# Dissecting Image Crops

Basile Van Hoorick
Columbia University
New York, NY, USA
basile@cs.columbia.edu

Carl Vondrick
Columbia University
New York, NY, USA
vondrick@cs.columbia.edu

## Abstract

*The elementary operation of cropping underpins nearly every computer vision system, ranging from data augmentation and translation invariance to computational photography and representation learning. This paper investigates the subtle traces introduced by this operation. For example, despite refinements to camera optics, lenses will leave behind certain clues, notably chromatic aberration and vignetting. Photographers also leave behind other clues relating to image aesthetics and scene composition. We study how to detect these traces, and investigate the impact that cropping has on the image distribution. While our aim is to dissect the fundamental impact of spatial crops, there are also a number of practical implications to our work, such as detecting image manipulations and equipping neural network researchers with a better understanding of short-cut learning. Code is available at* https://github.com/basilevh/dissecting-image-crops.

## 1. Introduction

The basic operation of cropping an image underpins nearly every computer vision paper that you will be reading this week. Within the first few lectures of most introductory computer vision courses, convolutions are motivated as enabling feature invariance to spatial shifts and cropping [44, 28, 2]. Neural networks rely on image crops as a form of data augmentation [25]. Computational photography applications will automatically crop photos in order to improve their aesthetics [42, 10]. Predictive models extrapolate pixels out from crops [49, 47]. Even the latest self-supervised representations depend on crops for contrastive learning to induce rich visual representations [11, 40].

This core visual operation can have a significant impact on photographs. As Oliva and Torralba told us twenty years ago, scene context drives perception [39]. Recently, image cropping has been at the heart of media disinformation. Figure 1 shows two popular photographs where the photographer or media organization spatially cropped out part of the context, altering the message of the image. Twitter's auto-



Figure 1: We show two infamous image crops, visualized by the red box. (**left**) An Ugandan climate activist had been cropped out of the photo before it was posted in an online news article, the discovery of which sparked controversy [14]. (**right**) A news network had cropped out a large stick being held by a demonstrator during a protest [12]. Cropping dramatically alters the message of the photographs.
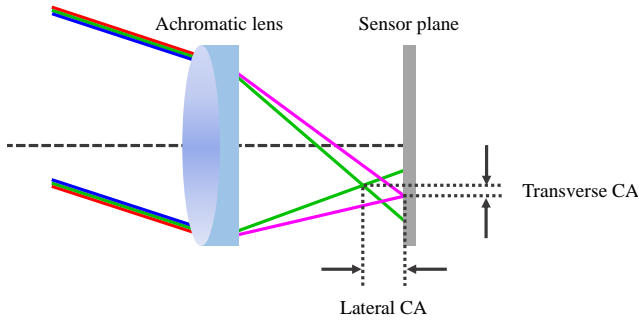
crop feature relied on a saliency prediction network that was racially biased [9].

The guiding question of this paper is to understand the traces left behind from this fundamental operation. What impact does image cropping have on the visual distribution? Can we determine when and how a photo has been cropped?

Despite extensive refinements to the manufacturing process of camera optics and sensors, nearly every modern camera pipeline will leave behind subtle lens artefacts onto the photos that it captures. For example, vignetting is caused by a lens focusing more light at the center of the sensor, creating images that are slightly brighter in the middle than near its borders [32]. Chromatic aberration, also known as purple fringing, is caused by the lens focusing each wave length differently [5]. Since these artefacts are correlated with their spatial position in the image plane, they cause image crops to have trace signatures.

Physical aberrations are not the only traces left behind during the operation. Photographers will prefer to take photos of interesting objects and in canonical poses [45, 4, 19]. Aesthetically pleasing shots will have sensible compositions that respect symmetry and certain ratios in the scene. Violating these principles leaves behind another trace of the cropping operation.

Both of these traces are very subtle, and the human eye often cannot detect them, which makes studying and characterizing them challenging. However, neural networks are

(a) Lens with transverse and longitudinal chromatic aberration. In this illustration, the red and blue channels are aligned (hence the magenta rays), but green-colored light is magnified differently in addition to having a separate in-focus plane.



(b) Close-up of two photos, revealing visible transverse chromatic aberration (TCA) artefacts.

Figure 2: Origin and examples of chromatic aberration.

excellent at identifying these patterns. Indeed, extensive effort goes into preventing neural networks from learning such shortcuts enabled by image crops [13, 38].

In this paper, we flip this around and declare that these shortcuts are not bugs, but instead an opportunity to dissect and understand the subtle clues left behind from image cropping. Capitalizing on a large, high-quality collection of natural images, we train a convolutional neural network to predict the absolute spatial location of a patch within an image. This is only possible if there exist visual features that are *not* spatially invariant. Our experiments analyze the types of features that this model learns, and we show that it is possible to detect traces of the cropping operation. We can also use the detected artefacts, along with semantic information, to recover where the crop was positioned in the original sensor plane.

While the aim of this paper is to analyze the fundamental traces of image cropping, we believe our investigation could have a large practical impact as well. Historically, asking fundamental questions has spurred significant insight into core computer vision problems, such as invariances to scale [8], asymmetries in time [41], and visual chirality [30]. For example, insight into image crops could enable detection of image manipulation attacks, or spur developments to mitigate shortcut learning.

## 2. Background and Related Work

**Optical aberrations.** No imaging device is perfect, and every step in the imaging formation pipeline will leave traces behind onto the final picture. The origins of these signatures range from the physics of light in relation to the camera hardware, to the digital demosaicing and compression algorithms used to store and reconstruct the image. Lenses typically suffer from several aberrations, including chromatic aberration, vignetting, coma, and radial distortion [23, 5, 29, 21]. As shown in Figure 2a, chromatic aberration is manifested in two ways: *transverse* (or *lateral*) chromatic aberration (TCA) refers to the spatial discrepancies in focus points across color channels perpendicular to the optical axis, while *longitudinal* chromatic aberration (LCA) refers to shifts in focus along the optical axis instead [20, 21]. TCA gives rise to color channels that appear to be scaled slightly differently relative to each other, while LCA causes the distance between the focal surface and the lens to be frequency-dependent, such that the degree of blurring varies among color channels. Chromatic aberration can be leveraged to extract depth maps from defocus blur [17, 46, 21], although the spatial sensitivity of these cues is often undesired [13, 37, 38, 36]. To the best of our knowledge, TCA has not yet been regarded as a useful feature in predicting the original position of an image region relative to the lens, or even as a means to expose cropped images whatsoever. Consequently, our contribution is to leverage TCA to fulfill both of these tasks.

**Patch localization.** While one of the first major works in self-supervised representation learning focused on predicting the *relative* location of two patches among eight possible configurations [13], it was also discovered that the ability to perform *absolute* localization seemed to arise out of chromatic aberration. For the best-performing 10% of images, the mean Euclidean distance between the ground truth and predicted positions of single patches is 31% lower than chance, and this gap narrowed to 13% if every image was pre-processed to remove color information along the green-magenta axis. Although there are reasons to believe that modern network architectures might perform better, these rather modest performance figures suggest a priori that the attempted task is a difficult one. Note that the learnability of absolute location is often regarded as a bug; treatments used in practice include random color channel dropping [13], projection [13], grayscale conversion [37, 38], jittering [37], and chroma blurring [36].

**Visual crop detection.** Almost all existing forensics research has centered around 'hard' tampering such as splicing and copy-move operations. We argue that some forms of 'soft' tampering, notably cropping, are also worth investigating. While a few papers have addressed image crops, they are tailored toward specific types of scenes only. For
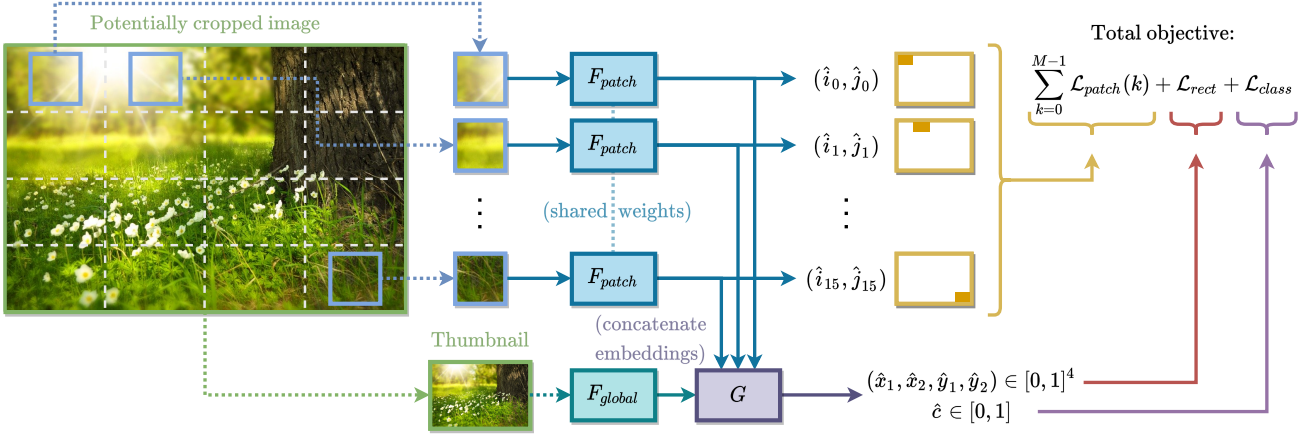
Figure 3: **Full architecture of our crop detection model.** We first extract $M = 16$ patches from the centers of a regularly spaced grid within the source image, a priori not knowing whether it is cropped or not. The patch-based network $F_{patch}$ looks at each patch and classifies its absolute position into one out of 16 possibilities, whereby the estimation is mostly guided by low-level lens artefacts. The global image-based network, $F_{global}$ instead operates on the downscaled source image, and tends to pick up semantic signals, such as objects deviating from their canonical pose (*e.g.* a face is cut in half). Since these two networks complement each other's strengths and weaknesses, we integrate their outputs into one pipeline via the multi-layer perceptron $G$. Note that $F_{patch}$ is supervised by all three loss terms, while $F_{global}$ only controls the crop rectangle $(\hat{x}_1, \hat{x}_2, \hat{y}_1, \hat{y}_2)$ and the final score $\hat{c}$.

example, [7] exploits the regularity of JPEG block structures to reveal cropping, which can easily be bypassed by using another format, or by resizing to dimensions that are multiples of 8 pixels. Both [35] and [15] instead rely heavily on structured image content in the form of vanishing points and lines, which works only if many straight lines (*e.g.* man-made buildings or rooms) are prominently visible. To the best of our knowledge, our approach is the first to tackle unconstrained images, including natural scenes.

## 3. Dataset

The natural clues for detecting crops are subtle, and we need to be careful to preserve them when constructing a dataset. Our underlying dataset has around $700,000$ high-resolution photos from Flickr, which were scraped during the fall of 2019. We impose several constraints on the training images, most importantly that they should not already have been cropped and that they must maintain a constant, fixed aspect ratio and resolution. Appendix A describes this selection and collection process in detail.

We generate image crops by first defining the *crop rectangle* $(x_1, x_2, y_1, y_2) \in [0, 1]^4$ as the relative boundaries of a cropped image within its original camera sensor plane, such that $(x_1, x_2, y_1, y_2) = (0, 1, 0, 1)$ for unmodified images. We always maintain the aspect ratio and pick a random size factor $f$ uniformly in $[0.5, 0.9]$, representing the relative length in pixels of any of the four sides compared to the original photo: $f = x_2 - x_1 = y_2 - y_1$. Every crop must stick to an edge; the rectangle essentially has 1 in 4 chance of touching either the top, right, bottom, or left

border of to the sensor plane.

After cropping exactly half of all incoming photos this way, we intend to give our model access to both patches and global context. We select a patch size of $96 \times 96$, which is large enough to allow the network to get a good idea of the local texture profile, while also being small enough to ensure that neighboring patches never overlap. We then downscale the whole image to a $224 \times 149$ thumbnail, such that it remains accessible to the model in terms of its receptive field and computational efficiency.[1]

Interrelating contextual, semantic information to its spatial position within an image might turn out to be crucial for exposing crops. We therefore add coordinates as two extra channels to the thumbnail, similarly to [31]. Lastly, several shortcut fuzzing procedures had to be used to ensure that the learned features are generalizable; see Appendix B for an extensive description.

## 4. Methodology

We describe our approach and the challenges associated with revealing whether and how a variably-sized single image has been cropped. First, we construct a neural network that can trace image patches back to their original position relative to the center of the lens. Then, we use this novel network to expose and analyze elementary manipulations using an end-to-end trained crop detection model, which

---

[1] The reason we care about receptive field is because, even though high-resolution images are preferable when analyzing subtle lens artefacts, a ResNet-$L$ with $L \leq 50$ has a receptive field of only $\leq 483$ pixels [3], which pushes us to prefer *lower* resolutions instead.

also incorporates the global semantic context of an image in a way that can easily be visualized and understood. Figure 3 illustrates our method.

## 4.1. Predicting absolute patch location

One piece of the puzzle towards analyzing image crops is a neural network called $F_{patch}$, which discriminates the original position of a small image patch with respect to the center of the lens. We frame this as a classification problem for practical purposes, and divide every image into a grid of $4 \times 4$ evenly sized cells, each of which represents a group of possible patch positions. Since this pretext task can be considered to be a form of self-supervised representation learning, with crop detection being the eventual downstream task, we call $F_{patch}$ the *pretext model*.
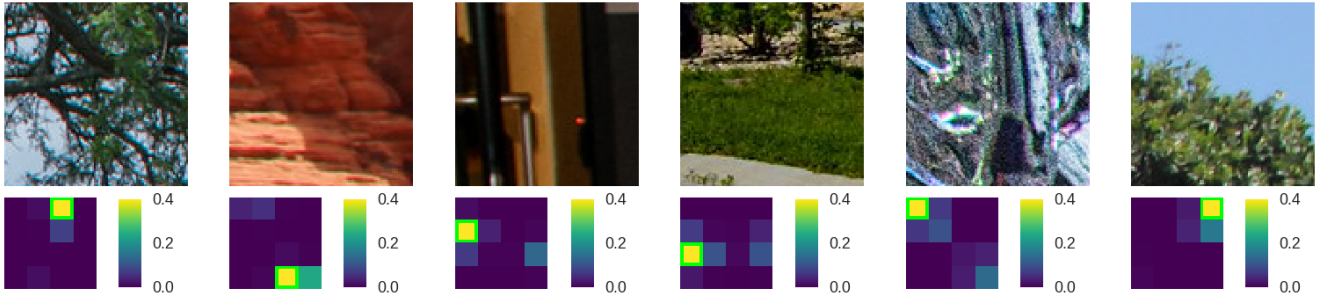
But before embarking on an end-to-end crop detection journey that simply integrates this module into a larger system right from the beginning, it is worth asking the following questions: When exactly does absolute patch localization work well in the first place, and how could it help in exposing manipulated images in an interpretable manner? To this end, we trained $F_{patch}$ in isolation by discarding $F_{global}$ and forcing the network to decide based on information from patches only. The 16-way classification loss term $\mathcal{L}_{patch}$ is responsible for pretext supervision, and is applied onto every patch individually.
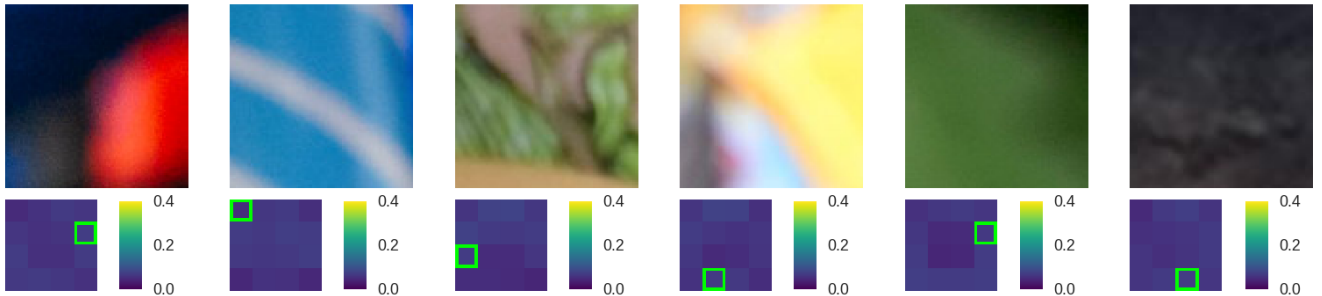
Intriguing patterns emerge when discriminating between different levels of confidence in the predictions produced by $F_{patch}$. Although the accuracy of this localization network is not that high ($\sim 21\%$ versus $\sim 6\%$ for chance) due to the inherent difficulty of the task, Figure 4 shows that it works quite well for some images, particularly those with a high degree of detail coupled with apparent lens artefacts. On the flip side, blurry photos taken with high-end cameras tend to make the model uncertain. This observation suggests that chromatic aberration has strong predictive power for the original locations of patches within pictures. Hence, it is reasonable to expect that incorporating patch-wise, pixel-level cues into a deep learning-based crop detection framework will improve its capabilities.

## 4.2. Architecture and objective

Guided by the design considerations laid out so far, Figure 3 shows our main model architecture. $F_{patch}$ is a ResNet-18 [18] that converts any patch into a length-64 embedding, which then gets converted by a single linear layer on top to a length-16 probability distribution describing the estimated location $(\hat{i}_k, \hat{j}_k) \in \{0 \ldots 3\}^2$ of that patch.



(a) Selecting for **high confidence** yields samples biased toward highly textured content with many edges, often with visible chromatic aberration. The pretext model is typically more accurate in this case.



(b) Selecting for **low confidence** yields blurry or smooth samples, where the lack of detail makes it difficult to expose physical imperfections of the lens. The pretext model tends to be inaccurate in this case.

Figure 4: **Absolute patch localization performance.** By leveraging classification, an uncertainty metric emerges for free. Here, we display examples where the pretext model $F_{patch}$ performs either exceptionally well or badly at recovering the patches' absolute position within the full image. The output probability distribution generated by the network is also plotted as a spatial heatmap (□ = ground truth).

4

$F_{global}$ is a ResNet-34 [18] that converts the downscaled global image into another length-64 embedding. Finally, $G$ is a 3-layer perceptron that accepts a 1088-dimensional concatenation of all previous embeddings, and produces 5 values describing (1) the crop rectangle $(\hat{x}_1, \hat{x}_2, \hat{y}_1, \hat{y}_2) \in [0,1]^4$, and (2) the actual probability $\hat{c}$ that the input image had been cropped. By simultaneously processing and combining aggregated patch-wise information with global context, we allow the network to draw a complete picture of the input, revealing both low-level lens aberrations and high-level semantic cues. The total, weighted loss function is as follows (with $M = 16$):

$$\mathcal{L} = \frac{\lambda_1}{M} \sum_{k=0}^{M-1} \mathcal{L}_{patch}(k) + \frac{\lambda_2}{4}\mathcal{L}_{rect} + \lambda_3 \mathcal{L}_{class} \quad (1)$$

Here, $\mathcal{L}_{patch}(k)$ is a 16-way cross-entropy classification loss between the predicted location distribution $\hat{l}(k)$ of patch $k$ and its ground truth location $l(k)$. For an uncropped image, $l(k) = k$ and $(i_k, j_k) = (k \mod 4, \lfloor k/4 \rfloor)$, although this grid alignment obviously does not necessarily hold for cropped images. Second, the loss term $\mathcal{L}_{rect}$ encourages the estimated crop rectangle to be near the ground truth in a mean squared error sense. Third, $\mathcal{L}_{class}$ is a binary cross-entropy classification loss that trains $\hat{c}$ to state whether or not the photo had been cropped. More formally:

$$\mathcal{L}_{patch}(k) = \mathcal{L}_{CE}(\hat{l}(k), l(k)) \quad (2)$$
$$\mathcal{L}_{rect} = [(\hat{x}_1 - x_1)^2 + (\hat{x}_2 - x_2)^2$$
$$+ (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2] \quad (3)$$
$$\mathcal{L}_{class} = \mathcal{L}_{BCE}(\hat{c}, c) \quad (4)$$

Note that the intermediate outputs $(\hat{i}_k, \hat{j}_k)$ and $(\hat{x}_1, \hat{x}_2, \hat{y}_1, \hat{y}_2)$ exist mainly to encourage a degree of interpretability of the internal representation, rather than to improve the accuracy of the final score $\hat{c}$. Specifically, the linear projection of $F_{patch}$ to $(\hat{i}_k, \hat{j}_k)$ should make the embedding more sensitive to positional information, thus helping the crop rectangle estimation.

### 4.3. Training details

In our experiments, all datasets are generated by cropping exactly 50% of the photos with a random crop factor in $[0.5, 0.9]$. After that, we resize every example to a uniformly random width in $[1024, 2048]$ both during training and testing, such that the image size cannot have any predictive power. We train for up to 25 epochs using an Adam optimizer [24], with a learning rate that drops exponentially from $5 \cdot 10^{-3}$ to $1.5 \cdot 10^{-3}$ at respectively the first and last epoch. The weights of the loss terms are: $\lambda_1 = 2.4$, $\lambda_2 = 3$, and $\lambda_3 = 1$.

## 5. Analysis and Clues

We quantitatively investigate the model in order to dissect and characterize visual crops. We are interested in conducting a careful analysis of what factors the network might be looking at within every image. For ablation study purposes, we distinguish three variants of our model:

- **Joint** is the complete patch- and global-based model from Figure 3 central to this work;
- **Global** is a naive classifier that just operates on the thumbnail, *i.e.* the whole input downscaled to $224 \times 149$, using $F_{global}$;
- **Patch** only sees 16 small patches extracted from consistent positions within the image, using $F_{patch}$.
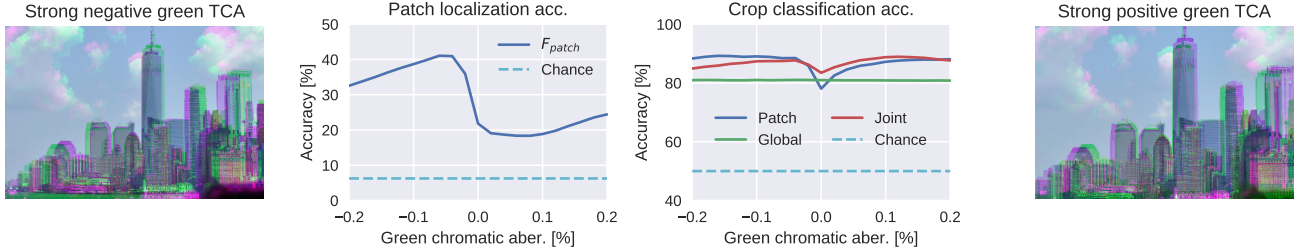
We classify the information that a model uses as evidence for its decision into two broad categories: **(1) characteristics of the camera or lens system**, and **(2) object priors**. While (1) is largely invariant of semantic image content, (2) could mean that the network has learned to leverage certain rules in photography, *e.g.* the sky is usually on top, and a person's face is usually centered.

To gain insight into what exactly our model has discovered, we first investigate the network's response to several known lens characteristics by artificially inflating their corresponding optical aberrations on the test set, and computing the resulting performance metrics. Next, we measure the changes in accuracy when the model is applied on datasets that were crafted specifically as to have divergent distributions over object semantics and image structure. We expect both lens flaws and photographic conventions to play different but interesting roles in our model.
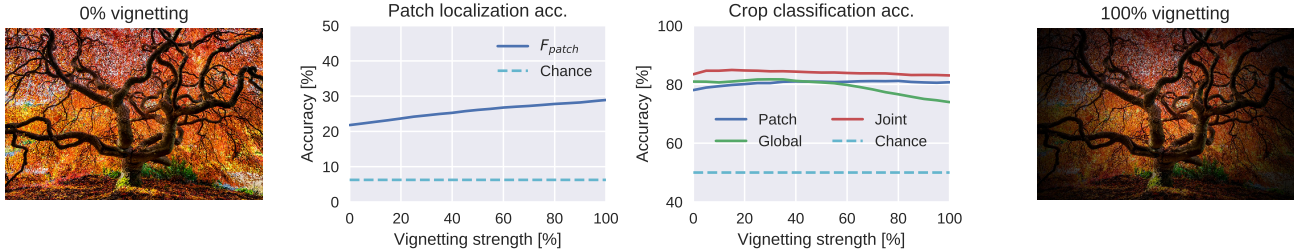
A discussion of chromatic aberration expressed along the green channel, vignetting, and photography patterns follows; see Appendix C for the effect of color saturation, radial lens distortion, and chromatic aberration of the red and blue channels.

### 5.1. Effect of chromatic aberration

A common lens correction to counter the frequency-dependence of the refractive index of glass is to use a so-called *achromatic doublet*. This modification ensures that the light rays of two different frequencies, such as the red and blue color channels, are aligned [23]. Because the remaining green channel still suffers from TCA and will therefore be slightly downscaled around the optical center, this artefact is often visible as green or purple fringes near edges and other regions with contrast or texture [6]. Figure 2b depicts real examples of what chromatic aberration looks like. Note that the optical center around which radial magnification occurs does not necessarily coincide with the image center due to the complexity of multi-lens systems [50], although both points have been found to be

5

(a) **Green transverse chromatic aberration** in the negative (inward) direction considerably boosts performance for patch localization, although asymmetry is key for crop detection. The *global* model remains unaffected since it is unlikely to be able to see the artefacts. (We show examples with excessive distortion for illustration; the range used in practice is much more modest.)



(b) **Vignetting** also contributes positively to the pretext model's accuracy. Interestingly, the crop detection performance initially increases but then drops slightly for strong vignetting, presumably because the distorted images are moving out-of-distribution.

Figure 5: **Breakdown of image attributes that contribute to features relevant for crop detection.** In these experiments, we manually exaggerate two characteristics of the lens on 3,500 photos of the test set, and subsequently measure the resulting shift in performance.

very close in practice [20]. Furthermore, chromatic aberration can vary strongly from device to device, and is not even present in all camera systems. Many high-end, modern lenses and/or post-processing algorithms tend to accurately correct for them, to the point that it becomes virtually imperceptible.

Nonetheless, our model still finds this spectral discrepancy in focus points to be a distinctive feature of crops and patch positions: Figure 5a (left plot) demonstrates that artificially downscaling the green channel significantly improves the pretext model's performance. This is because the angle and magnitude of texture shifts across color channels can give away the location of a patch relative to the center of the lens. Moreover, the downstream task of crop detection (right plot) becomes easier when TCA is introduced in either direction. Horizontally mirrored plots were obtained upon examining the red and blue channels, confirming that the green channel suffers an inward deviation most commonly of all in our dataset. It turns out that the optimal configuration from the perspective of $F_{patch}$ is to add a little distortion, but not too much — otherwise we risk hurting the realism of the test set.

## 5.2. Effect of vignetting

A typical imperfection of multi-lens systems is the radial brightness fall-off as we move away from the center of the image, seen in Figure 5b. Vignetting can arise due to mechanical and natural reasons [32], but its dependence on the position within a photo is the most important aspect in this context. We simulate vignetting by multiplying every pixel value with $\frac{1}{g(r)}$, where:

$$g(r) = 1 + ar^2 + br^4 + cr^6 \qquad (5)$$
$$(a, b, c) = (2.0625, 8.75, 0.0313) \qquad (6)$$

$g(r)$ is a sixth-grade polynomial gain function, the parameters $a, b, c$ are assigned typical values taken from [32], and $r$ represents the radius from the image center with $r = 1$ at every corner. The degree of vignetting is smoothly varied by simply interpolating every pixel between its original (0%) and fully modified (100%) state.

Figure 5b shows that enhanced vignetting has a positive impact on absolute patch localization ability, but this does not appear to translate into noticeably better crop detection accuracy. While the gradient direction of the brightness across a patch is a clear indicator of the angle that it makes with respect to the optical center of the image, modern cameras appear to correct for vignetting well enough such that the lack of realism of the perturbed images hurts $F_{global}$'s performance more so than it helps.

6

| Flickr | Flickr (no humans) | Upright | Tilted | Vanishing points | Texture | Smooth |

Figure 6: **Representative examples of the seven test sets.** The first two are variants of Flickr, one unfiltered and one without humans or faces, and the remaining five are custom photo collections we intend to measure various other kinds of photographic patterns or biases with. These were taken at various locations throughout New York, Boston, and San Francisco Bay Area, and every category contains between 18 and 148 pictures.

| Dataset | Joint | Global | Patch | Human |
|---|---|---|---|---|
| Flickr | 86% | 79% | 77% | 67% |
| Flickr (no humans) | 81% | 75% | 73% | - |
| Upright | 78% | 68% | 73% | - |
| Tilted | 67% | 59% | 69% | - |
| Vanish | 83% | 80% | 73% | - |
| Texture | 61% | 50% | 59% | - |
| Smooth | 50% | 52% | 50% | - |

Table 1: **Accuracy comparison between three different crop detection models on various datasets.** All models are trained on Flickr, and appear to have discovered common rules in photography to varying degrees. These results reveal intriguing and sometimes surprising facts about the visual distributions that we work with almost every day.

## 5.3. Effect of photography patterns and perspective

The desire to capture meaningful content implies that not all images are created equal. Interesting objects, persons, or animals will often intentionally be centered within a photo, and cameras are generally oriented upright when taking pictures. Some conventions, *e.g.* grass is usually at the bottom, are confounded to some extent by the random portrait-to-landscape rotations during training, although there remain many facts to be learned as to what constitutes an appealing or sensible photograph. One clear example of these so-called *photography patterns* in the context of our model is that when a person's face that is cut in half, this might reveal that the image had been cropped. This is because, intuitively speaking, it does not conform to how photographers typically organize their visual environment and constituents of the scene.

The structure of the world around us not only provides high-level knowledge on where and how objects typically exist within pictures, but also gives rise to perspective cues,

for example the angle that horizontal lines make with vertical lines upon projection of a 3D scene onto the 2D sensor, coupled with the apparent normal vector of a wall or other surface. Measuring the exact extent to which all of these aspects play a role is difficult, as no suitable dataset exists. The ideal baseline would consist of photos without any adherence to photography rules whatsoever, taken in uniformly random orientations at arbitrary, mostly uninteresting locations around the world.

We constructed and categorized a small-scale collection of such photos ourselves, using the Samsung Galaxy S8 and Google Pixel 4 smartphones, spanning the 5 right-most columns in Figure 6. Columns 3 and 5 depict photos that are taken with the camera in an upright, biased orientation. Column 5 specifically encompasses vanishing point-heavy content, where perspective clues may provide clear pointers. Columns 4, 6, and 7 contain pictures that are unlikely to be taken by a normal photographer, but whose purpose is instead to measure the response of our system on photos with compositions that make less sense.

Quantitative results are shown in Table 1. On the Flickr test set, the crop classification accuracy is 79% for the thumbnail-based model, 77% for the patch-based model, and 86% for the joint model. For comparison purposes, we also asked 16 people to classify 100 random Flickr photos into whether they look cropped or not, resulting in a human accuracy of 67%. This demonstrates that integrating information across multiple scales results in a better model than a network that only sees either patches or thumbnails independently, in addition to having a significant performance margin over humans.

Our measurements also indicate that the model tends to consistently perform better on sensible, upright photos. Flickr in particular seems to exhibit a high degree of photographic conventions involving persons, so we also tested a manually filtered subset of 100 photos that do not contain humans or faces, resulting in a modest drop in accuracy.
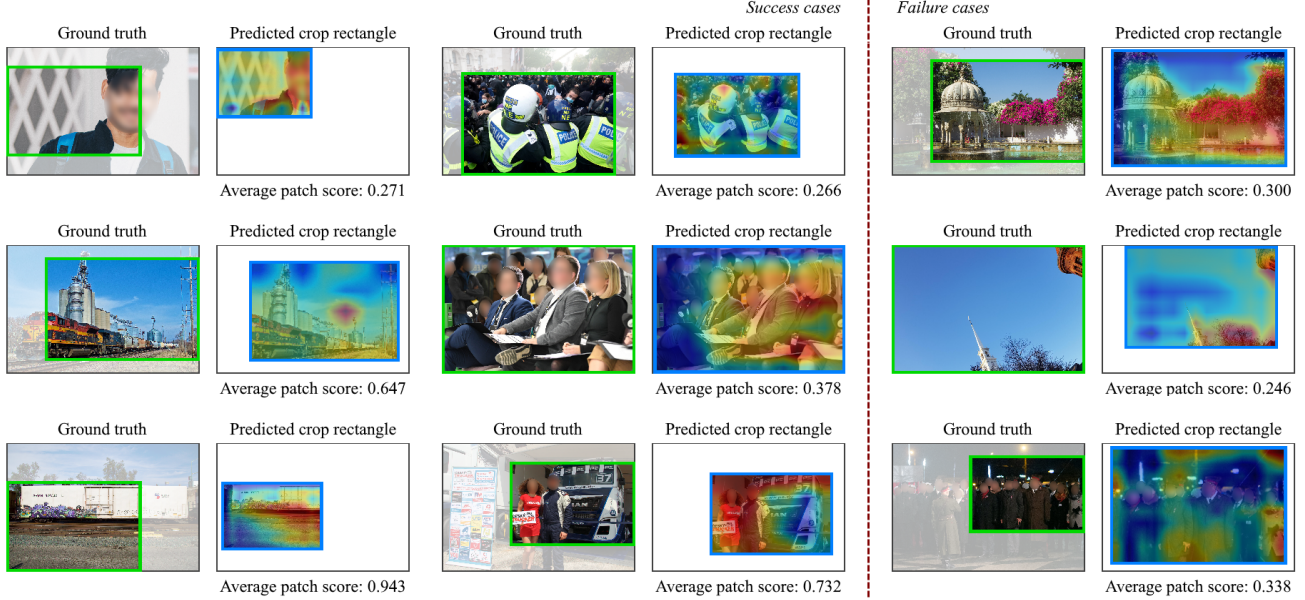
Figure 7: **Qualitative examples and interpretation of our crop detection system.** High-level cues such as persons and faces appear to considerably affect the model's decisions. Note that images don't always *look* cropped, but in that case, patches can act as the giveaway whenever they express lens artefacts. Regardless, certain scene compositions are more difficult to get right, such as in the failure cases shown on the right. (Faced blurred here for privacy protection.)

Interestingly, the patch-based model outperforms all other networks on *tilted*, suggesting that global context can sometimes substantially confuse the model if the photo is taken in an abnormal way. Fully textured or white-wall images appear to be even more out-of-distribution. However, most natural imagery predominantly contains canonical and appealing arrangements, where our model displays a stronger ability to distinguish crops.

# 6. Visualizing Image Crops

## 6.1. Embeddings

In order to depict the changing visual distribution as images are cropped to an increasingly stronger extent, we look at the output embeddings produced by the thumbnail network $F_{global}$. We first apply Principal Component Analysis (PCA) to transform the data points from 64 to 24 dimensions, and subsequently apply t-SNE [33] to further reduce the dimensionality from 24 to 2. The result is shown in Figure 8.

## 6.2. Attribution

As discussed in the previous sections, there could be many reasons as to why the model predicts that a certain photo appears or does not appear to be cropped. However, to explain results obtained from any given single input, we can also apply the Grad-CAM technique [43] onto the global image. This procedure allows us to construct a
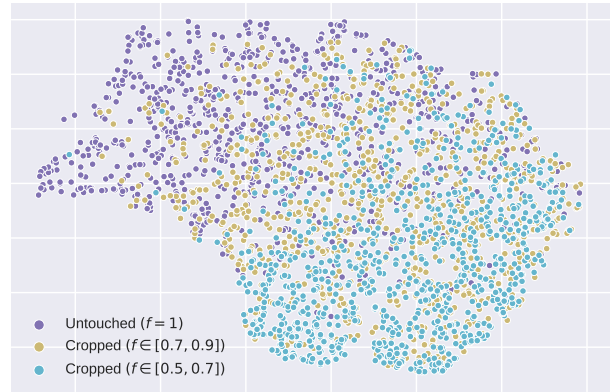


Figure 8: **Dimensionality-reduced embeddings generated by** $F_{global}$ **on Flickr.** Here, the size factor $f$ stands for the fraction of one cropped image dimension relative to the original photo. The model is clearly able to separate untampered from strongly cropped images, although lightly cropped images can land almost anywhere across the spectrum as the semantic signals might be less pronounced and/or less frequently present.

heatmap that attributes decisions made by $F_{global}$ and $G$ back to the input regions that contributed to them. In addition, relying on the earlier observation that the pretext model $F_{patch}$ usually knows when it is correct and when it is not, we note the average output score of the top 4 patches. We investigate and showcase a few examples in Figure 7.

Here, since we have access to the original images, we create cropped variants (by the green ground truth rectangle) and feed them into the network to visualize its prediction. The model is often able to *uncrop* the image, using semantic and/or patch-based clues, and produce a reasonable estimate of which spatial regions are missing (if any). For example, the top left image clearly violates routine principles in photography. The top or bottom images are a little harder to judge by the same measure, though we can still recover the crop frame thanks to the absolute patch localization functionality.

## 7. Discussion

We found that image regions clearly contain information about their spatial position relative to the lens, breaking established assumptions about translational invariance [27]. Our network has automatically discovered various relevant clues, ranging from subtle lens flaws to photographic priors. These features are likely to be acquired to some extent by many self-supervised representation learning methods, such as contrastive learning, where cropping is an important form of data augmentation [11]. Although they are often treated as a bug, there are also compelling cases where the clues could prove to be useful. For example, we believe that our crop detection and analysis framework has significant implications for digital forensics, especially in the increasingly relevant and important fight against misleading or fake news. We also hope that our work inspires further research into how the traces left behind by image cropping, and the altered distributions over images that it gives rise to, can be leveraged in other interesting ways.

## References

[1] Opencv: Geometric image transformations. 12

[2] Alexander Amini and Ava Soleimany. Mit 6.s191: Introduction to deep learning, spring 2020. 1

[3] André Araujo, Wade Norris, and Jack Sim. Computing receptive fields of convolutional neural networks. *Distill*, 2019. https://distill.pub/2019/computing-receptive-fields. 3

[4] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In *Advances in Neural Information Processing Systems*, pages 9453–9463, 2019. 1

[5] Steven Beeson and James W Mayer. *Patterns of light: chasing the spectrum from Aristotle to LEDs*. Springer Science & Business Media, 2007. 1, 2

[6] David Brewster and Alexander Dallas Bache. *A Treatise on Optics...: First American Edition, with an Appendix, Containing an Elementary View of the Application of Analysis to Reflexion and Refraction*. Carey, Lea, & Blanchard, 1833. 5

[7] AR Bruna, Giuseppe Messina, and Sebastiano Battiato. Crop detection through blocking artefacts analysis. In *International Conference on Image Analysis and Processing*, pages 650–659. Springer, 2011. 3

[8] Peter Burt and Edward Adelson. The laplacian pyramid as a compact image code. *IEEE Transactions on communications*, 31(4):532–540, 1983. 2

[9] Katie Canales. Twitter is making changes to its photo software after people online found it was automatically cropping out black faces and focusing on white ones, Oct 2020. 1

[10] Jiansheng Chen, Gaocheng Bai, Shaoheng Liang, and Zhengqin Li. Automatic image cropping: A computational complexity study. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 507–515, 2016. 1

[11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020. 1, 9

[12] London Broadcasting Company. Bbc criticised for cropping out weapon in black lives matter protest photo, Jun 2020. 1

[13] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 2, 12

[14] Sahar Esfandiari and Will Martin. Greta thunberg slammed the associated press for cropping a black activist out of a photo of her at davos, Jan 2020. 1

[15] Marco Fanfani, Massimo Iuliani, Fabio Bellavia, Carlo Colombo, and Alessandro Piva. A vision-based fully automated approach to robust image cropping detection. *Signal Processing: Image Communication*, 80:115629, 2020. 3

[16] Alex Franz and Thorsten Brants. All our n-gram are belong to you, Aug 2006. 11

[17] Josep Garcia, Juan Maria Sanchez, Xavier Orriols, and Xavier Binefa. Chromatic aberration and depth extraction. In *Proceedings 15th International Conference on Pattern Recognition. ICPR-2000*, volume 1, pages 762–765. IEEE, 2000. 2

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4, 5

[19] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. *arXiv preprint arXiv:1907.07174*, 2019. 1

[20] Sing Bing Kang. Automatic removal of chromatic aberration from a single image. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007. 2, 6

[21] Masako Kashiwagi, Nao Mishima, Tatsuo Kozakaya, and Shinsaku Hiura. Deep depth from aberration map. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4070–4079, 2019. 2

[22] Josh Kaufman. Github: first20hours/google-10000-english, Aug 2019. 11

[23] Michael J Kidger. Fundamental optical design. In *Fundamental optical design*. SPIE Bellingham, 2001. 2, 5

[24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[25] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. 1

[26] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Tom Duerig, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982*, 2018. 11

[27] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015. 9

[28] Fei-Fei Li, Ranjay Krishna, and Danfei Xu. Stanford cs231n: Convolutional neural networks for visual recognition, spring 2020. 1

[29] Xufeng Lin and Chang-Tsun Li. Image provenance inference through content-based device fingerprint analysis. In *Information Security: Foundations, Technologies and Applications*, pages 279–310. IET, 2018. 2, 15

[30] Zhiqiu Lin, Jin Sun, Abe Davis, and Noah Snavely. Visual chirality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12295–12303, 2020. 2

[31] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. In *Advances in Neural Information Processing Systems*, pages 9605–9616, 2018. 3

[32] Laura Lopez-Fuentes, Gabriel Oliver, and Sebastia Massanet. Revisiting image vignetting correction by constrained minimization of log-intensity entropy. In *International Work-Conference on Artificial Neural Networks*, pages 450–463. Springer, 2015. 1, 6

[33] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008. 8

[34] Francesco Marra, Diego Gragnaniello, Luisa Verdoliva, and Giovanni Poggi. A full-image full-resolution end-to-end-trainable cnn framework for image forgery detection. *IEEE Access*, 8:133488–133502, 2020. 11

[35] Xianzhe Meng, Shaozhang Niu, Ru Yan, and Yezhou Li. Detecting photographic cropping based on vanishing points. *Chinese Journal of Electronics*, 22(2):369–372, 2013. 3

[36] T Nathan Mundhenk, Daniel Ho, and Barry Y Chen. Improvements to context based self-supervised learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9339–9348, 2018. 2

[37] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016. 2

[38] Mehdi Noroozi, Hamed Pirsiavash, and Paolo Favaro. Representation learning by learning to count. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5898–5906, 2017. 2, 12

[39] Aude Oliva and Antonio Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International journal of computer vision*, 42(3):145–175, 2001. 1

[40] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 1

[41] Lyndsey C Pickup, Zheng Pan, Donglai Wei, YiChang Shih, Changshui Zhang, Andrew Zisserman, Bernhard Scholkopf, and William T Freeman. Seeing the arrow of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2035–2042, 2014. 2

[42] A Samii, R Měch, and Zhe Lin. Data-driven automatic cropping using semantic composition search. In *Computer graphics forum*, volume 34, pages 141–151. Wiley Online Library, 2015. 1

[43] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 8

[44] James Tompkin. Brown cs231n: Csci 1430: Introduction to computer vision, spring 2020. 1

[45] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528. IEEE, 2011. 1

[46] Pauline Trouvé, Frédéric Champagnat, Guy Le Besnerais, Jacques Sabater, Thierry Avignon, and Jérôme Idier. Passive depth estimation using chromatic aberration and a depth from defocus approach. *Applied optics*, 52(29):7152–7164, 2013. 2

[47] Basile Van Hoorick. Image outpainting and harmonization using generative adversarial networks. *arXiv preprint arXiv:1912.10960*, 2019. 1

[48] Todd Vorenkamp. Understanding crop factor, 2016. 11

[49] Yi Wang, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Wide-context semantic image extrapolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1399–1408, 2019. 1

[50] Reg G Willson and Steven A Shafer. What is the center of the image? *JOSA A*, 11(11):2946–2955, 1994. 5

## Supplementary material

## A. Dataset constraints, collection, and description

### A.1 Lack of cropping

As a starting point, any sufficiently large collection of untampered photos suffices. In order to simulate scenarios where a user only has access to pixels but not the metadata (which commonly happens when downloading photos from *e.g.* social media), no labels are needed. Training and testing data can be retrieved 'for free' by extracting patches and thumbnails from any dataset consisting of real-world images, where the only important constraint is the lack of tampering. However, it turns out that cropping, as well as various other kinds of 'soft tampering', is a natural part of the digital editing process. Because these operations are mostly harmless and probably happen more often than we realize, it becomes almost impossible to know to what extent a given database really is untampered.

### A.2 Sufficiently high resolution

Acknowledging the fact that the dataset might be noisy to some degree, we proceed with adding a resolution constraint. Image datasets for deep learning are often down-scaled such that the maximal dimension lies around 500 to 1,000 pixels[2], presumably because the benefit of an even finer level of detail for recognizing object semantics rarely outweighs the extra computational cost. However, in order to better pick up lens flaws that are typically exhibited in subtle pixel-level features, we prefer to keep the resolution higher and closer to the original photo. This matches the observation in image forensics that resizing should be avoided because it tends to damage high-frequency details [34]. We decided to settle for (*i.e.* download images with) a maximal dimension of 2,048 pixels for each sample, which is deemed high enough to detect optical imperfections, but also low enough to avoid exceeding realistic dimensions of photos that may be shared online.

### A.3 Inter-device variation considerations

Every lens and sensor is different, and this variation in standards might make what exactly constitutes a 'crop' less precise. For example, if a full-frame lens is coupled with a crop sensor (*i.e.* the film frame width is less than 35mm) as in Figure 9, every resulting picture can be thought of

---
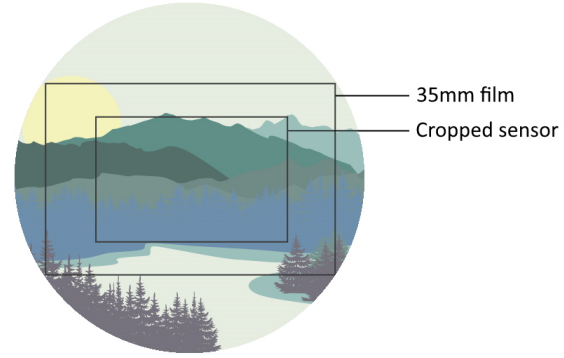[2]For example, every sample in Open Images V5 [26] has at most 1,024 pixels on its longest side.



Figure 9: Comparison of a full-frame sensor versus a crop sensor with respect to the lens circle. (Adjusted and reprinted from [48] with permission.)

as inherently cropped because the light captured by the sensor does not fully cover the lens circle. Mobile phones have an especially large crop factor, since their sensors are typically much smaller than those used in professional DSLR camera systems. In fact, there is a vast number of possible configurations, and trying to take all of them into account would become impracticable. We thus clear confusion by defining a 'cropped image' to be any deviation from what was originally captured by the imaging sensor at the time of shooting. Since our method is camera make and model-blind, we rely on the learning-based approach to discover modal values within this combinatorial space of configurations in the dataset, such that our network will learn to take the diversity among devices and settings into account automatically.

### A.4 Scraping and dataset bias

We scraped Flickr by querying the API with 10,000 different search terms and downloading up to 500 photos for every tag. The keywords were gathered from an online list of the 10,000 most commonly used words in English, which was in turn generated by performing N-gram frequency analysis on the Google Trillion Word Corpus [16, 22]. The resulting database has around 1.3 million images, which would have been 5 million if the search results did not overlap due to many entries having multiple tags. Note that Flickr seems to be biased toward photos (1) depicting persons, (2) of somewhat professional quality, and (3) taken using expensive cameras, but we view neither of these aspects as a drawback considering the relevance of our project to photography patterns and even photojournalism.
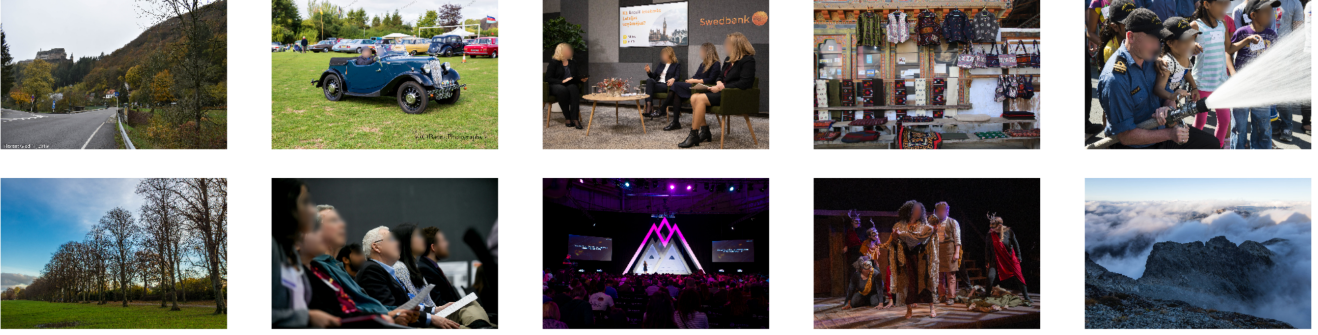
Figure 10: Random examples of the Flickr dataset. (Faced blurred for privacy protection.)

## A.5 Aspect ratio

Mixing training examples with different aspect ratios together also changes the shapes of the grid cells of $F_{patch}$, which should clearly be avoided. Otherwise, patches that have the same absolute position with respect to the lens circle, might be assigned different labels depending on the aspect ratio of the sensor within said lens circle. Most digital camera systems have a sensor size of 36mm×24mm, corresponding to an aspect ratio of 1.5. We therefore fix the aspect ratio to 1.5 and enforce landscape-only photos (by rotating portrait images either left or right) to further enhance consistency, which shrinks the pool of files meeting all discussed criteria down to 700,000 files.

## A.6 Dataset split

Lastly, we perform a 3-way train / validation / test set split distributed as 90% / 5% / 5%. A few samples of the test set are shown in Figure 10.

## B. Shortcut mitigation

Convolutional neural networks have been shown to be surprisingly adept at finding and leveraging often irrelevant shortcuts [13, 38]. Here, we present our approach to ensure that the models learn useful features.

## B.1 Image patch extraction

Patches are extracted from the centers of a regularly sized $4 \times 4$ grid within every image (cropped or not), but we also apply random jittering of $\pm 8$ pixels in both dimensions. This way, we discourage $F_{patch}$ from learning low-level image processing-related shortcuts, for example JPEG block artefact alignment.

## B.2 Resizing global images

Since $F_{global}$ uses a downscaled variant of the incoming image with fixed dimensionality $224 \times 149$, but cropping an image also changes its raw dimensions, we were obliged to employ some tricks in order to prevent the model from learning glitches that are unrelated to physical imaging aberrations, notably resampling factor detection. Resampling shortcuts have occurred in various previous works [38], and are typically an undesired factor. For example, a neural network is able to trivially distinguish images that have been downsized starting from $2048 \times 1365$ as opposed to starting from $1536 \times 1024$ based on pixel-level resampling artefacts, even if the interpolation method is randomized [38]. To work around this issue, we perform random resizing in multiple stages to make the original dimensions nearly impossible to recover, without noticeably damaging the image contents.

Given the potentially cropped source image of size $W \times H$, we first resize 3 times to a random $W' \times H'$ where $W'$ is uniformly distributed in $[1024, 2048]$, and $H'|W'$ is conditionally uniformly distributed in $[0.8W'/A, 1.2W'/A]$, with $A$ the aspect ratio. Note that the interpolation method itself is also random, and is chosen from one of {NEAREST, LINEAR, AREA, CUBIC, LANCZOS4} as provided by the OpenCV library [1]. Finally, the whole image is downscaled to $224 \times 149$, and from now on it should be nearly impossible to tell what its original resolution was.

Indeed, if we replace the cropping operation with a rescaling to the same dimensions that the cropped image would otherwise have, the accuracy of our global model drops to chance (50%). This suggests that only altered image contents play a role, while input resolution does not anymore.

Note that the way in which patches are extracted remains unaltered by this procedure; only thumbnails must be treated to ensure that $F_{global}$ predominantly looks at semantically meaningful content.

## B.3 Joint model

Another, more sophisticated shortcut arose which occurs only when the model has access to both patches and thumb-
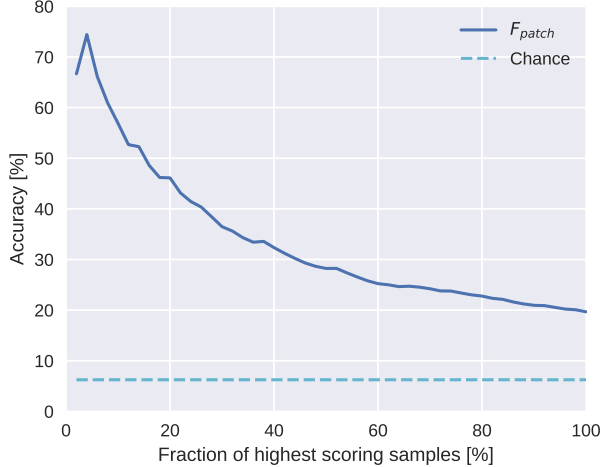
Figure 11: **Sample selectivity versus patch localization performance.** The accuracy improves significantly once we discard more and more predictions that $F_{patch}$ is uncertain about.
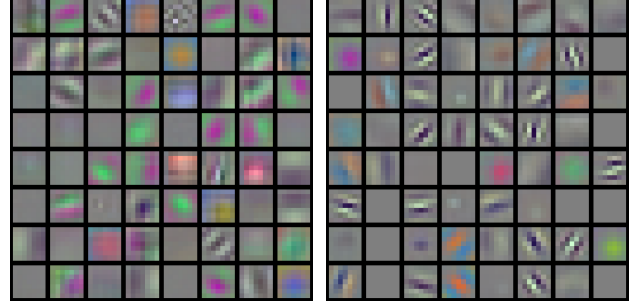
nails simultaneously. Even if the original dimensions of a global image cannot be inferred, the integrated network could still learn to measure how 'large' the patches are in comparison to the thumbnail, since they are extracted from a 'smaller' image if the input is cropped. To alleviate this issue, we perform an extra random resizing step *before* extracting patches but *after* cropping, where the width is uniformly distributed in $[1024, 2048]$ and the height is chosen proportionally such that the aspect ratio is retained. This guarantees that the fraction of the thumbnail that is being covered by patches loses its predictive power, discouraging $G$ from trying to exploit low-level correlations among the outputs produced by $F_{patch}$ and $F_{global}$. This approach serves the additional purpose of enforcing our ignorance about both the crop rectangle and the sequence of resizes that images at test time could have undergone; hence, during our evaluations, we also randomize input resolutions the same way.

## C. Patch localization accuracy versus confidence

Figure 11 plots the accuracy of $F_{patch}$ as a function of the response rate, where moving to the left on the horizontal axis means that an increasingly smaller fraction of only the patches with the highest scores are considered. This supports the earlier claim that the maximum value in the output distribution correlates positively with the correctness of the pretext model.

## D. Convolutional filter visualization

We display and compare the values of the convolution operations applied by the very first layers of both $F_{patch}$ and a regular ImageNet classifier in Figure 12.



(a) $F_{patch}$ (ResNet-18).  (b) ImageNet-trained ResNet-34.

Figure 12: **First convolutional layer filter visualization.** At the lowest level, the absolute patch localization model is clearly more sensitive to alternations between green and magenta (*i.e.* lack of green) pixel values in various directions, as compared to a vanilla ImageNet-trained neural network.

| Model | Color | Grayscale | Chance |
|---|---|---|---|
| $F_{patch}$ (patch loc.) | 21% | 15% | 6% |
| Joint (crop det.) | 86% | 81% | 50% |
| Global (crop det.) | 79% | 78% | 50% |
| Patch (crop det.) | 77% | 72% | 50% |

Table 2: **Accuracies with or without color.** Removing all color information on the test set decreases the model's performance, but only considerably so when a model relies on patches.

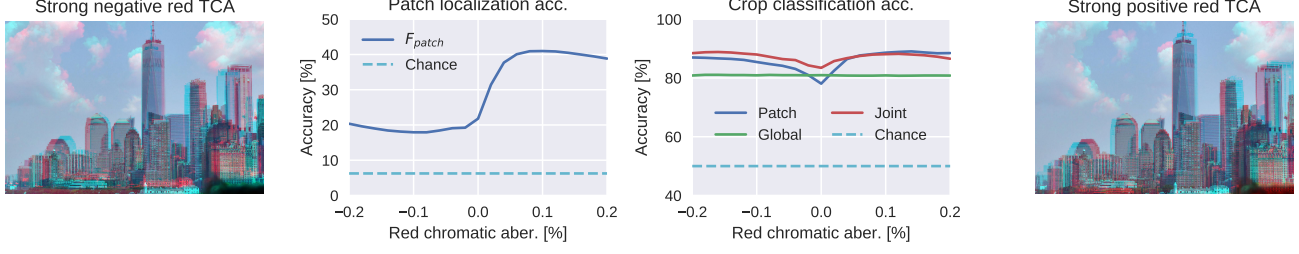## E. Additional experiments for lens-related clues

### E.1 Effect of red and blue chromatic aberration

As shown in Figures 13a and 13b, the patch localization accuracy plots appear horizontally flipped with respect to Figure 5a. This indicates that the modal value of purple fringing in our dataset corresponds to the green channel being scaled toward one preferred direction more often than in the other direction. (Inward green TCA is visually the same as a combination of outward red and blue TCA.)
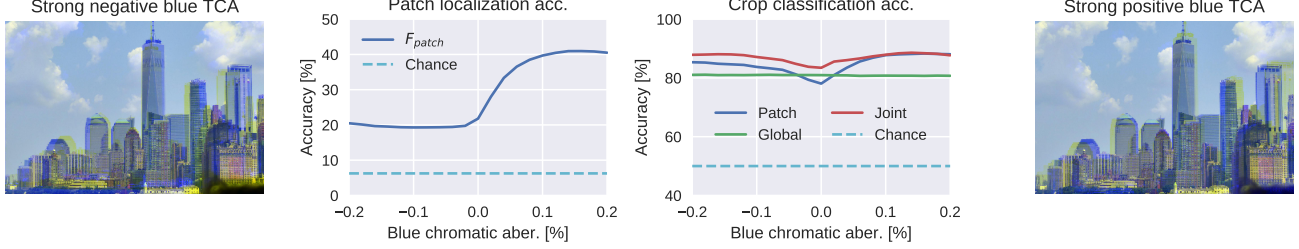
### E.2 Effect of color saturation and grayscale

In order to quantify the significance of color information in general beyond just chromatic aberration, it may be instructive to control the saturation of the test set. A saturation factor of 0% is equivalent to grayscale imagery, 100% is identity, and larger numbers represent exaggerated colors. The result is shown in Figure 13c. This feature does not depend on the location of a patch, therefore it is not unexpected that the best performance corresponds with untampered images. Any other value simply moves the images away from the expected distribution.
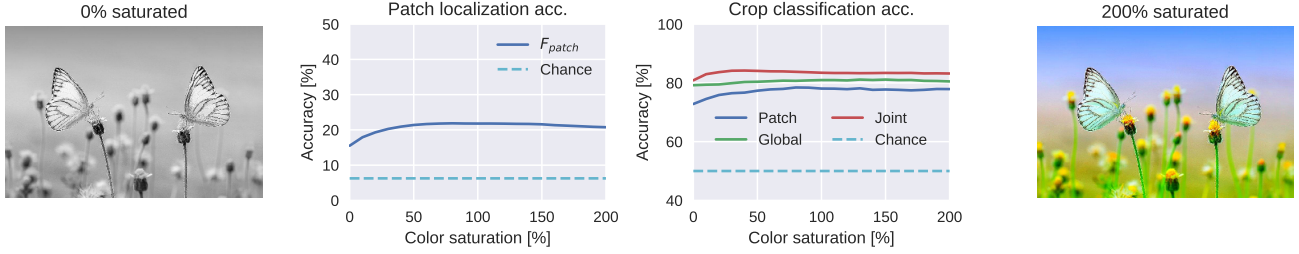
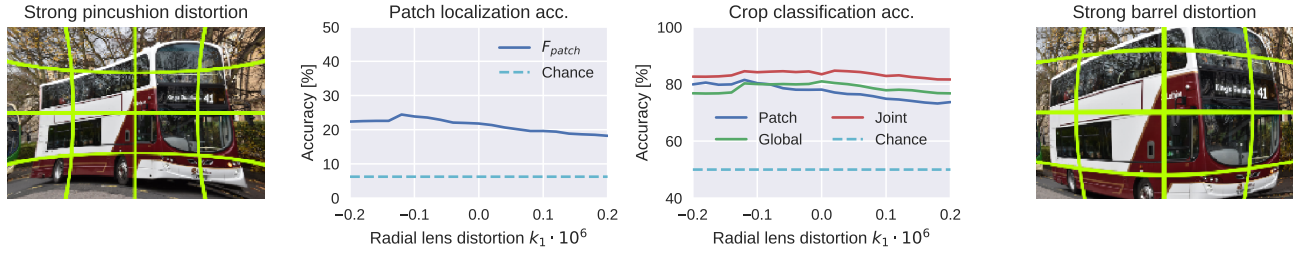Table 2 also compares the performance of the model

(a) **Red transverse chromatic aberration** in the positive (outward) direction boosts performance.



(b) **Blue transverse chromatic aberration** in the positive (outward) direction boosts performance.



(c) Adjusting **color saturation** away from 100% (= identity) slightly degrades performance.



(d) The degree of **radial lens distortion** in our dataset may be too subtle to substantially affect the integrated crop detection model, although due to the noisy results, this is inconclusive.

Figure 13: **Extended breakdown of image attributes.** See Figure 5 for the main results.

when tested on grayscale and regular color images. Although color information clearly constitutes a respectable gain to the network's correctness relative to chance levels, there is a large residual gap that does not rely on color. Apart from vignetting, we hypothesize this is mostly related to photography patterns and object priors, which we discussed in Section 5.3. Moreover, the only model that is

likely unable to perceive lens aberrations in the first place (*global*) seems to care the least about color information, suggesting that the object priors involved in revealing crops can be learned with minimal dependence on color.

### E.3 Effect of radial lens distortion

Pincushion or barrel distortion, illustrated left and right respectively in Figure 13d, arises from the fact that the magnification of a scene through a lens does not stay constant across the image plane, but depends on the radius $r = \sqrt{x^2 + y^2}$ from the optical center [29]. We replicate this distortion by applying a geometric coordinate transformation with a simple square law that scales every destination pixel $(x_d, y_d)$ relative to its source $(x_s, y_s)$ as follows:

$$d = 1 + k_1 r^2 \tag{7}$$
$$(x_d, y_d) = (dx_s, dy_s) \tag{8}$$

Figure 13d shows the effect of inflating lens distortion on the test set according to Equation (7-8).