b><math>83.1 \pm 2.0</math></b>	<b><math>89.4 \pm 1.3</math></b>	92.0	<b>94.1</b>

[Schwarz, Schulz and Behnke. ICRA 2015]

# Amazon Bin-Picking Challenge

- Picking a large variety of objects
- Placing them on a shelf or packing boxes
- NimbRo team came in 2nd
- Computer vision challenge



[Schwarz et al. ICRA 2017]

# Object Capture and Scene Synthesys

- Capture data with a turn table
- Rendered realistic scenes



# Object Detection and Segmentation

---

- RefineNet architecture [1]
- Trained on rendered data



[1] Lin et al. CVPR 2016

# Soccer Robots

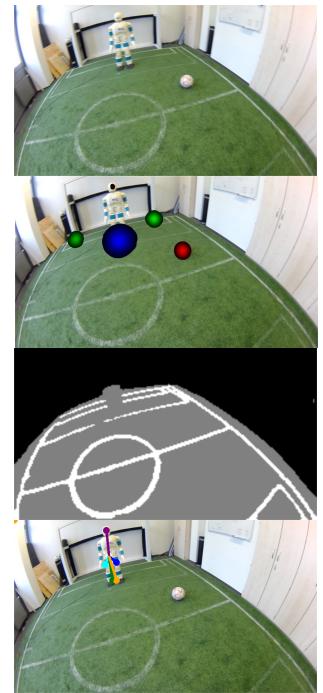
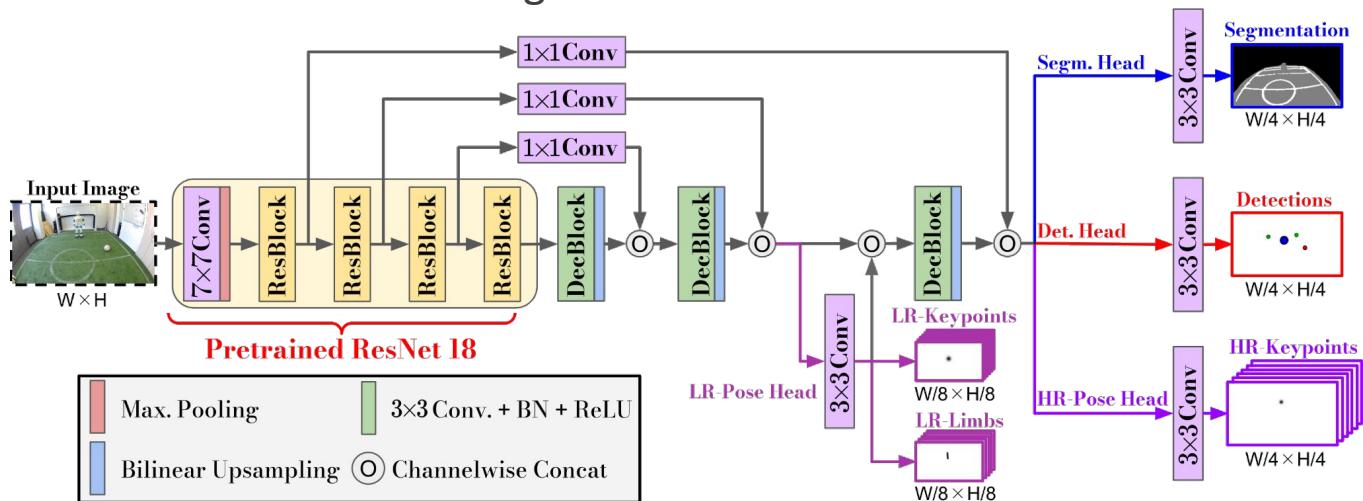
---

- NimbRo participates in humanoid soccer robot competitions (RoboCup)
- Challenging perception scenario



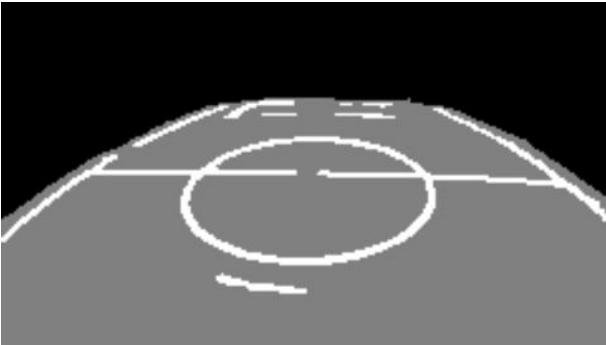
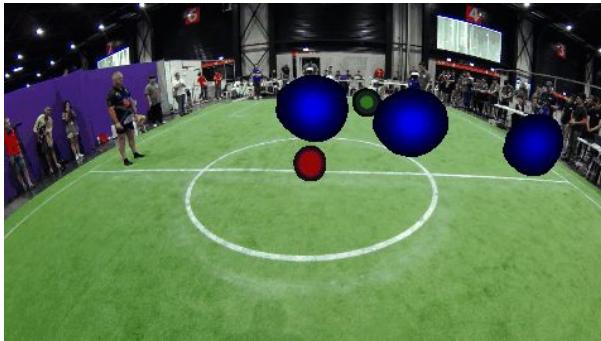
# Scene Understanding

- End-to-end convolutional model
- Robot and ball detection
- Soccer field segmentation



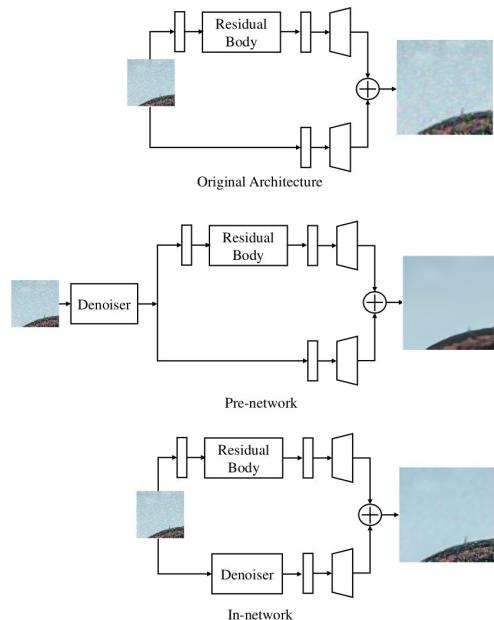
[Villar-Corrales et al. RoboCup 2023]

# Scene Understanding



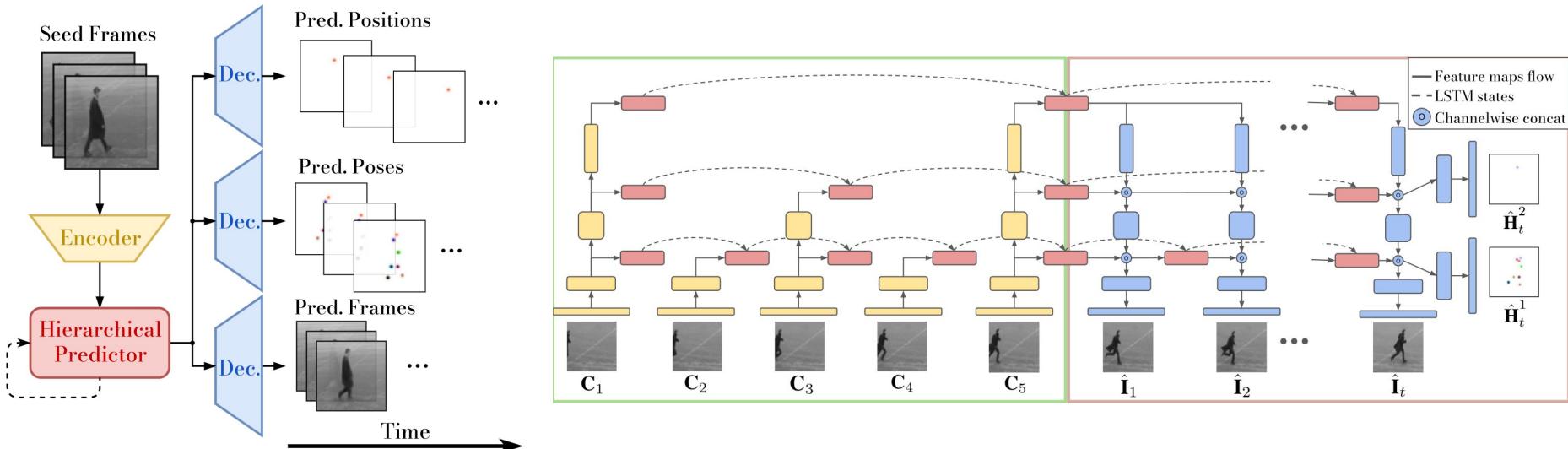
# Denoising and Super-Resolution

- Architectural designs for joined denoising and super-resolution



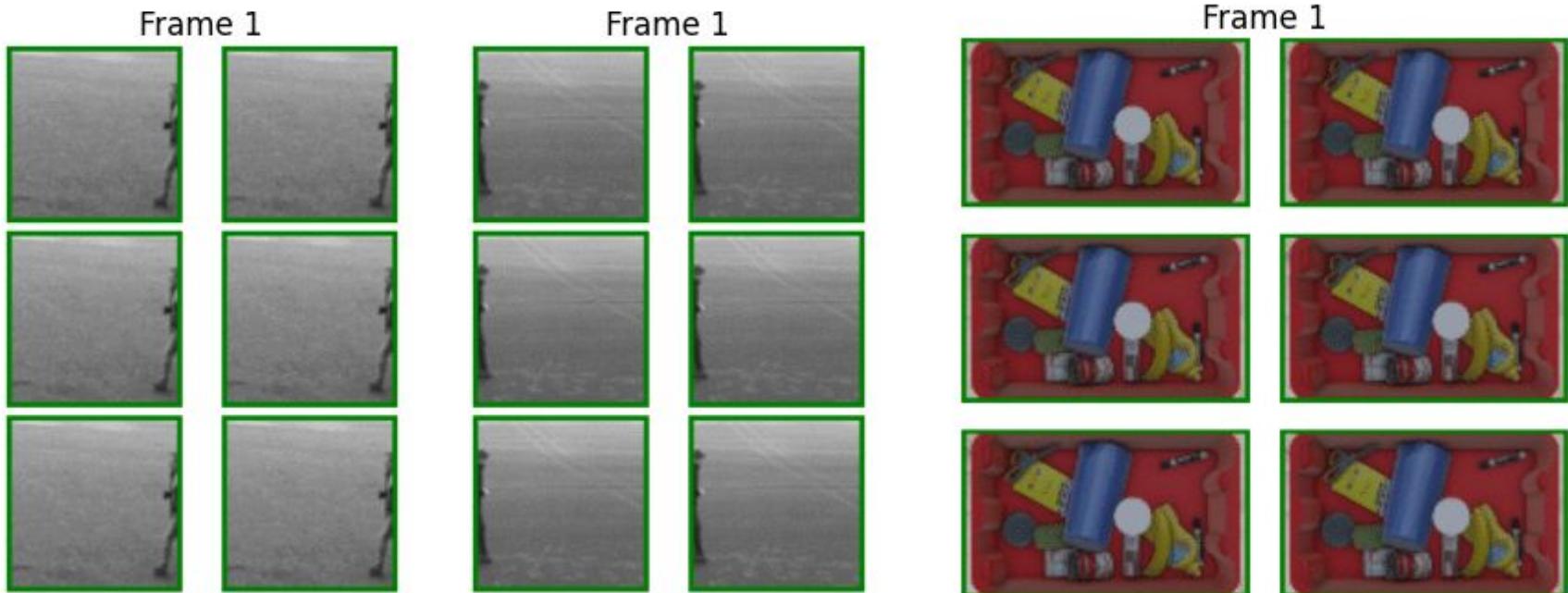
[Villar-Corrales et al. ICASSP 2021]

# Future Frame Video Prediction



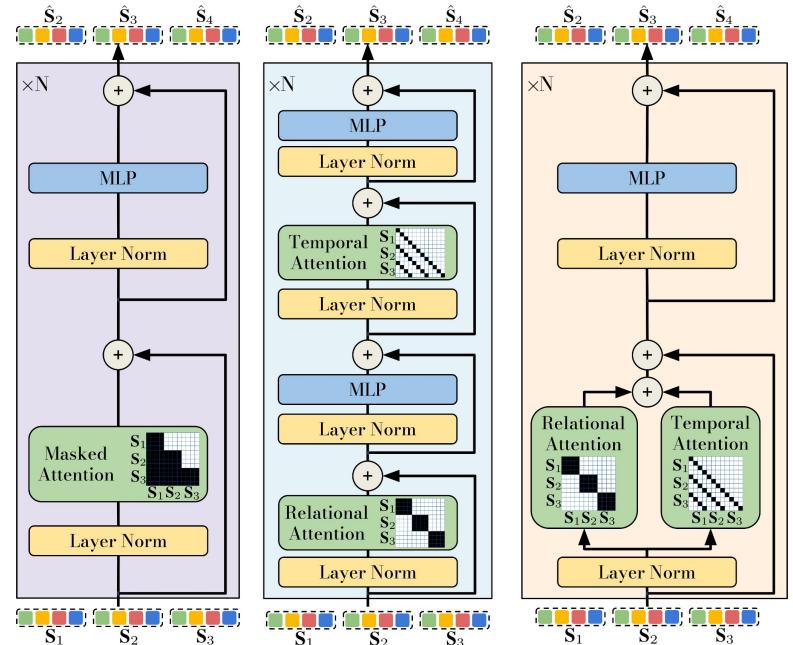
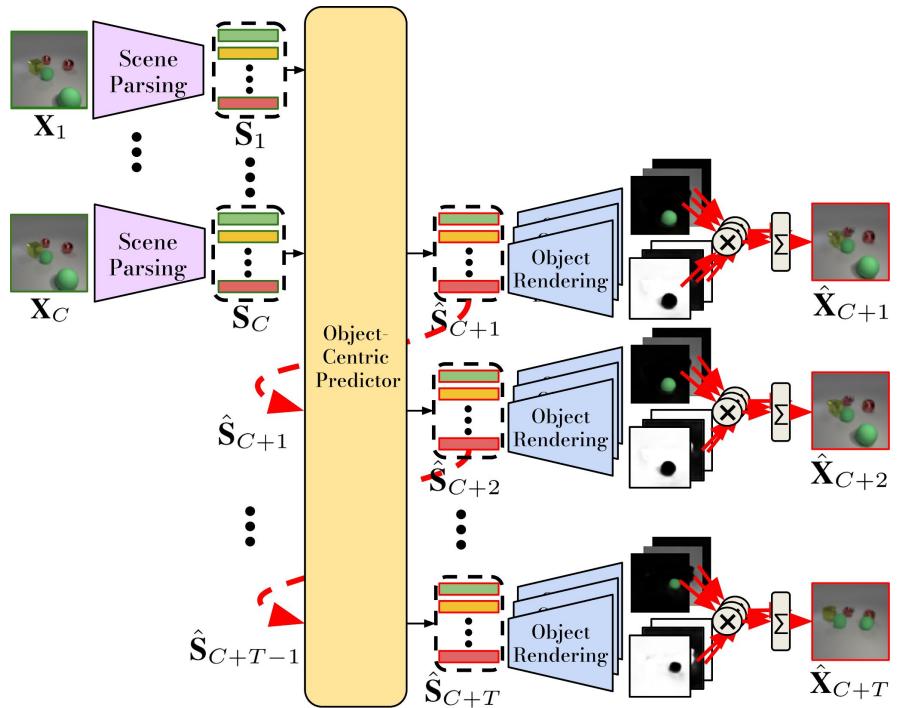
[Villar-Corrales et al. BMVC 2022]

# Future Frame Video Prediction



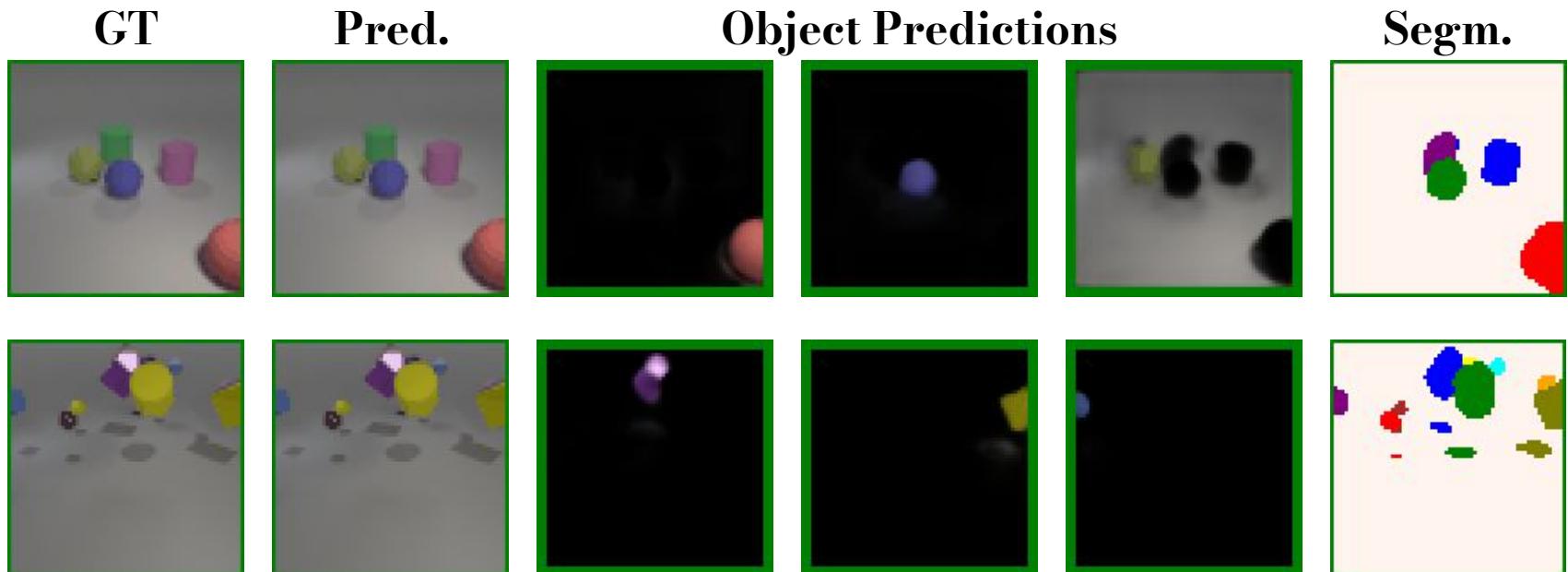
[Villar-Corrales et al. BMVC 2022]

# Object-Centric Video Prediction



[Villar-Corrales et al. ICIP 2023]

# Object-Centric Video Prediction



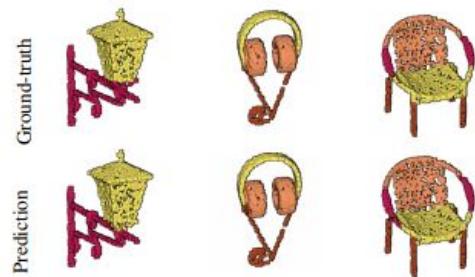
[Villar-Corrales et al. ICIP 2023]

# And much more...

UAV Perception



3D Deep Learning



Neural Representations & Rendering

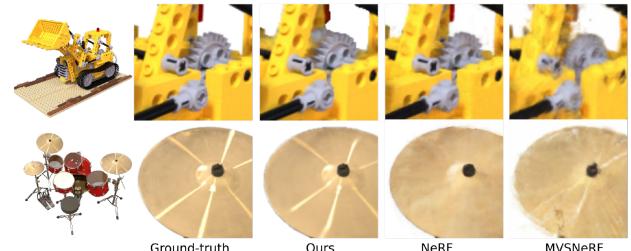
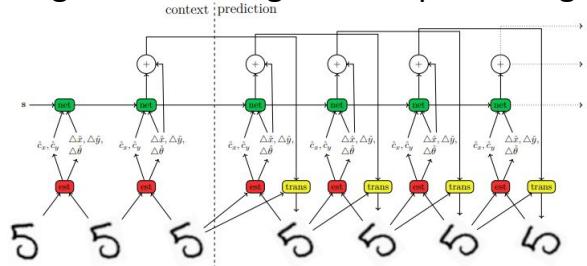


Fig. 8: ShapeNet [39] results of our method.

Signal Processing and Deep Learning



Scene Synthesis



# Once Again...

# Slot Assignment Selection

---

- 16 slots for students
  - Assigned at random

# Registration

---

- 16 slots assigned at random
- Please fill this form before 18.10 (select the time slots that work for you)  
<https://forms.gle/AZ7CiXx3SeMV8wAu8>
- Registration (contact me first):
  - **Uni Bonn:** in BASIS
  - **Bonn-Rhein-Sieg & Others:** Contact me

	Monday	Tuesday	Wednesday	Thursday	Friday
<b>9-11 or 10-12</b>					
<b>11-13 or 12-14</b>					
<b>13-15 or 14-16</b>					
<b>15-17</b>					

# Questions?

