

Enabling ISP-less Low-Power Computer Vision

Gourav Datta, Zeyu Liu, Zihan Yin, Linyu Sun, Akhilesh R. Jaiswal, Peter A. Beerel
University of Southern California, Los Angeles, USA

{gdatta, liuzeyu, zihanyin, linyusun, akhilesh, pabeerel}@usc.edu

Abstract

Current computer vision (CV) systems use an image signal processing (ISP) unit to convert the high resolution raw images captured by image sensors to visually pleasing RGB images. Typically, CV models are trained on these RGB images and have yielded state-of-the-art (SOTA) performance on a wide range of complex vision tasks, such as object detection. In addition, in order to deploy these models on resource-constrained low-power devices, recent works have proposed in-sensor and in-pixel computing approaches that try to partly/fully bypass the ISP and yield significant bandwidth reduction between the image sensor and the CV processing unit by downsampling the activation maps in the initial convolutional neural network (CNN) layers. However, direct inference on the raw images degrades the test accuracy due to the difference in covariance of the raw images captured by the image sensors compared to the ISP-processed images used for training. Moreover, it is difficult to train deep CV models on raw images, because most (if not all) large-scale open-source datasets consist of RGB images. To mitigate this concern, we propose to invert the ISP pipeline, which can convert the RGB images of any dataset to its raw counterparts, and enable model training on raw images. We release the raw version of the COCO dataset, a large-scale benchmark for generic high-level vision tasks. For ISP-less CV systems, training on these raw images result in a $\sim 7.1\%$ increase in test accuracy on the visual wake works (VWW) dataset compared to relying on training with traditional ISP-processed RGB datasets. To further improve the accuracy of ISP-less CV models and to increase the energy and bandwidth benefits obtained by in-sensor/in-pixel computing, we propose an energy-efficient form of analog in-pixel demosaicing that may be coupled with in-pixel CNN computations. When evaluated on raw images captured by real sensors from the PASCALRAW dataset, our approach results in a 8.1% increase in mAP. Lastly, we demonstrate a further 20.5% increase in mAP by using a novel application of few-shot learning with thirty shots each for the novel PASCALRAW dataset, constituting 3 classes.

1. Introduction

Modern high-resolution cameras generate huge amount of visual data arranged in the form of raw Bayer color filter arrays (CFA), also known as a mosaic pattern, as shown in Fig. 1, that need to be processed for downstream CV tasks [43, 1]. An ISP unit, consisting of several pipelined processing stages, is typically used before the CV processing to convert the raw mosaiced images to RGB counterparts [20, 42, 26, 29]. The ISP step that converts these single-channel CFA images to three-channel RGB images is called demosaicing. Historically, ISP has been proven to be extremely effective for computational photography applications, where the goal is to generate images that are aesthetically pleasing to the human eye [29, 8]. However, is it important for high-level CV applications, such as face detection by smart security cameras, where the sensor data is unlikely to be viewed by any human? Existing works [42, 20, 26] show that most ISP steps can be discarded with a small drop in the test accuracy for large-scale image recognition tasks. The removal of the ISP can potentially enable existing in-sensor [31, 10, 2] and in-pixel [5, 27, 12, 13, 14] computing paradigms to process CV computations, such as CNNs partly in the sensor, and reduce the bandwidth and energy incurred in the data transfer between the sensor and the CV system. Moreover, most low-power cameras with a few MPixels resolution, do not have an on-board ISP [3], thereby requiring the ISP to be implemented off-chip, increasing the energy consumption of the total CV system.

Although the ISP removal can facilitate model deployments in resource-constrained edge devices, one key challenge is that most large-scale datasets, that are used to train CV models, are ISP-processed. Since there is a large covariance shift between the raw and RGB images (please see Fig. 1 where we show the histogram of the pixel intensity distributions of RGB and raw images), models trained on ISP-processed RGB images and inferred on raw images, thereby removing the ISP, exhibit a significant drop in the accuracy. One recent work has leveraged trainable flow-based invertible neural networks [44] to convert raw to RGB images and vice-versa using open-source ISP datasets. These networks have recently yielded SOTA test

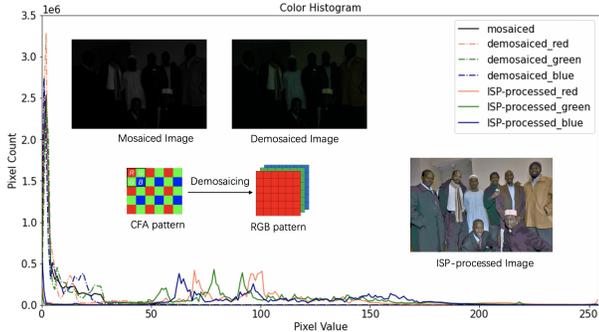


Figure 1. Difference in frequency distributions of pixel intensities between mosaiced raw, demosaiced, and ISP-processed images.

performance in photographic tasks, which we propose to modify to invert the ISP pipeline, and build the raw version of any large-scale ISP processed database for high-level vision applications, such as object detection. This raw dataset can then be used to train CV models that can be efficiently deployed on low-power edge devices without any of the ISP steps, including demosaicing. To further improve the performance of these ISP-less models, we propose a novel hardware-software co-design approach, where a form of demosaicing is applied on the raw mosaiced images inside the pixel array using analog summation during the pixel read-out operation, i.e., without a dedicated ISP unit. Our models trained on this demosaiced version of the visual wake words (VWW) lead to a 8.2% increase in the test accuracy compared to standard training on RGB images and inference on raw images (to simulate the ISP removal and the in-pixel/in-sensor implementation). Even compared to standard RGB training and inference, our models yield 0.7% (1.6%) higher accuracy (mAP) on the VWW (COCO) dataset. Lastly, we propose a novel application of few-shot learning to improve the accuracy of real raw images captured directly by a camera (which has limited number of annotations) with our generated raw images constituting the base dataset.

The key contributions of our paper can be summarized as follows.

- Inspired by the energy and bandwidth benefits obtained by in-sensor computing approaches and the removal of most ISP steps in a CV pipeline, we present and release a large-scale raw image database that can be used to train accurate CV models for low-power ISP-less edge deployments. This dataset is generated by reversing the entire ISP pipeline using the recently proposed flow-based invertible neural networks and custom mosaicing. We demonstrate the utility of this dataset to train ISP-less CV models with raw images.
- To improve the accuracy obtained with raw images, we propose a low-overhead form of in-pixel demosaicing that can be implemented directly on the pixel array alongside other CV computations enabled by recent

paradigms of in-pixel/in-sensor computing approaches and that also reduces the data bandwidth.

- We present a thorough evaluation of our approach with both *simulated* (our released dataset) and *real* (captured by a real camera) raw images, for a diverse range of use-cases with different memory/compute budgets.
- To improve the accuracy of real raw images, we propose a novel application of few-shot learning, with the simulated raw images having a large number of labelled classes constituting the base dataset.

2. Related Works

2.1. ISP Reversal & Removal

Since most ISP steps are irreversible, and depend on the camera manufacturer’s proprietary color profile [6], it is difficult to invert the ISP pipeline. To mitigate this challenge, a few recent works [25, 32, 46] proposed learning-based methods, but they result in large losses and the recovered RAW images may be significantly different from the originals captured by the camera. To reduce this loss, a more recent work [44] used a stack of k invertible and bijective functions $f = f_1 \cdot f_2 \cdot \dots \cdot f_k$ to invert the ISP pipeline. For a raw input x , the RGB output y and the inverted raw input x is computed as $y = f_1 \odot f_2 \odot \dots \odot f_k(x)$ and $x = f_k^{-1} \odot f_{k-1}^{-1} \odot \dots \odot f_1^{-1}(y)$.

The bijective function f_i is implemented through affine coupling layers [44]. In each affine coupling layer, given a D dimensional input m and $d < D$, the output n is

$$n_{1:d} = m_{1:d} + r(m_{d+1:D}) \quad (1)$$

$$n_{d+1:D} = m_{d+1:D} \odot \exp(s(m_{1:d})) + t(m_{1:d}) \quad (2)$$

where s and t represent scale and translation functions from R^d to R^{D-d} that are realized by neural networks, \odot represents the Hadamard product, and r represents an arbitrary function from R^{D-d} to R^d . The inverse step is

$$m_{d+1:D} = (n_{d+1:D} - t(n_{1:d})) \odot \exp(-s(n_{1:d})) \quad (3)$$

$$m_{1:d} = n_{1:d} - r(m_{d+1:D}) \quad (4)$$

The authors then utilize invertible 1×1 convolution, proposed in [23], as the learnable permutation function to revert the channel order for the subsequent affine coupling layer.

Recent works have also investigated the role of the ISP in image classification and the impact of its’ removal/trimming on accuracy for energy and bandwidth benefits. For example, [20] demonstrated that removal of the whole ISP during edge inference results in a $\sim 8.6\%$ loss in accuracy with MobileNets [36] on ImageNet [15], which can mostly be recovered by using just the tone-mapping stage. Another work [42] attempted to integrate the ISP and CV processing using tone mapping and feature-aware downscaling blocks that reduce both the number of bits

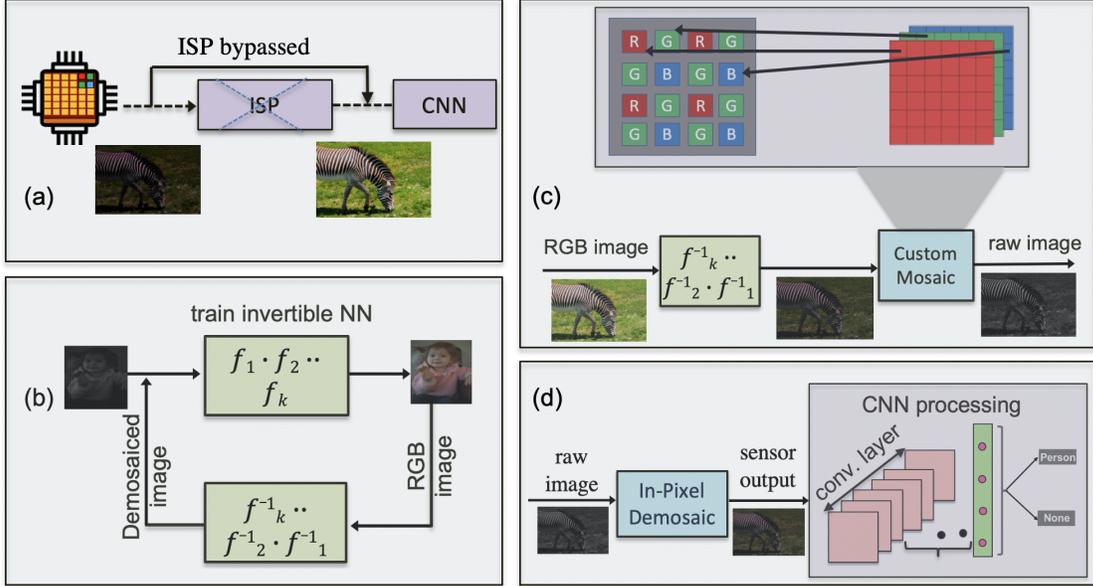


Figure 2. (a) Proposed ISP-less CV system, (b) Invertible NN training on demosaiced raw image, without any white balance or gamma correction, (c) Generation of raw images using the trained inverse network and custom mosaicing, and (d) Application of in-pixel demosaicing and training of the ISP-less CV models. Note the *In-Pixel Demosaic* implementation in the pixel array is illustrated in Fig. 3.

per pixel and the number of pixels per frame. A more recent work [37] used knowledge distillation on an ISP neural network model to align the logit predictions of an off-the-shelf pretrained model for raw images with that for ISP-processed RGB images.

2.2. Few-Shot Object Detection

In recent years, few-shot object detection (FSOD) has gained significant traction as ML accuracy in low-data scenarios continues to improve. There are two mainstream training paradigms in FSOD, meta-learning and finetune-based methods. Meta-learning methods attempt to capture aggregated information from multiple annotated data-rich support datasets. Thus, when required to train on a dataset with novel classes and less data, the model can leverage the prior knowledge learned from support datasets to generalize to new classes. For example, [22] used a re-weighting module to adjust coefficients of the query image meta features by capturing global features of the support images to be suitable for novel object detection. Authors in [45] proposed a Predictor-head Remodeling Network (PRN) module to generate class-attentive vectors to provide aggregated features between support and query images for the meta-learner predictor head. Additionally, [17] introduced an attention-based region proposal network to match the candidate proposals with the support images and a multi-relation detector which can measure the similarity between proposal boxes from the query and the support objects. Compared to meta-learning which requires a complicated training process, the finetune-based methods have

a simpler pipeline. For instance, [41] proposed the two-stage fine-tune based approach (TFA) which only finetunes the bounding box classification and regression parts on a class-balanced training set, but outperforms many meta-learning methods. Moreover, to mitigate misclassifying novel instances as confusing base classes, [39] introduced contrastive learning into the FSOD pipeline, that helps the learned target features represent high intra-class similarity and inter-class variability.

3. Inverting ISP Pipeline

Similar to [44], we propose to generate the raw demosaiced images from ISP-processed RGB images using the affine coupling layers described in Section 2.1. However, [44] models the ISP pipeline on the demosaiced, white-balanced and gamma-corrected raw image, and hence, the invertible ISP pipeline does not generate the raw image that are captured directly by a camera. The authors apply gamma correction on the RAW data (i.e. without storing on disk) to compress the dynamic range for faster convergence. Hence, for ISP-less in-sensor CV systems, the naive application of the invertible ISP pipeline proposed in [44] will require performing these operations in the sensor. This is challenging due to the limited compute/memory footprint available in the pixel array and the periphery. In particular, traditional demosaicing involves matrix operations that involve interpolation (nearest neighbour, bi-linear, bi-cubic, etc.) techniques which scale with the input resolution. Moreover, white balancing involves a variable gain amplification for each pixel location which requires complex control logic, and gamma correction involves logarithm-

mic computation which is challenging to process using analog logic in advanced high-density pixels.

For these reasons, we propose to train an invertible network on the demosaiced images from the MIT-Adobe 5K dataset [7]. Despite our focus on classification/detection tasks, we propose to use this photographic dataset to train the invertible ISP because we do not have large-scale ground truth raw-RGB image pairs for those tasks. We train using demosaiced images because the input size of the invertible neural network must be equal to its output size. Once trained, we use this network to obtain the raw demosaiced images from the ISP-processed RGB images from the large-scale classification/detection datasets.

We then invert the demosaicing i.e., perform the mosaicing operation by selecting the appropriate pixel color corresponding to each location, as shown in Fig. 2. For example, to generate the red pixel in a particular mosaiced RGGB patch, we select the pixel intensity of the red channel in the same location as in the demosaiced image. Although this final mosaiced image is obtained after inverting the entire ISP pipeline, it might still be slightly different from the raw image captured by a camera. This is partially because we do not explicitly model the latent distribution of the different ISP steps to stabilize the training of the invertible network. We mitigate this concern using few-shot learning.

4. Requirement for Demosaicing

Although training on raw images in the Bayer CFA format increases the test accuracy of ISP-less CV applications, it might lack the representation capacity that multiple colors spanning different spectral bands might provide for each pixel location. Hence, a natural question is can we increase this capacity without an additional ISP unit? Since demosaicing is the ISP technique that yields separate RGB channels from the raw CFA format, one intuitive idea is to implement demosaicing directly in the pixel array, and then process the CV computations required for CNNs using in-pixel/in-sensor computing. However, as explained above, traditional demosaicing approaches involve complex operations which are hard to map on the pixel array, especially when the pixel array needs to process the initial CNN layers in the in-pixel computing paradigms. Hence, we propose a low-overhead custom in-pixel demosaicing approach that significantly increases the test accuracy for our benchmarks compared to inference on raw images.

5. Proposed Demosaicing Technique

We propose to implement a simple but effective custom demosaicing operation inside the analog pixel array. Let us consider a demosaiced RGB image with shape $X \times Y \times 3$, to be processed for CV applications. Then, our custom demosaicing technique requires the input mosaiced raw image to have a shape $2X \times 2Y$, such that each 2×2 RGGB patch produces the corresponding 3 channels for a single pixel,

thereby yielding a 25% reduction in data dimensionality. Functionally, the custom demosaicing copies the red and blue pixel intensities from the camera to the demosaiced RGB channel output, while the two green pixels from the RGGB patch of the camera pixel array are *averaged* to produce one effective value for green pixel intensity. While the summation is performed by analog computation inside the pixel array described below, the division is performed in the digital domain after the analog to digital converter (ADC) in the pixel periphery by a simple logical right shift operation.

The proposed implementation of the pixel array to accomplish this custom demosaicing functionality is shown in Fig. 3(a). We propose to include two select lines for each row of the pixel array - the first set of select lines called 'Row-Select' are connected to the select transistors of the red and blue pixels, while the second set of select lines called 'Green-Select' are connected only to the green pixels. Essentially, the pixels in RGGB Bayer pattern are connected in an interleaved manner to the two select lines. Therefore, the read-out of the red and blue pixels are controlled by the 'Row-Select' lines, while the 'Green-Select' lines control the read-out of the green pixels. Now consider activating two rows of 'Row-Select' lines (Row-Select-1 and Row-Select-2) in the 2×2 pixel array of Fig. 3(a). This would result in read-out of the red and blue pixels on 'Column-Line-1' and 'Column-Line-2', respectively. The two green pixels would remain deactivated as the 'Green-Select' lines are kept at low voltage. In a subsequent cycle, the two 'Row-Select' lines are kept at low voltage and the two 'Green-Select' lines are activated by pulling them to high voltage. Further, the two 'Column-Lines' are connected together by closing the 'Column-Switch', shown in Fig. 3(a). Consequently, the voltage on the now connected 'Column-Lines' represents the accumulated response of the two green pixels, which are fed to column-ADCs for analog to digital conversion. Note, the proposed scheme is similar to pixel binning approaches [4, 40, 38, 21], except that in this case binning is selectively performed only for the two green pixels in each patch of Bayer RGGB pattern using interleaved connections to 'Row-Select' and 'Green-Select' lines. In summary, in two cycles, wherein two rows of 'Row-Select' and 'Green-Select' lines are activated in each cycle, the proposed scheme can generate demosaiced red, blue, and green pixels. Note, since we are able to read two rows (consisting of RGGB pixels) in two cycles, the proposed scheme does not incur any overhead in terms of read-out speed (or frame-rate) of the camera.

In yet another approach, we propose to combine the custom demosaicing and the computations of the first-layer of CNN inside pixel array using the P²M (Processing-in-pixel-in-memory) paradigm proposed in [12], as shown in Fig. 3(b). Modifying the P²M pixel array of [12], Fig. 3(b) presents a novel pixel array that can combine demosaicing and convolution computations using memory-embedded pixels. Essentially, the CNN weights associated

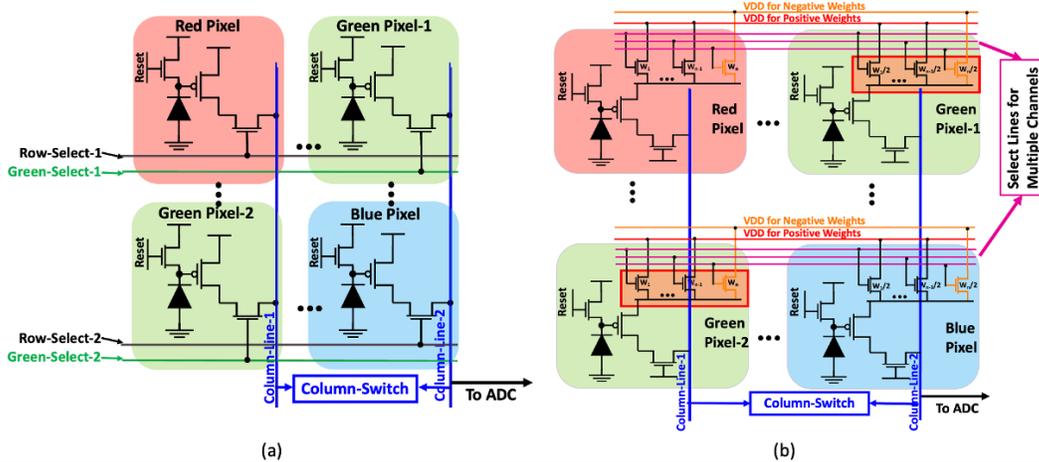


Figure 3. Implementation of the proposed (a) demosaicing and (b) demosaicing coupled with in-pixel convolution for ISP-less CV.

with two green pixels in a single patch of Bayer RGB pattern are kept the same. This is achieved by keeping the sizes of the weight transistors the same across the two green pixels. Further, these transistor weights are kept such that each set of weight transistors for a single pixel has half of the effective algorithmic weight value associated with the green channel in the input convolutional layer. This ensures that the resultant analog dot product obtained from the P^2M scheme [12] involve effective averaging of the intensities of the green pixels and then multiplying it with the corresponding weights associated with the convolutional layer.

While the in-pixel convolution on the demosaiced image can lead to significantly higher bandwidth reduction [12] (quantified later in Section 8.5), the analog non-idealities involved in the multiply-and-accumulate operation and weight mismatch in the green pixels can lead to large errors, require re-training the entire CNN network, and introduce manufacturing challenges, which might require non-trivial changes to the design pipeline of sensors.

6. Few-Shot Learning

Compared to the abundance of RGB image datasets, it is difficult to obtain large-scale annotated raw images. For example, to the best of our knowledge, the only raw image database for classification/detection tasks, PASCALRAW, contains only 4,259 annotated images, with 3 object classes, which are not enough to train a deep CV model. Even with a pre-trained model on a large-scale RGB dataset, it might be difficult to fine-tune on this small-scale raw dataset (due to the co-variance shift) and yield satisfactory performance.

As described in Section 2.2, recent works have proposed a plethora of few-shot learning approaches that achieve great performance on datasets with some novel classes, and a few images per class. Our problem is not exactly same as a typical few-shot learning setup, given we can find a large-scale annotated RGB image database, having the same classes as the raw dataset. For example, the Microsoft COCO dataset consists of 80 classes, which can cover ob-

jects from a range of applications such as autonomous driving, aerial imagery recognition, and can be used for fine-tuning on a raw image database with a subset of their classes. We propose to leverage TFA [41] (see Section 2.2) for this fine-tuning process, and to the best of our knowledge, this is the first application of few-shot learning in improving the accuracy/mAP of raw images.

However, naively applying TFA with COCO as the base dataset can only bring limited improvement in accuracy, due to the co-variance shift between the RGB and raw images. Note that in typical few-shot learning setups, such as TFA, the images in the base and novel datasets are assumed to have similar intensity distributions [41]. Hence, we propose a novel application of few-shot learning which leverages our simulated raw COCO dataset as the base class to increase the test mAP on the real raw dataset. We choose a class-balanced subset of the real raw dataset as the samples with ‘novel class’ to perform TFA on the model pretrained on our raw COCO dataset to further improve the mAP.

7. Experimental Setup

7.1. Implementation Details

We evaluate our proposed method on three CNN backbones/frameworks with varying complexities and use-cases as described below. For object detection experiments, we use the [9] framework, while for few-shot learning, we use the *mmfewshot* [28] and *FsDet* [41] framework. Our training details are provided in the supplementary materials. *MobileNetV2* [36]: A lightweight depthwise convolution neural network that has gained significant traction for being deployed on resource-constrained edge devices, such as mobile devices. In this work, we use a lower complexity version of MobileNetV2, namely MobileNetV2-0.35x [35], which shrinks the output channel count by $0.35\times$ to satisfy the compute budget of 60M floating point operations (FLOPs) representing standard micro-controllers, where ISP-less CV may be the most relevant.

Faster R-CNN [34]: A two-stage object detection framework that consists of a feature extraction, a region proposal, and a RoI pooling module. For our experiments, we use ResNet101 as the backbone network for feature extraction, since MobileNetV2 significantly degrades the test mAP compared to SOTA.

YOLOv3 [33]: YOLOv3 (You Only Look Once, Version 3) is a real-time object detection framework that identifies specific objects in videos or images. We use MobileNetV2 as the backbone network for feature extraction in YOLOv3.

7.2. Dataset Details

We evaluate our proposed approaches on the simulated raw versions of the Visual Wake Words (VWW) and COCO datasets, and a real raw dataset captured by a real camera introduced in [16]. The details of the datasets are below.

VWW [11]: The Visual Wake Words (VWW) dataset consists of high resolution images that include visual cues to “wake-up” AI-powered resource-constrained home assistant devices that require real-time inference. The goal of the VWW challenge is to detect the presence of a human in the frame (a binary classification task with 2 labels) with very little resources-close to 250KB peak RAM usage and model size, which is only satisfied by MobileNetV2-0.35x, and hence, used in our experiments.

Microsoft COCO: To evaluate on the multi-object detection task, we use the popular Microsoft COCO dataset [24]. Specifically, we use an image resolution of 1333×800 for the Faster-RCNN framework, and 416×416 for the YoloV3 framework [33], the same as used in [33]. We use the 80 available classes used for our experiments. We evaluate the performance of each method using mAP averaged for IoU $\in \{0.5, 0.75, [0.5 : 0.05 : 0.95]\}$, denoted as mAP@0.5, mAP@0.75 and mAP@[0.5, 0.95], respectively. Note that we also report the individual mAPs for small (area <32² pixels), medium (area between 32² and 96² pixels), and large (area >96² pixels) objects.

PASCALRAW: This RAW image database was developed to simulate the effect of algorithmic hardware implementations such as embedded feature extraction at the image sensor or readout level on end-to-end object detection performance. The annotations of this dataset were made in accordance with the original PASCAL VOC guidelines [16]. For the few-shot learning experiments, we choose 29 images containing the class ‘bicycle’, 25 images containing the class ‘car’ and 21 images containing the class ‘person’ to construct a balanced training set where each class has 30 annotated objects (i.e., 30-shot), and use the remaining 4178 images as the test dataset.

8. Experimental Results

8.1. VWW Results

For VWW, we compare the accuracy of the tinyML-based MobileNetV2-0.35x model with our proposed de-

Table 1. Evaluation of our approach on ISP-less CV systems with MobileNetV2-0.35x on VWW dataset. Demosaiced¹ denotes traditional demosaicing, while demosaiced² denotes our in-pixel demosaicing. WB, GC, and IPC denotes white balance, gamma correction, and in-pixel computing. Also, note that models trained on mosaiced images can only be tested with mosaiced images.

Method	Test Acc. (%)			
	Inference	Mosaiced	demosaiced ²	IPC
Training				
Mosaiced		87.47	-	-
demosaiced ¹		-	88.84	88.04
demosaiced ²		-	89.92	89.07
demosaiced ¹ +WB		-	86.47	86.23
demosaiced ¹ +WB+GC		-	82.70	81.45
ISP-processed		-	81.97	81.43

Table 2. mAP on different versions of the COCO raw dataset to emulate ISP-less CV systems using a Faster R-CNN framework with ResNet101 backbone.

model	mean average precision					
	0.5:0.95	0.5	0.75	S	M	L
baseline	33.8	50.5	37.0	16.6	36.6	46.7
demosaiced ¹	42.8	64.1	47.1	25.6	46.9	55.0
mosaiced	29.4	45.7	31.8	12.7	32.1	42.9
demosaiced ²	37.8	57.7	39.8	20.2	48.6	53.2

¹ ‘baseline’ indicates testing on our proposed COCO raw dataset with model pretrained on ISP-processed COCO dataset

² ‘demosaiced¹’ indicates training and testing on our proposed COCO raw dataset

³ ‘mosaiced’ indicates training and testing on mosaiced images obtained from the COCO raw dataset from our invertible ISP

⁴ ‘demosaiced²’ means training and testing on our in-pixel demosaiced images

mosaicing and in-pixel computing technique against inference on mosaiced raw and RGB images in Table 1. We also compare our approach with traditional demosaicing (*opencv* library in Python), white balancing (*rawpy* package in Python), and gamma correction. Note that, as shown in Table 1, using identically distributed images during training yields the best accuracy during inference.

Table 1 further illustrates that using the off-the-shelf model pre-trained on ISP-processed images yields an accuracy of 81.97% when deployed on an ISP-less CV system with our in-pixel demosaicing, which is 7.32% lower compared to the ISP-processed inference. Note that we cannot avoid the demosaicing step, because the pre-trained model is trained with 3 channel input images. With the generated mosaiced image database from our invertible pipeline, the accuracy gap (training and testing both on the mosaiced image) reduces to 2.82%. Additionally, with our in-pixel demosaicing on this mosaiced image, we yield an accuracy of 89.92%, which is even 0.63% higher than the RGB test accuracy. Appending the first layer convolution inside the pixel, coupled with the demosaicing results in a little lower accuracy of 89.07%.

Table 3. Comparison of our proposed approach on PASCALRAW dataset.

Framework	Method	mAP					
		@[0.5,0.95]	@0.5	@0.75	@small	@medium	@large
Yolov3	ISP-processed	2.7	8.2	1.2	0.2	2.2	4.4
	ISP-processed+few-shot*	5.2	15.4	2.4	0.6	5.5	7.9
	ISP-processed+few-shot**	6.2	17.0	3.3	0.2	3.9	11.2
	demosaiiced raw	13.4	38.5	5.4	0.9	12.2	22.9
	demosaiiced raw+few-shot*	16.9	40.6	10.9	0.5	17.8	26.3
	demosaiiced raw+few-shot**	20.8	47.4	14.5	0.9	17.3	30.4
Faster RCNN	ISP-processed	1.2	4.2	0.2	0.0	1.3	3.5
	ISP-processed+few-shot*	5.9	14.8	3.3	0.0	3.8	8.6
	ISP-processed+few-shot**	9.5	26.0	4.2	0.0	6.6	15.0
	demosaiiced raw	9.3	29.9	2.2	1.7	10.5	19.5
	demosaiiced raw+few-shot*	27.4	52.8	25.7	6.9	27.1	37.3
	demosaiiced raw+few-shot**	29.8	58.1	28.0	8.0	28.1	40.6

* The experiments apply few-shot learning with 30 shots of both base classes and novel classes.

** The experiments apply few-shot learning with 30 shots of only the novel classes.

8.2. COCO raw Results

The detailed results on COCO raw dataset are summarized in Table 2. Our experiments indicate that direct inference on the COCO demosaiced raw dataset using the model pre-trained on COCO ISP-processed RGB dataset yields an mAP of 33.8%, which is 7.2% lower compared to ISP-processed inference. Note that the mAP of small objects reduces significantly by nearly 35%. However, with finetuning on our COCO demosaiced raw dataset, the mAP increases to 42.8%. Unlike in VWW, where models can be accurately trained from scratch, training and testing on the COCO mosaiced raw image leads to a reduced mAP of 29.4%. This reduction might be because the pre-trained model (where the backbone is also pre-trained on ImageNet) cannot be leveraged because of the difference in the number of input channels. Lastly, applying our proposed in-pixel demosaicing on the mosaiced raw dataset yields an mAP of 37.0%, which is 5.0% lower than ISP-processed inference, unlike VWW. This might be because our demosaicing reduces the spatial resolution of the image, which might be detrimental for the complex object detection task. Interestingly, our approach is effective in detecting medium sized objects, and achieves the highest mAP of 48.6%.

8.3. PASCALRAW Results

8.3.1 YOLOv3

Table 3 shows the performance of six different methods with YOLOv3 on the PASCALRAW dataset. Direct inference on this dataset with models pre-trained on ISP-processed COCO dataset yields only 2.7% mAP due to the significant co-variance shift between the two datasets. Using the ISP-processed base dataset, we compare two different few-shot learning approaches, one where we use 30 shots for both the base and novel classes, and the other where we use 30 shots only for the novel classes. We observe the latter leads to 1.0% higher mAP compared to the

former, which might be because the former may underfit to the three target classes due to its improved generalization. Note, due to the difference in the dataset distributions, few-shot learning fails to significantly increase the mAP, as observed from the modest mAP improvement from 2.7% to 6.2%. On the other hand, after finetuning on our custom demosaiced COCO raw dataset (without any few-shot learning), the mAP increases by more than 4 \times to 13.4%. This strongly demonstrates the effectiveness of our large-scale raw database. Lastly, applying few-shot learning with this base raw dataset further increases the mAP to 20.8%.

8.3.2 Faster R-CNN

We perform a series of similar experiments with Faster R-CNN model with ResNet101 backbone on the PASCALRAW dataset. As we can see in Table 3, the results are consistent with that from the YOLOv3 model, except that there is no mAP increase with fine-tuning on the COCO raw dataset compared to applying few-shot learning with the ISP-processed base dataset. This might be because our demosaicing approach that incurs 4 \times spatial down-sampling might not be that competitive compared to the faster R-CNN framework with ultra-high input resolution. Applying few-shot learning with 30 shots of only the novel classes on our custom demosaiced COCO raw dataset yields an mAP of 29.8%, which is 28.6% higher compared to directly using the model pretrained on the ISP-processed COCO dataset.

8.4. Comparison with Prior Works

We compare the test accuracy and mAP obtained by our ISP-less CV models with existing similar works on the VWW and COCO dataset respectively in Fig. 4(a-b). As we can see, we yield similar performance on average compared to using the invertible ISP pipeline proposed in [44], while providing bandwidth and energy reduction quantified in Section 8.5. Even compared to testing on ISP-processed RGB images which require the entire ISP pipeline, we

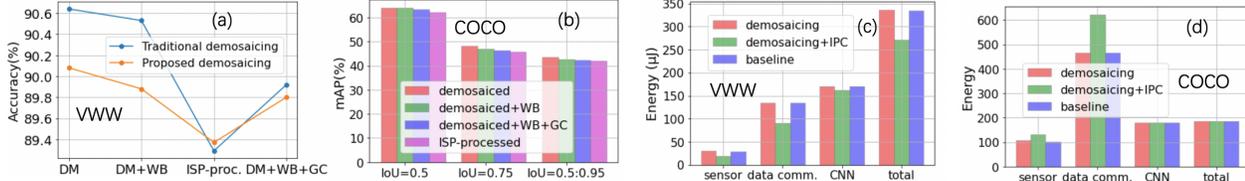


Figure 4. Comparison of the (a) accuracy and (b) mAP of our proposed demosaicing method with different ISP pipelines on COCO dataset with Faster-RCNN framework with ResNet101 backbone and VWW dataset with MobileNetV2-0.35x respectively, where DM denotes our proposed demosaicing technique, WB and GC denote white balancing and gamma correction respectively. The energy consumptions of our approaches are compared with the normal pixel read-out in (c) and (d) on VWW and COCO respectively, where IPC denotes in-pixel computing. Note, for (d), the energy unit is μJ for ‘sensor’ & ‘data comm.’, and $100\mu J$ for ‘CNN’ & ‘total’.

obtain 0.63% (1.6%) increase in accuracy (mAP) on the VWW (COCO) dataset. It is difficult to directly compare our approach with other works [20, 42], as they do not release the ISP model, and evaluate the impact of the removal of different ISP stages on in-the-wild datasets, such as ImageNet [15] and KITTI [18], which may not be a relevant use-case of ISP-less low-power edge deployment.

8.5. Bandwidth & Energy Benefits

Removing the entire ISP pipeline, and applying the proposed in-pixel demosaicing operation directly on the raw images can lead to significant energy and bandwidth savings, thereby aiding the deployment of CNN models on ultra low-power edge devices. The complete image captured by the sensors is transmitted to a down-stream SoC processing the ISP and CV units typically through energy-hungry MIPI interfaces, which cost significant bandwidth [19]. As explained in Section 5, the demosaicing operation leads to a dimensionality reduction of $\frac{4}{3}$, which implies a 25% reduction in bandwidth. Quantizing the demosaiced outputs to 8-bits using custom ADCs (inputs to modern CNNs have unsigned 8-bit representation) leads to a $(\frac{12}{8})$ or 50% reduction in bandwidth, assuming the raw image has a bit-depth of 12 [30]. Lastly, appending the first convolution layer inside the sensor yields a $3\times$ increase in bandwidth for MobileNetV2-0.35x. This is convolutional layer has a stride of 2, which implies a $4\times$ dimensionality reduction, while there is a $(\frac{8}{3})$ dimensionality increase due to the 3 channels in the input demosaiced image and 8 output channels in the first convolutional layer. In summary, the total bandwidth/data transmission energy reduction due to our proposed demosaicing operation is 75%, while for the in-pixel computing approach (on the proposed demosaiced image as illustrated in Section 5) is $12\times$.

Note that this energy benefit is in addition to the energy savings obtained by removing the ISP operations in an SoC, and transferring the ISP output to a CV processing unit. It is difficult to accurately quantify this saving as it depends on the underlying hardware implementation and dataflow, as well as the proprietary implementation of ISP. That said, we compare the sensor (pixel+ADC), data communication, and the CNN energy consumption of our demosaicing and in-pixel computing approaches with normal

pixel read-out in Fig. 4(c-d). While Fig. 4(c) represents the tinyML use-case on VWW using MobileNetV2-0.35x, Fig. 4(d) represents the more difficult use-case on COCO using YoloV3. We compute the pixel energies using our in-house circuit simulation framework, while the ADC, data communication, and CNN energies are obtained from [12]. While our demosaicing approach incurs a sensor energy overhead of $\sim 5\%$ on average, the proposed in-pixel implementation reduces (increases) the sensor energy by 33% (23%) on VWW (COCO) with MobileNetV2-0.35x (YoloV3). The energy increase is due to the increased number of convolutional output channels (first layer) in the MobileNet backbone of YoloV3.

9. Discussions

In this work, we propose an ISP-less computer vision paradigm to enable the deployment of CNN models on low-power edge devices that involve processing close to the sensor nodes with limited compute/memory footprint. Our proposal has two significant benefits: 1) We release a large-scale RAW image database that can be used to train and deploy CNNs for a wide range of vision tasks (including those related to photography) and 2) Our hardware-software co-design approach leads to significant bandwidth savings compared to traditional CV pipelines. To the best of our knowledge, this is the first work to address the widely overlooked ISP pipeline in near-sensor and in-sensor processing paradigms while also proposing novel in-pixel schemes for custom demosaicing, coupled with convolution computation. Our proposed approach increases the test accuracy (mAP) of a tinyML (generic object detection) application by 7.32% (7.2%) compared to direct deployment of the off-the-shelf pre-trained models on ISP-less CV systems. Our approach, coupled with few-shot learning, has been shown to be effective in detecting real raw objects captured directly by a camera from the PASCALRAW dataset.

10. Acknowledgements

We would like to acknowledge the DARPA HR00112190120 award and the NSF CCF-1763747 award for supporting this work. The views and conclusions contained herein are those of the authors and should not be

interpreted as necessarily representing the official policies or endorsements of DARPA or NSF.

References

- [1] Scaling CMOS Image Sensors. <https://semiengineering.com/scaling-cmos-image-sensors/>, 2020. Accessed: 04-20-2020.
- [2] Sony to Release World’s First Intelligent Vision Sensors with AI Processing Functionality. <https://www.sony.com/en/SonyInfo/News/Press/202005/20-037E/>, 2020. Accessed: 12-01-2022.
- [3] AP0201AT: Image Signal Processor, 2 MP. <https://www.onsemi.com/products/sensors/image-signal-processors-isps/ap0201at>, 2021.
- [4] Nikolai E Bock. Methods for pixel binning in an image sensor, July 22 2008. US Patent 7,402,789.
- [5] Laurie Bose, Piotr Dudek, Jianing Chen, Stephen J. Carey, and Walterio W. Mayol-Cuevas. Fully embedding fast convolutional networks on pixel processor arrays. In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXIX*, volume 12374, pages 488–503. Springer, 2020.
- [6] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T. Barron. Unprocessing images for learned raw denoising, 2018.
- [7] Vladimir Bychkovsky, Sylvain Paris, Eric Chan, and Fredo Durand. Learning photographic global tonal adjustment with a database of input / output image pairs. In *CVPR 2011*, volume 1, pages 97–104, 2011.
- [8] Prashant Chaudhari, Franziska Schirmacher, Andreas Maier, Christian Riess, and Thomas Köhler. Merging-ISP: Multi-exposure high dynamic range image signal processing. In Christian Bauckhage, Juergen Gall, and Alexander Schwing, editors, *Pattern Recognition*, pages 328–342, Cham, 2021. Springer International Publishing.
- [9] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open MMLab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- [10] Zhe Chen, Huifeng Zhu, Erxiang Ren, Zheyu Liu, Kaige Jia, Li Luo, Xuan Zhang, Qi Wei, Fei Qiao, Xinjun Liu, and Huazhong Yang. Processing near sensor architecture in mixed-signal domain with CMOS image sensor of convolutional-kernel-readout method. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 67(2):389–400, 2020.
- [11] Aakanksha Chowdhery, Pete Warden, Jonathon Shlens, Andrew Howard, and Rocky Rhodes. Visual wake words dataset. *arXiv preprint arXiv:1906.05721*, 2019.
- [12] Gourav Datta et al. P2M: A processing-in-pixel-in-memory paradigm for resource-constrained TinyML applications. *arXiv preprint arXiv:2203.04737*, 2022.
- [13] Gourav Datta et al. Toward efficient hyperspectral image processing inside camera pixels. *arXiv preprint arXiv:2203.05696*, 2022.
- [14] Gourav Datta, Souvik Kundu, Zihan Yin, Joe Mathai, Zeyu Liu, Zixu Wang, Mulin Tian, Shunlin Lu, Ravi T. Lakkireddy, Andrew Schmidt, Wael Abd-Almageed, Ajey P. Jacob, Akhilesh R. Jaiswal, and Peter A. Beerel. P2M-DeTrack: Processing-in-pixel-in-memory for energy-efficient and real-time multi-object detection and tracking, 2022.
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, volume 1, pages 248–255, 2009.
- [16] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>, 2012.
- [17] Qi Fan, Wei Zhuo, Chi-Keung Tang, and Yu-Wing Tai. Few-shot object detection with attention-RPN and multi-relation detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4013–4022, 2020.
- [18] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, volume 1, pages 3354–3361, 2012.
- [19] Jorge Gomez, Saavan Patel, Syed Shakib Sarwar, Ziyun Li, Raffaele Capocchia, Zhao Wang, Reid Pinkham, Andrew Berkovich, Tsung-Hsun Tsai, Barbara De Salvo, and Chiao Liu. Distributed on-sensor compute system for AR/VR devices: A semi-analytical simulation framework for power estimation. *arXiv preprint arXiv:2203.07474*, 2022.
- [20] Patrick Hansen, Alexey Vilkin, Yury Khrustalev, James Imber, David Hanwell, Matthew Mattina, and Paul N. Whatmough. ISP4ML: Understanding the role of image signal processing in efficient deep learning vision systems. *arXiv preprint arXiv:1911.07954*, 2019.
- [21] Xiaodan Jin and Keigo Hirakawa. Analysis and processing of pixel binning for color image sensor. *EURASIP Journal on Applied Signal Processing*, 2012(1), Dec. 2012.
- [22] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell. Few-shot object detection via feature reweighting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [23] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [25] Yu-Lun Liu, Wei-Sheng Lai, Yu-Sheng Chen, Yi-Lung Kao, Ming-Hsuan Yang, Yung-Yu Chuang, and Jia-Bin Huang. Single-image HDR reconstruction by learning to reverse the

- camera pipeline. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [26] Ekdeep Singh Lubana, Robert P. Dick, Vinayak Aggarwal, and Pyari Mohan Pradhan. Minimalistic image signal processing for deep learning applications. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 4165–4169, 2019.
- [27] Lukas Mennel, Joanna Karolina Symonowicz, Stefan Wachter, Dmitry K. Polyushkin, Aday J Molina-Mendoza, and T. Mueller. Ultrafast machine vision with 2D material neural network image sensors. *Nature*, 579:62–66, 2020.
- [28] mmfewshot Contributors. OpenMMLab few shot learning toolbox and benchmark. <https://github.com/open-mmlab/mmfewshot>, 2021.
- [29] Ali Mosleh, Avinash Sharma, Emmanuel Onzon, Fahim Mannan, Nicolas Robidoux, and Felix Heide. Hardware-in-the-loop end-to-end optimization of camera image processing pipelines. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [30] ON Semiconductor. *CMOS Image Sensor, 1.2 MP, Global Shutter*, 3 220. Rev. 10.
- [31] Reid Pinkham, Andrew Berkovich, and Zhengya Zhang. Near-sensor distributed DNN processing for augmented and virtual reality. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 11(4):663–676, 2021.
- [32] Abhijith Punnappurath and Michael S. Brown. Learning raw image reconstruction-aware deep image compressors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4):1013–1019, 2020.
- [33] Joseph Redmon and Ali Farhadi. YOLOv3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [34] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015.
- [35] Oindrila Saha, Aditya Kusupati, Harsha Vardhan Simhadri, Manik Varma, and Prateek Jain. RNNPool: Efficient non-linear pooling for RAM constrained inference. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 20473–20484. Curran Associates, Inc., 2020.
- [36] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.
- [37] Eli Schwartz, Alex Bronstein, and Raja Giryes. ISP distillation. *arXiv preprint arXiv:2101.10203*, 2021.
- [38] SMA Sharif, Rizwan Ali Naqvi, and Mithun Biswas. Beyond joint demosaicking and denoising: An image processing pipeline for a pixel-bin image sensor. *arXiv preprint arXiv:2104.09398*, 2021.
- [39] Bo Sun, Banghuai Li, Shengcai Cai, Ye Yuan, and Chi Zhang. FSCE: Few-shot object detection via contrastive proposal encoding. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2021.
- [40] Samet Taspinar, Manoranjan Mohanty, and Nasir Memon. Effect of video pixel-binning on source attribution of mixed media. In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, volume 1, pages 2545–2549, 2021.
- [41] Xin Wang, Thomas E. Huang, Trevor Darrell, Joseph E Gonzalez, and Fisher Yu. Frustratingly simple few-shot object detection. July 2020.
- [42] Chyuan-Tyng Wu, Leo F. Isikdogan, Sushma Rao, Bhavin Nayak, Timo Gerasimow, Aleksandar Sutic, Liron Ainkedem, and Gilad Michael. VisionISP: Repurposing the image signal processor for computer vision applications. In *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, sep 2019.
- [43] Jiyang Xie, Yixiao Zheng, Ruoyi Du, Weiyu Xiong, Yufei Cao, Zhanyu Ma, Dongpu Cao, and Jun Guo. Deep learning-based computer vision for surveillance in ITS: Evaluation of state-of-the-art methods. *IEEE Transactions on Vehicular Technology*, 70(4):3027–3042, 2021.
- [44] Yazhou Xing, Zian Qian, and Qifeng Chen. Invertible image signal processing. In *CVPR*, 2021.
- [45] Xiaopeng Yan, Ziliang Chen, Anni Xu, Xiaoxi Wang, Xiaodan Liang, and Liang Lin. Meta R-CNN: Towards general solver for instance-level low-shot learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9577–9586, 2019.
- [46] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. CycleISP: Real image restoration via improved data synthesis. In *CVPR*, 2020.