# Formulating Camera-Adaptive Color Constancy as a Few-shot Meta-Learning Problem

Steven McDonagh[*], Sarah Parisot[*], Fengwei Zhou, Xing Zhang,
Ales Leonardis, Zhenguo Li, Gregory Slabaugh
Huawei Noah's Ark Lab

{steven.mcdonagh, sarah.parisot, zhoufengwei, zhang.xing1
ales.leonardis, li.zhenguo, greg.slabaugh}
@huawei.com

## Abstract

*Digital camera pipelines employ color constancy methods to estimate an unknown scene illuminant, in order to re-illuminate images as if they were acquired under an achromatic light source. Fully-supervised learning approaches exhibit state-of-the-art estimation accuracy with camera-specific labelled training imagery. Resulting models typically suffer from domain gaps and fail to generalise across imaging devices. In this work, we propose a new approach that affords fast adaptation to previously unseen cameras, and robustness to changes in capture device by leveraging annotated samples across different cameras and datasets. We present a general approach that utilizes the concept of color temperature to frame color constancy as a set of distinct, homogeneous few-shot regression tasks, each associated with an intuitive physical meaning. We integrate this novel formulation within a meta-learning framework, enabling fast generalisation to previously unseen cameras using only handfuls of camera specific training samples. Consequently, the time spent for data collection and annotation substantially diminishes in practice whenever a new sensor is used. To quantify this gain, we evaluate our pipeline on three publicly available datasets comprising 12 different cameras and diverse scene content. Our approach delivers competitive results both qualitatively and quantitatively while requiring a small fraction of the camera-specific samples compared to standard approaches.*

## 1. Introduction

The colors of an image captured by a digital camera are always affected by the prevailing light source color in the scene. Accounting for the effect of scene illuminant to produce images of canonical appearance (as if captured under an achromatic light source) is an essential component of digital photography pipelines, and is of great importance for many practical high-level computer vision applications including image classification, semantic segmentation and

---

* Authors contributed equally



(a) Input image
(b) Ground-truth
(c) Standard fine tuning of a pre-trained model with access to only 10 images from a test camera.
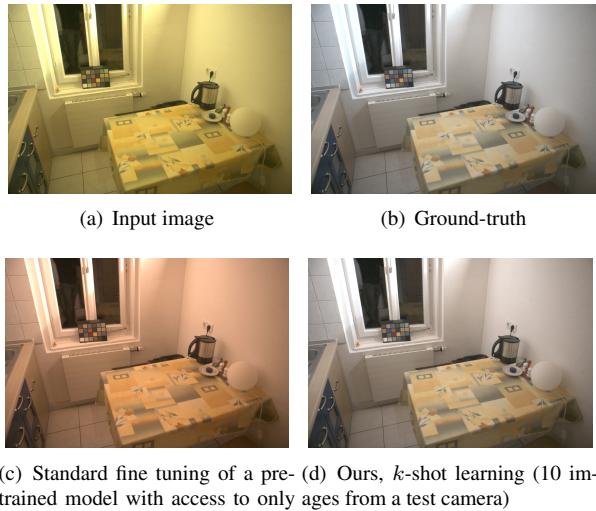(d) Ours, $k$-shot learning (10 images from a test camera)

Figure 1: An example image before and after color constancy correction. Our approach can quickly adapt to new unseen camera sensors using few samples where a pre-trained model, fine-tuned naively, fails to adapt well.

machine vision quality control [34, 44, 19]. Such applications commonly require that input images are device independent and illuminant color-unbiased. Extraction of the intrinsic color information from scene surfaces by compensating for scene illuminant color is commonly referred to as "Color Constancy" (CC) or "Automatic White Balance" (AWB). The process of computational CC can be defined as the transformation of the source image, captured under an unknown illuminant, to a target image representing the same scene under a canonical illuminant. CC algorithms typically consist of two stages; first, estimation of the scene illuminant color and second, transformation of the source image, accounting for the illuminant, such that the resulting image illumination appears achromatic.

The perceived color of surfaces are determined by the intrinsic surface reflectance properties of objects in the scene, the spectral power distribution of the light(s) illuminating
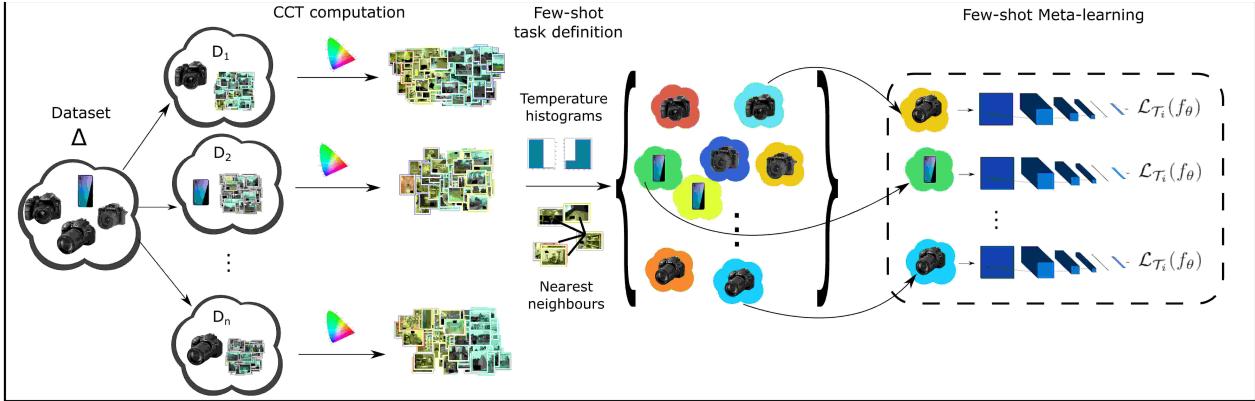
1

Figure 2: Overview of the proposed strategy defining task distribution $\mathcal{T}_i \sim p(\mathcal{T})$. Considering a set of cameras and camera specific images, we separate images into subtasks based on illuminant color. This is done by computing color temperature for each image, and building a CCT histogram for each camera. Images in the same task are defined as images captured using the same camera and belonging to the same CCT histogram bin.

them and physical capture device properties encompassing image sensor, camera spectral sensitivity (CSS) and lens effects. This combination of properties makes the problem locally underdetermined. In practice, if no device calibration prior is available, we can only observe a product of these factors, as measured in the digital image. More formally, we model a tri-chromatic photosensor response in the standard way such that:

$$\rho_k(X) = \int_\Omega E(\lambda)S(\lambda, X)R_k(\lambda)d\lambda \quad k \in \{R, G, B\}. \tag{1}$$

where $\rho_k(X)$ is the intensity of color channel $k$ at pixel location $X$, the wavelength of light $\lambda$ such that $E(\lambda)$ represents the spectrum of the illuminant, $S(\lambda, X)$ the surface reflectance at pixel location $X$ and $R_k(\lambda)$ the CSS for channel $k$, considered over the spectrum of visible wavelengths $\Omega$. The goal of computational CC then becomes estimation of the global illumination color $\rho_k^E$ where:

$$\rho_k^E = \int_\Omega E(\lambda)R_k(\lambda)d\lambda \quad k \in \{R, G, B\}. \tag{2}$$

Given that there exist infinitely many combinations of illuminant color and surface reflectance that result in identical (pixel value) observations $\rho_k(X)$, the problem of finding $\rho_k^E$ is (locally) ill-posed.

Modern supervised learning techniques can be used to infer this global image illuminant color and currently provide state-of-the-art estimation accuracy [9]. However, such approaches are typically CSS specific (i.e. consistent $R_k(\lambda)$) and therefore require, for each camera considered, acquisition of large sets of manually labelled images comprising a variety of scenes and illumination colors. This poses a barrier preventing such tools from providing highly accurate and robust illuminant estimation for new, previously unseen cameras in a manner that can be regarded as both quick and cheap.

In this paper we propose a new approach that removes the expensive, yet necessary for standard approaches, requirement of large amounts of labelled, sensor-specific image acquisition by decomposing the illuminant estimation problem such that it is robust with respect to variation in capture device. Using the concept of color temperature to infer the nature of scene light source, we frame the CC problem as a set of related yet distinct few-shot regression tasks, where each task is camera and illuminant specific. This enables us to exploit small image datasets, captured from disparate sources, and construct models with a capacity to learn camera-specific color biases quickly and cheaply using only a handful of target-device labelled images. Integrating our task definition approach within a meta-learning framework [17], we are able to train a joint model capable of quickly adapting to new unseen capture devices and report performance competitive with the fully-supervised, state of the art using a handful of camera-specific training samples. An overview of the proposed approach is depicted in Fig. 2.

The main contributions of this work are:

1. Our work constitutes the first few-shot learning approach for color constancy and enables the use of order(s) of magnitude fewer device-specific training images in comparison to contemporary work.

2. We introduce color temperature to the CC problem, demonstrate how it allows to estimate the type of light source from a photograph, and use it to frame CC as a set of simpler physically intuitive problems.

3. We provide extensive experiments on three public datasets and provide a comparative analysis of three meta-learning algorithm variants on a real-world image regression problem.

## 2. Related Work

Our contributions are closely related to previous learning based color constancy work, inter-camera considerations and few-shot learning techniques. We now provide brief review of these topics.

**Fully supervised methods.** Prior work can broadly be divided into statistics-based and learning-based methods [25]. Classical methods utilise low-level statistics that are fast and typically contain few free parameters. However, performance is highly dependent on strong scene content assumptions and these methods falter in cases where assumptions fail to hold. Early learning-based work [20, 46, 45, 36, 22] comprised of combinational and direct approaches, typically relying on hand-crafted image features which limited their overall performance.

Recent fully supervised convolutional CC work now offers state-of-the-art estimation accuracy. Both local patch-based [10, 40, 11, 27] and full image [8, 32, 9] input have been considered, investigating different model architectures [10, 11, 40] and the use of semantic information [27, 32]. Barron [8, 9] alternatively frames computational CC as a 2D spatial localisation problem. He represents image data using log-chroma histograms for which a single convolutional layer learns to evaluate illuminant color solutions in the chroma plane. Despite strong performance, fully supervised deep-learning techniques require large amounts of calibrated and hand-labelled sensor specific data to learn robust models for each target device [3]. This makes collection and calibration of imagery for data driven color constancy both restrictive and costly, commonly requiring placement of physical calibration objects in a large variety of scenes and illuminants, and subsequent manual segmentation to measure ground-truth illuminants.

Image augmentation, *eg.* synthetic relighting [11], and transfer-learning [10] using models pre-trained for alternative tasks, *eg.* ImageNet classification [30], have been previously employed to mitigate lack of available data. The former strategy commonly struggles with synthetic-data domain gap issues and may not generalise well to real-world image manifolds at inference time, while the misalignment between object classification and computational CC likely results in learning features invariant to appearance attributes of critical importance for CC, limiting the performance of the latter strategy.

**Inter-camera and unsupervised approaches.** Few color constancy works have attempted to mitigate the costs of sensor-specific data collection, calibration and image labelling. The work of [21] learns a transformation between pairs of camera CSS, but requires a priori knowledge of the cameras' CSS and is limited to pairs of sensors. Early unsupervised work [42] introduces a linear statistical model learned on a single sensor from video frame observations. Banic et al. [5] use classical statistical approaches to learn parameters that approximate the unknown ground-truth illumination of the training images, avoiding calibration and image labelling. Despite promising inter-camera performance [5], unsupervised techniques still require the collection of a large amount of unlabelled images under varying light sources, and yield subpar performance when compared to fully-supervised methods.

Our approach proposes to strongly restrict the data requirements via a few-shot formulation that is robust to variations in camera sensor, and bridges the gap between fully-supervised and unsupervised performances.

**Few-shot Learning.** Few-shot learning problems consist of learning a new task or concept using only a handful of data points (typically 1-10 samples per task) and have recently received considerable attention [43, 35, 41, 17]. This promises a number of advantages with regard to efficient model building for new tasks; reducing the need for data acquisition and labelling by order(s) of magnitude, and decreasing effort spent on fine-tuning and adaptation of existing models to novel problems. A popular meta-learning strategy consists of finding model initialisations that allow fast adaptation to new, previously unseen tasks [17]. The strategy has since been widely adopted for classification tasks [29] and several recently proposed extensions report increases to efficiency [33] and performance [31, 2]. While a natural separation of few-shot tasks exists for image classification problems, in contrast, problems framed as a regression (e.g. image illuminant estimation), require a careful task definition process so as to provide distinct yet homogeneous tasks that aid fast and accurate model adaption to new problem instances using limited training data.

To take advantage of such tools for color constancy, an important research question emerges, namely, how to decompose our problem in a set of few-shot tasks? This is crucially important when camera specific data is sparse.

## 3. Camera-Adaptive Color Constancy

An overview of our proposed Meta-AWB method is shown in Fig. 2. We consider a set of datasets $\Delta = \{\mathcal{D}_j\}$, where each dataset $\mathcal{D}_j = \{C_j, \{I_i\}_j\}$ comprises images acquired using a single camera $C_j$, representing various scenes under varying illumination conditions. Our objective is to provide a color constancy framework that is robust to variations in capture device, so as to leverage all the data available in $\Delta$ in order to adapt to previously unseen cameras with limited new training samples. We propose a physically intuitive model that casts CC as a set of simple illuminant and camera specific regression tasks. We use
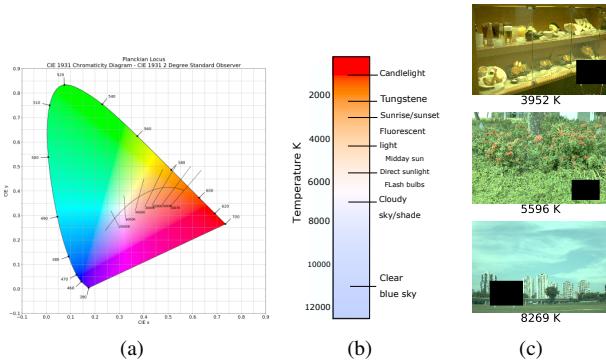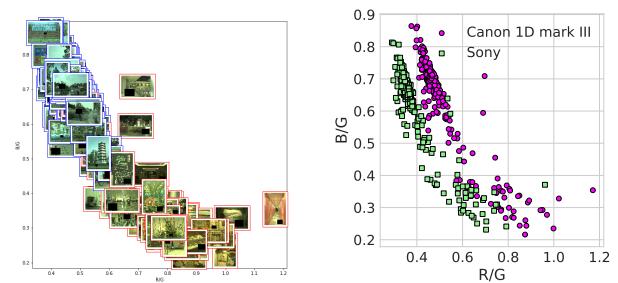
(a)    (b)    (c)

Figure 3: (a) Chromaticity space with Planckian locus. (b) Color temperature chart and types of light source associated with specific temperatures (c) Examples of images and corresponding color temperature $K$.

the concept of color temperature to approximate the type of light source illuminating each image in our datasets, allowing us to separate images based on light source in an unsupervised way. As tasks are illuminant specific, estimated illuminant corrections have limited variability and can be evaluated accurately with limited training samples. This allows us to frame each task as a few-shot regression problem, which can be addressed using recent few-shot learning strategies, such as meta-learning. Section 3.1 reviews the concept of color temperature and its computation from images before introducing our few-shot task definition strategy in Section 3.2. Section 3.3 details how our novel formulation is integrated in a meta-learning framework.

### 3.1. Correlated Color Temperature

Color Temperature (CT) is a common measurement in photography, often used in high-end camera software to describe the color of the illuminant for setting white balance [28]. By definition, CT measures, in degrees Kelvin, the temperature that is required to heat a Planckian (or black body) radiator to produce light of a particular color. A Planckian radiator is defined as a theoretical object that is a perfect radiator of visible light [38]. The Planckian locus, illustrated in Fig. 3(a), is the path that the color of a black body radiator would take in chromaticity space as the temperature increases, effectively illustrating all possible CTs.

In practice, the chromaticity of most light sources is off the Planckian locus, so the Correlated Color Temperature (CCT) is computed. CCT is the point on the Planckian locus closest to the non-Planckian light source [37, 38]. Intuitively, CCT describes the color of the light source and can be approximated from photos taken under this light. As shown in Fig. 3(b), different temperatures can be associated with different types of light [28]. For each image, we



(a) Pre-white balance images (Canon1D camera), marking their individual ground-truth gain corrections in $[\frac{r}{g}, \frac{b}{g}]$ space. Image frame border color (red, blue) indicates corresponding image temperature histogram bin membership.

(b) Ground-truth gain corrections for images observing identical scenes under similar illumination yet captured with distinct cameras (Canon1D, Sony; NUS-9 [14]).

Figure 4: Our meta-task definition is conditioned on both image temperature and camera, motivated by the expected image light source separability and CSS distribution shifts.

can compute CCT using standard approaches that map CIE 1931 chromaticities $x$ and $y$ to CCT [26]. Chromaticities $x$, $y$ are coordinates in the chromaticity space which can easily be estimated from the image's RGB values [37].

### 3.2. Illuminant and camera-specific tasks

Device-specific Camera Spectral Sensitivities (CSS) ($R_k(\lambda)$ in Eq. 2) affect the color domain of captured images and the recording of scene illumination. Images captured by different cameras can therefore exhibit ground-truth illuminant distributions that occupy differing regions of the chromaticity space [21], as can be observed in Fig. 4(b). Intuitively, this means that two images of the same scene and illuminant will have different illuminant corrections if taken by different cameras. In this context, a standard approach is to treat each camera dataset as an independent regression task. However, we expect to observe large variability in illuminant correction within one camera dataset, due to both scene and light source diversity. Achieving good performance and efficient generalisation to unseen cameras using tasks that contain too wide *intra-task* diversity may be difficult in a setting where camera specific data is sparse. Gamut based color constancy methods [18, 16, 6] assume that the color of the illuminant is constrained by the colors observed in the image. We make a similar hypothesis and aim to regroup images with similar dominant colors in the same task.

As a result, we decompose the inter-camera color constancy problem in a set of regression problems that comprise images acquired *with the same camera and with similar CCT* (i.e. similar dominant color). Our intuition is that, for each of these problems, illuminant corrections are clus-

tered such that good performance can be obtained quickly with only a handful of training samples.

We propose two strategies to separate images based on color temperature values. Our first approach is to compute a histogram $H_s$ for camera $s$ containing $M$ bins of CCT values and define each task as containing the set of images in each histogram bin. As a result, we define a task $T(\mathcal{D}_s, m) \in \mathcal{T}$ as: $T(\mathcal{D}_s, m) = \{ I \mid a_s^m \leq CCT(I) \leq b_s^m, Cam(I) = C_s\}$ where $Cam(I)$ is the camera used to acquire image $I$, and $a_s^m$, $b_s^m$ are the edges of bin $m$ in histogram $H_s$. Intuitively, images within the same temperature bin will have a similar dominant color, and therefore one could expect them to have similar illuminant corrections. Figure 3(b) highlights that a large variety of light sources are defined by relatively low temperatures. Accounting for this non-uniform distribution, we define bin edges of $H_s$ as a partition of temperature values on a logarithmic scale. In particular, when setting $M = 2$, we expect to separate images under a *warm* light source from images under a *cold* light source (*eg.* indoor images vs. outdoor images). This effect is illustrated in Fig. 4(a) where images are plotted marking their respective ground-truth gain correction in $[\frac{r}{g}, \frac{b}{g}]$ space and image frame border colors indicate temperature bin membership. This low granularity decomposition yields a few-shot scenario such that only 10 to 20 training images will be required to adapt to a previously unseen camera. A second, more granular approach consists of sampling K-nearest neighbour images in terms of temperatures, where K is the number of images comprising the regression task. Such a setting allows to separate the types of illuminants more precisely, but conversely requires more illuminant specific training images at test time.

### 3.3. Meta-learning formulation

Using the task formulation of Section 3.2, we can frame camera-adaptive illuminant estimation as a few-shot learning problem where tasks are used to define learning episodes. One way of approaching this type of problem is to use meta-learning techniques such as the popular MAML algorithm [17], where the strategy is to learn an optimal neural network *initialisation* capable of achieving strong performance on a new unseen task in only a few gradient updates, using only a small number of training samples.

Each regression task instance $\mathcal{T} : \hat{\boldsymbol{\rho}}_\theta = f_\theta(I)$ aims to estimate a global illuminant correction vector $\boldsymbol{\rho} = [r, g, b]$ for an input image $I$, where $f_\theta$ is a nonlinear function described by a neural network model. The model's parameters $\theta$ are learned by minimising the angular error loss:

$$\mathcal{L}_\mathcal{T}(\hat{\theta}) = \arccos(\frac{\hat{\boldsymbol{\rho}}_\theta}{\parallel \hat{\boldsymbol{\rho}}_\theta \parallel} \cdot \frac{\boldsymbol{\rho}}{\parallel \boldsymbol{\rho} \parallel}). \qquad (3)$$

Angular error provides a standard metric sensitive to the inferred orientation of the $\hat{\boldsymbol{\rho}}_\theta$ vector, with respect to the ground-truth $\boldsymbol{\rho}$ yet agnostic to its magnitude, providing independence to the brightness of the illuminant.

MAML is an iterative algorithm that learns a global set of parameters $\theta^*$ across tasks by optimising fine-tuning performance on each training task. Each iteration comprises an *inner update* which consists of fine-tuning global parameters $\theta$ to be task specific on a set of training images via $n$ gradient descent steps with learning rate $\alpha$. The second step, the *outer update*, updates $\theta$ as:

$$\theta^* = \theta - \beta \nabla_\theta \sum_i \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i}), \qquad (4)$$

where $\beta$ is the meta-learning rate parameter and $\mathcal{L}_{\mathcal{T}_i}(f_\theta)$ is the regression loss function as described in Eq. 3, computed using task specific fine-tuned parameters on a new set of (previously unseen) meta-test images. At test time, parameters are fine-tuned for a new unseen task for $n$ gradient updates and $K$ training samples. We finally compute the illuminant correction for each test image $I$ as $\boldsymbol{\rho}_{\theta_i} = f_{\theta_i}(I)$.

## 4. Results

**Datasets and preprocessing.** Three public color constancy datasets: Gehler-Shi [39, 22], NUS-9 [14], and Cube [5] are combined to investigate the capabilities of the proposed methodology. We use a total of 4128 images captured by 12 different cameras. For each dataset, we make use of the provided 'almost-raw' PNG images for all experimental work that follows in Section 4. Ground-truth illumination is measured by Macbeth Color Checker (MCC) in each dataset except for the 'Cube' database that alternatively uses a SpyderCube [1] calibration object. Illuminant ground-truth information (respective calibration objects) are masked in all images (using provided MCC coordinates and Cube+ using the fixed SpyderCube image location with mask value RGB $= [0, 0, 0]$), during both learning and inference. The **Gehler-Shi** dataset [39, 22] contains 568 images of indoor and outdoor scenes. Images were captured using Canon 1D and Canon 5D cameras. The **NUS-9**-Camera dataset [14] consists of 9 subsets of $\sim$210 images per camera providing a total of $(1736 + 117^1) = 1853$ images. All subsets comprise images representing the same scene, highlighting the influence of the camera sensor. The **Cube** dataset [5] contains 1365 images and consists of predominantly outdoor imagery (Canon EOS 550D camera). It was recently updated to include an additional 342 indoor images (renamed **Cube+** dataset). We use all Cube+ data at train time. At inference time, we limit evaluation to the Cube dataset, in order to be directly comparable to previous work.

Camera specific black-level corrections are applied in keeping with offsets specified in the dataset descriptions.

---

[1]The NUS dataset has recently been updated to include 117 additional images from a ninth camera. During training we use all nine cameras.

We apply a standard gamma correction ($\gamma = 2.2$) and normalize network input to $[0, 1]$. Input images are converted to 8-bit where required, providing bit depth consistency across all datasets. Impoverished input (low bit-depth) has been shown to make the color constancy problem more difficult [9] yet also provides a good real-world test bed as specialized camera hardware typically performs illuminant estimation using small, low bit-depth images.

**Implementation.** For each camera, we train a model using random image crops of variable size ($128 \times 128$ to original image size) from the remaining 11 cameras, spatially resized to $128 \times 128$. Due to the discussed limits on typical dataset size per camera and assumed simplicity of our regression tasks, we adopt a simple architecture comprising four convolutional layers (all sized $3 \times 3 \times 64$), an average pooling layer and two fully connected layers (sizes $64 \times 64$, $64 \times 3$), with ReLU activations on all hidden layers. Due to the makeup of our proof-of-concept amalgamated dataset, the majority of cameras capture similar scene content (NUS-9 imagery). Since we aim to optimise generalisation between cameras without overfitting to particular scenes, we choose the Gehler-Shi (Canon 5D) images as a validation set. This allows for optimisation of model hyperparameters using imagery containing unique scene content.

We evaluate three variants of MAML, characterised by different definitions of the learning rate $\alpha$. *MAML* [17], uses a constant $\alpha$; *metaSGD* [31] learns an $\alpha$ value per parameter in the network, allowing the direction and magnitude of the gradient descent to be learned; *LSLR* [2] substantially reduces the number of trainable parameters, learning a $\alpha$ single parameter per *layer* in the network for each inner gradient update. We train our CNN model for $25k$ iterations, using a meta-batch size of 10, number of training images per batch $K_{train} = 10$, learning rate $\beta = 0.001$ (with exponential decay). The inner-update learning rate $\alpha$ is set (or initialised for *metaSGD* and *LSLR*) to 0.001. We use layer normalisation on all convolutional layers. At inference time, for each test image, we randomly select $K_{test}$ training samples (from the test image task) and fine-tune the model for 10 iterations. To evaluate the statistical robustness of our method to variation in the selection of the $K_{test}$ images, we independently repeat 10 draws for each test image. We report the median angular error over all images, averaged over all draws. As a baseline, we train a model with standard back-prop and leave-one-out cross validation on camera using network architecture matching our introduced base-learner. At test time we report both with (*Baseline - fine tune*) and without (*Baseline - no fine tune*) $K$-shot fine tuning. Baselines are trained for $25k$ iterations using the same parameters as our Meta-AWB model.

**Histogram-based task definition.** To explore the validity of our learning task formulation, we plot CCT histograms per camera and ground-truth $[r, g, b]$ gain correc-



(a) Cube    (b) Canon600D (NUS)    (c) Canon 600D (NUS) ground-truth illuminants in RGB space
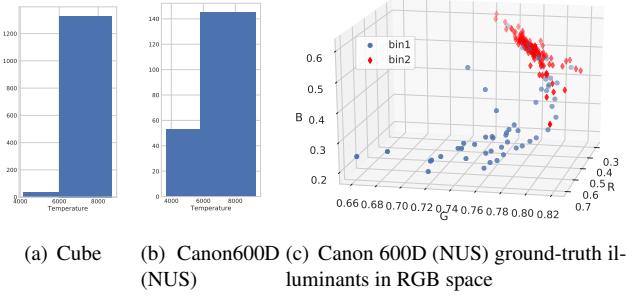
Figure 5: Task definition: Temperature histograms and corresponding separation in RGB space.

tions per image in RGB space, with CCT bin assignment indicated. Figure 5(b) and 5(c) provide examples of these respectively for the NUS Canon 600D image set ($M = 2$ bins). Correlating bin edges (as defined in Section 3.2) to the temperature chart found in Figure 3(b), we see that images can be separated based on light type, and more specifically, indoor vs outdoor light sources. This observation is confirmed in Fig. 5(a) where the CCT histogram, obtained using the original subset of Cube images (all images contain outdoor scenes), essentially contains all images in one bin. Figure 5 also illustrates that our CCT-histogram strategy generates homogeneous learning tasks since ground-truth illuminant corrections belonging to the same bin are well clustered in RGB space. This setup assigns images to a task, conditioned on both camera sensor and CCT bin, resulting in $M \cdot |\{C_j \in \Delta\}|$ valid learning tasks assuming that each CCT bin is non-empty. In the remaining experiments, we set $M = 2$ as discussed in Sec. 3.2.

**Parameter and method analysis.** Using our validation camera we evaluate the influence of key parameters and meta-learning strategy. For each method considered (*MAML*, *metaSGD* and *LSLR*), we train models for variable numbers of inner gradient updates $n_{train} \in \{1, 5, 10\}$. While $K_{train} = 10$ images is fixed during meta-model training, we evaluate the influence of available $k$-shot image count at test time, computing performance for $K_{test} \in \{5, 10, 20\}$. We also report results for inner updates $n_{test} \in \{1, 5, 10\}$. Since *LSLR* learns a different learning rate per inner update, we set $\alpha_i = \alpha_n, \forall i \geq n$ when $n_{test}$ is set to a value larger than that used during training.

We report results in Fig. 6(a) and Table 1, with the best results for each $K_{test}$ reported in bold. We observe a substantial improvement in performance when increasing the number of inner updates from 1 to 5, but note that 10 updates do not improve performance. *LSLR* appears to offers a compromise between simplicity and flexibility, yielding the best results when $n_{train} = 5, 10$. However, interestingly, *LSLR* and *metaSGD* perform poorly compared to *MAML*
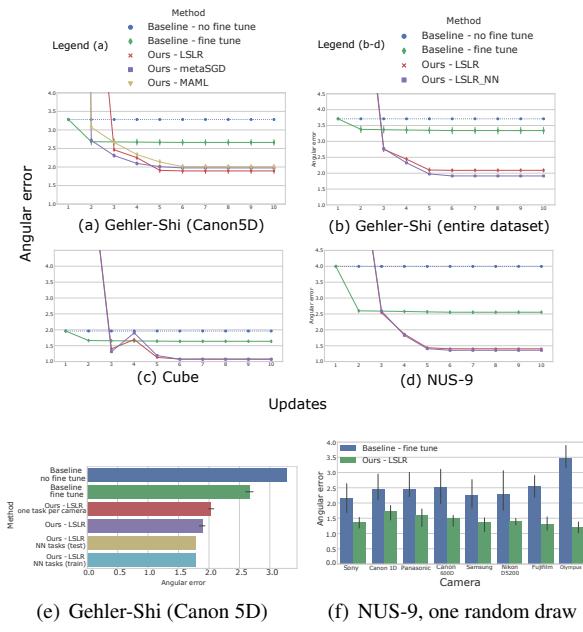
Figure 6: Median angular-error with respect to the number of fine-tuning updates. (a,e) Parameter study results, (b-d) Dataset specific results compared to the baselines, (f) Per camera angular-error after 10 updates, all NUS-9 cameras. Error bars in (a-e) report inter-draws variance.

| | | MAML | | | metaSGD | | | LSLR | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | \multicolumn{9}{c}{Training updates $u = n_{train}$} | | | | | | | | |
| | | $u{=}1$ | $u{=}5$ | $u{=}10$ | $u{=}1$ | $u{=}5$ | $u{=}10$ | $u{=}1$ | $u{=}5$ | $u{=}10$ |
| | $K_{test}=5$ | | | | | | | | | |
| | 1 update | 2.18 | 3.08 | 3.11 | 2.12 | 2.80 | 3.10 | 2.42 | 9.07 | 4.68 |
| | 5 updates | 2.06 | 2.07 | 2.08 | 2.27 | 2.11 | 2.07 | 2.31 | **2.00** | 2.76 |
| | 10 updates | 2.06 | 2.07 | 2.08 | 2.28 | 2.11 | 2.06 | 2.31 | **2.00** | 2.05 |
| | $K_{test}=10$ | | | | | | | | | |
| | 1 update | 2.13 | 3.05 | 3.12 | 2.00 | 2.71 | 3.09 | 2.37 | 9.05 | 4.52 |
| | 5 updates | 2.02 | 2.01 | 2.00 | 2.09 | 1.98 | 1.94 | 2.23 | **1.87** | 2.50 |
| | 10 updates | 2.02 | 2.00 | 2.00 | 2.09 | 1.97 | 1.93 | 2.23 | **1.87** | 1.92 |
| | $K_{test}=20$ | | | | | | | | | |
| | 1 update | 2.10 | 3.05 | 3.09 | 1.99 | 2.70 | 3.05 | 2.31 | 9.01 | 4.40 |
| | 5 updates | 1.98 | 1.94 | 2.00 | 2.06 | 1.89 | 1.91 | 2.18 | **1.81** | 2.41 |
| | 10 updates | 1.98 | 1.94 | 1.95 | 2.06 | 1.89 | 1.91 | 2.18 | **1.81** | 1.84 |

(Testing parameters — row group label on the left)

Table 1: Our meta-learning hyper-parameter investigation and method analysis. Median angular-error. Best results for each $K$-shot configuration are reported in bold.

when training with only 1 gradient update. Predictably, we observe an increase in performance as we increase $K_{test}$ reaching our overall best performance with a median angular error of 1.81 degrees and $K_{test} = 20$ samples. Finally, we can see that all methods benefit from $n_{test} > 1$ updates, but tend to plateau when $n_{test} > 5$. Considering our experimental observations, we use $n_{train} = 5$ and $n_{test} = 10$ with our *LSLR* variant for the remainder of our experimental work. We set $K_{test} = 10$ unless otherwise specified.

**Influence of task definition strategy.** Using our validation camera, we further evaluate the influence of different task definition approaches. As shown in Fig. 6(e), we compare our $M = 2$ bins histogram based approach (*Ours - LSLR*) to 1) $M = 1$, corresponding to the naive approach of setting one camera dataset $D_s$ as a task (*Ours - LSLR - one task per camera*), 2) defining tasks as temperature nearest neighbours at test time (*Ours - LSLR - NN tasks (test)*) and both and train and test time (*Ours - LSLR - NN tasks (train)*). Results are compared to the baseline method for context. We observe that results improve as task definition methods get more granular (i.e. provide better illuminant separation), while our two experiments on nearest neighbour tasks show that testing on more granular tasks (with respect to train tasks) can equally improve performance and allows to make use of potential larger datasets at test time.

**Comparisons to the state of the art** Figure 6 shows the evolution of the median angular error, with respect to the number of gradient updates, for all datasets under both baselines and our approach (with and without nearest neighbour tasks at test time). We observe a significant gap in performance for all datasets. Fig. 6(f) shows per camera performance and highlights that, unlike the baseline which particularly struggles for some cameras (*eg.* Olympus), our method provides a consistent and better performance on each NUS-9 test camera individually, in addition to average performance. We provide a visual example in Fig. 7.

Finally, we compare our results on all datasets with recent state of the art approaches. We report results on NUS-8 (without the recently added Nikon D40 camera) to provide a fair and accurate comparison. Quantitative results are shown in Tables 2 (NUS-8), 3 (Gehler-Shi) and 4 (Cube) where we report standard angular-error statistics (Tri. is trimean and G.M. geometric mean). We obtain results that are competitive with the state of the art and fully supervised methods, despite using an order of magnitude less camera specific training data. We achieve good generalisation on all datasets, in particular with the NUS-8 and Cube datasets, where we outperform most state of the art methods. The superior performance on NUS-8 can be linked to the fact that the NUS scene content is repeatedly seen during training.

## 5. Conclusion

In this paper, we propose a novel formulation of the computational color constancy problem that adapts and generalises quickly to a large variety of camera sensors. We exploit the concept of color temperature to approximate the type of light source from images, so as to decompose the problem into a set of simpler regression tasks, each associated with a camera sensor and type of light source. The simplified nature of the obtained regression tasks allows

| Algorithm | Mean | Median | Tri. | Best 25% | Worst 25% | G.M. |
|---|---|---|---|---|---|---|
| Low-level statistics-based methods | | | | | | |
| White-Patch [12] | 9.91 | 7.44 | 8.78 | 1.44 | 21.27 | 7.24 |
| Gray-world [13] | 4.59 | 3.46 | 3.81 | 1.16 | 9.85 | 3.70 |
| Edge-based Gamut [24] | 4.40 | 3.30 | 3.45 | 0.99 | 9.83 | 3.45 |
| Natural Image Statistics [23] | 3.45 | 2.88 | 2.95 | 0.83 | 7.18 | 2.81 |
| Fully-supervised learning | | | | | | |
| Bayesian [22] | 3.50 | 2.36 | 2.57 | 0.78 | 8.02 | 2.66 |
| Cheng et al. 2014 [14] | 2.93 | 2.33 | 2.42 | 0.78 | 6.13 | 2.40 |
| SqueezeNet-FC4 [27] | 2.23 | 1.57 | 1.72 | 0.47 | 5.15 | 1.71 |
| CCC [8] | 2.38 | 1.48 | 1.69 | 0.45 | 5.85 | 1.74 |
| Cheng et al. 2015 [15] | 2.18 | 1.48 | 1.64 | 0.46 | 5.03 | 1.65 |
| Shi et al. 2016 [40] | 2.24 | 1.46 | 1.68 | 0.48 | 6.08 | 1.74 |
| FFCC [9] (thumb, 8bit input) | 2.06 | 1.39 | 1.53 | 0.39 | 4.80 | 1.53 |
| FFCC [9] | 1.99 | 1.31 | 1.43 | 0.35 | 4.75 | 1.44 |
| {Unsupervised, Few-shot} learning | | | | | | |
| Color Tiger [5] | 2.96 | 1.70 | 1.97 | 0.53 | 7.50 | 2.09 |
| Meta-AWB $K = 10$ | 1.93 | 1.38 | 1.49 | 0.47 | 4.37 | 1.52 |
| Meta-AWB $K = 20$ | 1.89 | 1.34 | 1.44 | 0.45 | 4.28 | 1.47 |

Table 2: Performance on the NUS-8 dataset [14]. We follow the same format as [9], reporting average performance (geometric mean) over the 8 original NUS cameras. Results outperformed by ours are marked in gray.

| Algorithm | Mean | Median | Tri. | Best 25% | Worst 25% | G.M. |
|---|---|---|---|---|---|---|
| Low-level statistics-based methods | | | | | | |
| Gray-world [13] | 6.36 | 6.28 | 6.28 | 2.33 | 10.58 | 5.73 |
| White-Patch [12] | 7.55 | 5.86 | 6.35 | 1.45 | 16.12 | 5.76 |
| Edge-based Gamut [24] | 6.52 | 5.04 | 5.43 | 1.90 | 13.58 | 5.40 |
| Fully-supervised learning | | | | | | |
| Bayesian [22] | 4.82 | 3.46 | 3.88 | 1.26 | 10.49 | 3.86 |
| Cheng et al. 2014 [14] | 3.52 | 2.14 | 2.47 | 0.50 | 8.74 | 2.41 |
| Bianco CNN [10] | 2.63 | 1.98 | 2.10 | 0.72 | 3.90 | 2.04 |
| Cheng et al. 2015 [15] | 2.42 | 1.65 | 1.75 | 0.38 | 5.87 | 1.73 |
| CCC [8] | 1.95 | 1.22 | 1.38 | 0.35 | 4.76 | 1.40 |
| SqueezeNet-FC4 [27] | 1.65 | 1.18 | 1.27 | 0.38 | 3.78 | 1.22 |
| DS-Net [40] | 1.90 | 1.12 | 1.33 | 0.31 | 4.84 | 1.34 |
| FFCC [9] (thumb, 8bit input) | 2.01 | 1.13 | 1.38 | 0.30 | 5.14 | 1.37 |
| FFCC [9] | 1.61 | 0.86 | 1.02 | 0.23 | 4.27 | 1.07 |
| Few-shot learning | | | | | | |
| Meta-AWB $K = 10$ | 3.07 | 2.08 | 2.28 | 0.56 | 7.31 | 2.26 |
| Meta-AWB $K = 20$ | 2.99 | 2.02 | 2.18 | 0.55 | 7.19 | 2.20 |
| Meta-AWB NN $K = 10$ | 2.66 | 1.91 | 1.99 | 0.49 | 6.20 | 1.98 |
| Meta-AWB NN $K = 20$ | 2.57 | 1.84 | 1.94 | 0.47 | 6.11 | 1.92 |

Table 3: Performance on the Gehler-Shi dataset [39, 22]. Previous methods as reported by [9]. Results outperformed by our best method are marked in gray.

| Algorithm | Mean | Median | Tri. | Best 25% | Worst 25% | G.M. |
|---|---|---|---|---|---|---|
| Low-level statistics-based methods | | | | | | |
| White-Patch [12] | 6.58 | 4.48 | 5.27 | 1.18 | 15.23 | 4.88 |
| Gray-World [13] | 3.75 | 2.91 | 3.15 | 0.69 | 8.18 | 2.87 |
| General Gray-World [7] | 2.50 | 1.61 | 1.79 | 0.37 | 6.23 | 1.76 |
| Smart Color Cat [4] | 1.49 | 0.88 | 1.06 | 0.24 | 3.75 | 1.04 |
| {Unsupervised, Few-shot} learning | | | | | | |
| Color Tiger [5] | 2.94 | 2.59 | 2.66 | 0.61 | 5.88 | 2.35 |
| Meta-AWB $K = 10$ | 1.63 | 1.08 | 1.20 | 0.31 | 3.89 | 1.17 |
| Meta-AWB $K = 20$ | 1.59 | 1.02 | 1.15 | 0.30 | 3.85 | 1.16 |

Table 4: Performance on the Cube dataset. Previous methods as reported by [5]. Results outperformed by ours are marked in gray.

us to cast color constancy as a few-shot learning problem that we address using meta-learning. Extensive experiments



(a) input image



(b) Ground-truth solution



(c) Meta-AWB, angular error: 5.42°



(d) Baseline fine-tuned, angular error: 18.76°

Figure 7: A challenging Gehler-Shi [22, 39] test image from our worst 25%. Images are shown in sRGB space and clipped at the 97.5 percentile. See supplementary material for additional qualitative results.

across three benchmark datasets and 12 different camera sensors result in performance competitive with the fully-supervised state-of-the-art, using only a small fraction of camera specific data at test time. We presented and studied the influence of several variants of our technique, including task definition approaches. We show improved learning ability over standard fine-tuning, resulting in efficient use of only few training samples. Meta-AWB has the ability to generalise quickly and learns to solve the computational color constancy problem in a camera agnostic fashion. This provides the potential for high accuracy performance as new sensors become available yet mitigates arduous and time-consuming calibration of training imagery, required for fully-supervised approaches. One would expect to reach better generalisation performance with more imaging content variability per camera. Future work will investigate different base learner components, alternatives to meta-learning to address the few-shot learning problem and more diverse task definition approaches (e.g. scene content).

# References

[1] Datacolor SpyderCube. http://www.datacolor.com/photography-design/product-overview/spydercube/. Accessed: 2018-11-05. 5

[2] A. Antoniou, H. Edwards, and A. Storkey. How to train your MAML. *arXiv preprint arXiv:1810.09502*, 2018. 3, 6

[3] Ç. Aytekin, J. Nikkanen, and M. Gabbouj. A data set for camera-independent color constancy. *IEEE Transactions on Image Processing*, 27(2):530–544, 2018. 3

[4] N. Banić and S. Lončarić. Using the red chromaticity for illumination estimation. In *Image and Signal Processing and Analysis (ISPA), 2015 9th International Symposium on*, pages 131–136. IEEE, 2015. 8

[5] N. Banic and S. Loncaric. Unsupervised learning for color constancy. *CoRR*, abs/1712.00436, 2017. 3, 5, 8

[6] K. Barnard. Improvements to gamut mapping colour constancy algorithms. In *European conference on computer vision*, pages 390–403. Springer, 2000. 4

[7] K. Barnard, V. Cardei, and B. Funt. A comparison of computational color constancy algorithms. i: Methodology and experiments with synthesized data. *IEEE transactions on Image Processing*, 11(9):972–984, 2002. 8

[8] J. T. Barron. Convolutional color constancy. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 379–387, 2015. 3, 8, 11

[9] J. T. Barron and Y.-T. Tsai. Fast fourier color constancy. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, 21–26 July 2017*, 2017. 2, 3, 6, 8

[10] S. Bianco, C. Cusano, and R. Schettini. Color constancy using cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 81–89, 2015. 3, 8

[11] S. Bianco, C. Cusano, and R. Schettini. Single and multiple illuminant estimation using convolutional neural networks. *IEEE Transactions on Image Processing*, 26(9):4347–4362, 2017. 3

[12] D. H. Brainard and B. A. Wandell. Analysis of the retinex theory of color vision. *JOSA A*, 3(10):1651–1661, 1986. 8

[13] G. Buchsbaum. A spatial processor model for object colour perception. *Journal of the Franklin institute*, 310(1):1–26, 1980. 8

[14] D. Cheng, D. K. Prasad, and M. S. Brown. Illuminant estimation for color constancy: why spatial-domain methods work and the role of the color distribution. *JOSA A*, 31(5):1049–1058, 2014. 4, 5, 8, 11

[15] D. Cheng, B. Price, S. Cohen, and M. S. Brown. Effective learning-based illuminant estimation using simple features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1000–1008, 2015. 8

[16] G. D. Finlayson. Color in perspective. *IEEE transactions on Pattern analysis and Machine Intelligence*, 18(10):1034–1038, 1996. 4

[17] C. Finn, P. Abbeel, and S. Levine. Model-Agnostic Meta-Learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*, 2017. 2, 3, 5, 6

[18] D. A. Forsyth. A novel algorithm for color constancy. *International Journal of Computer Vision*, 5(1):5–35, 1990. 4

[19] B. Funt, K. Barnard, and L. Martin. Is machine colour constancy good enough? In *European Conference on Computer Vision*, pages 445–459. Springer, 1998. 1

[20] B. Funt and W. Xiong. Estimating illumination chromaticity via support vector regression. In *Color and Imaging Conference*, volume 2004, pages 47–52. Society for Imaging Science and Technology, 2004. 3

[21] S.-B. Gao, M. Zhang, C.-Y. Li, and Y.-J. Li. Improving color constancy by discounting the variation of camera spectral sensitivity. *JOSA A*, 34(8):1448–1462, 2017. 3, 4

[22] P. V. Gehler, C. Rother, A. Blake, T. Minka, and T. Sharp. Bayesian color constancy revisited. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008. 3, 5, 8, 11

[23] A. Gijsenij and T. Gevers. Color constancy using natural image statistics and scene semantics. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(4):687–698, 2011. 8

[24] A. Gijsenij, T. Gevers, and J. Van De Weijer. Generalized gamut mapping using image derivative structures for color constancy. *International Journal of Computer Vision*, 86(2-3):127–139, 2010. 8

[25] A. Gijsenij, T. Gevers, J. Van De Weijer, et al. Computational color constancy: Survey and experiments. *IEEE Transactions on Image Processing*, 20(9):2475–2489, 2011. 3

[26] J. Hernandez-Andres, R. L. Lee, and J. Romero. Calculating correlated color temperatures across the entire gamut of daylight and skylight chromaticities. *Applied optics*, 38(27):5703–5709, 1999. 4

[27] Y. Hu, B. Wang, and S. Lin. Fc4: Fully convolutional color constancy with confidence-weighted pooling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR17)*, pages 4085–4094, 2017. 3, 8

[28] R. Jacobson, S. Ray, G. G. Attridge, and N. Axford. *Manual of Photography*. Taylor & Francis, 2000. 4

[29] M. A. Jamal, G.-J. Qi, and M. Shah. Task-agnostic meta-learning for few-shot learning. *arXiv preprint arXiv:1805.07722*, 2018. 3

[30] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 3

[31] Z. Li, F. Zhou, F. Chen, and H. Li. Meta-SGD: Learning to learn quickly for few shot learning. *arXiv preprint arXiv:1707.09835*, 2017. 3, 6

[32] Z. Lou, T. Gevers, N. Hu, M. P. Lucassen, et al. Color constancy by deep learning. In *BMVC*, pages 76–1, 2015. 3

[33] A. Nichol, J. Achiam, and J. Schulman. On first-order meta-learning algorithms. *CoRR, abs/1803.02999*, 2018. 3

[34] R. Ramanath, W. E. Snyder, Y. Yoo, and M. S. Drew. Color image processing pipeline. *IEEE Signal Processing Magazine*, 22(1):34–43, 2005. 1

[35] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. 2017. 3

[36] C. Rosenberg, A. Ladsariya, and T. Minka. Bayesian color constancy with non-gaussian models. In *Advances in neural information processing systems*, pages 1595–1602, 2004. 3

[37] J. Schanda. *Colorimetry: understanding the CIE system.* John Wiley & Sons, 2007. 4

[38] E. F. Schubert. *Light-emitting diodes*. E. Fred Schubert, 2018. 4

[39] L. Shi and B. Funt. Re-processed version of the gehler color constancy dataset of 568 images. *http://www. cs. sfu. ca/˜ color/data/*, 2000. 5, 8, 11

[40] W. Shi, C. C. Loy, and X. Tang. Deep specialized network for illuminant estimation. In *European Conference on Computer Vision*, pages 371–387. Springer, 2016. 3, 8

[41] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, pages 4077–4087, 2017. 3

[42] K. Tieu and E. G. Miller. Unsupervised color constancy. In *Advances in neural information processing systems*, pages 1327–1334, 2003. 3

[43] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems*, pages 3630–3638, 2016. 3

[44] M. Vrhel, E. Saber, and H. J. Trussell. Color image generation and display technologies. *IEEE Signal Processing Magazine*, 22(1):23–33, 2005. 1

[45] N. Wang, D. Xu, and B. Li. Edge-based color constancy via support vector regression. *IEICE transactions on information and systems*, 92(11):2279–2282, 2009. 3

[46] W. Xiong and B. Funt. Estimating illumination chromaticity via support vector regression. *Journal of Imaging Science and Technology*, 50(4):341–348, 2006. 3

## Appendix A. Additional qualitative results

In Figure 8 we provide additional qualitative results in the form of test images from the Gehler-Shi dataset [22, 39]. For each image we show the input image produced by the camera and a white-balanced image corrected using the ground-truth illumination. We also show the output of our model ("Meta-AWB"), and that of the baseline fine-tuning approach reported in the paper. Color checker boards are visible in the images, however the relevant areas are masked prior to inference. Images are shown in sRGB space and clipped at the 97.5 percentile.

In similar fashion to [8], we adopt the strategy of sorting the test images by the combined mean angular-error of the two evaluated methods. We present images of increasing average difficulty however images to report were selected by instead ordering from *"hard"* to *"easy"* and sampling with a logarithmic spacing, providing a greater number of samples that proved challenging, on average.

## Appendix B. Color temperature distributions for learning task formulation

In our paper, we decompose the inter-camera color constancy problem into a set of regression tasks such that each task comprises images acquired from the same camera and under similar illumination settings. One approach we propose is to compute Correlated Color Temperature (CCT) histograms and define the images of each task as those that belong to a given histogram bin.

In the main paper, we limit the number of bins to $M = 2$, as this allows the broad separation of images into indoor and outdoor illuminants. Here we provide additional examples of CCT histograms and their corresponding clustering of ground-truth (GT) illuminant corrections in RGB space. Results for the $M = 3$ bins case, for two cameras (Canon 600D and Nikon D40 from the NUS-9 dataset [14]) are shown in Fig. 9.

We first observe that for smaller sized datasets (e.g. Nikon D40 data), there are not enough images ($< 10 = K$) available in the first bin as shown in 9(c). Furthermore, while ground-truth illuminants retain reasonable cluster separability, bin edges between different cameras may be misaligned. Finally, increasing the number of histogram bins also increases training data requirements as each bin requires 10 to 20 training images. These results motivate the choice of $M = 2$, as shown in Figure 4 of the main paper, thus providing distinctive tasks with sufficient images in each bin of the CCT histogram.

## Appendix C. Interpretation of camera specific results

We visualise the distributions of ground truth RGB corrections per camera in Fig.11 where datasets with larger median angular errors are shown with warmer colormaps, as well as their respective median value in Fig.11(a). Ground-truth distributions are linked to differences in scene content and camera (CSS, lens and sensor effects). In particular, camera latent effects can be considered dominant for the NUS dataset cameras due to the matching inter-camera scene content of images. Distribution difference for other cameras can be linked to both scene content and the considered hardware effects and is therefore less comparable.

Linking distributions in Fig. 11 to the results obtained, we observe that best results are obtained for cameras with compact distributions (ie. similar scene content, as observed for the Cube+ dataset). Furthermore, the very different nature of the Gehler-Shi (Canon 1D) and Nikon D40 distributions suggest that these cameras are more difficult to adapt to.

In particular, comparative results on the NUS-9 dataset between the fine-tuned baseline (blue) and our approach (green) are shown in Fig. 11(b). Cameras are ordered from left to right from lower to larger gap between baseline and our method. This shows that our method adapts well to cameras with distributions that are shifted with respect to the majority of NUS-9 cameras (e.g. Fujifilm, Olympus). These cameras exhibit a larger increase in performance with respect to our baseline. Similarly, we observe one of the larger gaps in performance for the Nikon D40 camera, with images that were acquired later and comprises different scenes and illuminations.
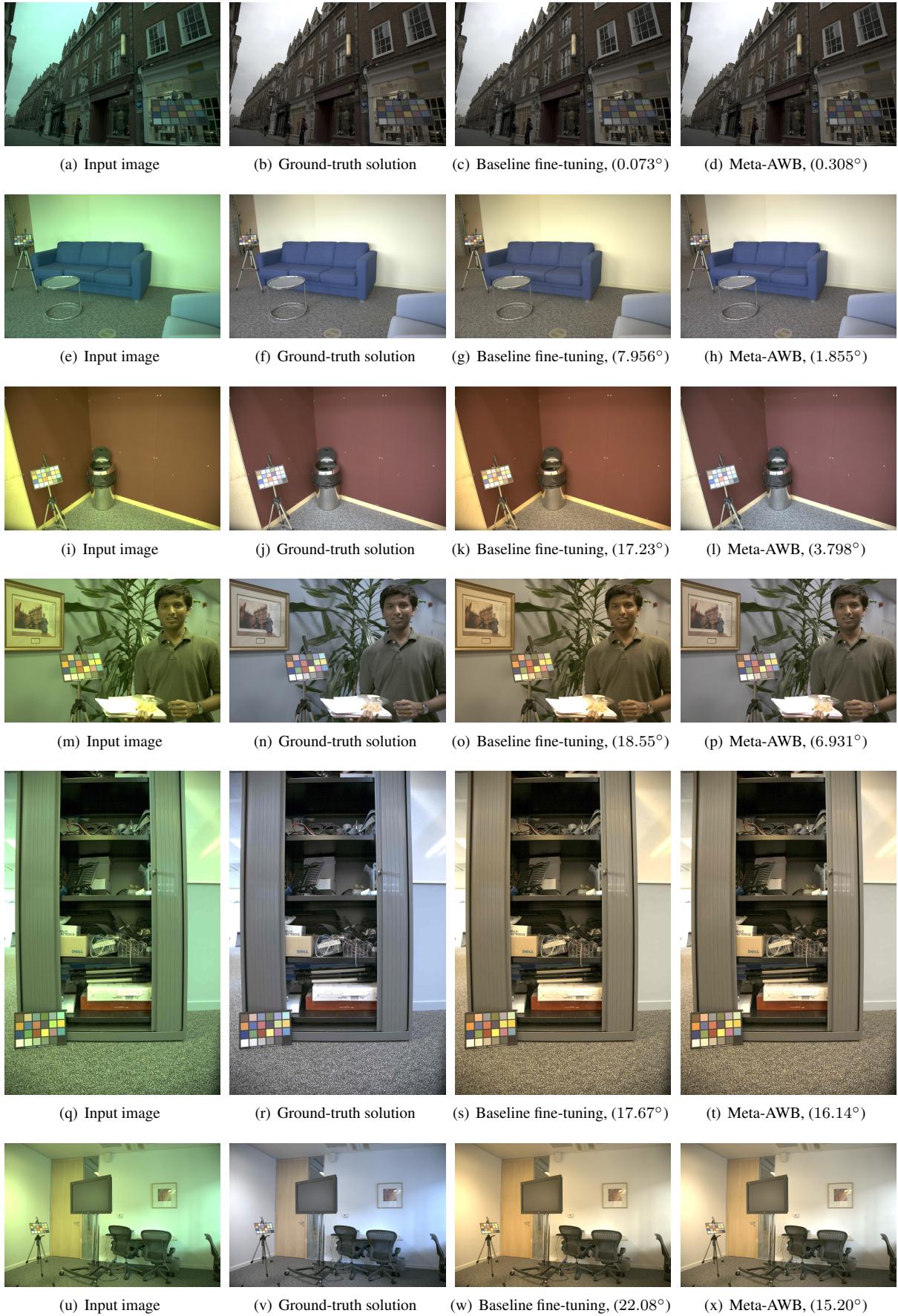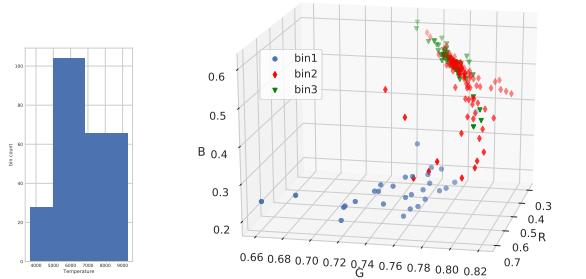
(a) Input image      (b) Ground-truth solution      (c) Baseline fine-tuning, (0.073°)      (d) Meta-AWB, (0.308°)

(e) Input image      (f) Ground-truth solution      (g) Baseline fine-tuning, (7.956°)      (h) Meta-AWB, (1.855°)

(i) Input image      (j) Ground-truth solution      (k) Baseline fine-tuning, (17.23°)      (l) Meta-AWB, (3.798°)

(m) Input image      (n) Ground-truth solution      (o) Baseline fine-tuning, (18.55°)      (p) Meta-AWB, (6.931°)

(q) Input image      (r) Ground-truth solution      (s) Baseline fine-tuning, (17.67°)      (t) Meta-AWB, (16.14°)

(u) Input image      (v) Ground-truth solution      (w) Baseline fine-tuning, (22.08°)      (x) Meta-AWB, (15.20°)
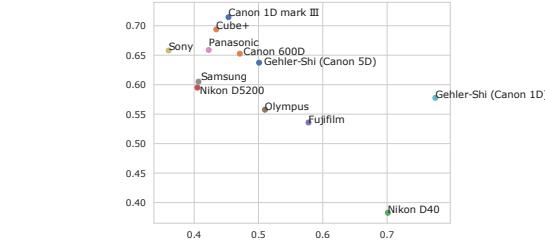
Figure 8: For each scene we present the input image produced by the camera alongside the ground-truth white-balanced image. Images are shown in sRGB space and normalized to the 97.5th percentile. For both our "Meta-AWB" approach and the baseline algorithm we show the white-balanced image, as well as the angular-error of the estimated illumination in degrees, with respect to the ground-truth.

(a) Canon 600D histogram

(b) Canon 600D GT illuminants in RGB space



(c) Nikon D40 histogram

(d) Nikon D40 GT illuminants in RGB space

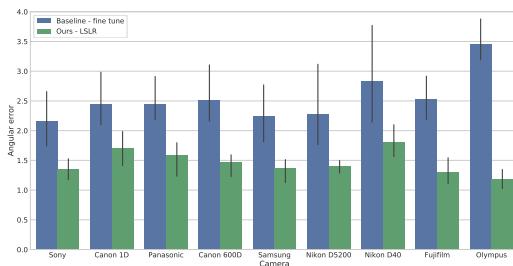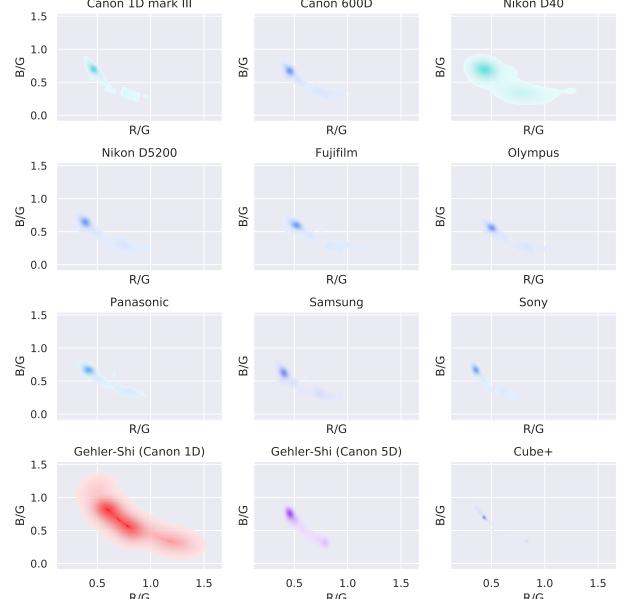Figure 9: CCT histograms $H_s$, 3 bin example.



Figure 10: Median angular error for each NUS-9 camera, our method (green) and fine-tuned baseline (blue). Cameras are ordered from left to right with increasing gap in performance between baseline and ours.



(a) Median ground-truth RGB correction for all cameras.



(b) Ground-truth RGB distributions for all cameras. Shifts in distributions are linked to camera effects and scene content. The color of each dataset distribution is defined by a colormap linked to the angular errors obtained over all images and cameras. Warmer colors correspond to larger median dataset errors.

Figure 11: Visualisation of variability between distributions of ground-truth RGB illuminant corrections for all cameras.