



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

| | |
|-----------------------------|---|
| Title | Deep learning for consumer devices and services 2-AI gets embedded at the edge |
| Author(s) | Corcoran, Peter; Lemley, Joseph; Costache, Claudia; Varkarakis, Viktor |
| Publication Date | 2019-09-02 |
| Publication Information | Corcoran, P., Lemley, J., Costache, C., & Varkarakis, V. (2019). Deep Learning for Consumer Devices and Services 2—AI Gets Embedded at the Edge. IEEE Consumer Electronics Magazine, 8(5), 10-19. doi: 10.1109/MCE.2019.2923042 |
| Publisher | IEEE |
| Link to publisher's version | https://dx.doi.org/10.1109/MCE.2019.2923042 |
| Item record | http://hdl.handle.net/10379/15482 |
| DOI | http://dx.doi.org/10.1109/MCE.2019.2923042 |

Downloaded 2024-04-19T18:33:46Z

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



Deep Learning for Consumer Devices & Services 2 – AI gets Embedded at the Edge

Peter Corcoran, Joe Lemley, Claudia Costache & Viktor Varkarakis.

| | |
|--|----------------------------|
| 1. INTRODUCTION | 1 |
| 1.1. AI MOVES TO THE EDGE | 2 |
| 2. EXAMPLE CE USE CASES AND DL SOLUTIONS | 3 |
| 2.1. EYE-GAZE SYSTEMS | 3 |
| 2.2. BIOMETRICS AND DEVICE AUTHENTICATION | 5 |
| 2.3. IMMERSIVE AUDIO FOR MIXED REALITY HEADSETS | 7 |
| 2.4. IMAGE SIGNAL PROCESSING PIPELINE IN A CAMERA | 8 |
| 3. CHALLENGES FOR AI DEPLOYMENTS IN CONSUMER ELECTRONICS? | 9 |
| 3.1. THE PROBLEM-SPECIFIC NATURE OF AI SOLUTIONS | 10 |
| 3.2. DEVICE-SPECIFIC ASPECTS | 10 |
| 3.3. THE DATA BOTTLENECK | 11 |
| 4. WHAT IS NEXT FOR <i>EDGE-AI</i>? | 12 |
| 4.1. EMERGING OPPORTUNITIES FOR <i>EDGE-AI</i> | 12 |
| 4.2. CHALLENGES FOR <i>EDGE-AI</i> | 13 |
| | CONCLUDING THOUGHTS |
| | 14 |
| 5. | 14 |
| 6. BIBLIOGRAPHY | 15 |

1. Introduction

This article follows-on from our earlier publication [1], “Deep Learning for Consumer Devices and Services: Pushing the limits for machine learning, artificial intelligence, and computer vision”. In that earlier work we discussed the emerging importance of deep learning (DL) techniques for the next generation of consumer electronic devices and services.

Our vision at that time was that new embedded hardware solutions would enable advanced new devices and services that would incorporate convolutional neural network (CNN) based artificial intelligence across a broad range of new CE devices and services. This article

introduced many of the basics of deep learning and the supporting tools and methodologies and we referred to the growth of new device categories, relying on cloud-based AI, such as smart-speakers and mobile voice assistants.

Since that time the explosive growth in new AI research has continued and there has been a growing interest and investment by industry into moving key element of the AI away from the cloud towards the sensors and the embedded devices themselves. I originally wrote about this new trend in editorials for two special issues of CE Magazine [2], [3] towards the end of 2016. At that time industry was focussed on OpenFog, an initiative to define a new generation of low-latency services for the Internet of Things. But more recently we have seen the reach of AI moving onto the device itself with many companies and researchers focussing on developing FPGA based solutions [4] and most recently embedded AI hardware accelerators [5], [6].

1.1. AI Moves to the Edge

Most of the large semiconductor manufacturers are currently working on a new generation of AI accelerator chipsets so it is only a matter of time before we begin to see widespread deployments of neural networks in the device itself, independent of any network connection or services. In fact such technology is already incorporated into the latest generation of mobile devices, high-end television panels, professional digital cameras and many new automotive subsystems.

I like to refer to this embedding of AI into the device itself as *Edge-AI*, and distinguish it from edge AI services provided over a local network link, such as envisaged by the proponents of Fog Computing solutions. There is a role for both networked based and device based AI, but it is the recent emergence of AI implementations *on-device* that I find most exciting as a CE engineer.

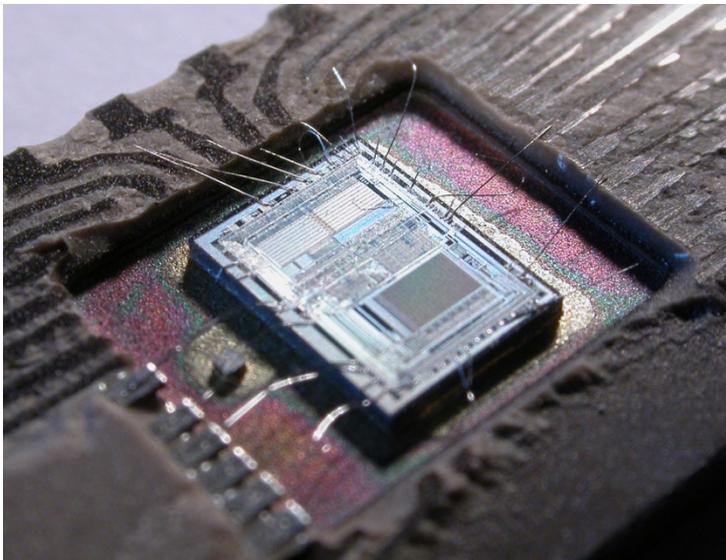


Figure 1: A system-on-chip solution can incorporate dedicated CNN accelerator circuitry to enable low-power AI in a consumer device. (license info at: https://commons.wikimedia.org/wiki/File:Intel_8742_153056995.jpg)

Improvements in the energy efficiency of dedicated AI accelerators over today's GPU based solutions will provide significant improvements in the energy required to power Edge-AI, and this in turn makes it an ideal solution for the challenging data-processing problems introduced

by a new generation of lightweight, battery-powered wearable and internet-of-things (IoT) devices.

In this article we take a look at some of the progress that we have seen in *Edge-AI* over the past two years and describe some examples of practical problems that have been tackled with deep neural networks (DNNs) over that time. Some of these are contributions made within our research team, other are drawn from the literature, but each of them illustrates how DNNs are being used today and how DNN-based *Edge-AI* will be at the core of many new devices, systems and services that emerge over the next decade. We will also take a look at some of the challenges posed by this migration of AI onto the device itself and comment on several of the new research challenges we anticipate over the next 4-5 years.

2. Example CE Use Cases and DL Solutions

A good starting point for our discussion is how can DL solutions be used beneficially in consumer devices? Where can edge-AI improve performance and operational efficiency to a point where the benefits outweigh the costs of incorporating an inference engine or platform into the system.

Looking at the research literature and the evolution of consumer electronics products over the past decade once clear area where DL can add value is for computer vision applications, in particular for wireless, internet-of-things (IoT) devices. The cost of a complete VGA camera module has dropped below \$1.00 and the cost of a CMOS image sensor is almost negligible in today's devices. Thus many interesting uses of DL inference at the edge will focus on providing advanced image analysis capabilities to low-cost CMOS sensors. Let us begin by considering some examples of new CE applications that *Edge-AI* can enable.

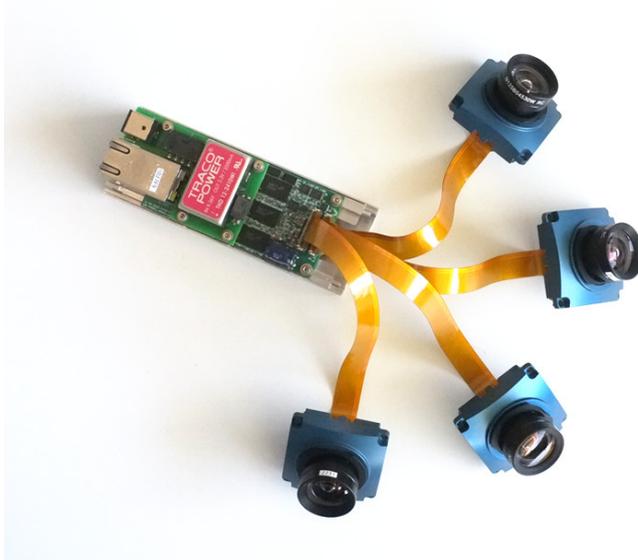


Figure 2: CMOS based camera modules are now so inexpensive that embedded devices often include multiple units. (licence info at:

https://commons.wikimedia.org/wiki/File:NC393_multisensor_camera_kit_for_developers.jpeg)

2.1. Eye-Gaze Systems

Eye tracking has begun to find uses in a variety of consumer applications. A range of works based on the use of gaze information for consumer platforms like automotive (for driver monitoring), augmented and virtual reality (for foveated rendering and immersive experiences), smartphones and TV (for gaze based menu selection and navigation) have been described in a recent review [7]. Recent works in this field include [8] where gaze duration and

patterns are used to assess how easily users can navigate an interface of a connected self-injection system. Gaze data is used as an indicator to study the usability, efficiency and ease of use of the drug delivery device. Similarly, the design effectiveness of observation charts in a hospital is evaluated in [9] by comparing viewing patterns of users derived from eye tracking data. Eye movements and pupillary response are used as indicators of cognitive load while users answered mathematical questions in [10] and for studying cognitive processes and learning aspects in [11]. In [12], the influence of emotions on the visual acuity of users was studied which showed that eye movements like fixations and saccades clearly respond to levels of stress.

Eye gaze estimation is a “must have” feature for the latest driver monitoring systems and is crucial for the functioning of AR/VR systems. For these applications, deep learning can either be applied as an end-to-end solution [13] or as a component of a more traditional gaze estimation pipeline, such as iris segmentation [14], [15].



Figure 3: State of art Driver Monitoring System – screenshots courtesy of Xperi, Inc.

End-to-end approaches to eye gaze estimation are state of the art for uncontrolled environments where the camera is at a distance, and such methods are closing the gap in AR/VR settings [16]. Real world Edge-AI must balance accuracy requirements with power availability and speed. Typical Edge products have strict limits on the number of MAC operations per second, and these limits often preclude the use of the largest and most popular networks without significant optimization. Similar approaches can have vastly different performance. Many papers report performance on GPUs but it is important to consider their speed on the low cost commodity processors that are typical of Edge Devices.

For example, one published gaze estimation solution operates at approximately 120 FPS on a modern GPU, but 0.0043 FPS on an ARM processor (approximately one frame every 4 minutes) while another network with similar accuracy runs at 92.59 FPS on the same ARM processor, or for a $\sim 1.5\%$ increase in accuracy, at 19 FPS [13].

2.2. Biometrics and Device Authentication

Biometrics have a huge potential to solve many pressing problems related to security and privacy for today's CE devices [17]. This potential has been realized to some extent on today's mobile devices, but the real breakthrough will come when biometric technology can be incorporated into lower power devices such as wearables and internet-of-things (IoT) peripherals.

A key aspect here is that biometric acquisition and processing should occur on the device or peripheral in order to secure the privacy of the biometric [18]–[20]. Realistically this can only be achieved through leveraging *Edge-AI* to enable more accurate and verifiable acquisition and more energy-efficient authentication of the biometric.

Let us next consider a few examples of emerging use cases for biometrics.

Automotive biometrics offer a use case where *Edge-AI* offers a practical solution. Our cars already have cameras to monitor the driver and detect drowsiness so it is quite straightforward to incorporate technology to verify the identity of the driver once they sit into the vehicle. External cameras, typically used to detect pedestrians or replace physical mirrors can be repurposed to provide an initial authentication based on facial recognition and supplement today's keyless entry systems.

A key challenge here is acquiring the biometric in unconstrained conditions which can be challenging [21], [22]. While facial biometrics offer a reasonable level of security and verifiability they are prone to a range of attacks [23]–[25] and even iris biometrics can be challenged [26], [27].

Incorporating all the elements of a biometric authentication chain – unconstrained acquisition, liveness verification and, finally, authentication of the registered biometric – requires significant computational resources and energy if we rely on conventional processing approaches. But each of these elements can be tackled independently, and in parallel, by neural network based solutions. As examples, large pose 3D facial segmentation [28], facial segmentation and alignment [29] and iris segmentation [14] can all be tackled with CNNs; liveness detection is also receiving much attention from researchers [30]–[33].

Biometrics are also important for wearable devices. Consider the emerging category of augmented and mixed reality (AR/MR) headsets where there is no keyboard, yet these devices

will inevitably be connected to an online account. So how can the device authenticate the user in a frictionless manner?

Fortunately these AR/MR headsets need to track the user's eye-gaze in order to accurately render objects onto the real-world scene. In order to track eye-gaze a user-facing camera is required to analyse the eye-region, which also includes monitoring and segmenting the eye-iris and pupil to determine the direction of gaze. It now becomes a simple task to authenticate the segmented iris region [34].



Figure 4: User-Facing Cameras for a Mixed Reality Headset – courtesy of pupil-labs.com [35](Check with Viktor about use agreement from pupil-labs).

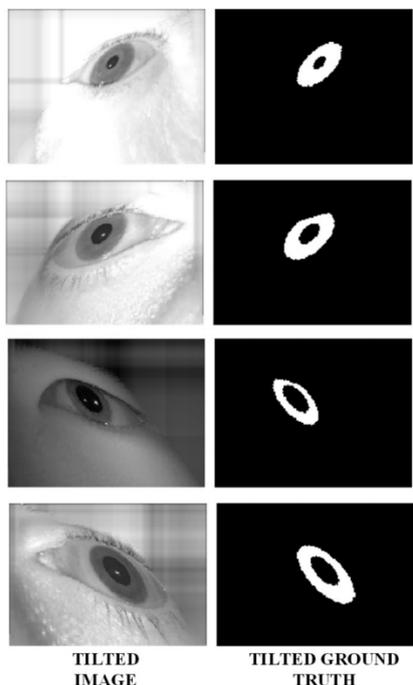


Figure 5: Off-axis iris regions; accurate, per-pixel, segmentation is essential for practical user authentication. [15]

2.3. Immersive Audio for Mixed Reality Headsets

As we've mentioned, the emerging category of augmented and mixed reality (AR/MR) headsets has become a driver of new innovation. This is another interesting area for deep networks research. Mixed reality implies adding a mixture of new visual elements onto the visual field of view, but it also implies adding non-visual immersive elements and audio can provide even more compelling immersion than visual cues. If you doubt that, then consider the music score of a horror movie which creates and maintains the underlying atmosphere of impending doom, and then leads you cleverly into those 'big scares' that are the cornerstone of the genre.

But adding audio to the perceived environment, particularly when you are seeking to maintain an illusion of the real-world, is quite challenging. Each of us has a unique ear canal and our brains process and perceive audio in a highly personalized way. And our audio senses are also attuned to visual cues in our environment, so that we anticipate changes in the environmental acoustics. Thus if you move into a large cathedral you expect that your footsteps and voice will echo more. In a room with carpets and soft furnishings the acoustics are more muffled.

This presents an interesting challenge in that acoustic cues should be adapted to match the surroundings of the wearer of a headset, or the illusion of immersion is lost. But analysing and evaluating the perceived environment with conventional image processing would not be possible on a wearable headset where energy usage is even more critical than on a smartphone. And thus another excellent CE use-case for advanced neural networks presents itself – scene analysis [36], [37] and materials recognition [38] combined with depth [39], [40] can help build a detailed analysis of the surrounding acoustic environment.

And importantly, with a state-of-art *Edge-AI* neural accelerator we can run multiple neural networks in parallel at a fraction of the power budget of a GPU based computational unit. Some exciting new developments in immersive multimedia experiences are to be expected, in the near future.

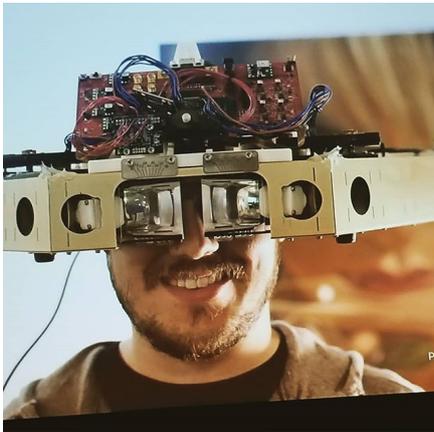


Figure 6: Multimedia, Mixed Reality Headset early Magic Leap Prototype

<https://www.flickr.com/photos/egoant/40067827415> (Peter to try get permission & hi-res version from Magic Leap)



Figure 6 (alternative image): Multimedia, Mixed Reality Headset

<https://commons.wikimedia.org/wiki/File:Ar-vs-vr.jpg>



Figure 6 (alternative): Mixed Reality with a VR Headset, [courtesy of Pierre Faure

https://commons.wikimedia.org/wiki/File:Mixed_Reality_with_a_Virtual_Reality_Headset.png]

2.4. Image Signal Processing Pipeline in a Camera

Processing the raw Bayer data from an image sensor is a classic example where camera engineers and photographic experts have devoted hundreds of man-years of effort to create a highly specialized image processing pipeline. Bryce Bayer's original patent [41] filed in 1976 is a classic example of an engineering compromise that has stood the test of time. It employs a color filter array with twice as many green pixels as red or blue pixels to match the color sensitivity of the human visual system (HVS). And it works so well that it has become the basis of modern digital imaging.

Now as computer rendered images use an equal number of red, green and blue pixels it becomes necessary to convert the HVS-compatible sensor data to RGB images which is where the camera pipeline comes in. But given that image sensing can occur in many different lighting conditions and is subject to other environmental conditions it is actually a very complex process to achieve a high-quality final RGB image from the raw sensor data.

Until recently this 'magic' happened within a complex set of image analysis and processing algorithms known as the *image processing pipeline* (IPP) [42]. In fact many consumer imaging devices feature a dedicated *image signal processor* (ISP) – a dedicated hardware component – to handle this conversion process.

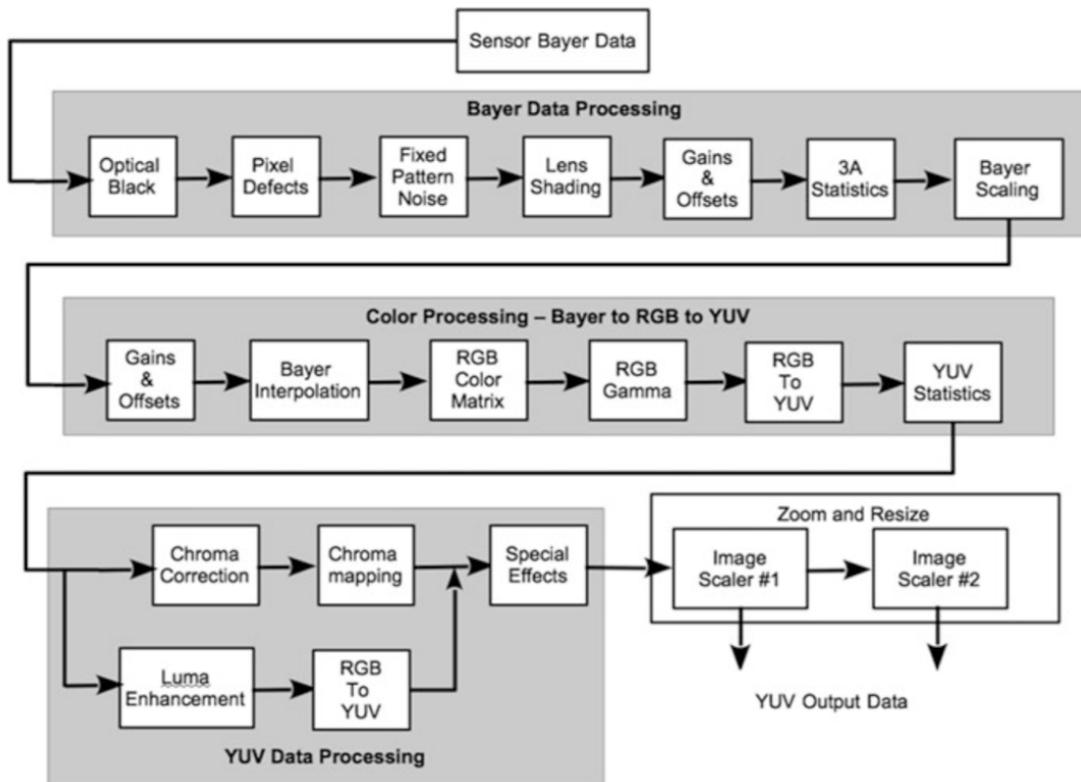


Figure 7: Detailed view inside the Image Processing Pipeline (IPP) of a typical Digital Camera.

(From my chapter in ref [42])

But now, thanks to the magic of deep learning methodologies and the potential to implement the corresponding CNNs in hardware it has become possible to consider replacing the IPP in an imaging system. This work began with studying the replacement of the de-mosaicking step of image conversion [43], [44] quickly followed by the idea to replace the two key steps of de-noising and de-mosaicking from the IPP with a single CNN network [45]–[47].

Other authors have further extended this idea to completely replace the traditional IPP [48] or alternatively to learn the detailed camera model embodied in an existing IPP [49], [50] which has potential both to identify the source of processed images, but could also lead to reprogrammable IPPs based on a CNN hardware accelerator. Imagine that you can completely reprogram how your camera captures & develops raw images ‘on the fly’. Well this approach is beginning to make its way into actual products so expect to see some exciting new features in higher-end digital cameras over the next couple of years!

And once you start to think about replacing the IPP with a convolutional network new ideas will emerge such as adapting the camera to facilitate “Learning to see in the Dark” [51]. In this particular example researchers have used a CNN to solve a key challenge for today’s smartphone cameras – that of capturing images in low-illumination conditions.

3. Challenges for AI deployments in Consumer Electronics?

There has been a lot of very rapid progress with AI technologies since our last article, but there remain many challenges that are specific to *Edge-AI* and the implementation of solutions in

consumer devices. Our recent work on a number of quite typical example of CE-device problems have highlighted these challenges. ...

3.1. The Problem-Specific Nature of AI Solutions

Every practical problem that we solve with *Edge-AI* is typically a part of a larger problem set. Focussing on a specific problem typically allows the research engineer to accurately define the data characteristics and the criteria required to solve that particular problem. Let's consider the task of iris authentication where there is a processing pipeline as shown in *Figure 8*.

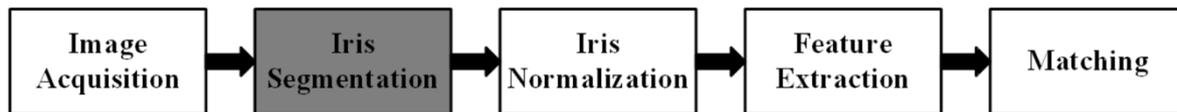


Figure 8: Iris Authentication processing pipeline. [15]

Note that in this sequence of tasks some tasks can employ proven techniques such as iris-matching which has been deployed and verified for more than 2 decades by the biometrics research community. However the adoption of the authentication on mobile and wearable devices has now introduced new processing tasks such as the segmentation of unconstrained iris regions and the normalization of these to serve as input to established feature extraction and matching algorithms.

These new task represent more recent challenges and it is now well appreciated that iris segmentation is the predominant source of errors in mobile and wearable devices [14], [21].

It is important for the research engineer to appreciate this need to break down problems in this way. While deep learning is a powerful tool and can often achieve very impressive levels of accuracy it is important to appreciate that neural networks are also susceptible to adversarial data and if appropriate training data is not provided they can easily learn the wrong features from a poorly designed set of data samples. This is a challenge we'll return to in some of the subsequent discussion and it is a very open-ended challenge.

3.2. Device-Specific Aspects

A unique aspect of applying deep learning techniques to Consumer Electronics is the device-specific nature of consumer data.

The complexity of the processing pipeline in a digital camera was previously discussed and it is well-known by imaging experts that every production camera has unique characteristics. Indeed recent research has shown that these characteristics can be learned and images can be uniquely associated with a particular image processing pipeline [49], [50]. This observation will apply across most sensing capabilities of consumer devices. Thus the collection of data from any particular consumer device, but it video, audio, motion or other forms of data collection will invariably exhibit some unique characteristics.

This leads us to the observation that to achieve the more accurate and optimal performance from an deep learning solution the network should be trained on device-specific datasets. This is, in fact, well-known in the CE industry where more traditional algorithms have always been tailored to individual models of production devices. Indeed sometimes the same device, but manufactured in a different facility, can perform differently due to variations in the manufacturing process, the local environment or the calibration procedures applied.

Thus, when applying deep learning methodologies to CE problems it is important to bear in mind that optimal performance will be achieved by tuning an AI network to device specific data. However the corollary to this is that sometime a network that is tuned to a specific device may not perform well on other devices. While we have not explored this phenomenon across

enough different problems cases our current experience suggests that a two-step approach makes sense. At stage one a network should be designed and tuned on a generic dataset representative of a range of broadly similar data acquisition systems (e.g. data acquired from multiple f2.0, 12 MP cameras). Once the performance of this network is tuned to an acceptable level then stage two should involve additional tuning of the network on data from a specific production stream of camera modules. Our experience shows that additional performance can be achieved, but at the expense of a loss of performance for the other streams of camera module.

3.3. The Data Bottleneck

This brings us to what is undoubtedly the greatest challenge for *Edge-AI* – that of obtaining the large datasets that are typically needed to achieve convergence of the training process for a deep neural network. Data acquisition is time consuming, but more importantly every problem also needs a *ground truth* to train against!

Taking the iris segmentation problem as an example, every iris image has to be marked-up to obtain the ground truth. Even where some form of automated mark-up is possible there should be a manual check to detect mark-up failures. This is very time-consuming and the costs can quickly mount-up when large datasets are involved.

This data bottleneck is a problem in general for deep learning researchers, but it is even more so for CE engineers who may need to adapt and re-train networks for multiple device models or in new use-case geometries such as off-axis iris authentication.

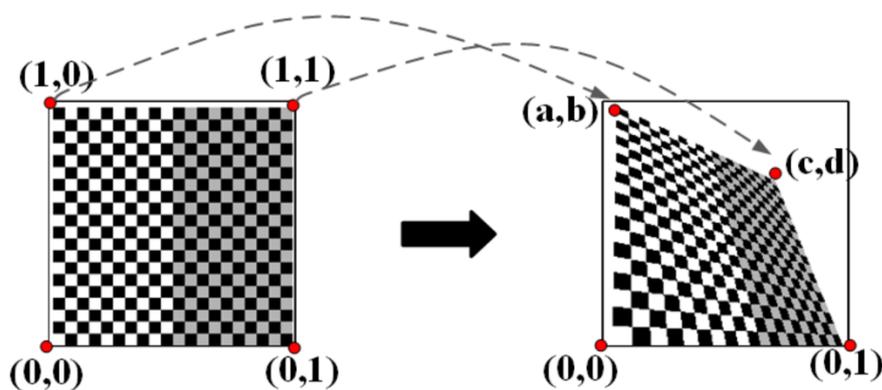


Figure 9: Novel data augmentation strategies can help grow the available training dataset for a specific problem; in this example we show how iris data samples can be transformed to solve off-axis iris segmentation [15].

There are some approaches that can help here such as data-augmentation where a seed-dataset is modified and transformed in particular ways to match the underlying problem as is done in the case of iris segmentation in [14], [52]. It has also been shown that we can train a network to learn how to make ‘new’ data by combining existing samples in its convolutional layers [53]. And another technique that pairs two deep-learning networks in a configuration known as a generative adversarial network (GAN) enables researchers to train a *data generator* that learns the key features of an existing dataset and can then make ‘new’ random data samples that match these [54].

Ultimately data is the biggest challenge for *Edge-AI* and its successful deployment in consumer devices. There is a need for improved approaches to build the large datasets that are needed and to determine and record the corresponding ground truth associated with each individual acquisition. Engineers need improved tools and data management methodologies to increase the efficiency and accuracy of dataset acquisition and ground truth estimation. This will,

without doubt, become a key focus area for future research and I hope to write in more detail on this topic in an upcoming special issue of CE Magazine.

4. What is Next for *Edge-AI*?

When we last considered the state-of-art for machine learning, artificial intelligence, and computer vision in the context of consumer electronics systems it was clear that new hardware accelerators for embedded devices were becoming available and, at that time, the AI-stick from *Movidius* was discussed as a practical example. Since that time this company was acquired by Intel and a 2nd generation of the AI-stick is now available.

But there are now several other AI accelerators from mainstream players such as Nvidia's Jetson family of devices and in mid-2018 Google introduced an "Edge" version of its Tensor Processing Unit (TPU) which can integrate with the well-known deep learning platform, *Tensorflow*. Outside of the mainstream players there are many start-ups, spin-outs and in-house projects working to deliver new low-power neural network accelerators.

Over the next 1-2 years it is likely that many of you will be developing systems and products on these new AI platforms, if you are not already doing so. *Edge-AI* with the promise of intelligent devices that have minimal power requirements – some able to run on a coin sized battery for months without replacement - allows devices to have full or partial functionality that only last year would have required a large, power-hungry GPU or an always-on connection to the cloud.

Our group has had the opportunity to build some interesting prototypes with these hardware accelerators and some of you may have attended our "hands-on" workshop at IEEE GEM 2018 which was a great success – you can view a Twitter "moment" of the conference at <https://twitter.com/i/moments/1061174127342559232>. In our current graduate program lab classes we help students build a handheld computer-vision terminal that can implement the Yolo one-shot object detector. This runs happily at 30 fps on a Raspberry PI coupled with the Intel AI-stick.

It a fair statement that the age of embedded AI is now a practical reality. The open questions are how it will impact on today's technology and where are the new challenges and opportunities that a broader adoption of *Edge-AI* will bring?

4.1. Emerging Opportunities for *Edge-AI*

One area where *Edge-AI* will have enormous impact is on personal privacy. It is difficult for people to feel comfortable about personal privacy when video data is constantly being uploaded and processed in the cloud. By processing data on the device, even within the sensor module we can avoid transmitting raw data over networks, or storing data in cloud repositories where it offers an attractive target for cyber criminals.

One example where *Edge-AI* will have significant impact are *driver monitoring systems* (DMS) which will be mandated by the EU in 2022. These are already deployed in most high-end vehicles, and require devices that can instantly and intelligently react in cases where a driver is distracted or impaired. They offer a stepping stone to autonomous vehicles, but also pose a significant design challenge in the context of the EU's *General Data Protection Regulation* (GDPR). While 5G technology can arguably perform the advanced processing required by DMS this approach also requires moving data off-vehicle with associated data security and privacy issues. In contrast *Edge-AI* enables data processing to occur within the DMS and can even be restricted to a secure compute unit within the sensing sub-system. Placing computation as close as possible to the sensor allows reduction in latency and cost while increasing privacy and usability.

Other examples of new opportunities lie in new wearable devices and smart-cities. With the introduction of a new generation of smart-glasses (e.g. Magic Leap and Hololens) and a new market in wearable audio enhancement devices, known as hearables, we are seeing a wave of devices that can directly modify our perception of the surrounding environment. This is known as *Mixed Reality* (MR), and as a concept it has been discussed since the 1960's but only recently have researchers had access to devices that can actually realize MR effects. But the computational requirements to achieve real-time perceptual analysis followed by a realistic blending of additional visual and acoustic elements into the user experience are beyond the capabilities of today's embedded-GPU solutions. The answer lies, as you might have expected, with the new generation of *Edge-AI* hardware accelerators which can achieve the required real-time data processing rates with levels of energy efficiency that are orders of magnitude lower.

In smart-cities we have an urban environment permeated with ubiquitous networks of sensors and services, but this poses some significant challenges. How can we authenticate individuals in such an environment to validate their access to services and, more importantly, how do we guard the privacy of individuals when their every move is tracked by a multitude of cameras and sensing technologies?

Again, *Edge-AI* can offer new solutions. Biometric processing can be implemented within devices so that registered users can be authenticated without a global sharing of their biometric data [55]. Once we have authenticated individuals they can be linked with a global-ID that is independent of the local device authentication using techniques such as Blockchain or Zero-Knowledge Proof (ZKP) [56]. And once the individual is globally authenticated, then they can be flagged with a 'do not track' marker and compliance with regulations such as GDPR can be explicitly recorded.

4.2. Challenges for Edge-AI

Without a doubt the biggest challenge for *Edge-AI* is that of data acquisition, annotation and curation. The training of deep neural networks requires large datasets and an accurately annotated ground truth. There are challenges in acquiring the data in the first instance, as often the data is of a personal nature – e.g. facial images. Then the annotation of a large dataset is time-consuming and costly. And often an absolute ground truth is not available.

Beyond these basic challenges there is that of data curation – if data samples are not carefully chosen to match the problem at hand then neural networks can easily learn incorrect features from the training dataset. These networks are very powerful, but they are only as good as the data that is used to train them. For AI applications in CE devices this dependency of the solution on the training data is both a challenge and an opportunity.

In a nutshell, the better aligned the training data is with the original sensor system the more robust and accurate the trained network will be. For many large online datasets the training data is typically gathered by many devices and images are harvested from various online sources. If we consider the variety of video cameras used to create a large collection of *youtube* videos, for example, it is easy to see that networks trained on such datasets won't be able to take account of device-specific characteristics.

In solving a problem for a CE device we have the potential advantage of training on device-specific data, or more likely of using transfer learning to map a generic network onto device-specific data. While this can improve the accuracy and robustness of a particular AI solution for that devices, it also requires that we can gather, or generate a large dataset of relevant data directed at the unique characteristics of that particular device model with the corresponding challenge to annotate and curate this dataset. The flip side of this is that such a solution may

no longer work if you change the sensor, or the signal processing pipeline for the device and it will be necessary to train a new solution for these modifications.

Thus the biggest challenge for *Edge-AI* is that of data. We need improved tools and methodologies to better support how we acquire, annotate and curate our training dataset. This is a fascinating topic and we plan to write a follow-up article to address it in more detail.

5. Concluding Thoughts

A lot has happened in the last two years. Our research group has continued to work on a number of fascinating problems, leveraging deep learning techniques to achieve state-of-art solutions. In parallel many other researchers have been working on and solving a broad range of problems in computer-vision, machine learning, signal-processing, data analytics and advanced sensor fusion all of which have relevance to new and emerging CE devices and services.

There has been explosive growth in the use of deep learning and advanced neural network methodologies and much of this research can be leveraged into new CE solutions. The challenge for CE engineers and researchers is how to pick and chose across this vast array of possibilities and deliver practical and useful solutions that can meet the needs of consumers. We hope this article has helped to open your eyes to some of the potential of *Edge-AI* and in some follow-up articles we'll explore some of the associated opportunities and challenges in more detail.

6. Bibliography

- [1] J. Lemley, S. Bazrafkan, and P. Corcoran, "Deep Learning for Consumer Devices and Services: Pushing the limits for machine learning, artificial intelligence, and computer vision.," *IEEE Consum. Electron. Mag.*, vol. 6, no. 2, pp. 48–56, Apr. 2017.
- [2] P. Corcoran and S. K. Datta, "Mobile-Edge Computing and the Internet of Things for Consumers: Extending cloud computing and services to the edge of the network.," *IEEE Consum. Electron. Mag.*, vol. 5, no. 4, pp. 73–74, 2016.
- [3] P. Corcoran, "Mobile-Edge Computing and Internet of Things for Consumers: Part II: Energy efficiency, connectivity, and economic development," *IEEE Consum. Electron. Mag.*, 2017.
- [4] S. Venieris, A. Kouris, C. B.-A. C. S. (CSUR), and undefined 2018, "Toolflows for mapping convolutional neural networks on fpgas: A survey and future directions," *dl.acm.org*.
- [5] A.-A. Erofei, C.-F. Druta, and C. Daniel Căleanu, "Embedded Solutions for Deep Neural Networks Implementation," in *2018 IEEE 12th International Symposium on Applied Computational Intelligence and Informatics (SACI)*, 2018, pp. 425–430.
- [6] M. Kotlar, D. Bojic, M. Punt, and V. Milutinovic, "A Survey of Deep Neural Networks: Deployment Location and Underlying Hardware," in *2018 14th Symposium on Neural Networks and Applications (NEUREL)*, 2018, pp. 1–6.
- [7] A. Kar, S. Member, and P. Corcoran, "A Review and Analysis of Eye-Gaze Estimation Systems , Algorithms and Performance Evaluation Methods in Consumer Platforms," pp. 16495–16519, 2017.
- [8] Q. Lohmeyer, A. Schneider, C. Jordi, and J. Lange, "Expert Opinion on Drug Delivery Toward a new age of patient centricity ? The application of eye-tracking to the development of connected self-injection systems," *Expert Opin. Drug Deliv.*, vol. 16, no. 2, pp. 163–175, 2019.
- [9] L. Cornish, A. Hill, M. S. Horswill, S. I. Becker, and M. O. Watson, "Eye-tracking reveals how observation chart design features affect the detection of patient deterioration : An experimental study," *Appl. Ergon.*, vol. 75, no. September 2017, pp. 230–242, 2019.
- [10] M. Shojaeizadeh, S. Djamasbi, R. C. Paffenroth, and A. C. Trapp, "Detecting task demand via an eye tracking machine learning system," *Decis. Support Syst.*, vol. 116, no. June 2018, pp. 91–101, 2019.
- [11] J. Choi, T. Oh, S. Members, and I. S. Kweon, "Human Attention Estimation for Natural Images : An Automatic Gaze Refinement Approach," pp. 1–12.
- [12] J. Przybyło, E. Kańtoch, and P. Augustyniak, "Eyetracking-based assessment of affect-related decay of human performance in visual tasks," *Futur. Gener. Comput. Syst.*, vol. 92, pp. 504–515, 2019.
- [13] J. Lemley, A. Kar, A. Drimbarean, and P. Corcoran, "Convolutional Neural Network Implementation for Eye-Gaze Estimation on Low-Quality Consumer Imaging Systems," *IEEE Trans. Consum. Electron.*, vol. 2019, no. 1, pp. 1–10, 2019.
- [14] S. Bazrafkan, S. Thavalengal, and P. Corcoran, "An end to end Deep Neural Network for iris segmentation in unconstrained scenarios," *Neural Networks*, vol. 106, pp. 79–95, Oct. 2018.
- [15] V. Varkarakis, S. Bazrafkan, and P. Corcoran, "Deep Neural Network and Data Augmentation Methodology for off-axis iris segmentation in wearable headsets," Mar. 2019.
- [16] J. Lemley, A. Kar, and P. Corcoran, "Eye Tracking in Augmented Spaces: A Deep Learning Approach," in *2018 IEEE Games, Entertainment, Media Conference, GEM 2018*, 2018.
- [17] P. Corcoran, "Biometrics and Consumer Electronics: A Brave New World or the Road to Dystopia?," *Consum. Electron. Mag. IEEE*, vol. 2, no. 2, pp. 22–33, 2013.
- [18] P. Corcoran, "The Battle for Privacy In Your Pocket [Notes from the Editor]," *IEEE Consum. Electron. Mag.*, vol. 5, no. 3, pp. 3–36, Jul. 2016.
- [19] P. M. Corcoran, "A privacy framework for the Internet of Things," in *2016 IEEE 3rd World Forum on Internet of Things, WF-IoT 2016*, 2017.
- [20] P. Corcoran and C. Costache, "Biometric Technology and Smartphones: A consideration of the practicalities of a broad adoption of biometrics and the likely impacts.," *IEEE Consum. Electron. Mag.*, vol. 5, no. 2, pp. 70–78, 2016.
- [21] S. Thavalengal, P. Bigioi, and P. Corcoran, "Iris authentication in handheld devices - considerations for constraint-free acquisition," *IEEE Trans. Consum. Electron.*, vol. 61, no. 2, pp. 245–253, May 2015.
- [22] J. J. Lozoya-Santos, R. A. Ramirez-Mendoza, S. Savaresi, J. C. Tudon-Martinez, and V. Sepúlveda-Arróniz, "Survey on Biometry for Cognitive Automotive Systems," *Cogn. Syst. Res.*, 2019.
- [23] A. Mohammadi, S. Bhattacharjee, and S. Marcel, "Deeply vulnerable: a study of the robustness of face recognition

- to presentation attacks,” *IET Biometrics*, vol. 7, no. 1, pp. 15–26, Jan. 2018.
- [24] S. Q. Liu, P. C. Yuen, X. Li, and G. Zhao, “Recent progress on face presentation attack detection of 3D mask attacks,” in *Advances in Computer Vision and Pattern Recognition: Handbook of Biometric Anti-Spoofing*, 2019, pp. 229–246.
- [25] S. Bhattacharjee, A. Mohammadi, A. Anjos, and S. Marcel, “Recent Advances in Face Presentation Attack Detection,” 2019, pp. 207–228.
- [26] S. Thavalengal, T. Nedelcu, P. Bigioi, and P. Corcoran, “Iris liveness detection for next generation smartphones,” *IEEE Trans. Consum. Electron.*, vol. 62, no. 2, pp. 95–102, May 2016.
- [27] A. F. Sequeira, S. Thavalengal, J. Ferryman, P. Corcoran, and J. S. Cardoso, “A realistic evaluation of iris presentation attack detection,” in *2016 39th International Conference on Telecommunications and Signal Processing, TSP 2016*, 2016, pp. 660–664.
- [28] A. S. Jackson, A. Bulat, V. Argyriou, and G. Tzimiropoulos, “Large Pose 3D Face Reconstruction from a Single Image via Direct Volumetric CNN Regression,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [29] Y. Zhao, F. Tang, W. Dong, F. Huang, and X. Zhang, “Joint face alignment and segmentation via deep multi-task learning,” *Multimedia Tools and Applications*, 2018.
- [30] Y. A. U. Rehman, L. M. Po, and M. Liu, “LiveNet: Improving features generalization for face liveness detection using convolution neural networks,” *Expert Syst. Appl.*, 2018.
- [31] A. Sengur, Z. Akhtar, Y. Akbulut, S. Ekici, and U. Budak, “Deep Feature Extraction for Face Liveness Detection,” in *2018 International Conference on Artificial Intelligence and Data Processing (IDAP)*, 2018, pp. 1–4.
- [32] Y. A. U. Rehman, L.-M. Po, M. Liu, Z. Zou, W. Ou, and Y. Zhao, “Face liveness detection using convolutional-features fusion of real and deep network generated face images,” *J. Vis. Commun. Image Represent.*, vol. 59, pp. 574–582, 2019.
- [33] L. Wu, Y. Xu, M. Jian, X. Xu, and W. Qi, “Face liveness detection scheme with static and dynamic features,” *Int. J. Wavelets, Multiresolution Inf. Process.*, vol. 16, no. 02, p. 1840001, Mar. 2018.
- [34] V. Varkarakis, S. Bazrafkan, and P. Corcoran, “A Deep Learning Approach to Segmentation of Distorted Iris Regions in Head-Mounted Displays,” in *2018 IEEE Games, Entertainment, Media Conference, GEM 2018*, 2018, pp. 402–406.
- [35] M. Kassner, W. Patera, and A. Bulling, “Pupil: an open source platform for pervasive eye tracking and mobile gaze-based interaction,” in *Proceedings of the 2014 ACM international joint conference on pervasive and ubiquitous computing: Adjunct publication*, 2014, pp. 1151–1160.
- [36] M. Cimpoi, S. Maji, and A. Vedaldi, “Deep filter banks for texture recognition and segmentation,” in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2015.
- [37] F. Husain, H. Schulz, B. Dellen, C. Torras, and S. Behnke, “Combining Semantic and Geometric Features for Object Class Segmentation of Indoor Scenes,” *IEEE Robot. Autom. Lett.*, 2017.
- [38] S. Bell, P. Upchurch, N. Snaveley, and K. Bala, “Material recognition in the wild with the Materials in Context Database,” in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2015.
- [39] S. Bazrafkan, H. Javidnia, and J. Lemley, “Semiparallel deep neural network hybrid architecture: first application on depth from monocular camera,” *J. Electron. Imaging*, 2018.
- [40] C. Hazirbas, L. Ma, C. Domokos, and D. Cremers, “FuseNet: Incorporating depth into semantic segmentation via fusion-based CNN architecture,” in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2017.
- [41] B. E. Bayer, “Color Imaging Array,” US Pat. 3,971,065, 1976.
- [42] P. Corcoran and P. Bigioi, “Consumer Imaging I -- Processing Pipeline, Focus and Exposure,” in *Handbook of Visual Display Technology*, J. Chen, W. Cranton, and M. Fihn, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014, pp. 1–25.
- [43] P. Amba, D. Alleysson, and M. Mermillod, “Demosaicing using Dual Layer Feedforward Neural Network,” *Color Imaging Conf.*, vol. 2018, no. 1, pp. 211–218, 2019.
- [44] F. Kokkinos and S. Lefkimiatis, “Deep image demosaicking using a cascade of convolutional residual denoising networks,” in *Lecture Notes in Computer Science (Lecture Notes in Artificial Intelligence)*, 2018, vol. 11218 LNCS, pp. 317–333.
- [45] M. Gharbi, G. Chaurasia, S. Paris, and F. Durand, “Deep joint demosaicking and denoising,” *ACM Trans. Graph.*, vol. 35, no. 6, pp. 1–12, 2016.
- [46] W. Dong, M. Yuan, X. Li, and G. Shi, “Joint Demosaicing and Denoising with Perceptual Optimization on a

Generative Adversarial Network,” *arXiv Prepr. 1802.04723*, 2018.

- [47] A. Buades and J. Duran, “Joint Denoising and Demosaicking of Raw Video Sequences,” in *25th IEEE International Conference on Image Processing (ICIP)*, pp. 2172–2176.
- [48] E. Schwartz, R. Giryes, and A. M. Bronstein, “DeepISP: Toward learning an end-to-end image processing pipeline,” *IEEE Trans. Image Process.*, vol. 28, no. 2, pp. 912–923, Jan. 2019.
- [49] A. Tuama, F. Comby, and M. Chaumont, “Camera model identification with the use of deep convolutional neural networks,” in *8th IEEE International Workshop on Information Forensics and Security, WIFS 2016*, 2017, pp. 1–6.
- [50] L. Bondi, L. Baroffio, D. Guera, P. Bestagini, E. J. Delp, and S. Tubaro, “First Steps Toward Camera Model Identification with Convolutional Neural Networks,” *IEEE Signal Process. Lett.*, vol. 24, no. 3, pp. 259–263, 2017.
- [51] C. Chen, Q. Chen, J. Xu, and V. Koltun, “Learning to See in the Dark,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3291–3300.
- [52] V. Varkarakis, S. Bazrafkan, and P. Corcoran, “Deep Neural Network and Data Augmentation Methodology for off-axis iris segmentation in wearable headsets,” *arXiv Prepr. 1903.00389*, 2019.
- [53] J. Lemley, S. Bazrafkan, and P. Corcoran, “Smart Augmentation Learning an Optimal Data Augmentation Strategy,” *IEEE Access*, vol. 5, pp. 5858–5869, 2017.
- [54] S. Bazrafkan, H. Javidnia, and P. Corcoran, “Latent space mapping for generation of object elements with corresponding data annotation,” *Pattern Recognit. Lett.*, no. 116, pp. 179–186, 2018.
- [55] Z. Rui and Z. Yan, “A Survey on Biometric Authentication: Toward Secure and Privacy-Preserving Identification,” *IEEE Access*, vol. 7, pp. 5994–6009, 2019.
- [56] S. Grzonkowski and P. Corcoran, “Sharing cloud services: user authentication for social enhancement of home networking,” *IEEE Trans. Consum. Electron.*, vol. 57, no. 3, pp. 1424–1432, Aug. 2011.