



Provided by the author(s) and University of Galway in accordance with publisher policies. Please cite the published version when available.

Title	Edge-AI: Rise of the neural accelerators
Author(s)	Corcoran, Peter
Publication Date	2019-08-29
Publication Information	Corcoran, Peter. (2019). Edge-AI: Rise of the neural accelerators, Paper presented at the EMAI 2019 Emerging Memory and Artificial Intelligence Workshop (Plenary presentation), Stanford Center for Magnetic Nanotechnology, Stanford, 29 August
Publisher	EMAI 2019 Emerging Memory and Artificial Intelligence Workshop
Link to publisher's version	https://emai19.sites.stanford.edu/workshop-presentations
Item record	http://hdl.handle.net/10379/15489

Downloaded 2024-05-19T09:42:13Z

Some rights reserved. For more information, please see the item record link above.



<http://bit.ly/neural2019>

EMAI 2019

EMERGING MEMORY AND ARTIFICIAL INTELLIGENCE WORKSHOP

Prof. Peter Corcoran, National University of Ireland Galway



Center for Cognitive, Connected and Computational Imaging
College of Engineering, Science & Informatics, NUI Galway

WHO AM I?

- IEEE Volunteer (Electronic & ICT Engineer)
 - Board Member of IEEE Consumer Electronics Society (6 years)
 - Editor-in-Chief of IEEE Consumer Electronics Magazine (2010-2016)
 - IEEE Fellow in 2010 (Contributions to Digital Camera Technology)
 - IEEE Distinguished Lecturer, Conference Chair, Editor & Reviewer
- Day Job(s):
 - University Professor & Former Vice-Dean (H-Index 85; 20k citations)
 - Active Researcher (currently 8 PhD & 3 PostDoctoral researchers)
 - Entrepreneur, Inventor & Technologist; (300+ patents)
 - Industry Consultant
- Contact Information
 - E-Mail: dr.peter.corcoran@ieee.org
 - peter.corcoran@nuigalway.ie
 - Twitter: @pcor LinkedIn:
 - Google Scholar:
<https://scholar.google.com/citations?hl=en&user=J6YWBB4AAAAJ>



WHAT IS IN THIS TALK?

- 1) How Big Data became Fool's Gold ...
- 2) ... and Artificial Intelligence is moving to the Edge
- 3) What's inside a Camera? (<\$1)
Today's camera tech provides good examples of where Edge-AI is headed ...
- 4) AI + Camera < \$5 → Disruptive Edge Tech
- 5) But new, disruptive applications will need Training Data!
Data Acquisition is complex, difficult to get right and expensive to collect
Driver Monitoring System as an Example (Face Pose, Eye-Gaze, Gestures)
- 6) Solutions to the Data Problem
- 7) Thoughts & Take-Aways on Storage & Data Bandwidth

<http://bit.ly/neural2019>

RISE OF THE NEURAL ACCELERATORS

Why AI at the Edge will drive the need for Storage & Bandwidth



NUI Galway
OÉ Gaillimh



C3I

Cognitive, Connected & Computational Imaging

Center for Cognitive, Connected and Computational Imaging
College of Engineering, Science & Informatics, NUI Galway

THERE WAS A TIME WHEN COMPUTERS WERE SIMPLE ...

1960's



1980's



THEN EVERYTHING CHANGED ...

2000's – Virtual Machines



A computer was no longer a computer ...

By mid-2000's – The Cloud



And data started to disappear into the network ...

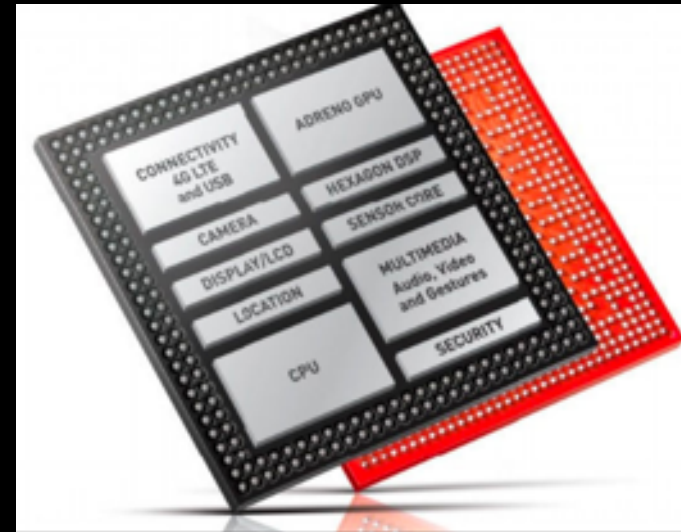
AND CHANGED SOME MORE ...

2010's – Phones got Smart ...



And started to take over our daily lives ...

& became very, Sophisticated ...



Driving the cutting edge of real-time sensing & data analytics ...

EVERYTHING MOVED INTO THE CLOUD ...

2015+ – Data got BIG ...

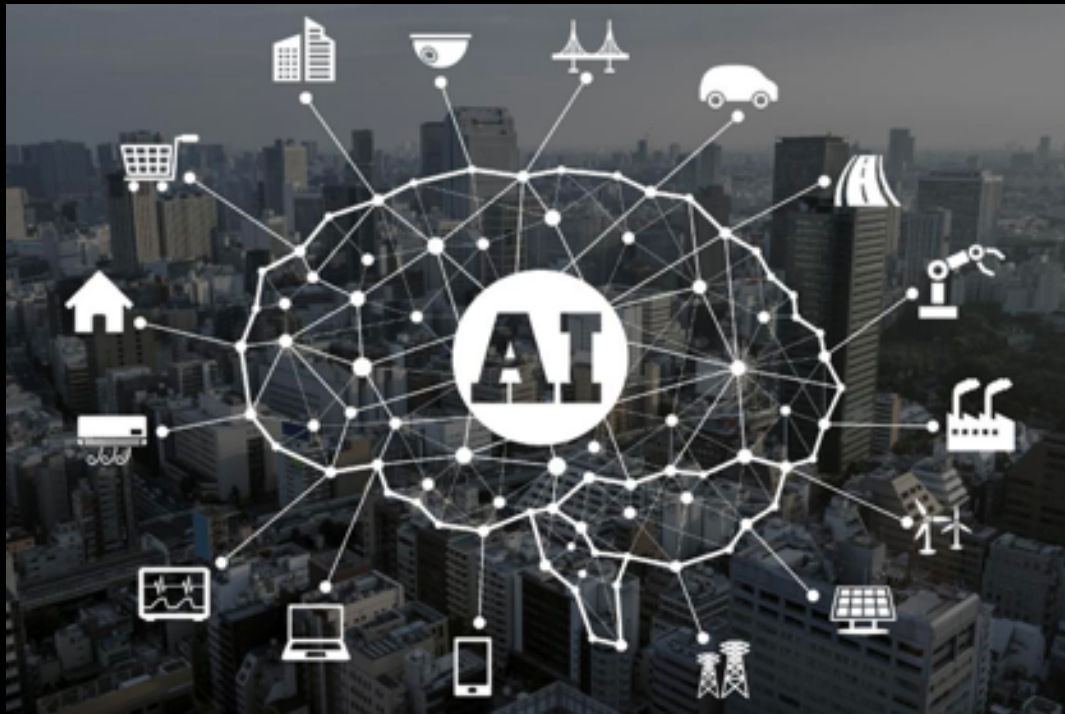


& Speakers got Smart ...



Driving a new wave of Artificial Intelligence ...

AND ARTIFICIAL INTELLIGENCE TURNED DATA INTO GOLD



COMPANIES GOT FAT ON DATA....



The Observer

www.observer.co.uk | Monday 18 March 2013 | £1.00

Revealed: 50m Facebook files taken in record data breach

The Cambridge Analytica Files

Full interview
Whistleblower Christopher Wylie lifts the lid

Cover story
New Review

Like or dislike
The algorithm that reveals all about you

Report, page 9

Facebook
How its destructive ethos imperils democracy

Observer
Comment, 61



Exclusive

- Whistleblower tells of bid to influence votes
- Tech giant suspends controversial data firm

Candice Cadwalladr
@ Emma Graham-Harrison

The data analysis firm that worked with Donald Trump's election team and the winning Brexit campaign harvested millions of Facebook profiles of US users, in one of the web giant's biggest ever data breaches, and used them to build a powerful software program to predict and influence choices at the ballot box.

A whistleblower has revealed to the Observer how Cambridge Analytica

Christopher Wylie, who worked with a Cambridge University academic to obtain the data, told the Observer: "We exploited Facebook to harvest millions of people's profiles. And built models to exploit what we knew about them and target their inner demons. That was the basis the entire company was built on."

Documents seen by the Observer, and confirmed by a Facebook statement, show that by late 2015 the company had found out that information had been harvested on an unprecedented scale. However, at the time it failed to alert users and took only limited steps to recover and remove the private information of more than 50 million individuals.

The New York Times is reporting that copies of the data harvested for Cambridge Analytica could still be found online; its reporting team had viewed some of the raw data.

The data was collected through an



NUI Galway
OÉ Gaillimh



C3I

Cognitive, Connected & Computational Imaging

THEN "PRIVACY"
HIT THE FAN ...

Center for Cognitive, Connected and Computational Imaging
College of Engineering, Science & Informatics, NUI Galway

EUROPEAN SOCIAL RESPONSIBILITY (AKA. REGULATORS) FOUGHT BACK ...

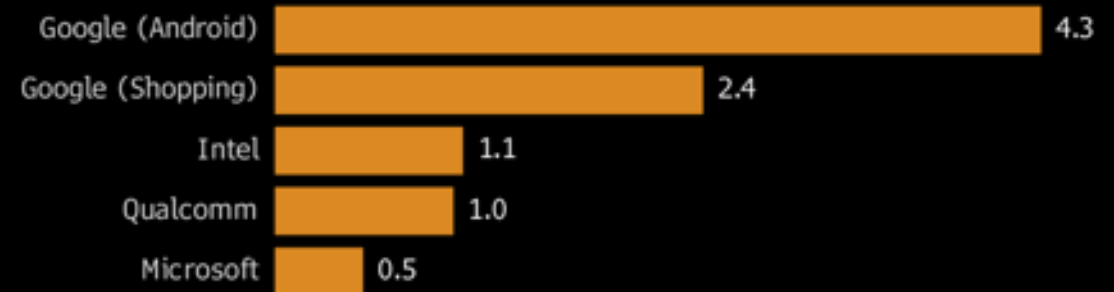
€4.3B fine



Another Record Fine for Google

The EU's Five Biggest Antitrust Penalties to Date

■ Fines in billions of euros



Source: European Commission

Bloomberg

POINT #1 FROM TODAY'S TALK:
THE NEW "AGE OF PRIVACY" HAS ARRIVED
– GATHERING CENTRALIZED DATA JUST
BECAME A FOOL'S ERAND



BUT TECHNOLOGY MARCHES ON ...

- Neural Accelerators are here ...
- Now you can analyze data where & when it is created ...



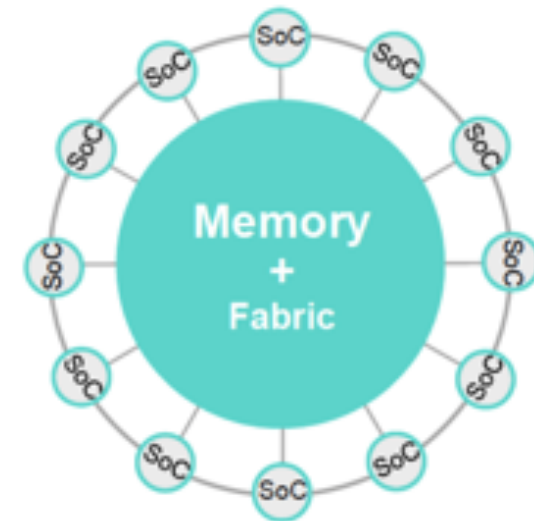
- At the Edge

THIS ISN'T THE END OF THE CLOUD ...

- But it is an important 'saddle point' for ICT Technology
- This is Important because it enables "on-chip" Memory-Driven computing performed IN the Memory Fabric!

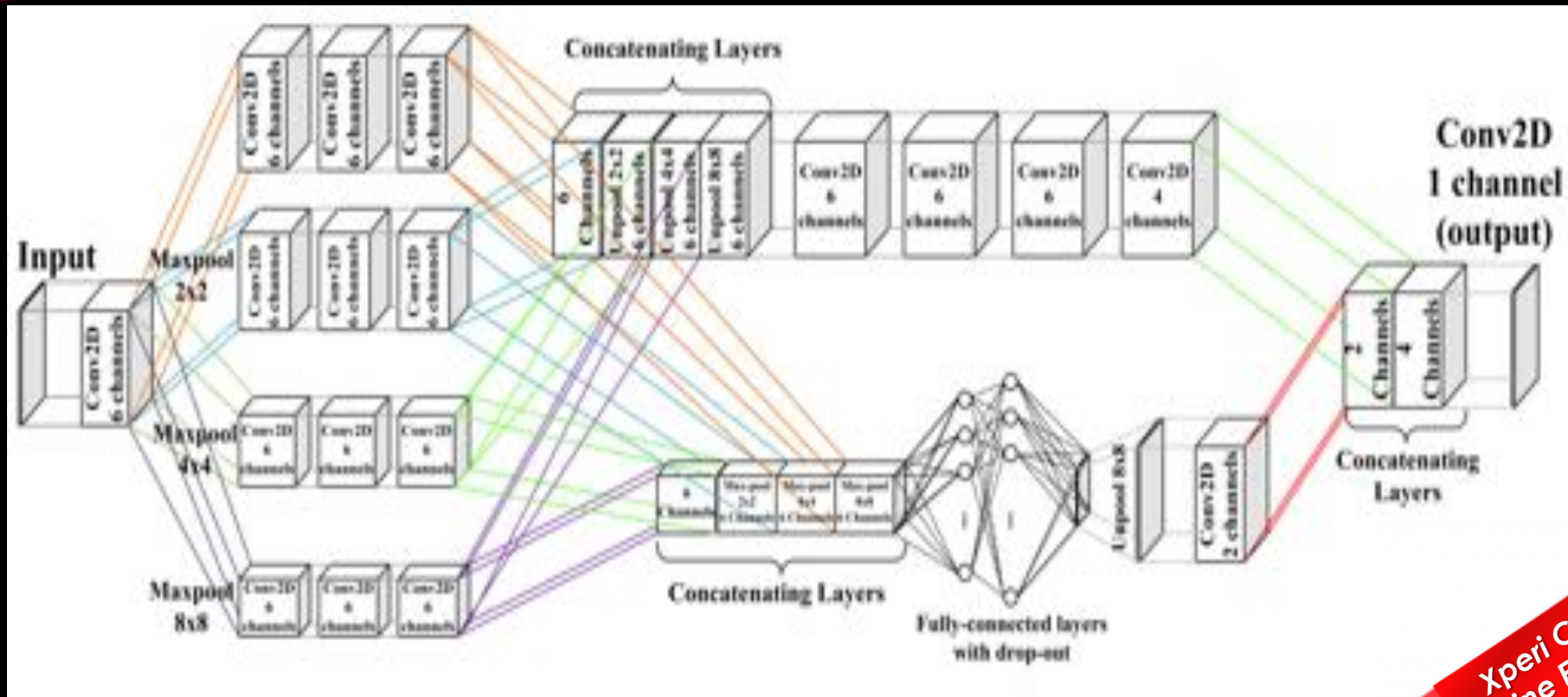


From Processor-Centric Computing...



...to Memory-Driven Computing

NOW POSSIBLE TO PROGRAM AI NETWORKS INTO SD CARD FOOTPRINT WITH 2 ORDERS OF MAGNITUDE LOWER ENERGY USE THAN GPUS ...



POINT #2 FROM TODAY'S TALK:
AI HAS STARTED TO MOVE TOWARDS THE
EDGE OF THE NETWORK AND ABANDON
THE CENTRALIZED CLOUD PARADIGM

WHAT IS INSIDE A CAMERA?

- **& Why are they everywhere in new Consumer Technologies & Use Cases?**
 - Smartphones
 - AR Headsets (user-facing cameras)
 - Driver Monitoring Systems
 - Smart-City Applications
 - IoT Devices (Security, Elderly Monitoring, etc)



NUI Galway
OÉ Gaillimh



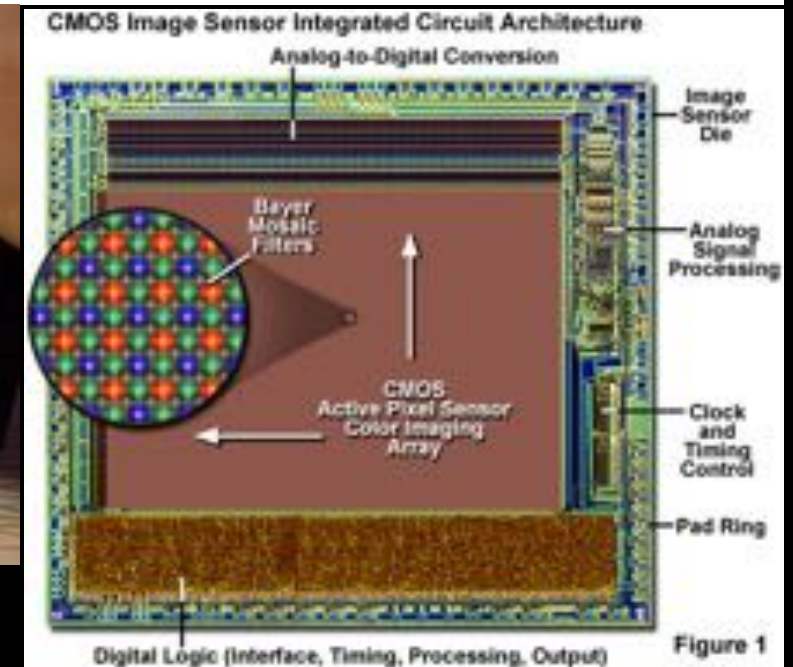
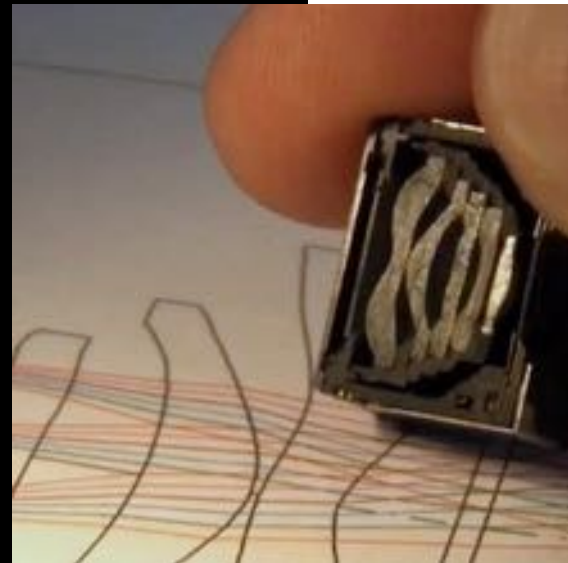
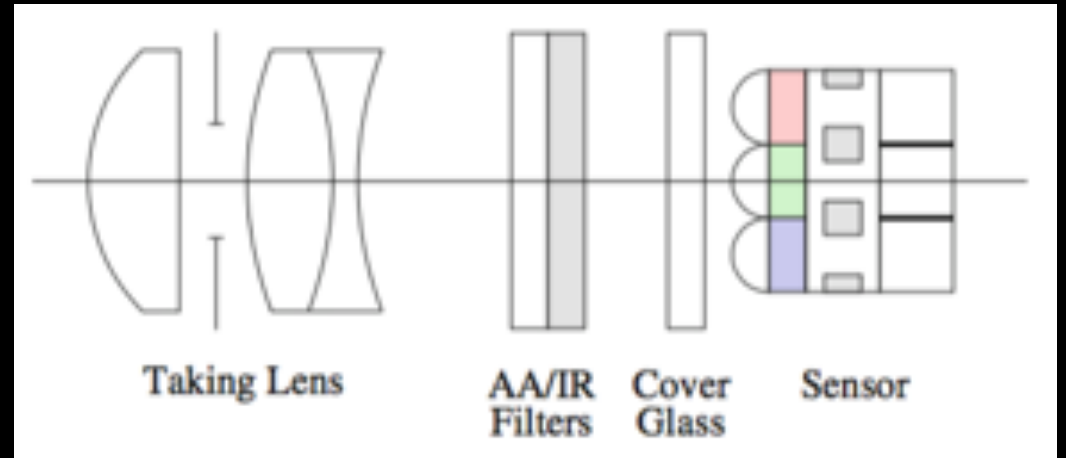
C3I

Cognitive, Connected & Computational Imaging

DIGITAL CAMERA TECHNOLOGIES #1

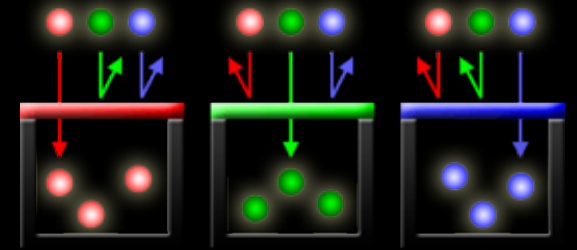
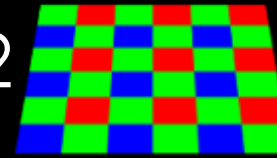
BASICS – THE OPTICAL IMAGE PATH #1

- Multi-Element Lens
 - Typically at least 5-element
 - Telecentric (see reading #1)
 - Small Point-Spread Function (PSF)
- Anti-Aliasing Filter
 - Removes High-Frequency (Spatial) Artifacts
- Infrared Cutoff Filter
 - Silicon is sensitive to NIR
 - NIR focus is different to Visible
- Sensor
 - Bayer Color Filter Array (CFA)
 - Back-Illuminated

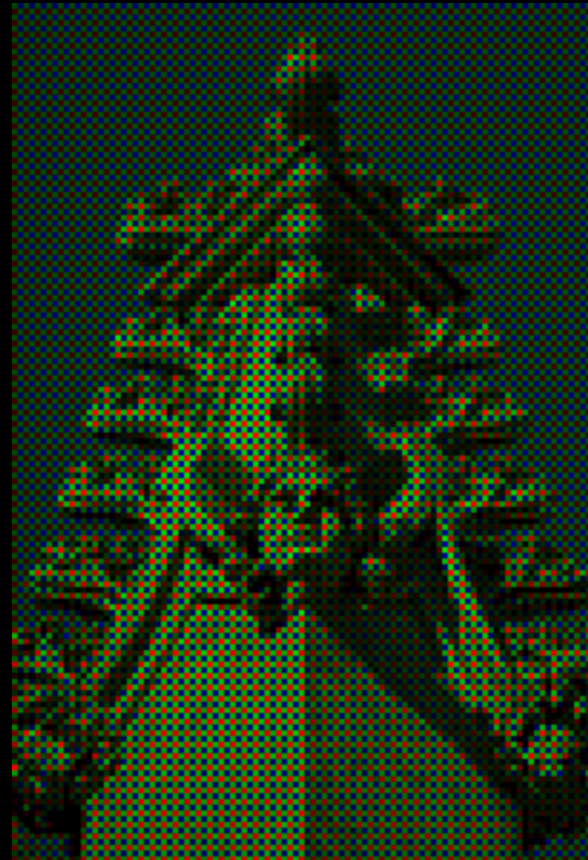


DIGITAL CAMERA SENSORS #2

BAYER IMAGE



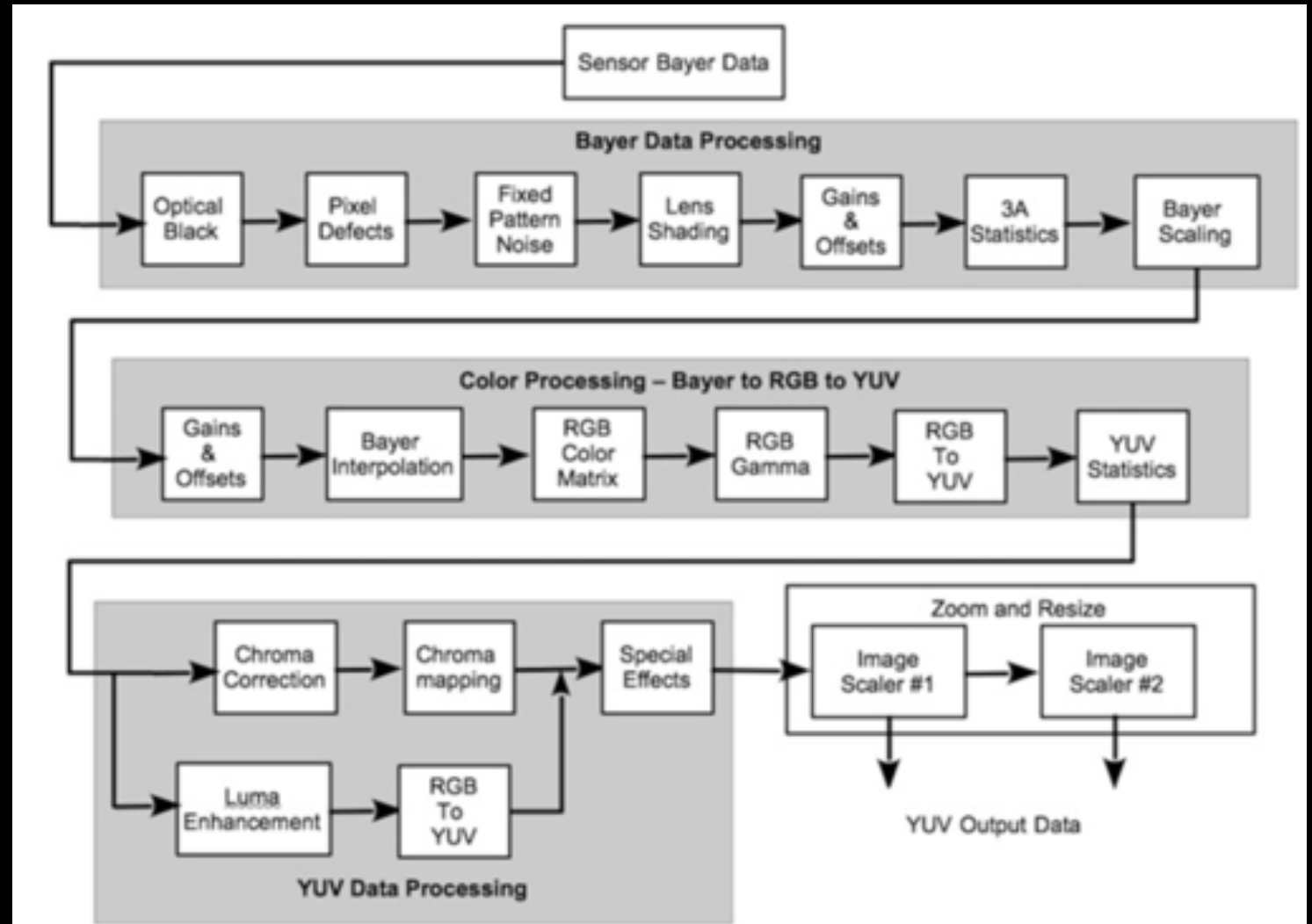
- A Bayer array consists of alternating rows of red-green and green-blue filters.
- Notice how the Bayer array contains twice as many green as red or blue sensors.
- Each primary color does not receive an equal fraction of the total area because the human eye is more sensitive to green light than both red and blue light.
- Redundancy with green pixels produces an image which appears less noisy and has finer detail than could be accomplished if each color were treated equally.
 - Noise in the green channel is less than for the other two primary colors simply because there are twice as many pixels.
- Bayer's technique is > 30 years old – clearly a robust engineering approximation!



DIGITAL CAMERA TECHNOLOGIES #3

BASICS – THE IMAGE PROCESSING PIPELINE #1

- To fully understand the complexity of what happens in a modern digital camera, we need to illustrate the concept of the **image processing pipeline** (IPP) – the sequence of digital manipulations of the original image data to get to the image that you see on the main camera screen.



DIGITAL IMAGE & COMPRESSION BASICS #6

COMPRESSION - JPEG #1

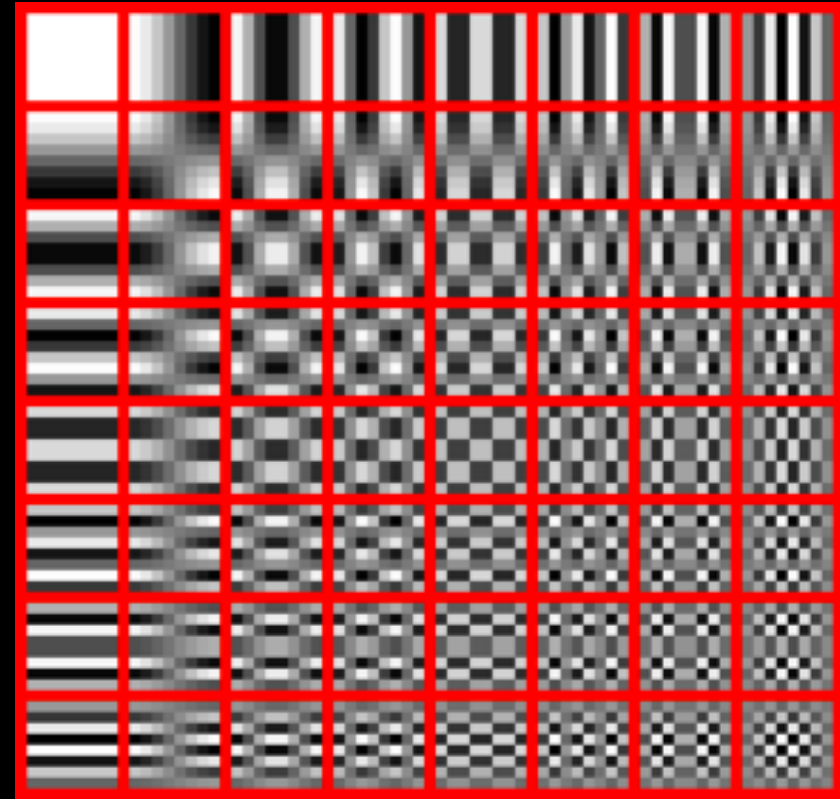
- **JPEG** is a commonly used method of lossy compression for digital images, particularly for those images produced by digital photography.
- The degree of compression can be adjusted, allowing a selectable tradeoff between storage size and image quality.
 - JPEG typically achieves 10:1 compression with little perceptible loss in image quality.
- **JPEG/Exif** is the most common image format used by digital cameras and other photographic image capture devices; along with **JPEG/JFIF**, it is the most common format for storing and transmitting photographic images on the World Wide Web.
- The term "JPEG" is an acronym for the **Joint Photographic Experts Group**, which created the standard.
- As the typical use of JPEG is a lossy compression method, which somewhat reduces the image fidelity, it should not be used in scenarios where the exact reproduction of the data is required (such as some scientific and medical imaging applications and certain technical image processing work).



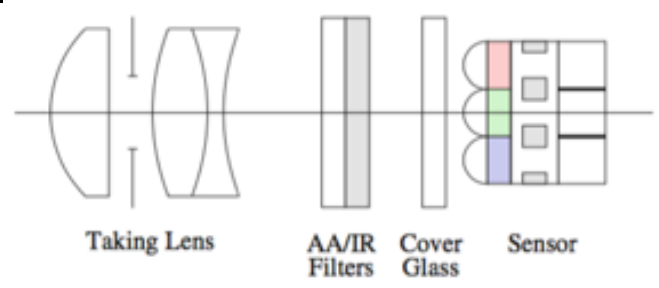
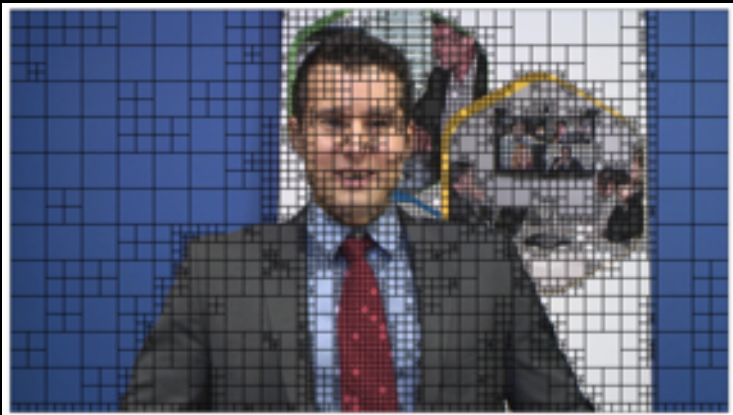
DIGITAL IMAGE & COMPRESSION BASICS

COMPRESSION - JPEG #3

- **JPEG encoding example:** Although a JPEG file can be encoded in various ways, most commonly it is done with JFIF encoding. The encoding process consists of several steps:
 - **Color Space Transformation:** the representation of the colors in the image is converted from RGB to Y'CBCR. (This step is sometimes skipped.)
 - **Chroma Downsampling:** the resolution of the chroma data is reduced, usually by a factor of 2 or 3. This reflects the fact that the eye is less sensitive to fine color details than to fine brightness details.
 - **Block Splitting & DCT:** The image is split into blocks of 8×8 pixels, and on each block, each of the Y, CB, and CR data undergo a discrete cosine transform (**DCT**). A DCT is similar to a Fourier transform in the sense that it produces a form of **spatial frequency spectrum**.
 - **Quantization:** the amplitudes of frequency components are quantized - human vision system (HVS) is more sensitive to small variations in color or brightness over large areas than to high-frequency (edge) variations. Thus, the magnitudes of high-frequency components are stored with lower accuracy than low-frequency components.
 - The quality setting of the encoder affects to what extent the resolution of each frequency component is reduced. If a very low quality setting is used, the high-frequency components may be discarded altogether.
 - **Entropy Encoding:** The resulting data for all 8×8 blocks is further compressed with a lossless algorithm, a form of **Huffman encoding**.
- The decoding process reverses these steps, except the quantization because it is irreversible. Also, modern devices with larger image sensors may use 16×16 or larger DCT blocks.
 - A detailed example is given at: <https://en.wikipedia.org/wiki/JPEG>

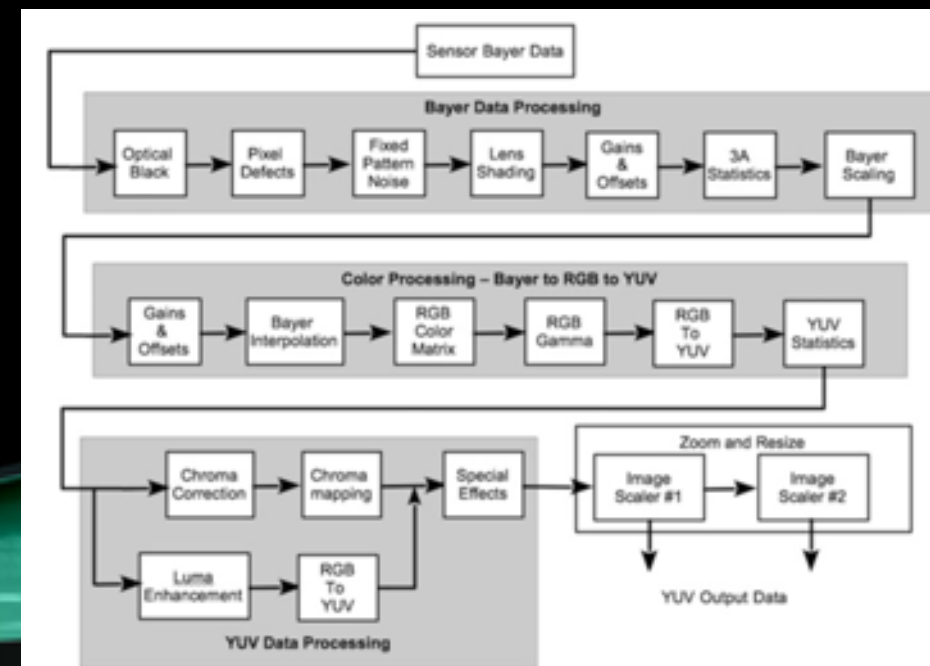
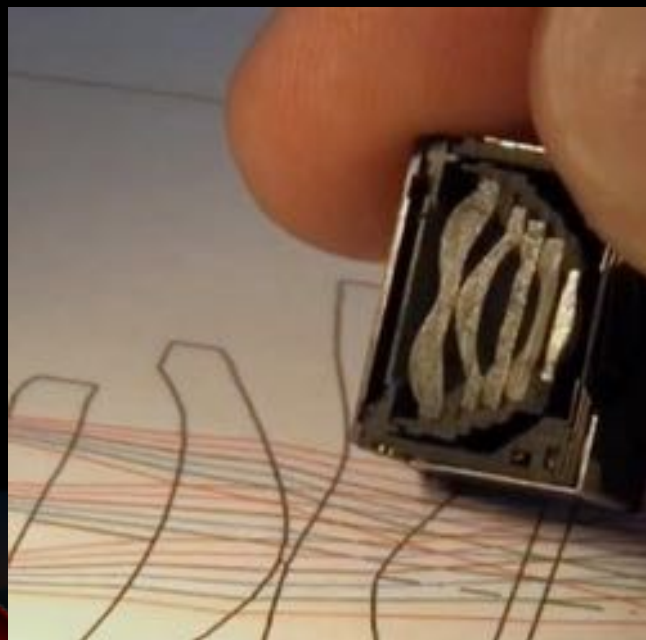
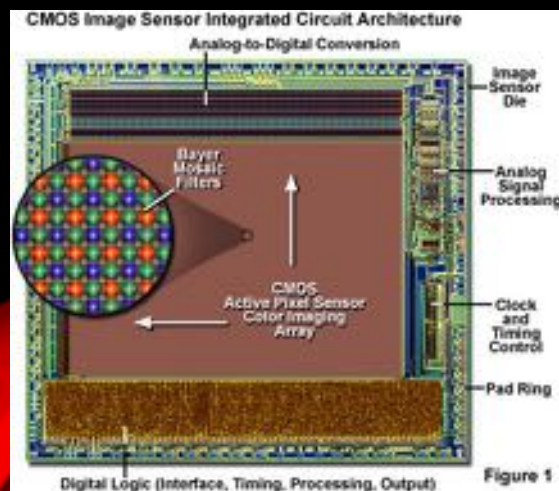


The DCT transforms an 8×8 block of input values to a linear combination of these 64 patterns. The patterns are referred to as the 2D **DCT basis functions**, and the output values are **transform coefficients**. The horizontal index is u and the vertical index is v .

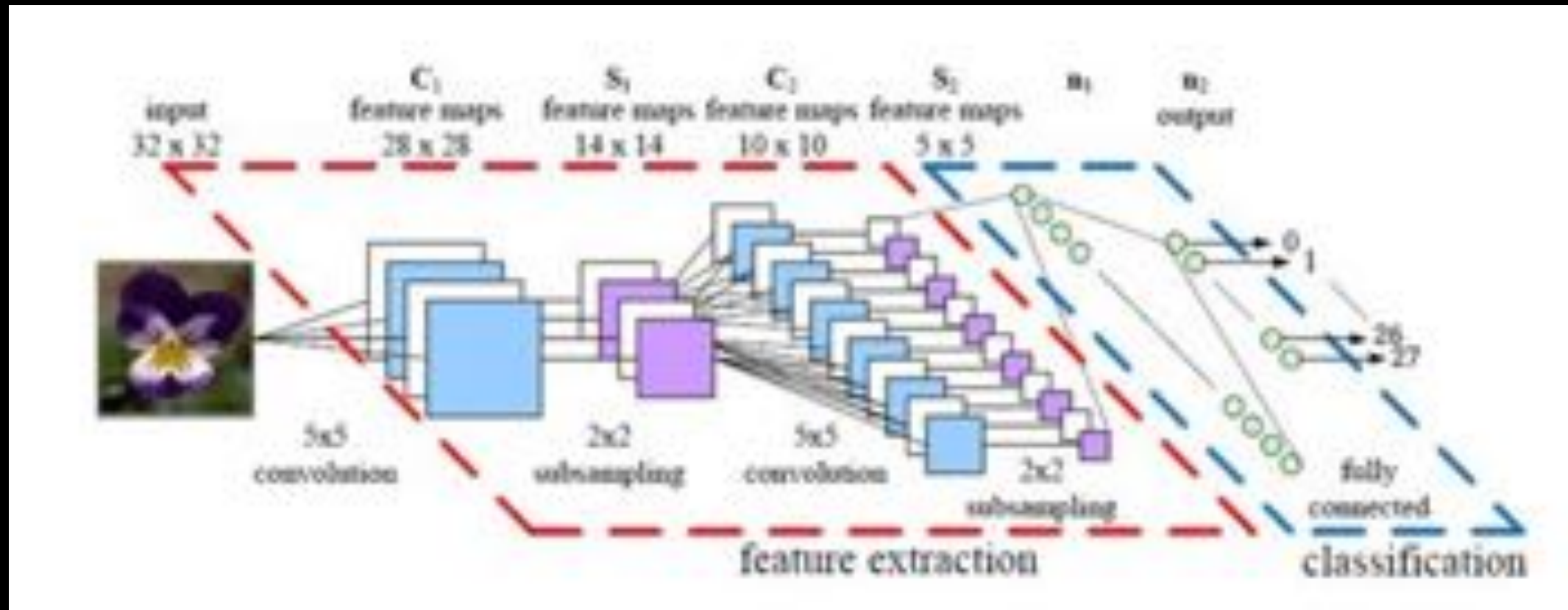


< \$1

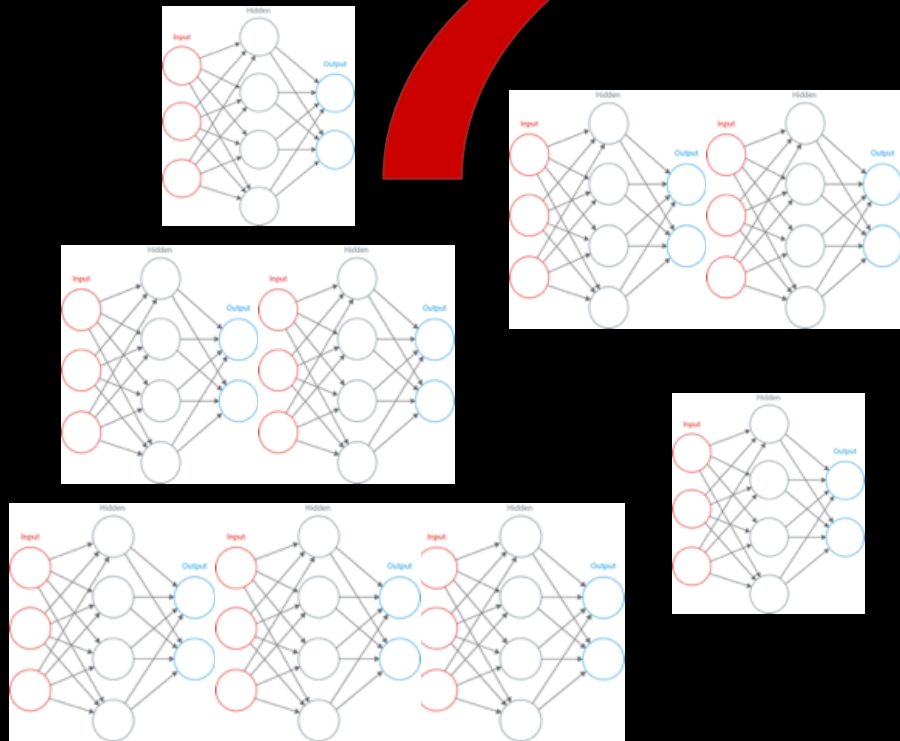
POINT #3 FROM TODAY'S TALK:
CAMERA TECH IS INCREDIBLY SOPHISTICATED, BUT
TODAY IT IS A VERY **LOW-COST SENSING COMMODITY!**



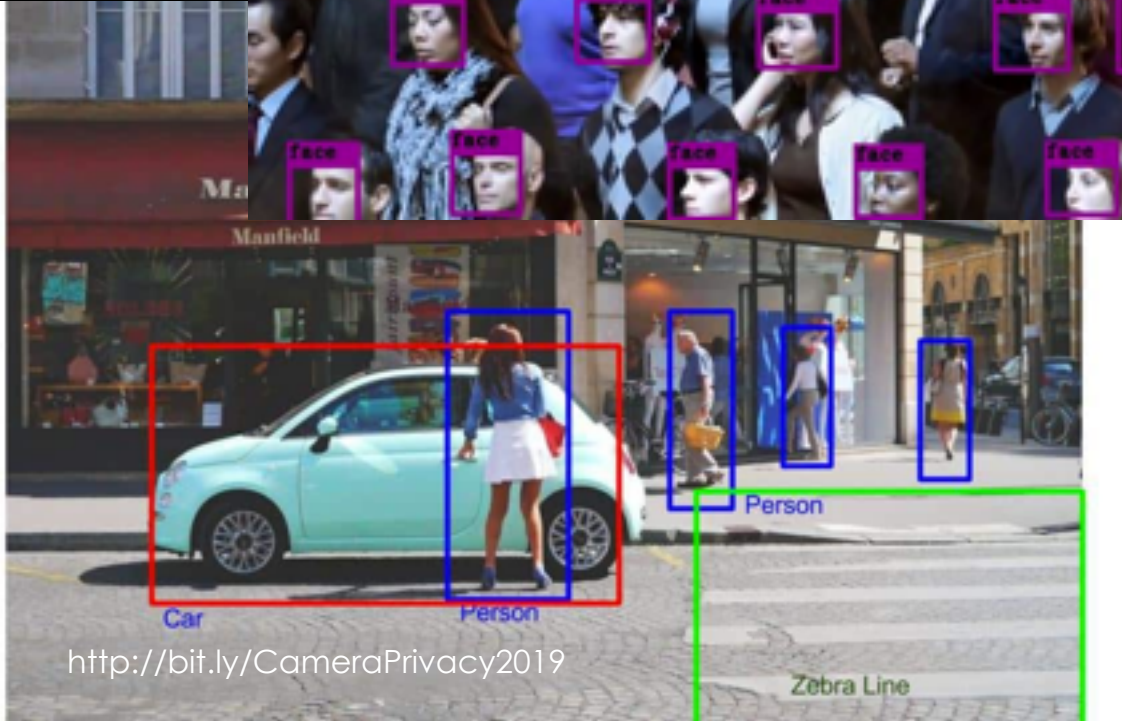
THE NEXT GENERATION OF CAMERAS? CONVOLUTIONAL NEURAL NETWORKS (CNNs) ...



... EMBEDDED INSIDE THE CAMERA!



TO MAKE A REALLY "SMART" CAMERA ...

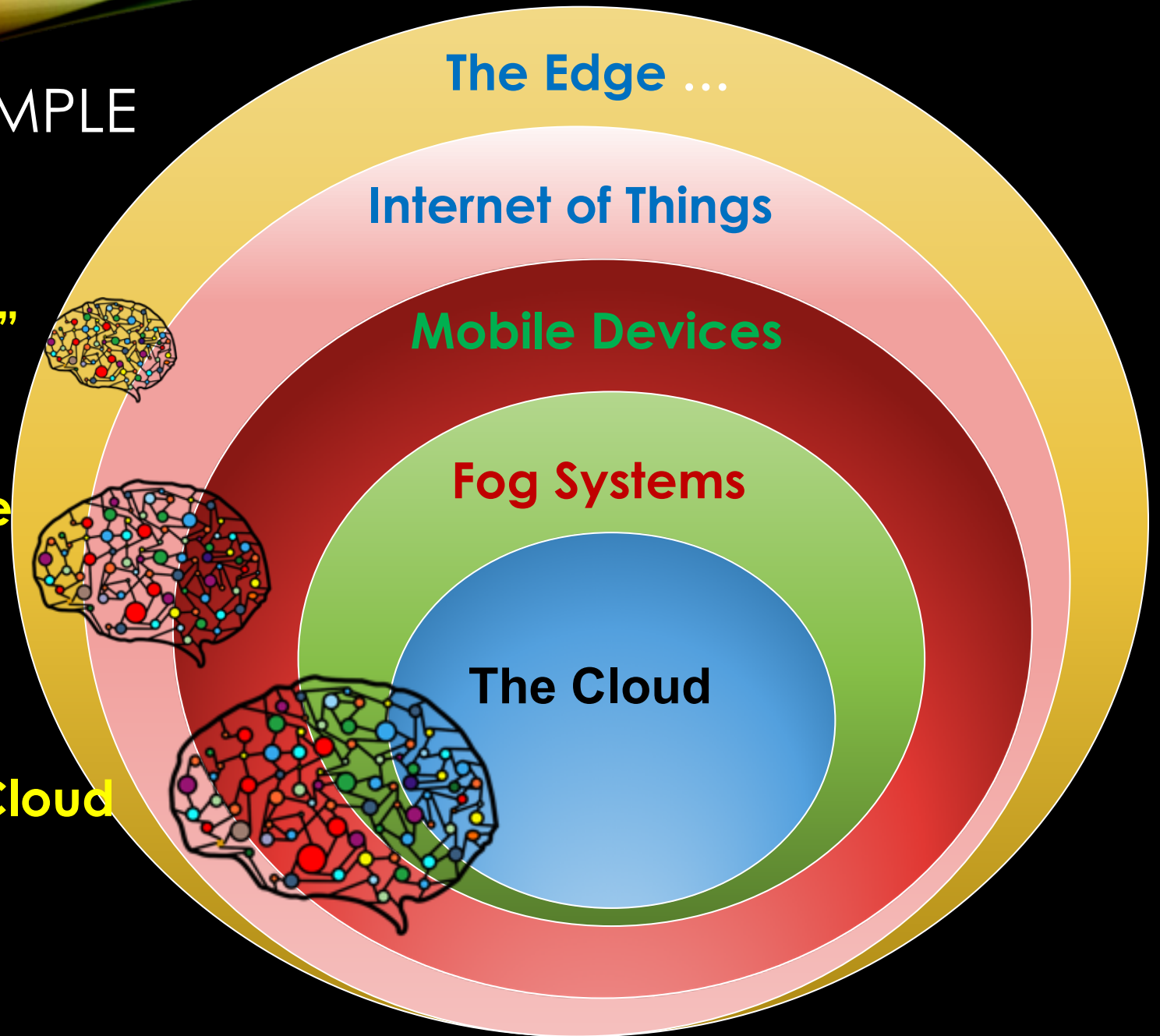


THIS IS A GOOD EXAMPLE
OF WHAT I REFER TO
AS "EDGE AI" ...

"Edge AI"

"Mobile/Edge
AI"

"Edge to Cloud
AI"



IN FACT WE LIKE EDGE-AI SO MUCH WE BUILT A HAND-HELD DEVICE TO SHOW WHAT IT IS CAPABLE OF ...



POINT #4 FROM TODAY'S TALK:
CHEAP CAMERA TECH + ARTIFICIAL INTELLIGENCE
LEADS TO MANY NEW USE CASES & APPLICATIONS –
HIGHLY DISRUPTIVE!

And this is just one example of AI disruption at the edge!

WHAT DOES THIS ALL MEAN FOR MEMORY TECHNOLOGY?

- The AI Chips themselves feature novel, memory-centric architectures
 - New design opportunities & challenges
 - Some architectures will become big winners and create new 'memory standards'
- But, IMHO, this is the thin edge of the wedge for Memory Technology opportunities

AI HAS TO BE TRAINED!

- It needs a ton of data to get good results!
- And real-world data is complex and expensive to obtain
- Lets consider a simple example – suppose I want to obtain data to train a Driver Monitoring System (DMS) for an Automotive Manufacturer
 - At minimum the DMS will have:
 - (i) a facial pose estimator,
 - (ii) an eye-gaze tracker

DRIVER MONITORING SYSTEMS REQUIRED IN EU

FROM 2021



NUI Galway
OÉ Gaillimh



C3I

Cognitive, Connected & Computational Imaging

DATA WE NEED TO TRAIN EXAMPLE AI

- (i) Capture Video Data of Subjects Face (while driving)
- For each video frame we also need to measure:
 - (ii) **Head distance** from the camera; also eyes & other facial key-points
 - (iii) **Head pose** relative to the camera position
 - (iv) Direction of **eye-gaze** (two eyes)
 - (v) **Lighting conditions** (ambient & directional - e.g. sun, car headlights, etc)
- & Ideally we need data from **100s of subjects** – variations in ethnic origin, gender, face & body sizes, glasses, facial hair, etc ...

& WHAT HAPPENS WHEN, FOR EXAMPLE, THE CAMERA LOCATION CHANGES?

- A New Cabin design or Different Model Vehicle?
- Gather Data all over again?
- 3-4 Engineers working for > 1 month with 100+ subjects, data acquisitions, post-processing
- [Industry team I work with carries datasets around on **10TB HDDs !!!!**]

THE “NEW” SD CARD FOR AI ENGINEERS?



POINT #5 FROM TODAY'S TALK:
GATHERING TRAINING DATASETS TAKES A LOT OF
TIME, EXPERTISE, SUBJECTS AND IS VERY, VERY
COSTLY!

> \$100,000

BUT, THINKING TANGENTIALLY

VIRTUAL REALITY & ANIMATION TOOLS ARE
NOW GOOD ENOUGH TO SIMULATE 'REAL
DATA' ...

ULTRA-REALISTIC FACIAL ANIMATIONS

<https://youtu.be/TxErDzslIdKI>



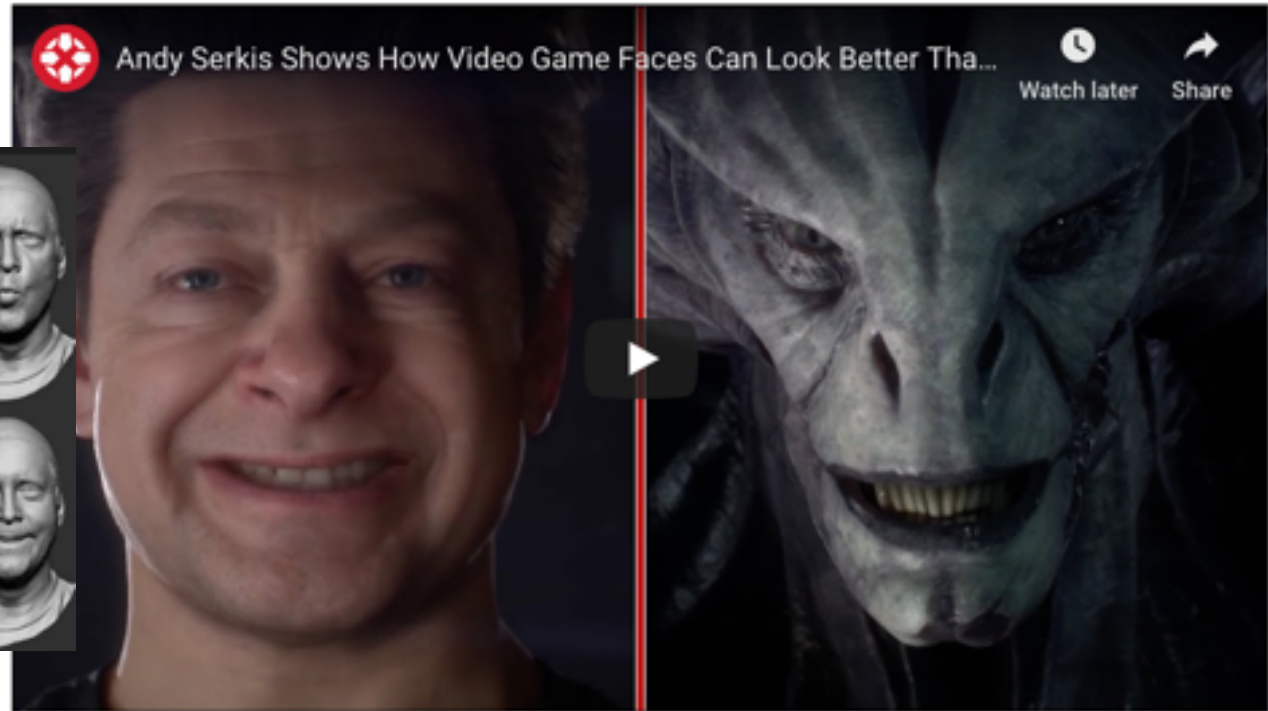
>\$2.7B



Digital Version of Andy Serkis Recites 'Macbeth' to Show the Future of Performance Capture in Gaming

by Justin Page at 10:16 AM on March 23, 2018

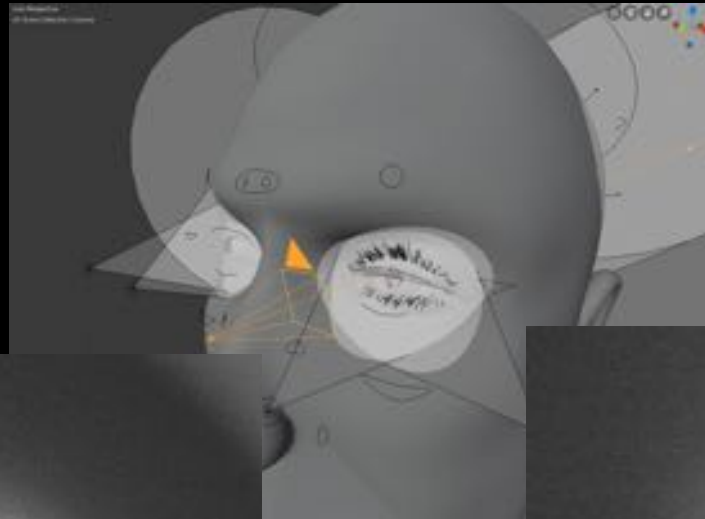
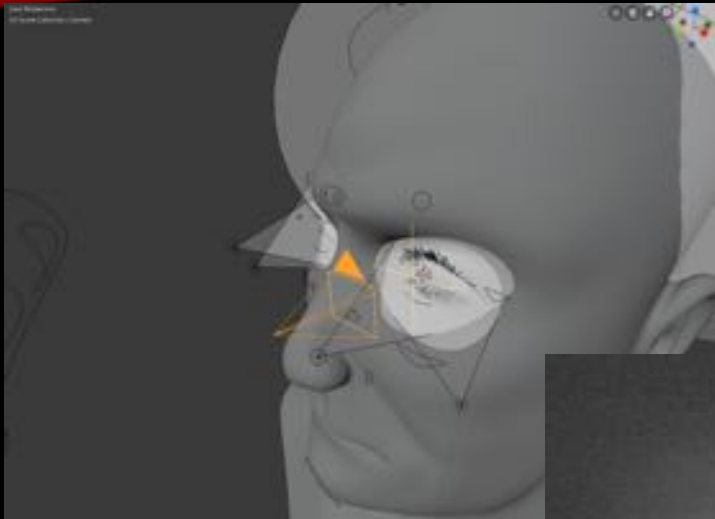
Facebook Twitter Flipboard Pinterest Tumblr More



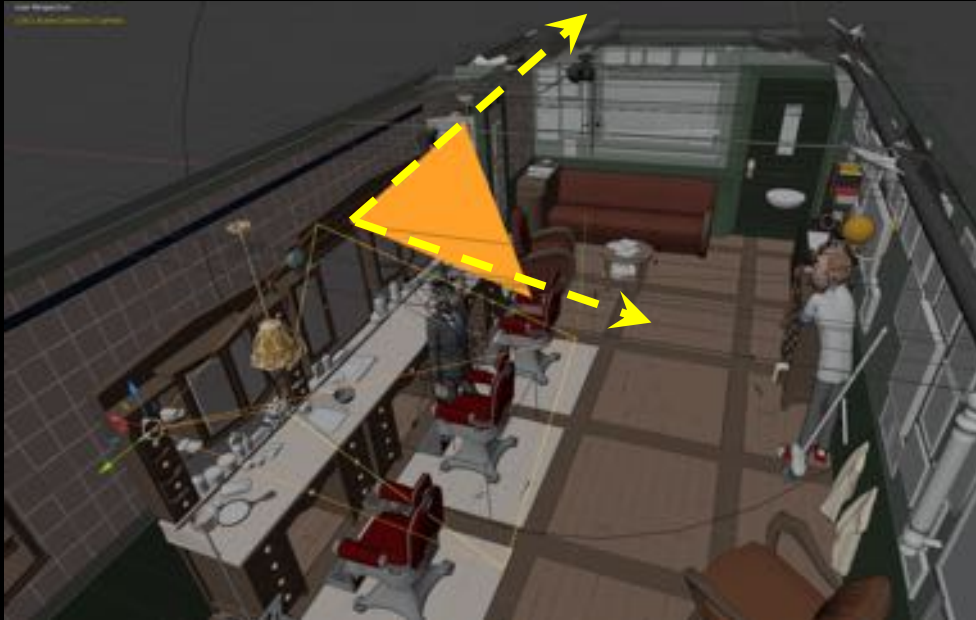
FULL BODY CHARACTERS & ANIMATIONS



SOPHISTICATED FACE & EYE MODELS



COMPLEX 3D SCENES, ADJUSTABLE CAMERA FOV & HIGH QUALITY RENDERING ENGINES



MANY BENEFITS OF 'FAKE DATA' ...

- Provides a **more accurate ground truth** (*depth, pose angle, object dimensions*) than real data, and of as many 3D points as needed ...
- Large numbers of **2D viewpoints** can be rendered from a single 3D scene ...
- **Camera models, locations & paths** can be controlled and a 3D scene re-rendered with new cameras & locations
- Most **annotations can be automated** as part of the rendering process ...
- Data is essentially **free of noise & blur** (but these can be simulated if needed) ...

POINT #6 FROM TODAY'S TALK:
**FAKE DATA IS BETTER, EASIER AND CHEAPER
THAN REAL DATA**



BUT THERE ARE SOME CHALLENGES!

- **Uncompressed Data Rates** for generated image frames
 - 16-bit uncompressed 1080p24 4:2:2 file will have the size = $(1920 \times 1080 \times 16 \times 3 \times 1.05) / (1024 \times 1024 \times 8) = 12.46$ MB per frame.
 - The data rate for such a file = $12.46 \times 24 \text{ fps} \times 0.667 = 200$ MB/s
- For a **Driver Monitoring System**:
 - 2 fixed camera viewpoints = 400 MB/s
 - + take into account additional data-elements – a **16 bit depth map** adds another c.100 MB/s; other metadata (**face & eye metadata; lighting info**) might add another 30-50 MB/s so lets ballpark at 600 MB/s
 - + assume a **3 minute random head & eye-motion cycle**, with **200 subjects** will generate a training dataset of (600x180x500 MB) = **22 TB dataset**
 - About **0.1 TB per subject** – data bandwidth challenges for Training System!

THE “NEW” SD CARD FOR AI ENGINEERS?



**POINT #7 FROM TODAY'S TALK:
'NEW' STORAGE & DATA BW CHALLENGES FROM
TRAINING WITH LARGE VOLUMES OF SYNTHETIC DATA**

....

THIS IS STILL VERY MUCH A WORK IN PROGRESS!

<http://bit.ly/neural2019>

??? QUESTIONS ???

SOME ARTICLES TO CONSIDER ...

- Privacy, Smartphones & Internet of Things
 - P. Corcoran, "The Battle for Privacy In Your Pocket" [Notes from the Editor], IEEE Consumer Electronics Magazine. **2016** Jul;5(3):3-36.
 - P. Corcoran, "Privacy in the Age of the Smartphone". IEEE Potentials. **2016** Sep;35(5):30-35.
 - P. Corcoran, "A privacy framework for the Internet of Things", In Internet of Things (WF-IoT), **2016** IEEE 3rd World Forum on 2016 Dec 12 (pp. 13-18). IEEE.
- Biometrics & Personal Authentication
 - P. Corcoran, "Biometrics and consumer electronics: A brave new world or the road to dystopia?" [Soapbox]. IEEE Consumer Electronics Magazine. **2013** Apr;2(2):22-33.
 - P. Corcoran, C. Costache, "Biometric Technology and Smartphones: A consideration of the practicalities of a broad adoption of biometrics and the likely impacts", IEEE Consumer Electronics Magazine, 5 (2), pp. 70–78, **2016**.
 - P. Corcoran, C. Costache, "Smartphones, Biometrics, and a Brave New World", IEEE Technology and Society Magazine. **2016** Sep;35(3):59-66.

MORE ARTICLES TO CONSIDER ...

- Mobile Edge, IoT & Edge-AI
 - Corcoran P, Datta SK., Mobile-Edge Computing and Internet of Things for Consumers: Part II: Energy efficiency, connectivity, and economic development. IEEE Consumer Electronics Magazine. 2017 Jan;6(1):51-2..
 - Corcoran P, Datta SK. Mobile-edge computing and the Internet of Things for consumers: Extending cloud computing and services to the edge of the network. IEEE Consumer Electronics Magazine. 2016 Oct;5(4):73-4.
 - Corcoran P. The Internet of Things: why now, and what's next?. IEEE consumer electronics magazine. 2016 Jan;5(1):63-8.
- Deep Learning & Consumer Electronics use cases
 - Bazrafkan S, Corcoran PM. Pushing the AI envelope: merging deep networks to accelerate edge artificial intelligence in consumer electronics devices and systems. IEEE Consumer Electronics Magazine. 2018 Mar;7(2):55-61..
 - Lemley J, Bazrafkan S, Corcoran P. Deep Learning for Consumer Devices and Services: Pushing the limits for machine learning, artificial intelligence, and computer vision. IEEE Consumer Electronics Magazine. 2017 Apr;6(2):48-56.
 - Bazrafkan S, Javidnia H, Lemley J, Corcoran P. Depth from monocular images using a semi-parallel deep neural network (SPDNN) hybrid architecture. arXiv preprint arXiv:1703.03867. 2017 Mar 10.

AND EVEN MORE ARTICLES TO CONSIDER ...

- Deep Learning & Biometric use cases
 - Bazrafkan S, Thavalengal S, Corcoran P. An end to end deep neural network for iris segmentation in unconstrained scenarios. *Neural Networks*. 2018 Oct 1;106:79-95.
 - Varkarakis V, Bazrafkan S, Corcoran P. A Deep Learning Approach to Segmentation of Distorted Iris Regions in Head-Mounted Displays. In *2018 IEEE Games, Entertainment, Media Conference (GEM) 2018 Aug 15 (pp. 1-9)*. IEEE..
 - Ungureanu AS, Thavalengal S, Cognard TE, Costache C, Corcoran P. Unconstrained palmprint as a smartphone biometric. *IEEE Transactions on Consumer Electronics*. 2017 Aug;63(3):334-42.
 - Bazrafkan S, Nedelcu T, Filipczuk P, Corcoran P. Deep learning for facial expression recognition: A step closer to a smartphone that knows your moods. In *2017 IEEE International Conference on Consumer Electronics (ICCE) 2017 Jan 8 (pp. 217-220)*. IEEE.
 - Lemley J, Kar A, Drimbarean A, Corcoran P. Efficient CNN Implementation for Eye-Gaze Estimation on Low-Power/Low-Quality Consumer Imaging Systems. *arXiv preprint arXiv:1806.10890*. 2018 Jun 28.
 - Lemley J, Kar A, Drimbarean A, Corcoran P. Convolutional Neural Network Implementation for Eye-Gaze Estimation on Low-Quality Consumer Imaging Systems. *IEEE Transactions on Consumer Electronics*. 2019 Feb 15.