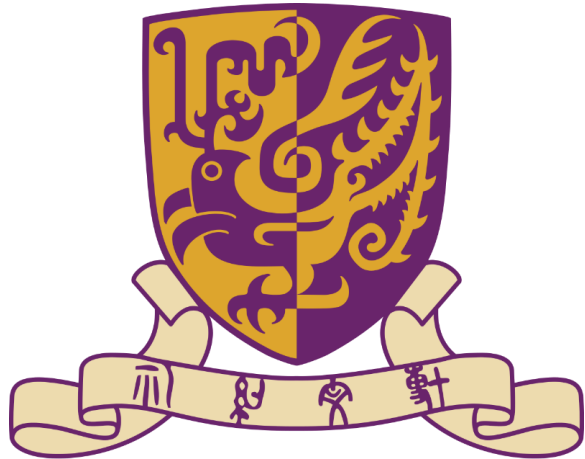# How to make a smart camera pipeline

Tianfan Xue

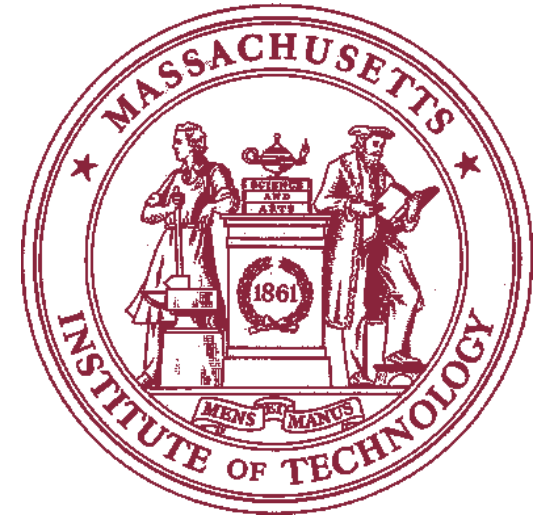The Chinese University of Hong Kong

# About myself



The Chinese Univ. of HK

Assist Prof.

Google Research, Gcam

Staff Eng.

Manager: David Saleson

2017-2022

MIT

Ph.D.

Advisor: William T. Freeman

2012-2017

# Smart camera = ML algorithms applied images?


Object detection


Face beautification
[Leyvand et al., 2006]

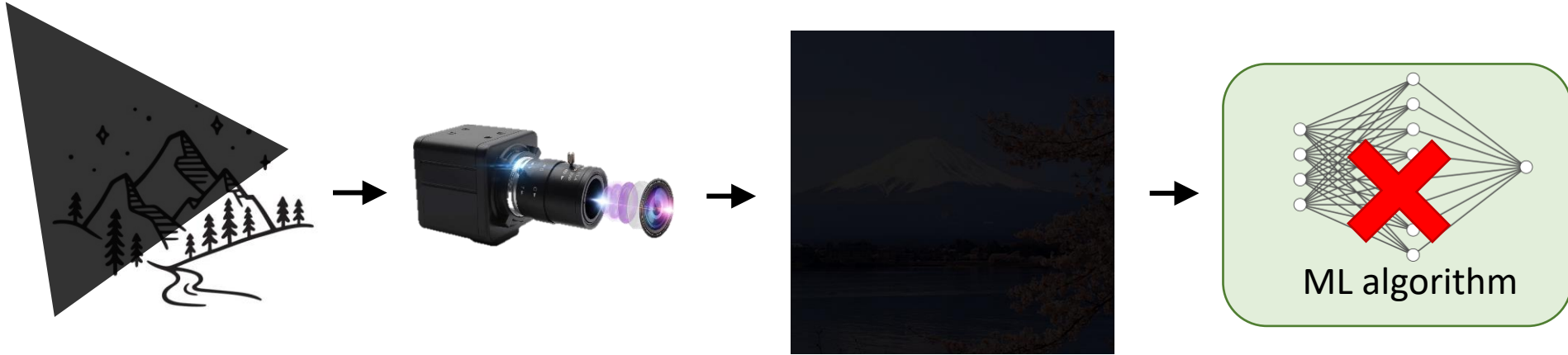# Camera may not capture visual signal for ML system
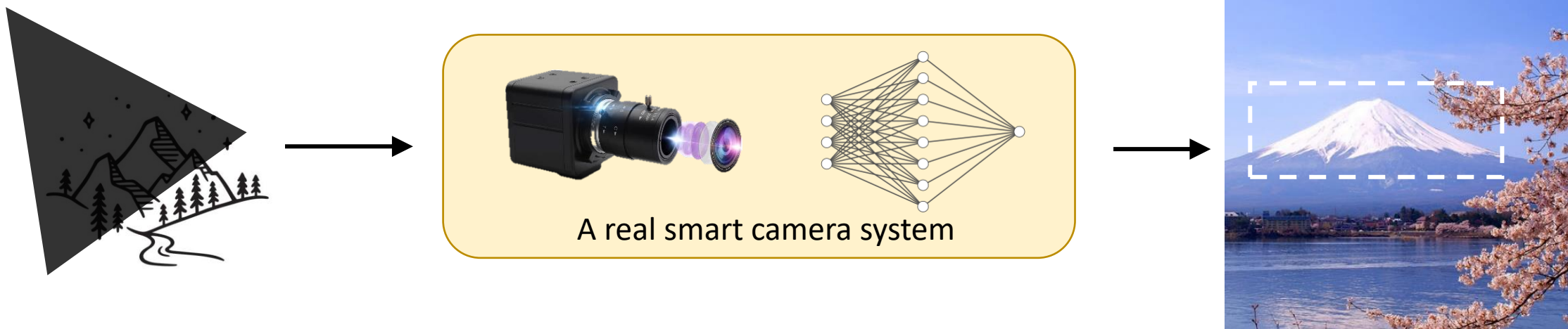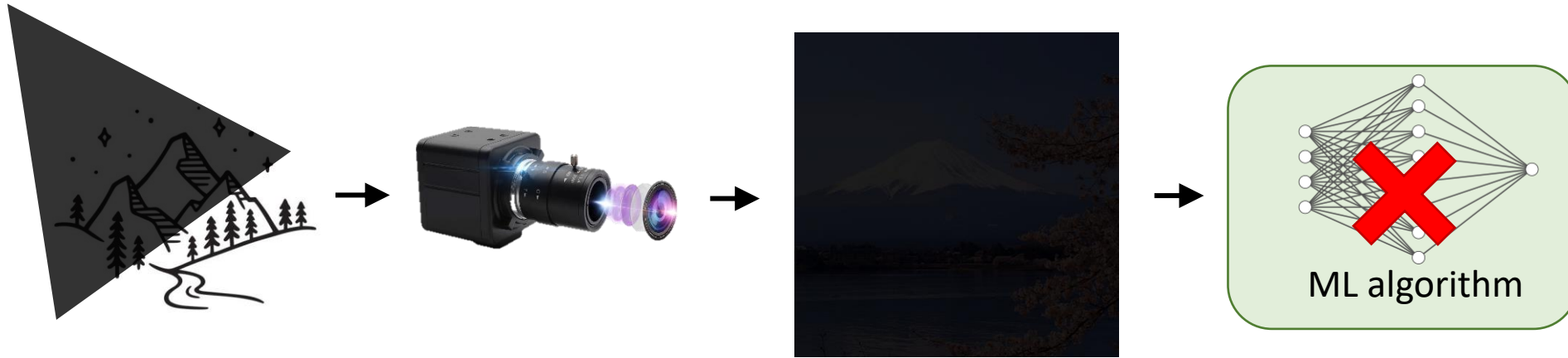


One frame

A frame 2s later

**Images are too dark** for ML to detect this biking person

*Video from a fatal car crash*

# Separate design may fail
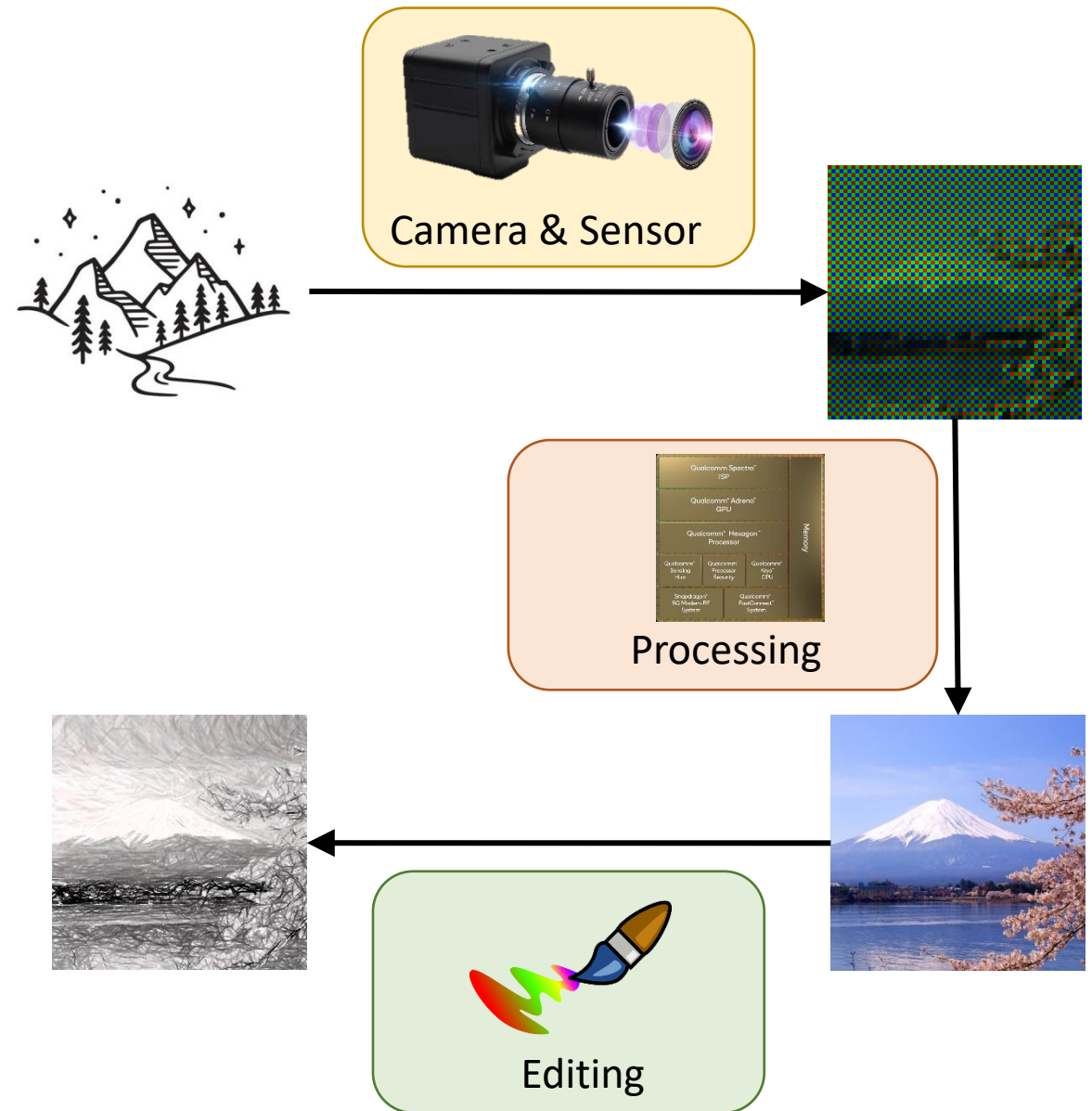


ML algorithm

# Machine learning embedded in the camera

# Overview

- **Capturing**: multiple source fusion

- **Processing & editing**

  - **Training data:** synthetic data

  - **Network**: combine classic image processing algorithm and machine learning

# Overview

- **Capturing**: multiple source fusion

- **Processing & editing**

  - **Training data:** synthetic data

  - **Network**: combine classic image processing algorithm and machine learning
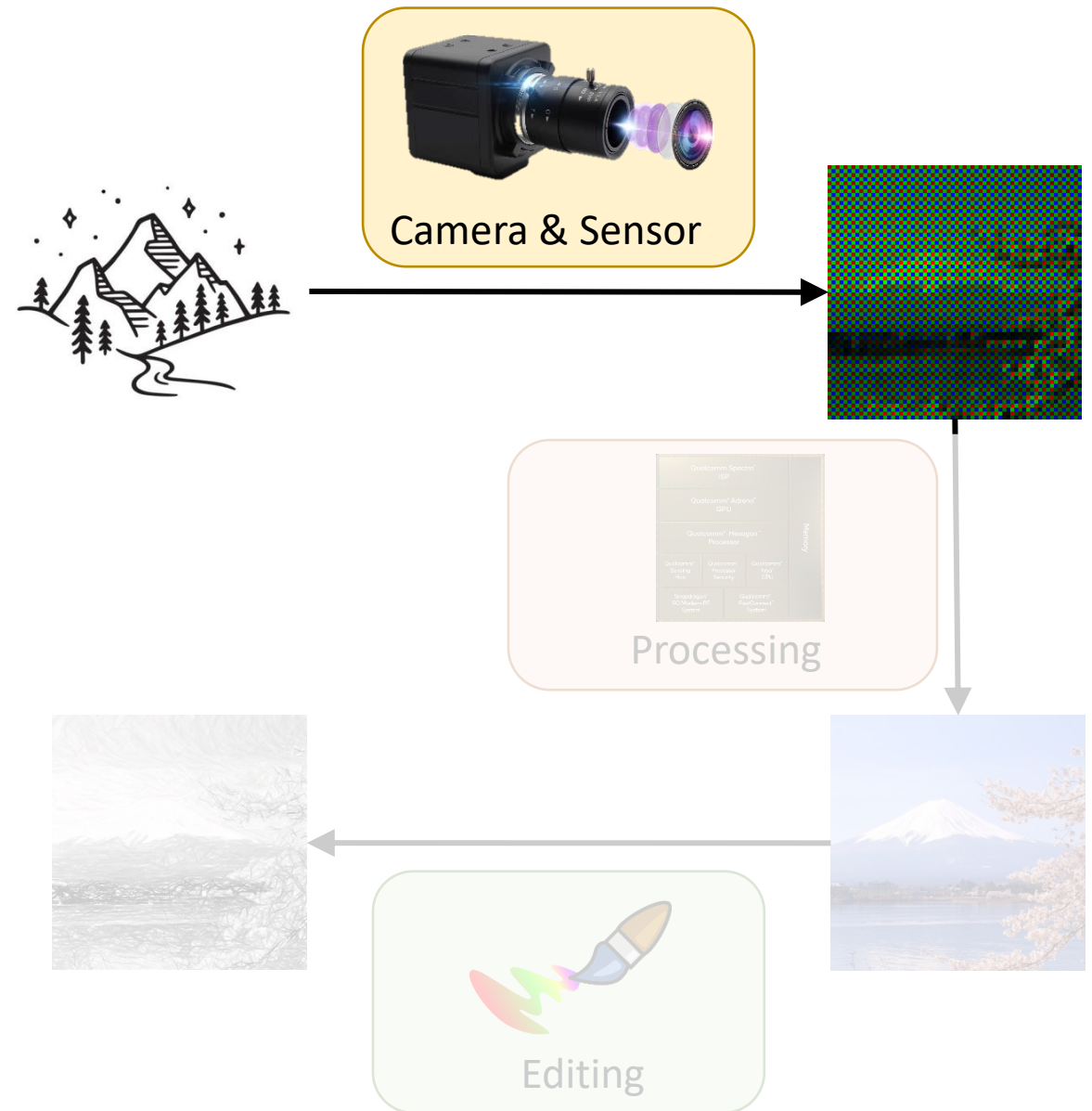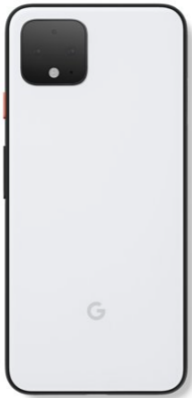


Camera & Sensor

Processing
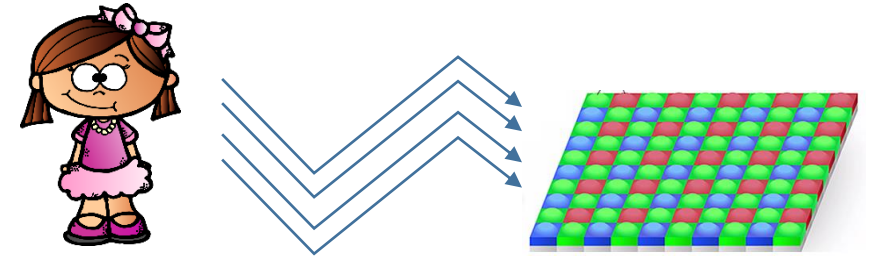
Editing

What is this object?

Image in lowlight

"Handheld Mobile Photography in Very Low Ligh", SIGGRAPH Asia 2019

Flower
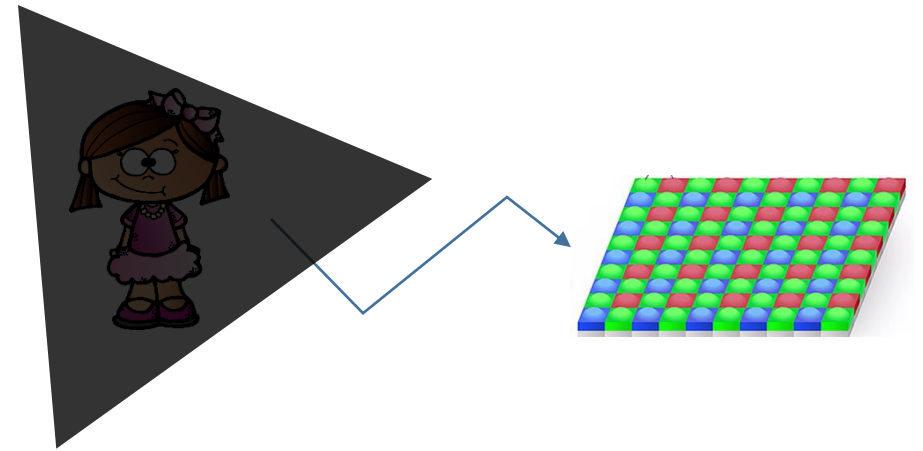
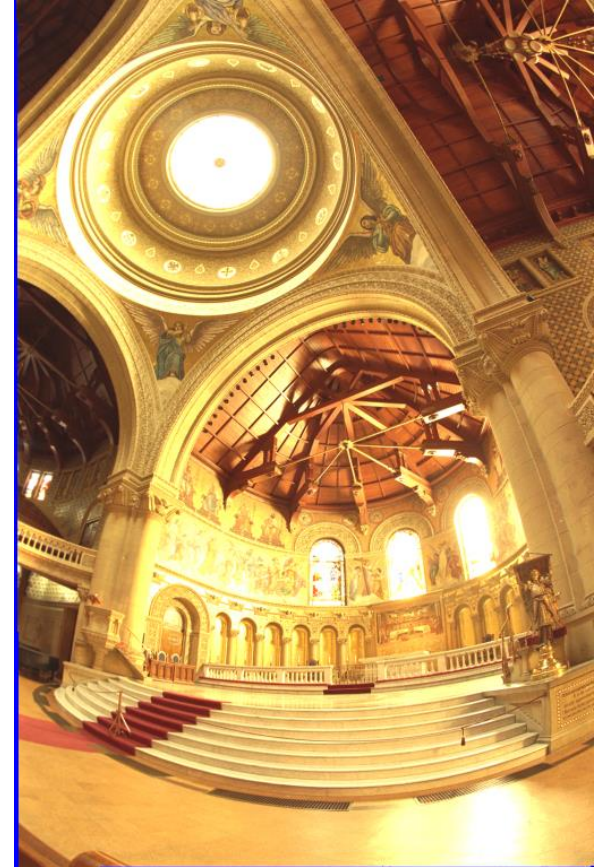by our night sight algorithm
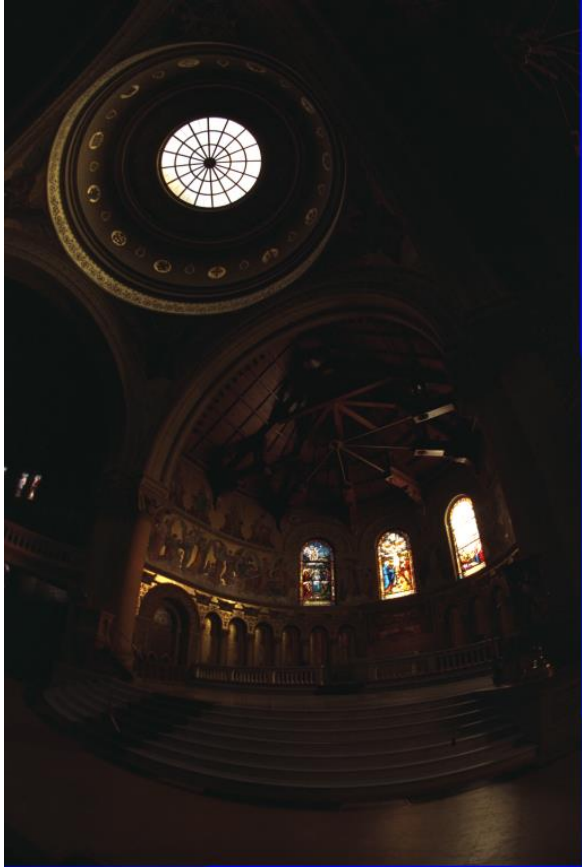
# Not enough photons in lowlight



Good lighting

Lowlight

# Exposure bracketing



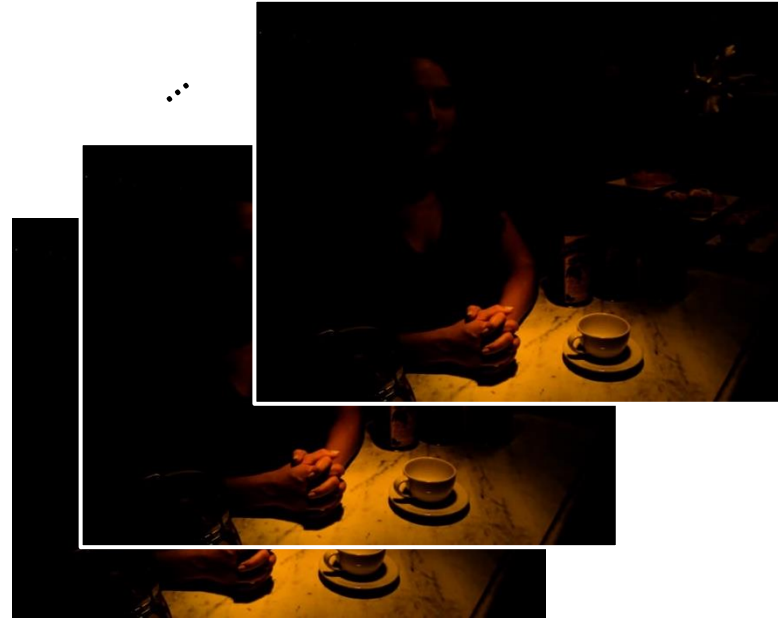[Debevec et al., 2007]
[Gallo and Sen 2016]

# From long exposure to burst photography



Long-exposure

Motion blur

Burst photography
**burst of short exposed frames**

# HDR+



S. Hasinoff et al, "Burst photography for high dynamic range and low-light imaging on mobile cameras ", SIGGRAPH Asia 2016.

# Night sight



O. Liba, et al., "Handheld Mobile Photography in Very Low Ligh", SIGGRAPH Asia 2019.

# Multiple captures also helps to remove reflection



Images with reflection

Reflection-free image

T. Xue, et al., "A Computational Approach for Obstruction-Free Photography", SIGGRAPH, 2015.

# Reflection removal using stereo input



2 frames from stereo camera → Output

S. Niklaus et al., "Learned dual-view reflection removal", WACV, 2021.
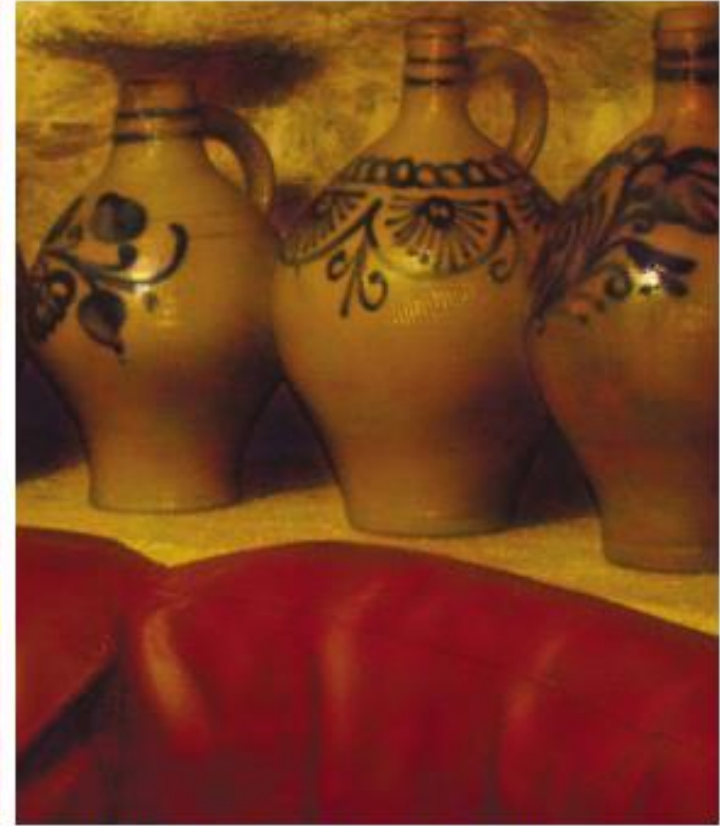
# Flash / Non-flash Photography



Flash          No-Flash          Detail Transfer with Denoising

G. Petschnigg et al., "Digital Photography with Flash and No-Flash Image Pairs", SIGGRAPH 2004.

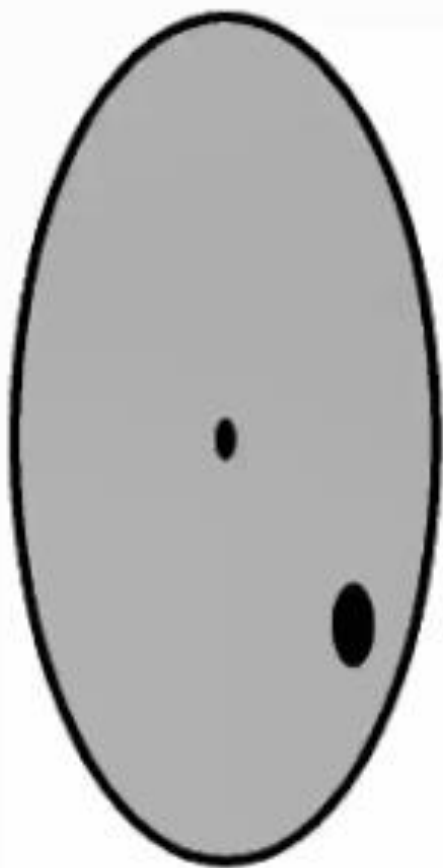# Stereoscopic Dark Flash for Low-light Photography



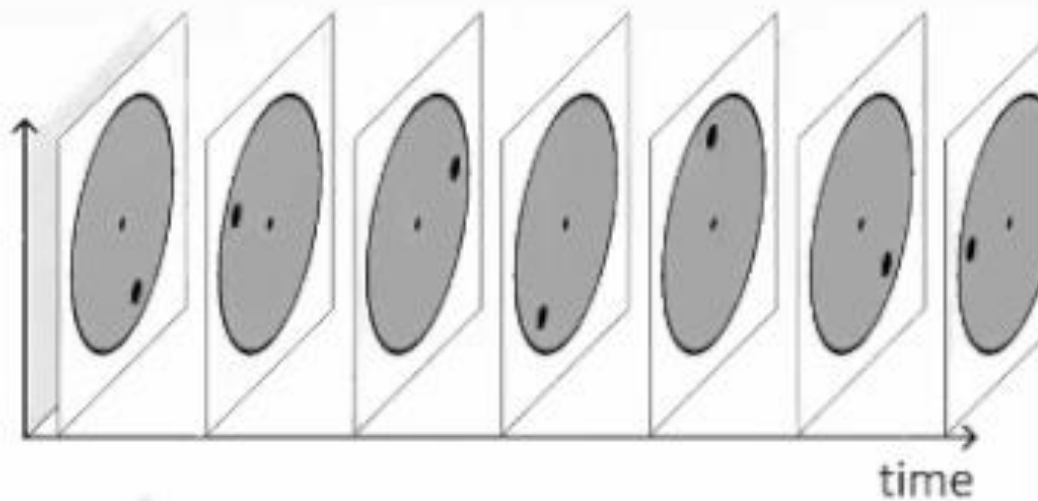RGB        IR (infrared)

Merged result

Jian Wang, Tianfan Xue, Jonathan T. Barron, Jiawen Chen
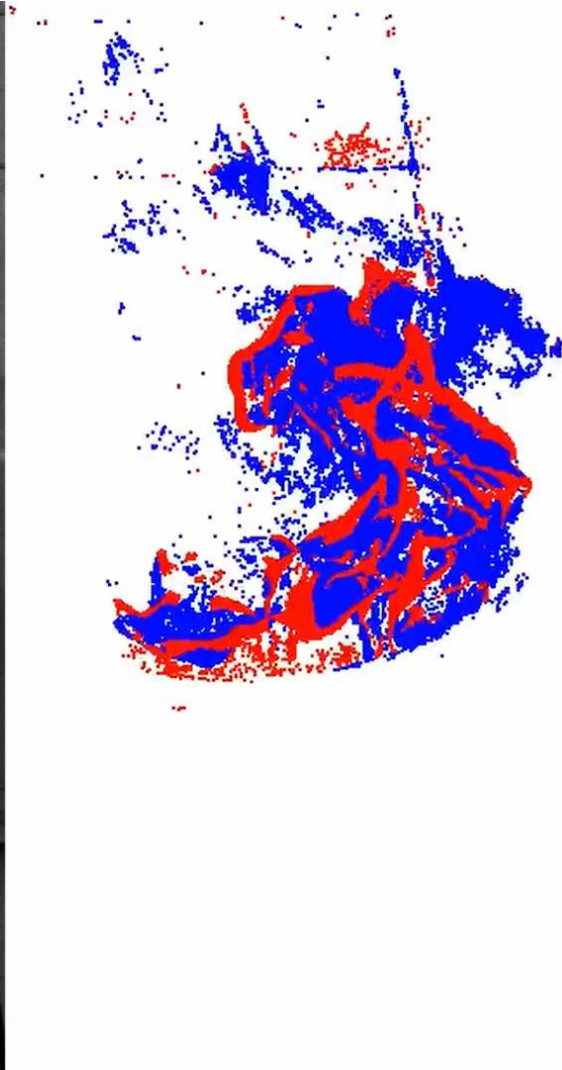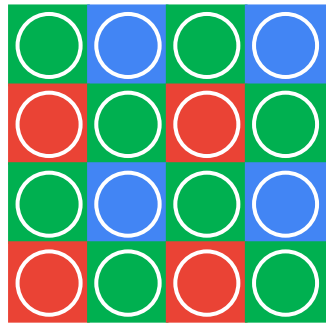
**ICCP 2019**

# Event camera

# Use event camera to recover high-speed motion
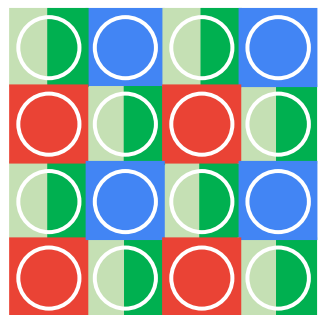


S. Tulyakov et al., "Time Lens: Event-based Video Frame Interpolation", CVPR 2021.
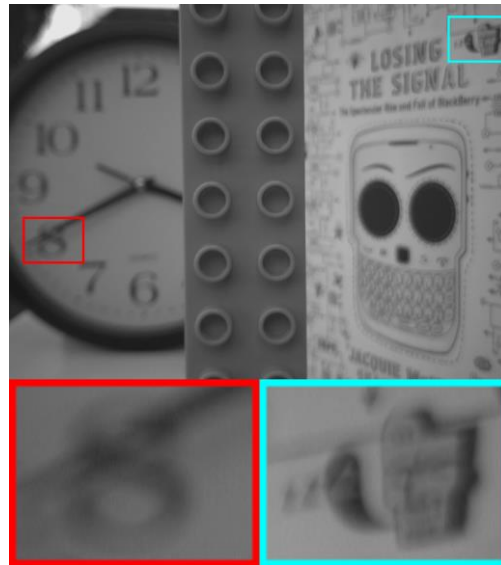
# Depth and debluring from DP images



Regular sensor

DP sensor

Input DP image

Depth map

All-in-focus image

Far

Near

Defocus Map Estimation and Deblurring from a Single Dual-Pixel Image

S. Xin, N. Wadhwa, T. Xue, J. T. Barron, P. P. Srinivasan, J. Chen, I. Gkioulekas, R. Garg

**ICCV 2021**

# Overview

- **Capturing**: multiple source fusion

- **Processing & editing**

  - **Training data:** synthetic data

  - **Network**: combine classic image processing algorithm and machine learning



Camera & Sensor

Processing

Editing

# An input/output pair is needed for ML training



Noisy input

Neural network

Denoised output

# Ground truth output are hard to label



Labeling detection is **easy**.
few seconds / image

Image Credit: https://github.com/tzutalin/labelImg



Label denoising is **hard**:
few hours / image

Image Credit: Nik Collection

# Capturing ground truth requires a lot of manual efforts



Collecting ground truth for denoising (<100/day)

# Devices differences

DSLR  Mobile camera  Webcam
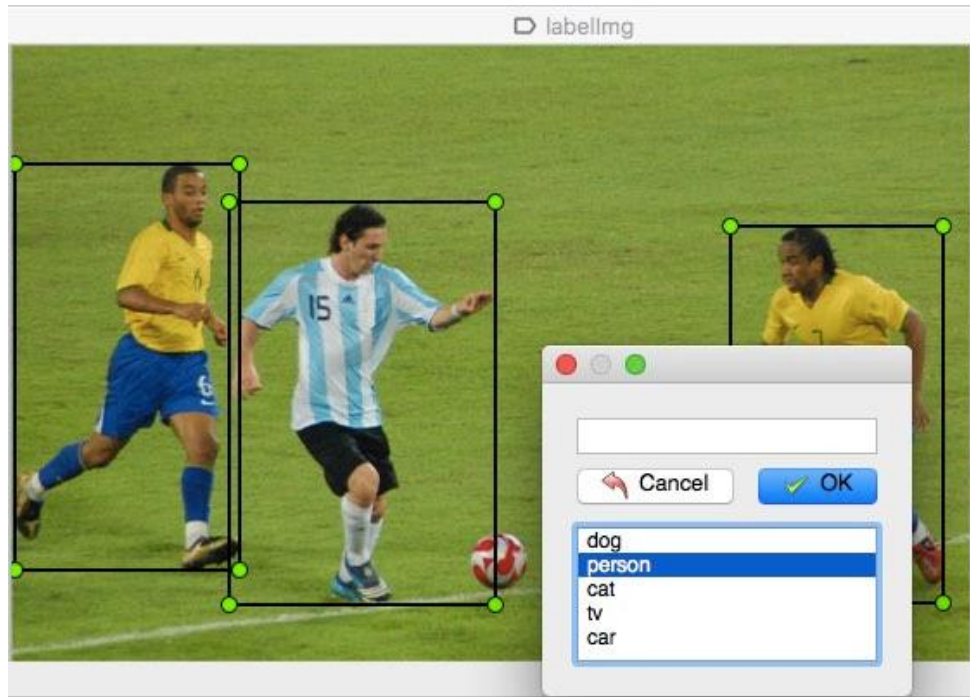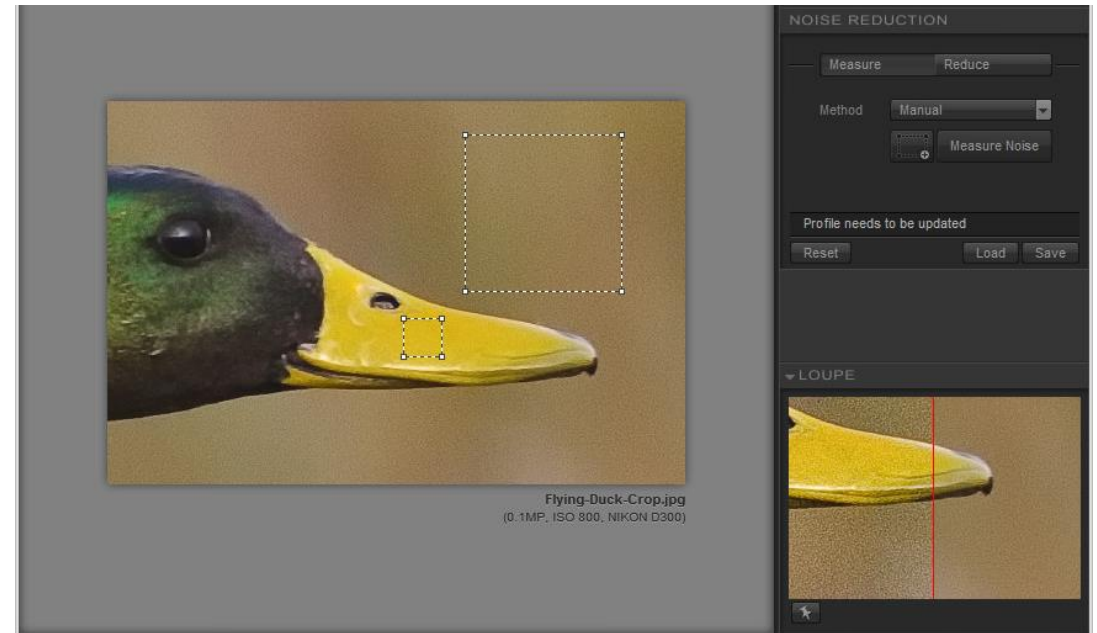
## List of Android smartphones

From Wikipedia, the free encyclopedia

This is a list of devices that run on the Android operating system.

*This is a dynamic list and may never be able to satisfy particular standards for completeness. You can help by adding missing items with reliable sources.*

### References  [ edit ]

1. ^ "Asus PadFone - Full phone specifications". *GSMArena.*
2. ^ "Asus PadFone 2 - Full phone specifications". *GSMArena.*
3. ^ "Asus PadFone Infinity - Full phone specifications". *GSMArena.*
4. ^ "Asus PadFone Infinity 2 - Full phone specifications". *GSMArena.*
5. ^ "Asus PadFone mini - Full phone specifications". *GSMArena.*
6. ^ "Asus PadFone E - Full phone specifications". *GSMArena.*
7. ^ "Asus PadFone Infinity Lite - Full phone specifications". *GSMArena.*
8. ^ "Asus Zenfone 5 A500CG (2014) - Full phone specifications". *GSMArena.*
9. ^ "Asus Zenfone 4 (2014) - Full phone specifications". *GSMArena.*
10. ^ "Asus Zenfone 6 A600CG (2014) - Full phone specifications". *GSMArena.*
...

271. ^ "Honor 8X - Full phone specifications". *GSMArena.*
272. ^ "Honor 8X Max - Full phone specifications". *GSMArena.*
273. ^ "Honor 8C - Full phone specifications". *GSMArena.*
...

553. ^ "Motorola ATRIX HD MB886 - Full phone specifications". *GSMArena.*
...

804. ^ "Oppo A72 - Full phone specifications". *GSMArena.*
...

1058. ^ "Samsung Galaxy A8+ (2018) - Full phone specifications". *GSMArena.*
...

1317. ^ "Xiaomi Mi Max - Full phone specifications". *GSMArena.*
1318. ^ "Xiaomi Mi 5s - Full phone specifications". *GSMArena.* Retrieved 2020-06-06.
...

# Can we use images on the web



Exposure: 30s   Exposure: 5s   Exposure: 0.5s

<100 images / day



No. of images uploaded to internet:
3,000,000,000,000 images / day

*by Leon Seibert, Unsplash*

# Apply degeneration to images on the web

Input/output pairs



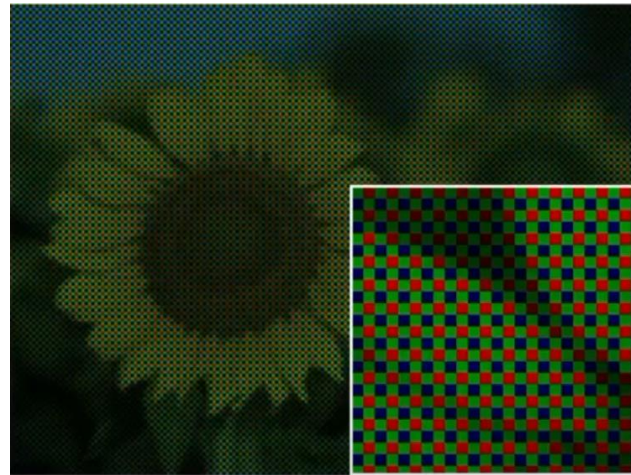Clean images from web

Synthesize degeneration

Degenerated images

Processing network

Recovered clean image

**How to generate realistic degeneration?**
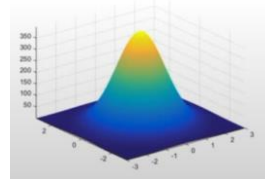
# Real noise does not directly apply to sRGB



Raw image

Camera pipeline

sRGB image

Noise

$$y \sim \mathcal{N}(\mu = x, \sigma^2 = \lambda_{read} + \lambda_{shot}x)$$

Noise model

# Synthesize raw from sRGB



T. Brooks, B. Mildenhall, **T. Xue**, J. Chen, D. Sharlet, J. T. Barron, "Unprocessing images for learned raw denoising", CVPR, 2019

# Unprocess improves the image quality



Noisy input                    N3Net                    Ours

# Simulate realistic rain drops



(a) Heavy rain.          (b) Rain acccumulation.

Yang et al., "Deep Joint Rain Detection and Removal from a Single Image", 2017

# Sometimes, it is important to understand 3D geometry in the simulation



RGB image   Syn. Depth   Syn. Mesh

| Symbol | Definition |
|--------|------------|
| $T$ | Front scene image |
| $R$ | Back scene image |
| $\tilde{R}$ | Back scene image reflected by a glass |
| $X^*$ | Predicted image of $X$ |
| $T/R$ | $T$ or $R$ |

(1) $I$    (2) $T$    (3) $\tilde{R}$    (4) $R$

Kim et al., "Single Image Reflection Removal with Physically-Based Training Images", CVPR 2020

# We can even resort to rendering engine



S. Niklaus et al., "Learned dual-view reflection removal", WACV, 2021.

# We can even resort to rendering engine



S. Niklaus et al., "Learned dual-view reflection removal", WACV, 2021.

# Lens flare

# Flare formation



Wu et al,. "How to Train Neural Networks for Flare Removal", ICCV 2021

# Overview

- **Capturing**: multiple source fusion

- **Processing & editing**

  - **Training data:** synthetic data

  - **Network**: combine classic image processing algorithm and machine learning

# Representing point operation: 3D LUT



Green Axis

Red Axis

Blue Axis

https://www.bromptontech.com/what-is-a-3d-lut/

# Learning to enhance -> Learning 3D LUT



H. Zeng et al., "Learning Image-adaptive 3D Lookup Tables for High Performance Photo Enhancement in Real-time", T-PAMI, 2020

# Learning 3D LUT significantly reduces the time cost

| Resolution | 1920×1080 | 3840×2160 | 6000×4000 |
|---|---|---|---|
| Pix2Pix [49] | 1.2e2 | N.A. | N.A. |
| CycGAN [50] | 5.6e2 | N.A. | N.A. |
| DPE [7] | 8.6e1 | N.A. | N.A. |
| White-Box [9] | 5.0e3 | 9.1e3 | 2.0e4 |
| Dis-Rec [8] | 2.5e4 | 1.1e5 | 3.3e5 |
| UIE [11] | 1.0e4 | 2.0e4 | 3.3e4 |
| HDRNet [2] | 4.5e1 | 2.1e2 | 5.9e2 |
| UPE [10] | 4.5e1 | 2.1e2 | 5.9e2 |
| **Ours** | **0.64** | **1.66** | **3.76** |

# Image Enhancement: Using color tran. and global curve



Song et al, "StarEnhancer: Learning Real-Time and Style-Aware Image Enhancement", ICCV 2021

# Direct prediction is expensive

[PhotoWCT: Li et al., ECCV 2018]
[WCT2: Yoo et al., ICCV 2019]
[LST: Li et al., CVPR 2019]

All of them are OOM when applied to 4MP image



Style image



Input

Neural network

Stylized output

# Can we approximate it using tone curves?



Input

Stylized output

# Model a set of tone curves as bilateral grid



Bilateral grid

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \end{bmatrix} \begin{bmatrix} r \\ g \\ b \\ 1 \end{bmatrix}$$

J. Chen, A. Adams, N. Wadhwa, S. Hasinoff, "Bilateral guided upsampling", 2017
M. Gharbi, J. Chen, J. Barron, S. Hasinoff, F. Durand, " Deep Bilateral Learning for Real-Time Image Enhancement", SIGGRAPH 2017

# Style transfer using a set of tone curves



Tone curves baked in bilateral grid

**Low resolution: 256x256**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Full resolution: 4Kx3K**

Apply

Output

Input

X. Xia, M. Zhang, **T. Xue**, Z. Sun, H. Fang, B. Kulis, J. Chen, "Joint bilateral learning for real-time universal photorealistic style transfer", ECCV 2020

# Performance

| Image Size | PhotoWCT | LST | WCT$^2$ | Ours |
|---|---|---|---|---|
| 512 × 512 | 0.68s | 0.25s | 3.85s | < 5 ms |
| 1024 × 1024 | 1.51s | 0.84s | 6.13s | < 5 ms |
| 1000 × 2000 | 2.75s | OOM | 10.94s | < 5 ms |
| 2000 × 2000 | OOM | OOM | OOM | < 5 ms |
| 3000 × 4000 | OOM | OOM | OOM | < 5 ms |

Latency

| Mean Score | PhotoWCT | LST | WCT$^2$ | Ours |
|---|---|---|---|---|
| Photorealism | 2.02 | 2.89 | **4.21** | 4.14 |
| Stylization | 3.10 | 3.19 | 3.24 | **3.49** |
| Overall quality | 2.23 | 2.84 | 3.60 | **3.79** |

User study of visual quality

Results on 12MP image

# HDRnet tonemapping



12 megapixel 16-bit linear input
(tone-mapped for visualization)

tone-mapped with HDR+
**400 – 600 ms**

processed with our algorithm
**61 ms**, PSNR = **28.4 dB**

M. Gharbi, J. Chen, J. Barron, S. Hasinoff, F. Durand, " Deep Bilateral Learning for Real-Time Image Enhancement", SIGGRAPH 2017

# Used by Google Tensor Chip

# Denoising using spatially varying kernels



Input       MalleConv       Selected kernel 1       Selected kernel 2

Jiang et al, "Fast and High-quality Image Denoising via Malleable Convolutions", ECCV 2022

# We can even learn to reorder different modules



Basic camera modules

Learn a task-specific pipeline

K. Yu et al, "ReconfigISP: Reconfigurable Camera Image Processing Pipeline", ICCV 2021

# We can even learn to reorder different modules



Basic camera modules

Different task may need different pipelines

K. Yu et al, "ReconfigISP: Reconfigurable Camera Image Processing Pipeline", ICCV 2021

Future smart cameras research

# Simulation is important to collect training data



Image credit: Tesla AI Day

# Is there Isaac Gym for computational photography?



Issac Gym by NVIDIA, for robotic algorithm design

# Computational Photography and Hardware

# Joint lens and algorithm design



Tseng et al., "Neural nano-optics for high-quality thin lens imaging", 2021

# Camera is not only for better selfies



First black hole image
*[image credit: NASA]*



VLBI image formation

K. Bouman, "Computational Imaging for VLBI Image Reconstruction", CVPR 2016

# How VR may impact us

# Can we capture more 3D content using our cameras?



B. Mildenhall et al., "NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis", ECCV, 2020

# Can cameras be as good as our eyes?



Optical see through displays

Video see through displays

camera

# Can our algorithms be as fast as to deal with 8K60fps?

# AIGC and Computational Photography

# Where is the boundary between **editing** and **synthesis**



Which one is real, which one is fake?

# Where is the boundary between **editing** and **synthesis**



Robin et al., "High-Resolution Image Synthesis with Latent Diffusion Models", CVPR 22.

# Where is the boundary between **editing** and **synthesis**

# Where is the boundary between **editing** and **synthesis**



| input texture | Portilla & Simoncelli [17] | Xu et.al. [21] | Wei & Levoy [20] | Image Quilting |

A. Efros and W. T. Freeman, "Image quilting for texture synthesis and transfer", CG 2001

# Do we even need a powerful lens?

Only less powerful camera hardware is needed in the future?
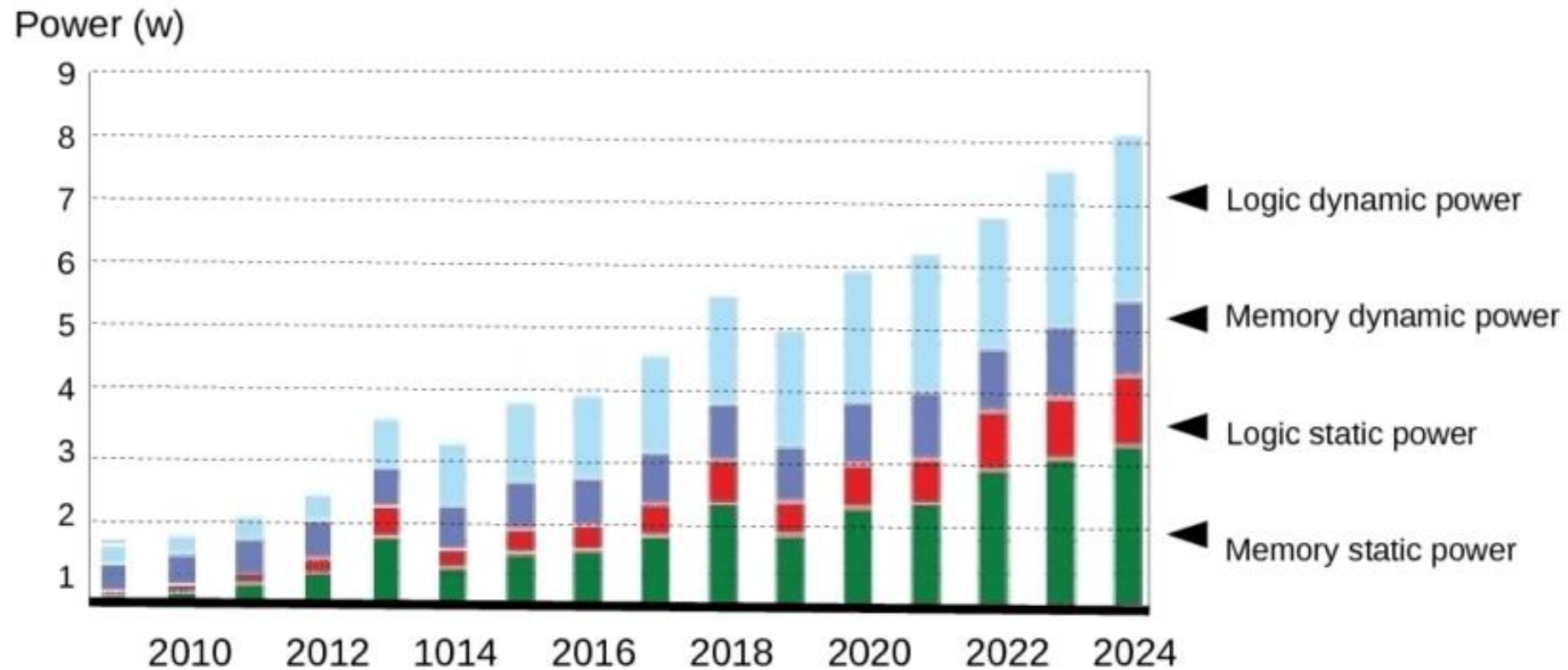


Cameras become **more powerful** in past 10y

# Gap between academic research and industrial design

# Industrial not only care quality, but also speed & power



SoC consumer portable power consumption

Yahia Benmoussa, "Performance and Energy Consumption Characterization and Modeling of Video Decoding on Multi-core Heterogenous Mobile SoC and their Applications"

# High PSNR != Better Image
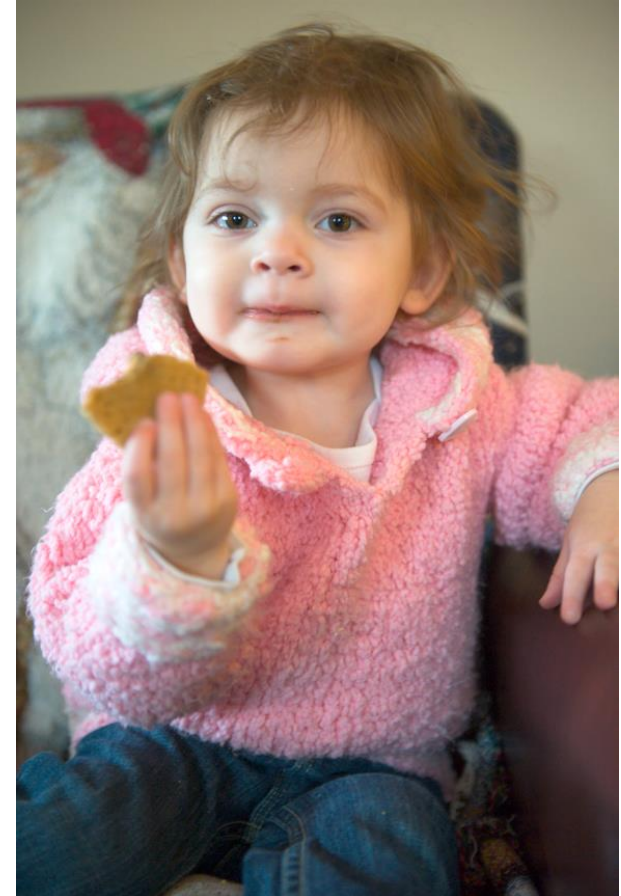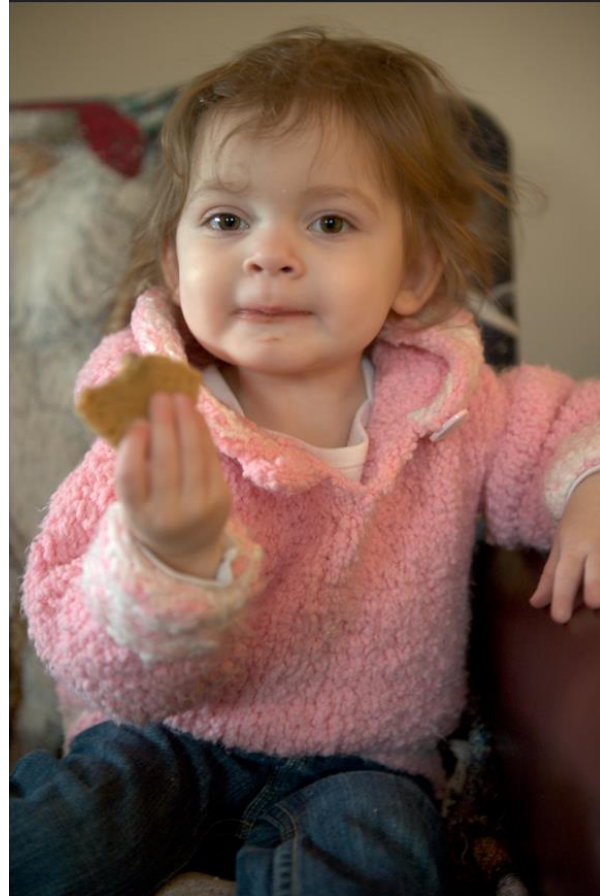


bicubic (21.59dB/0.6423)     SRResNet (23.53dB/0.7832)     SRGAN (21.15dB/0.6868)

L., Christian, et al. "Photo-realistic single image super-resolution using a generative adversarial network." *CVPR*. 2017.

# Diff users may even have diff preferences



3 experts gave very different tunings

# It is even hard to describe what is best



Some photographers don't like over-smoothed image, and call it "**like oil painting**"



Some photographers don't like HDR image, and call it "**Cartoon-look**"

http://barney-streit.squarespace.com/blog/2013/6/5/good-hdr-bad-hdr

# Collaborators and acknowledgement



Bill Freeman · Fredo Durand · Antonio Torralba · Josh Tenenbaum · Michael Geosele · Rick Szeliski · Daniel · Ce Liu · Orly Liba

Miki Rubinstein · Joseph Lim · Yuandong Tian · Hossein Mobahi · Jiajun Wu · Katie Bouman · Neal Wadhwa · Chengkai Zhang · Yun-Ta Tsai
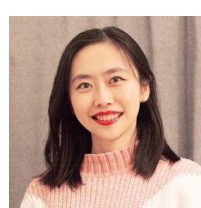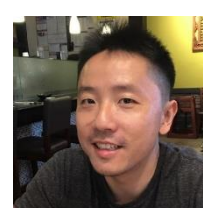
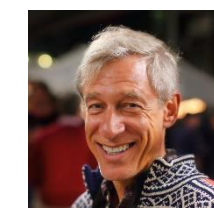Jon Barron · Sam Hasinoff · Dillon Sharlet · Jiawen Chen · Xide Xia · Zheng Sun · Kiran Murthy · Marc Levoy · Shumian Xin
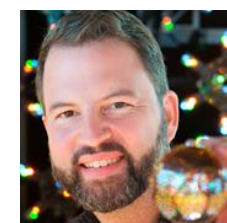
Pratul Srinivasan · Ioannis Gkioulekas · Rahul Garg · Jian Wang · Qiurui He · Brian Kulis · Simon Niklaus · Paul Debevec · Xiuming Zhang

Q&A