

A New Deep Learning-based Food Recognition System for Dietary Assessment on An Edge Computing Service Infrastructure

Chang Liu, Yu Cao, *Senior Member, IEEE*, Yan Luo, *Member, IEEE*, Guanling Chen, *Member, IEEE*, Vinod Vokkarane, *Senior Member, IEEE*, Yunsheng Ma, Songqing Chen, *Member, IEEE*, Peng Hou

Abstract—Literature has indicated that accurate dietary assessment is very important for assessing the effectiveness of weight loss interventions. However, most of the existing dietary assessment methods rely on memory. With the help of pervasive mobile devices and rich cloud services, it is now possible to develop new computer-aided food recognition system for accurate dietary assessment. However, enabling this future Internet of Things-based dietary assessment imposes several fundamental challenges on algorithm development and system design. In this paper, we set to address these issues from the following two aspects: (1) to develop novel deep learning-based visual food recognition algorithms to achieve the best-in-class recognition accuracy; (2) to design a food recognition system employing edge computing-based service computing paradigm to overcome some inherent problems of traditional mobile cloud computing paradigm, such as unacceptable system latency and low battery life of mobile devices. We have conducted extensive experiments with real-world data. Our results have shown that the proposed system achieved three objectives: (1) outperforming existing work in terms of food recognition accuracy; (2) reducing response time that is equivalent to the minimum of the existing approaches; and (3) lowering energy consumption which is close to the minimum of the state-of-the-art.

Index Terms— Mobile Applications, Object Recognition, Deep Learning, Edge Computing, Food Recognition

1 INTRODUCTION

In the US, more than one-third (34.9% or 78.6 millions) of adults are obese and approximately 17% (or 12.7 millions) of children and adolescents aged 2 to 19 years are obese [1]. There were more than 1.9 billion adults, 18 years and older, were overweight on earth in 2014 [2]. Documenting dietary intake accurately is crucial to help fight obesity and weight management. Unfortunately, most of the current methods for dietary assessment (for example, 24 hour dietary recall [3] and food frequency questionnaires [4]) must rely on memory to recall foods eaten.

In the last few years, we have witnessed an explosive increase of mobile and wearable computing devices (e.g., the smart watch and smart phone) in the consuming electronics market. One common characteristic of these devices is that many of them have inexpensive, unobtrusive and multimodal sensors. These sensors enable us to collect

multimedia data (e.g., video and audio) in natural living environments. Due to the ubiquitous nature of mobile and wearable devices, it is now possible to use these devices to develop pervasive, automated solutions for dietary assessment [5-11]. One example of such solutions is to use mobile devices as a pervasive food journal collection tool and to employ cloud service as a data analysis platform. The combination of mobile device and cloud service could contribute to improving the accuracy of dietary assessment. As a result, in the last few years, we have seen several mobile cloud software solutions [12-14] to improve the accuracy of dietary intake estimation. One common issue among these solutions is that the users of the software must enter what they have eaten manually. To address this issue, visual-based food recognition algorithms and systems have been proposed [6-11]. A recent review by Martin et al. [15] also indicated that using digital imaging techniques for food recognition is superior to many other methods of dietary assessment techniques. Some advantages of visual-based food recognition systems include: reduced burden for users to recall the food, improved accuracy and efficiency of dietary recall.

While promising, one of the major barriers of adopting automatic dietary assessment system into practice is how to design and develop effective and efficient algorithms and system to derive the food information (e.g., food type) from food images. Considering the limited computation resources and low battery life on mobile device, it is more challenging to develop such a system within the mobile

- Chang Liu, Yu Cao, Guanling Chen, Peng Hou are with the University of Massachusetts, Lowell, MA, 01854. E-mail: ycao@cs.uml.edu.
- Yan Luo and Vinod Vokkarane are with the Department of Electrical and Computer Engineering, University of Massachusetts, Lowell, MA, 01854.
- Yunsheng Ma is with the Department of Medicine, University of Massachusetts Medical School, Worcester, MA, 01601.
- Songqing Chen is with the Department of Computer Science, George Mason University, Fairfax, VA, 22030.

***Please provide a complete mailing address for each author, as this is the address the 10 complimentary reprints of your paper will be sent

Please note that all acknowledgments should be placed at the end of the paper, before the bibliography (note that corresponding authorship is not noted in affiliation box, but in acknowledgment section).

Digital Object Identifier 10.1109/TSC.2017.2662008

cloud computing paradigm. We have carefully investigated this problem and have identified two major challenges. The first major challenge is how to design effective and efficient analytics algorithms to achieve optimal recognition accuracy. The second major challenge is how to develop a system that can minimize energy consumption and response time.

To address the first issue (recognition accuracy), we plan to develop new deep learning-based algorithms. Deep learning [16, 17] (also known as representation learning, feature learning, deep structured learning, or hierarchical learning) is a new area of machine learning research. It allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction [18]. In the last five years, these techniques have improved the state-of-the-art in speech recognition, computer vision, natural language processing, and many other domains. Our extensive experiments in this paper have shown that, compared with traditional hand engineered features (e.g., SIFT [19]) and shallow learning-based classification algorithms (e.g., Support Vector Machine (SVM)), our proposed deep learning-based classification algorithms could improve the recognition accuracy substantially. We also developed other image analysis algorithms to enhance the food image quality for data analysis. All these algorithms have been integrated into an edge computing-based real-time computing system, which is discussed in the next paragraph.

To address the second issue (energy consumption and response time), we aim to design and employ a real-time food recognition system employing edge computing service paradigm. The proposed system distributes the data analytics throughout the network by splitting the food recognition task between the edge devices (close to end users) and the servers (in the cloud). Edge computing refers to the enabling technologies that allow computation to be performed at the edge of the network in a stream fashion. Edge computing is a non-trivial extension of cloud computing from the core network to the edge network [20-26]. The proposed edge computing service infrastructure is particularly useful for our application because most of the mobile devices have limited computation capacity and battery life. Hence, it is difficult for them to support computational-intensive tasks. At the same time, our proposed food image analysis algorithms usually involve heavy computation and may require much more computation resources.

In this paper, we focus on two major research efforts. The first research effort aims to develop new food recognition algorithms, including new food image recognition algorithms based on deep learning and image pre-processing and segmentation algorithms to enhance the quality of food image. The second research effort aims to design a real-time food recognition system for dietary assessment. The proposed system employs edge computing service paradigm and distributes the data analytics throughout the network. Specifically, the proposed system will split the food recognition tasks between the edge devices (which is physically close to the user) and the server (which is usually located in the remote cloud). For example, in our system, the edge devices (e.g., user's smart

phone) can perform light-weight computation on food image for food recognition. Then, our system will transmit the food images (after the light-weight computation at edge device) to the server in the cloud to perform more accurate recognition tasks. By distributing the analytics throughout the network, our system can achieve significant improvement in the recognition accuracy, while minimizing the response time and energy consumption. In this project, we implemented a prototype system to verify our hypothesis and evaluate the proposed algorithms. Our prototype runs on both edge device (Xiaomi Note, running Android 6.0.1 "marshmallow") and server (an in-house GPU cluster). We also conducted extensive experiments with real-world data. The results show that our system achieves very impressive results on the following three aspects. First, to the best of our knowledge, the food recognition accuracy using our proposed approach outperformed all other reported results. Second, the response time of the proposed system is equivalent to the minimum of the existing approaches. Last but not the least, the energy consumption of the proposed system is close to the minimum of the state-of-the-art.

The rest of the paper is organized as follows. In Section 2, we introduce related work in computer-aided dietary assessment, visual-based food recognition, deep learning, and edge computing. In Section 3, we present the architecture, components, and algorithms for the proposed system based on deep learning and edge computing. In Section 4, we describe the implementation details of our system. Section 5 presents the evaluation results, which include recognition accuracy, power consumption, response time, etc. Section 6 discusses the system limitations. In Section 7, we make concluding remarks.

2 RELATED WORK

Estimating dietary intake accurately with a high-quality food journal is crucial for managing weight loss [27]. Unfortunately, due to many technical barriers, how to improve the accuracy of dietary intake estimation is still an open question. In this paper, we aim to develop a systematic approach as a first step to address this issue. We envision that there are four most relevant research areas, listed as below.

The first related research area is to enhance the accuracy of diet assessment with computer-aided solutions. Due to the recent advances in electronics, it is now possible to develop computer-aided solutions to transform healthcare from reactive and hospital-centered to preventive, proactive, evidence-based, person-centered. Dietary assessment is one such area that has gained a lot of attentions from both academia and industry. Among thousands of existing mobile cloud health software and hardware, we have seen many of them (e.g., MyFitnessPal [12], MyNetDiary [13], and FatSecret [14]) are dedicated for improving the accuracy of dietary estimates. However, all these applications require the user to enter everything they ate manually. To address this issue, several applications have been developed to improve the level of automation. For example, a

recent App entitled “Meal Snap” [28] aims to reduce human efforts by asking the user to take a picture, enter some quick information such as whether user is eating breakfast or lunch, and add a quick text annotation if the user wants to. Unfortunately, the accuracy of calorie estimation is heavily dependent on the accuracy of the manually entered text from user. Therefore, the accuracy is very unstable. Another example of such application is named “Eatly” [29]. This application requires the user to take the food image and then rates the food into one of the three categories (“very healthy”, “it’s O.K.”, and “unhealthy”). However, the actual rating is performed manually by the community, which consists of the users of this App. In this paper, we propose new algorithms and system that can recognize the food images (captured by the user with their mobile devices) automatically. This automation reduces the user’s burden substantially.

The second related research area is to perform dietary analysis using food images and/or videos. In one paper [6], researchers proposed an approach to combine a learning method (manifold ranking-based techniques) and a statistics method (co-occurrence statistics between food items) to recognize multiple food items. In another study [7], the authors proposed a method for fast food detection by researching the relative spatial relationships of local features of the ingredients and a feature fusion technique. NIH also funded a project named “Technology Assisted Dietary Assessment (TADA)” [11]. Researchers under this project have investigated different aspects of computer-aided dietary assessment, such as food item recognition, mobile interface design, and data development for food images. They have published several papers on food image recognition [8-10]. Most of the existing visual-based food recognition algorithms employed traditional signal processing with hand-engineered features (e.g., SIFT [19], HOG [30]) and shallow machine learning algorithms (e.g., SVM). Only very recently, with the striking success of deep learning, people started to research the application of deep learning for food image recognition [31]. Deep learning has the potential to address one main issue associated with existing techniques, which is that the hand engineered features may be useful for screening a few categories of food item but are unable to generalize to other food types. The proposed approach in this paper is also based on recent advances in deep learning. Related work in deep learning is introduced in the next paragraph.

The third related field is deep learning, which is a branch of machine learning. It allows the computers to learn from experience and understand the world in terms of a hierarchy of concepts using a deep graph with multiple processing layers. Each concept is defined in terms of its relation to simpler concepts [32]. Essentially, deep learning is trying to solve the central problem in representation learning by introducing representations that are expressed in terms of other, simpler representations [32]. It has already been proven useful in many disciplines, such as computer vision, speech recognition, natural language processing, bioinformatics, etc. There are two main classes of deep learning techniques. The first class is purely supervised learning algorithms, such as Deep Convolutional

Neural Network (CNN). The second class is unsupervised and semi-supervised learning algorithms, such as Denoising Auto-encoders and Deep Boltzmann Machines. In this paper, we focus on deep Convolutional Neural Network (CNN) [33]. Our proposed approach is rooted from CNN and it belongs to the category of supervised learning algorithms. CNNs are biologically-inspired [34] (animal visual cortex) variants of Multilayer Perceptrons (MLPs). It is consisted of neurons that have learnable weights and biases. Compared with MLPs, CNN has several distinct features. First, by enforcing a local connectivity pattern between neurons of adjacent layers, CNN could exploit spatially local correlation. Second, each filter in CNN is replicated across the entire visual field, which share the same parameters (e.g., weight vector and bias). Third, the neurons are arranged in three dimensions (width, height, and depth). Furthermore, a feature map can be generated by repeated application of a function across sub-regions of the whole image. Early implementation of CNNs, such as LeNet-5 [35], has been successfully applied to hand writing digital recognition. However, due to the lack of large scale labeled data and limited computation power, CNNs failed to address more complex problems. With the help of large-scale and well-annotated dataset like ImageNet [36], new computing hardware such as graphics processing unit (GPU), and several algorithms advancements such as Dropout [37], it is now possible to train large scale CNNs for complex problems. Recently, many research, such as VGGNet [38], ZFNet [39], GoogLeNet [40], Residual Network [41], has been proposed to address the issue of limited abilities of feature representation. One common strategy is to make the network deeper and avoid saturation issues. Our proposed approach was directly inspired by CNN work from LeNet-5 [42], AlexNet [33], and GoogLeNet [40]. The LetNet-5 [42] is a 7-layer network structure with 32x32 grey-scale image as input for hand written digital recognition. It includes three convolutional layers (C1, C3 and C5), two sub-sampling layers marked as (S2, S4), one fully connected layers (F6), and one output layer. LeNet-5 generates a feature map and feed the feature map into the two fully-connected layers. After that, a 10-class output is generated. A receptive field (a.k.a. “fixed-size patch” or “kernel”) is chosen during the convolutional layer to compute convolution with the same size patch in the input image. A stride is defined to make sure every pixel in the original image will be covered. The system will perform convolution operation first, followed with a sub-sampling with the feature map. The goal of sub-sampling is for dimension reduction. Then, we will move to the fully connected layers, which are used to join the multi-dimension feature maps. Finally, we will generate a ten-class output, each of which represents one digital (from zero to nine). Please note, at each layer, the parameters (e.g., weight vector and bias) are trainable. Recent progresses in CNN have focused on enhancing the object representation with more complex models. For example, AlexNet [33] is a seven-layer model which includes five convolution layers and two fully connected layers. It outperformed the state-of-the-art object recognition techniques in 2012 ImageNet [36] challenges with large margin (over 10%). Later on, we

witnessed many new models with increased layers, increased layer size, more complex neurons, as well as sophisticated computation units and layer structures. Drop-out and ReLU were proposed to address the issue of over-fitting and saturation, respectively. Some excellent examples include VGG net [38], ZFNet[39], GoogLeNet[40], Residual Network [41]).

The last (but not the least) related research area is edge computing service infrastructure [22-26]. Under this infrastructure, part of the data processing tasks may be pushed to the edge of the network. One of the core ideas is called “Collaborative Edge” [22], which refers to the architecture that connects the edges of multiple stakeholders. These stakeholders may be geographically distributed and they may have distinct physical location and network structure. Under this infrastructure, the cloud paradigm is extended to the edge of the network. Therefore, such an edge computing service infrastructure offers a unifying paradigm for cloud-based computing and Internet of Things (IoT)-based computing. It has the potential to address the issues of delayed response time, reduced battery life, limited bandwidth, and data security and privacy. However, most of the existing use cases of edge computing-based digital health applications [43-45] are relatively simple examples with small data sets. Novel user cases and intriguing applications with more challenging tasks, such as larger data sets and sophisticated computation, are needed for evaluating the efficacy and effectiveness of edge computing in digital health. Our proposed application, which focuses on food image recognition for dietary assessment, employ very complicated computation tasks (e.g., image pre-processing, image segmentation, and deep learning) with large image data sets (in the size of GB). This application scenario provides an excellent playground to evaluate the efficacy and effectiveness of edge computing in digital health.

3 SYSTEM DESIGN

3.1 Overview

Our food recognition system employs visual sensors to capture food images as the source data. Due to the recent advances of electronics, visual sensors are now available in many Internet-of-Things(IoT) devices, such as smart phones. To simplify the design, we utilized the camera on smartphones for visual sensing. Besides the smartphone for sensing and image capturing, the recognition is done in a collaborative manner between the edge device (e.g., smartphone) and servers (e.g., servers in the cloud). As shown in **Figure 1**, our system includes end user layer (left most of **Figure 1**), edge layer (middle of **Figure 1**), and cloud layer (right most of **Figure 1**), together form a three-layer service delivery model. In our proposed system, data and computation are kept close to end users at the edge of network. Also, the end user’s device can passively record the geological location. Hence, the system could provide low latency, reduced energy consumption, and location-awareness for end users. The computations are distributed throughout the network, including both the edge devices and servers in the cloud. Please note, in our system, the

recognition is done in a collaborative manner.

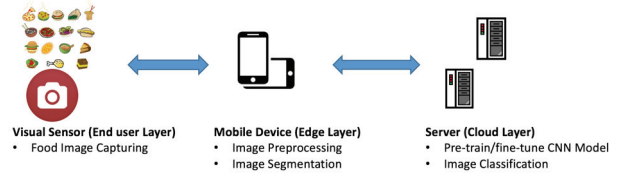


Figure 1: Overview architecture of our proposed “Deep Food on the Edge” system

The system design and related components are shown in **Figure 2**. As shown in this figure, our system consists of the following three major modules:

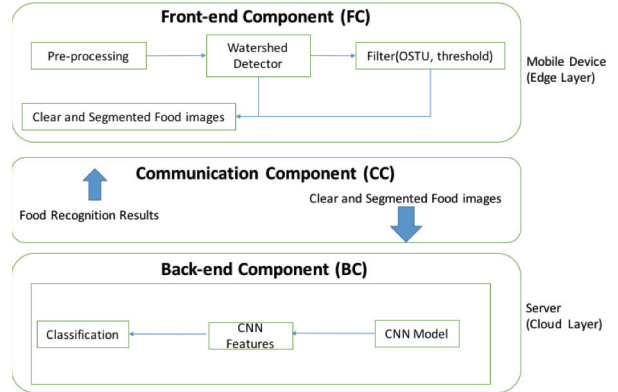


Figure 2. System design and components

Front-end Component (FC): we deploy the FC module on the edge device (smartphone). As shown in the top box in **Figure 2**, it’s consisted of three submodules, which are image pre-processing (e.g., blurry image detection), watershed detectors, and the filters (OTSU or threshold)-based segmentation. After the image pre-processing module, an original clear image is generated for segmentation. Next, the watershed detector, combined with different filters (e.g., OTSU-based threshold) to segment the original image. After segmentation, we can generate the clear and segmented image. These images will be transferred to the server via the Communication Module (introduced below) for further classification.

Communication Component (CC): CC provides two channels for communication between the FC and the Back-end Component (BC), which will be introduced in more detail in the next paragraph. It transfers the image data from the FC to the BC via Input Channel, and it also passes the detection results from the BC to the FC via Output Channel.

Back-end Component (BC): the BC module runs on the cloud server, which is configured to use Caffe [46] (an open source deep learning framework) for CNN model training and testing. We use pre-trained GoogLeNet by ImageNet and fine-tune it on our food dataset (UEC-256 and Food-101). Then the trained model is deployed on the server and used for classifying the image. More specifically, the segmented image is first passed through our proposed CNN model (which is rooted from GoogLeNet model [40]), then the features are generated from the model. furthermore, a

softmax classifier is used with these features to generate the probability of each category. Here we use the top-5 and top-1 probability as our prediction/classification of the food image. Our evaluation of accuracy is also based on these criteria.

3.2 Food Recognition Algorithms

In this section, we will introduce our proposed food recognition algorithms, which runs on the FC and BC. Essentially, our system is a multiple-stage food recognition system that distributes the analytics throughout the network.

3.2.1 Food Image Analysis Algorithms Running at FC

Once the food images are captured, we will conduct two types of computations at mobile device in the Front-end Component (FC) (a.k.a., Edge Layer): image pre-processing and image segmentation.

The main objective of the first computation (image pre-processing) is to identify if the image being captured is blurry or not. While many cameras on mobile devices have features such as optical zoom or auto focusing, in real-world practice, when users take the pictures of food, they may have very limited time to do so due to their busy schedule and their photo taking action may be interrupted by other matters. Hence, the chances of device shaking and other interruptions while taking pictures are high. An automatic image blurry detection algorithm running at the mobile device is needed to give a real-time alarm to reminder user to re-take the picture if the image is blurry. We define an out-of-focus image as blurry image. Our goal is to develop a light weight and effective blurry image detection algorithms running at the mobile device. In literature, image restoration has been proposed to handle blurry images. Unfortunately, these existing methods are not applicable to our case because these techniques need a reference image to compute the quality of the test image. In our applications, we may only have test images. Followed our previous research [47], we propose a simple-feature (such as “edginess” of the image) and threshold-based method to divide the images captured into two groups (i.e., the clear image group and the blurry image group). The “edginess” of the image is defined as the number of edge pixels (e.g., detected by Sobel operator) divided by the total number of image pixels. The rationale behind this method is that the percentage of edge pixels for clear image (with clear object of interests) is much higher than the percentage of edge pixels for blurry image. In our previous research [47], we also noticed that there are different patterns between the frequency spectrums of clear image and blurry image. The Fourier spectrum of a blurry image usually shows prominent components along the certain degree (e.g., 45 degree) directions that correspond to the corners of the image. This is because the blurry image usually does not contain clear object information except the four strong edges at the corners of the image running at certain degree relative to the sides. On the other hand, the clear image usually has a lot of clear edge information so that its spectrum does not show prominent components along certain

degree directions because it has a wider range of bandwidth from low to high frequencies. Based on the aforementioned observation, we first employ texture analysis algorithms on the frequency spectrum image. Then we extract different types of texture features (e.g., entropy, contrast, correlation, homogeneity) from each image. Once the features are extracted, we employ different types of classifiers to classify the images into two categories (blurry image or clear image). Similar to our previous work [47], we employ a two-step K-means clustering algorithms, the details is illustrated in the **Algorithm 1**.

```

Data: A set of Images  $Set : \{I_1, I_2, \dots, I_n\}$ 
Result: Two clusters  $Set_b : \{b_1, b_2, \dots, b_p\}$ ,  $Set_c : \{c_1, c_2, \dots, c_q\}$ 
initialization, set  $i$  to 1;
while  $i$  is no more than  $n$  do
    Read one image  $I_i$ ;
    Extract texture features  $T_i$  from frequency spectrum;
    Apply entropy feature extraction from  $T_i$  as  $S_1$ ;
    Apply contrast feature extraction from  $T_i$  as  $S_2$ ;
    Apply correlation feature extraction from  $T_i$  as  $S_3$ ;
    Apply homogeneity feature extraction from  $T_i$  as  $S_4$ ;
    Use binary classifier for  $S_1, S_2, S_3, S_4$  separately;
    Combine classification result using majority vote;
    if blurry then
        group  $I_i$  into the  $Set_b$ ;
    else
        group  $I_i$  into the  $Set_c$ ;
    end
    go to next iteration;
end

```

Algorithm 1. Image Pre-processing in the Front-end Component(FC)

The main objective of the second computation (image segmentation) is to segment the image into two parts: foreground (which contains the actual food) and background. Based on the size of foreground, we could crop the image by removing some portion of the background that does not overlap with foreground. According to our own experiments and other people research results, when using deep learning-based model (which is the main algorithms used in server) for image analysis and object detection, if we could reduce the background information, the object detection and recognition accuracy could be improved. Inspired by this observation, we employ watershed [48] segmentation algorithm to preprocessing the image at FC. In this process, we first pre-process the image by image segmentation. Then we generate a new cropped image and send the updated image to the server in the cloud for further processing. By doing so, we can achieve the following performance improvements: (1) the volume of data transfer over the network may be reduced substantially. It also reduces the power consumption caused by network transferring; (2) The response time may be reduced by shorter transmission time, which will improve the user experiences; (3) The system uses much less network flow consumption, which is very helpful when the network connection is unreliable, or when the user is connected to the server via cellular network and/or he or she has limited data plan with the mobile device; (4) More importantly, the cropped image will eliminate the abundant information

and further improve the accuracy for classification. In theory, the watershed algorithm is based on the following observations: any grayscale image can be viewed as a topographic surface, in which the high intensity indicates peaks and hills while low intensity represents valleys. The watershed algorithm starts filling every isolated valley with different colored water. When water rises, water from different valleys with different colors will start to merge. We could avoid this by building barriers in the locations where water merges. The algorithm continues to fill water and build barriers until all the peaks are under water. Finally, the barriers the system created are the segmentation result. **Algorithm 2** illustrates the details of the algorithm.

Algorithm 2 Watershed Algorithm using topographical distance

INPUT: The lower complete grey scale Image (V, E, im) , which is defined from original image with a lower boundary

OUTPUT: a sequence of labels on V , representing background or foreground

```

1: WATERSHED  $\leftarrow 0$   $\triangleright$  (The label of watershed for every pixel)
2: Init Label with a minima and MASK for other pixels
3:  $U \leftarrow \{p \in V | \exists q \in N_G(p): im[p] \neq im[q]\}$ 
4: while not empty( $U$ ) do
5:   select point  $p$  from  $U$  with minimal grey value
6:   remove  $p$  from  $U$ 
7:   for all  $q$  steeper than  $p$  do  $\triangleright$  (pixel value is greater in the neighbour)
8:     if  $label[q] == MASK$  then
9:        $label[q] \leftarrow label[p]$ 
10:    else
11:       $label[q] \leftarrow WATERSHED$ 
12:    end if
13:  end for
14: end while  $\triangleright$  (Label array represents the boundary)

```

Algorithm 2. Image Segmentation in the Front-end Component(FC)

3.2.2 CNN-based Food Image Analysis Algorithms Running at BC

After the image pre-processing and segmentation at FC, we will further analyze these images at BC. Our proposed approach running at BC is based on the recent advances on deep learning, which aims to learn multiple levels of representation and abstraction that help infer knowledge from data such as images, videos, audio, and text.

Our proposed approach was directly inspired by CNN work from LeNet-5 [42], AlexNet [33], and GoogLeNet [40], and it employs a new module called “Inception Module”, which is motivated by recent advances named “Network-in-Network” [49]. This is also similar to the one used in GoogLeNet [40]. In this “Inception Module”, an additional 1x1 convolutional layers are added to the original AlexNet [33] network architecture. This additional layer undoubtedly increases the depth of the network. However, this addition could also substantially reduce the feature map’s dimension. Therefore, this module could help to remove the computation bottlenecks. Specifically, we use feature map as the input for the “Inception Module”. After that, we apply multiple levels of convolutional layers and max-pooling layers. The kernel size of the convolutional layer varies from 1x1 to 3x3 and 5x5. At each layer, different outputs are generated and are concatenated to form the new feature map, which is used as input for the convolution and pooling operation for next layer. In order to perform dimension deduction, an optimized convolution is proposed based on the “Inception Module”. Please note, instead of feeding the input directly into the convolutional layer, an additional convolutional layer with size 1x1 is added to reduce the input dimension. In addition, the output from the 3x3 max-pooling layer is sent into an additional convolutional layer with size 1x1. These new designs enable the new architecture to reduced dimension even the depth of the network is increased. Not surprisingly, our experiments have demonstrated that, even under constrained computational complexity, this new network structure is able to enhance the ability to capture more visual information. **Figure 3** illustrates the improved inception module. The network structure in the left (Figure a) is the original structure in regular CNN, such as AlexNet [33]. The right figure (Figure b) is the snapshot of the new network architecture with “Inception Module”. As shown in Figure b, the three added 1x1 convolutional layers are annotated with dotted rectangle and green color. While the number of layers in Figure b is four (which is one layer more than the number of layers in Figure a), the total dimension of the output (at feature concatenate layer) in Figure b is still

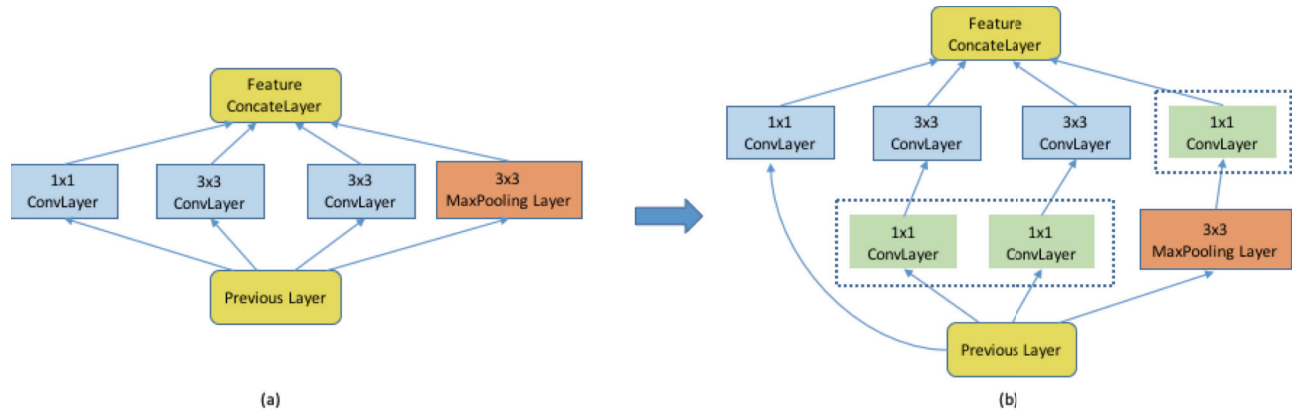


Figure 3: Illustration of the “Inception Module”. Figure (a) in the left is the snapshot of the original network architecture in regular CNN, such as AlexNet. Figure (b) in the right is the snapshot of the new network architecture with “Inception Module”. The figure is best viewed in color.

smaller than the output dimension of Figure a.

The next step after forming the “Inception Module” is to employ multiple modules to form the network (similar to GoogLeNet). In this step, we will connect the two modules using one additional max pooling layer. The output from the previous module will be used as the input for the next module. Specifically, the concatenated features (output) from the previous module are fed into the newly added max pooling layer. The output from the max pooling layer is used as input for next module. **Figure 4** illustrates this architecture. This figure includes two “Inception Module”, one (figure “a”) is located on the top of **Figure 4** and another (figure “b”) is located in the bottom of the **Figure 4**. These two components (figure “a” and figure “b”) are connected via a 3x3 max-pooling layer. Essentially, the new network architecture becomes a hierarchical level step by step. In order to address the issue of increased time complexity associated with the increased network layers, we resort to the lessons learned in recent paper [50], which offer some insights for designing the network architectures by balancing factors such as depth, numbers of filters, filter sizes, etc. In this study, we design a network structure with 22 layers (similar to the one used in GoogLeNet) with different kernel size and stride. We have found that, in our study, using an input size of 224x224 with three channels (RGB), combined with “1x1”, “3x3” and “5x5” convolutions, produces the best results. The 22 layers are layers with parameters. We design the pooling layer whose filter size is 5x5. The convolutional layer is 1x1 and includes 128 filters and ReLU (rectified linear activation). The dimension of the fully-connected layers is 1024. During pre-train-

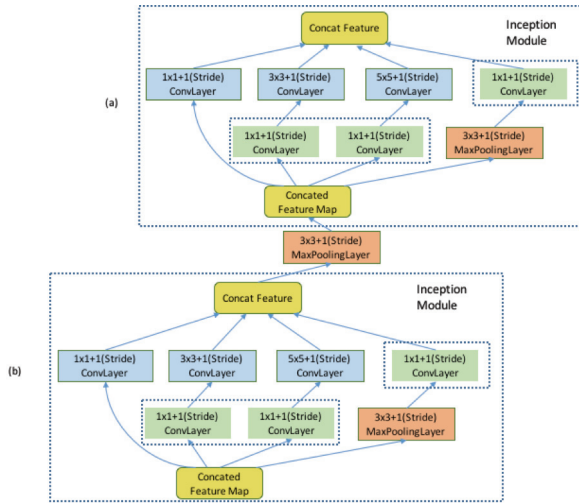


Figure 4. Illustration of module connection (best viewed in color)

ing stage, it is mapped into a 1,000-class output, similar to the ImageNet data set [36]. We use a 70% dropout rate to address the overfitting issue. Softmax is used for final classifier. Please note, based on the actual food categories, we will need to adjust the output class number during the fine-tuning stage. The proposed approach is implemented on top of open source deep learning framework Caffe [46]. CentOS 7.0 is chosen as our host system. We also use

NVidia Tesla K40 GPUs for model training. The model definition is adjusted in prototxt file in Caffe. We will introduce the implementation details in Section 4.

4 SYSTEM IMPLEMENTATION

In order to verify the efficacy and effectiveness of the proposed system, we implemented a prototype system for food recognition. Specifically, the front-end component (FC) is implemented on Android 6.0.1 (Marshmallow). The back-end component (BC) is implemented using server equipped with CentOS 7.0. The implementation of communication component (CC) includes two part. The first part is on the smartphone where we use Apache HttpClient to communicate with server. The second part is on the server we employed Django web development framework [51] and the associated RESTful web service. In this section, we present implementation details.

4.1 Implementation of Front-end Component (FC)

We develop an Android application for the front-end module. It runs on Xiaomi Note running Android 6.0.1 marshmallow. The image pre-processing algorithm, the watershed segmentation algorithm, and the threshold filter are also implemented in this application. The watershed algorithm runs on the local mobile devices and it is implemented using OpenCV [52] on Android devices. Several pre-defined markers are first constructed, the algorithm treats each pixel as a local topography, and then it fills the basins from the markers, until the basins meet on watershed lines. Here we set the markers as the local minimal of the food image, so that we can start from the bottom to fill the basins. We use OpenCV 3.10 and port the java SDK into the android studio project, which supports the OpenCV for Android SDK and also involves the image processing class.

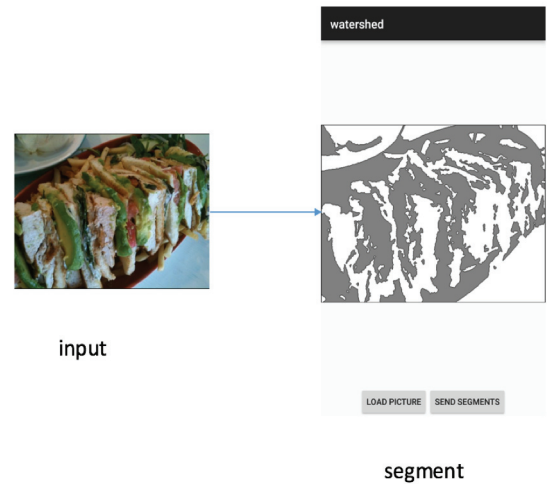


Figure 5: Screenshots showing image segmentation implementation in FC module

The App we implemented has an UI for processing and loading the image. A screenshot of the UI is shown in **Figure 5**. There is a background thread for preprocessing the

image. After it finishes, the App will display the segmented image in the application's mainframe. While in the background, the thread does several tasks when pre-processing the image, that includes: (1) rescaling the image if it's exceed 1024x786, since too large image will increase the computing time and energy consumption; (2) converting the RGB image to grey level image for further image processing, the grey image is more easily computed when there're many channels and pixels; (3) creating the watershed class and watershed threshold for dividing the image into segments and non-segments; (4) saving and generating a unified image segments for future transferring.

4.2 Implementation of Communication Component (CC)

There are two implementations for communication between the Android device and cloud server. For the Android application, we use Apache HTTP Client and construct the HTTP POST request to send the segmented image into the cloud server. First, a connection bound to the server is established, and then we construct the necessary HTTP header, and fill the content with image file. Then we send the POST request to the cloud server to finish the transmission. On the cloud server, we deploy a RESTful web server using Django [53], which supports the file transferring (image, audio, video) using HTTP requests. When the server is up and deployed, it will listen to the port and save the requested file into the pre-configured destination. Our server will store all the necessary segments for the classification task using trained-well CNN models.

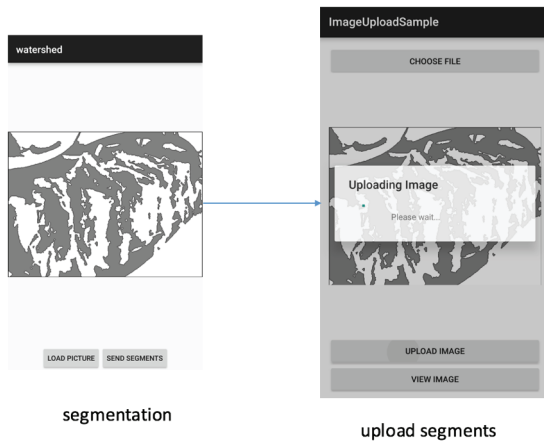


Figure 6: Screenshot showing segmented images being uploaded to the server in CM module

4.3 Implementation of Back-end Component (BC)

Our back-end system is mainly used for classification when we receive the images from the mobile device. Before testing, we used pre-trained GoogLeNet model from ImageNet, and then fine-tuned on public food data set like Food-101 and UEC-100/UEC-256. After these steps, a fine-grained model is generated which can be used for specifically food image classification. We use Caffe to train and tune the model. And our deployment of model is also based on Caffe's python interface. We first load the model

into memory, when the test food image is fed into the Convolutional neural network as the input, CNN features are extracted, with max-pooling and rectified linear-unit (ReLU) layers for dimension reduction and accelerating the convergence of computing.

5 PERFORMANCE EVALUATION

5.1. Experiment Setup and Evaluation Data Set

In all the following experiments, we use Xiaomi Note running Android 6.0.1 "Marshmallow" as the front-end to install the FC of our system. This smartphone uses Qualcomm MSM8974AC Snapdragon 801 featuring Quad-core 2.5 GHz Krait 400 and an Adreno 330 GPU. It also has a 64 GB of internal storage and 3 GB of RAM. In the back-end, we use an in-house GPU server. This server is a SuperServer 4027GR-TR from SuperMicron. It has two Intel Xeon processor E5-2600 with 512GB RAM. This server is also equipped with four NVIDIA Tesla K40 GPU.

In order to evaluate the effectiveness and efficiency of our system, we implemented two other systems running the state-of-the-art visual-based food recognition algorithms for comparisons. The first one, entitled as C-System, employs different types of computer vision algorithms using hand engineered features (e.g., SIFT [19], SURF [54], HOG [55], Cascade [56]) running at the mobile device for food image recognition, without relying on any algorithms running in the server. These algorithms (e.g., the Cascade algorithm) have been used in many embedded computer vision systems. We also implement the second system, called D-System, for comparisons. The D-system mainly relies on using the state-of-the-art deep learning algorithms running in the server, without using any image analysis and/or pre-processing computation at mobile device. Both systems are evaluated against our proposed system, and the performance metrics we use include response time, energy consumption, and detection accuracy.

In our experiment, we use two publicly available and challenging data sets, which are UEC-256/UEC-100 [57] and Food-101 [58]. As shown in the sub-sections below, the results of our proposed approach outperformed all the existing techniques in terms of accuracy. At the same time, the response time and energy consumption of our system are close to the minimum of the existing approaches.

5.2 Experimental Results on UEC-256/UEC-100 Dataset

As we have introduced before, we employ two data sets for our experiments. We will introduce our experimental results for the first category in this section, which is UEC dataset [57]. It was first developed by DeepFoodCam project [59] and the majority of the food items in UEC dataset is Asian food. This data set includes two sub-data sets: UEC-100 and UEC-256. The first sub-data set (UEC-100) includes 100 food categories with a total of 8643 images. There are around 90 images in each category. The second sub-data set (UEC-256) includes 256 categories with a total of 28375 images. There are around 110 images in each category. The researchers have added correct annotations for each image, including food category and bounding box (used to indicate the positions of the food). We

use UEC-256 as the baseline dataset since we prefer to have large scale training data. We divided the images into 5 groups (5 folds). 3 groups (out of 5 groups) were used for training and the rest of the images were used for testing.

In our experiments, the publicly available, 1000-class category from ImageNet dataset was used as the pre-trained model. This model (pre-trained model) was generated by training over 1.2 million images and testing over 100,000 images. Once we have the pre-trained model, we fine-tuned this model with the UEC-256 dataset. We fine-tuned the model with a base-learning rate of 0.01, a momentum of 0.9 and 100,000 iterations. The results are shown below in **Table 1**. If we compare the results in **Table 1** with the results in our previous publication [60], we could make two discoveries. First, our detection accuracy in this paper is slightly better. Second, the number of iterations when we reach the best performance is less than our previous paper. These two discoveries indicate that, due to the adaption of the proposed new system and algorithms, both the accuracy and the time complexity have been slightly reduced.

Table 1: Comparison of accuracy on UEC-256 at different iterations using UEC-256

# of Iterations	Top-1 accuracy	Top-5 accuracy
4,000	46.0%	77.5%
24,000	51.0%	78.8%
56,000	51.3%	79.6%
84,000	53.3%	80.7%
92,000	54.5%	81.8%

We also compared our results with both the C-System and the D-System. As we introduced before, the D-system is employing different sophisticated deep learning-based food image recognition algorithms, including the algorithms from the DeepFoodCam papers [57, 59]. To make a fair comparison, we used the same dataset as original papers, which is UEC-100, as well as the same strategy of dividing image dataset, the result is shown in the **Table 2**. Please note, there are five “C-system” in this table because we tried different types of computer vision algorithms using hand engineered features. Each sub-category of “C-system” (the first five rows in **Table 2**) represents one type of hand engineered feature. From this table, we can tell that our proposed method outperformed all existing methods using the same dataset:

Table 2: Comparison of accuracy between our proposed approach and existing approaches using the same data set (UEC-100)

Method	Top-1	Top-5
C-System (SURF-BoF+ColorHistogram)	42.0%	68.3%
C-System (HOG Patch-FV+Color Patch-FV)	49.7%	77.6%
C-System (HOG Patch-FV+Color Patch-FV(flip))	51.9%	79.2%
C-System (MKL)	51.6%	76.8%
C-System (Extended HOG Patch-FV+Color Patch-FV(flip))	59.6%	82.9%
D-System (DeepFoodCam(ft))	72.26%	92.0%
Proposed Approach in this paper	77.5%	95.2%

Table 3 shows the corresponding energy consumption of the three systems upon each food image. This table shows that the energy consumption of our system is very close to the energy consumption of the both C-system and D-system. Please note, in **Table 3**, we computed the energy consumption for both image analysis (on mobile device) and the image transferring (from the mobile device to the server). However, we did not compute the energy consumption if the computation is performed at the server in the cloud. Therefore, the D-system’s energy consumption for image analysis is zero because D-system does not include any computation on mobile device. On the other hand, the energy consumption for image transferring for C-system is zero. Because in C-system, there is no need for data uploading since all the recognition tasks have been done on the mobile device.

Table 3: Energy consumption (Joule) from different systems

Method	Energy Consumption (Joule) Per Image for Image Analysis	Energy Consumption (Joule) Per Image for Image Transferring
C-System	1.01	0
D-System	0	0.98
Proposed Approach in this paper	0.51	0.57

As of the computation and response time, let’s first discuss the computing time. Indeed, our algorithms is based on deep learning and training a large deep learning model requires a large amount of time. For example, on a NVidia Tesla K40 GPU, it takes 2 to 3 seconds per image for forward-backward pass using our proposed architecture. Since large dataset like ImageNet and Microsoft COCO [61] contains so many images, it may not be wise to train the model from scratch. One practical strategy is to use the pre-trained model in model zoo from existing implementation (e.g., Caffe [46]), which is public for all researchers. In our own experiment, the training time is largely impacted by the computation capacity of the server (e.g., the types of CPU and GPU), how large the image candidate is, how many iterations we choose, and what value we choose for learning rate, etc. According to the rough estimation, if we use the pre-trained GoogLeNet model, then fine-tune on the UEC-100, UEC-256, Food-101 dataset, it roughly takes 2 to 3 days nonstop for a server equipped with Nvidia K40 GPU to train the model. Once the model is trained, we can directly apply the model for classifying the image. On average, it takes 50 seconds for recognition for one image. Therefore, the average response time (the time duration between capturing the image and getting the food recognition results) is 1 minute per image for our proposed approach, which include time for image pre-processing on mobile device, the time to uploading the image to server, and the time for recognition in the server. As a comparison, the response time for C-system is usually around 35 to 55 seconds (depends on what hand engineered features we use). For example, the average computation time for a SIFT-like feature extraction and analysis algorithm on a

mobile device (Xiaomi Note) is 50 seconds. On the other hand, in the D-system, the response time (the time duration between capturing the image and getting the food recognition results) is 70 seconds per image in our experiments. This is mainly because in D-system, the image being processed is the raw image without pre-processing. Hence, we could conclude that the response time of our proposed approach is very close to the minimal response time of existing approach.

5.2 Experimental Results on Food-101

In addition to the first data set (UEC data set), we use the second data set, Food-101 data set [58], in our experiment. This dataset includes a total of 101 categories. For each food category, there are around 1000 images. We used around three-quarters (75%) of these images for training and the rest of the images are used for testing. Altogether, there are over 100,000 images in this data set. One thing about this data set is that this data set does not provide any bounding box information (which can be used to indicate the food location in the image). Instead, this data set offers food type information for each image. Different from the UEC data set, most of the images in this data set are popular western food images.

For this data set, we used a similar implementation as the one used in Section 4.1. The parameters were adjusted to fit for this new data set. We used a base learning rate of 0.01, a momentum of 0.9. Similar to the methods we used in Section 4.1, we fine-tuned the model on Food-101 dataset. **Table 4** below shows the accuracy (both top-1 accuracy and top-5 accuracy are listed as below). Again, if we compare the results in **Table 4** with the results in our previous publication [60], we can find that, due to the new system and algorithms in this paper, both the accuracy and the time complexity have been slightly reduced.

Table 4: Comparison of accuracy on Food-101 at different iteration

# of Iterations	Top-1 accuracy	Top-5 accuracy
5,000	65.6%	88.7%
10,000	70.7%	91.2%
20,000	73.4%	92.6%
60,000	77.0%	94.0%

We also compared our experimental results with the results of both the C-System and the D-System using the same data set (Food-101 datasets). As shown in **Table 5**, our proposed method is better than all existing work using the same dataset and division.

Table 5: Comparison of accuracy using different method on Food-101

Method	top-1	top-5
C-System (RFDC-based Approach from Lukas et.al[58])	50.76%	NA
D-System (CNN-based Approach from Lukas et.al[58])	56.40%	NA
Proposed Approach in this paper	77.0%	94%

From the above table, we can see that pre-trained model

with domain specific fine-tuning can boost the classification accuracy significantly. And fine-tuning strategy improves the accuracy comparing with non-fine-tuning method. The “NA” value in the “top-5” column means “not available”, as we used the original experiment data from their paper[58], and they don’t provide the top-5 result in it.

As of the energy consumption and response time, we have similar results reported as our previous data set (UEC-256/UEC-100), as introduced in the last paragraph of Section 5.1. Due to the space limitation, we did not report the exact numbers here.

5.3 The Employment of Bounding Box

As shown in both Section 5.1 and Section 5.2, the detection accuracy of our proposed approach is better than all existing approaches. We believe that one of the reasons we could achieve such performance boost is because in our proposed approach, image pre-processing and image segmentation are performed at the mobile device before analyzing these images in the server. To verify this hypothesis, we conducted a simple experiment. Our goal is to demonstrate that even very simple pre-processing can help improve the recognition performance. For example, we can use a simple bounding-box strategy to reduce the image size without analyzing the image content fully.

Specifically, we first employed the bounding box to crop the raw image. After this processing, only the food image part is remained for training and testing. Then, we performed similar experiment on UEC-256 dataset.

Table 6: Comparison of accuracy of proposed approach using bounding box on UEC-256

Method	top-1	top-5
Proposed approach (no bounding box)	53.7%	80.7%
Proposed approach (with bounding box)	63.6%	87.0%

We also conduct the experiment on UEC-100, as follows:

Table 7: Comparison of accuracy of proposed approach using bounding box on UEC-100

Method	top-1	top-5
Proposed approach (no bounding box)	54.8%	81.4%
Proposed approach (with bounding box)	76.3%	94.6%

As we can see from the two tables (**Table 6** and **Table 7**), the employment of bounding box could boost the classification accuracy substantially. A simple explanation for this is that the abundant information in the raw image is removed after the images were cropped using bounding-box. Therefore, a more accurate and clear image candidate for training can be generated. Please note, these results are valid only if we assume the majority of food image have the foreground centered on the image. Using this simple cropping-based approach will not work well if the food is scattered on different parts of the image. In this case, our

proposed approach, which conducts image pre-processing and image segmentation based on the image content, is certainly necessary to improve the recognition accuracy.

6 DISCUSSIONS

Our findings indicated that our system achieves very high detection accuracy, as shown in previous sections. However, the response time, while very close the minimal of existing systems, is still around 5% slower than the best performer. While this is not surprising since deep learning-based algorithms are usually very time-consuming, we believe that more research should be devoted to further improving the speed. In particularly, we plan to investigate new deep learning algorithms that can be executed in mobile devices. There are some recent papers that have started to explore this area with some preliminary results [62, 63], which further motivate us to pursue this route in the future. While pushing the deep learning-based computation further to the edge device sounds like a good idea in the initial look, we will have to consider the energy consumption if we execute the deep learning algorithms at the edge device. We believe much more research is needed in the area of distributed deep learning-based analytics in the era of edge computing.

7 CONCLUSION

In this paper, we aimed to develop a practical deep learning based food recognition system for dietary assessment within the edge computing service infrastructure. The key technique innovation in this paper includes: the new deep learning-based food image recognition algorithms and the proposed real-time food recognition system employing edge computing service paradigm. Our experimental results on two challenging data sets using our proposed approach have demonstrated that our system has achieved the three major objectives: (1) it outperforms the results from all existing approaches in terms of recognition accuracy; (2) it develops a real-time system whose response time is close to the minimal of existing techniques; and (3) it saves the energy by keep the energy consumption equivalent to the minimum of the existing approaches. In the future, we plan to continue improving performance of the algorithms (in terms of detection accuracy) and system (in terms of response time and energy consumption). We also plan to integrate our system into a real-world mobile devices and edge/cloud computing-based system to enhance the accuracy of current measurements of dietary caloric intake estimate. As our research is related to the biomedical field, much larger data sets are needed to provide convincing evidence to verify the efficacy and effectiveness of our proposed system. Backed by several major federal grants from NSF and NIH, we are in the process of collaborating with UMass Medical School and the University of Tennessee, College of Medicine to deploy our system in the real-world clinical practice.

ACKNOWLEDGMENT

The paper described is supported in partial by National

Science Foundation (NSF) of the United States (Award No. 1547428, 1541434, 1440737, and 1229213). Points of view or opinions in this document are those of the authors and do not represent the official position or policies of the U.S. NSF.

REFERENCES

- [1] C. L. Ogden, M. D. Carroll, B. K. Kit, and K. M. Flegal, "Prevalence of childhood and adult obesity in the United States, 2011-2012," *Jama*, vol. 311, pp. 806-814, 2014.
- [2] (2016). *World Health Organization: fact sheets on obesity and overweight*, Available at <http://www.who.int/mediacentre/factsheets/fs311/en/>.
- [3] G. H. Beaton, J. Milner, P. Corey, V. McGuire, M. Cousins, E. Stewart, et al., "Sources of variance in 24-hour dietary recall data: implications for nutrition study design and interpretation," *American Journal of Clinical Nutrition*, vol. 32, 1979.
- [4] J. Cade, R. Thompson, V. Burley, and D. Warm, "Development, validation and utilisation of food-frequency questionnaires—a review," *Public health nutrition*, vol. 5, pp. 567-587, 2002.
- [5] R. Steele, "An overview of the state of the art of automated capture of dietary intake information," *Critical Reviews in Food Science and Nutrition*, 2013.
- [6] Y. Matsuda and K. Yanai, "Multiple-food recognition considering co-occurrence employing manifold ranking," in *Pattern Recognition (ICPR), 2012 21st International Conference on*, 2012, pp. 2017-2020.
- [7] S. Yang, M. Chen, D. Pomerleau, and R. Sukthankar, "Food recognition using statistics of pairwise local features," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, 2010, pp. 2249-2256.
- [8] F. Zhu, M. Bosch, I. Woo, S. Kim, C. J. Boushey, D. S. Ebert, et al., "The use of mobile devices in aiding dietary assessment and evaluation," *Selected Topics in Signal Processing, IEEE Journal of*, vol. 4, pp. 756-766, 2010.
- [9] B. L. Daugherty, T. E. Schap, R. Ettienne-Gittens, F. M. Zhu, M. Bosch, E. J. Delp, et al., "Novel technologies for assessing dietary intake: evaluating the usability of a mobile telephone food record among adults and adolescents," *Journal of medical Internet research*, vol. 14, 2012.
- [10] C. Xu, Y. He, N. Khannan, A. Parra, C. Boushey, and E. Delp, "Image-based food volume estimation," in *Proceedings of the 5th international workshop on Multimedia for cooking & eating activities*, 2013, pp. 75-80.
- [11] TADA: Technology Assisted Dietary Assessment at Purdue University, West Lafayette, Indiana, USA, available at <http://www.tadaproject.org/>.
- [12] MyFitnessPal.com: Free Calorie Counter, Diet & Exercise Tracker, available at <http://www.myfitnesspal.com/>.
- [13] MyNetDiary: the easiest and smartest free calorie counter and free food diary for iPhone, iPad, Android, and BlackBerry applications, available at <http://www.mynetdiary.com/>.
- [14] FatSecret: All Things Food and Diet, available at <http://www.fatsecret.com/>.
- [15] C. K. Martin, T. Nicklas, B. Gunturk, J. B. Correa, H. R. Allen, and C. Champagne, "Measuring food intake with digital

- photography," *J Hum Nutr Diet*, vol. 27 Suppl 1, pp. 72-81, Jan 2014.
- [16] Y. Bengio, "Learning deep architectures for AI," *Foundations and trends® in Machine Learning*, vol. 2, pp. 1-127, 2009.
 - [17] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, pp. 1527-1554, 2006.
 - [18] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436-444, 2015.
 - [19] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision (IJCV)*, vol. 60, pp. 91-110, November 2004.
 - [20] M. Satyanarayanan, P. Bahl, R. Caceres, and N. Davies, "The case for vm-based cloudlets in mobile computing," *Pervasive Computing, IEEE*, vol. 8, pp. 14-23, 2009.
 - [21] S. Yi, Z. Hao, Z. Qin, and Q. Li, "Fog Computing: Platform and Applications," in *Hot Topics in Web Systems and Technologies (HotWeb), 2015 Third IEEE Workshop on*, 2015, pp. 73-78.
 - [22] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge Computing: Vision and Challenges," *IEEE Internet of Things Journal*, 2016.
 - [23] M. Satyanarayanan, Z. Chen, K. Ha, W. Hu, W. Richter, and P. Pillai, "Cloudlets: at the leading edge of mobile-cloud convergence," 2014, pp. 1-9.
 - [24] W. Shi and S. Dustdar, "The Promise of Edge Computing," *IEEE Computer Magazine*, vol. 29, pp. 78-81, 2016.
 - [25] S. Yi, C. Li, and Q. Li, "A survey of fog computing: concepts, applications and issues," in *Proceedings of the 2015 Workshop on Mobile Big Data*, 2015, pp. 37-42.
 - [26] P. Garcia Lopez, A. Montresor, D. Epema, A. Datta, T. Higashino, A. Iamnitchi, *et al.*, "Edge-centric computing: Vision and challenges," *ACM SIGCOMM Computer Communication Review*, vol. 45, pp. 37-42, 2015.
 - [27] R. R. Wing and S. Phelan, "Long-term weight loss maintenance," *The American journal of clinical nutrition*, vol. 82, pp. 222S-225S, 2005.
 - [28] *Meal Snap: Magical Meal Logging for iPhone*, available at <http://mealsnap.com/>.
 - [29] *Eatly: Eat Smart (Snap a photo of your meal and get health ratings)*, available at <https://itunes.apple.com/us/app/eatly-eat-smart-snap-photo/id661113749>.
 - [30] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," 2005, pp. 886-893.
 - [31] A. Meyers, N. Johnston, V. Rathod, A. Korattikara, A. Gorban, N. Silberman, *et al.*, "Im2Calories: towards an automated mobile vision food diary," 2015, pp. 1233-1241.
 - [32] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*: MIT Press, 2016.
 - [33] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *NIPS*, 2012, p. 4.
 - [34] D. H. Hubel and T. N. Wiesel, "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex," *The Journal of physiology*, vol. 160, pp. 106-154, 1962.
 - [35] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, pp. 2278-2324, 1998.
 - [36] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database" in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, 2009, pp. 248-255.
 - [37] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, pp. 1929-1958, 2014.
 - [38] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
 - [39] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Computer vision-ECCV 2014*, ed: Springer, 2014, pp. 818-833.
 - [40] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, *et al.*, "Going Deeper With Convolutions," in *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on*, 2015.
 - [41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *arXiv preprint arXiv:1512.03385*, 2015.
 - [42] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE* vol. 86, pp. 2278-2324, 1998.
 - [43] V. Stantchev, A. Barnawi, S. Ghulam, J. Schubert, and G. Tamm, "Smart Items, Fog and Cloud Computing as Enablers of Servitization in Healthcare," *Sensors & Transducers*, vol. 185, p. 121, 2015.
 - [44] Y. Cao, S. Chen, P. Hou, and D. Buhl-Brown, "FAST: A Fog Computing Assisted Distributed Analytics System to Monitor Fall for Stroke Mitigation," in *Prof. of 10th IEEE International Conference on Networking, Architecture, and Storage (NAS 2015) (Best Paper Award)*, Boston, MA, U.S.A, 2015.
 - [45] Y. Cao, C. Liu, B. Liu, M. J. Brunette, N. Zhang, T. Sun, *et al.*, "Improving Tuberculosis Diagnostics using Deep Learning and Mobile Health Technologies among Resource-poor and Marginalized Communities," in *IEEE Conference on Connected Health: Applications, Systems and Engineering Technologies*, 2016.
 - [46] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, *et al.*, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the ACM International Conference on Multimedia*, 2014, pp. 675-678.
 - [47] J. Oh, S. Hwang, Y. Cao, W. Tavanapong, D. Liu, J. Wong, *et al.*, "Measuring objective quality of colonoscopy," *IEEE Transactions on Biomedical Engineering*, vol. 56, pp. 2190-2196, September 2009.
 - [48] R. C. Gonzalez and R. E. Woods, *Digital Image Processing*: Prentice Hall, 2008.
 - [49] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.
 - [50] K. He and J. Sun, "Convolutional neural networks at constrained time cost," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5353-5360.
 - [51] A. Holovaty and J. Kaplan-Moss, *The definitive guide to Django: Web development done right*: Apress, 2009.
 - [52] G. Bradski and A. Kaehler, *Learning OpenCV: Computer vision with the OpenCV library*: "O'Reilly Media, Inc.", 2008.
 - [53] *Django: The Web framework for perfectionists with deadlines*, available at <https://www.djangoproject.com/>.
 - [54] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, "Speeded-up robust features (SURF)," *Computer Vision and Image Understanding*, vol. 110, pp. 346-359, 2008.
 - [55] N. Dalal and B. Triggs, "Histograms of oriented gradients for

- human detection," in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, 2005, pp. 886-893.
- [56] P. Viola and M. Jones, "Rapid Object Detection using a Boosted Cascade of Simple Features," in *Proc. of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, Kauai Marriott, Hawaii, USA, 2001, pp. 511-518.
 - [57] Y. Kawano and K. Yanai, "Foodcam: A real-time food recognition system on a smartphone," *Multimedia Tools and Applications*, pp. 1-25, 2015.
 - [58] L. Bossard, M. Guillaumin, and L. Van Gool, "Food-101—mining discriminative components with random forests," in *Computer Vision—ECCV 2014*, ed: Springer, 2014, pp. 446-461.
 - [59] Y. Kawano and K. Yanai, "Food image recognition with deep convolutional features," in *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing: Adjunct Publication*, 2014, pp. 589-593.
 - [60] C. Liu, Y. Cao, Y. Luo, G. Chen, V. Vokkarane, and Y. Ma, "DeepFood: Deep Learning-based Food Image Recognition for Computer-aided Dietary Assessment," in *Prof. of The 14th International Conference on Smart homes and Health telematics (ICOST 2016)*, Wuhan, China, 2016.
 - [61] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, *et al.*, "Microsoft coco: Common objects in context," in *Computer Vision—ECCV 2014*, ed: Springer, 2014, pp. 740-755.
 - [62] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized Convolutional Neural Networks for Mobile Devices," *arXiv preprint arXiv:1512.06473*, 2015.
 - [63] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, *et al.*, "MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems," *arXiv preprint arXiv:1512.01274*, 2015.