

ImageBeacon: Broadcasting Color Images over Connectionless Bluetooth LE Packets

Chong Shao

Department of Computer Science

University of North Carolina at Chapel Hill
Chapel Hill, North Carolina 27599
cshao@cs.unc.edu

Shahriar Nirjon

Department of Computer Science

University of North Carolina at Chapel Hill
Chapel Hill, North Carolina 27599
nirjon@cs.unc.edu

ABSTRACT

This paper describes the first ‘image beacon’ system that is capable of broadcasting color images over a very long period (years, as opposed to days or weeks) using a set of cheap, low-power, memory-constrained Bluetooth Low Energy (BLE) beacon devices. We design an image processing pipeline that takes into account the background and foreground information of an image and then applies an adaptive encoding method which prioritizes more important regions of an image during encoding, in order to achieve the best quality image under a very strict size limit. We test our system with different types of RGB images that contain indoor and outdoor objects, buildings, and road signs. We empirically determine the tradeoffs between the system lifetime and the quality of broadcasted images, and determine an optimal set of parameters for our system, under user-specified constraints such as the number of available beacon devices, maximum latency, and life expectancy. We develop a smartphone application that takes an image and user-requirements as inputs, shows previews of different quality output images, writes the encoded image into a set of beacons, and reads the broadcasted image back. Our evaluation shows that a set of 2–3 beacons is capable of broadcasting high-quality images (70% structurally similar to original images) for a year-long continuous broadcasting, and both the lifetime and the image quality improve when more beacons are used.

CCS CONCEPTS

•Hardware → Signal processing systems;

KEYWORDS

Bluetooth beaconing, image processing

ACM Reference format:

Chong Shao and Shahriar Nirjon. 2017. ImageBeacon: Broadcasting Color Images over Connectionless Bluetooth LE Packets. In *Proceedings of The 2nd ACM/IEEE International Conference on Internet-of-Things Design and Implementation, Pittsburgh, PA USA, April 2017 (IoTDI 2017)*, 12 pages.
DOI: <http://dx.doi.org/10.1145/3054977.3054985>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IoTDI 2017, Pittsburgh, PA USA

© 2017 ACM. 978-1-4503-4966-6/17/04...\$15.00
DOI: <http://dx.doi.org/10.1145/3054977.3054985>

1 INTRODUCTION

In this modern age of the Internet of Things (IoT), it is now possible to literally glue tiny computers to everyday objects, so that they can sense, react, and tell their own stories. The IoT community has embraced wireless standards such as Bluetooth Low Energy (BLE) and developed programmable ‘beacon’ devices that periodically broadcast a small amount of preloaded data, while lasting for multiple years on a coin-cell battery. Broadcast messages from beacon devices typically contain information about an object, a location, a web-resource, or just an arbitrary string. This connectionless mode of BLE does not require a receiver to pair/bond or connect to a sender, and hence, there is no overhead of connection setup and no inconvenience of requiring a user to enter pins and passwords. These broadcast messages are received by a BLE capable mobile device to obtain relevant information just-in-time and on-the-spot. Emerging applications of beacon devices include advertising merchandise in retail stores, identifying late passengers at the airports, authorizing people at the hospitals, smarter signage, indoor navigation, and tracking moving platforms like airline cargo containers, computers on wheels, museum artworks, or even humans.

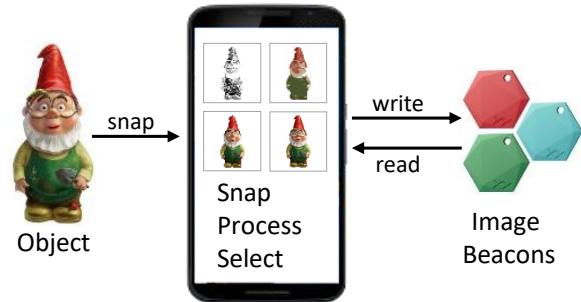


Figure 1: An image beacon system.

The enabling technology behind these applications is the ability of a beacon to simply broadcast a few bytes of data (usually a URL or a UUID) as BLE 4.0 advertisement packets. The bound in data rate comes from the lifetime requirement of these devices. Such a tight budget on payload size and the maximum data rate have limited a beacon’s capability to only be able to broadcast an identifier or a small amount of text (effectively ~30 bytes). The next generation BLE 5.0 beacon is expected to have an 8X increase in broadcasting capacity (~256 bytes). Such an increase opens up the possibility to design beacons that can serve larger assets, e.g., an image, carried

by connectionless BLE advertisement packets. However, even a simple 72×72 PNG image, such as the Android launcher icon, has a size of over 3KB. To store and broadcast this image, either we require to use a dozen of BLE 5.0 beacons, or we will have to accept a very long image transmission and loading time.

Image compression is a natural way to deal with this problem. Existing image compression algorithms, however, fail to achieve the desired compression ratio for an image to be broadcasted over BLE. Hence, a fundamental challenge toward realizing an image beacon is to devise an algorithm that efficiently represents an image using as few bits as possible, while taking into account the application-driven limits on the number of usable beacons per image, broadcast message size, data rate, latency, and lifetime. In an earlier work [14], we devised an image beacon system that broadcasts binary images of a few limited categories (e.g., handwritten characters) only. This paper is a continuation to that line of work, but this time, we have taken a harder challenge, i.e., to develop a beacon system that works for color images, e.g., images taken with a mobile phone.

Being able to broadcast images from beacons enables more powerful and feature rich applications than the ones supported by today's beacons. We envision that like the web has evolved from serving hypertexts to streaming multimedia contents, the natural successor of today's beacon devices would be the ones that broadcast images. Applications of image beacons would be in scenarios where there is no Internet connectivity but there is a need for storing and broadcasting information that can be best described by an image. For example, coordinating rescue workers in disaster areas, creating a bread-crumb system for adventurous hikers and mountaineers, remote surveillance (when coupled with a camera), or even a simple system just to let someone know that 'We were here'. Recently, Google started to experiment with an idea called '*Fat Beacons*', where they are looking into broadcasting html pages over BLE. However, for lack of a suitable image compression technique, the pages do not support images.

In this paper, we chase this seemingly impossible goal of creating an image beacon system that efficiently broadcasts color images, carried by BLE broadcast messages, over an extended period of time. We propose a self-contained system that stores and broadcasts actual image contents as opposed to IDs, links, or URLs of an image. We assume availability of no additional information on the broadcasted image from any other sources – globally (on the web) or locally (on a user's smartphone that receives the broadcast).

The crux of the system is an algorithm that analyzes an image to identify its 'important' semantic regions (as defined by the user or the use case) and then encodes them differently than the rest of the image to reduce the overall image size. The image data are written to and read from the image beacon system using a smartphone application, which runs the proposed compression and rendering algorithms. We use the term 'beacon system' instead of 'a beacon', since a compressed image may still require more than one physical beacons to ensure its acceptable quality. Allowing multiple beacons per image makes the system flexible. It widens our scope for optimizations and helps satisfy users who are willing to dedicate more beacons for better results. Besides, until BLE 5.0 is available, we need to simulate its broadcast capacity with multiple BLE 4.0 devices anyways.

We have developed a prototype of an image beacon system using a set of commercially available Estimote beacons [2], and developed an Android application that takes images of an object of interest along with user-specified requirements and constraints on broadcasting the image as inputs, generates previews of the image to be written, writes the image representation into a set of beacons, and reads the broadcasted image back. Figure 1 shows an example scenario where a user snaps photos of a gnome statue which he is interested in broadcasting. The smartphone application performs image processing on the phone to produce multiple versions of broadcast image. The user selects one of these compressed images that satisfies his requirements (e.g. available beacons, image quality, lifetime, and image loading latency). The user is allowed to change his requirements and the app immediately shows options for the best possible compressed images under those constraints. The application writes the image data into the beacon system and the image is broadcasted by the beacons. A reader application reads the broadcasted image and displays it on the phone.

We perform an in-depth evaluation of the beacon system. We describe a set of results showing the tradeoffs between system lifetime and image quality, when the image type and the number of beacons are varied. We also deploy an image beacon system indoors, and perform a user study in a real-world scenario in order to have a subjective measure of the quality of the received images, where a group of 20 participants are asked to identify objects from their beaconed images of various resolutions, and locate it among a set of similar looking objects in the real-world.

The main contributions of this paper are as follows:

- (1) To the best of our knowledge, we are the first to propose an image beacon system that uses multiple BLE beacons to broadcast color images over the BLE advertisement messages.
- (2) We have devised an image approximation algorithm that is tailored to the need of an image beacon system. We quantify the tradeoffs between the image quality and the device lifetime, and determine the best set of parameters, under the user-specified constraints on the number of beacons, latency, and expected system lifetime.
- (3) We have developed and evaluated a prototype of an image beacon system that broadcasts color images of various types (e.g., near-distance indoor and outdoor objects, road signs, and buildings). Our evaluation shows that one BLE 5.0 beacon would be capable of broadcasting good-quality images (70% structurally similar to original images) for a year-long continuous broadcasting, and both the lifetime and the image quality improve when more beacons are used.

2 PROBLEM FORMULATION

2.1 Generic Problem Setting

The problem is formally stated as: given an image x (where each pixel is represented by b bit) having the dimensions of $N \times M$, the number of available beacon devices K , the payload size of each beacon packet C bytes, the maximum allowable broadcast rate of R packets/sec, and the maximum allowable latency for an image T , the objective is to find an approximate representation of the image

\hat{x} so that the lifetime τ of the beacon system is maximized while the approximation ratio $\lambda(x, \hat{x}) \in [0, 1]$ of the image is high ($\lambda = 1$ means no distortion). Now, for a single beacon, the broadcast rate:

$$R = \left(\frac{bNM}{8C} \right) \frac{1}{T} \quad (1)$$

For K beacons, considering $\log K$ overhead bits for addressing the beacons, and K times more payload capacity:

$$R = \left(\frac{bNM + \log K}{8CK} \right) \frac{1}{T} \quad (2)$$

Both (1) and (2) are for undistorted images.

The lifetime τ of a BLE device depends on its inter packet interval and in general, $\tau \propto \frac{1}{R}$. Replacing R and incorporating approximation ratio λ into (2):

$$\frac{1}{\tau} \propto \left(\frac{\lambda bNM + \log K}{8CK} \right) \frac{1}{T} \quad (3)$$

The above equation relates the lifetime of an image beacon system and the approximation ratio of any image compression algorithm. Ideally, we look for an image approximation algorithm that achieves a sufficiently large λ for a reasonably high lifetime of the system.

2.2 Broadcast Capacity of Bluetooth LE

According to the BLE 4.0 specification, the maximum payload size C available in beacons is 30 bytes. The expected lifetime of BLE beacons depends on the inter-packet interval [5]. For example, a BLE 4.0 beacon would last up to 3.5 years, if a packet is sent at every second (i.e. $R = 1$). Therefore, for a beacon system to last for 3.5 years, its broadcast bandwidth cannot exceed 30 bytes/sec.

Recently, BLE 5.0 has been announced [1] to offer an 8X increase in broadcast capacity and a 2X increase in transmission speed. It is scheduled to be released in early 2017. In coming days when BLE 5.0 capable devices will be widespread, we expect to have a 128 byte sized payload and about 256 bytes/sec broadcast bandwidth.

2.3 The Case for Loss-Less Image Broadcast

The size of a typical 72×72 PNG image can be anywhere between 3–13 KB. Therefore, to transmit such an image, a BLE 4.0 beacon would require 191 – 832 broadcast packets, or alternatively, we would require up to $K = 832$ beacons to simultaneously broadcast different slices of an image. The latency of a complete image transmission cycle would be up to $T = 13.9$ minutes for a single beacon, or 1 second for a set of 832 beacons.

When BLE 5.0 beacons will replace 4.0, the transmission latency will drop to 52 seconds for one beacon, or 1 second 52 of them. Therefore, without compressing the image content, even the new BLE 5.0 beacons will not be able to support a fast image beacon system with a reasonably small number of beacons.

2.4 The Case for Compressed Image Broadcast

If standard image compression algorithms could generate compressed images that meet the size and quality requirements of an image beacon system, the problem would have been already solved. But the fact is, even the best of existing image compression methods, such as JPEG/JPEG2000 and PNG, are not capable of optimizing

for both quality and size at the same time. Figure 2 illustrates that JPEG/JPEG2000 generates extremely poor quality images given a size requirement of 300 bytes even for a very low-resolution (64×64 pixels) image. On the other hand, to have a compressed image of acceptable quality (having a minimal useful visual information to the viewer), JPEG/JPEG2000 takes about 2K bytes.

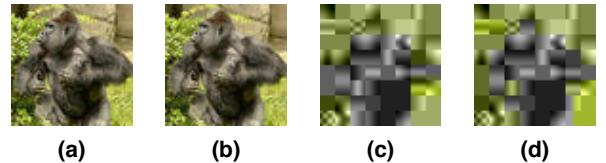


Figure 2: A 64×64 resolution image compressed in high/low quality settings using JPEG/JPEG2000: (a) JPEG high quality, 1963 bytes (b) JPEG2000 high quality, 2026 bytes (c) JPEG lowest possible quality, 738 bytes (d) JPEG2000 lowest possible quality, 391 bytes.



Figure 3: Two types of 64×64 resolution image compressed in PNG (a) from natural scene, 12112 bytes (b) JPEG2000 high quality, 1012 bytes. PNG is good for handling images with large uniform color regions.

PNG and Vector Graphics image, on the other hand, have the potential to generate a smaller compressed image that *may* fit our constraints. However, these codecs generate smaller images only if the input image is of a specific type – such as an image containing a few regions of uniform colors like a cartoon drawing, or when the shape is not complicated. This is illustrated in Figure 3. In general, PNG and Vector Graphics image encoding do not meet the requirements of an image beacon system that broadcasts color images taken by a smartphone user.

3 IMAGE BEACON SYSTEM OVERVIEW

In a typical usage scenario of the proposed image beacon system, a user at first takes pictures of an object with his smartphone's camera. A smartphone application analyzes the image (which may contain objects, portraits, scenes, shapes, signs, and/or text), identifies semantic regions on it, and processes each region differently to produce a compressed version that satisfies the beacon system's requirements such as the number of available beacon devices, maximum allowable loading time, and lifetime. The user is also shown an interactive preview of the image so that he can verify it, as well as relax/constrain the system requirements. Finally, when he is satisfied with the preview, the image is written into the image beacon system. The beacon system would then broadcast the image periodically over BLE, and any other smartphone user would be able to receive that broadcast and see the image on their phones.

3.1 System Design Choices

The design choices we made in developing an image processing algorithm for the proposed image beacon system are as follows:

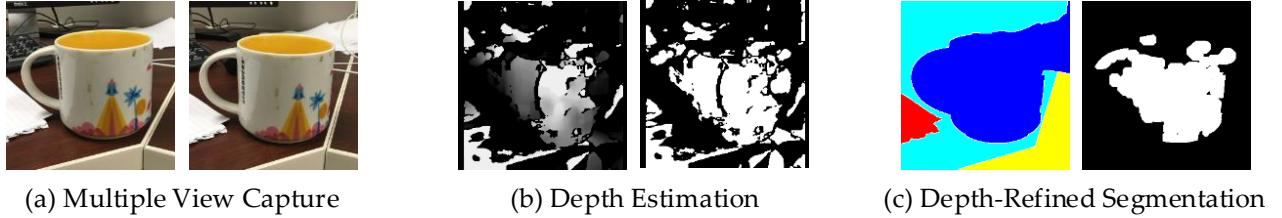


Figure 4: Multiple views of a scene are used to estimate the depth map. Combined with standard image segmentation, this can identify the pixels of an image that may be of more interest than the rest, e.g. a foreground object.

- (1) First, our custom image compression technique is designed to work for images taken with a smartphone. We assume that the phone has an on-board IMU in it. The final compressed image will be a color image with a lower resolution, such as 64×64 pixels.
- (2) Second, we make a reasonable assumption that the image to be compressed is linked to a real-word “thing” like a near-distance object, a road sign, or a building – which has one or more regions of interest that a person who took the picture wants to preserve with a higher priority than the rest. By exploiting this, we design an image encoding technique that prioritizes foreground information preservation during image compression, so that the most important information in the image is delivered under a given size constraint.
- (3) Third, under a very tight budget for the final image size, any image compression algorithm would distort the original image – which is reflected in different ways such as lacking boundary details, increased noise, changed colors, or removal of texture. We introduce the concept of *adaptive encoding* that applies different encodings to different regions of an image based on the image content (e.g., a road sign vs. a t-shirt), image regions (foreground vs. background), and what a user would prefer to preserve (e.g., texture or true color). The proposed compression algorithm should employ an adaptive approach that applies the most suitable encoding technique for different image types and region types, so that an optimal compression strategy is chosen for a given image based on its content.
- (4) Fourth, since both capturing an image and writing the compressed version into the beacon devices involve the smartphone user in the loop, we provide an interactive user interface in order to guide the user in taking pictures, and to preview and select the desired image under a given set of system constraints.

3.2 Image Processing Pipeline Overview

For a given set of user-defined beacon system requirements, the overall image processing and compression pipeline (Figure 5) consists of four basic stages: multiple view capture, depth estimation, depth-refined segmentation, and image compression. These steps are briefly described in this section, and elaborated in detail in the subsequent sections.

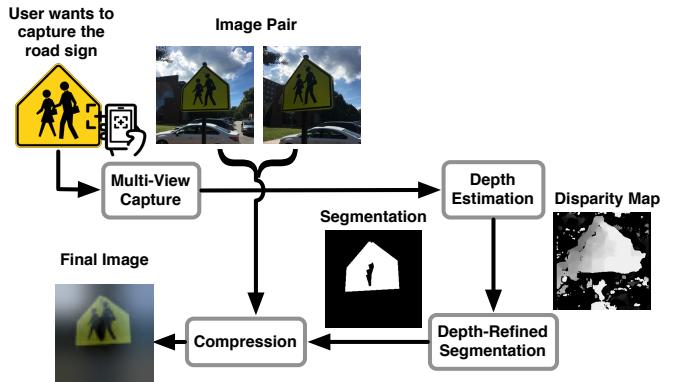


Figure 5: Image processing stages.

- (1) *Multiple View Capture:* The proposed system requires a user to capture two or more views of an object – which helps at a later stage when the depth map is generated. Estimating pixel depths using a pair of images takes about 2 seconds on a mobile device. Because processing too many images would be time consuming, a careful selection of views (e.g. images having adequate overlaps) makes a difference. Figure 4(a) shows two views of a mug that have enough overlap to create a depth map. To guide the user and to select the best pair of images for depth estimation, we leverage IMUs of the smartphone. The algorithm is described in Section 4.
- (2) *Depth Estimation:* Depth of each pixel is estimated by finding and matching corresponding ‘feature points’ (e.g., corners and edges) in two or more images. The matched points are then used to generate the camera relative geometry, so that the depth of every pixel can be estimated. Figure 4(b) shows two depth maps of the same image. The left one is the computed depth map, and the right one is thresholded to separate the background from the foreground pixels. However, due to lack of enough views, low resolution, and inaccuracies in estimation, depth map alone is not sufficient to segment semantic regions in an image. Depth estimation from multiple views is discussed as part of Section 5.1.
- (3) *Depth-Refined Segmentation:* Like depth map, color/textured-based image segmentation algorithms often fail to identify semantically different/similar regions in an image. For example, the left image in Figure 4(c) is the result of applying marker-controlled watershed segmentation [12] on

the original image. When we overlay this with the depth map, we obtain a better segmentation, which performs a much better job in isolating the mug from the rest. This step is inspired by one of our earlier work [10] that used RGB and depth images from Kinect sensors. In this work, we use only images to estimate the depth (previous step) and then apply this step to get the final segmentation. The details of this step are in Section 5.2.

- (4) *Image Compression:* The image compression stage takes both an image (for texture and content information) and its segmentation map (for semantic region information), and produces the best quality image under the user-specific constraints of the beacon system. Until the resultant image size satisfies the system requirements, the algorithm gracefully degrades the quality of different semantic regions, starting from the least important one (e.g., the background). This step is described in detail in Section 6.

4 MULTIPLE VIEW CAPTURE

4.1 Need for Multiple Views

The first step of our proposed image processing pipeline is to guide the user in capturing two or more views of an object of interest. Further down the pipeline, these images are used to estimate the depth information of each pixel, so that an image can be segmented into background and foreground regions, prior to applying appropriate region-specific encodings.

An alternative to using multiple views is to apply standard image segmentation algorithms [7, 11, 12, 19] on a single image. These algorithms group adjacent pixels of an image based on information derived from pixel intensity in various ways. However, in order to obtain a sufficiently accurate segmentation for the proposed system, we require computationally expensive algorithms, such as convolutional/recurrent neural networks [19], which are not suitable for running on smartphones and does not produce results in real-time.

4.2 Challenges with Multiple Views

Even though smartphones cannot run state-of-art image segmentation algorithms in real-time, many other computer vision techniques, including depth estimation from two views, can be implemented on them. In order to get a sense of their real-time performance, we used OpenCV library for Android to compute the depth map for a pair of 400×400 pixel images. It took about two seconds for the algorithm to finish on a Nexus 5 phone. This gives us the lower limit for depth map computation, which may only happen if the user is well-trained and knowledgeable to know which views or camera poses would produce the most effective depth map.

A good pair of images is critical in generating a good depth map. However, the finding of a good pair of images depends on many factors. The most important of which is a suitable difference in view angles. It also depends on the distance between the object/scene and the smartphone. Furthermore, there are other factors such as lighting, texture and shapes of the image.

If a real-time depth estimation system could display the current depth map as the user takes images, it would be easier for him

to generate a good pair of images for depth estimation. However, depth estimation does not run in true real-time on most smartphones. In absence of a real-time feedback, a user has to take the trial-and-error approach, i.e., he has to take two images, look at the result after two seconds, and then repeat the entire process until the result looks good. This may lead to a very long time in just taking the right photos, and result in a non-smooth user experience.

4.3 IMU Assisted View Capture

To address this problem, we designed a method to make use of the inertial measurement unit (IMU) of the phone to shorten the image capture time and to improve the user experience. We adopt a machine-learning based approach.

In the offline training phase, we use a smartphone to capture video/image of multiple indoor and outdoor objects. We also keep record of the IMU values for each captured image. The IMU data consists of δx_r , δy_r and δz_r , which represent the components of the difference vector between the rotation vectors between a pair of images. The IMU data also contains x_a , y_a and z_a components of acceleration when taking an image. After this, we run depth estimation for all pairs of images and estimate its accuracy by comparing the result against a manually generated ground-truth segmentation.

The segmentation accuracy is measured in terms of *intersection over union* (IoU), where *intersection* is defined as the area of intersections between foreground regions of two segmentations, and *union* is defined as the set of pixels either marked as foreground in the testing segmentation or in the ground truth segmentation. For a segmentation that is identical to the ground truth, IoU equals to 1.

Using IoU values as the variable Y, and the IMU data for the corresponding pair of images as the variable X, where $X = [\delta x_r, \delta y_r, \delta z_r, x_a, y_a, z_a]$, we create a data set for many pairs of images, and then train a regression tree model to learn the relationship between the change in IMU values between a pair of images and an expected quality of depth segmentation.

During the online phase, when the user is taking images for depth estimation, the trained regression tree model keeps track of current IMU readings and gives hints about if current view is a good choice, given the already taken photo(s).

5 DEPTH ESTIMATION AND SEGMENTATION

The second and third stages of the proposed image processing pipeline are described together in this section.

5.1 Depth Estimation

In this step, the depth of each pixel is estimated from a pair of images. We use a standard algorithm [6] that at first estimates the ‘disparity’ between the corresponding points on two images, and then estimated depth from disparity. For example, if a point P_1 on the first image and a point P_2 on the second image correspond to the same point P on the actual physical world object, then $(P_1 - P_2)$ is called the disparity between them. Depth of a pixel is, in general, inversely proportional to its disparity. This is based on the principle that points in the scene that are closer to the camera will have larger disparity, and points that are very far away will be effectively at

the same or very close location on both images. Hence, finding the depth map is essentially equivalent to finding the disparity map.

The disparity map is generated by the *stereo matching* algorithm described in [6]. The goal of the algorithm is to find matching pixel blocks in a pair of images. This is also called the *correspondence problem* as it looks for the pixel coordinates on the image pair that corresponds to the same world point. The disparity map is computed based on the matching result. The disparity map is a gray-scale map, where the intensity directly corresponds to depth.

5.2 Depth-Refined Segmentation

The depth estimation algorithm groups pixels purely based on depth. It may not group pixel regions even if the regions share a common appearance pattern. As a result, even an accurate depth map tends to contain holes in foreground regions and isolated, incorrectly marked, small, bright regions in the background. Therefore, it is necessary to introduce other types of information derived from the image to obtain a cleaner and better segmentation. The refinement process is described as follows:

Thresholding: At first, a binary segmentation of the depth map is obtained by applying a threshold on depth values.

Combining with Watershed Segmentation: We combine the depth-based segmentation map with another image segmentation method which groups pixels into several connected large regions and is computationally inexpensive to run on a mobile device. Watershed segmentation algorithm fits these requirements. However, the traditional watershed segmentation tends to generate an over-segmented result. Hence, we adopt the marker-controlled watershed segmentation, which uses mathematical morphology operations to pre-process an input image to avoid over-segmentation [6].

At first, the input image is converted into grayscale. Then we run the marker-controlled watershed segmentation on the grayscale image. A successful segmentation contains more than one segmented regions to separate foreground from background in the image. To determine which region belongs to the foreground, we apply a voting approach: the region that includes the highest number of common pixels with the foreground region in the depth-based segmentation map is labeled as foreground. Here we denote the number of common pixels as N , where N is defined as:

$$N = \max_i C(W_i, D) \quad (4)$$

Here, i is the index of a region, and W_i is the corresponding region. D represents the foreground region in the depth-based segmentation. $C(,)$ computes the number of common pixels between two regions. If there is another region W_j for which, $C(W_j, D) \geq 0.8 \times N$, then it is also considered as foreground. This procedure iterates until no more regions can be added to the foreground. Having two result segmentation maps, we produce a final map by labeling pixels that are considered foreground in both maps as foreground.

Final Refinement: Finally, we perform a pixel-level refinement process. We remove all connected foreground pixel regions with size less than 1000 pixels because regions of this size tend to be a background region. Then we apply two common mathematical morphology operations erosion and dilation, to clean out any remaining bright pixel islands in the background and to expand the

foreground regions, respectively. Lastly, we enforce the accuracy of the segmentation boundary by combining the result with the labeled watershed foreground region.

6 IMAGE COMPRESSION

The last stage of the proposed image processing pipeline is the image compression step. Using the segmentation information from the previous stage, this step encodes different segments of an image using different encoding techniques. The overall goal is to make sure that the resultant image fits the storage requirement of a beacon system, while making sure that the foreground regions are the least affected by during the compression process.

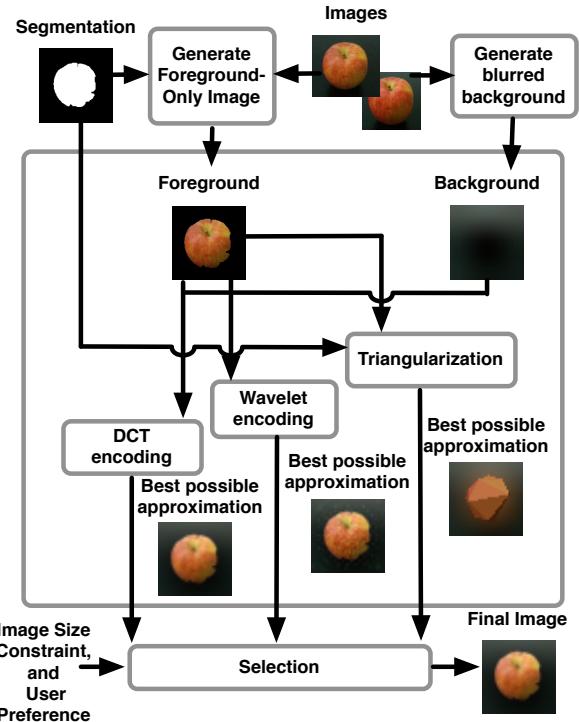


Figure 6: Image compression details.

We propose three encoding options for image compression – *discrete cosine transform* (DCT) and coefficient reduction, *wavelet transform* with coefficient reduction, and *foreground texture triangulation*. All three are applied on the input image and finally the one that produces the best quality image is chosen as the output. The effect of applying different encodings is illustrated in Figure 7.

6.1 Discrete Cosine Transform Encoding

Discrete Cosine Transform (DCT) is a widely adopted image encoding technique. We integrate DCT encoding into our compression system as a baseline encoding option. The benefit of using DCT is that – by reducing low frequency coefficients of an image (in the DCT transformed space), the resultant compressed image's *appearance* details are removed first, while its *global shape* is preserved. This is useful in usage scenarios when a user wants to preserve the shape of an object in the image more than its detailed appearance.

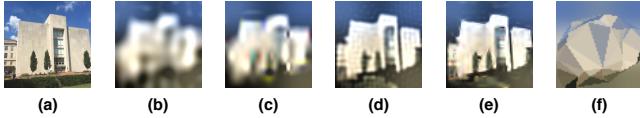


Figure 7: 64x64 resolution building image compressed in high/low quality settings using our customized DCT/Wavelet/Triangle encoding: (a) Original image, (b) DCT 342 bytes, (c) Wavelet high 360 bytes, (d) DCT 1114 bytes, (e) Wavelet 1098 bytes, and (f) triangularization 366 bytes. For a similar compressed image size, DCT preserves less details than Wavelet method. But for low quality settings (about 350 bytes), Wavelet-encoded images have strange color block defects. Triangularization failed to preserves the information in the original image.

At the beginning of encoding, we generate a foreground image by using the segmentation map to set the background pixels' intensity of the input image to zero. The image is then down sampled to 64×64 pixels. To further reduce the size, we make use of the fact that the quality of an image depends more on its brightness information than its color information. During the encoding, at first, an image is transformed into the YUV space, where Y represents the brightness information and U/V represents the image color information. We further down sample U and V channels into 32×32 pixels, while keeping the resolution of Y channel intact. Then DCT is applied on all three channels, followed by a data-reduction step that sets high frequency components to zero. Finally, the resultant frequency coefficients are compressed using gzip. The data is sent to the beacon along with the blurred background image, which is also encoded using DCT. The size of a DCT compressed image, S_{DCT} can be expressed as:

$$S_{DCT} = B(g(d)) + B(g(b)) \quad (5)$$

Where, $B()$ denotes the bit length, $g()$ denotes the gzip encoded data size, d is the DCT transformed (reduced coefficient version) foreground information, and b is the coefficient of the DCT transformed (blurry) background image data.

At the receiving end, a broadcast image is recovered by superimposing the foreground image and the background image. Note that, the background image needs to have the foreground pixel intensities set to zero, before it is down sampled. This is done to make sure that the foreground pixel intensities are not added up twice.

A limitation of this above approach is that, for an image with a uniform dark background and a foreground having more details, DCT may yield a ringing artifact close to the sharp edges, especially in a low quality setting. This problem can be addressed by switching to using wavelet.

6.2 Wavelet Encoding

Wavelet is the second image encoding method that we integrate in the adaptive image encoding process. When compared to DCT, wavelet tends to better handle images whose backgrounds have an uniform intensity. Similar to DCT, the best information reduction parameters are sent to the wavelet encoding module in order to generate the highest quality image under a given storage limit. We adopt global thresholding of the wavelet coefficients and Huffman encoding, based on the method described in [13].

Similar to DCT, prior to encoding the foreground image, its background pixel intensities are set to zero to obtain the wavelet data. The down sampled background image (with zero intensity foreground) is also sent along with the wavelet data. At the receiving end, wavelet coefficients are inverse-transformed to generate the foreground image and then superimposed on the background image to render the final image.

The size of a wavelet compressed image S_w is as follows:

$$S_w = B(H(w)) + B(g(b)) \quad (6)$$

where, $H()$ denotes the Huffman encoded data, and w denotes the reduced wavelet coefficients on the wavelet transformed foreground image. All other symbols carry the same meaning as discussed in the previous section.

The weakness of wavelet encoding is that, for limited storage requirements (< 500 bytes), a wavelet encoded image may have unrealistic texture patches after decoding.

6.3 Triangularization-Based encoding

Both DCT and wavelet encodings blend the geometric information and texture information of the foreground image. However, for some cases, an accurate texture information and a fine-grained boundary representation are not necessary. For example, a “football” image’s foreground texture and shape could be decoupled. For a viewer to understand that the image is about a “football”, a repeated patch of a football’s surface texture along with an approximate “football” shape information would suffice. Both DCT and wavelet would encode too much redundant information for such an image. To address this, we designed a triangularization-based image encoding method, which consists of three stages: triangularization, triangle reduction, and color/texture filling. These are illustrated in Figure 8, and are described as follows:

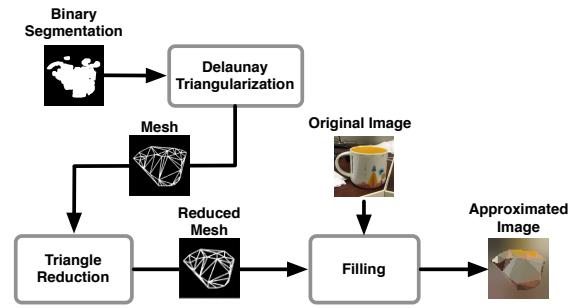


Figure 8: The process of Triangularization-based encoding.

- (1) *Triangularization:* Given an input image with the foreground/background segmentation, the first step is to generate a binary *boundary map* from the segmentation map, in which, only the pixels on the segmentation boundary have an intensity of 1. A Delaunay triangularization [4] is performed on the boundary map. The parameters are chosen to produce a high number of triangles to capture boundary details.

- (2) *Triangle Reduction*: Having a set of triangles, the next step is to reduce the number of vertices, iteratively, one vertex at a time. For this, we compute the sum of distances for every vertex from its nearest 3 neighbors, and then remove the one with the minimum sum of distances. The intuition behind this process is that – a region of vertices group together densely because of the non-smooth boundary in the boundary image. Since the goal of removing vertices is to reduce the details and preserve the general shape information, we should pick a vertex from dense regions.
- (3) *Color/Texture Filling*: We provide two options for filling a triangle – with texture or with a single color. For texture, we choose to fill all triangles with a limited set of textures which are derived from regions surrounded by each triangle. To reduce the number of textures, we take an average of textures from different triangles. Fig. 9 shows the process. For each triangle, we transform it to a fixed-size triangle by an affine transform, and then compute one or two average textures for all triangles. For the case of two textures, we apply k-means algorithm.

The size of the compressed image using triangularization encoding with color filling S_{Tc} , and with texture filling S_{Tt} are as follows:

$$S_{Tc} = B(g([v, c, f]) + B(g(b)) \quad (7)$$

$$S_{Tt} = B(g([v, c, i]) + B(g(t)) + B(g(b)) \quad (8)$$

Here, v denotes the location of the vertices, c denotes the connectivity list, f denotes RGB color values, i denotes the texture index, and t denotes the reduced DCT coefficients on the texture patch transformed using DCT. All other symbols carry the same definition as in previous sections.

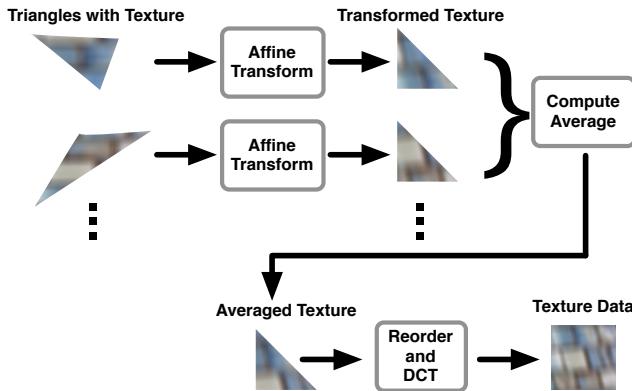


Figure 9: Triangle texture averaging process.

7 EMPIRICAL EVALUATION

In this section, we describe a series of empirical evaluations. First, we evaluate the performance of IMU-guided multiple view capture and depth-refine segmentation. Then our image compression approach is compared with JPEG encoding. After that, we describe a set of results that quantifies the trade-offs between the beacon

system lifetime and image quality, when the image type and number of beacons are varied. We also perform a full system evaluation involving real users, which is described in Section 8.

7.1 Experimental Setup

In all of our experiments, we have used Estimote model REV.F2.3 Radio Beacon [2] having a 32-bit ARM Cortex M0 CPU, 256 KB flash memory, 4 dBm output power, 40 channels (3 for advertising), and 2.4-2.4835 GHz operating frequency. We vary the BLE broadcast interval for a beacon between 100 ms to 2,000 ms. However, an encoded image (broadcasted from multiple beacons) reaches a user's device in less than 1 second. The transmission power is set to -12 dBm, which limits the range of each beacon to about 30 meters. The image writing and reading application runs in a Nexus 5 smartphone having a 2.26GHz quad-core Qualcomm Snapdragon 800 processor, 2 GB RAM, BLE v4, and runs Andriod 6. We mimic BLE 5.0 broadcast packets by a set of rolling BLE 4.0 packets. The rolling mechanism is implemented by configuring the Estimote Location beacons to broadcast customized advertising packets. The customized data is received from an Android compatible LightBlue Bean device [3] via the beacon's GPIO, configured as an UART interface.

We use four types of images in our experiments: images containing road signs, common indoor and outdoor objects, and buildings. Examples of these images are shown in Figure 10. Indoor object images are taken from a 50cm - 150cm distance. Outdoor objects and signs are taken from 2m to 5m distance. Building images are taken from far. All images are cropped into square shaped images, and are down-sampled to 288 × 288 pixels prior to compression for a fast disparity map and watershed result computation. Each image is down sampled to 64 × 64 pixels before writing into the beacon system.

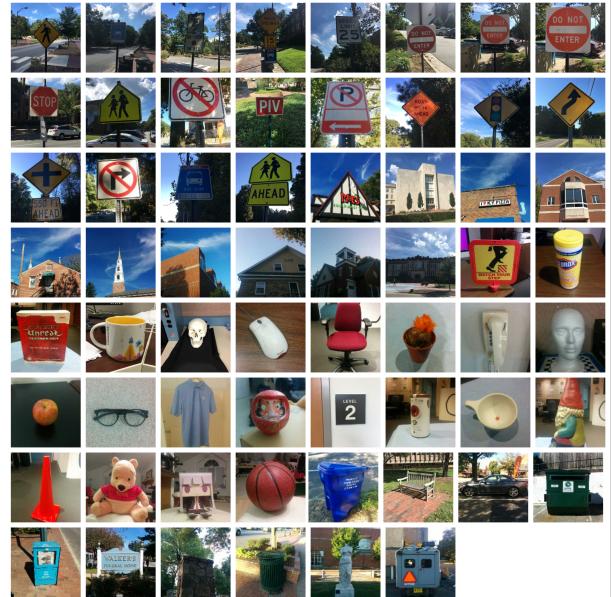


Figure 10: Test images used in the empirical evaluation.

The two main metrics that are used in the experiments are structural similarity (SSIM) scores, and device lifetime in months. We

measure these two under different conditions and show their trade-offs. The structural similarity scores are used to measure the quality of the produced images when compared to the original ones. The device lifetime is estimated from its relation to a beacon's transmission frequency. Before each experiment, we program the beacons to set a transmission frequency and use the corresponding estimated device lifetime (as reported by the Estimote beacon API) in our experiments.

7.2 Performance of IMU-Guided Multiple View Capture

We evaluate the accuracy and time to capture multiple views with and without the guidance of IMU. Recall that, the regression tree model takes IMU data and predicts if the current smartphone positioning is good for taking an image for depth estimation, given an already taken image. In the evaluation, for each test, a set of images and their corresponding IMU readings are recorded. If the predicted image produces the best segmentation result when compared to other images in the set, we record this test result as a ‘hit’, otherwise, a ‘miss’. The accuracy of prediction is determined by the ratio of ‘hits’ to total tests. A total of four tests are performed separately for indoor and outdoor objects. For indoors, there is 1 miss over 4 tests, whereas for outdoors, the model correctly identifies the best image pair for all tests. The model for outdoor objects is more robust since the smartphone’s rotation angle shows smaller variations and usually the phone is held vertical. Indoors, users often take pictures of an object from its above, from its below, or from the same height, which results in a large variation in angles.

Figure 11 shows the expected time for finding a good pair of images with and without the assistance IMU. We observe that the average time to obtain the depth map from a pair of images takes about 2 seconds, which includes the time to focus, raw image processing, and disparity map computation. Since the IMU-guided system predicts a good pair in real time, it significantly decreases the time from 8 seconds to 2-2.5 seconds.

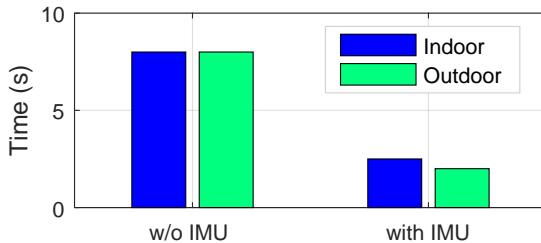


Figure 11: Performance of IMU-guided view capture.

7.3 Performance of Depth-Refined Segmentation

We compare our depth-refined segmentation method against a segmentation approach that is purely based on the depth map. We use intersection over union (IoU) as the evaluation metric. We tested our segmentation method on 20 images of four types: road signs, buildings, indoor and outdoor objects. For each set, we manually

generate the ground truth segmentation. Figure 12 shows the result. We observe that our segmentation technique, which combines depth information and watershed segmentation information, outperforms depth-only approach by up to 26%. Our method performs better in cases where the object and the background has larger difference in depths, e.g., signs and indoor objects.

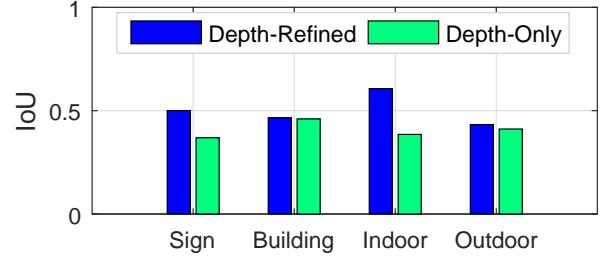


Figure 12: Performance of depth-refined segmentation.

7.4 Comparison with JPEG

In this experiment, we compare our proposed image encoding method with JPEG. For an in-depth illustration, we use the picture of an apple as our test image. We compare four different options of our adaptive image encoding methods in the comparison, i.e. DCT-based, wavelet-based, triangularization with color filling, and texture filling. We measure the quality of a compressed image using Structure Similarity (SSIM) [16]. A SSIM score ranges from 0 to 1, and two identical images have the best SSIM score of 1. We compute the SSIM value between every compressed image and the original image. We plot the SSIM versus compressed image size in bytes in Figure 13.

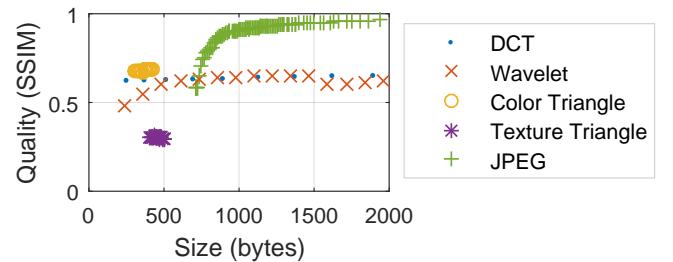


Figure 13: Image quality versus image size for different encoding methods.

The result suggests that JPEG is able to generate a compressed image with a higher quality (SSIM close to 1). However, such a high quality image has the size larger than 2 KB. When JPEG is set to compress an image into lower than 1KB, the image quality drops sharply. On the other hand, our image compression method is based on foreground/background separation. The background in a compressed image is always blurred. This makes our method impossible to get a high SSIM score larger than 0.8. But our method is able to allocate bits more efficiently under a tight space constraint.

This makes our method (when using wavelet/DCT encoding) outperform JPEG when the compressed image size is about 800 bytes. The plot also shows that our method can compress an image into as low as 240 bytes (left most point on the wavelet method curve), which is impossible with JPEG, even in its lowest quality setting.

Triangularization with color filling option performs the best for this test case, and the compressed image size stays less than 500 bytes, while triangularization with texture filling encoding fails to generate a good compressed image.

Besides JPEG, we also studied several other image compression techniques [9][17][18]. But none of these yield suitably small sized images.

7.5 Effect of Image Type

In this experiment, we test how our image encoding method's performance changes as we vary the type of input images. We consider four categories of images, i.e., road signs, buildings, indoor objects, and outdoor objects. For each image, our algorithm selects the best compressed image from different versions of the adaptive encoded images based on SSIM score. Then we compute the average image quality for a given lifetime for each category of images. We simulate a BLE 5.0 beacon in this experiment to store and read the images. We limit our system to deliver the image data within 0.5 second from one beacon.

The result shown in Figure 14 suggests that, on average, indoor object images achieves the best quality over all types. The building images are best approximated if the beacon system broadcasts packets at a higher frequency, sacrificing the system life time. The road sign images have relatively lower quality than other images, but they tend to last for up to 28 months on a single battery.

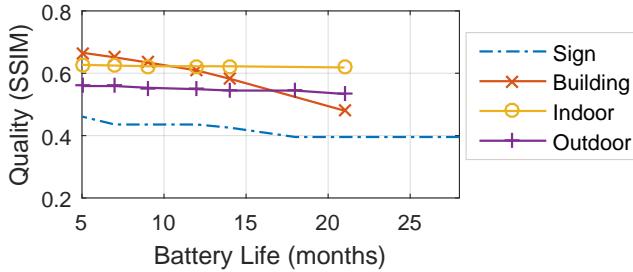


Figure 14: Image quality versus beacon battery life for different image types.

7.6 Effect of Number of Beacon Devices

We explore the impact of the number of beacons on image quality. Since each beacon in a system of beacons can be set to broadcast different parts of an image, the more beacons we have, the better quality images we can generate by utilizing the additional space. This is based on the assumption that the image loading time and the device lifetime requirements are fixed.

In this experiment, we vary the number of beacons, and record the SSIM of the best quality compressed image generated by our system. The best image is chosen from the adaptive-encoding results with the highest SSIM score.

We plot the SSIM scores for various expected system lifetime in Figure 15. The experiment results suggest that as the number of beacons is increased from 1 to 3, for the beacon system to have an expected lifetime of e.g., 32 months, the quality of produced images also increases from 0.45 to 0.69.

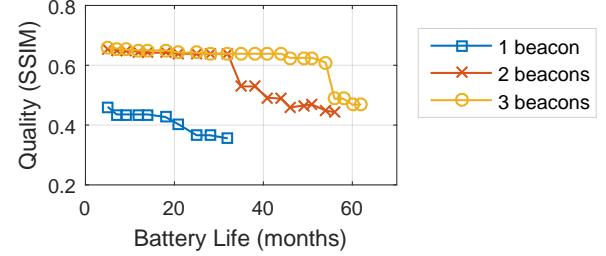


Figure 15: Image quality versus device lifetime for various number of beacons.

8 REAL DEPLOYMENT

We deploy an image beacon system in an indoor environment to evaluate the ability of image beacons in delivering visual information of various real-world objects. An approximated image stored in the beacon system contains an object's shape as well as its texture information. The goal of this deployment experiment was to understand how an image beacon system delivers these two kinds of information of an image.

To have a subjective measure of the performance of an image beacon system, we conduct a user study involving 20 participants. Each participant is given a smartphone that receives image broadcasts from four different beacons placed inside a room. The goal of the user is to identify the objects in all four images, and then locate them in the room.

The four broadcasted images are of an apple, a chair, a text book and a computer mouse. For each image, we compress it under three size constraints: 256 bytes, 512 bytes, and 768 bytes, and obtain the best quality image produced by our compression algorithm under the constraint. The images compressed in three levels along with the original image are shown in Figure 16. We choose these sizes to mimic 1, 2, and 3 BLE 5.0 beacons, respectively; but our actual implementation used a rolling mechanism with BLE 4.0 (connected with an Arduino), as described earlier. Each user is progressively shown a better quality image until he is able to identify the object in the room.

We make sure that there are at least 3 objects in the room that are similar to the one that the user is looking at on his phone. To test if an object's texture details are preserved, for the apple test case, we put another apple having a green/red mixed color and an orange next to it. For the book test case, we put two other similar sized books next to it that have different covers. To test if an object's shape details are preserved, for the chair test case, we add another two chairs of the same color. For the mouse test case, we add an iPhone and a mac mouse. Figure 17 shows photos of the objects along with the objects that we added to introduce confusions.

The experiment results are shown in Figure 18. For each size limit, we plot the number of correct guesses by our participants

for various categories of images. We observe that, as expected, when the image size limit is larger, participants tend to perform better. Even with the lowest size, at about 50% images were always guessed correctly by the participants. Therefore, with 3 or more beacons, the image quality of our system is high enough to let people distinguish very similar objects.

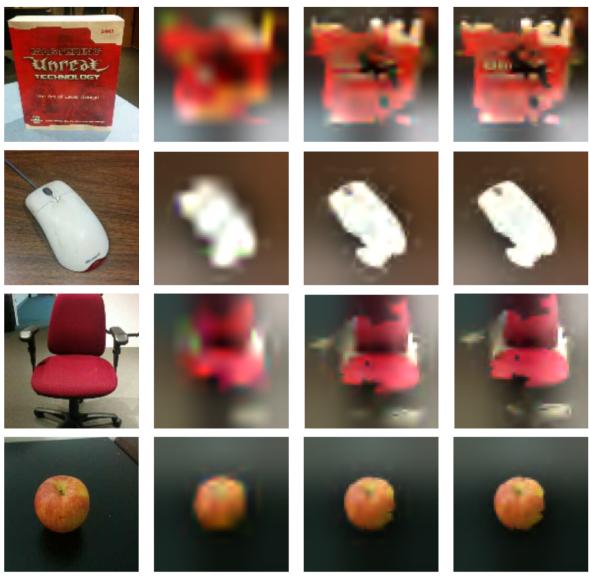


Figure 16: Test images compressed in three quality levels.



Figure 17: Photos of the object used in the experiment.

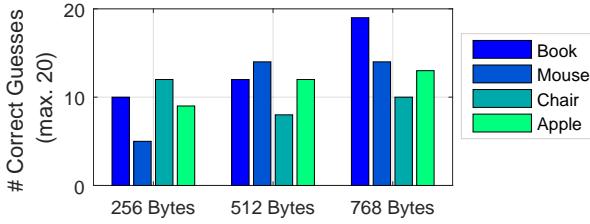


Figure 18: Responses from user study.

9 DISCUSSION

Our proposed image beacon system only considers stationary objects. This is an inherent problem of any depth estimation technique. In such case, we have to resort to texture or color based segmentation.

Our IMU-based prediction algorithm uses a regression tree model. Its prediction accuracy is lower indoors than outdoors. A robust model, may train a separate regression tree for different cases, such as one model for taking images on objects below the phone, one model for objects on the same height to the phone.

A property of marker-controlled watershed segmentation algorithms is that they always generate a clean segmentation result. Our combined segmentation method does not fully exploit this feature. We could further enhance the power of combining depth estimation and watershed results by deploying a smarter foreground region growing method.

Animated images (e.g., GIFs) are not supported by the system. However, with additional optimization for time-domain redundancy, we believe it is possible to develop image beacons that supports animation.

10 RELATED WORK

Previously, a binary image beacon system [14] was developed to enable BLE beacons to store and broadcast binary images. The limitation of the system is that it only supports a limited category of binary images such as hand-written characters, shapes, and symbols.

A rich set of image segmentation methods exists in the literature. In the past decade, more modern techniques involving machine learning have been invented. State of art neural-network based method [19] can achieve a very accurate result (highest score 90.4 on IoU evaluation on airplane type testing data). The method is based on neural-network and conditional random field. However, this method is computationally expensive to run on a cellphone. Otsu's method [11] is based on finding a separation on image pixel intensity histograms, which does not take care of local image structure. Gabor filter segmentation is based on finding edges in an image using Gabor filtering. Marker-controlled watershed methods [12] use mathematical morphology to pre-process the data, followed by a watershed segmentation. This avoids over segmentation, which is a weakness of traditional watershed method.

The notion of foreground / background information can be obtained by disparity estimation with semi-global matching [6]. The method enforces smoothness in the neighbor matching process to reduce matching errors.

Shapiro [15] developed an image encoding technique named embedded zerotree wavelet (EZW) encoding, which is computationally expensive and slow. Said and Pearlman [13] developed a better wavelet-based image encoding method based on set partitioning in hierarchical trees. This method gives similar image compression performance regarding quality and size, and the same time it achieves a faster computation. We have used this method in our system.

Lu et al. [8] introduced a piece-wise linear image encoding method using surface triangulation. Their triangulation algorithm fits the image surface in a top-down manner. The idea is to apply a constrained resource planning to allocate the least amount of triangles while achieving a small image approximation error. Their experiment result shows that the triangulation method compresses images with a compact code length with a guaranteed error bound. But their method's target is to achieve

near lossless compression, while our ImageBeacon system is a lossy compression method in general.

11 CONCLUSION

In this paper, we described a system involving beacon devices and smartphones with BLE receiving functionality. Our system allows a user to generate an approximation of an input color image, write the approximated image into the very limited beacon device storage, and receive compressed images from beacon devices' broadcast packets. The main contribution of this work is the overall system construction, the adaptive encoding image compression method, the evaluation of various parameters of the system, and quantifying the trade-off between image quality and beacon battery lifetime for our image compression method. Our work widens the usage of the energy efficient, long-lasting beacon devices by allowing easy storage and access of custom image data in scenarios where there is no Internet connection.

Our system will perform at its best with the beacons that adopts the upcoming Bluetooth 5.0 standard. A future work is to evaluate the different aspects of the system performance as the Bluetooth 5.0 is released and gets popular. Moreover, a ‘fat-beacon’ standard is under development at Google, that allows an even higher broadcast transmission capacity for BLE beacons. The goal of that standard is to equip beacon devices with the ability to broadcast basic web contents to smartphones in absence of the Internet connectivity. It will be meaningful to study the application of our image beacon system combined with a fat beacon.

REFERENCES

- [1] 2016. Bluetooth 5.0 Press Release. <http://tinyurl.com/hvrosf5>. (2016).
- [2] 2016. Estimote Beacons. (2016). <http://estimote.com/>.
- [3] 2016. Lightblue Bean. punchthrough.com/bean. (2016).
- [4] Mark De Berg, Marc Van Kreveld, Mark Overmars, and Otfried Cheong Schwarzkopf. 2000. Computational geometry. In *Computational geometry*. Springer, 1–17.
- [5] Artem Dementyev, Steve Hodges, Stuart Taylor, and Joshua Smith. 2013. Power consumption analysis of Bluetooth Low Energy, ZigBee and ANT sensor nodes in a cyclic sleep scenario. In *Wireless Symposium (IWS), 2013 IEEE International*. IEEE, 1–4.
- [6] Heiko Hirschmuller. 2005. Accurate and efficient stereo processing by semi-global matching and mutual information. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, Vol. 2. IEEE, 807–814.
- [7] Jonathan Long, Evan Shelhamer, and Trevor Darrell. 2015. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3431–3440.
- [8] Tianyu Lu, Zisheng Le, and DYY Yun. 2000. Piecewise linear image coding using surface triangulation and geometric compression. In *Data Compression Conference, 2000. Proceedings. DCC 2000*. IEEE, 410–419.
- [9] Sherif A Mohamed and Moustafa M Fahmy. 1995. Binary image compression using efficient partitioning into rectangular regions. *Communications, IEEE Transactions on* 43, 5 (1995), 1888–1893.
- [10] Shahriar Nirjon and John A Stankovic. 2012. Kinsight: Localizing and tracking household objects using depth-camera sensors. In *2012 IEEE 8th International Conference on Distributed Computing in Sensor Systems*. IEEE, 67–74.
- [11] Nobuyuki Otsu. 1975. A threshold selection method from gray-level histograms. *Automatica* 11, 285–296 (1975), 23–27.
- [12] K Parvati, Prakasa Rao, and M Mariya Das. 2009. Image segmentation using gray-scale morphology and marker-controlled watershed transformation. *Discrete Dynamics in Nature and Society* 2008 (2009).
- [13] Amir Said and William A Pearlman. 1996. A new, fast, and efficient image codec based on set partitioning in hierarchical trees. *IEEE Transactions on circuits and systems for video technology* 6, 3 (1996), 243–250.
- [14] Chong Shao, Shahriar Nirjon, and Jan-Michael Frahm. 2016. Years-Long Binary Image Broadcast using Bluetooth Low Energy Beacons. In *Proceedings of the International Conference on Distributed Computing in Sensor Systems (DCOSS 2016)*.
- [15] Jerome M Shapiro. 1993. Embedded image coding using zerotrees of wavelet coefficients. *IEEE Transactions on signal processing* 41, 12 (1993), 3445–3462.
- [16] Zhou Wang, Alan Conrad Bovik, Hamid Rahim Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *Image Processing, IEEE Transactions on* 13, 4 (2004), 600–612.
- [17] Saif Zahir, Kal Dhou, and BC Prince George. 2007. A new chain coding based method for binary image compression and reconstruction. *PCS, Lisbon, Portugal* (2007), 1321–1324.
- [18] Saif Zahir and Mehmood Naqvi. 2005. A new rectangular partitioning based lossless binary image compression scheme. In *Electrical and Computer Engineering, 2005. Canadian Conference on*. IEEE, 281–285.
- [19] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. 2015. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*. 1529–1537.