

# Item Popularity Prediction in E-commerce Using Image Quality Feature Vectors

Stephen Zakrewsky  
Drexel University  
sz372@drexel.edu

Kamelia Aryafar  
Etsy  
karyafar@etsy.com

Ali Shokoufandeh  
Drexel University  
ashokouf@cs.drexel.edu

**Abstract**—Online retail is a visual experience- Shoppers often use images as the first order information to decide if an item matches their taste. Image characteristics such as color, scene composition, texture, style, aesthetics and overall quality play a crucial role in making a purchase decision, clicking on or liking a product listing. In this paper we use a set of image features that indicate quality to predict product listings popularity on a major e-commerce website, Etsy<sup>1</sup>. We first define listing popularity through search clicks, favoriting and purchase activity. Next, we infer listing quality from the pixel-level information of listed images as quality features. We then compare our findings to text-only models for popularity prediction. Our initial results indicate that a combined image and text modeling of product listings on Etsy outperforms text-only models in popularity prediction.

## I. INTRODUCTION

The informative presentation of product listings through text and images is the foundation of modern e-commerce. Shoppers often have a specific style or visual preference for many of the available items such as jewelry, clothing, home decor, etc. Images provide the first order information for product listings. Users usually use images in combination with other data modalities such as textual description, price, ratings and etc. to decide if an item is a suitable match for what they need and have in mind. The selection of proper high quality images is then an important step in listing a successful product. In this paper we examine the role of image quality in a listing popularity on a major e-commerce website, Etsy<sup>1</sup>.

Etsy is an online marketplace for artisans selling unique handcrafted goods, and vintage wares that couldn't be found elsewhere. Etsy caters to the long tail of online retail [1], [2]. With more than one million sellers, 35 million unique product listings and nearly a hundred million images, Etsy is uniquely positioned to answer some interesting questions about the role of images as a rich visual experience in an e-commerce setting. Each Etsy listing is composed of text information such as title, tags, item description, shop and seller name and complementary images. For a product listing to stand out, high-quality images

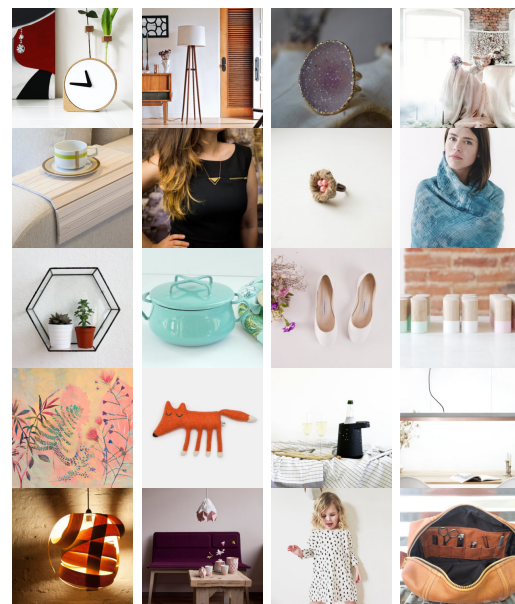


Fig. 1: Sample Etsy listing images are shown with different lighting, scene composition, and quality.

describing the content of the product listing is a necessity [3], [4].

In this paper we introduce a mechanism for product listings popularity prediction from the images representing those listings. We then explore the correlation between listings' popularity and user interaction with what is for sale. Because sales are rare in comparison to the number of items available on a large site such as Etsy, we look into an alternative mechanism for interaction, the "favorite." Favorites on Etsy are similar to any number of "like" mechanisms available online, the most familiar of which is Facebook's ubiquitous "thumbs-up." By considering what users explicitly express interest in, we are able to form relationships between user preferences and popularity-based image features.

The remainder of the paper is organized as follows: Section II discusses background on predicting image popularity. Section III describes the features

<sup>1</sup>[www.etsy.com](http://www.etsy.com)

we used to predict popularity. We examine the performance of popularity features in predicting favorite listings classification in section IV-A. Finally, we conclude this paper in section V and propose future research directions.

## II. BACKGROUND

Early work in the literature has defined popularity as quality [5] or aesthetics [6] and use data from photography rating websites where users who have interest in photography upload their photos and rate others. Popularity has also been defined as memorability [7], and interestingness [8], [9]. More recent work has directly tackled popularity. In [10], popularity is defined as the number of views on Flickr, and [2] uses favorited listings on Etsy.

Popularity tends to be predicted using SVM classification or regression [6] [10] [11] [12]. Datta et. al. [6] uses a two class SVM classifier with a forward selection algorithm to find good feature sets. By using elastic net to rank feature relevance to aesthetics, and a best first algorithm to find feature sets that minimize the RMSE cross validation error, [12] are able to achieve a 30.1% improvement compared to [11]. A few have explored other machine learning techniques. In [5] a naive Bayes classifier is used, not SVM. Aryafar et. al [2] studied the significance of color in favorited listings on Etsy using logistic regression, perceptron, passive aggressive and margin infused relaxed algorithms.

The features used in popularity prediction model the same qualities professional photographers use such as light, color, rule of thirds, texture, smoothness, blurriness, depth of field, scene composition [5] [6] [11] [12]. Most of these features are unsupervised, but some such as the spacial edge distribution and color distribution features of [5] require all of the labeled training data. Some recent work has looked at semantic object features. [10] used the popular CNN ImageNet to detect the presence of 1000 difference object categories in the image. The presence/absence of these categories is used as the feature.

With more than 30 million active listings and over 90 million listing images, Etsy provides a unique visually enticing experience for users. Because images are uploaded by users of the site, representing the myriad items for sale, these images are composed of different items, presented with various lighting conditions, scene geometries and background selections. One of the key components of e-commerce websites is efficient image search and color filtering methods. Presence of occluded backgrounds and highly textured material can hinder the accuracy of color detection algorithms. In our work, we define popularity as listings that have been favorited, clicked on, or purchased, and we show that unsupervised image

popularity features are statistically significant when combined with traditional text meta-data features in predicting popularity.

## III. FEATURES

The quality features extracted from images are composed of a set of hand-crafted features. The implementation of this features is made publicly available <sup>2</sup>.

### A. Simplicity

High quality photos are typically simpler than others. They often have one subject placed deliberately in the frame. Sometimes the background is out of focus to emphasize the subject. Poor quality photographs tend to have cluttered backgrounds and it may be difficult to distinguish the subject of the scene. We used the four measures of simplicity from [5], spatial edge distribution, hue count, contrast and lightness, and blur.

1) *Spatial Edge Distribution*: Spatial edge distribution measures how spread out sharp edges are in the image. A single subject is expected to have a small distribution while an image with a cluttered background would have a large distribution. Edges are detected by applying a 3x3 Laplacian filter and taking the absolute value. The filter is applied to each RGB channel independently and the final image is computed as the mean across all three channels. The Laplacian image is resized to 100x100 and normalized to sum to 1. Then, the edges are projected onto the x and y axis independently. Let  $w_x$ , and  $w_y$  be the width of 98% of the projected edges respectively. The image quality feature  $f = 1 - \frac{w_x w_y}{100}$  is the percent of area outside the majority of edges. Figure 2 shows the edges detected from two different images and their respective feature value.

2) *Hue Count*: Professional photographs look more colorful and vibrant, but actually tend to have less distinct hues because cluttered scenes contain many heterogeneous objects. We use a hue count feature by filtering an image in HSV color space such that V is in the range of [0.15, 0.95] and S is greater than 0.2. A 20 bin histogram is computed on the remaining H values. Let  $m$  be the maximum value of the histogram and let  $N = \{i | H(i) > \alpha m\}$ , be the set of bins values greater than  $\alpha m$ . The quality feature  $f = 20 - ||N||$  is 0 when there are a many different hues and larger as the number of distinct hues in the image goes down. We used  $\alpha = 0.05$  as in [5].

<sup>2</sup>We make our feature extraction pipeline for image quality features available at: <https://github.com/szakrewsky/quality-feature-extraction>



Fig. 2: The Laplacian image for computing spacial edge distribution for two images. The feature for figure a. is 0.013 and for b. is 0.30.

3) *Contrast and Lightness*: Brightness is a well known variable that professional photographers are trained to understand and adjust. We use the average brightness feature [5], [11] computed from the L channel of the Lab color space. Contrast is similar, and is the ratio of maximum and minimum pixel intensities. We sum the RGB level histograms, and normalize it to sum to 1. We use the width of the center 98% mass of the histogram [5].

### B. Blur

Blurry images are almost always considered to be of poor quality. We use the blur features of [5] and [13]. In [5] blur is modeled as  $I_b = G_\sigma * I$  where  $I_b$  is the result of convolving a Gaussian filter with an image. The larger the  $\sigma$  the more high frequencies are removed from the image. Assuming the frequency distribution of all  $I$  is approximately the same, then the maximum frequency  $\|C\|$  can be estimated as  $C = \{(u, v) \mid \|FFT(I_b)\| > \Theta\}$ . The feature is  $f = \|C\| \sim 1/\sigma$ , after normalizing by the image size.

In [13], blur estimation is done based on changes in the edge structures. The blur operation will cause gradual edges to lose sharpness. Assuming that most

images have gradual edges that are sharp enough, the blur is measured as the ratio of gradual edges that have lost their sharpness.

### C. Rule of Thirds

The rule of thirds is an important composition technique. Thirds lines are the horizontal and vertical lines that divide an image into a 3x3 grid of equal sized cells. The rule of thirds states that subjects placed along these lines are aesthetically more pleasing and more natural than subjects centered in the photograph. In order to segment the subject of the image from the background, we use the Spectral Residual saliency detection algorithm [14]. The feature is a 5x5 map where each cell is the average saliency value [15]. Let  $w_p$  be the saliency value of the pixel and  $A(W_i)$  is the area of the cell, then the value of each cell is

$$w_i = \frac{\sum_{p \in W_i} w_p}{A(W_i)}. \quad (1)$$

To compute the feature, the image is divided into a 5x5 grid with emphasis on the thirds lines; the horizontal and vertical regions centered on the thirds lines are 1/6 of the image size. Figure 3 shows the

saliency detection with the 5x5 grid overlay, and the thirds map feature for an image.

#### D. Texture

A smooth image may indicate blur or out-of-focus, and the lack of which may indicate poor film, or too high an ISO setting. In contrast, texture in the scene is an important composition skill of a photographer. Smoothness may indicate the lack of texture. Texture and smoothness are some of the most statically correlated features for quality/popularity [12] and [10]. We use three smoothness/texture features from these.

A three level wavelet transform is applied to the L channel of the Lab color space. We only use the bottom level of the pyramid. The result is squared to indicate power. Let  $b = \{HH, HL, LH\}$  be the bottom level of a wavelet transform, the feature is

$$f = \frac{1}{3MN} \sum_{m=1}^M \sum_{n=1}^N \sum_b w^b(m, n) \quad (2)$$

where  $w$  is the square of the wavelet value. Because the Laplacian is often used as a pyramid of different scales, another feature

$$f = \frac{1}{MN} \sum_{m=1}^M \sum_{n=1}^N l(m, n) \quad (3)$$

is also used. This time  $l$  is the second level from the bottom of a Laplacian pyramid.

Another texture feature is computed using local binary pattern (LBP). Then a pyramid of histograms are computed as in [16]. Figure 4 shows the similarities of LBP features and the three channels of Daubechies db1 wavelet.

#### E. Depth of Field

Depth of field is the distance between the nearest and farthest objects that appear in sharp focus. A technique of professional photographers is to use low depth of field to focus on the photographic subject while blurring the background. We used the feature [6] of the ratio of high frequency detail in center regions of the image compared to the entire image. Let  $w$  be the bottom level of a wavelet transform, the feature is

$$f = \frac{\sum_{(x,y) \in M_6 \cup M_7 \cup M_{10} \cup M_{11}} w(x, y)}{\sum_{i=1}^{16} \sum_{(x,y) \in M_i} w(x, y)}, \quad (4)$$

where  $M_i | 1 \leq i \leq 16$  are the cells of a 4x4 grid. The same feature is also reapplied using the Laplacian pyramid  $l$  instead of  $w$  [12]. These features only look at the center region of the image. A third feature [12] looks at the spacial distribution of high frequency details. Let  $l$  be the bottom layer of a Laplacian

pyramid and  $c_{row}, c_{col}$  are the center of mass, the feature is

$$f = \frac{1}{MN} \sum_{m=1}^M \sum_{n=1}^N l(m, n) \sqrt{(m - c_{row})^2 + (n - c_{col})^2}. \quad (5)$$

Figure 5 visualizes how these features are computed for an image.

#### F. Experimental

Maximally Stable Extremal Regions (MSER) [17] can be used to detect text because characters are typically single solid colors with sharp edges that stand out from the background [18]. Additionally, texture patterns are also often detected by MSER, like bricks on a wall. We used the experimental feature the count of the number of MSER regions. We would like to continue this experiment into other features based on text in images.

### IV. POPULARITY PREDICTION

Once we have extracted all the quality features

Once an item is listed on Etsy, the users can favorite a listing which allows them to bring all the items they like in one place. We first examine the listings dominant colors to predict if a listings is likely to be favorited or not. Then we explore the entropy of dominant colors among users favorites to indicate the color variations among favorited listings. Two datasets are used for classification of favorited listings and color entropy experiments. The classification dataset consists of 2.73 million unique listings that have been created within the last month on Etsy. The listings images are tagged with the top three dominant colors and labeled as positive if they have been favorited by a user in that period. To collect the entropy dataset, a set of 11235 active users with more than 20 and less than 2000 favorited listings within the past six months are selected. The entropy dataset then consists of 2.32 million unique listings that candidate users have favorited over the last 6 months on Etsy. These listings images are also tagged with top three dominant colors. These two datasets contain more than 5.05 million listing images and dominant color tags and are available through Etsy's API<sup>3</sup>.

#### A. Classification

Once the classification dataset has been tagged with top three dominant colors, we extract textual information from the listings. These textual features consist of the tokenized listings titles unigrams and bigrams and tokenized listings tags unigrams. We then represent each listing with a feature vector including textual features and color unigrams. A binary

<sup>3</sup>[www.etsy.com/developers](http://www.etsy.com/developers)



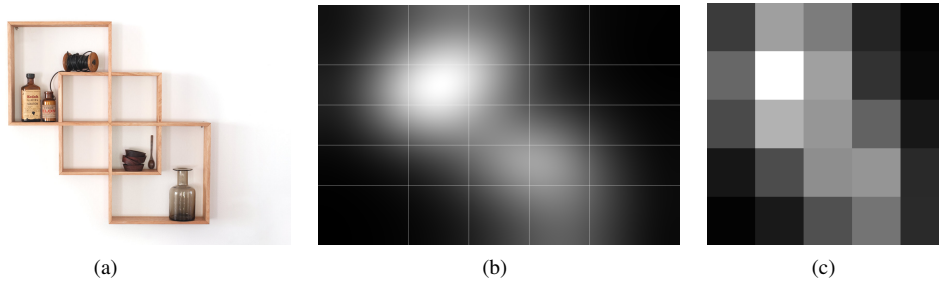


Fig. 3: Example of Rule of Thirds feature. Figure b. shows the SR saliency detection, and c. shows the thirds map feature.

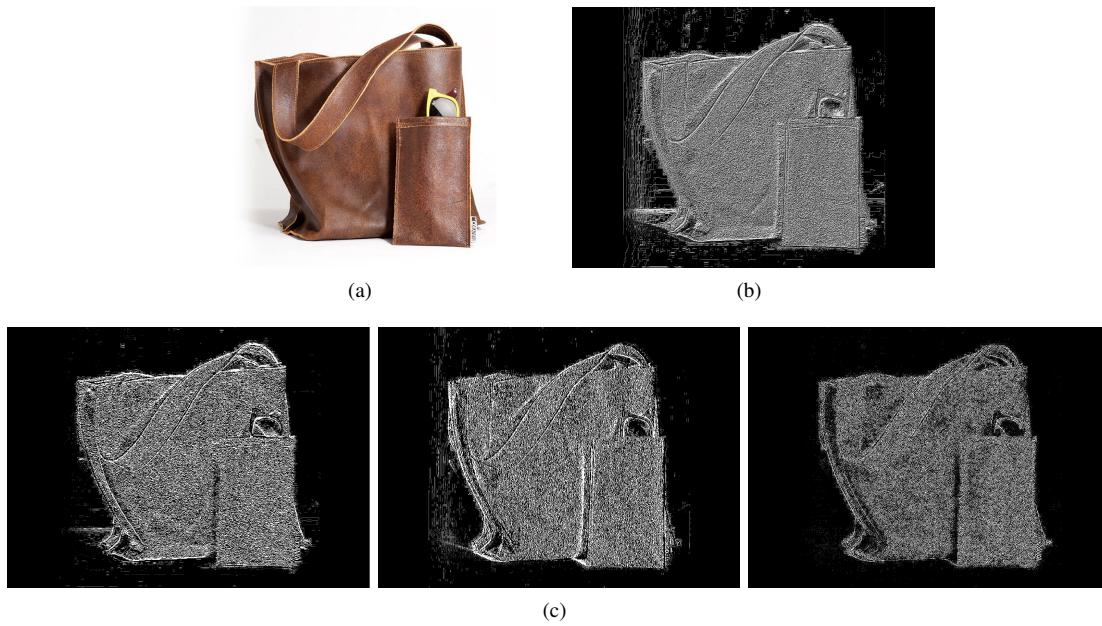


Fig. 4: Smoothness and texture features. Figure b. shows Local Binary Pattern (LBP) feature image, and c. shows the 3 channels of the DB1 wavelet transform.

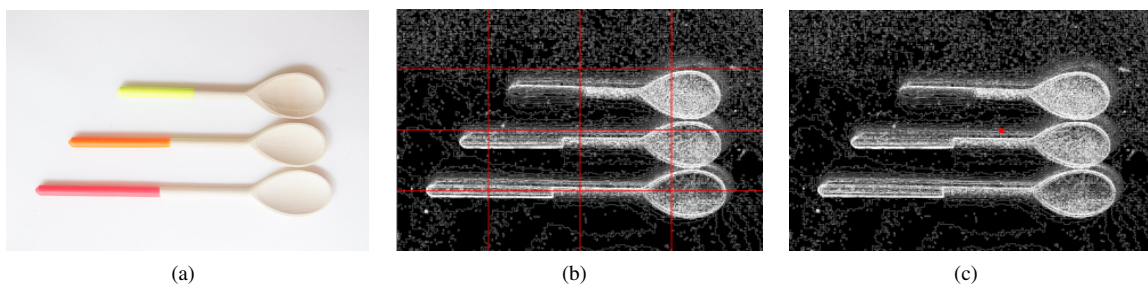


Fig. 5: Figure b. shows the Low Depth of Field features in the center grid region for the Laplacian image. Figure c. shows the same image with its center of mass.

Feature	Dimension
'Ke06-qa': spatial edge distribution	1
'Ke06-qh': hue count	1
'Ke06-ql': blur	1
'Ke06-tong': blur tong etal	1
'Ke06-qct': contrast	1
'Ke06-qb': brightness	1
'-msr count': msr count	1
'Mai11-thirds map': thirds map	25
'Wang15-f1': avg lightness	1
'Wang15-f14': wavelet smoothness,	1
'Wang15-f18': laplacian smoothness	1
'Wang15-f21': wavelet low dof	1
'Wang15-f22': laplacian low dof	1
'Wang15-f26': laplacian low dof swd	1
'Khosla14-texture': texture	5120

TABLE I: Feature Dimensions

Classification method	Average accuracy rate	AUC
logistic regression	0.5512	0.5694
perceptron	<b>0.5600</b>	<b>0.5906</b>
passive aggressive	0.5329	0.5240
MIRA	0.5232	0.5240

TABLE II: Average classification accuracy rate lift are reported using text features and multimodal feature vectors.

classification is then performed to predict if the test listings are favorited by users. We report the average classification accuracy rate and the area under the curve (AUC) with four different classifiers. Logistic regression, passive-aggressive classifier, perceptron and margin infused relaxed algorithm (MIRA) are used as the learning models. Table II shows the results of four rounds of 5-fold classification on this dataset. The textual information and color unigrams do not indicate a strong improvement in favoriting behaviour prediction. It will be interesting to observe the effects of image quality, memorability, aesthetics and interestingness for the similar problem on this unique dataset.

## V. CONCLUSION

This work represents a initial study on understanding how images, specifically color-based image features, can be used to represent items in an e-commerce setting, thereby providing better user understanding and a better overall shopping experience. To facilitate this understanding, this work proposed an empirical method to estimate the dominant colors of images representing product listings on Etsy, using object localization.

We used this dominant colors to filter listings in a conventional text-based e-commerce search, according to user input. Moreover we explored the color distribution among candidate users favorited listings. We also examined the impact of color unigrams on

users favoriting behaviour. This work represents the tip of the iceberg of understanding how visual cues influence user action in an e-commerce setting.

## REFERENCES

- [1] C. Anderson, *The Long Tail: Why the Future of Business Is Selling Less of More*. Hyperion, 2006. 1
- [2] K. Aryafar, C. Lynch, and J. Attenberg, "Exploring user behaviour on etsy through dominant colors," in *2014 22nd International Conference on Pattern Recognition (ICPR)*. IEEE, 2014, pp. 1437–1442. 1, 2
- [3] Y. J. Wang, M. S. Minor, and J. Wei, "Aesthetics and the online shopping environment: Understanding consumer responses," *Journal of Retailing*, vol. 87, no. 1, pp. 46–58, 2011. 1
- [4] P. Obrador, X. Anguera, R. de Oliveira, and N. Oliver, "The role of tags and image aesthetics in social image search," in *Proceedings of the first SIGMM workshop on Social media*. ACM, 2009, pp. 65–72. 1
- [5] Y. Ke, X. Tang, and F. Jing, "The design of high-level features for photo quality assessment," in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, vol. 1. IEEE, 2006, pp. 419–426. 2, 3
- [6] R. Datta, D. Joshi, J. Li, and J. Z. Wang, "Studying aesthetics in photographic images using a computational approach," in *Computer Vision–ECCV 2006*. Springer, 2006, pp. 288–301. 2, 4
- [7] P. Isola, J. Xiao, A. Torralba, and A. Oliva, "What makes an image memorable?" in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 145–152. 2
- [8] S. Dhar, V. Ordonez, and T. L. Berg, "High level describable attributes for predicting aesthetics and interestingness," in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 1657–1664. 2
- [9] M. Gygli, H. Grabner, H. Riemenschneider, F. Nater, and L. Van Gool, "The interestingness of images," 2013. 2
- [10] A. Khosla, A. Das Sarma, and R. Hamid, "What makes an image popular?" in *Proceedings of the 23rd international conference on World wide web*. ACM, 2014, pp. 867–876. 2, 4
- [11] M. Chen and J. Allebach, "Aesthetic quality inference for online fashion shopping," in *IS&T/SPIE Electronic Imaging*. International Society for Optics and Photonics, 2014, pp. 902 703–902 703. 2, 3
- [12] J. Wang and J. Allebach, "Automatic assessment of online fashion shopping photo aesthetic quality," in *Image Processing (ICIP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 2915–2919. 2, 4
- [13] H. Tong, M. Li, H. Zhang, and C. Zhang, "Blur detection for digital images using wavelet transform," in *Multimedia and Expo, 2004. ICME'04. 2004 IEEE International Conference on*, vol. 1. IEEE, 2004, pp. 17–20. 3
- [14] X. Hou and L. Zhang, "Saliency detection: A spectral residual approach," in *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*. IEEE, 2007, pp. 1–8. 3
- [15] L. Mai, H. Le, Y. Niu, and F. Liu, "Rule of thirds detection from photograph," in *Multimedia (ISM), 2011 IEEE International Symposium on*. IEEE, 2011, pp. 91–96. 3
- [16] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, vol. 2. IEEE, 2006, pp. 2169–2178. 4

- [17] J. Matas, O. Chum, M. Urban, and T. Pajdla, “Robust wide-baseline stereo from maximally stable extremal regions,” *Image and vision computing*, vol. 22, no. 10, pp. 761–767, 2004. 4
- [18] H. Chen, S. S. Tsai, G. Schroth, D. M. Chen, R. Grzeszczuk, and B. Girod, “Robust text detection in natural images with edge-enhanced maximally stable extremal regions,” in *Image Processing (ICIP), 2011 18th IEEE International Conference on*. IEEE, 2011, pp. 2609–2612. 4