

Synergy of physics and learning-based models in computational imaging and display

Zihao Wang
PhD defense, July 23, 2020

Committee:
Prof. Oliver Cossairt (advisor),
Prof. Aggelos Katsaggelos,
Prof. Jack Tumblin,
Dr. Nathan Matsuda

Computational imaging

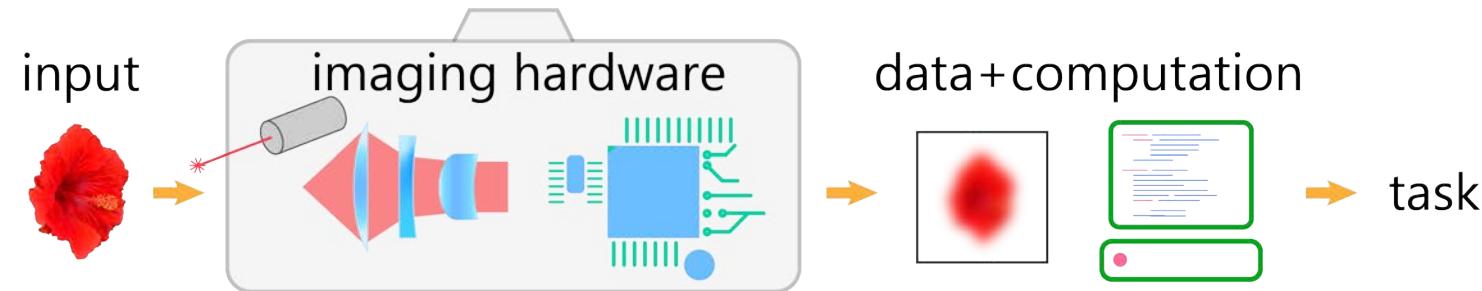


Image formation models

Physics-based
models

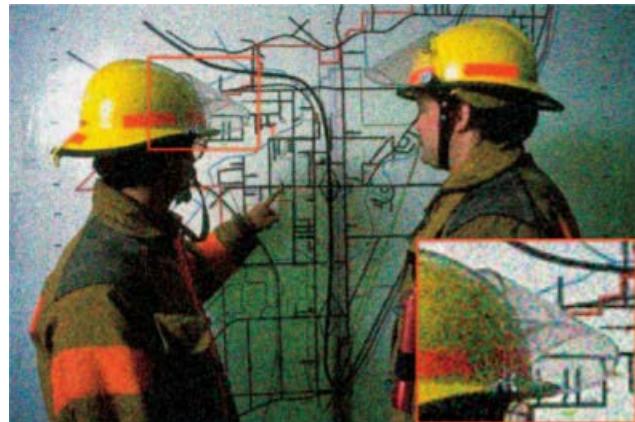
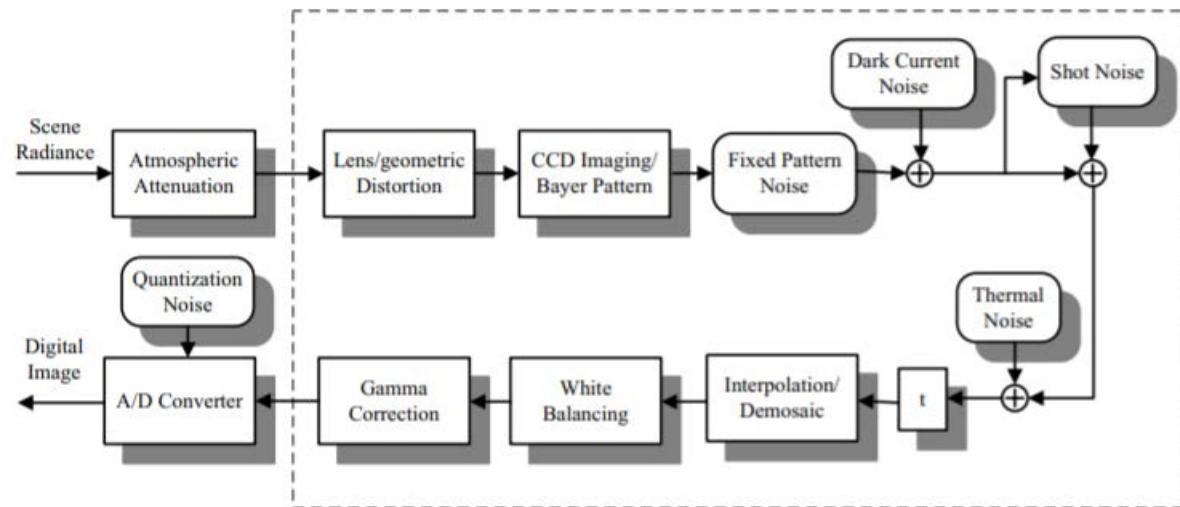
Learning-based
models



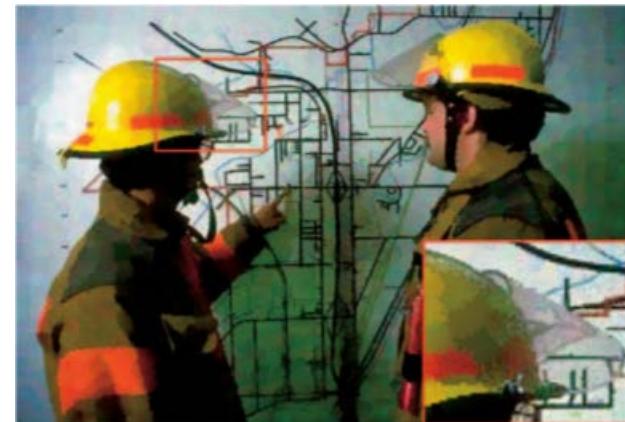
Example: image denoising

Liu et al. CVPR'06, PAMI'07

In-camera ISP pipeline



noisy image



denoised image

Noise modeling

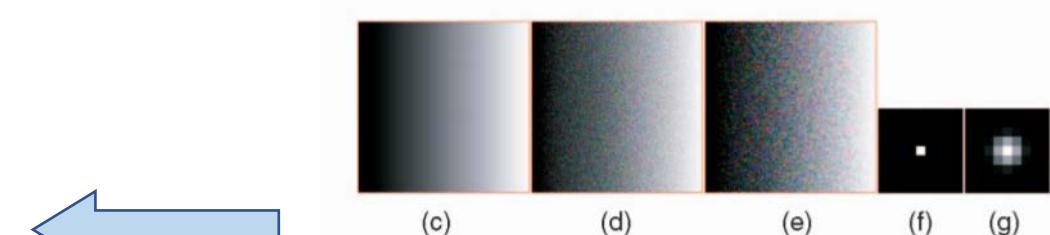
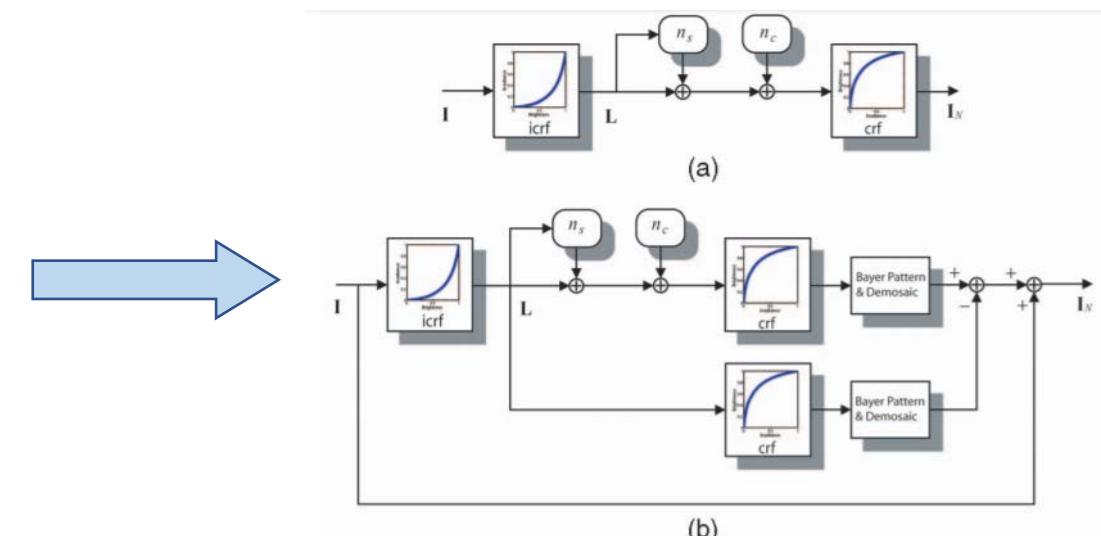
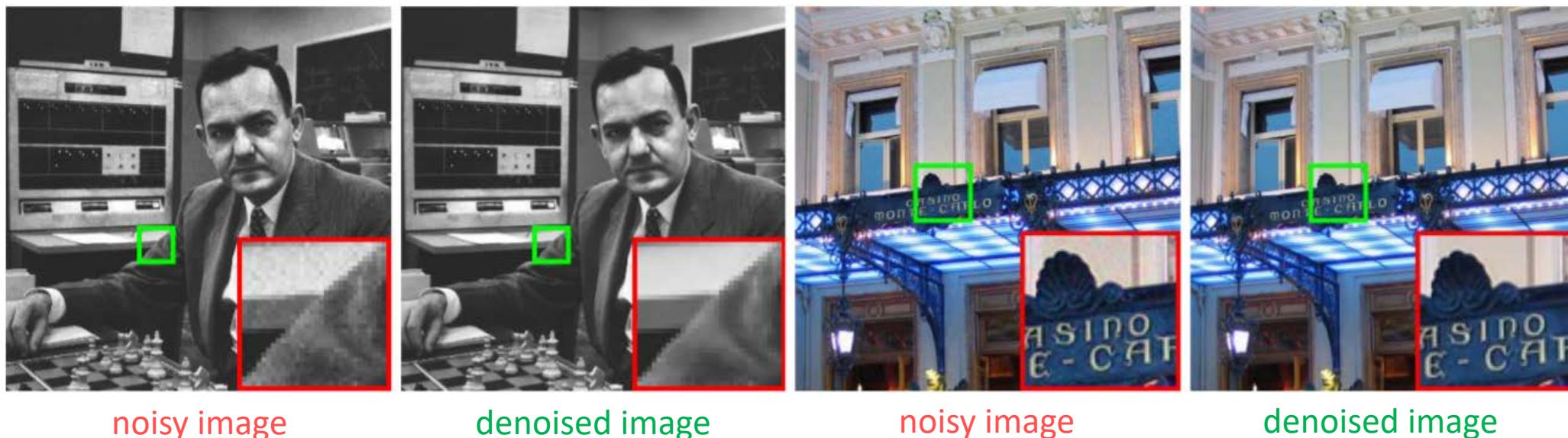
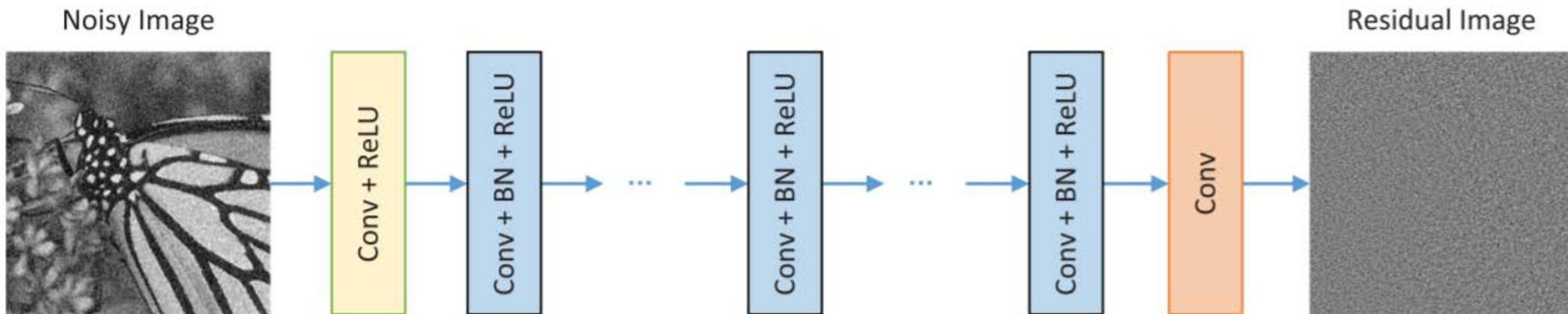


Fig. 4. Block diagrams showing noise simulations for color camera images. (a) shows independent white noise synthesis and (b) adds CCD color filter pattern sensing and demosaicing to model spatial correlations in the camera noise [28]. (c) Test pattern. (d) and (e) The synthesized images of (a) and (b). (f) and (g) The corresponding autocorrelation.

Example: image denoising (10 years later)

Zhang et al. TIP'17

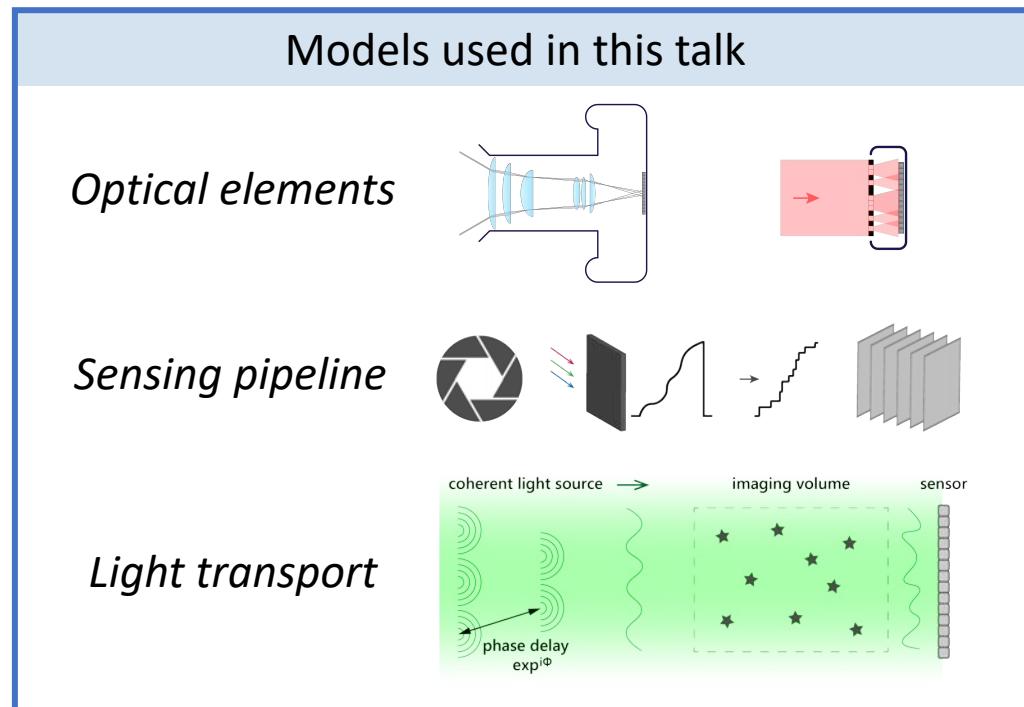
DnCNN



Non-universal definition

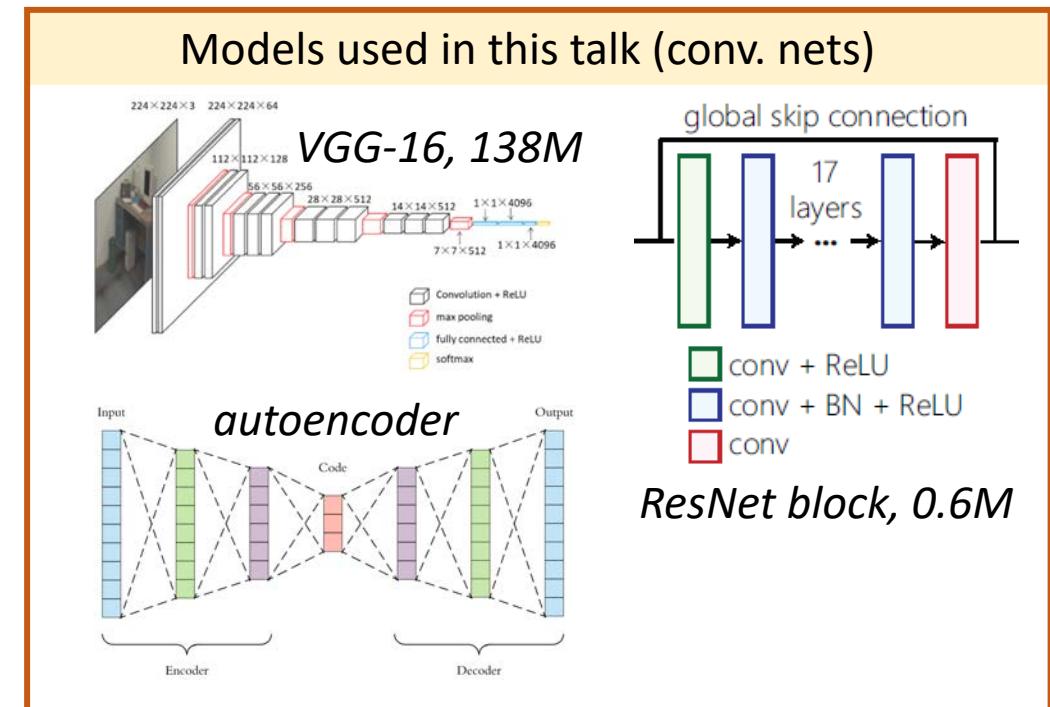
Physics-based models

*Analytical models from the image formation (hardware)
few parameters*



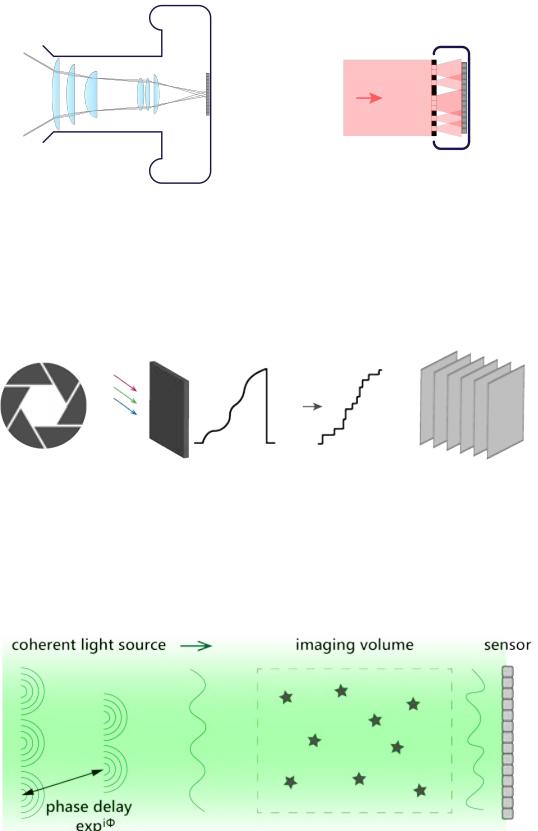
Learning-based models

*Statistical priors from lots of data
millions of parameters*



Overview

Physics-based models



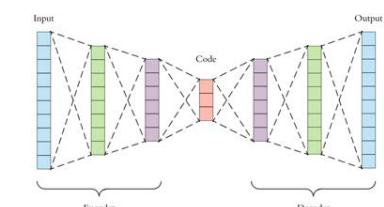
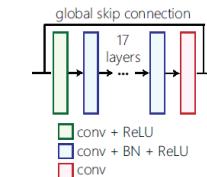
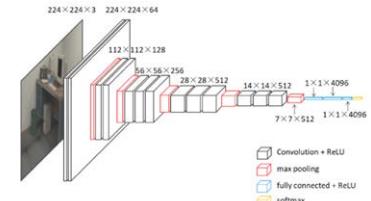
Synergistic models

1. Privacy preserving action recognition, *CVPRW'19*

2. High frame-rate video frame synthesis, *ICCVW'19*

3. 3D holographic display design, image compression

Learning-based models



Why are learning-based models not enough?

Imaging hardware

Datasets + neural networks

task

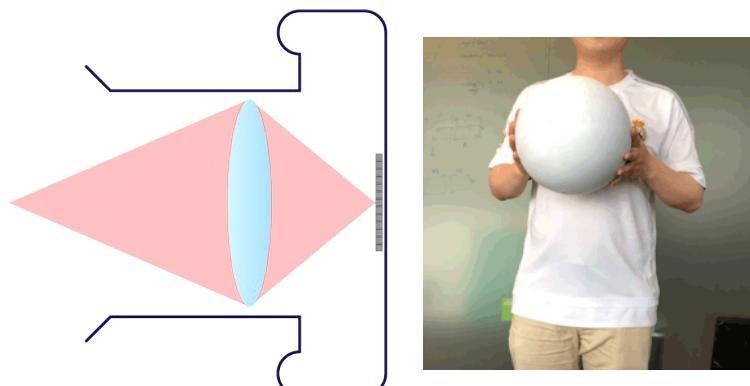
Synergistic models

Image formation

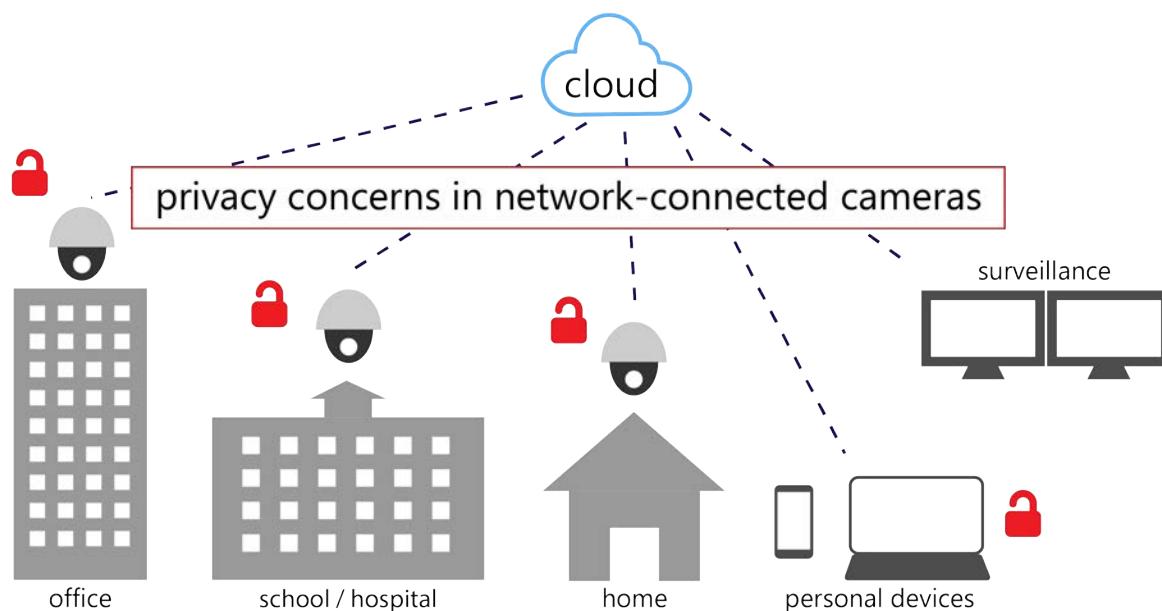
Neural networks

Pre-capture privacy preservation

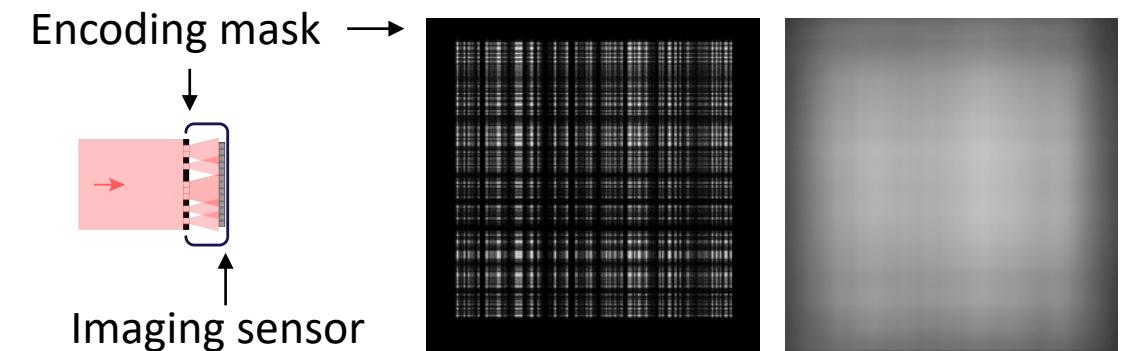
Lens-based cameras



[This is a video]



Lens-free coded aperture cameras



[This is a video]

Restoration can be done with mask info., not without masks.



Privacy preserving action recognition:
Mask-free & restoration-free

Privacy preserving action recognition

	(%)	Training	Validation
<i>Overfitting!</i>	CA	50 th : 79.06	50 th : 63.21
<i>Info. is still there!</i>	grayscale	50 th : 99.56	50 th : 94.39



Need to remove the mask effect
Design mask-independent features

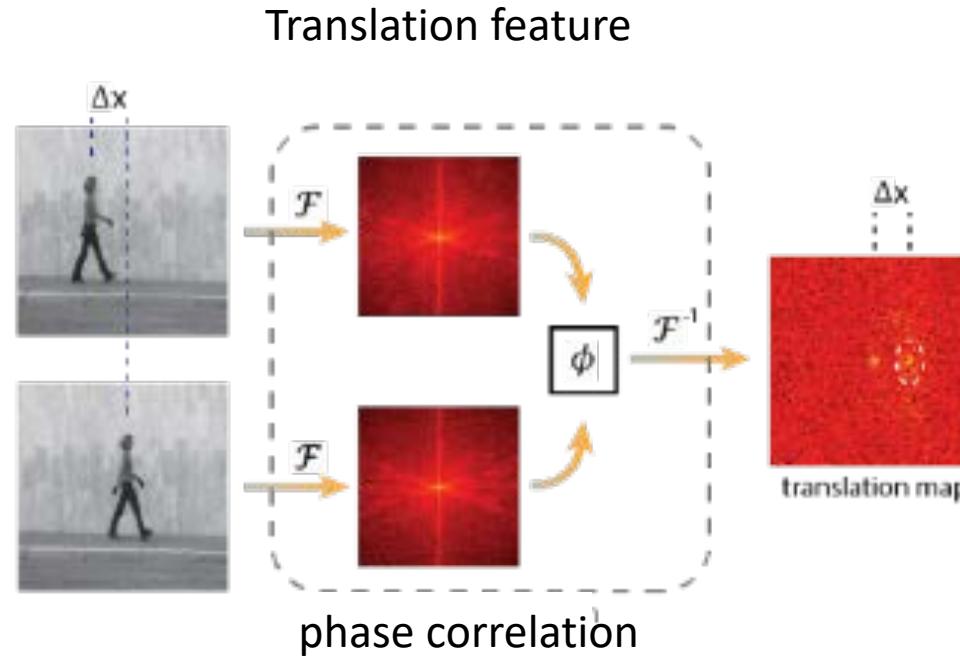
Task: 5-class classification. (RGB2gray vs. synthetic CA images)

“writing on board”, “Wall pushups”, “blowing candles”, “pushups”, “mopping floor” from UCF-101

3 frames for each sample

Classifier: VGG-16

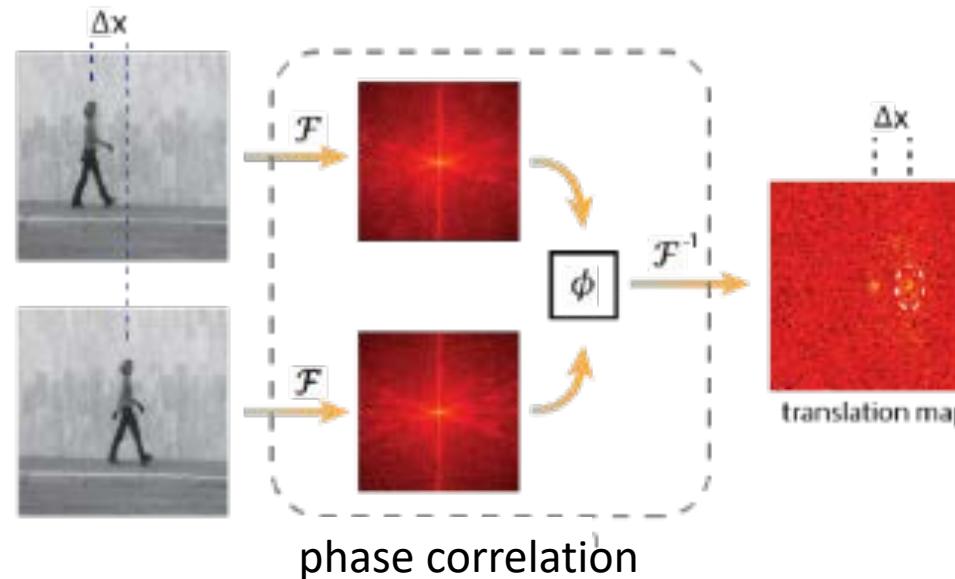
Mask-independent motion features



- $o_2(\mathbf{p}) = o_1(\mathbf{p} + \Delta\mathbf{p})$
- $O_2(\mathbf{v}) = \phi(\Delta\mathbf{p})O_1(\mathbf{v})$ Fourier transform
- $C(\mathbf{v}) = \frac{O_1 \cdot O_2^*}{|O_1 \cdot O_2^*|} = \phi^* \frac{O_1 \cdot O_1^*}{|O_1 \cdot O_1^*|} = \phi(-\Delta\mathbf{p})$ Cross power spectrum
- $c(\mathbf{p}) = \delta(\mathbf{p} + \Delta\mathbf{p})$ Inverse Fourier transform

Mask-independent motion features

Translation feature



Mask-independent

- $d_2(\mathbf{p}) = o_2(\mathbf{p}) * a = d_1(\mathbf{p} + \Delta\mathbf{p})$ a : coded aperture
- $D_2(\mathbf{v}) = O_2 \cdot A = \phi \cdot O_1 \cdot A$ Fourier transform
- $C_d(\mathbf{v}) = \frac{D_1 \cdot D_2^*}{|D_1 \cdot D_2^*|} = \phi^* \frac{O_1 \cdot \mathbf{A} \cdot \mathbf{A}^* \cdot O_1^*}{|O_1 \cdot \mathbf{A} \cdot \mathbf{A}^* \cdot O_1^*|} \approx \phi^*$ Cross power spectrum
- $c_d(\mathbf{p}) = \delta(\mathbf{p} + \Delta\mathbf{p})$ Inverse Fourier transform

Mask-independent motion features

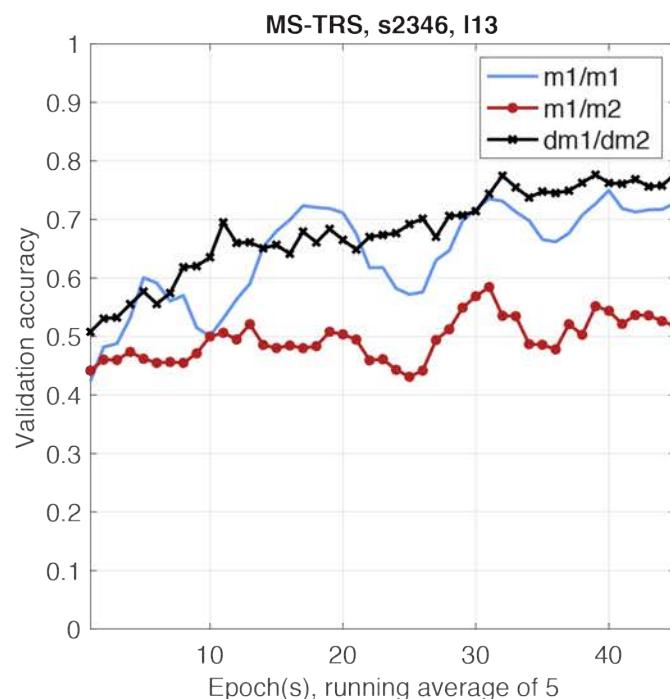
Translation features

phase correlation

Mask-independent



Augmentation: change masks during training



Rotation & Scale features

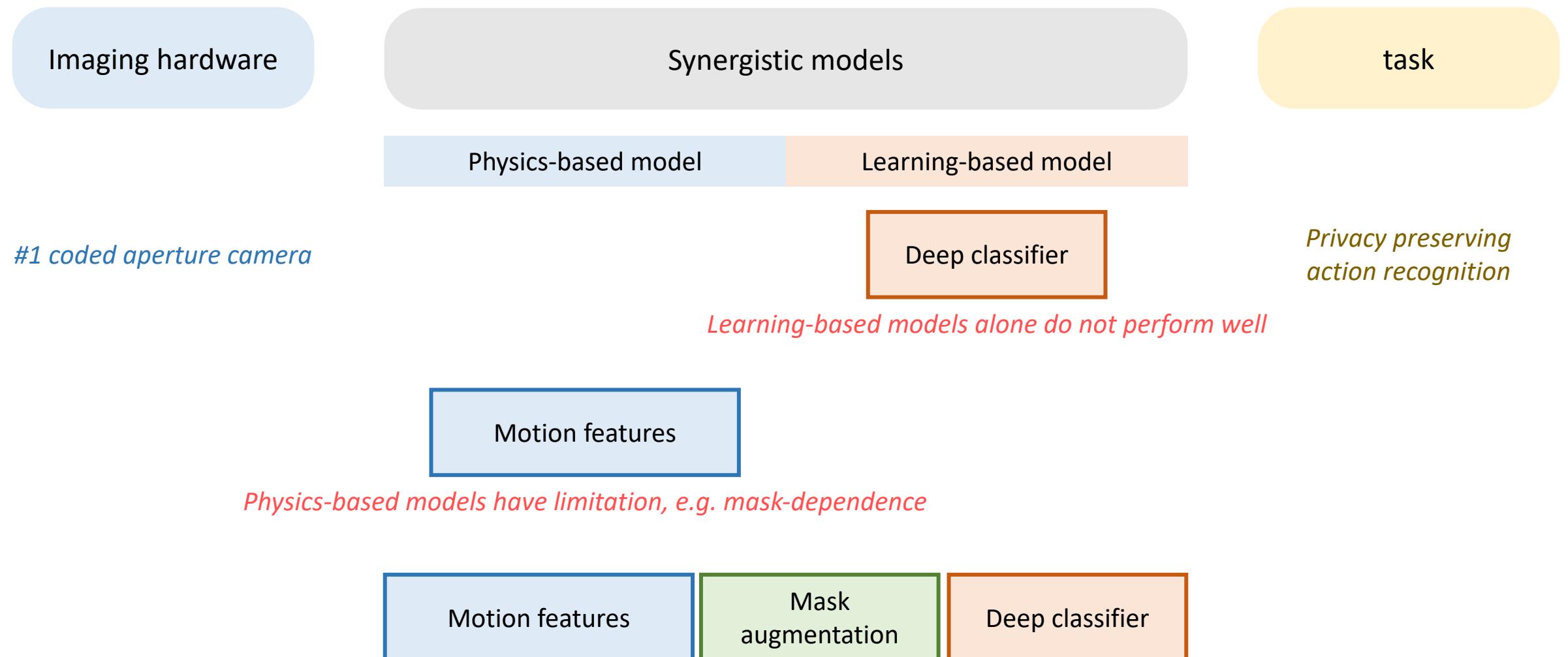
*Fourier-Merlin transform (Cartesian to Polar)
+ phase correlation*

Not mask-independent

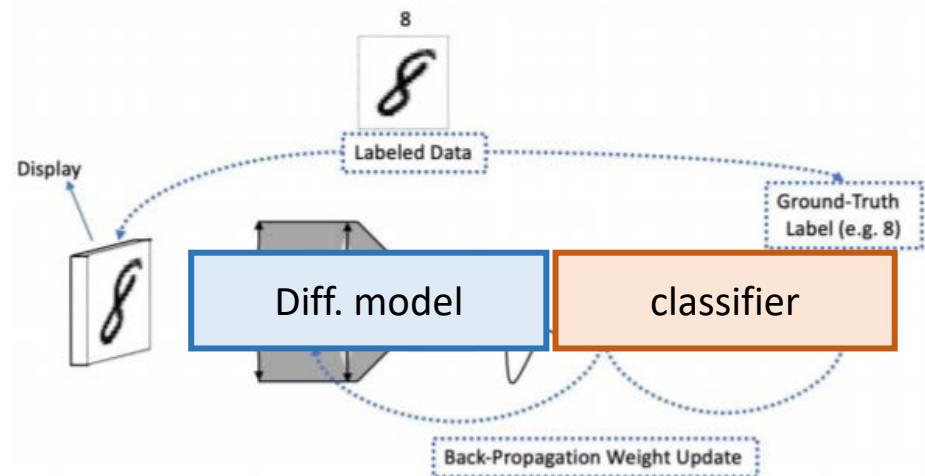
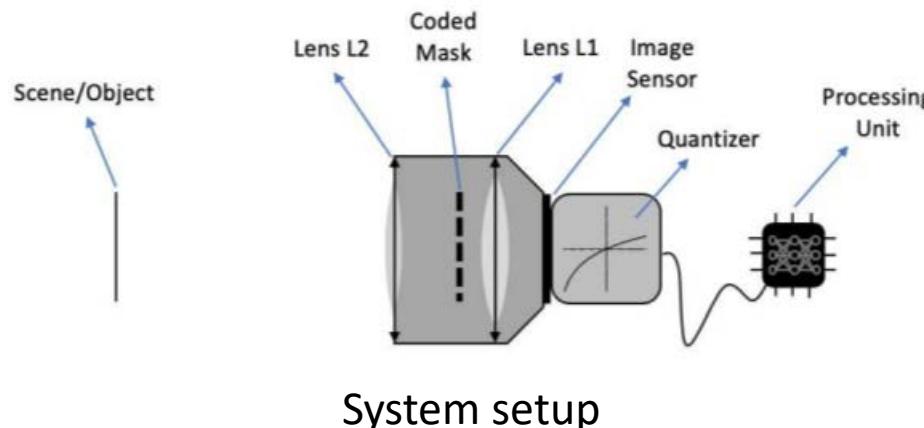
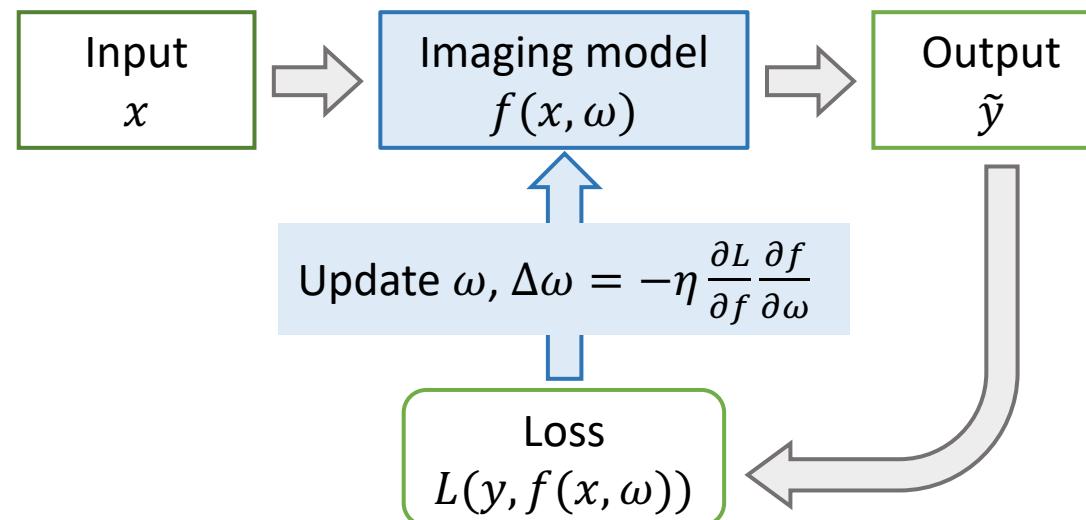
Training w/ randomly varying masks
vs.
w/ single fixed mask, test w/ a different mask

More results are in the paper

Why synergistic models?



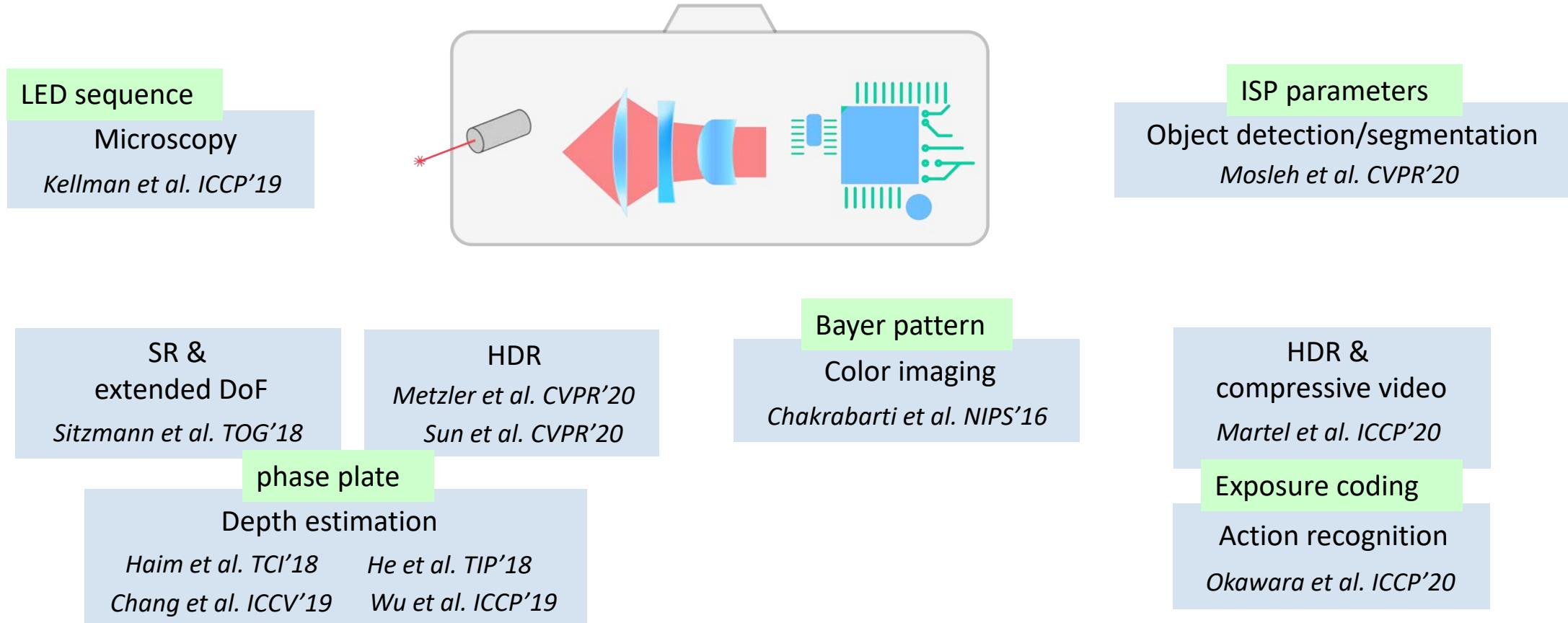
Differentiable imaging model



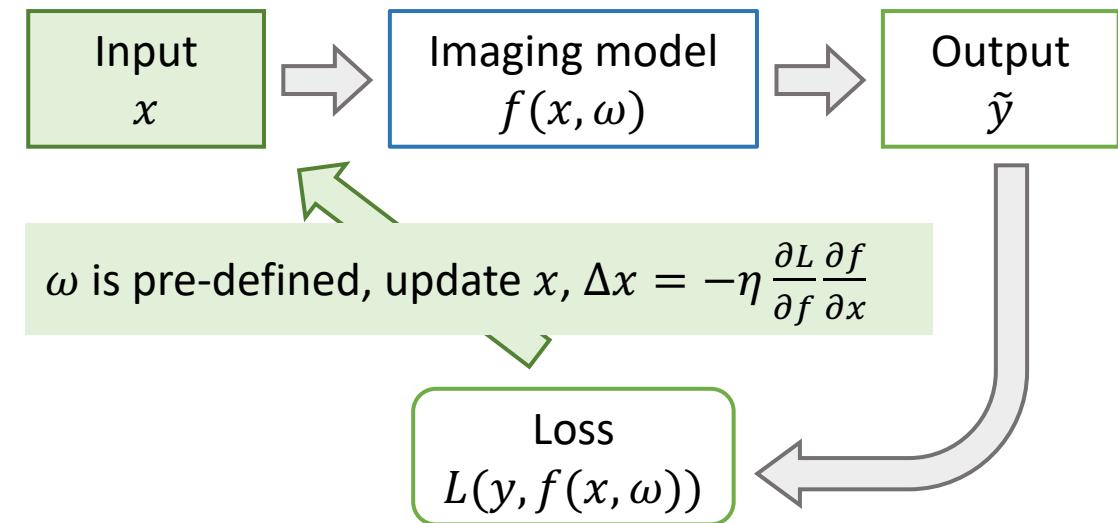
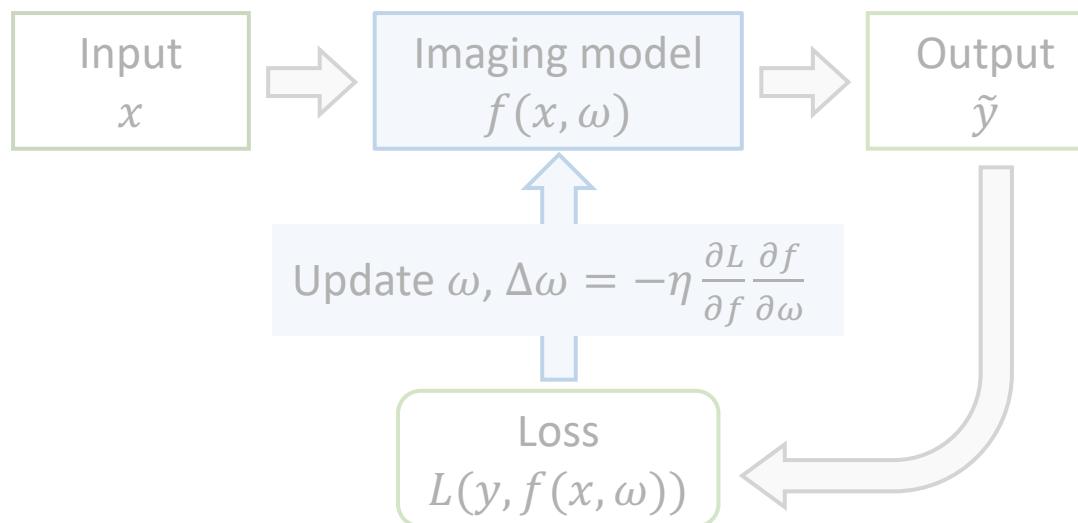
The coded mask is the first layer of the ConvNet

Task-specific differentiable imaging

Part of the learnable parameters ω has physical counterparts



Task-specific differentiable imaging



Differentiable inverse solver

Optical tomography

Kamilov et al. Optica'15

Neural 3D rendering

Kato et al. CVPR'18 ...

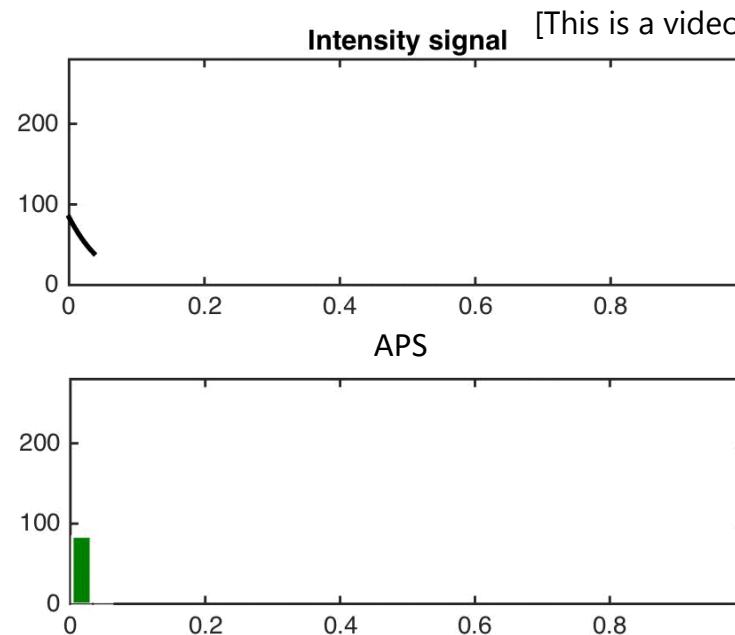
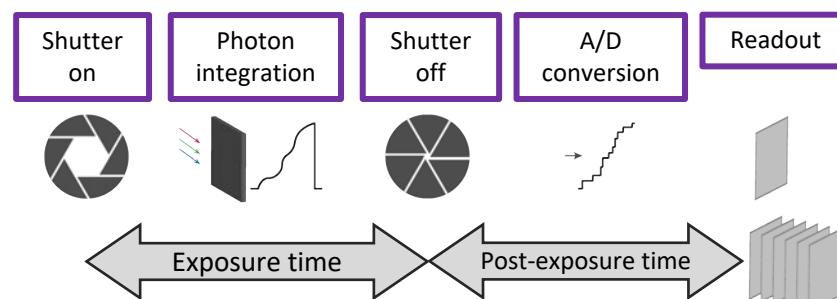
High frame-rate video frame synthesis
ICCVW'19

Holographic display

Chakravarthula et al. TOG'19

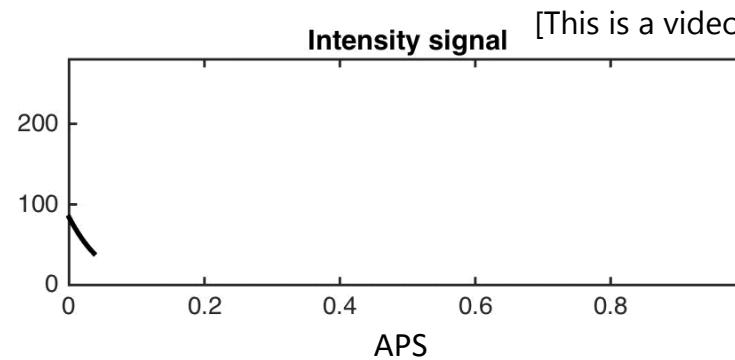
Background: what's making a camera slow?

Active pixel sensors (APS):
agnostic to scene motion

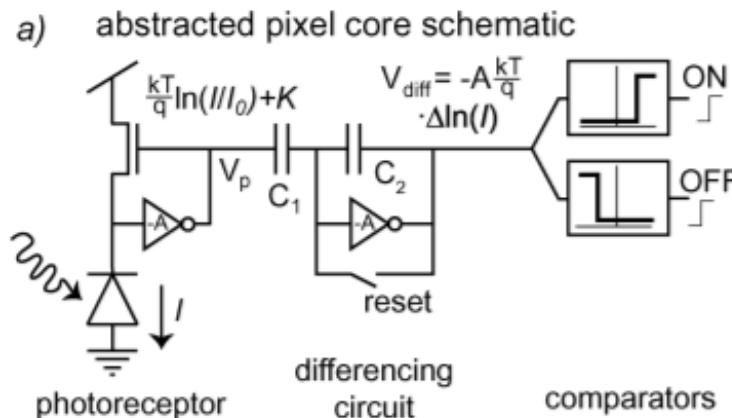


Background: what's making a camera slow?

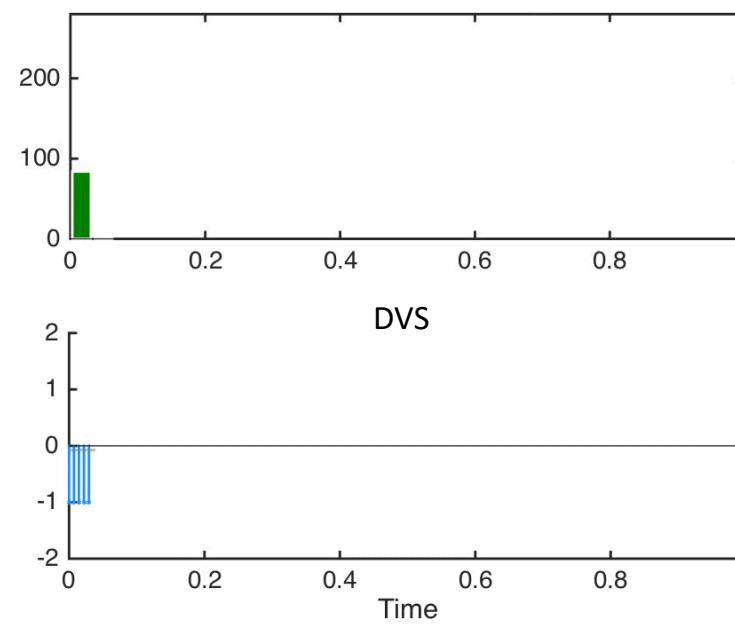
Active pixel sensors (APS):
agnostic to scene motion



Dynamic vision sensors (DVS):

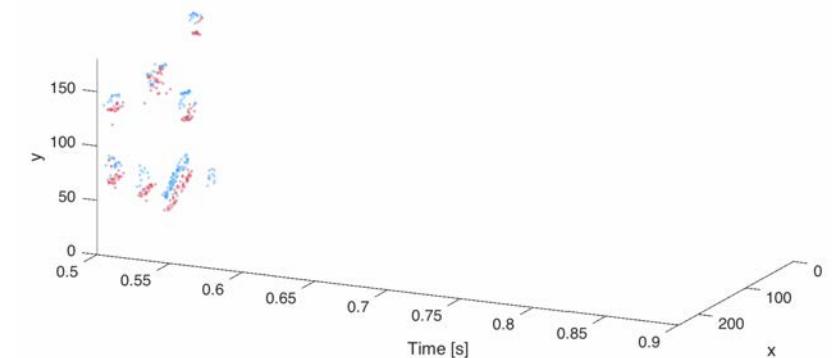


Lichtsteiner et al. IEEE SSC 2008

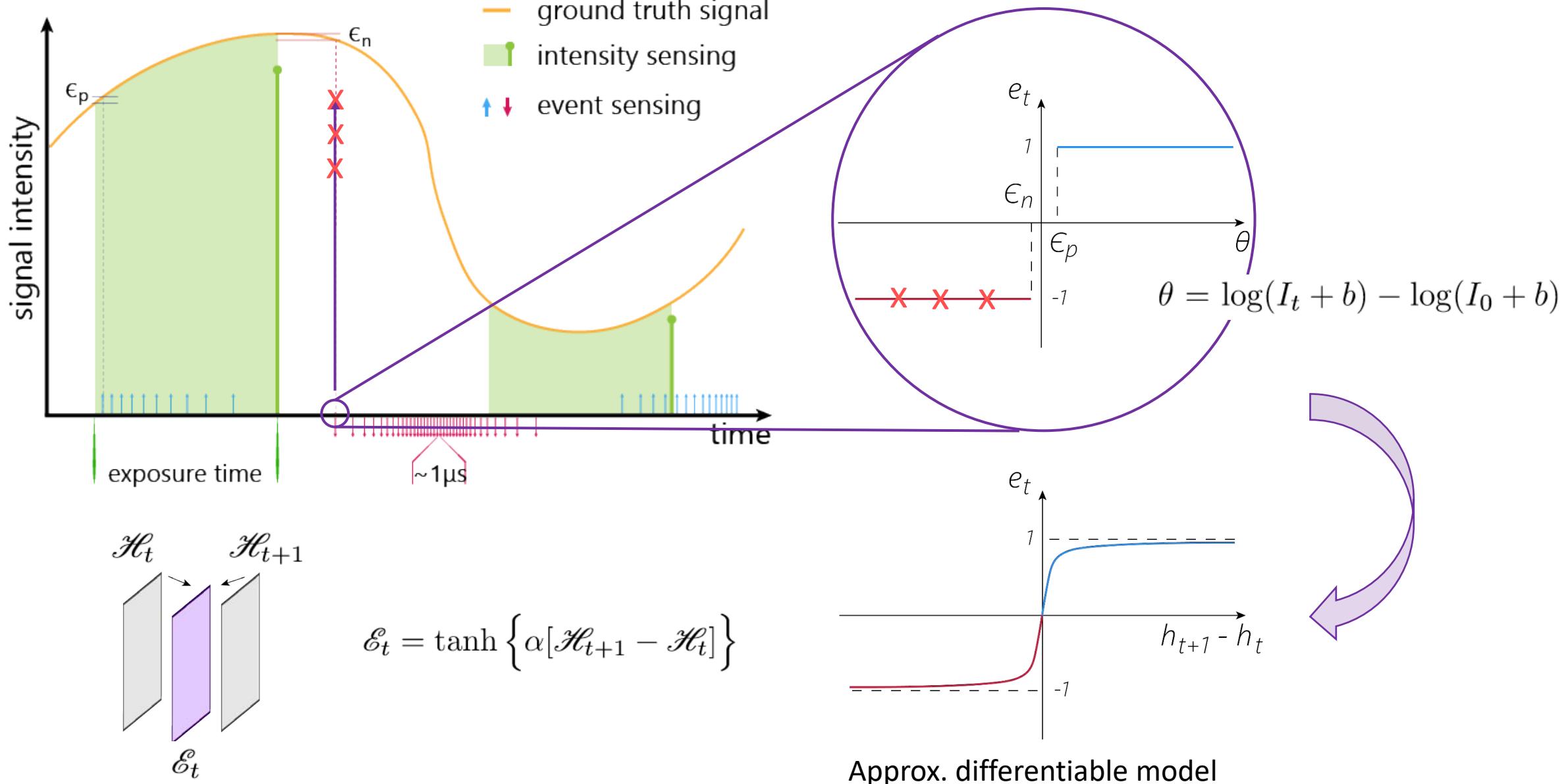


Fast (~1us latency)

But only binary data; also noisy



Synthesizing dense frames from intensity (APS) and events (DVS)

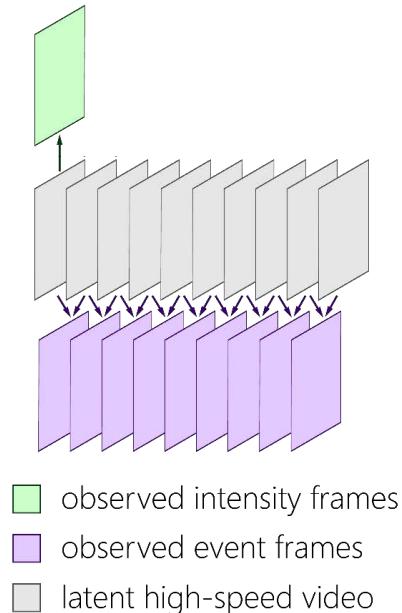


Differentiable inverse solver for video frame synthesis

Fusion setting:

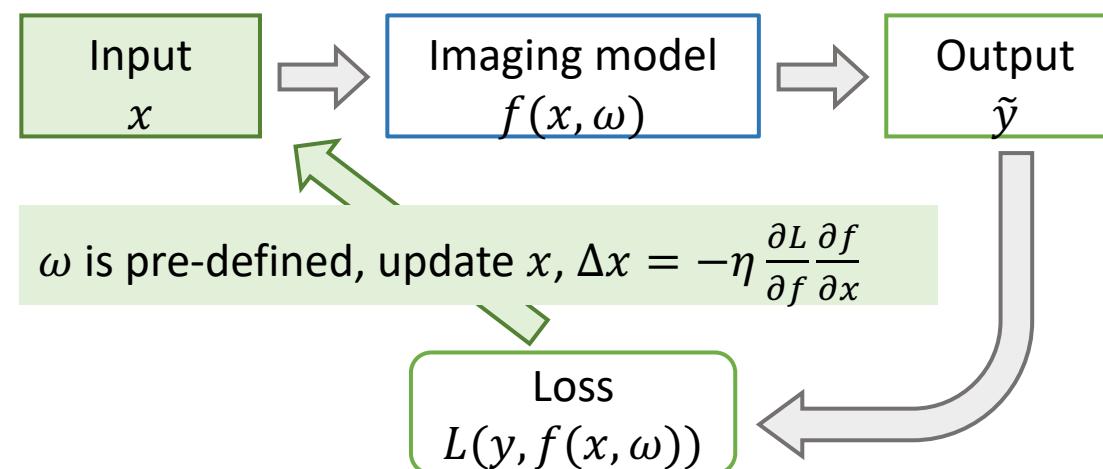
Known first intensity frame

Known subsequent event frames



■ latent high-speed video

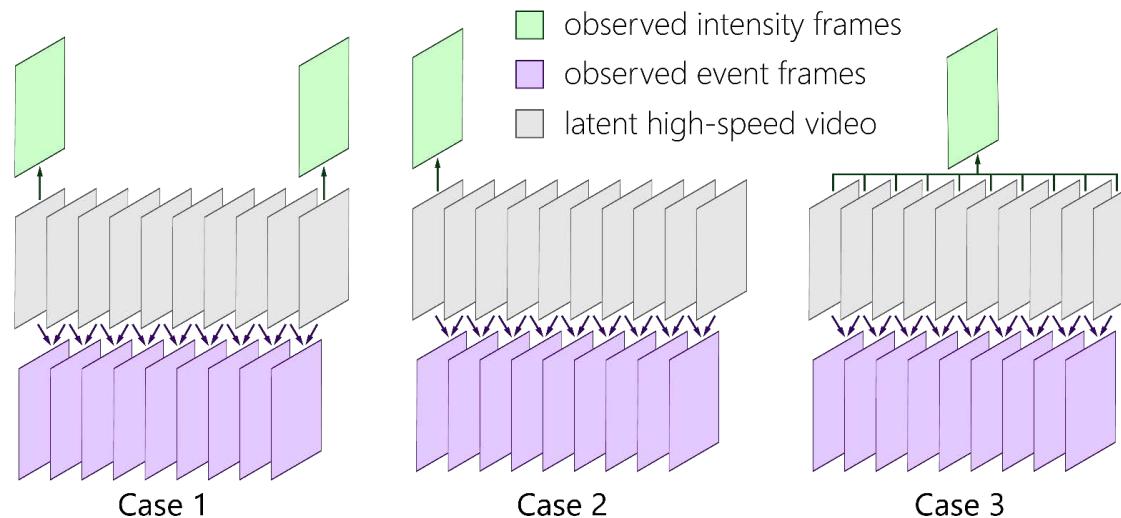
■ observed intensity frames
■ observed event frames



Differentiable inverse solver

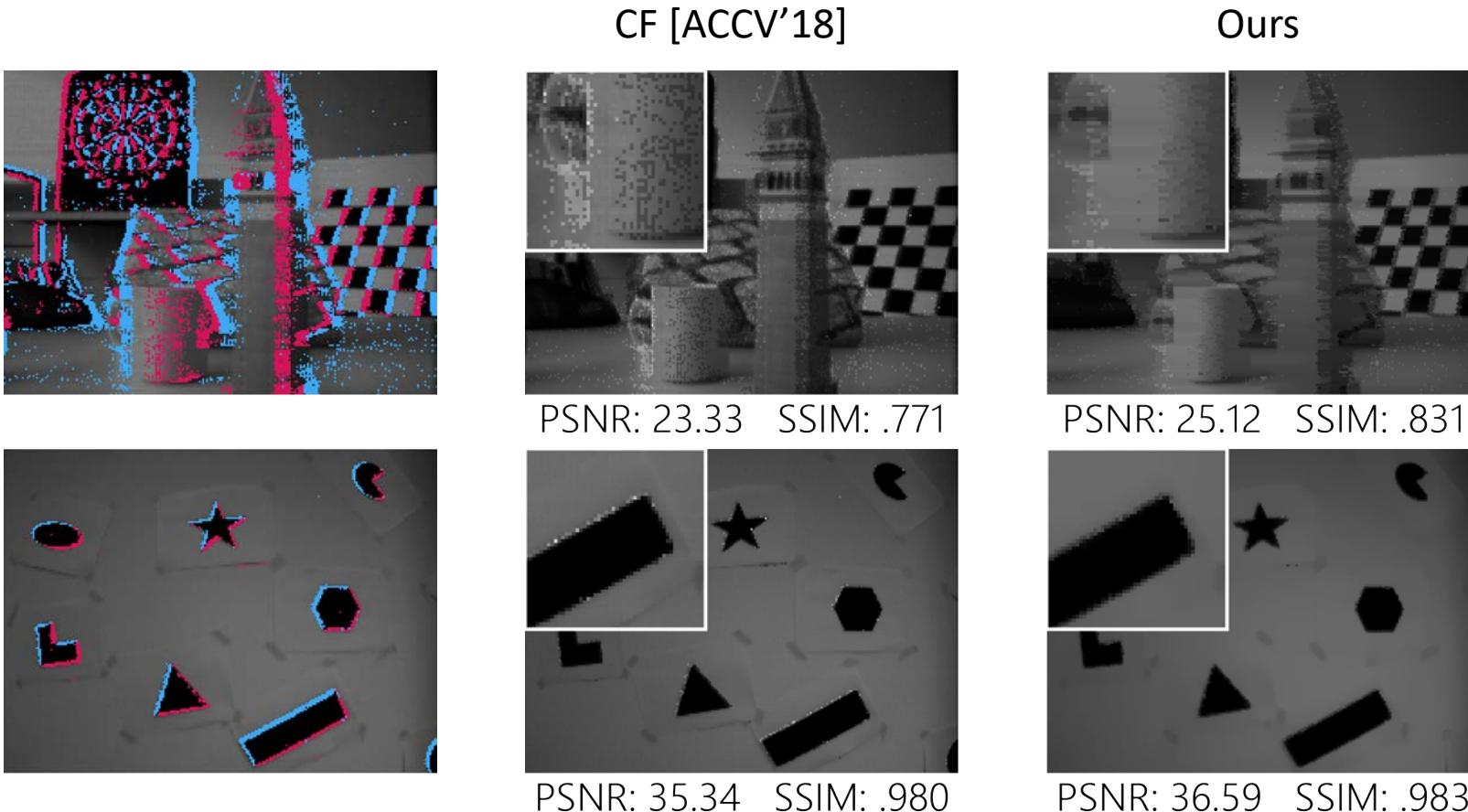
Variants of diff. inverse solvers

Differentiable inverse solvers are flexible to solve different fusion settings



Results

- Prediction case
 - Given start frame and future events, recover future frames



Results

- Interpolation case
 - Given start & end frames + events in-between, recover intermediate frames

Low-speed intensity frames
(2 frames)



[This is a video]

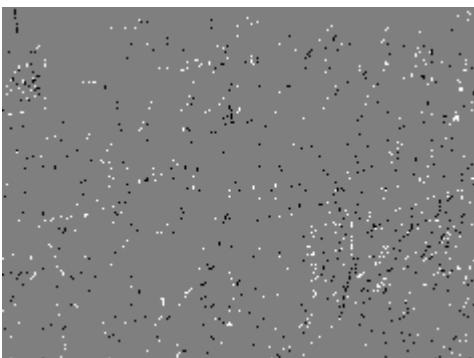
High-speed video
(21 frames)



[This is a video]

The middle frame is withheld for evaluation

[This is a video]

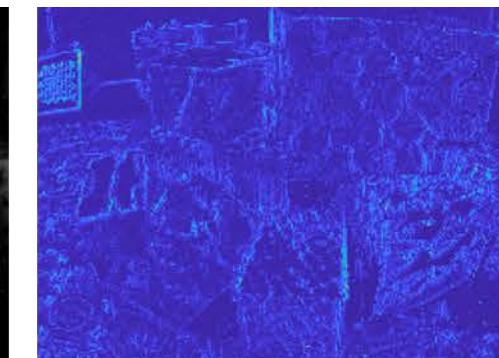


Event frames (20 frames)



PSNR: 33.41 SSIM: .955

Frame #10



0.0 0.25

Error map of Frame #10

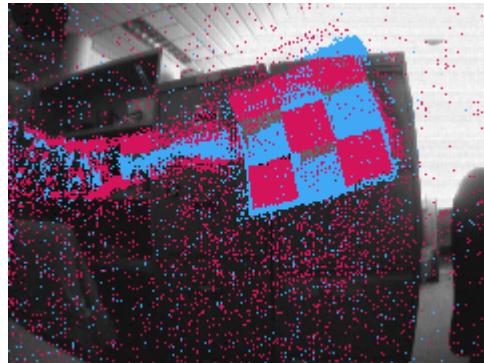
Results

- Motion deblur case
 - Given a blurry image + events in-exposure, recover intermediate sharp frames.

Blurry images



Events during exposure



EDI [CVPR'19]



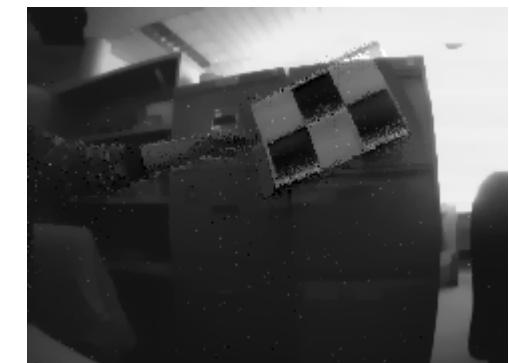
Ours



Reconstruction



[This is a video]



[This is a video]

Synergy: diff. inverse solver + residual learning

Imaging hardware

Synergistic models

task

Physics-based model

Learning-based model

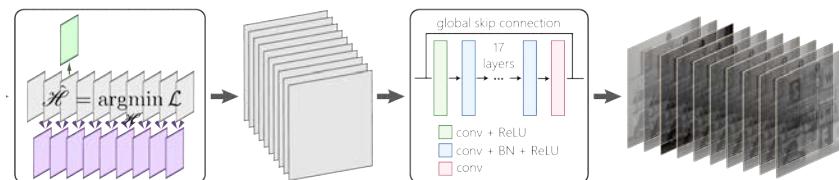
#1 APS + DVS camera

Differentiable inverse solver

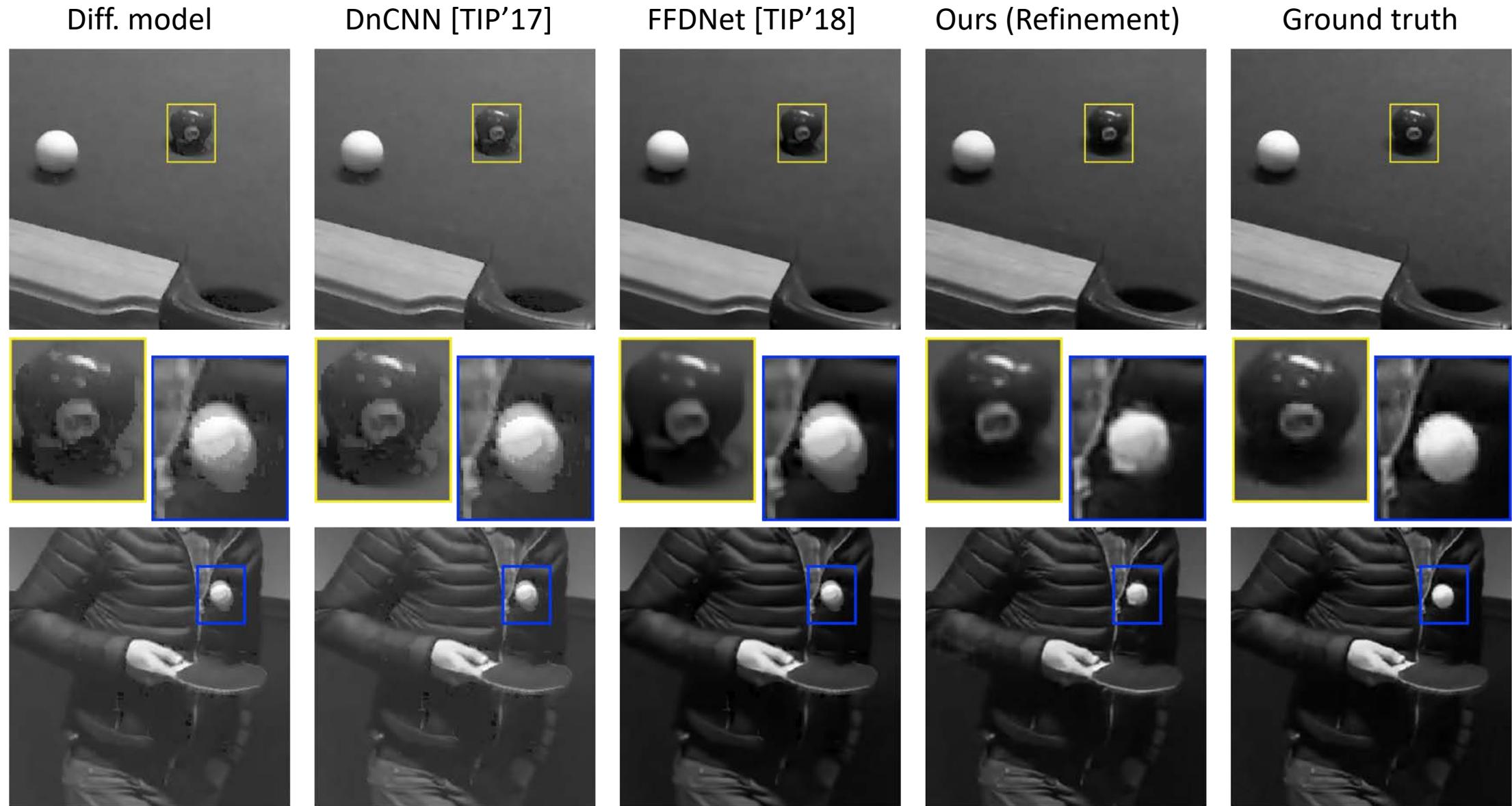
*Video frame synthesis:
Interpolation,
extrapolation,
motion deblur*

Augmentation

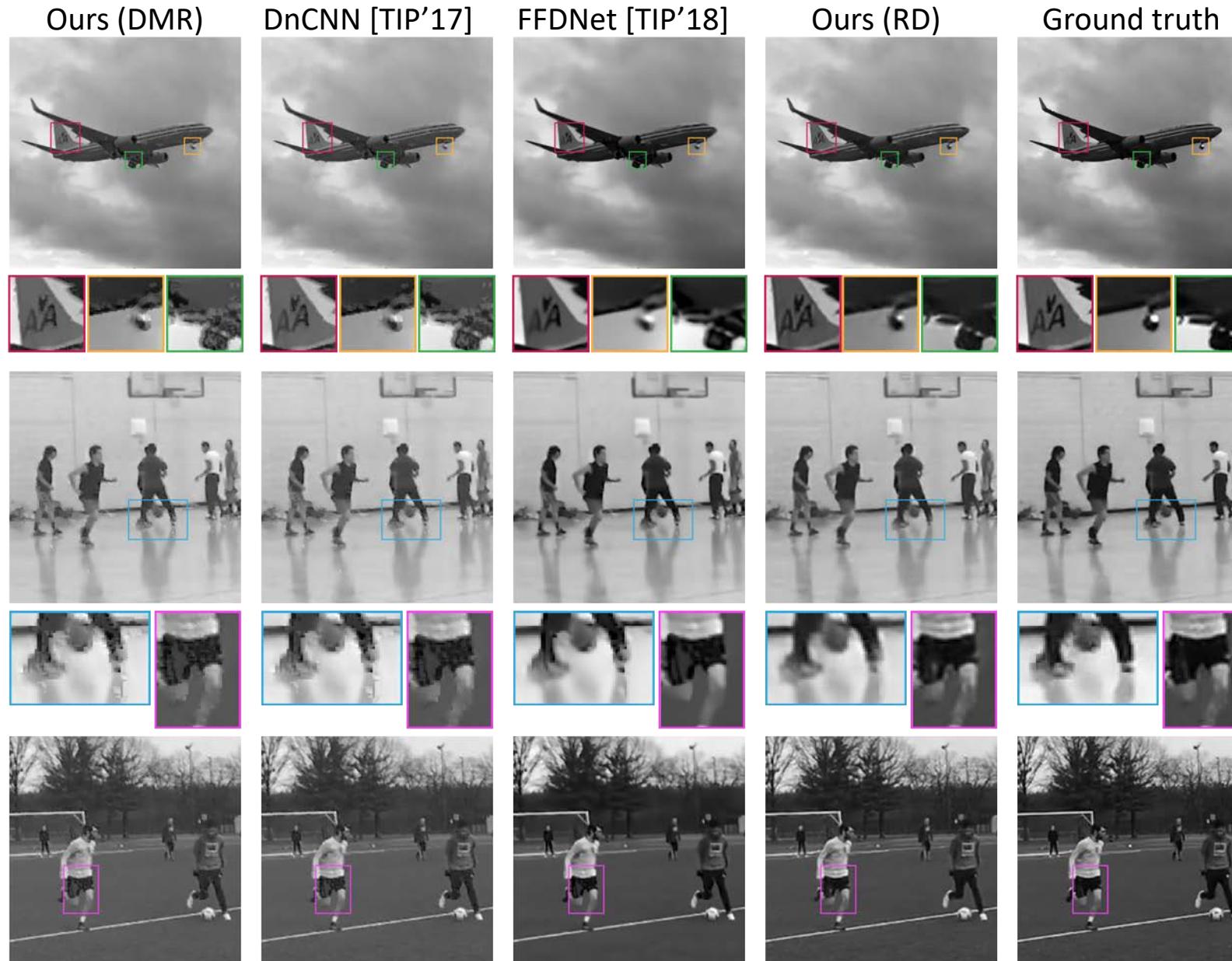
Residual refinement



Results (residual refinement)



More results (residual refinement)



clip name	metric	DMR	DnCNN	FFDNet	Ours
airplane	PSNR	30.91	31.10	30.92	31.38
	SSIM	0.975	0.982	0.976	0.982
basketball	PSNR	23.55	24.05	23.47	24.06
	SSIM	0.963	0.971	0.964	0.972
soccer	PSNR	29.96	31.08	30.13	31.29
	SSIM	0.961	0.974	0.962	0.975
billiard	PSNR	36.46	35.42	36.48	36.46
	SSIM	0.982	0.986	0.983	0.987
ping pong	PSNR	32.46	32.26	32.50	32.24
	SSIM	0.974	0.978	0.975	0.979

Results

Comparison with non-event-based frame interpolation approach.

Interpolation using APS-only

SepConv [CVPR'17]



[This is a video]

Ground truth



[This is a video]

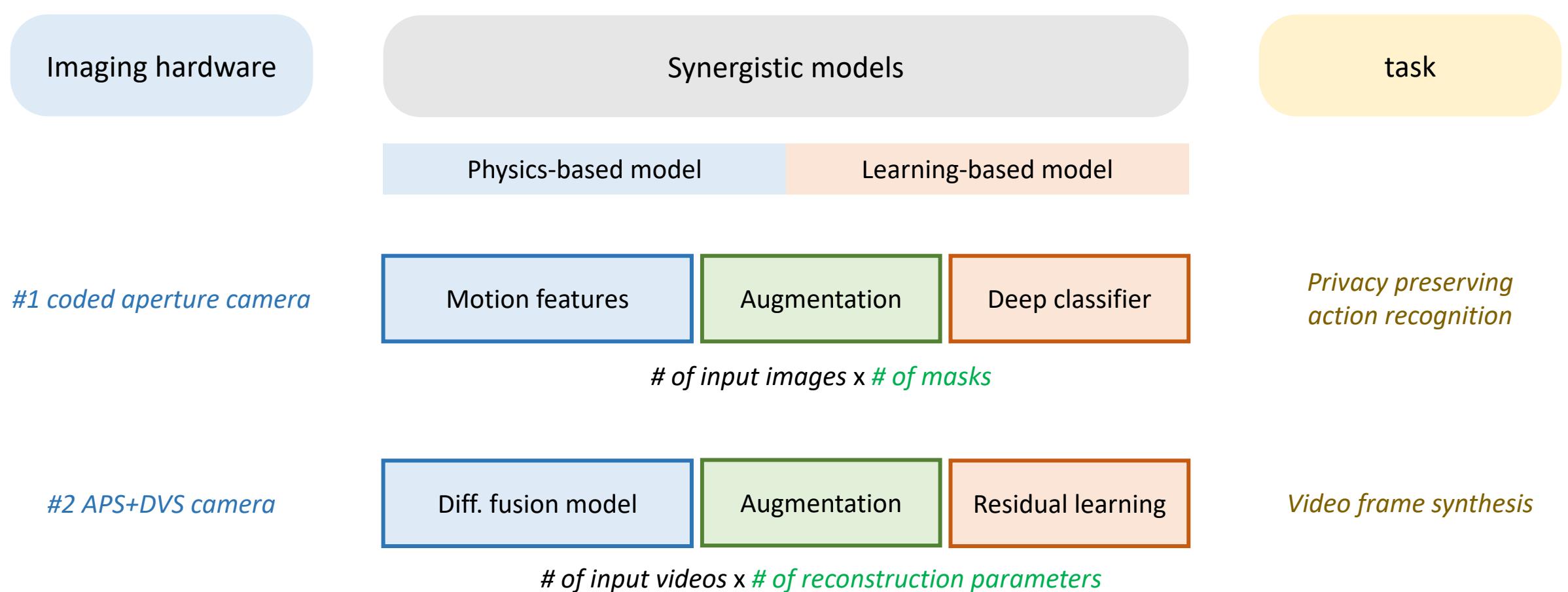
Interpolation using APS + DVS

Ours (DMR + Refinement)

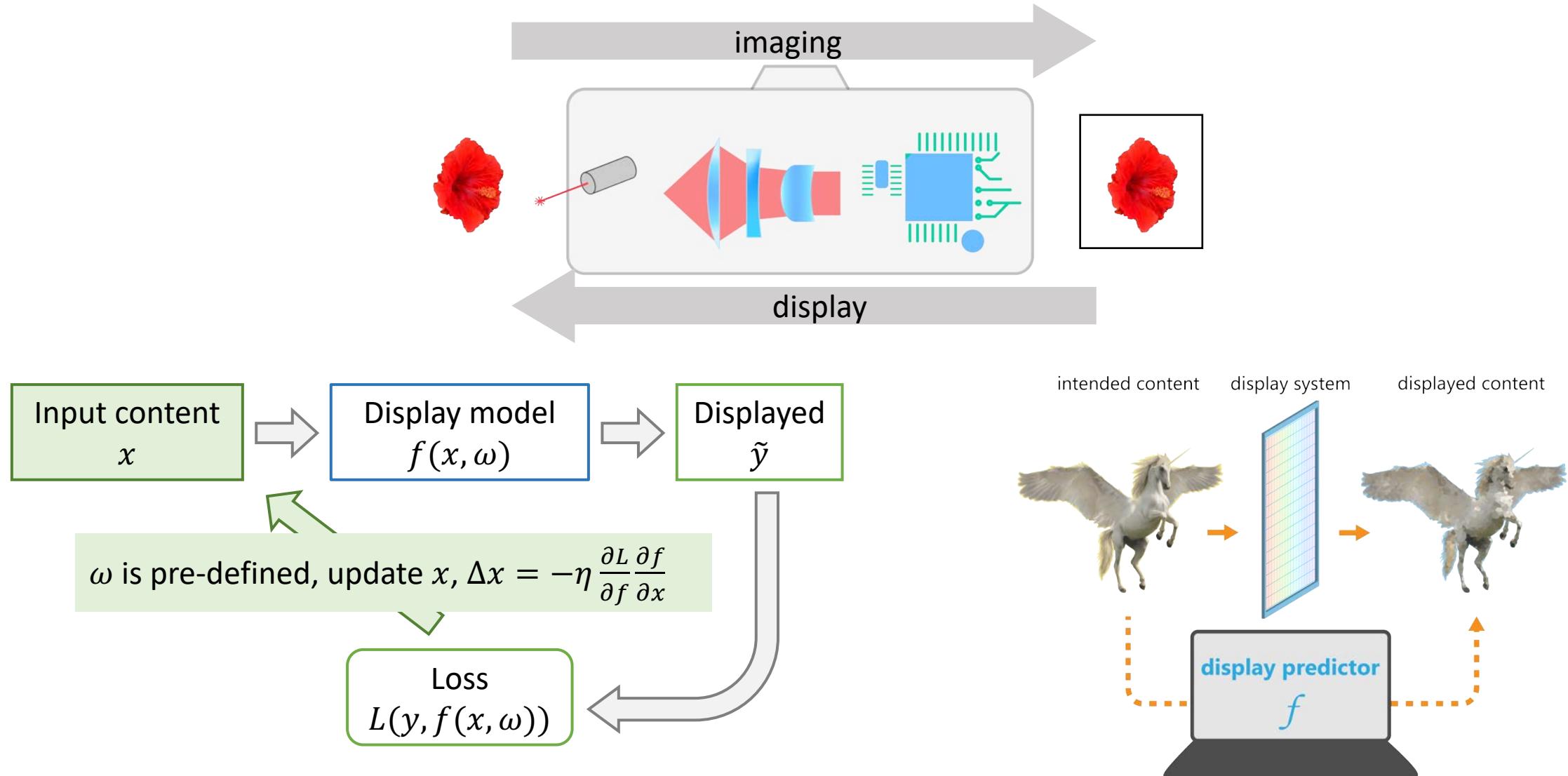


[This is a video]

Physics-based models as means of data augmentation

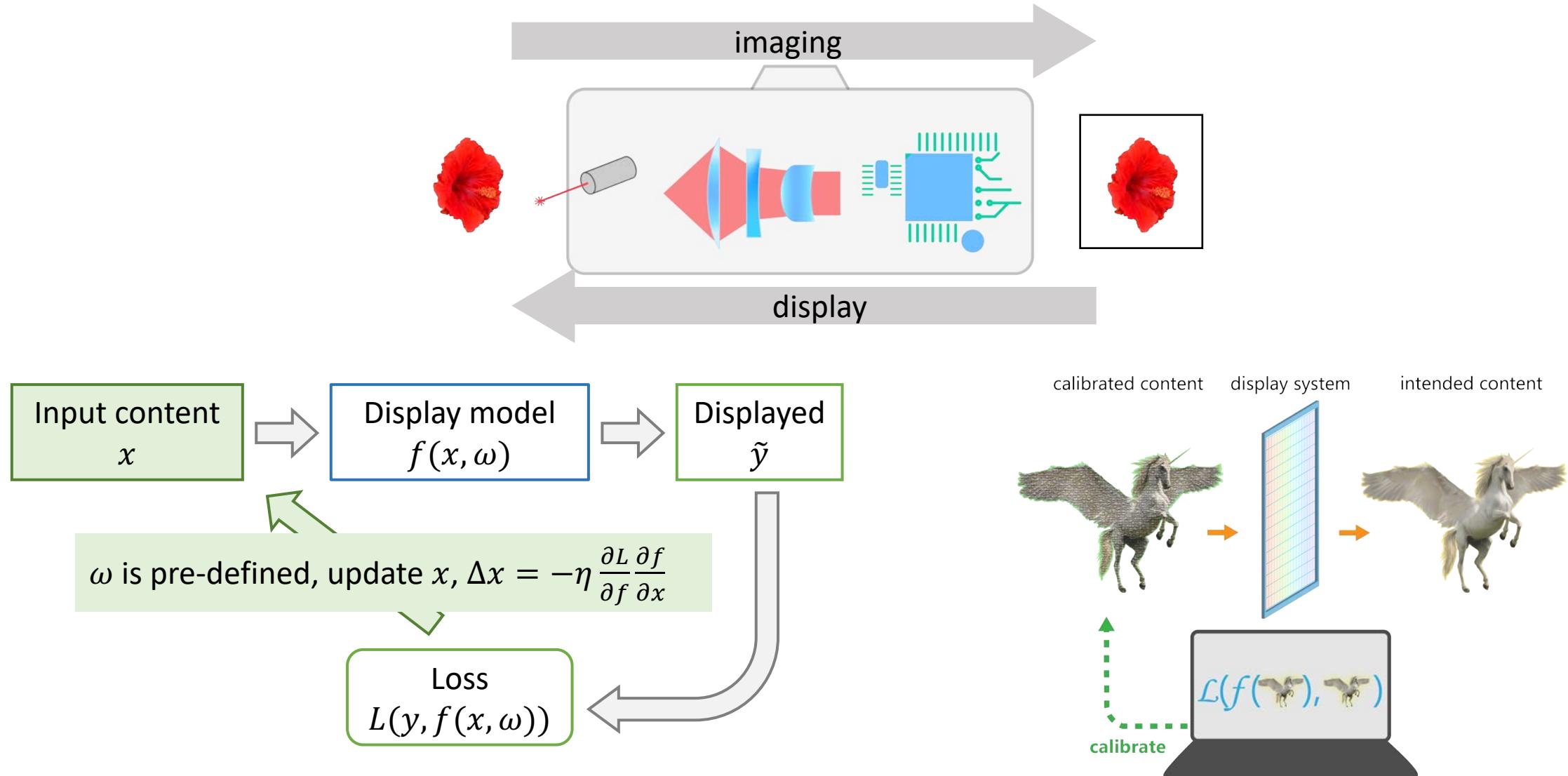


Computational display: a reverse process of imaging



Differentiable display model: optimize content to display

Computational display: a reverse process of imaging



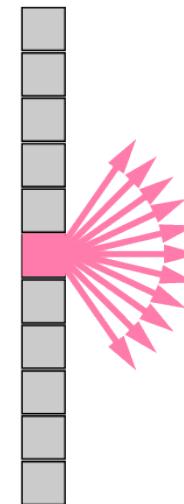
Differentiable display model: optimize content to display

What are we talking about when we talk about 3D display? – parallax



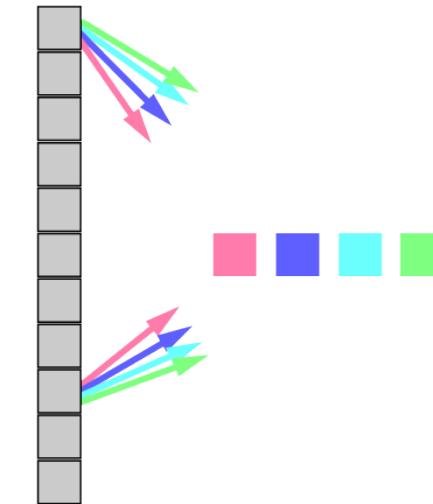
[video of a recorded volume hologram]

2D display



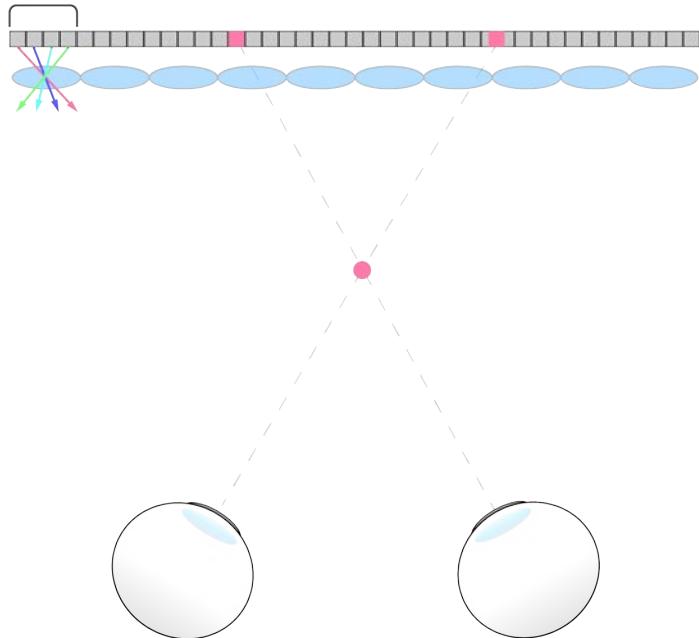
Program 2D pixels

3D display



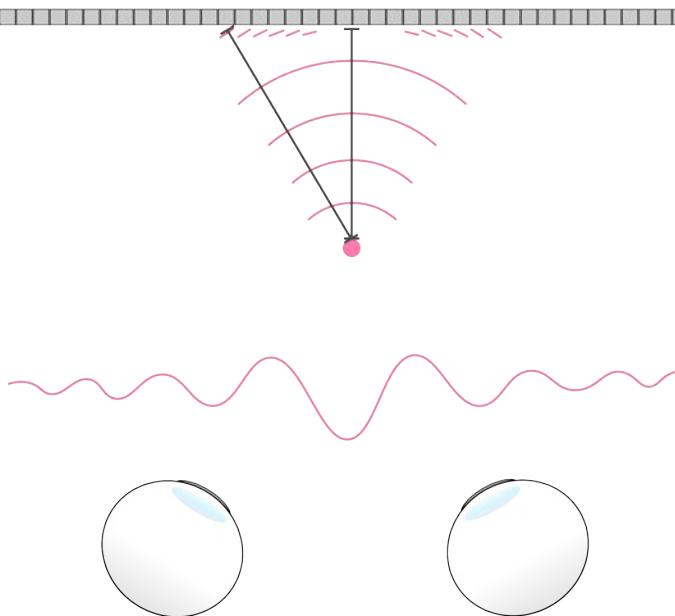
Program 4D rays
- 2D in spatial
- 2D in angular

Ray programming: light field vs. holographic display



Light field display:

- Use microlens/pinhole array to bend light



Holographic display:

- Use coherent light and phase-based Spatial Light Modulators

The total number of pixels (spatial-bandwidth product) remains unchanged:

- 2D display: $N \times N$ pixels
- 3D display: $(N/M \times N/M) \times (M \times M)$

M: angular resolution

Challenges:

- **computation:** $10^{12} - 10^{14}$ pixel/s for 3D iPad or 3D iMac, (not for HMD)
 - Currently can achieve $10^9 - 10^{11}$
- **hardware:** pixel pitch, form factor etc.

[Calculation adapted from Yamaguchi, M. JOSAA 2016]

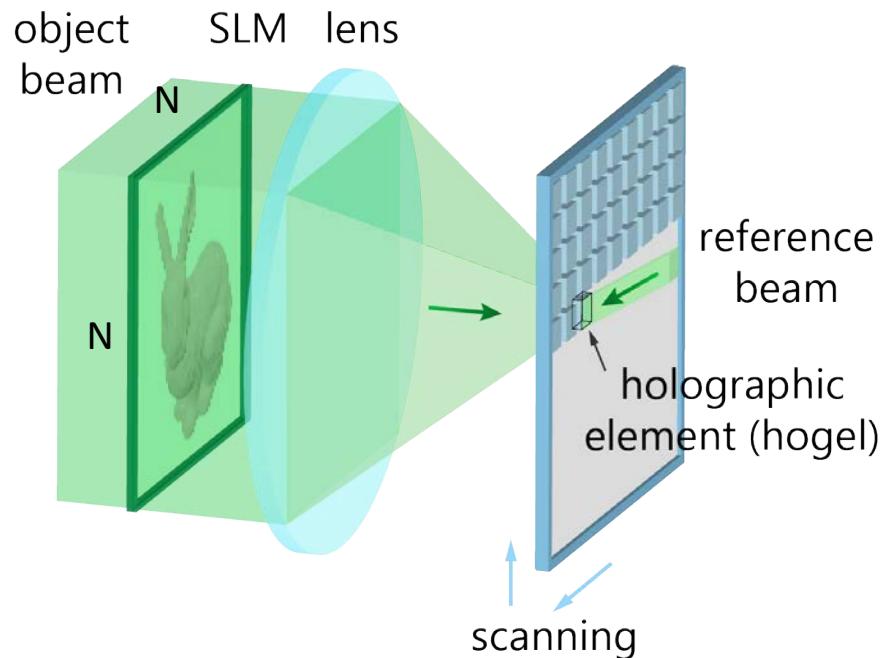
Can we program $N \times N \times M \times M$ rays using $N \times N$ pixels?



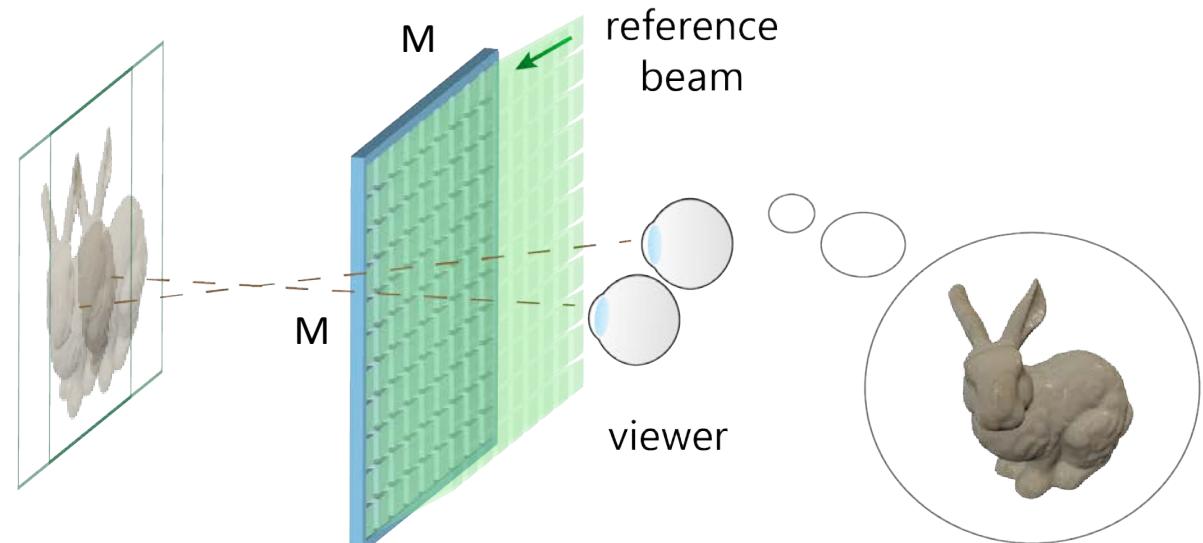
Holographic stereogram

Holographic stereograms

Recording (preparing the HS)



Playback (viewing the HS)



SLM has $N \times N$ pixels, and the stereogram has $M \times M$ hogels.

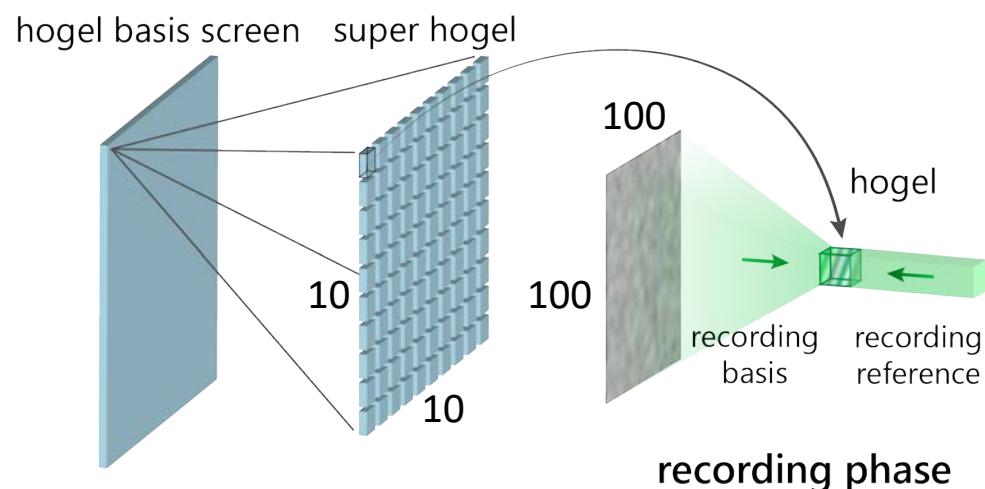
Total rays: $N \times N \times M \times M$

But this is static 3D image.

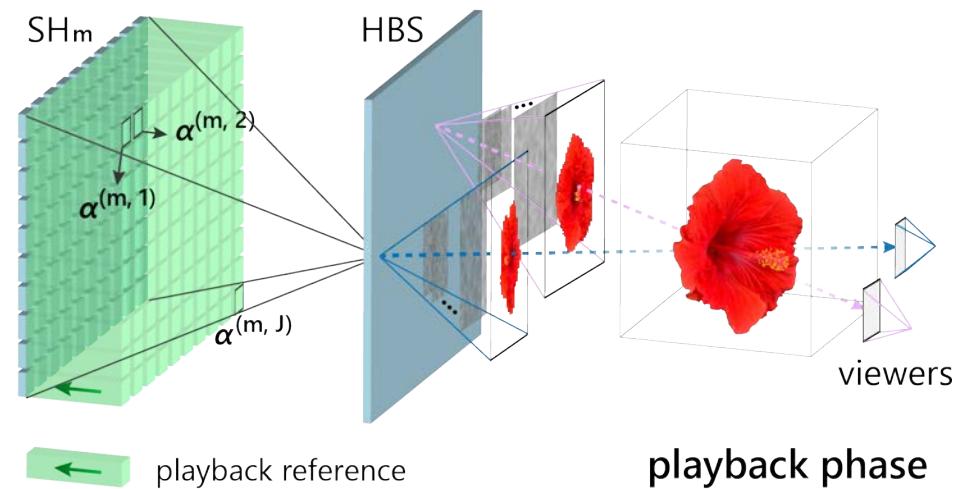


Add programmability to the ref. beam

Hogel basis display



HBS: 1000×1000 hogels = 100×100 super hogels (SH)
 Each hogel records a 100×100 basis pattern



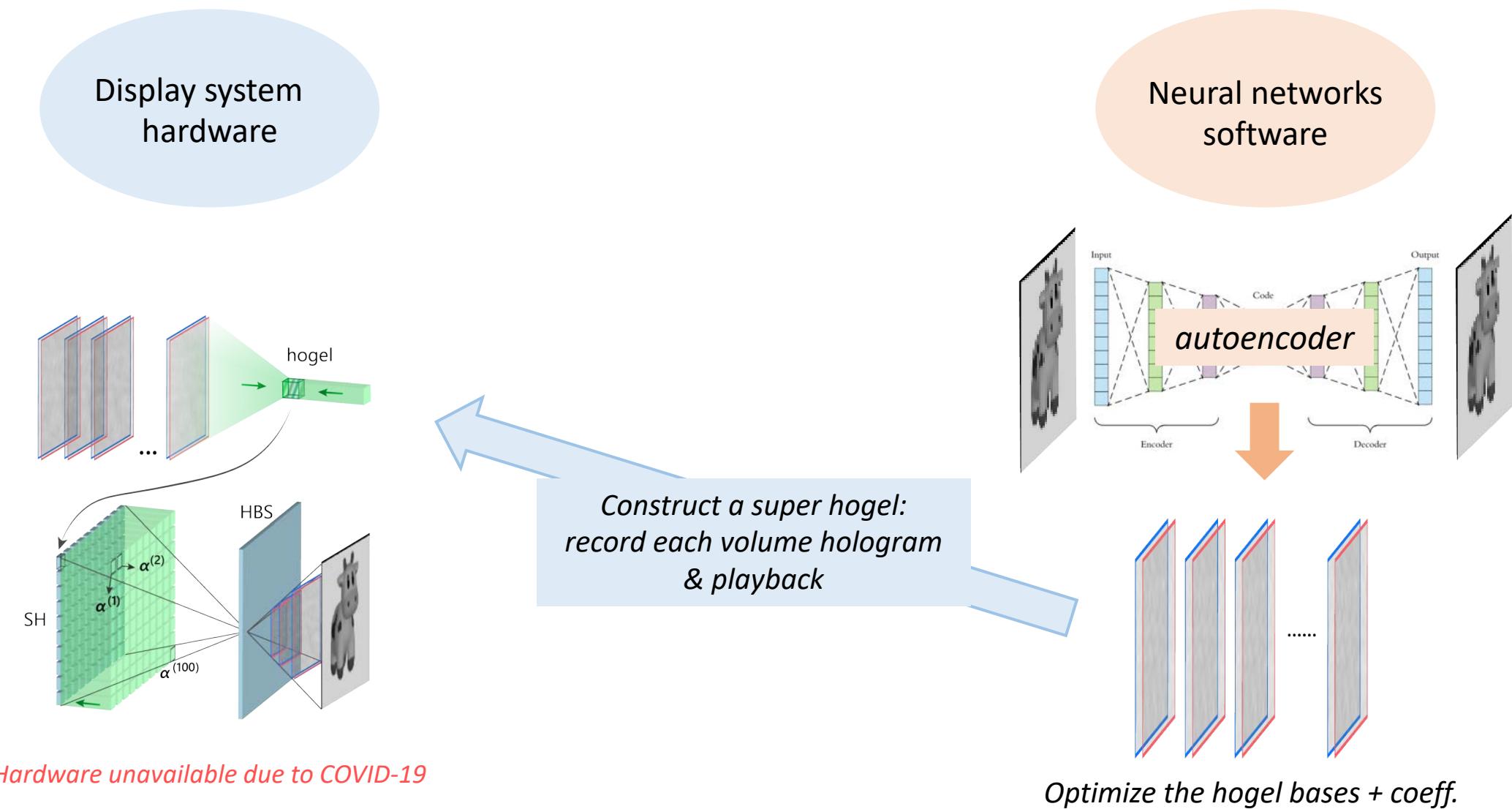
One SH produce one viewpoint image.

$$I_{SH_m} = \left| \sum_{j=1}^J \alpha^{(m,j)} B^{(j)} \right|^2$$

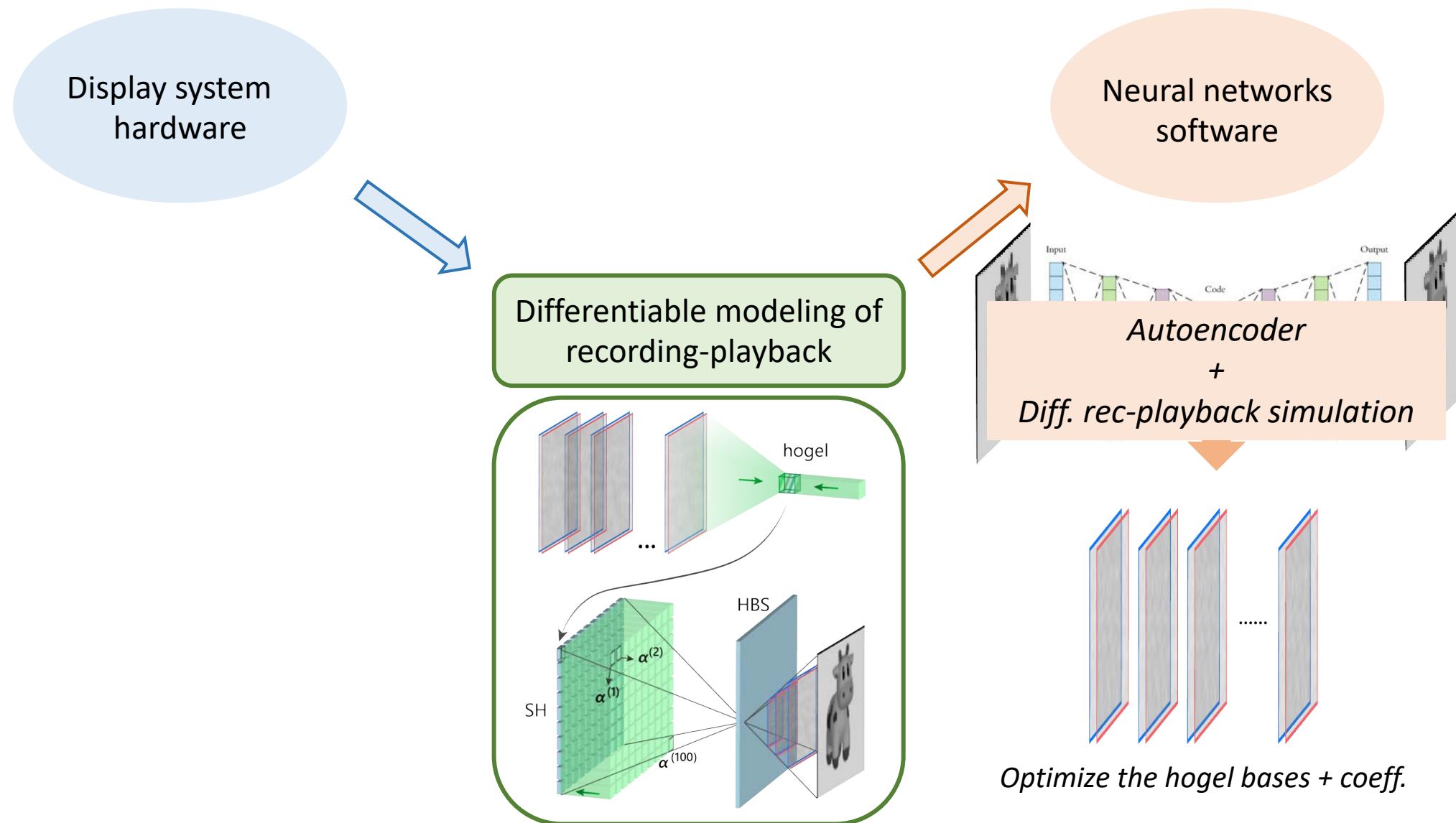
*Spatial averaging for SH:
 viewing distance should satisfy SH size within visual acuity.*

HBS: 1000×1000 ref. pixels, i.e. $\alpha^{(m,j)}$
 can produce 100×100 views;
 $100 \times 100 \times 100 \times 100$ rays.

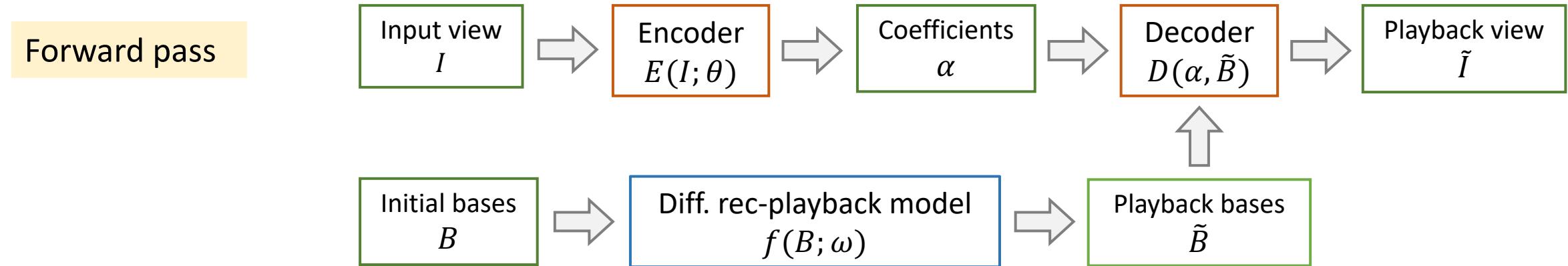
One super hogel (SH)



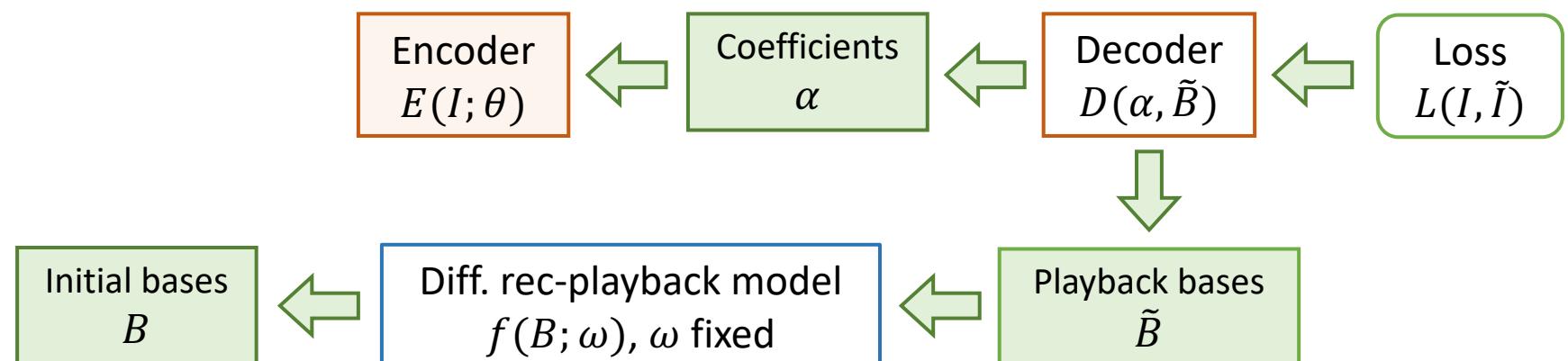
One super hogel (SH)



Synergistic model for (super) hogel basis learning



Backprop.

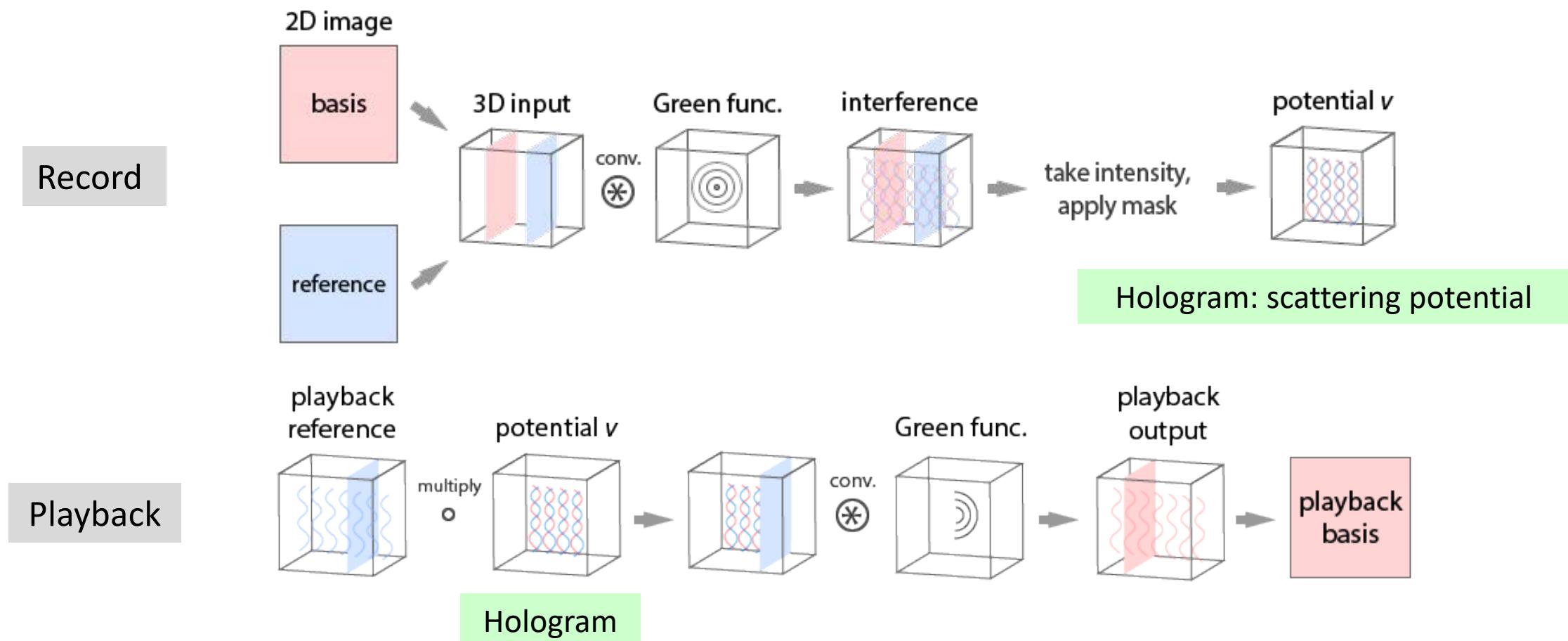


(Physics-in-the-loop autoencoder)

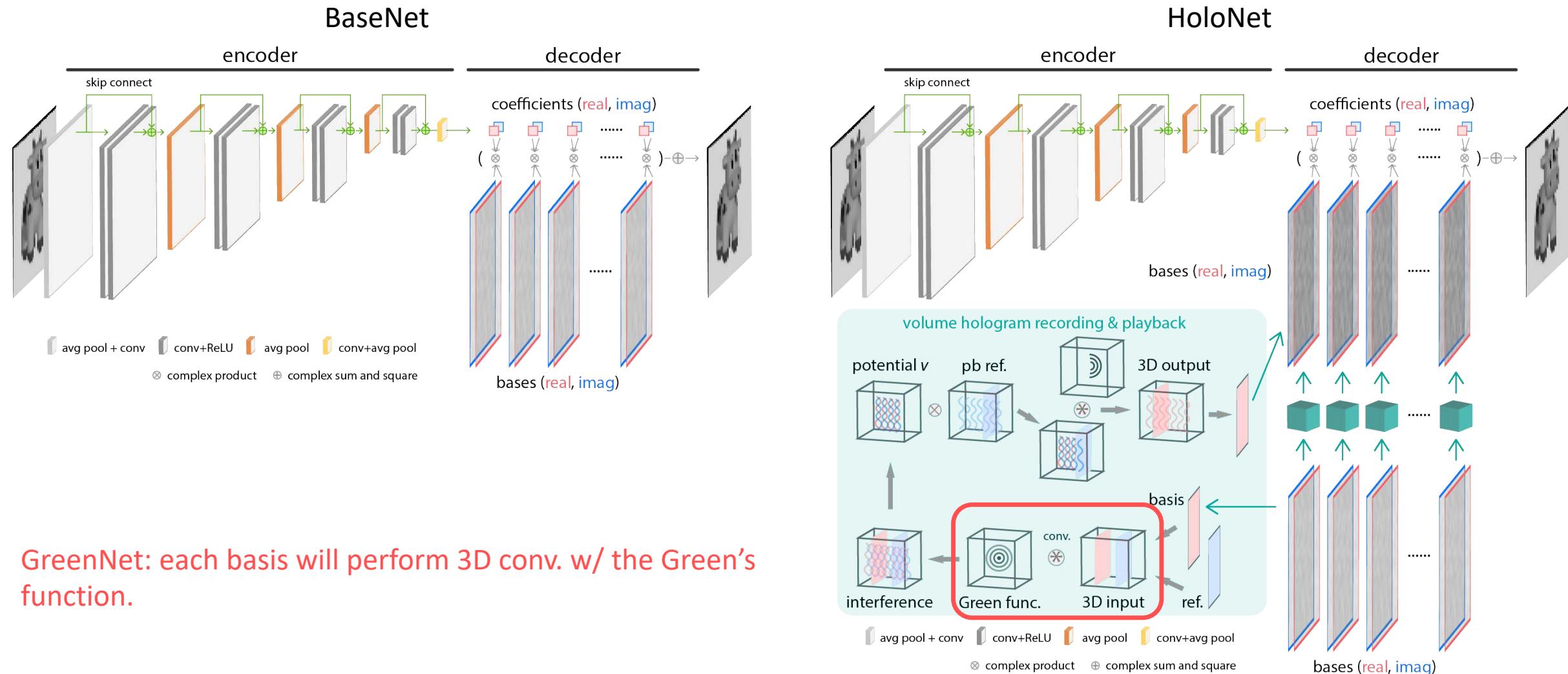
Physics-based (recording-playback) model

Solve the inhomogeneous Helmholtz equation, based on Born approx.

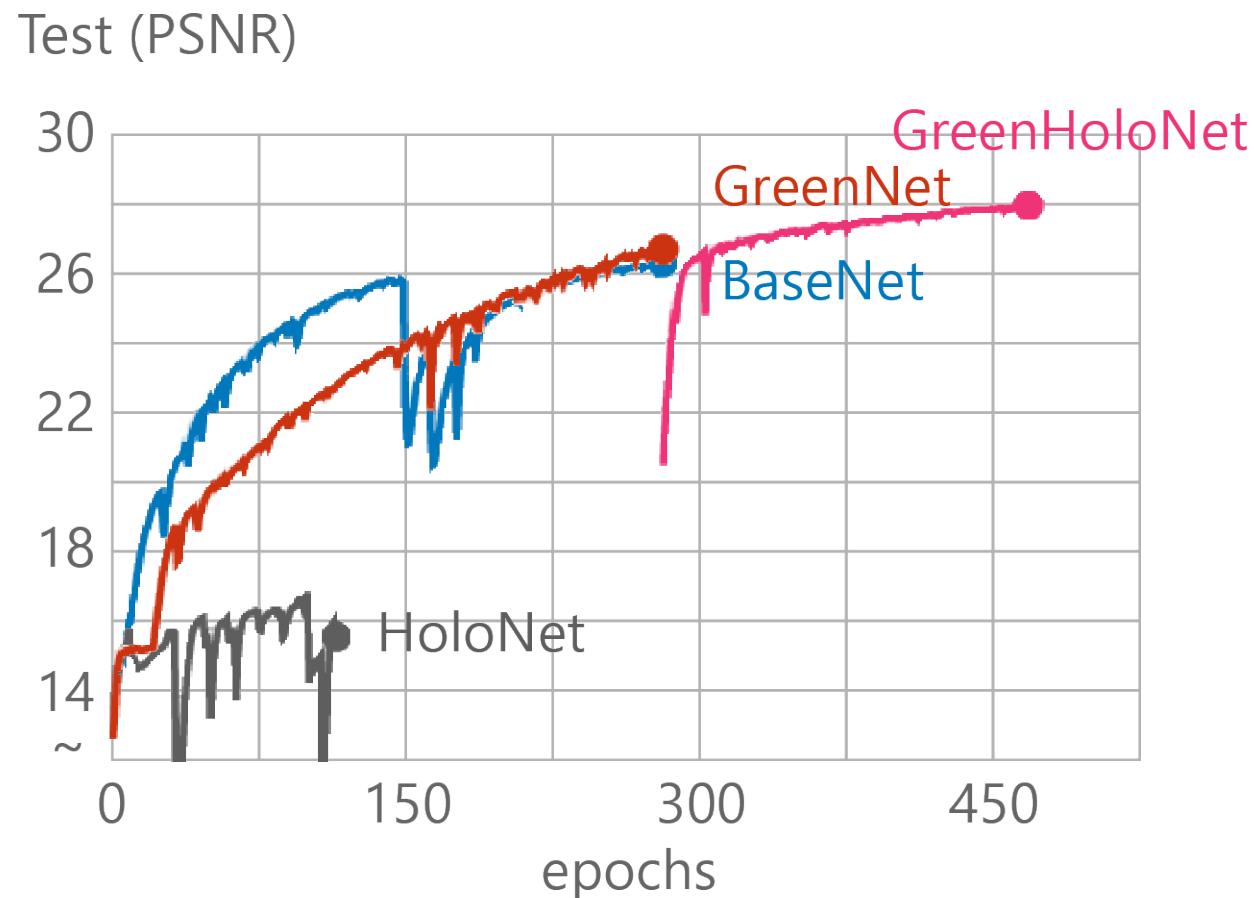
$$\text{3D convolution with the Green's function (PSF): } G(x) = \frac{\exp(ik_0|x|)}{|x|}$$



Learning-based model vs synergistic model

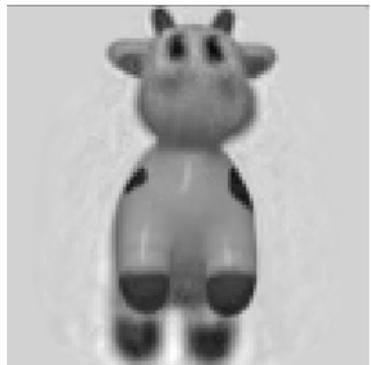


Results

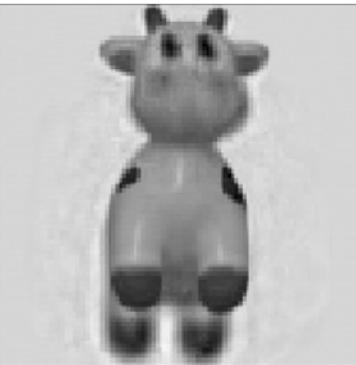


Results

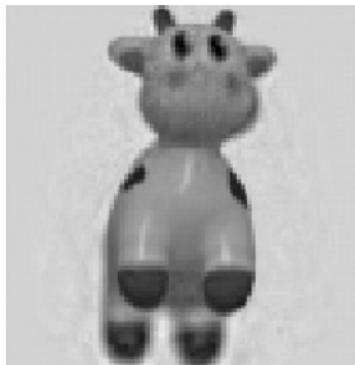
BaseNet



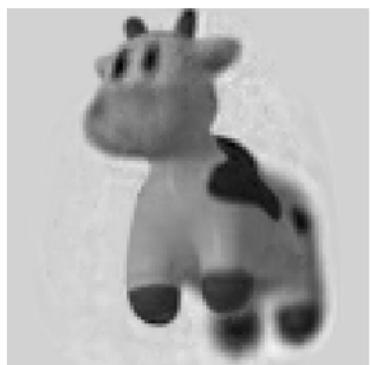
GreenNet



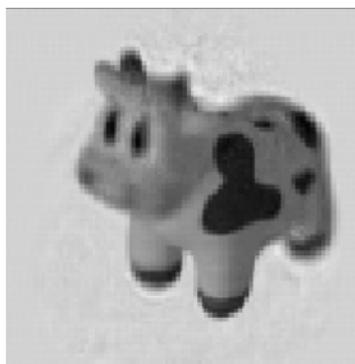
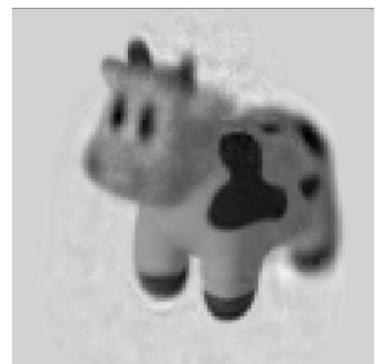
GreenHoloNet



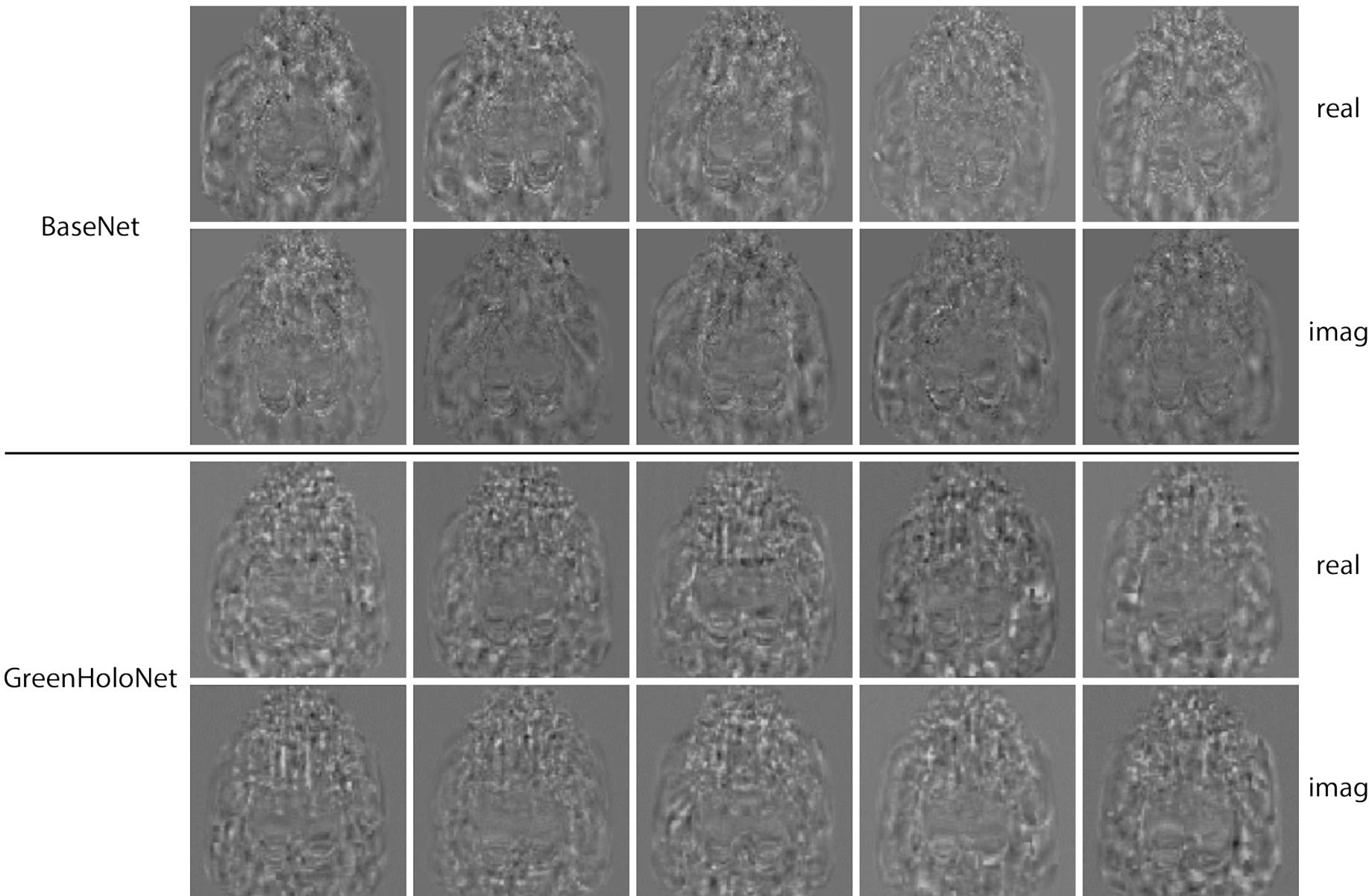
Ground truth



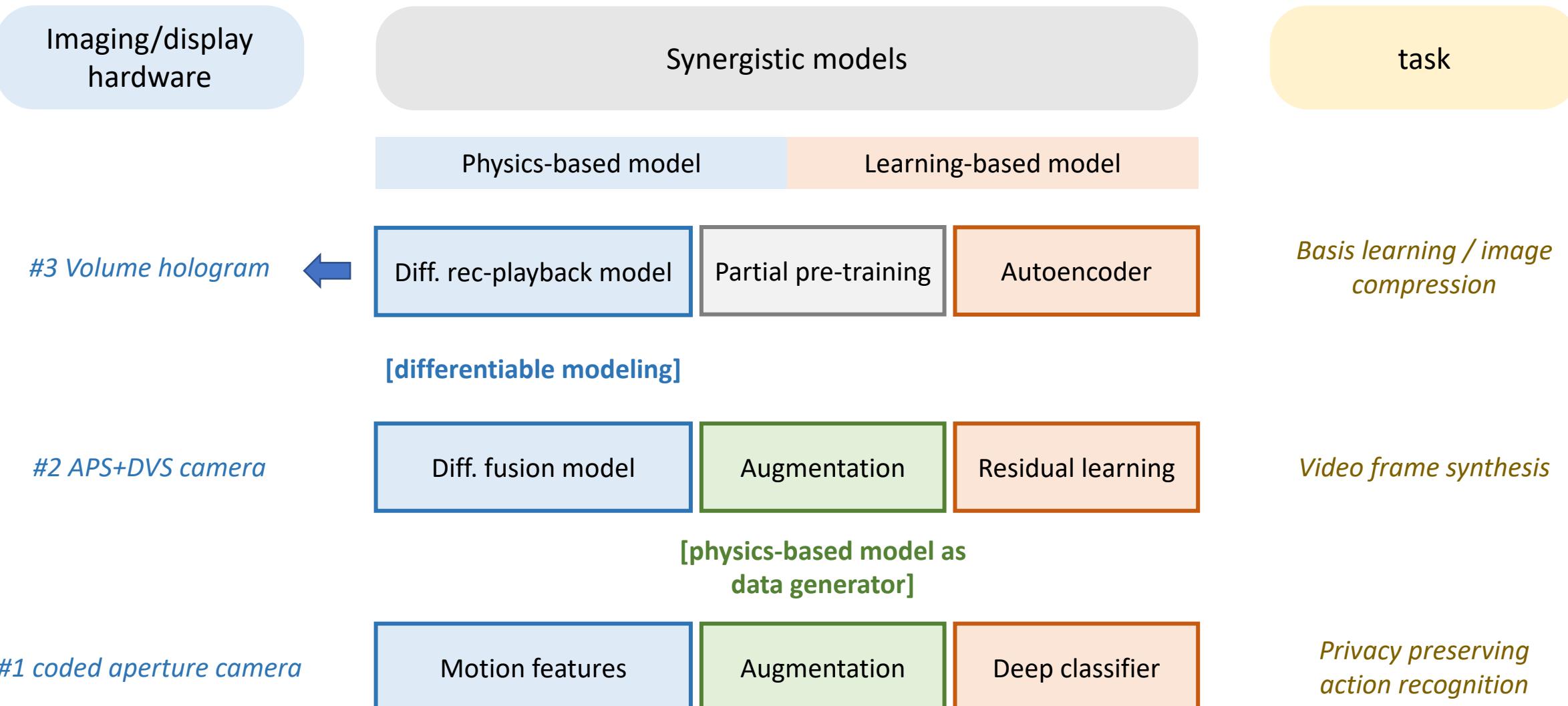
Network	PSNR (dB)
BaseNet	26.1
GreenNet	26.5
GreenHoloNet	27.9



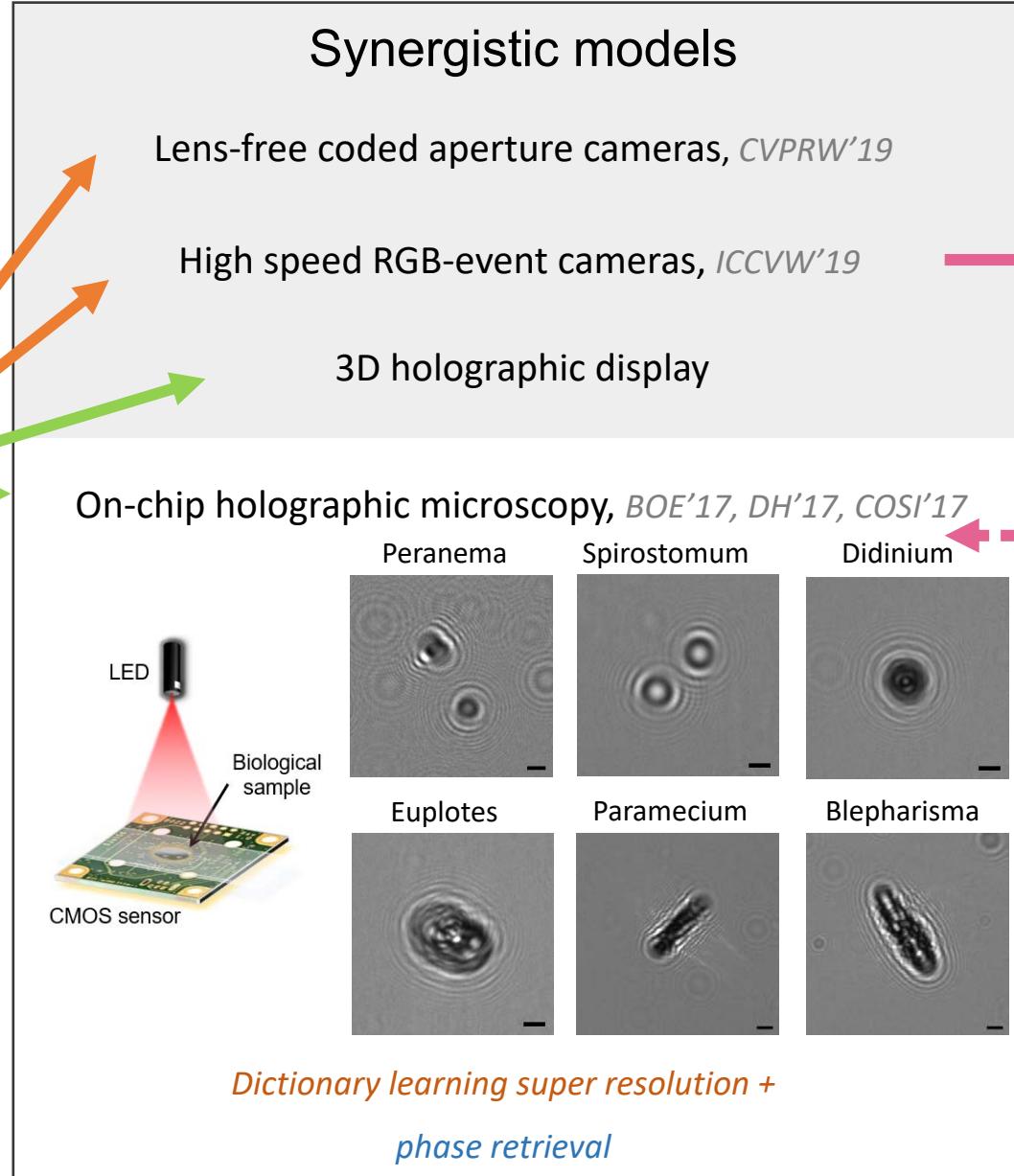
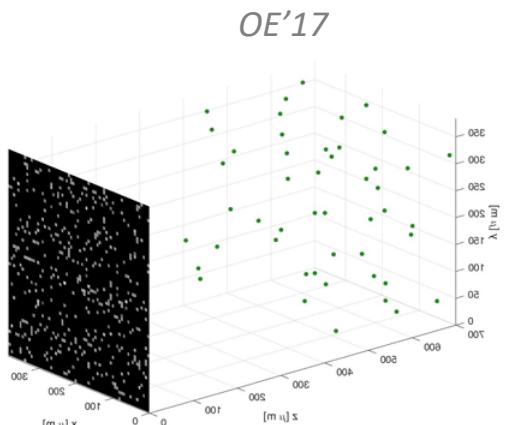
Visualization of example learned bases



Summary

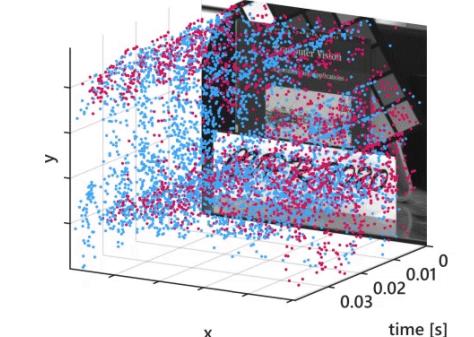


Holographic compressive video



Guided event filtering

CVPR'20 & journal extension
image+events



Event denoising and super resolution w/
HR image guidance



Lossy event compression

Srutarshi et al. in submission '20

Acknowledgement

Advisor: Ollie Cossairt,

Committee: Aggelos Katsaggelos, Jack Tumblin, Nathan Matsuda

Internship mentors: Dikpal Reddy (Light co.), Sing Bing Kang (Microsoft Research), Xiaokai Li (Apple Inc.)

Coauthors: Lei Tian (Boston U.), Roarke Horstmeyer (Duke U.), Boxin Shi (PKU), Leo Spinoulas, Kuan He, Donghun Ryu, George Chen, Qiqin Dai, Weixin Jiang, Francesco Pittaluga, Vibhav Vineet, Sudipta Sinha, Srutarshi Banerjee

Coauthors-to-be: Prasan Shedligeri, Hamid Hasani, Florian Willomitzer, Florian Schieffers

Comp-photo-lab members

My parents!