

Finding Good Composition in Panoramic Scenes

Yuan-Yang Chang and Hwann-Tzong Chen

Department of Computer Science, National Tsing Hua University
30013 Hsinchu, Taiwan

Abstract

We introduce a new problem of automatic photo composition, and present an effective technique for finding good views within a panoramic scene. Instead of applying heuristic rules of photo composition, we propose to imitate good composition presented in the artworks of professional photographers. Our approach tries to model the composition styles of professional photographs by analyzing the structural features and the layout of visual saliency. The task of finding good photo composition through a viewfinder is formulated as a search problem, and we present a stochastic search algorithm to look for good viewing configurations and to choose suitable reference images from the collection of masterpiece photographs. Given any initial location in the panoramic scene, our algorithm is able to suggest a better view that would often yield professional-like photo composition.

1. Introduction

The composition of a photo is one of the essential ingredients in the craft of photography. Good and balanced composition could make a photograph look more appealing even if the scene being shot is normal. Photographers often follow similar rules of composition that have been applied by painters for hundreds of years. However, a major difference between photography and painting is that in photography the composition is determined by observing and framing the scene through the viewfinder, while a painter may start with an empty canvas and modify the composition repeatedly until the outcome is satisfying, as reflected in the quote of famous photographer Edward Steichen: “Every other artist begins with a blank canvas, a piece of paper...the photographer begins with the finished product.”

Although sometimes experienced photographers might deliberately break the rules to create tension in a photograph through peculiar composition, in general being aware of basic rules of composition helps to produce more dynamic and visually pleasing photos. In photography, widely accepted principles such as *rule of thirds*, *leading lines*, *repeating*

patterns, *layering*, *horizon lines*, *relative scale* [11] are all useful compositional techniques for creating better photos. Amateur photographers might know well the rules of composition, but still find it difficult to apply the rules when they try to take pictures. It usually requires years of practice and experience to transform the rules into intuitions so that an instantly taken photo may have a better chance to resemble those great pictures in *National Geographic*. It will be a great benefit if the composition of masterpiece photographs can be modeled and be used as guidance when producing new photos. (See Fig. 1 for example.) In this paper we address the problem of automatic photo composition in a panoramic scene. This problem is new in computer vision, and we present a method that can construct graphically appealing images based on examples of good composition.

We build a collection of pictures taken by experienced photographers, and try to model the structural and salient elements presented in the pictures. These pictures are used as the *exemplars* of photo composition. We present a scheme that can automatically make suggestions of good composition by comparing the structure and saliency between the exemplars and the candidate views. Our goal is to search in the panoramic scene, by controlling certain parameters of the camera, to find a good way of constructing the photos. The process of search is guided by the exemplars, and the camera might move, pan, tilt, or zoom to find a favorable view with good composition for taking pictures of the panoramic scenery.

1.1. Related Work

Previous computational approaches to automatic photo composition focus on implementing simple heuristic guidelines, *e.g.* the rule of thirds, for the placement of the subject [3], [4]. We try to take a different approach to the task by introducing an auto-composition scheme, which is able to imitate good composition presented in the artworks of master photographers. The idea is similar to the work of Bae *et al.* [2] on tone management, in which they seek to improve the photographic quality of an image by transferring the photographic look of a model photograph to the image being edited.

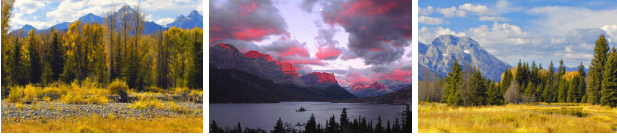


Figure 1. Left: A landscape photo with casual composition. Middle: A photograph taken by professional photographer William Neill. (From *Landscapes of the Spirit*, available on Flickr <http://www.flickr.com/photos/williamneill>) Right: Our method achieves good composition by imitating the artwork of William Neill. The composition presents the same placement of the horizon line and a similar layering effect as in the middle image.

To model the composition styles used in professional photography, we need an image representation that is able to capture the structural features and the spatial layout of salient content. With an appropriate image representation we could obtain a similarity measure to compare images based on the resemblance in their composition. In this work, we consider the low-dimensional global image representation derived from the GIST descriptor, which has been shown to be a good image representation for scene classification and scene matching [13], [17], as well as for depth estimation [19]. The GIST descriptor computes oriented edge filter responses at different scales aggregated into spatial bins of locations. The arrangement of structures within an image can be characterized by the GIST descriptor, and thus we can measure the similarity between the geometric patterns and scene structures presented in two images by comparing their GIST descriptors.

Besides the scene structures, we also wish to model the visual saliency in an image, since the placement of salient elements in the photo is also an important compositional issue. The layout of the salient elements needs to draw the eye into the photo and also has to be well balanced and pleasing to the eye. We try to assess the layout of salient elements by adopting the techniques developed for saliency detection and visual attention, which have long been studied in computer vision, *e.g.*, [9], [10], [14]. We combine the information provided by the GIST descriptor and the result of saliency detection, and use the coupled representation to model the composition of photos. The effectiveness of integrating gist and saliency has also been addressed by Siagian and Itti in their model of human vision for scene classification [16].

Our work is also related to the image editing methods for image warping, scene completion, and retargeting/resizing, *e.g.*, [1], [5], [7], [8], [20], which can be considered as alternative ways of modifying the photo composition. Gal *et al.* [7] present a method for feature-aware texture warping. Their method is able to retain the shape of the regions of foreground objects while changing the aspect ratio of the image. For the application of modifying the compo-

sition of images, their method may be used to adjust the relative scales of different parts in the image. Avidan and Shamir [1] present the *seam carving* algorithm for image retargeting. The algorithm aims to reduce the image size by removing less important seams of pixels from the original image such that the resized image preserves most of the perceptually significant parts. The *patch transform* algorithm presented by Cho *et al.* [5] divides an image into non-overlapping patches, and reorganizes the patches to form a new image, subject to user-specified constraints such as the spatial locations of patches. The patch transform can be applied to various image editing tasks, *e.g.*, image retargeting, or changing the location of foreground object. Hays and Efros [8] present a scene completion algorithm that, given an input image, looks for similar scenes in two million images using the GIST descriptor, and then fills the holes in the input image with good patches extracted from the selected similar scenes. Their algorithm can be used to produce contextually valid and visually pleasing composite images. The main limitation of applying these image editing methods to automatic photo composition is that the methods would tinker with the image content and inevitably change the real structures within the image. Moreover, the image editing methods still require suitable image representations and evaluation criteria for characterizing good composition.

Deselaers *et al.* [6] present the pan-zoom-scan method for automatic video cropping. Their method does not change the scene structures in images, and can find the best cropped viewing area for each image in a video sequence through panning and zooming. Santella *et al.* [15] describe an interactive method for cropping photos based on the information provided by eye tracking. Through user studies they show that their gaze-guided method performs better than fully-automatic cropping approach such as [18]. Unlike our goal of composition, the aforementioned methods focus on modeling visual saliency and attention in images, but do not take into account more complex compositional components. Lalonde *et al.* [12] propose the use of a physically-based sky model to analyze the information available in the visible sky. The model can be applied to the segmentation of the sky and cloud layers, and the bi-layered representation for sky and clouds is useful for data-driven sky matching. For landscape and outdoor photography, the bi-layered representation might provide helpful hints to modify visual balance in the composition.

2. Formulation of the Search Problem

For convenience sake, instead of using an active PTZ camera to take pictures directly, we simply try to simulate the process of observing a panoramic scene through a viewfinder and searching for good views at different zooms. Nevertheless, the approach presented in this paper for the virtual environment can be easily adapted to real control-

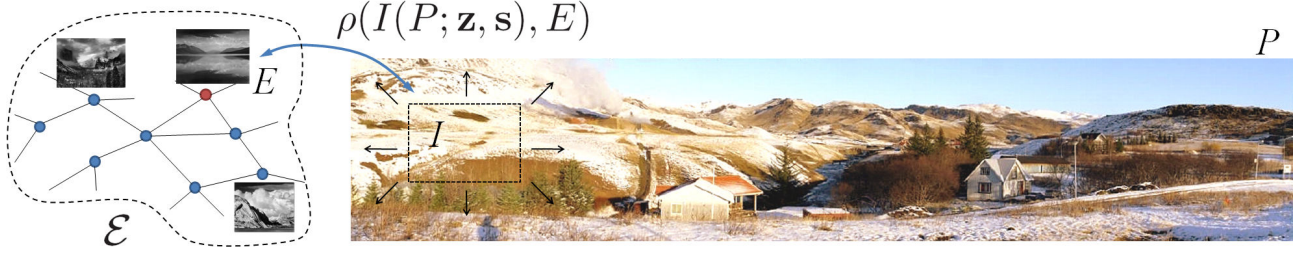


Figure 2. The scheme of finding good composition in a panorama. Given a panorama P to be explored, we can extract the view $I(P; \mathbf{z}, \mathbf{s})$, according to the current configuration of the viewfinder defined by position \mathbf{z} and scale \mathbf{s} . The distance measure $\rho(I(P; \mathbf{z}, \mathbf{s}), E)$ computes the dissimilarity in composition between the view $I(P; \mathbf{z}, \mathbf{s})$ and an exemplar E in the set \mathcal{E} .

lable camera systems, without modifying the search algorithm. In our formulation, we assume that a panorama of a scene is given, and only two camera motions are considered: translation and zoom. Furthermore, camera translation is approximated by a translation in the image plane, ignoring the parallax effects, and zoom is approximated by scaling the view.

We formulate the problem of finding good photographic composition within a panoramic scene as a search problem in a joint state-space. To begin with, we describe the configuration of the viewfinder by the 2D position of the viewing center, $\mathbf{z} = [x \ y]^T$, and the size of view, $\mathbf{s} = [w \ h]^T$. Given a panorama P to be explored, we can extract the image region $I(P; \mathbf{z}, \mathbf{s})$ within the panorama, according to the current configuration of the viewfinder defined by \mathbf{z} and \mathbf{s} . Suppose we have a set of exemplars, \mathcal{E} , for providing compositional guidance, we can express the search problem by

$$\{\mathbf{z}^*, \mathbf{s}^*, E^*\} = \arg \min_{\substack{\mathbf{z}, \mathbf{s} \\ E \in \mathcal{E}}} \rho(I(P; \mathbf{z}, \mathbf{s}), E), \quad (1)$$

where $\rho(\cdot, \cdot)$ is a distance measure to compute the dissimilarity in composition between an image region $I(P; \mathbf{z}, \mathbf{s})$ and an exemplar E . Fig. 2 illustrates the scheme of the search problem. In sum, the problem we are interested in solving is to select the most suitable exemplar for the panoramic scene, and at the same time, to find the best view that would look very similar to the exemplar.

Finding an optimal solution to the search problem would be time-consuming and impractical for real-world applications. Therefore we propose a stochastic-based algorithm to find approximate solutions in an active search setting, where each step of the algorithm just requires locally available information. Before going into the details of the search algorithm, we need to define the image representation and the dissimilarity measure ρ concerning photo composition.

3. Image Representation

We attempt to model the composition styles by analyzing the structural features, the geometric patterns, and the

spatial layout of salient elements. In particular, we use the GIST descriptor [13] to characterize the arrangement of structures and geometric patterns within an image. For the layout of salient elements, we introduce the *saliency descriptor*: we compute the saliency map using the spectral residual approach [9], and then divide the saliency map into spatial bins as in the computation of the GIST descriptor.

The GIST descriptor. The GIST descriptor measures the oriented edge responses at multiple scales, and aggregates the responses into spatial bins. As shown in Fig. 3a, the input image is converted into grayscale, and we use the 6×4 spatial bins to capture the structural elements in the input image. The descriptor is built from three coarse-to-fine scales with 8, 8, 4 filter orientations, and the aggregated descriptor is a 480 ($= (8 + 8 + 4) \times 24$) dimensional vector. We then extend the descriptor to 481 dimensions by adding a dummy dimension and assign a threshold value to the corresponding component. We normalize the 481-dimensional GIST descriptor to have unit norm. The purpose of appending the dummy component with threshold value is to prevent the noise in a featureless image being overemphasized.

The saliency descriptor. Given an input image, we use the spectral residual [9] to capture the statistical singularity in the frequency domain. The spectral residual is defined as the difference between the log-spectrum and the averaged spectrum of the input image. We can then derive the saliency map from the spectral residual through inverse Fourier transform. Larger intensity values in the saliency map correspond to more salient parts in the input image. The saliency map is also divided into 6×4 spatial bins. We use two scales 32 and 64 to construct the saliency maps, as shown in Figs. 3b & 3c, and compute the sum of the intensity values in each spatial bin for each saliency map, resulting in a 48 ($= 2 \times 24$) dimensional saliency descriptor. We also append a dummy component to the original saliency descriptor and then normalize the saliency descriptor to have unit norm, as is done for the GIST descriptor.

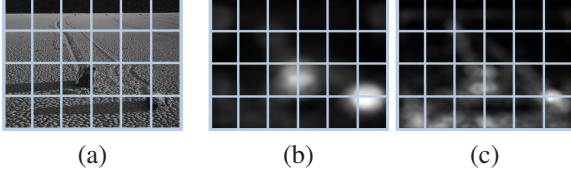


Figure 3. (a) The 6×4 spatial bins used to compute the GIST descriptor. (b) & (c) The saliency maps at two scales 32 and 64, also divided into 6×4 spatial bins.

Finally, we represent an image I using the GIST descriptor $G(I)$ and the saliency descriptor $S(I)$, and the dissimilarity in composition between two images can be measured by

$$\rho(I, I') = \|G(I) - G(I')\| + \lambda \|S(I) - S(I')\|, \quad (2)$$

where λ is a weighting parameter to determine the significance of the two properties in comparing photo composition: i) the structural and geometric patterns in the scene described by GIST, or ii) the layout of visual saliency represented by the saliency descriptor. We use $\lambda = 0.2$ in the experiments.

4. Algorithm

The neighborhood graph of exemplars. The search problem in (1) needs to find the most suitable exemplar for taking a picture within the panorama. It would be impractical to examine the entire set of exemplars if we hope to achieve an active searching performance. To make the search feasible, we pre-process the exemplar set and build a neighborhood graph on the exemplars. This graph structure helps the algorithm to explore locally in the set of exemplars for better candidates. For each exemplar E in the set \mathcal{E} , we find the k nearest neighbors of E according to the dissimilarity measure $\rho(E, \cdot)$, and connect them to E . We use an exemplar set containing more than a hundred photos. We construct the neighborhood graph with $k = 4$, and have found it to be effective for the problem. Fig. 4 shows the dissimilarity matrix and the neighborhood graph of the exemplars used in our experiments.

The stochastic search algorithm. Recall that the configuration of the viewfinder is defined by the position of the viewing center, \mathbf{z} , and the size of the viewing area, \mathbf{s} . We can use the configuration to extract the image $I(P; \mathbf{z}, \mathbf{s})$ of the corresponding view within the panorama, and then compare the extracted image with a chosen exemplar E by $\rho(I(P; \mathbf{z}, \mathbf{s}), E)$. Instead of finding the optimal solution to the problem in (1), we present a stochastic search algorithm to find approximate solutions as recommendations for shooting the scene. At the end of the search, we expect the

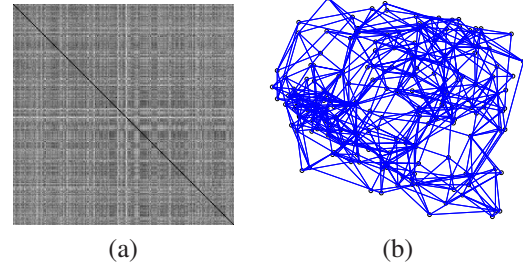


Figure 4. (a) The dissimilarity matrix of the exemplars \mathcal{E} used in our experiments. (b) The neighborhood graph of \mathcal{E} .

algorithm to suggest a candidate view that would yield good photo composition. The search algorithm is summarized as follows.

- **Input:** A panorama P and the neighborhood graph built from the exemplar set \mathcal{E} . A starting position of the center of view, $\mathbf{z}_0 = [x_0 \ y_0]^T$. The initial size of the view, $w_0 \times h_0$.

- **Initialization:** Let \mathbf{s}_0 denote the scale parameter as $\mathbf{s}_0 = [w_0 \ h_0]^T$. We randomly choose $(\lfloor \log_2 |\mathcal{E}| \rfloor + 1)$ exemplars from \mathcal{E} , and pick among them the exemplar E_0 that has the smallest dissimilarity to the initial view $I(P; \mathbf{z}_0, \mathbf{s}_0)$. The following steps are repeated for the initial configuration $\{\mathbf{z}_0, \mathbf{s}_0, E_0\}$.

- **For $t = 1, \dots, T$:**

1. **Find a better exemplar in the neighborhood graph.**

$$E_t = \arg \min_{E' \in \Omega} \rho(I(P; \mathbf{z}_{t-1}, \mathbf{s}_{t-1}), E'), \quad (3)$$

where $\Omega = \{E_{t-1}\} \cup \mathcal{N}(E_{t-1})$ is the union of E_{t-1} and its neighbors.

2. **Update the size of view.** Re-scale the size of the view by $\pm 5\%$ of the current size and compute the dissimilarity between the re-scaled view and the exemplar E_t . If the dissimilarities are not improved, then keep the current view size; otherwise update the scale parameter \mathbf{s}_t of the view according to the new size.
3. **Find a better position of the view.** Generate M random vectors, $\{\mathbf{v}^{(m)}\}_{m=1}^M$ (typically $M = 8$), from a two-dimensional normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$ with zero mean and identity covariance matrix. Solve

$$\mathbf{v}^* = \arg \min_{\mathbf{v}' \in \{\mathbf{v}^{(m)}\} \cup \{\mathbf{0}\}} \rho(I(P; \mathbf{z}_{t-1} + \mathbf{A}\mathbf{v}', \mathbf{s}_t), E_t), \quad (4)$$

where \mathbf{A} is diagonal scaling matrix for controlling the step size of a move. Update the position by $\mathbf{z}_t \leftarrow \mathbf{z}_{t-1} + \mathbf{A}\mathbf{v}^*$.

If the position \mathbf{z}_t does not change during the previous τ consecutive iterations, then do the next step for greedy local search; otherwise, continue the iteration.

4. **Greedy local search.** Keep s_t and E_t fixed, and find the target position z_t by taking the best moves along the path in a 4-connected pixel grid until no further improvement can be achieved. Terminate the iteration.

• **Output:** Obtain the final configuration $\{\hat{z}, \hat{s}, \hat{E}\}$ and suggest the corresponding view for taking photos.

5. Experimental Results

We collect more than a hundred landscape and outdoor photos taken by professional photographers, including Ansel Adams, Jay Dickman, Peter Essick, and William Neill. We use these photos as the exemplars. A neighborhood graph of the exemplars is constructed by linking each exemplar to its four nearest neighbors as shown in Fig. 4b. We consider only images in the *landscape* aspect ratio, and all exemplars are resized to 240×160 pixels. In our implementation with MATLAB, to find a candidate view in a typical panorama ($\sim 3000 \times 500$ pixels) takes about 20 seconds to one minute, and the number of iterations is set to $T = 25$. We have $\tau = 10$, *i.e.*, the algorithm turns to the greedy local search if it is stuck for 10 iterations. Our stochastic search algorithm is efficient: It would take more than two hours if an exhaustive search is instead performed to check 100×15 downsampled locations in the panorama at 13 different scales, using a fixed exemplar. An example of a view found by the exhaustive search is shown in Fig. 5, in comparison with a view found by our algorithm. The two views look quite similar, but the dissimilarity between the exhaustive search result and the exemplar is $\rho = 5.07$, which is a bit smaller than our solution, $\rho = 5.86$. Various experimental results are shown in Fig. 8. Our method performs well on choosing suitable exemplars for the scene, and the views found in the panorama present very similar layering effects and the arrangements of horizon lines, sky-lines, and textures.

User study. To further evaluate our approach, we conduct an assessment of the compositional quality of the views suggested by our algorithm. We choose 58 panoramas for the assessment. Inside each panorama we randomly pick 10 to 16 initial configurations with different positions and scales, and from each initial configuration we run our algorithm to find a better view. We keep the three images that correspond to the initial configuration, the final configuration, and the exemplar as a set of evaluation data. We have invited 12 people to do the assessment. These viewers are non-professionals, and only have basic knowledge about photography. The course of the assessment is divided into two sessions. In the first session, the viewer is required to evaluate 100 pairs of images. We create a UI for this task, as shown in Fig. 6a. In each of the 100 runs, the program

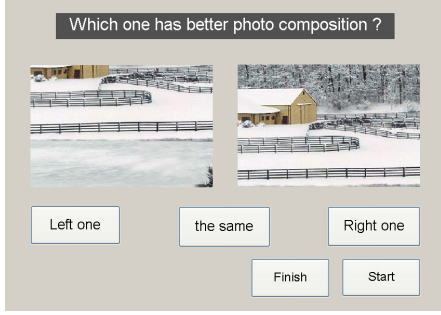


Figure 5. Left: A view found by our algorithm. Middle: The exemplar chosen by our algorithm. Right: A view found by exhaustive search using the same exemplar.

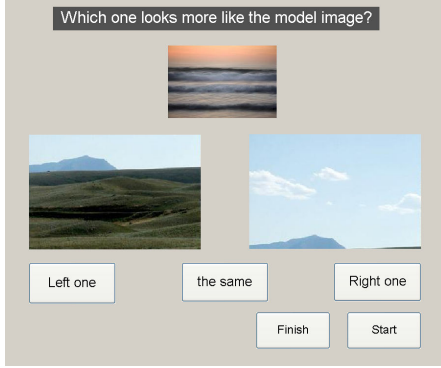
randomly selects a set of evaluation data, and shows the initial view in tandem with the suggest view to the viewer. We shuffle the two views and ask the viewer to choose the one that seems to have better photo composition. The viewer may click ‘left’ or ‘right’ as the choice, or click ‘the same’ if it is hard to judge which one is better. The setting for the second session is similar to the first session, but in each run we also show the selected exemplar in the set of evaluation data. The viewer is asked to decide which image looks more like the exemplar, in terms of photo composition. Fig. 6b illustrates an example of one run of the evaluation. The results of assessment are as follows. For the first session, we have 510 votes for the initial views, 526 votes for the suggested views, and 164 for ‘the same’. For the second session, we get 360 votes for the initial views, 575 votes for the suggested views, and 265 for ‘the same’. The vote distributions of 12 viewers are shown in Fig. 7. The result of the first session indicates that the preference on photo composition for non-experts is quite subjective, and the suggested views seem just slightly better than the initial views. Nevertheless, the results of the second session shows that our approach indeed finds views that are similar to the model images in terms of composition. A possible improvement of our approach is thus to allow the users to choose their own collection of preferred exemplars.

6. Conclusion

We have presented a method of finding photographic composition in a panorama for producing visually pleasing photos. It is possible to implement the stochastic search algorithm to control a PTZ camera, so the camera can automatically explore the scene and search for good views to take pictures. The algorithm might also be built in a point-and-shoot camera as an easy shooting mode for suggesting the user how to move the camera and zoom the view to get better composition of the photo. We are trying to build a much larger collection of exemplar photos, and to invite more people, particularly professional photographers, to do the assessment. We are also interested in extending our method to different themes in photography, especially for photographing people in a dynamic scene.



(a) First session



(b) Second session

Figure 6. User study. (a) The first session of the assessment: A random pair of images corresponding to the initial view and the suggested view are both presented to the viewer. The positions of the two images may be switched at random. The viewer is asked to determine which image has better photo composition. The viewer needs to vote for ‘left’ or ‘right’, or click the button ‘the same’ if it is hard to judge which one is better. Each viewer is required to evaluate 100 pairs of images. (b) The second session of the assessment: Similar to the first session, but this time we also show the viewer the model image (the selected exemplar), and ask the viewer which image looks more like the model image. The viewer can answer ‘left’, ‘right’, or ‘the same’.

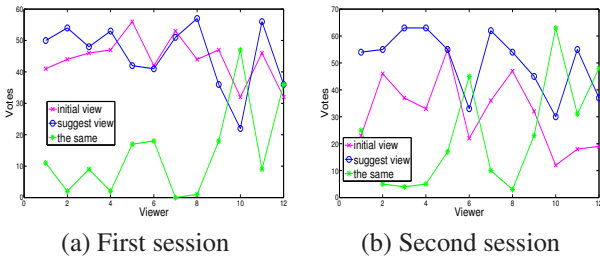


Figure 7. The vote distributions of the 12 viewers.

Acknowledgments. This research was supported in part by NSC grant 96-2221-E-007-132-MY2 and NTHU grant 98N2935E1. We thank the artists for giving us permission to use their photographs in our paper.

References

- [1] S. Avidan and A. Shamir. Seam carving for content-aware image resizing. *ACM Trans. Graph.*, 26(3):10, 2007.
- [2] S. Bae, S. Paris, and F. Durand. Two-scale tone management for photographic look. *ACM Trans. Graph.*, 25(3):637–645, 2006.
- [3] S. Banerjee and B. L. Evans. In-camera automation of photographic composition rules. *IEEE Transactions on Image Processing*, 16(7):1807–1820, 2007.
- [4] Z. Byers, M. Dixon, W. D. Smart, and C. M. Grimm. Say cheese! experiences with a robot photographer. *AI Mag.*, 25(3):37–46, 2004.
- [5] T. S. Cho, M. Butman, S. Avidan, and W. T. Freeman. The patch transform and its applications to image editing. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [6] T. Deselaers, P. Dreuw, and H. Ney. Pan, zoom, scan: Time-coherent, trained automatic video cropping. In *CVPR*, 2008.
- [7] R. Gal, O. Sorkine, and D. Cohen-Or. Feature-aware texturing. In *Proceedings of Eurographics Symposium on Rendering*, pages 297–303, 2006.
- [8] J. Hays and A. A. Efros. Scene completion using millions of photographs. *ACM Trans. Graph.*, 26(3):4, 2007.
- [9] X. Hou and L. Zhang. Saliency detection: A spectral residual approach. In *CVPR*, 2007.
- [10] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Trans. Pattern Anal. Mach. Intell.*, 20(11):1254–1259, 1998.
- [11] J. Kinghorn and J. Dickman. *Perfect Digital Photography*. McGraw-Hill/Osborne, 2005.
- [12] J.-F. Lalonde, S. G. Narasimhan, and A. A. Efros. What does the sky tell us about the camera? In *ECCV*, 2008.
- [13] A. Oliva and A. B. Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International Journal of Computer Vision*, 42(3):145–175, 2001.
- [14] A. Oliva and A. B. Torralba. Contextual guidance of eye movements and attention in real-world scenes: The role of global features on object search. *Psychological Review*, 113(4):766–786, 2006.
- [15] A. Santella, M. Agrawala, D. DeCarlo, D. Salesin, and M. F. Cohen. Gaze-based interaction for semi-automatic photo cropping. In *CHI*, pages 771–780, 2006.
- [16] C. Siagian and L. Itti. Rapid biologically-inspired scene classification using features shared with visual attention. *IEEE Trans. Pattern Anal. Mach. Intell.*, 29(2):300–312, 2007.
- [17] J. Sivic, B. Kaneva, A. Torralba, S. Avidan, and W. T. Freeman. Creating and exploring a large photorealistic virtual space. In *First IEEE Workshop on Internet Vision*, 2008.
- [18] B. Suh, H. Ling, B. B. Bederson, and D. W. Jacobs. Automatic thumbnail cropping and its effectiveness. In *UIST*, pages 95–104, 2003.
- [19] A. B. Torralba and A. Oliva. Depth estimation from image structure. *IEEE Trans. Pattern Anal. Mach. Intell.*, 24(9):1226–1238, 2002.
- [20] L. Wolf, M. Guttman, and D. Cohen Or. Non-homogeneous content-driven video-retargeting. In *ICCV*, 2007.

