# An Ensemble Mobile-Cloud Computing Method for Affordable and Accurate Glucometer Readout

Navidreza Asadi, and Maziar Goudarzi, *Senior Member, IEEE*

**Abstract**—Despite essential efforts towards advanced wireless medical devices for regular monitoring of blood properties, many such devices are not available or not affordable for everyone in many countries. Alternatively using ordinary devices, patients ought to log data into a mobile health-monitoring manually. According to medical specialists, it causes several issues: (1) due to the direct human intervention, it is prone to errors, and clients reportedly tend to enter unrealistic data; (2) typing values several times a day is bothersome and causes clients to leave the mobile app. Thus, there is a strong need to use now-ubiquitous smartphones, reducing error by capturing images from the screen of medical devices and extracting useful information automatically. Nevertheless, there are a few challenges in its development: (1) data scarcity has led to impractical methods with very low accuracy: to our knowledge, only small datasets are available in this case; (2) accuracy-availability tradeoff: one can execute a less accurate algorithm on a mobile phone to maintain higher availability, or alternatively deploy a more accurate and more compute-intensive algorithm on the cloud, however, at the cost of lower availability in poor/no connectivity situations. We present an ensemble learning algorithm, a mobile-cloud computing service architecture, and a simple compression technique to achieve higher availability and faster response time while providing higher accuracy by integrating cloud- and mobile-side predictions. Additionally, we propose an algorithm to generate synthetic training data which facilitates utilizing deep learning models to improve accuracy. Our proposed method achieves three main objectives: (1) $92.1\%$ and $97.7\%$ accuracy on two different datasets, improving previous methods by $\sim40\%$, (2) reducing required bandwidth by $45\times$ with $\sim1\%$ drop in accuracy, (3) and providing better availability compared to mobile-only, cloud-only, split computing, and early exit service models.

**Index Terms**—Mobile Computing, Ensemble Learning, Data Generation, Deep Learning, Smart Health

✦

## 1 INTRODUCTION

MANY mHealth/uHealth medical devices, especially those affordable in middle/low-income countries, show the measured quantity on a digital or seven-segment screen alongside other additional information such as date, time, diagrams and measurement units. In most commonly used mHealth services, particularly for diabetics, patients are required to manually type sensed values into their mobile app. As illustrated in Fig. 1(a), a client first reads a value from the medical device, opens the app, navigates to logging interface, and eventually logs the information through typing. These steps should be repeated every time they log a measurement.

### 1.1 Motivation

According to the reports [1], [2] as well as our own experience from our diabetes management application iDia [3], this procedure has a few drawbacks: (1) Manual logging multiple times a day, is deterring; it is bothersome and time-consuming, and based on the feedbacks we have received, leads users to eventually lose their interest in using the app. (2) It is prone to human errors. More importantly, the observations show that patients are tempted to enter fake information that are more acceptable and closer to the

normal values. This can have considerable negative effects on the whole process of prevention, control, and treatment.

Currently, there are two better approaches: (1) Some medical devices are able to transmit their data to mobile devices via distant communication technologies (e.g., Bluetooth), facilitating the logging procedure (Fig. 1(b)). Nevertheless, they are far more expensive and in some cases less accurate [4]. More importantly, most of them can only interact with their own software applications and do not permit third-party apps to receive data. In addition, different versions of transmission technologies cause incompatibility between different mobile and medical devices. Despite a potentially bright future, these devices currently hold less than one percent of the market [5]. This number is even lower in developing countries. (2) An alternative which has grown interest in academia, uses image capturing and computing capabilities of mobile devices (i.e., smartphones); they are available to almost everybody and are able to perform light-weight computing tasks. Fig. 1(c) illustrates this alternative: the mobile phone is used to capture an image of the medical device, and then image processing is applied to recognize the sensed values automatically. In this approach (Fig. 2), each digit is considered as an object having a region of interest (RoI) and an ordered sequence of digits (i.e., RoIs) forms a *read* string. A correct read means not only all digits are predicted correctly, but also in the same order as displayed on the device. In this paper we follow this methodology for its broader applicability especially in middle/low-income countries.

---

*N. Asadi is now with Computer Engineering Department, Technical University of Munich, Germany (navidreza.asadi@tum.de). He was with Computer Engineering Department, Sharif University of Technology, Iran while working on this project.*
*M. Goudarzi is with the Computer Engineering Department, Sharif University of Technology, Iran (goudarzi@sharif.edu).*
*Manuscript submitted to IEEE Transactions on Mobile Computing.*

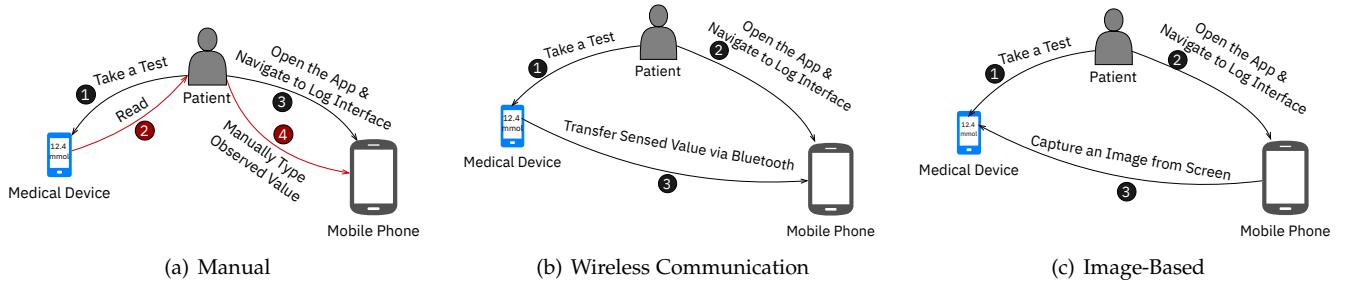(a) Manual         (b) Wireless Communication        (c) Image-Based

Fig. 1. Logging methods. Removing human interventions (red arrows) is preferred to avoid false logs.
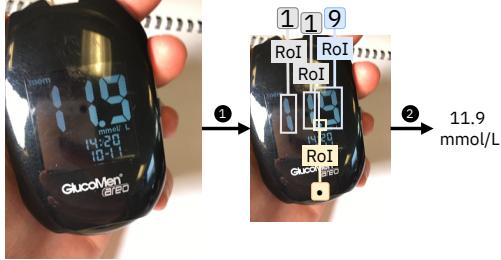


Fig. 2. Image-based method, showing Regions-of-Interest (RoIs) and the correct labels. The expected correct readout is "11.9 mmol/L".

## 1.2 Challenges

**Accuracy.** Although automated image-based reading may seem basic and easy, there are a few points that prove it a challenging problem in this specific case: (1) the current state-of-the-art [6], which has largely improved the previous algorithms, uses only two medical devices, yet achieves just **51.5%** accuracy, meaning it misreads almost half of the captured images. (2) as we discuss in Sections 2 and 3, the diversity of medical devices and a variety of structural and visual differences as well as noisy information on their screens (e.g., date and time) make it extremely difficult to read sensed value with business-as-usual image processing techniques.

**Data Scarcity.** The accuracy of deep learning algorithms relies on either big volume of annotated data to be trained in a supervised manner or a model pre-trained on a related task. To the best of our knowledge, there exists neither a big dataset nor a related pre-trained model in our particular task of imaged-based reading. Generative models (e.g., GANs) also would not help there because they require similar data to train, which is not applicable in our case. Additionally, it might be really difficult (if not impossible) to generate automatic annotation for our usecase.

**Resource Constraints.** Current state-of-the-art deep learning models are compute-intensive and need to be deployed on specialized cloud infrastructures. Thus, they introduce new challenges, including availability during poor network conditions, that is genuinely an expected issue in our target under-developed countries. Edge Computing [7], referring to every device near data source with some compute capacity (e.g., smartphones), is considered as a promising approach to improve availability and quality of service (QoS) by reducing delay. Mobile devices are usually resource-constrained, and therefore, can only run simpler deep neural network (DNN) models with much lower accuracy. Cloud Computing, however, is at the opposite side.

## 1.3 Main Contributions

We propose practical solutions to address the challenges above. We present an ensemble mobile-cloud computing architecture to get the best of both worlds: higher availability on the mobile, as well as higher accuracy and enhanced performability (i.e., a measure of level of performance/service-quality of the system) by integrating cloud and mobile modules (Fig. 3, described in §3). The followings items are our main contributions:

(1) A hybrid mobile-cloud service architecture, together with a compatible ensemble deep learning algorithm. This enables an accuracy-availability tradeoff based on network connectivity. We addressed the challenges of combining predictions of two separate models; e.g., differences and overlaps in the identified bounding boxes for each data element.

(2) A simple yet effective compression method. Combined with our ensemble model, this provides higher accuracy despite little communicated data.

(3) A high-fidelity data synthesizer algorithm, making utilizing deep learning models possible. This has basically turned the challenge into an opportunity; the variety of glucometer models, data formats, units, fonts, etc. is a challenge to conventional methods, but we used it in our data synthesizer mechanism to produce enough reasonable data to train high accuracy models that cover many varieties including those not seen before.

Our proposed method achieves $92.1\%$ and $97.7\%$ accuracy on two real-world datasets, and improves previously published results by more than $40\%$. It reduces the required bandwidth by $45\times$, and maintains higher availability compared to mobile-only, cloud-only, split computing, and early exit service models. Our proposal can be easily extended to other usecases with minor modifications.

The rest of this paper is organized as follows: In Section 2, we review related work. In Section 3, we present our proposed method. Section 4 explains our dataset generation algorithm. We evaluate our methods and algorithms in Section 5 and conclude in Section 6.

## 2 RELATED WORK

We separate related work into two different parts. The first one presents algorithms for reading sensed values or detection and recognition of digits on digital or seven-segment screens. The second one summarizes the studies that attempt to deploy part or whole of a deep learning model on the mobile.
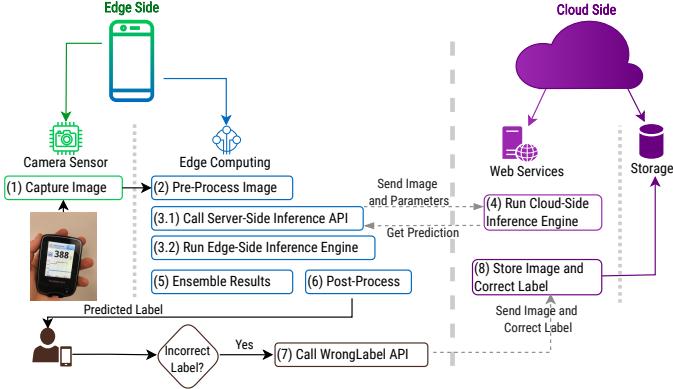
Fig. 3. Our mobile-cloud architecture.

## 2.1 Image-Based Automated Reading

Most of the proposed methods break the problem into multiple steps including image enhancement, localization of RoIs, detection and classification, and eventually ordering. We review related work within five criteria, as summarized in Table 1.

**Automated Localization.** Locating the RoIs is a crucial step and impacts the final accuracy. Many works try to simplify the problem while assuming the localization step is somehow already done, either manually by a client, or by fixing the device or using special markers [8], [9], [10], [11], [13], [15]. A few others [6], [12], [16], and our work take a more holistic approach, applying automated localization as well.

**Accuracy.** To our knowledge, no previous work achieves a reasonable readout accuracy. [6] with $51.5\%$ has by far the best performance. General OCR engines also are not helpful; our previous experiments along with the reports from [15], [16], and [6] confirm the state-of-the-art OCR engine, Tesseract's [14] poor performance ($<10\%$) in this specific task. On the other hand, our method can reach over $90\%$ accuracy on both datasets.

**Robustness.** Different methods, especially those using conventional algorithms are usually sensitive to different conditions such as skewness, various noises, camera perspective and angle, exposure, and illumination. Here, we define robustness as performing consistently in different conditions. That said, only few works [6], [16] try to adress it. We simulate all these variations in our data synthesizer.

Addtionally, the ensemble of two different models, mitigates such errors.

**Generalization.** There are a variety of medical devices each having unique characteristics, such as different font styles (including various types of seven-segments and/or digital styles), background and foreground colors, screen size and shape, units, backlit, etc. (Fig. 4). Nevertheless, the previous studies, except [16], use one or very few specific devices so that their proposed algorithms highly depend on the properties of those selected devices; hence, may not be considered as a general solution. For example, [13] is designed for a particular medical device with a blue backlit screen, [6] considers only seven-segment screens, and [15] assumes the largest contour as the screen, and hardcodes exact location of RoIs. In contrast, we cover a broad market, and illustrate this using an additional public dataset from Oxford [6].

**Respone Time.** Since larger portion of the previous studies use conventional and light-weight image processing methods, they can achieve reasonable response time. For instance, [13] is implemented on a Samsung Galaxy S i9000 mobile device, and can process 20 frames per second, or [12] is deployed on an N95 mobile device, and achieves five frames per second. The only exception is [15] which uses a deep convolutional neural network (CNN) for the digits classification step. It takes 10 seconds, that is unsatisfactory. Although we leverage CNN-based object detection models, we choose parameters so that the response is prepared in less than half a second. Besides, we design a simple compression technique when using the cloud-side engine in poor network conditions, and, therefore, reduce the end-to-end response time.

## 2.2 Deep Learning on Edge

In practice, an edge device can be any computing machine (indluding smartphones) that generates data or is near a data generation source. Edge devices are usually resource-constrained. So it is challenging to deploy big deep learning models on edge. A concise comparison of different methods is provided in Table 2. The efforts to overcome the limitations can be divided into three major directions.

The first direction is designing light-weight DNN models or optimizing existing ones. Several successful works have studied light-weight models including [17], [18], [19], [20], [21], and [30]. In general, related work in this category leverage a combination of using convolution blocks with lower parameters (e.g., separable convolutions), quantization, pruning, and model distillation. These techniques

TABLE 1
Summary of Literature Review on Image-Based Automated Reading.

| Work | Localization | Accuracy | Robustness* | Generalization | Response Time |
|---|---|---|---|---|---|
| [8] | ✗ | ✗ | ✗ | ✗ | ✓ |
| [9] | ✗ | ✗ | ✗ | ✗ | ? |
| [10] | ✗ | ✗ | ✗ | ✗ | ? |
| [11] | ✗ | ✗ | ✗ | ✗ | ? |
| [12] | ✓ | ✗ | ✗ | ✗ | ✓ |
| [13] | ✓ | ✗ | ✗ | ✗ | ✓ |
| [14] | ✓ | ✗ | ✗ | ✗ | ✓ |
| [15] | ✗ | ✗ | ✗ | ✗ | ✗ |
| [16] | ✓ | ✗ | ✓ | ✓ | ✓ |
| [6] | ✓ | ✗ | ✓ | ✗ | ✓ |
| Ours | ✓ | ✓ | ✓ | ✓ | ✓ |

\* Robustness to different environmental conditions.
? Information not available.

TABLE 2
Related Work on Deep Learning at Edge: A summary.

| Method | Works | Accuracy | Availability | No-Cloud* | Bandwidth** | Response Time | Deployment |
|---|---|---|---|---|---|---|---|
| D1 | [17], [18], [19], [20], [21] | ✗ | ✓ | ✓ | – | ✓ | E |
| D2 A1 | [22], [23], [24] | ✗ | ✓ | ✗ | ✓ | ✓ | E,C |
| D2 A2 | [25], [26], [27] | ✓ | ✗ | ✗ | ✓ | ? | E,C |
| D3 | [27], [28], [29] | ✗ | ✓ | ✓ | ✗ | ✓ | E or C |
| Ours | | ✓ | ✓ | ✓ | ✓ | ✓ | E,C |

\* Independence from a central cloud entity.
\*\* Bandwidth usage optimization.
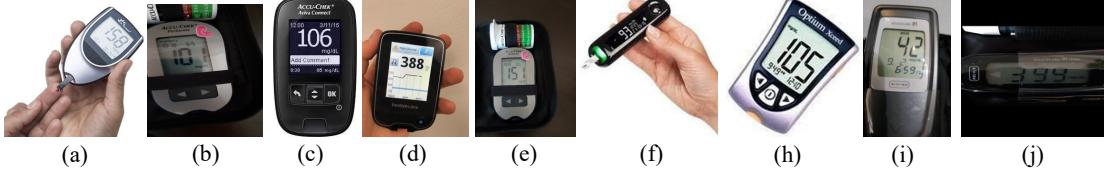D: Direction   A: Approach   E: Edge   C: Cloud

| (a) | (b) | (c) | (d) | (e) | (f) | (h) | (i) | (j) |

Fig. 4. Examples of images captured from medical devices.



(a) Cloud-Only     (b) Mobile-Only     (c) Split Computing
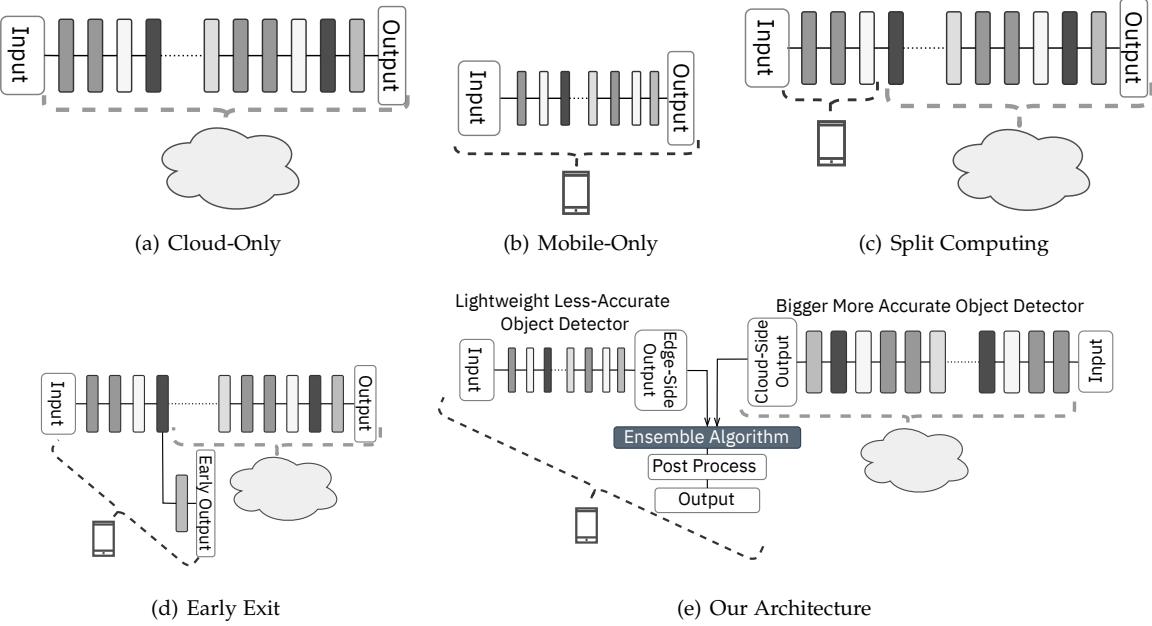
(d) Early Exit     (e) Our Architecture

Fig. 5. Deep Learning Serving Architectures.

mainly focus on optimizing response time and memory footprint, and thereby sacrifice accuracy (Fig. 5(b)).

The second direction, distributes computation across edge and central cloud, vertically (Fig. 5(c)). Some techniques aim to reduce the required computation and bandwidth by dropping insignificant frames of input stream at the edge, before sending them to the cloud, and depending on the nature of a task may follow different filtering policies [22], [23], [24]. Some distribute the inference model across edge and cloud (Split Computing) [25], [26], [27], and usually trade accuracy for better response time or bandwidth usage. These methods rely on central cloud; hence, in the case of network outages, they become unavailable.

Another recent direction determines one or more exit points within neural model (including pre-processing steps). Exit points are usually designed so that at least one of them stay on the edge. Depending on the task and its requirements, the model can exit early, sacrificing accuracy to meet delay constraints [27], [28], [29]. Despite their appealing results, they still are immature and are evaluated on simple tasks such as classification (Fig. 5(d)).

In comparison, the strength of our proposed mobile-cloud architecture (Fig. 5(e)) is its ability to take advantage of both worlds: it improves accuracy, and in general performability as well as availability, thanks to the independent nature of our ensemble models.

## 3 PROPOSED METHOD

Our service captures images from medical devices using phone's camera (Fig. 3(1)), then performs a pre-processing step on mobile (Fig. 3(2)). The mobile device concurrently sends the prepared image to cloud (Fig. 3(3.1, 4)), while executing light-weight inference engine locally (Fig. 3(3.2)). After receiving predictions of both sides, the mobile device runs our ensemble algorithm (Fig. 3(5)). It also performs a post-processing step, to correct the initial answer that is produced by the mobile, if required (Fig. 3(6)). If user recognizes a misprediction, they can send the image and its corresponding true reading (Fig. 3(7)) to cloud storage for future analysis and training iterations (Fig. 3(8)). The mobile model provides availability, while the cloud model improves performability (Refer to §5.5). Both models together improve accuracy.

We reduce the problem into object detection so that digits of the sensed value (and not other digits illustrating noisy information such as temperature, time and date) are objects of interest, together with a post-processing step to reorder objects and prepare the final response. DNN-based methods have shown remarkable results, but to get the most out of them, we need a well-annotated training dataset, which is not available for our problem. We design a better workflow and a data generation algorithm to automatically synthesize thousands of training images with well-aligned annotations. Our data generation workflow is described in detail in §4.

We use mobile smartphones to capture images from medical devices and to extract useful information. While inference on the mobile device provides low latency and high availability for users, it usually suffers from low accuracy due to resource constraints. On the other hand, inference on the cloud provides accuracy and performability, but

high latency and unavailability in the case of poor network connection are its big weaknesses. As depicted in Fig. 3, our proposed hybrid mobile-cloud serving architecture takes advantage of both mobile and cloud computing. Additionally, we designed a specialized ensemble model to further improve our desired performance metrics.

## 3.1 Deep Learning Models

We reduce the problem of image-based reading from screen of medical devices to an object detection and post-processing task. Each digit and decimal point belonging to the sensed value occupies a region of interest (RoI) and has a corresponding class (Fig. 2-1). In post-processing step, we convert a group of objects to a meaningful string (Fig. 2-2). For the object detection part, we design and train two convolutional models (CNNs) based on single-shot detector (SSD) architecture [31].

In general, we have powerful resources on cloud, but limitation at the mobile. Thus, we consider two proportionate backbones for our SSD architectures. One is highly optimized for smartphones and has fewer parameters; hence, lower accuracy. Another one is more accurate and consumes more resources so that can not be deployed on mobile. We prefer SSD-based networks for both sides, because they achieve better end-to-end latency [32].

### 3.1.1 Cloud-Side Model

We employ ResNet-50 [33] as the backbone of our SSD model on the cloud-side. It contains 16 convolutional blocks with shortcuts and one fully connected layer. It has more than $25M$ parameters, and requires $4$ billion multiplication-accumulation operations (MACs) per sample. We remove its last few layers including the classification head to use the remainder as a feature extractor backbone for the SSD architecture. Since the size of RoIs in input images varies, we utilize feature pyramid network (FPN) [34]. It leverages intrinsic pyramid structure of modern CNNs (i.e., ResNet) to generate multi-scale feature maps. FPN can improve accuracy of the model, and is much more compute efficient than feeding an image with different sizes multiple times.

### 3.1.2 Mobile-Side Model

Similar to the cloud-side, we stacked up a CNN-based classification model as a feature extraction backbone for our mobile-side SSD architecture. However, we employed edge-friendly architectures both for the backbone and the detection head. Our mobile-side model is directly inspired by the current state-of-the-art detection architecture for mobile devices, the MobileDets [35]. The reason of using depth-wise separable convolutional layers instead of regular convolutional layers was their fewer parameters and MACs, while hoping these metrics directly relate to less train and inference time. However, recent studies [36], [37] have shown network parameters and MACs may not be good proxies to model inference throughput and latency as they are not the only factors. Therefore, MobileDets expands the neural architecture search space by adding regular convolutions as well. After multiple experiments on different backbone architectures, we found the architecture proposed by [35]

for Mobile CPUs is the most appropriate for our mobile-side model. Since mobile devices are resource-constrained, we do not leverage FPN.

## 3.2 Ensemble Algorithm

As illustrated in Fig. 3, we design and deploy a light-weight engine at the mobile along with a more compute-intensive and accurate one on the cloud. We integrate predictions of our two deep learning models through our ensemble algorithm (Algorithm 1). Having regions of interest (RoIs), labels and confidence scores of the objects detected by both models, our ensemble algorithm first finds corresponding RoIs among mobile and cloud predictions (line 3). They must have identical labels (e.g., class of digit '5') and their distance should not be more than a tolerable amount $\epsilon$. The intuition behind is, the RoI coordinates predicted by different object detection models may not exactly match with each other. $\epsilon$ controls maximum distance between two corresponding predictions from two different models. There may exist more than one object with same labels. So if we do not check that, it will lead to misrecognition of RoIs. However, strict comparison is a bad idea because two quite different deep learning models may have small differences in RoI refinement after training. Hence, we add a tolerable distance to make it more suitable.

After finding, when both scores are higher than specified thresholds, the ensemble algorithm adds them up as the final confidence score (line 7). Otherwise, picks the highest score (line 9). In the case of only one prediction for a

---

**Algorithm 1** Ensemble Learning Algorithm

**Input:** ① maximum number of RoIs ($N$). ② RoIs matrices ($R$), and their corresponding vectors of ③ confidence ($\rho$) and ④ labels ($L$). ⑤ confidence thresholds ($T$). ⑥ tolerable distance ($\epsilon$).
  ▷ Each $R_l^K$ comprises four member points: left ($x_{min}$), right ($x_{max}$), up ($y_{min}$), down ($y_{max}$).
  ▷ Superscripts $C$, $M$ and $E$ stand for *Cloud*, *Mobile (Edge)* and *Ensemble*, respectively.

**Output:** A string equivalent to final prediction (e.g., "10.6").

1: $R^E \leftarrow \emptyset,\ \rho^E \leftarrow \emptyset,\ L^E \leftarrow \emptyset$
2: **for** $i \leftarrow 1$ **to** $N$ **do**
3:  Find $j$ such that:
   $-R_i^C[x_{min}] \approx R_j^M[x_{min}] \pm \epsilon$ **and**
   $-R_i^C[x_{max}] \approx R_j^M[x_{max}] \pm \epsilon$ **and**
   $-L_i^C = L_j^M$
4:  **if** $j \neq \emptyset$ **then**
5:   $L^E \leftarrow L_i^C \cup L^E,\ R^E \leftarrow R_j^M \cup R^E$
6:   **if** $\rho_i^C \geq T_i^C$ **and** $\rho_j^M \geq T_i^M$ **then**
7:    $\rho^E \leftarrow (\rho_j^M + \rho_i^C) \cup \rho^E$
8:   **else**
9:    $\rho^E \leftarrow \max(\rho_j^M, \rho_i^C) \cup \rho^E$
10:   **end if**
11:   $R^M \leftarrow R^M - R_j^M, \rho^M \leftarrow \rho^M - \rho_j^M, L^M \leftarrow L^M - L_j^M$
12:  **else**
13:   $L^E \leftarrow L_i^C \cup L^E,\ \rho^E \leftarrow \rho_i^C \cup \rho^E,\ R^E \leftarrow R_i^C \cup R^E$
14:  **end if**
15: **end for**
16: $R^E \leftarrow R^M \cup R^E, \rho^E \leftarrow \rho^M \cup \rho^E, L^E \leftarrow L^M \cup L^E$
17: $I \leftarrow \{i \in \mathbb{N} \mid i \leq |\rho^E|,\ \rho_i^E \geq T^E\}$
18: $L^E \leftarrow [L_{i \in I}^E \mid \forall(i,j) : i < j \Leftrightarrow R_i^E[x_{min}] \leq R_j^E[x_{min}]]$
19: **return** $\|_{i=1}^n L_i^E$   ▷ ($\|$ concats every element within $L^E$).

particular object on either mobile or cloud, simply adds it to our final results (lines 13, 16). The second elimination step happens in (line 17), where all remaining objects are compared via a higher universal threshold. Eventually, our algorithm reorders the remaining objects by their placement within image (line 18), and concatenates a sequence of corresponding labels to generate the output string (line 19).

### 3.3 Post-Processing Algorithm

To remove redundant objects more, and to improve accuracy even further, we add an additional post-processing algorithm right before the line 18 in Algorithm 1. We apply a modification of Non-Max Suppression (NMS) [38] to be compatible with our problem. Here, our objective is to find and remove those RoIs that significantly overlap with each other. It means that there are more than one object of interest in a region while it must not. Leveraging NMS reduces false positive (FP). Our modified NMS is defined in Algorithm 2. It executes iteratively. Having RoIs and their corresponding confidence scores and labels, it first sorts RoIs based on their confidence in a descending order (line 7). It then removes those objects that their overlap of the occupied area with another object is high enough ($> T_{nms}$), and their confidence is simultaneously lower (lines 13, 14). The overlap between two objects is calculated by intersection over union (IoU) of their corresponding RoIs. Although this works well, using a single threshold for all labels introduces a problem in our specific task. The *unit* label, representing the decimal point of sensed numbers, considerably overlaps with other objects, but it must not be removed. This also applies to the classes of '1' and '7'. To avoid that, we consider the *unit* label separately in our calculations within NMS procedure (lines 2-6).

---

**Algorithm 2** Post-Processing Algorithm

**Input:** ① overlap threshold ($T_{nms}$). ② RoIs matrix ($R$), and its corresponding vectors of ③ confidence ($\rho$) and ④ labels ($L$).
**Output:** reduced RoIs matrix, and its corresponding vectors of confidence and labels.

1: $R^M \leftarrow \emptyset,\ \rho^M \leftarrow \emptyset,\ L^M \leftarrow \emptyset$
2: $i \leftarrow j \in \mathbb{N} \mid j \geq |L|, L_j = '.'$     $\triangleright$ ('.' $\equiv$ decimal point)
3: **if** $i \neq \emptyset$ **then**
4:    $L^M \leftarrow L_i,\ \rho^M \leftarrow \rho_i,\ R^M \leftarrow R_i$
5:    $L \leftarrow L - L_i, \rho \leftarrow \rho - \rho_i, R \leftarrow R - R_i$
6: **end if**
7: $(L, \rho, R) \leftarrow [(L_i, \rho_i, R_i) \mid \forall (i,j) : i < j \Leftrightarrow \rho_i \geq \rho_j]$
8: **while** $R \neq \emptyset$ **do**
9:    $R^M \leftarrow R_1 \cup R^M,\ R \leftarrow R - R_1$
10:    $L^M \leftarrow L_1 \cup L^M,\ L \leftarrow L - L_1$
11:    $\rho^M \leftarrow \rho_1 \cup \rho^M,\ \rho \leftarrow \rho - \rho_1$
12:    **for** $i \leftarrow 1$ **to** $|R|$ **do**
13:      **if** $IoU(R_i, R_1^M) \geq T_{nms}$ **then**
14:        $R \leftarrow R - R_i,\ L \leftarrow L - L_i$
15:      **end if**
16:    **end for**
17: **end while**
18: **return** $R^M \cup R,\ L^M \cup L,\ \rho^M \cup \rho$

---

### 3.4 Bandwidth Optimization

Our target usecase is in underdeveloped countries, thus more often than not, some users may be located in areas or situations that do not have access to the internet or their connection is poor, e.g., due to limited available bandwidth or network congestion problems. Therefore, we prefer providing the highest availability at the cost of less accurate response, to increase performability, in general. One can still get the mobile answer in zero-connection situations. Nevertheless, user will experience accuracy degradation. For poor network conditions, we present a simple image compression technique which reduces the bandwidth usage when sending captured images (Algorithm 3). We downscale the image, transform it from RGB to HSV colorspace, and then only send the value ($V$) channel of the image to the cloud. On the cloud, $H$ and $S$ are filled with predefined constants and are up-scaled to the original dimension before performing inference.

---

**Algorithm 3** Simple Lossy Compression

**Input:** ① $Img$ ② output size ($\mathcal{H}\nu, \mathcal{W}\nu$) ③ filling constant ($\mathcal{K}$)
**Output:** At mobile : $Img_{comp.}$     At cloud : $Img_{decomp.}$

**At mobile**

1: $Img_{comp.} \leftarrow$ Resize $Img$ to ($3 \times \mathcal{H}\nu \times \mathcal{W}\nu$).
2: Transform $Img_{comp.}$ from $RGB$ colorspace to $HSV$.
3: $Img_{comp.} \leftarrow$ Drop channels $H$ (Hue) and $S$ (Saturation).
4: **return** $Img_{comp.}$

**At cloud**

5: $Img_{decomp.} \leftarrow$ Add channels $H$ and $S$, and fill pixels with $\mathcal{K}$.
6: Resize $Img_{decomp.}$ to original size of $Img$.
7: **return** $Img_{decomp.}$

---

### 3.5 Complexity of Proposed Algorithm

For every detected object we sweep the other objects. Similarly, this is done again in the post-processing algorithm, for fewer objects. The time complexity of our algorithm, therefore, is $O(N^2 + N'^2)$ where $N$ is the number of predicted RoIs and $N'$ is the number of objects to post-process. Since $N' < N$, the complexity can be written as $O(N^2)$. We use one-dimensional arrays to keep RoI boxes and confidence and label vectors, that their size depend on $N$. Consequently, the space complexity of the algorithm is $O(N)$.

## 4 DATA GENERATION

Data scarcity in image-based reading is a major reason for poor results in the previous works. The superpower of deep learning models (e.g., CNNs) is their ability to learn representations directly from data. Deep learning algorithms depend on either big volume of annotated data to be trained in a supervised manner or a model pre-trained on a related task; none are available for our task.

Related work commonly use conventional computer vision workflow which starts with data analysis and subsequently continues with data cleaning, pre-processing, feature extraction, feature selection and finally designing and deploying a suitable model (Fig. 6(a)). In contrast, we present a different workflow (Fig. 6(b)) which enables us to utilize deep learning models, outperforming previous algorithms by a large margin.
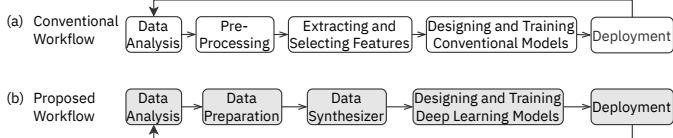
Fig. 6. (a) conventional computer vision workflow. (b) proposed workflow. Note that although some blocks are the same in both, their internal procedure might be completely different.

## 4.1 Data Preparation

After analyzing the images captured from displays of various glucometers, we discovered a number of common differentiating properties, summarized in Table 3. We learned that a lot of problems stem from the environmental conditions. For example, as depicted in Fig. 4, some images have been rotated (a, b, d, f, g), taken from different viewpoints (a, b, d, e, f), have poor contrast between the sensed quantity and its background (a, e, f), and have flashlight reflection (h), distorted and/or blurred (e, h). Some of the patterns from Table 3 are observable as well. For example, they may contain charts (d), have backlit (c, d, g), irregular screen shapes (g, h). And almost all of them are different in background color, display font style and type, etc.

This diversity in medical devices together with environmental variations introduce many difficulties for previous methods as they are based on feature engineering. Instead, we utilize such properties to synthesize a large training dataset with annotation. We gather 100 distinct images of medical devices and 150 open-source fonts (including seven-segment, dot-matrix, LCD, and LED styles) through the Internet. We then manually transform all images to obtain images with standard point of view. We also determine $\sim$20 different point coordinates per image corresponding to display screen corners and different items that can be on it (e.g., sensed value, measurement unit, date, etc.). We eventually feed all these preliminary information into our *Data Synthesizer* algorithm which can synthesize almost infinite distinct images of medical device with different sensed values.

## 4.2 Data Synthesizer

Our *Data Synthesizer* (Algorithm 4) receives the data prepared from the last step (line 1), namely a number of transformed images ($\mathcal{I}mages$) along with their corresponding point coordinates ($\mathcal{A}$) a set of fonts ($\mathcal{F}$), a degree-of-freedom

TABLE 3
Common differentiating properties and troublemakers in image-based readout from medical devices.

| # | Property | Pattern |
|---|---|---|
| 1 | Display | Different aspect ratios |
| 2 | | Different screen shapes |
| 3 | | LCD, LED, 7-segment, dot-matrix |
| 4 | | With or W/O backlit |
| 5 | | Regular or irregular screens |
| 6 | | Rounded or sharp corners |
| 7 | Value | Different measurement units ($mg/dL$, $mmol/L$) |
| 8 | (Label) | Usually, the biggest item on display |
| 9 | | No specific position on different displays |
| 10 | Additional *Items* | Include date, time, various signs and symbols, measurement unit, diagrams, etc. |
| 11 | | No specific position on different displays |
| 12 | | In some cases, fonts, styles, and colors differ from Value |

($\mathcal{D}$) for every item on the display and their properties, and maximum number of images to synthesize ($\mathcal{N}$). It then generates new images (line 4) and strict annotations (line 5) following Algorithm 4. Producing thousands of images from a small sample size, particularly to train deep learning models, is an important procedure, because any subtle noise or shift in distribution of training data leads to the overfitting issue: our models may reach to a reasonable accuracy on our synthesized training data, but perform poorly on real-world test data. Therefore, all artificially generated images must be quite similar to real images. On the other hand, to train the deep learning models, our training set should contain enough distinctive features. Also, our algorithm has to generate well-aligned annotations.

---

**Algorithm 4** Data Synthesizer

---

**Input:** ① source images $\mathcal{I}mages$, ② corresponding coordinates $\mathcal{A}$, ③ set of fonts $\mathcal{F}$, ④ degrees-of-freedom $\mathcal{D}$, and ⑤ maximum number of samples to generate $\mathcal{N}$.
**Output:** A set of $\mathcal{N}$ tuples each containing a synthesized image, its label (Value), and corresponding RoIs ($R^{\text{Value}}, R^{Items}$).

---

1: **procedure** SYNTHESIZER($\mathcal{I}mages^{Set}, \mathcal{A}, \mathcal{D}, \mathcal{F}, \mathcal{N}$)
2:   **for** $\iota \leftarrow 1$ **to** $\mathcal{N}$ **do**
3:     **repeat**
4:       $(Img, R^{Screen}) \leftarrow \text{Sample}_{Rand}(\ (\mathcal{I}mages^{Set}, \mathcal{A})\ )$
5:       $(\text{Value}, Items) \leftarrow \text{Generate}_{Rand}(Img, \mathcal{A}_{Img}, \mathcal{D}_{Img}, \mathcal{F})$
6:       $\forall \text{item } j \in Items \; \forall \text{object } k \in \text{item } j : \text{Calculate}(R_{j,k}^{Items})$.
7:       $\forall \text{object } i \in \text{Value} : \text{Calculate}(R_i^{\text{Value}})$.
8:     **until** $R_i^{\text{Value}} \cap R_{j,k}^{Items} = \emptyset$ **and**
          $\text{IoU}(R_i^{Items}, R_j^{Items}) < \epsilon, \; i \neq j$
9:     **repeat**
10:      $Img \leftarrow \text{Transform}(Img, \; Set)$
11:      $\forall R_i \in R^{\text{Value}} : \text{Calculate}(R_i)$     ▷ (Recalculates RoIs.)
12:     **until** $\forall R_i \in R^{\text{Value}} : \; R_i \cap (\mathbb{U} - R^{Screen}) = \emptyset$
13:     **store** $\{Img, R^{\text{Value}}, \text{Value}, R^{Items}\}$   ▷ Value $\equiv$ Label $L_{Img}$
14:   **end for**
15: **end procedure**

---

We use $\mathcal{D}$ to make sure the synthesized images and their corresponding annotations are both similar to real captured images and follow the patterns described in Table 3. We introduce controlled stochasticity to our synthesizing procedure. Randomness is a key enabler to increase mutual distinction between training samples, to prevent overfitting and improve eventual accuracy. Additionally, our algorithm randomly selects or generates different font styles, sizes, colors, background lights, time, date, and different indicators (e.g., temperature or blood drop symbol) (line 5). The font styles, font sizes, postions, and colors of additional information are selected from a range randomly as well. It validates the placement of each item afterwards (line 8). There must not be any overlap between the sensed Value and other additional information on screen ($Items$) such as time, date, temperature, etc. Also, there should not be any significant ($> \epsilon$) intersection between two different $Items$.

Consequently, it calls a transformation procedure (line 10) that modifies each generated image together with its annotation. Our transformation procedure consists of two major sections: *geometrical* and *visual*. In the geometrical part, we apply scaling, cropping, rotation, shearing, perspective transformation, and translation. In the visual part, we modify color, contrast, lightness, and sharpness. We insert noise, and arbitrarily drop some pixels to simulate

light reflections or scratches in display screen. The range and likelihood of each transformation depend on the generation set ($\mathcal{S}et$). In general, we designate a broader range and higher likelihood in the training set rather than validation set. The intuition behind is, we aim to generalize in the training time, while in validation, a set of data much more similar to real-world images is needed. Note that we never apply our *Data Synthesizer* on test set. In addition, to prevent any meaningful leakage, we make sure there is no similarity in medical devices in test set and training/validation set. After applying transformations, Value must still be within $Screen$ RoI ($R^{Screen}$). Otherwise, we apply transformation again until the condition is met. A few synthesized samples, generated using only one image from $\mathcal{I}mages$, is depicted in Fig. 7.

Value in synthesized images, ranges from 0 to 1000 and follows *discrete pseudo iid* distribution. Half of the numbers are integers. 35% are single decimal (e.g., 12.5), and remainders are double decimal (e.g., 1.25). The generated measurement unit must be compatible with its synthetic Value. Backlit and display background colors are not uniformly random, rather our algorithm selects colors that are more likely to appear in real-world devices with higher probability, but there is a probability to generate completely new colors in order to make the models robust against new devices or outliers. Each item $\in Items$ may randomly appear or disappear on $Screen$. The $\mathcal{I}mages$ are split into train ($\mathcal{I}mages^{train}$) and validation ($\mathcal{I}mages^{val}$) parts, in 90 and 10 shares, respectively. Each image comprises a different medical device.

## 5 EVALUATION

### 5.1 Experimental Setup

**Training Dataset.** We generated $1M$ distinct well-annotated training samples using our *Data Synthesizer*. The image size directly affects training throughput, inference latency, size of detection models, and the space needed to store the dataset. Therefore, we down-scale each image to $3 \times 320^2$ and save in JPEG format. Note that the input sizes of our CNN models are different. It occupies $45GB$ disk space and takes $\sim$11 hours to generate and store $1M$ samples.

**Test Dataset.** We collected 300 images directly captured by our clients from various glucometer devices without any modifications except down-sizing to the desired scale. As Fig. 8 depicts, the class of '1' has the most frequency of occurrence, and the class of '.' has the least among others. Additionally, we evaluate our method on another publicly available dataset from CameraLab from University of Oxford (referred to as Oxford for simplicity from now on) [6].

**Training.** We trained both our mobile and cloud neural models on a single server. We use an in-house GPU server (Table 6) to train our models. We designed and trained our models via TensorFlow framework. We reduce the time and resources needed for training the models by employing transfer learning from pre-trained weights on COCO dataset [39]. It helps our models extract better low-level features. We trained each model for 500-600k epochs ($\sim$3.5 days), and applied simple augmentations during training to prevent overfitting.

**Inference Setup.** For the proof of concept, we use the same server as the cloud side. For the mobile side, we use Galaxy Tab A (2019), a mediocre tablet (Table 5). We leveraged TensorFlow Lite for Android to convert, optimize, measure the performance of our mobile model. During the inference, input images get resized to $3 \times 416^2$ and $3 \times 350^2$ for the cloud and mobile models, respectively.

### 5.2 Accuracy Metric

We report the accuracy of our algorithm as formulated in (1). It is more illustrative in comparison with the previously reported formulas, namely Precision, Recall and F1-score for classification and localization. That said, it is the strictest as well, and is necessary due to medical nature of the values being read.

Given a prediction $\hat{Y}$ and its ground-truth $\Lambda$. Let:

$$\hat{Y} := \hat{y_1}\hat{y_2}\cdots\hat{y_m}$$
$$\Lambda := \lambda_1\lambda_2\cdots\lambda_n$$

$\hat{Y}$ is then considered as a *correct* prediction if:

$$m = n \quad \text{and} \quad \forall i : \hat{y}_i = \lambda_i,\ i \in \mathbb{N},\ i \leq n \tag{1}$$

A prediction is correct if ① all ground-truth RoIs are correctly detected, ② the labels are correctly predicted, ③ objects are in the same order as they are in the ground-truth, and ④ no additional object (e.g., from $Items$) is detected. For instance, assuming a ground-truth label *11.9* (Fig. 2), none of 119, 1.19, 111, 1149, 1109, 11, 19, 1.9, etc. are correct predictions, which makes it more difficult than Precision or Recall.

### 5.3 Accuracy on Our Test Set

We first evaluate our algorithms against 300 images directly captured by our clients from various glucometer devices without any modifications except down-sizing to the desired scale. As Fig. 8 depicts, the class of '1' has the most frequency of occurrence, and the class of '.' has the least among others. While cloud- and mobile-only predictions fluctuate around $89 - 90\%$ (Fig. 9), our ensemble model improves the accuracy by 7.3 percentage points. The reason behind is illustrated in Fig. 10. In each confusion matrix, every column represents objects in a predicted class while each row represents the objects in its actual class. Making matrices more informative, we added another row and column. The last column represents the objects in a true class that are ignored, and the last row shows the detected objects that do not belong to any classes. To get a better perspective, each matrix is scaled to one in rows. Our cloud model performs well on most classes except ''.' and '1'. On the other side, our mobile model performs poorly on '0', '4' and '9', but predicts '.' and '1' much better than the cloud. The last matrix confirms that our ensemble algorithm perform better than the single models. It also reduces false positives and false negatives. The mispredicted images usually contain some special objects/symbols on their screen. For example, the device shown in Fig. 12(f) displays an arrow on the screen which is pretty similar to '7', and sometimes misleads the models.
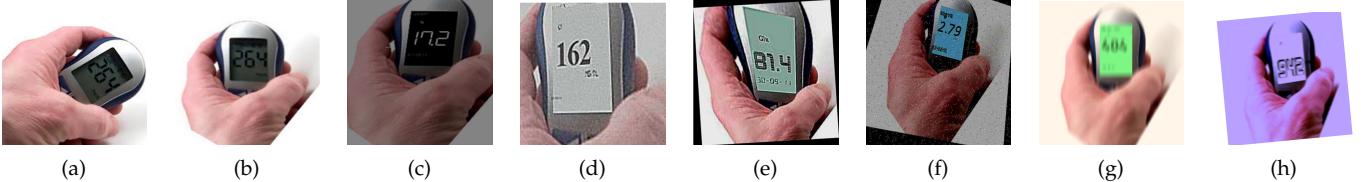
Fig. 7. A few examples of the images synthesized and transformed by Algorithm 4. (a) is the original image, and (b) is its manually transformed image. Note that for better understanding, only images generated from one source image are shown.
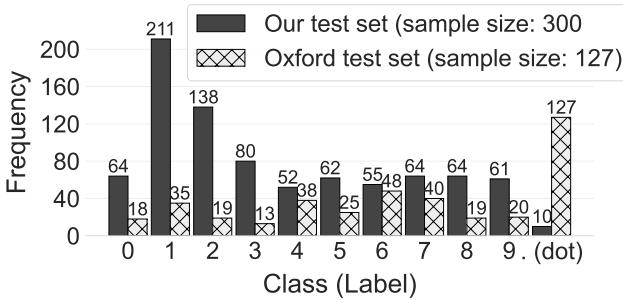


Fig. 8. Overview of the test sets used for evaluation. Our dataset contains 300 samples, and Oxford's [6] contains 127.
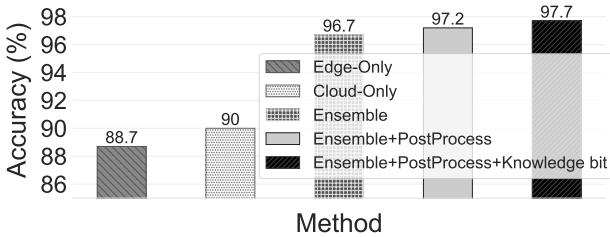


Fig. 9. Accuracy on our test set in different settings.

### 5.3.1 Post-Processing Impact

Our post-processing algorithm shows its positive effect on prediction of test samples. It increases accuracy by $0.5$ up to $1$ on our test set (Fig. 9). The main reason for fewer false positive errors is applying post-processing after the ensemble step (Fig. 10).

### 5.3.2 Device-Aware Prediction (Knowledge Injection)

The least accurate prediction in our proposed method belongs to the class '.' (decimal point) with $90\%$ accuracy. The current state-of-the-art [6] only detects objects and simply adds a decimal point manually. This is because their dataset contains two devices that both use $mmol/L$ as measurement unit, meaning that all values are floating point. In contrast, our dataset is more general and comprises $mg/DL$ as well. We can improve the accuracy of our algorithm by injecting one bit prior knowledge about the measurement unit of a medical device ($mg/DL$ or $mmol/L$). As shown in Fig. 9, it increases the accuracy by $0.5$. However, considering the displayed measurement unit on the screen as another object of interest can be used in future work.

### 5.4 Generalization: Accuracy on Oxford's Dataset

the Oxford's dataset [6] gives us the opportunity to evaluate robustness and generalization of our proposed algorithm. We can also have a direct comparison with the current state-of-the-art [6]. It consists of two different glucometer devices, each evenly divided into training and test sets. The two sets are quite similar (including the image capturing device),

except the values on the screens. Each set contains 127 samples, and all of them use mmol/L as the measurement unit on seven-segment displays (Fig. 8). The size of each image is $\sim 3 \times 2600 \times 4600$ pixels. We first report the accuracy of our models on the test set without learning the training set, to measure how robust our method is. We achieved **89.8%** accuracy in this configuration. Next, to get a fair comparison, we fine-tuned our models on the training set. In both cases, knowledge injection was not employed. All parameters except $T^E$ were fixed between the evaluations on the two datasets, we used $0.72$ for Oxford's and $0.82$ for ours. To prevent Class Imbalance Problem, the Oxford's training set was combined with a small part of our own training set. After fine-tuning the accuracy improved to **92.1%**, while the previous method achieved $51.5\%$ [6]. As the results suggest, our algorithm outperforms the state-of-the-art both with ($+40.6\%$) and without observing the training set ($+38.3\%$). Detailed assessment showed that in most misread images, devices are far from the camera so that in some cases it is difficult to read them correctly. The main reason is lower resolution of our models. Input resolution significantly impacts accuracy of RoI detection [32]. In general, reducing size of image by a factor of $4$ decreases accuracy by almost $16\%$ and lowers response time by $\sim 27\%$. The input resolution of our models are $<1\%$ of the images within Oxford's dataset. Thus, one may be able to achieve even lower error rates by increasing the input size of the neural models, but at the cost of longer training and response time. The memory and storage needed at the mobile-side are also the constraints that must be provisioned.

### 5.5 Availability and Performability

Our mobile device is able to prepare its prediction in $260ms$ and $360ms$ in GPU-enabled and only-CPU settings, respectively. For the situations that there is no internet connection, clients can still rely on the results of their smartphones. They will lose some accuracy, while higher availability can be achieved. When the connection quality is poor, users can enable our simple compression algorithm. It can reduce the bandwidth usage by **45×** while resulting $<\mathbf{1}\%$ degradation in accuracy, or by $90\times$ and $\sim 2\%$ accuracy loss. This is achieved by our *Data Synthesizer* algorithm. It forces our deep learning models to learn channel-invariable features. In addition, the ensemble approach inherently improves robustness. The response time breakdown in Fig. 11 depicts that by applying our compression technique, the transmission delay, which is the major contributor to total delay, collapses. Hence, users can still get a reasonable response in $<500ms$ with $512Kbps$ available bandwidth and $100ms$ RTT.

To compare related service architectures (Fig. 5), we assume a SLA in which *Performability* is described as *reaching*

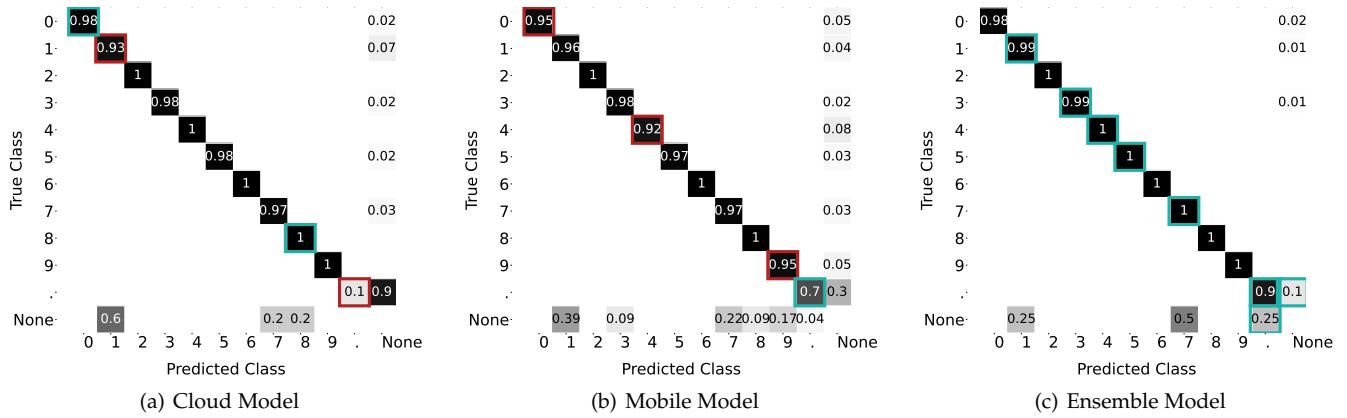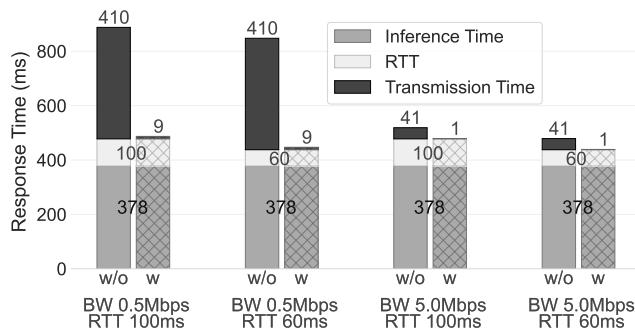(a) Cloud Model     (b) Mobile Model     (c) Ensemble Model

Fig. 10. Confusion matrix of different models. Matrices are row normalized.



Fig. 11. Cloud-side response time breakdown with (w) and without (w/o) applying compression, in two different bandwidth and RTT settings. RTT is the round-trip time, and can be measured by the command *ping*. Two RTT values of 100ms, 60ms assumed for illustration.

$L_k$ accuracy while being able to prepare response within $5s$. *Availability* can be achieved by responding in $<500ms$, as a soft deadline. Here, $L_1$ (✔✔) and $L_2$ (✔) stand for $>90\%$ and $>85\%$ accuracy, respectively. We evaluate each service in three different connection qualities: excellent, poor (e.g., due to congestion or limited bandwidth) and zero (no connectivity). As summarized in Table 4, our model performs better in all three cases. To the best of our knowledge, no comparable Split Computing or Early Exit method currently significantly compresses its intermediate data while losing little to no accuracy (rows 4 and 6).

## 6 CONCLUSION

We presented a mobile-cloud automated image-based glucometer reading system broadly applicable to various medical devices; our method uses the camera, the wireless communication, and the computing capabilities of the mobile phone to provide a low-cost alternative to expensive devices

available more in developed countries. Our deep learning-based ensemble algorithm together with our mobile-cloud service architecture achieves higher availability and performability compared with mobile-only, cloud-only, split computing, and early exit rival models. We proposed a data generation algorithm to address the data scarcity problem, and synthesized one million well-annotated samples. Note that the massive varieties in glucometer devices and their screen outputs, has conventionally posed challenges to applicability of existing techniques, but we instead took benefit from them in our data generation techniques to produce more training samples from existing photos, and thus to achieve a more robust model.

Our method is capable of proper readout even in imperfect conditions such as dark ambience, reflections on the screen, blurry and out-of-focus photos. Our ensemble algorithm efficiently combines results obtained from two separate DNN models. Specifically, we take into account the slight shifts in bounding boxes identified by the two models, as well as the special case of '.' symbol that, unlike other symbols, can significantly overlap with other detected objects.

We evaluated the accuracy of our method on two different real-world test sets, one collected from our users, and one from Camera Lab [6]. Our method achieved 97.7% on our test set, and 92.1% on Oxford's, outperforming the current state-of-the-art. Our results showed that our algorithm is robust and generalizes well on other datasets.

## REFERENCES

[1] J. E. Given, M. J. O'Kane, B. P. Bunting, and V. E. Coates, "Comparing patient-generated blood glucose diary records with meter memory in diabetes: a systematic review," *Diabetic medicine*, vol. 30, no. 8, pp. 901–913, 2013.

[2] D. Salvi, C. Velardo, L. Mackillop, and L. Tarassenko, "Algorithmic comparison of patient-reported blood glucose diary records with meters' memory in gestational diabetes," *Informatics in Medicine Unlocked*, vol. 20, p. 100397, 2020.

[3] "idia." [Online]. Available: https://idia.app/

[4] "Continuous Glucose Monitoring: Weighing the Pros and Cons." [Online]. Available: https://www.verywellhealth.com/continuous-glucose-monitoring-the-arrival-of-dexcom-5-3289566

[5] "Blood Glucose Meter Guide." [Online]. Available: https://diabetes.co.uk/diabetes_care/blood_glucose_monitor_guide.html

[6] E. Finnegan, M. Villarroel, C. Velardo, and L. Tarassenko, "Automated method for detecting and reading seven-segment digits from images of blood glucose metres and blood pressure monitors," *Journal of Medical Engineering and Technology*, vol. 43, no. 6, pp. 341–355, 2019.

TABLE 4
Overview of Availability (Avail.) and Performability (Perf.) in different methods and connections (Conn.).

| Service | Excellent Conn. | | Poor Conn. | | Zero Conn. | |
|---|---|---|---|---|---|---|
| | Avail. | Perf. | Avail. | Perf. | Avail. | Perf. |
| Mobile-Only | ✔✔ | ✔ | ✔✔ | ✔ | ✔✔ | ✔ |
| Cloud-Only | ✔✔ | ✔✔* | ⇓ | ✔✔* | ✗ | ✗ |
| Split Comp. | ✔✔ | ✔✔* | ✗ | ✗ | ✗ | ✗ |
| Split Comp.+ | ✔✔ | ✔✔* | ⇓ | ✔✔* | ✗ | ✗ |
| Early Exit | ✔✔ | ✔✔* | ✔✔ | ✗or✔ | ✔✔ | ✗or✔ |
| Early Exit+ | ✔✔ | ✔✔* | ✔✔ | ✔✔* | ✔✔ | ✗or✔ |
| **Ours** | ✔✔ | ✔✔ | ✔✔ | ✔✔ | ✔✔ | ✔ |

+ Assuming intermediate data were compressed by a lossless method.
⇓ Noticeable degradation.
* Poses extra cost.

[7] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," *IEEE internet of things journal*, vol. 3, no. 5, pp. 637–646, 2016.

[8] T. Morris, P. Blenkhorn, L. Crossey, Q. Ngo, M. Ross, D. Werner, and C. Wong, "Clearspeech: A display reader for the visually handicapped," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 14, no. 4, pp. 492–500, 2006.

[9] R. P. Ghugardare, S. P. Narote, P. Mukherji, and P. M. Kulkarni, "Optical character recognition system for seven segment display images of measuring instruments," in *TENCON 2009-2009 IEEE Region 10 Conference*. IEEE, 2009, pp. 1–6.

[10] M. Mariappan, V. Ramu, T. Ganesan, B. Khoo, and K. Vellian, "Virtual Medical Instrument for OTOROB based on LabView for acquiring multiple medical instrument LCD reading using optical charcater recognition," in *Proceedings from the International Conference on Biomedical Engineering and Technology (IPCBEE)*, vol. 11, 2011, pp. 70–74.

[11] S. Ghosh and S. Shit, "A low cost data acquisition system from digital display instruments employing image processing technique," in *2014 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*. IEEE, 2014, pp. 1065–1068.

[12] E. Tekin, J. M. Coughlan, and H. Shen, "Real-time detection and reading of LED/LCD displays for visually impaired persons," in *2011 IEEE Workshop on Applications of Computer Vision (WACV)*. IEEE, 2011, pp. 491–496.

[13] I. Rasines, P. Iriondo, and I. Díez, "Real-Time display recognition system for visually impaired," in *International Conference on Computers for Handicapped Persons*. Springer, 2012, pp. 623–629.

[14] "Tesseract Open Source OCR Engine." [Online]. Available: https://github.com/tesseract-ocr/tesseract

[15] M. K. Prakruthi, V. Kalyan, V. Suhas, M. G. Katti, and S. Pai, "Application of Convolutional Neural Networks in Mobile Devices for Inferring Readings from Medical Apparatus," *International Journal of Research and Scientific Innovation*, vol. IV, no. Vi, pp. 2321–2705, 2017.

[16] C. Liu, "Digits Recognition on Medical Device," Diss., University of Western Ontario, 2016.

[17] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. C. Chen, "MobileNetV2: Inverted Residuals and Linear Bottlenecks," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, jan 2018.

[18] X. Dai, P. Zhang, B. Wu, H. Yin, F. Sun, Y. Wang, M. Dukhan, Y. Hu, Y. Wu, and Y. Jia, "Chamnet: Towards efficient network design through platform-aware model adaptation," in *Proceedings of the IEEE Conference on computer vision and pattern recognition*, 2019, pp. 11 398–11 407.

[19] A. Howard, M. Sandler, G. Chu, L.-c. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan, G. Chu, L.-c. Chen, B. Chen, and M. Tan, "Searching for MobileNetV3 Accuracy vs MADDs vs model size," in *International Conference on Computer Vision*, 2019, pp. 1314–1324.

[20] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, "Mnasnet: Platform-aware neural architecture search for mobile," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2820–2828.

[21] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Enabling AI at the edge with XNOR-networks," *Communications of the ACM*, vol. 63, no. 12, pp. 83–90, 2020.

[22] Y. Li, A. Padmanabhan, P. Zhao, Y. Wang, G. H. Xu, and R. Netravali, "Reducto: On-Camera Filtering for Resource-Efficient Real-Time Video Analytics," in *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*, 2020, pp. 359–376.

[23] K. Hsieh, G. Ananthanarayanan, P. Bodik, S. Venkataraman, P. Bahl, M. Philipose, P. B. Gibbons, and O. Mutlu, "Focus: Querying large video datasets with low latency and low cost," in *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, 2018, pp. 269–286.

[24] J. Wang, Z. Feng, Z. Chen, S. George, M. Bala, P. Pillai, S.-W. W. Yang, and M. Satyanarayanan, "Bandwidth-efficient live video analytics for drones via edge computing," in *2018 IEEE/ACM Symposium on Edge Computing (SEC)*. IEEE, 2018, pp. 159–173.

[25] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," *ACM SIGARCH Computer Architecture News*, vol. 45, no. 1, pp. 615–629, 2017.

[26] D. Hu and B. Krishnamachari, "Fast and Accurate Streaming CNN Inference via Communication Compression on the Edge," in *2020 IEEE/ACM Fifth International Conference on Internet-of-Things Design and Implementation (IoTDI)*. IEEE, 2020, pp. 157–163.

[27] S. Laskaridis, S. I. Venieris, M. Almeida, I. Leontiadis, and N. D. Lane, "SPINN: synergistic progressive inference of neural networks over device and cloud," in *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking*, 2020, pp. 1–15.

[28] S. Teerapittayanon, B. McDanel, and H.-T. T. Kung, "Distributed deep neural networks over the cloud, the edge and end devices," in *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2017, pp. 328–339.

[29] L. Zhang, Z. Tan, J. Song, J. Chen, C. Bao, and K. Ma, "SCAN: A scalable neural networks framework towards compact and efficient models," in *Advances in Neural Information Processing Systems*, 2019, pp. 4027–4036.

[30] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.

[31] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.

[32] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy, "Speed/accuracy trade-offs for modern convolutional object detectors," *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, pp. 3296–3305, 2017.

[33] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, vol. 2016-Decem, 2016, pp. 770–778.

[34] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.

[35] Y. Xiong, H. Liu, S. Gupta, B. Akin, G. Bender, Y. Wang, P.-J. Kindermans, M. Tan, V. Singh, and B. Chen, "Mobiledets: Searching for object detection architectures for mobile accelerators," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3825–3834.

[36] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient Processing of Deep Neural Networks," *Synthesis Lectures on Computer Architecture*, vol. 15, no. 2, pp. 1–341, 2020.

[37] D. Stamoulis, "Hardware-Aware AutoML for Efficient Deep Learning Applications," 2020.

[38] A. Neubeck and L. Van Gool, "Efficient non-maximum suppression," in *18th International Conference on Pattern Recognition (ICPR'06)*, vol. 3. IEEE, 2006, pp. 850–855.

[39] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

# APPENDIX A
## MOBILE DEVICE AND CLOUD SPECIFICATIONS

TABLE 5
Mobile Device Specification

| Samsung Galaxy Tab A 8.0 with S Pen (2019, SM-P205) | |
|---|---|
| CPU | Octa-core big.LITTLE |
| | (2×1.8GHz Cortex-A73 and 6×1.6GHz Cortex-A53) |
| Chipset | Exynos 7 Octa 7904 (14 nm) – 64bit |
| GPU | ARM Mali-G71 MP2 |
| Memory | 3GB LPDDR4X |
| OS | Android 10 |
| Network | 2G/3G/4G(LTE) & Wi-Fi 802.11 |

TABLE 6
Cloud Server Specifications

| CPU | Intel Xeon E5-2630-v3 (x86_64) |
|---|---|
| | Frequency: 2.4GHz (1.2-3.2GHz) |
| | Cache: 32K-32K-256K-20480K |
| | #Cores: 24 (double threaded) |
| GPU | NVidia GeForce GTX 1080 Ti |
| | Total Memory: 11GB |
| | CUDA V9.0.176 , Compute Capacity 6.1 |
| Memory | 32GB (1600MHz) |
| OS | Gnu/Linux Ubuntu 18.04 |

# APPENDIX B
## RESPONSE TIME

Table 7 shows total response time on the mobile and cloud in different settings. For mobile inference, Tensorflow benchmarking tool for Android was used, and for the cloud, we measured the time it takes from sending a REST request to the server to get a response. The reported response time at the cloud side in TABLE 4 is the summation of inference time, typical transmission time (considering 5Mbps bandwidth), and typical round trip time (RTT equals 60ms).

TABLE 7
Response Time in Different Settings.

| Engine | Platform | Response Time (ms) |
|---|---|---|
| Mobile | CPU (1 Thread) | 494 ± 1 |
| | CPU (4 Threads) | 360 ± 25 |
| | **GPU (1 CPU Thread)** | **260 ± 1** |
| Cloud | CPU | **479 ± 20** |

# APPENDIX C
## BANDWIDTH OPTIMIZATION

Table 8 illustrates how our compression method affects bandwidth and accuracy with respect to different input resolutions.

TABLE 8
Bandwidth Usage Reduction Vs. Accuracy (%)

| Method | Original RGB | | Our Compression | |
|---|---|---|---|---|
| Image Size | Accuracy | BW Reduction | Accuracy | BW Reduction |
| 64×64 | 93.0% | 30× | 95.3% | 90× |
| 64×128 | 95.7% | 15× | 96.7% | 45× |
| 128×128 | 95.7% | 7.5× | 97% | 22.5× |
| 256×256 | 97.2% | 2.5× | 97.2% | 5.5× |
| 350×350 | 97.2% | 1.0× | 97.2% | 3.0× |

# APPENDIX D
## EXAMPLES OF MISRECOGNITION

We provide samples that our algorithm did not recognize them correctly. As shown in Fig 12, in some situations, additional elements on the screen mislead our algorithm. For example, in Fig. 12(f) and Fig. 12(c) an arrow and battery information are wrongly detected as '7' and '1', respectively. This can be corrected by adding similar samples to the training set in the next iterations.



(a) 53 → 5.3



(b) 61 → 67



(c) 5.81 → 5.8



(d) 99 → 399



(e) 94 → 9.4



(f) 1067 → 106

Fig. 12. Examples of mispredicted images. (prediction → ground-truth)