

Deep sensing -Jointly optimize imaging and processing-

Hajime Nagahara¹

¹ Institute for Dataability Science, Osaka University
2-8, Yamadaoka, Suita, Osaka, 565-0871, Japan
E-mail:nagahara@ids.osaka-u.ac.jp

Abstract Deep neural network (DNN) is a powerful tool for solving image processing and computer vision tasks such as image and video reconstructions, object recognition, and scene understanding, etc. However, DNN have been used for only digital domain in the imaging pipeline, such as the feature extractor and classifier models after an image is captured and digitized. In this research, we propose a new framework called “deep sensing.” The proposed framework also models the analog layer to the neural network model and jointly optimizes the parameters in optics and sensor designs of a camera, as well as reconstruction and classification models by the same training strategy.

Keywords: Deep neural network, Computational photography, Compressive sensing

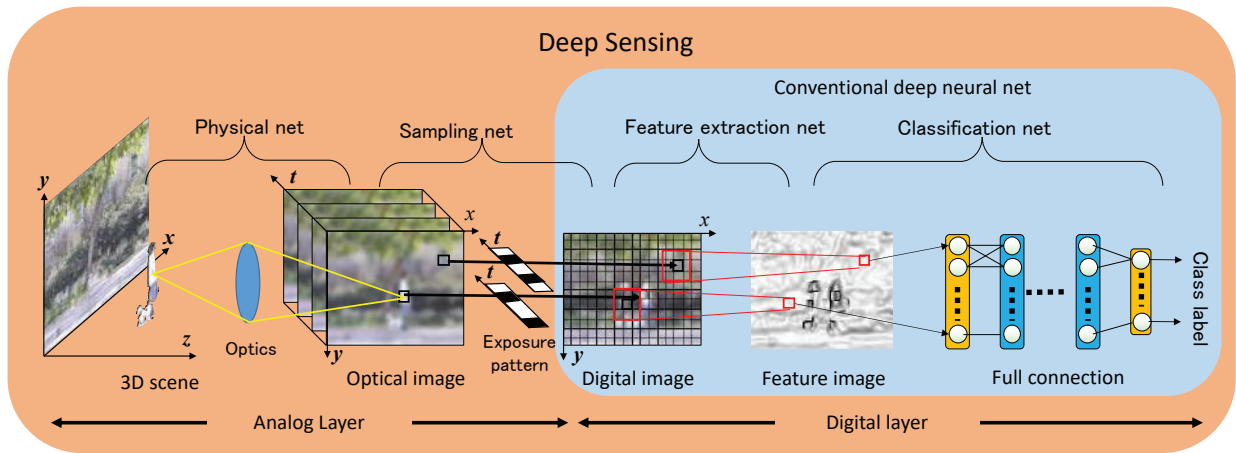


Figure 1: Conceptual figure of Deep sensing

1. Overview

Deep neural network (DNN) is a powerful tool for solving image processing and computer vision tasks such as image and video reconstructions, object recognition, and scene understanding, etc. It realizes to drastically improve the quality of reconstruction and the accuracy of recognition to the classical methods since feature extractor and classifier models are designed by training based on the target data. However, DNN have been used for only the digital domain in the imaging pipeline, such as the feature extractor and classifier models after an image is captured and digitized, as shown in the blue part of figure 1. On the other hand, optics and sensors in the analog layer still have been designed by hand based on theoretical or empirical analysis. It is not always grantee that the designs and hardware setting parameters are optimal to the applications and target tasks. In this research, we propose a new framework called “deep sensing,” as shown in figure 1. The proposed framework also models the analog layer to the neural network model and jointly optimizes the parameters in optics and sensor designs of a camera as well as reconstruction and classification models by the same training strategy. In this talk, I introduce the concept of deep sensing and show our work; compressive light field sensing [1,4,5], compressive video sensing [2], action recognition by a coded image [3], and privacy-preserving imaging. The details are in the corresponding papers.

References

- [1] Y. Inagaki, Y. Kobayashi, K. Takahashi, T. Fujii, and H. Nagahara, “Learning to capture light fields through a coded aperture camera,” European Conference on Computer Vision, 2018, pp. 418–434.
- [2] M. Yoshida, A. Torii, M. Okutomi, K. Endo, Y. Sugiyama, R. Taniguchi, and H. Nagahara, “Joint optimization for compressive video sensing and reconstruction under hardware constraints,” European Conference on Computer Vision, 2018, pp. 634–649.
- [3] T. Okawara, M. Yoshida, H. Nagahara, and Y. Yagi, “Action Recognition from a Single Coded Image,” in Proceedings of IEEE International Conference on Computational Photography, 2020.
- [4] K. Sakai, K. Takahashi, T. Fujii, H. Nagahara: “Acquiring Dynamic Light Fields through Coded Aperture Camera”, European Conference on Computer Vision 2020.
- [5] R. Mizuno, K. Takahashi, M. Yoshida, C. Tsutake, T. Fujii, H. Nagahara, “Acquiring a Dynamic Light Field through a Single-Shot Coded Image”, IEEE Conference on Computer Vision and Pattern Recognition, June, 2022.
- [6] S. Kumawat, T. Okawara, M. Yoshida, H. Nagahara, Yasushi Yagi, “Action Recognition From a Single Coded Imaging”, IEEE Transactions on Pattern Analysis and Machine Intelligence, pp. 1-14, Aug., 2022.
- [7] S. Kumawat, H. Nagahara, “Privacy-Preserving Action Recognition via Motion Difference Quantization”, European Conference on Computer Vision, Oct, 2022.

Deep sensing

- Jointly optimizing sensing and processing -

Hajime Nagahara
Institute for Datability Science, Osaka University

1

Concept of Deep sensing

H30-R1 科研挑戦の研究(萌芽)

Proposing a general framework to design a camera hardware by DNN optimization

2

Our research examples

- Learning to Capture Light Fields through A Coded Aperture Camera [Inagaki+ ECCV2018]
- Acquiring a Dynamic Light Field through a Single-Shot Coded Image [Mizuno+ CVPR2022]
- Joint optimization for Compressive video sensing [Yoshida+ ECCV2018]
- Action Recognition from a Single Coded Image [Okawara+ ICCP2020]

3

Compressive LF sensing

Presented at:
Y. Inagaki, Y. Kobayashi, K. Takahashi, T. Fujii, H. Nagahara:
"Learning to Capture Light Fields through A Coded Aperture Camera",
European Conference on Computer Vision (ECCV) 2018.

4

What is light field (LF)?

Dense set of multi-viewpoint images

Cameras

(s, t) Viewpoint
(u, v) Pixel position

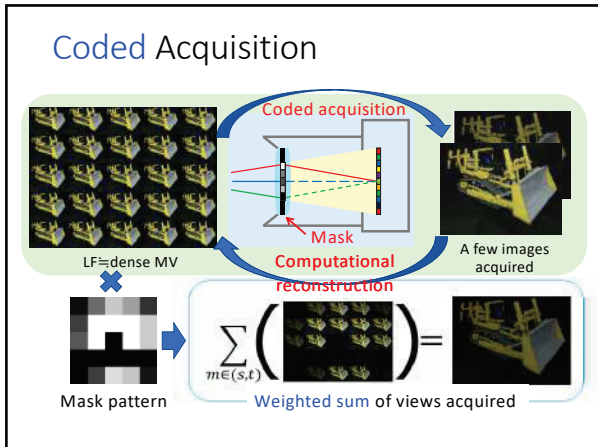
5

Naïve acquisition

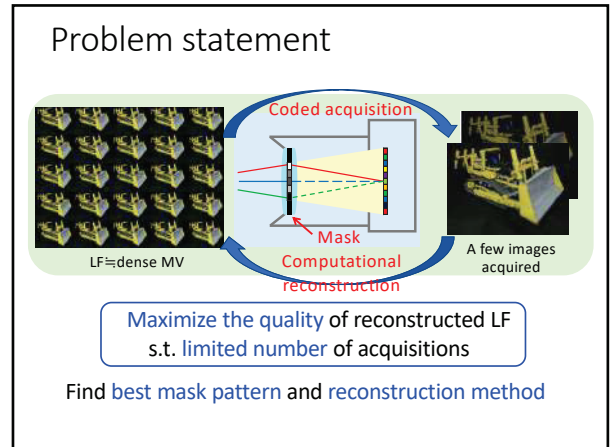
Mask pattern → View by view acquisition

- 25 acquisitions
- Noisy images

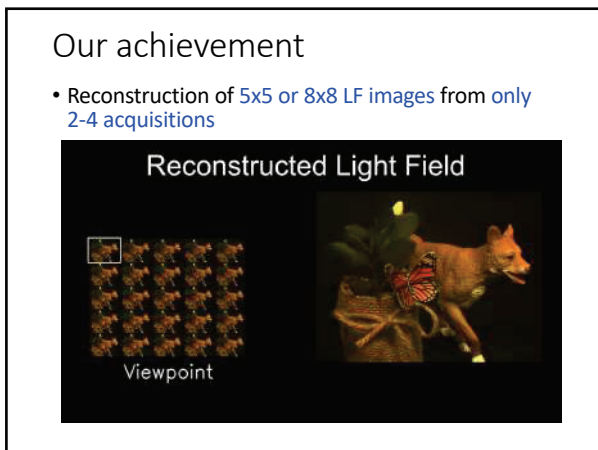
6



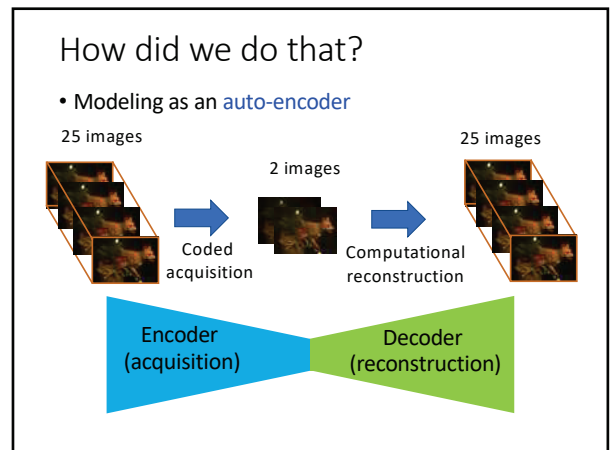
7



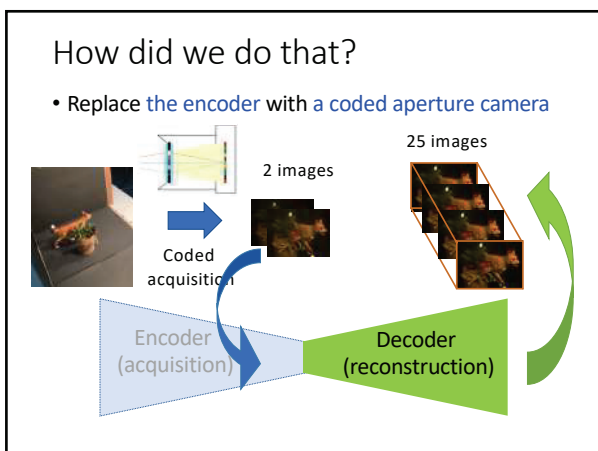
8



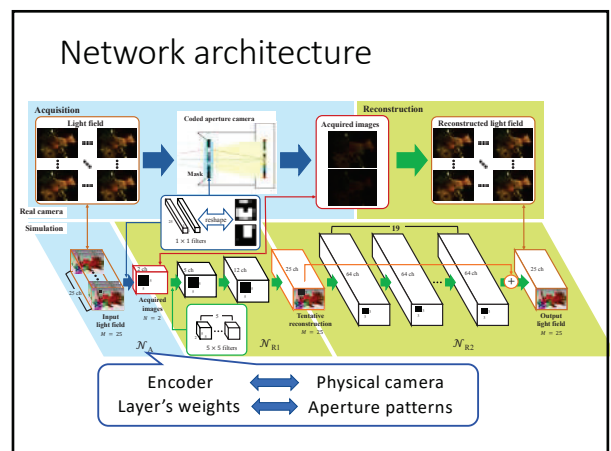
9



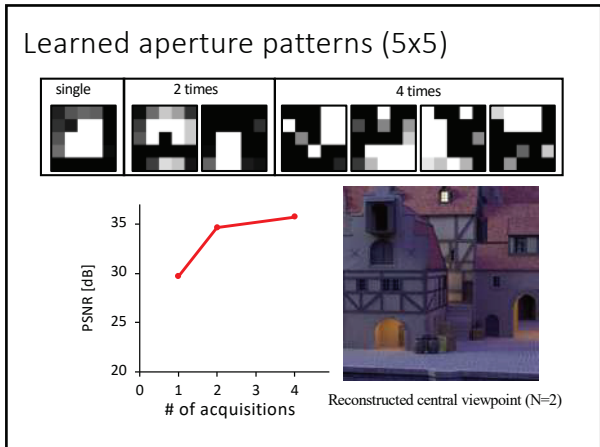
10



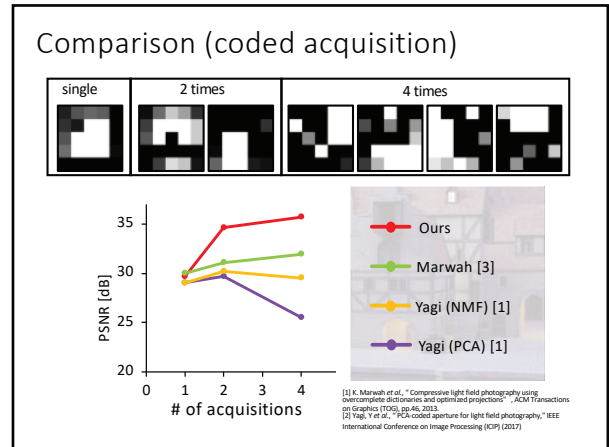
11



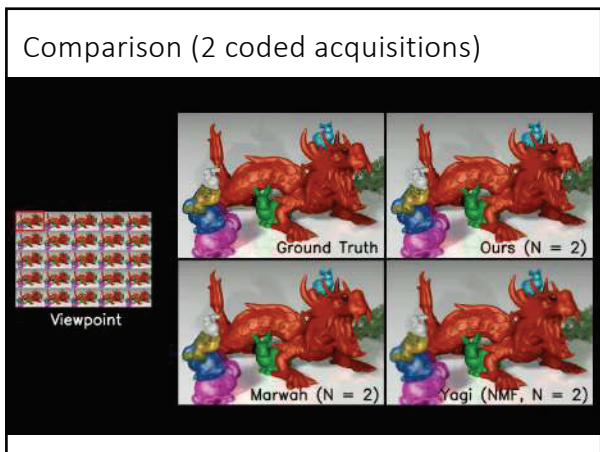
12



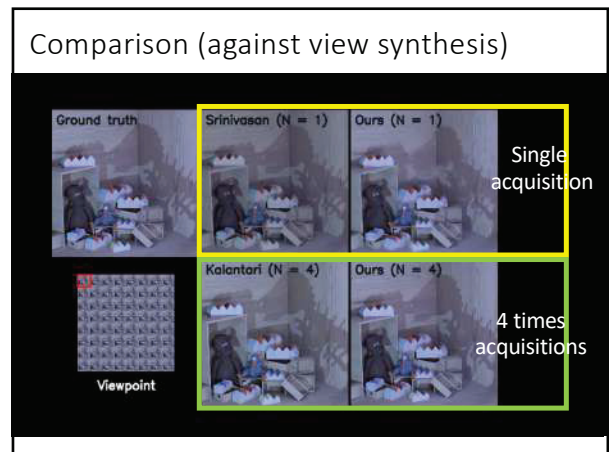
13



14



15

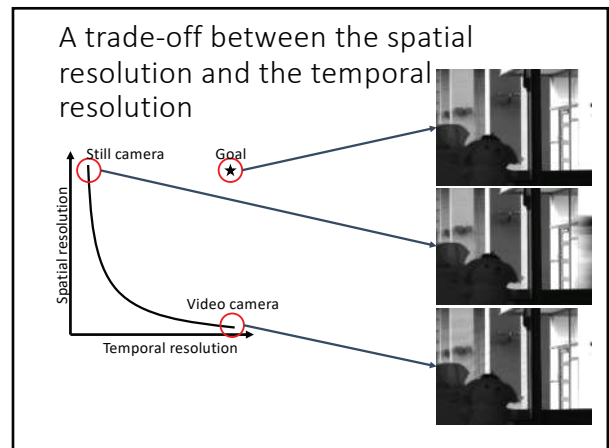


16

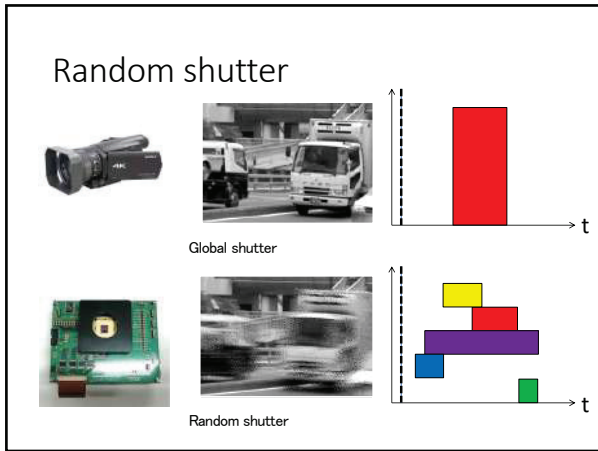
Joint optimization for Compressive video sensing

Michitaka Yoshida, Akihiko Torii, Masatoshi Okutomi,
Kenta Endo, Yukinobu Sugiyama, Rin-ichiro Taniguchi,
Hajime Nagahara
ECCV2018

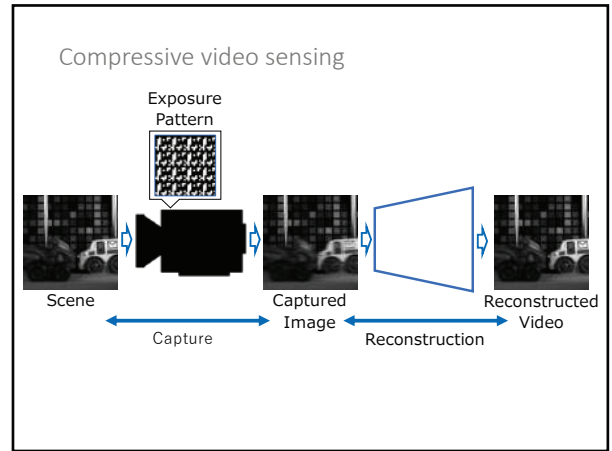
17



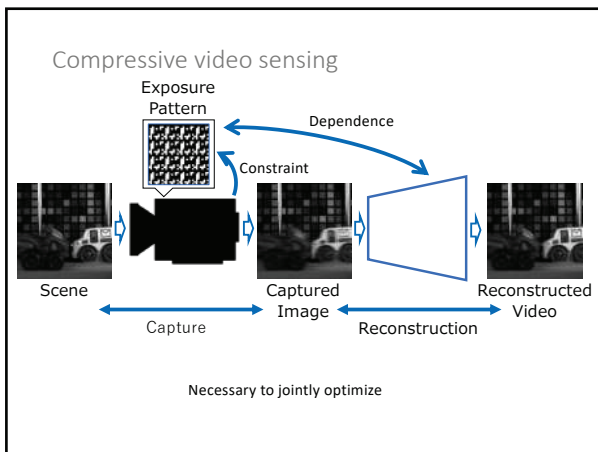
18



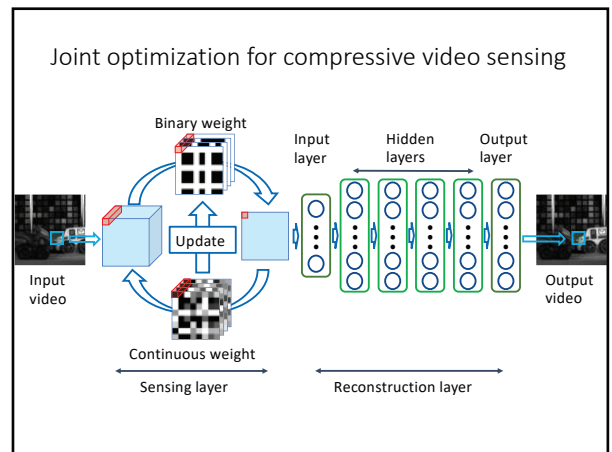
19



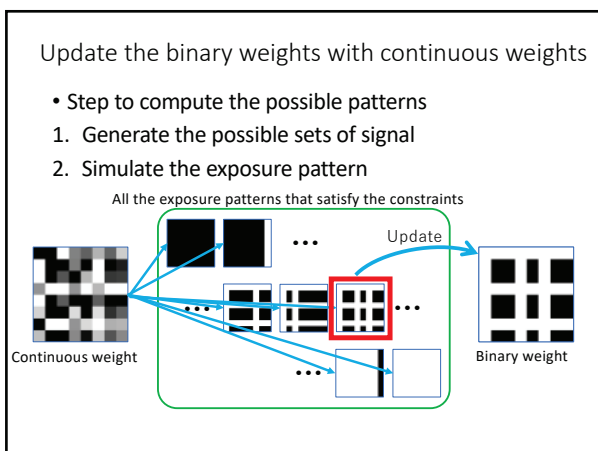
20



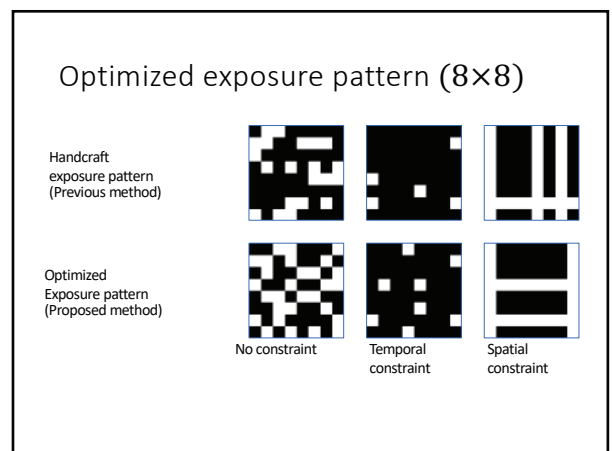
21



22



23



24

Simulation experiments

- No constraint
 - Sarhangnejad et al.

Original video(PSNR) Handcraft(29.77) Optimized(31.39)

14 videos Average	Handcraft	Optimize
PSNR(dB)	29.17	29.99

25

Real experiments

- Spatial constraint on exposure pattern
- Capturing coded images with 15 FPS
- Set 16 exposure patterns per frame

The reconstructed video was equivalent to 240 FPS

26

Real experiments

Captured Image Reconstructed Video

27

Compressive color video sensing

Scene → Capture → Captured Image → Reconstruction → Reconstructed Video

28

Network structure to optimize the color filter, exposure pattern and reconstruction

- Input : RGB Video ($8 \times 8 \times 3 \times 16$)
- Output : RGB Video ($8 \times 8 \times 3 \times 16$)

Input layer: 64 4 Hidden layers: 3072 Output layer: 3072

29

Color filter (8×8)

Bayer Filter Optimized Filter (100 epoch)

R:16 G:32 B:16 R:19 G:23 B:22

30

Exposure pattern (Different color filter)

Monochrome (500epoch) Bayer filter (100epoch) Optimized filter (100epoch)

31

Reconstruction results

Original Video

Color Filter	Bayer	Bayer	Optimize
Exposure Pattern	Fixed	Optimize	Optimize
Epoch	100	100	100
PSNR	24.18	23.92	24.34

Average 25 Videos	Only decoder	Exposure pattern	Filter + Pattern
PSNR	26.56	26.43	26.76

32

Acquiring a Dynamic Light Field through a Single-Shot Coded Image

Ryoya Mizuno, Keita Takahashi, Michitaka Yoshida, Chihiro Tsutake, Toshiaki Fujii, Hajime Nagahara
CVPR2022

33

Combination of aperture and pixel codings

Coded aperture camera
Coding viewpoint dimension to reconstruct Light field

Coded exposure camera
Coding space-time dimension to reconstruct video

5D Light field video can be captured and reconstructed

34

Coded aperture + Coded exposure

Dynamic LF

Lens

Coded exposure sensor

Aperture mask

Pixel mask

Coded aperture

Coded exposure

Coded captured image

Synchronized 4 aperture and 4 exposure patterns
Compressing 25 views x 4 frames into a single coded image

35

Simulation experiment (wakusei)

Ground Truth

提案 A+P

A_only

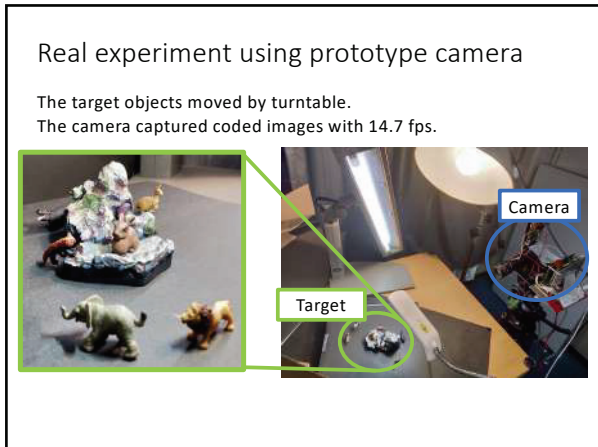
P_only

Normal

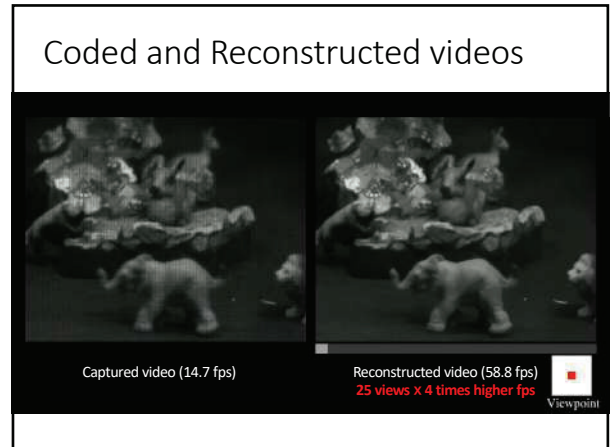
Viewpoint (5-5 views)

PSNR (dB)

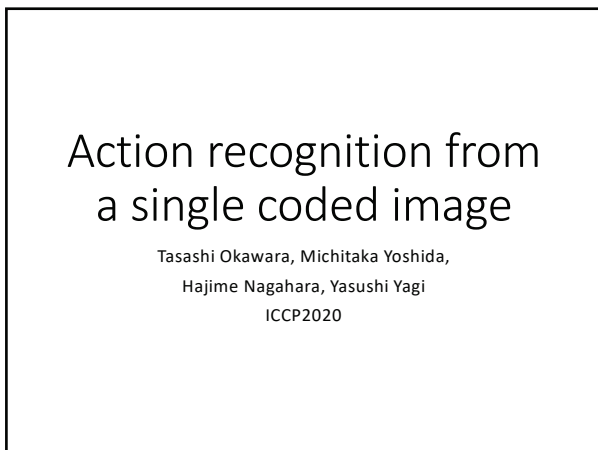
36



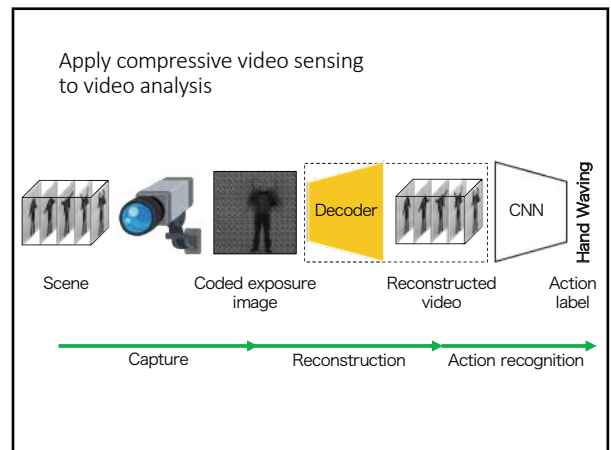
37



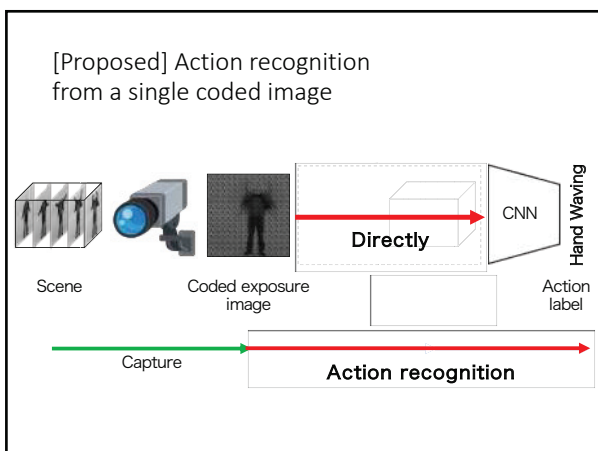
38



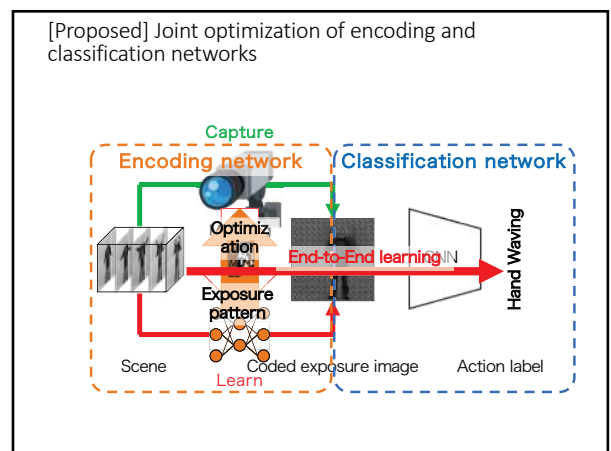
39



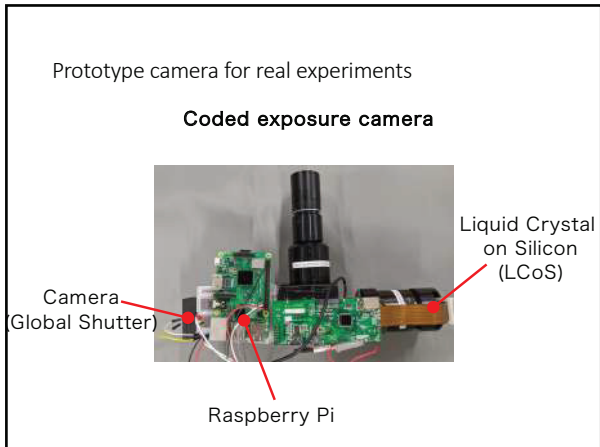
40



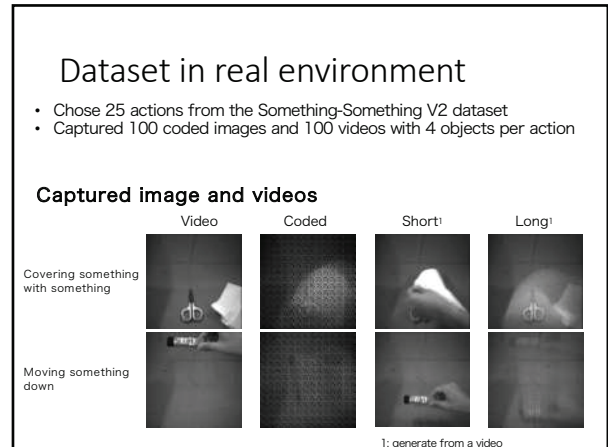
41



42



43

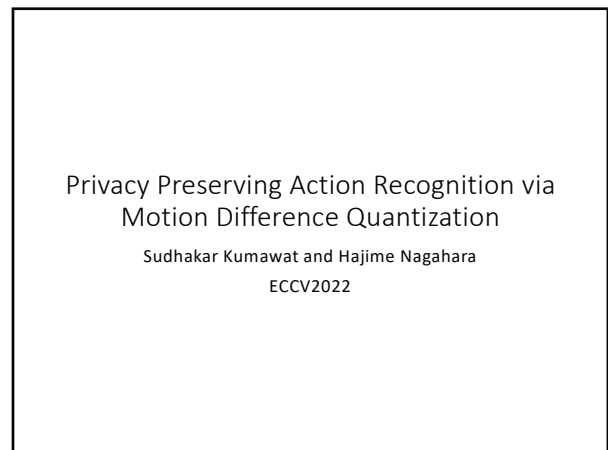


44

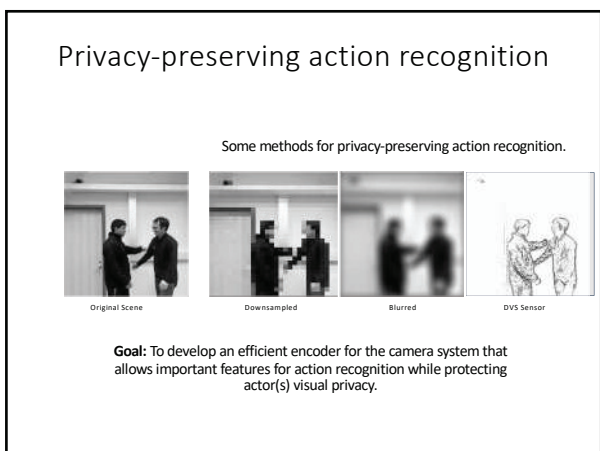
Results

Input		Model	Simulation			Real		
			Top1	Top3	Top5	Top1	Top3	Top5
Video (upper bound)		C3D	47.1	69.4	76.9	71.0	88.0	88.0
Single image	Coded exposure (Proposed)	SVC2D	41.6	58.9	67.2	72.0	84.0	88.0
	Long exposure	C2D	13.8	30.4	39.4	20.0	40.0	52.0
	Short exposure	C2D	14.6	32.5	40.5	21.0	47.0	60.0

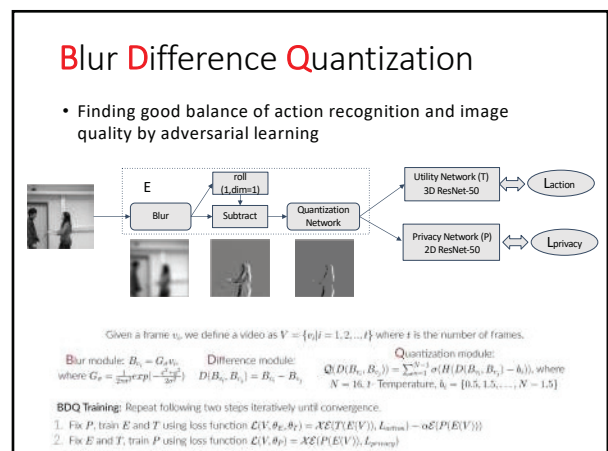
45



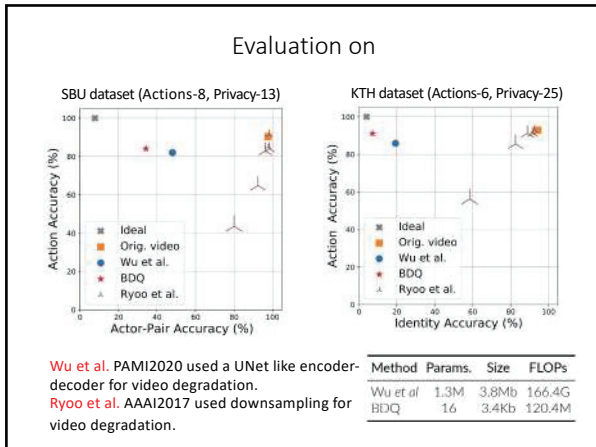
46



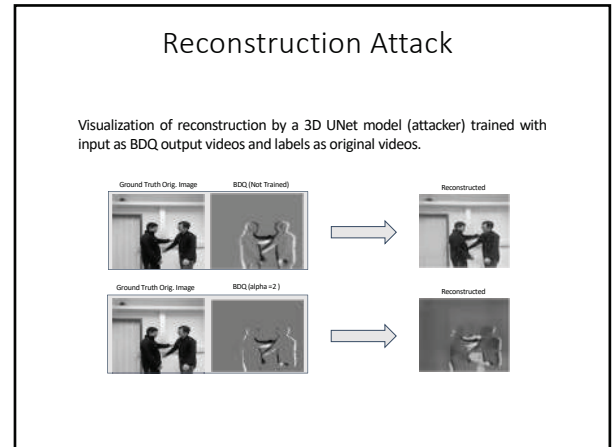
47



48



49



50

Insights of Deep Sensing

- Optimize the encoder to simulate the imaging devise.
- Finding the good 2D coded image as a bottleneck.

Reconstruction

- Important cues for reconstruction.
- Omitting redundant information which decoder can be interpolated.

Classification

- Distinguishable features btw classes.
- Filter out the redundant information.
- Rich information is not always better: Cost, data rate, and privacy.

51

Conclusion

- Proposing deep sensing
 - Parameters for optical designs can be jointly optimized with decoder by same deep learning frame work.
- Demonstrating two example cases:
 - Compressive light field acquisition (4D → 2D)
 - Compressive video sensing (3D → 2D)
 - Compressive LF video sensing (5D → 2D)
 - Action recognition directly from coded image (3D → 174 classes)
 - Balancing the accuracy and privacy.

52