# A Dietary Nutrition-aided Healthcare Platform via Effective Food Recognition on a Localized Singaporean Food Dataset

Kaiping Zheng[1], Thao Nguyen[1], Jesslyn Hwei Sing Chong[2], Charlene Enhui Goh[3]
Melanie Herschel[4], Hee Hoon Lee[5], Beng Chin Ooi[1], Wei Wang[6], James Yip[7]

[1]School of Computing, National University of Singapore
[2]Ng Teng Fong General Hospital, Dietetics and Nutrition Department, Singapore
[3]Faculty of Dentistry, National University of Singapore
[4]Universität Stuttgart, Germany
[5]Ng Teng Fong General Hospital, Allied Health Division, Singapore
[6]TikTok Singapore
[7]National University Heart Centre, Singapore

**Abstract**

Localized food datasets have profound meaning in revealing a country's special cuisines to explore people's dietary behaviors, which will shed light on their health conditions and disease development. In this paper, revolving around the demand for accurate food recognition in Singapore, we develop the `FoodSG` platform to incubate diverse healthcare-oriented applications as a service in Singapore, taking into account their shared requirements. We release a localized Singaporean food dataset `FoodSG-233` with a systematic cleaning and curation pipeline for promoting future data management research in food computing. To overcome the hurdle in recognition performance brought by Singaporean multifarious food dishes, we propose to integrate supervised contrastive learning into our food recognition model `FoodSG-SCL` for the intrinsic capability to mine hard positive/negative samples and therefore boost the accuracy. Through a comprehensive evaluation, we share the insightful experience with practitioners in the data management community regarding food-related data-intensive healthcare applications.

The `FoodSG-233` dataset can be accessed via: `https://foodlg.comp.nus.edu.sg/`.

## 1 Introduction

Healthcare analytics aims to conduct numerous analytic tasks with a broad range of healthcare data spanning cost and claim data, pharmaceutical data, research and development data, clinical data, and patients' behavior and sentiment data. Based upon these, healthcare analytics serves as a basis for different data-intensive healthcare applications, such as chronic disease progression modeling [37, 45, 47, 44], and disease diagnosis [15, 26, 43, 42]. Among the diverse data sources for healthcare analytics, large-scale food data plays a significant role in revealing the dietary information and knowledge of patients, posing a profound impact on humans' health and well-being [30, 11]. For instance, long-term consumption of unhealthy foods increases the risk of developing diseases such as diabetes, hypertension, and hyperlipidemia. Therefore, there is an increasing demand for food data for diet assessment and disease management.

Naturally, food computing emerges as a subfield of computing, being accorded for its importance. It involves computational methods to analyze heterogeneous food data for solving disparate food-related problems and thereby, contributes to diverse data-intensive
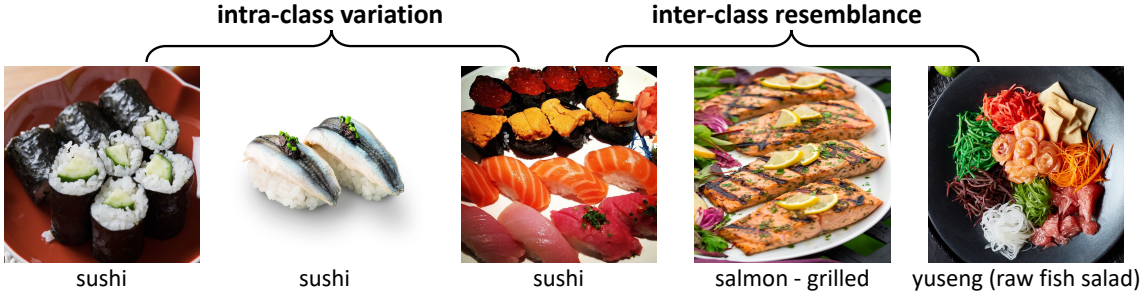
Figure 1: Two common issues in real-world food datasets.

healthcare applications [30]. Food recognition, which identifies food types from food images, is a cornerstone as it enables the effective assessment of people's diets whereby it provides much-needed intake information for disease prevention. Notably, combined with nutrition information, many other food-related applications become possible, such as estimating the intake of calories and nutrients.

Prior studies either make use of hand-crafted features for food recognition [17, 39, 12] or more recently turn to deep neural network models for boosted performance [22, 38], benefiting from the remarkable advancement of deep learning. As a consequence, they fuel the releasing of public food datasets [14, 16, 29, 23]. Although these datasets are highly useful in numerous studies, they lay emphasis on general food types or different geographical regions and hence, may not be applicable to more specific areas, in our case, Singapore. We note that localized food datasets are indispensable, as they manage to take into account each country's particular cuisines in ingredients, cooking styles, and varieties. Therefore, the localized food datasets could unveil valuable insights into the corresponding local human groups in dietary behaviors, which could ultimately cast light on their disease development.

Focusing on Singapore, we collect and curate a localized Singaporean food dataset in a systematic fashion, and the derived `FoodSG-233` dataset is of high quality in terms of volume and diversity compared with its counterparts. Built upon this `FoodSG-233` dataset, we collaborate closely with different healthcare sectors, clinicians, and dietitians on food recognition-based applications, and identify their shared requirements. We subsequently design and establish a full-fledged platform `FoodSG` driven by effective food recognition and dietary nutrition analysis, for supporting diverse data-intensive healthcare applications in Singapore. With the `FoodSG` platform as a service, we manage to bring huge benefits to healthcare management in Singapore from various aspects.

Probing into the `FoodSG-233` dataset, we observe two critical issues in this food dataset that are prevalent in real-world food datasets exemplified in Figure 1: (i) intra-class variation meaning that the food images in the same category may exhibit high dissimilarity in appearance, which gives rise to hard positive samples, and (ii) inter-class resemblance indicating that some dishes may look similar to each other despite being in different categories, which results in hard negative samples. The existence of such hard positive/negative samples tends to impede the food recognition approaches to achieve satisfactory performance. To address this issue, we introduce supervised contrastive learning (`SCL`) that encourages hard positive/negative sample mining into the food recognition task and propose the `FoodSG-SCL` model for recognizing food categories accurately. We summarize our main contributions below.

- Stemming from the shared requirements of diverse dietary management programs, we design and develop the `FoodSG` platform as a service for food-related analytics.

- `FoodSG` is currently serving a variety of healthcare-oriented applications and healthy living programs in Singapore, bringing benefits to the country and her people in the long run.

- We present real use cases to highlight the importance of localized datasets to the healthcare management of specific countries, and collect, clean, and release a localized Singaporean food dataset `FoodSG-233` with a systematic curation pipeline in order to foster data management research in food computing.

- We identify the particular challenges of food images, namely intra-class variation and inter-class resemblance, which result in hard positive/negative samples hindering the food recognition performance. To address these, we introduce `SCL` into the model design and devise the `FoodSG-SCL` model to facilitate the mining of hard positive/negative samples for boosted performance.

- We conduct an extensive experimental evaluation to validate the effectiveness of `FoodSG-SCL`, which confirms the accuracy of `FoodSG-SCL` and simultaneously offers practitioners fresh insights into data-intensive healthcare applications.

Despite our focus on Singaporean food, our design of the `FoodSG` platform for accommodating different applications, and our systematic data preparation pipeline can be used for other countries' food computing and related research, exhibiting enormous potential in promoting similar research directions and improving healthcare.

## 2   `FoodSG` Overview

In this section, we present an overview of `FoodSG`. The architecture of `FoodSG`'s platform as a service (PaaS) for supporting various applications is illustrated in Figure 2. `FoodSG` is designed for the cloud to enable adoption by local healthcare providers. We develop its backend with node.js, and adopt Redis for caching. We use PostgreSQL as `FoodSG`'s backend database system for users' diet and exercise information collection and storage, etc. We also rely on `FoodSG`'s backend server to schedule `FoodSG`'s core component - food recognition driven by our devised `FoodSG-SCL` model for analytics. Specifically, on our curated and released `FoodSG-233` dataset, we train `FoodSG-SCL` integrating the Data Augmentation Module, the Encoder Module, the Projection Module, and the Prediction Module in a two-stage manner to gain the contrastive power of discriminating food images and thus, generating accurate predictions. We shall elaborate on the `FoodSG-233` dataset and the `FoodSG-SCL` model in Section 4 and Section 5, respectively. Further, the data stream is managed by Apache Kafka, and we adopt ejabberd [7] to embed the instant messaging service in `FoodSG` for facilitating communication between patients and their clinicians or dietitians.

Upwards, `FoodSG` supports diverse applications through Angular [6], which provides appealing and informative user interfaces. In order to support different users' terminals covering browsers, iOS, and Android, we further introduce Ionic [10] for a consistent and enjoyable user experience. Such a design provides the functionalities shown in Figure 3, spanning dietary intake recording, nutrition information retrieval, communication within the community, etc.
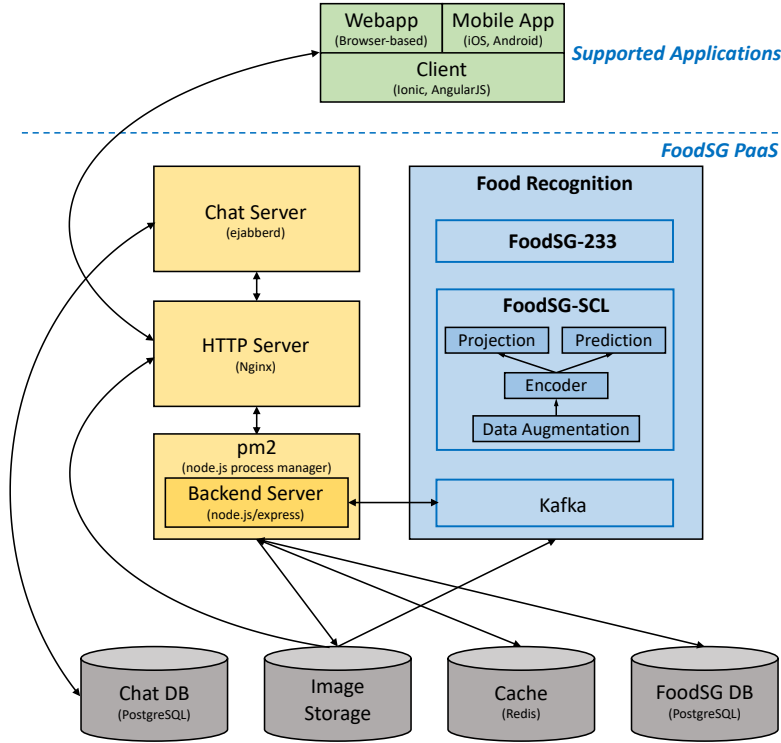
Figure 2: Overview of `FoodSG`'s architecture.

In summary, the `FoodSG` platform is equipped with the capability and flexibility of incubating diverse, multifunction applications with varied focuses and concerns.

## 3   Application Scenarios

In face of the growing prevalence of the diabetes epidemic and its resultant enormous burden on patients and economics, the Singapore Ministry of Health launched a campaign to combat diabetes in April 2016. The goal of this campaign is to facilitate a shared whole-of-nation effort to reduce the diabetes burden in Singapore [13, 32]. Afterward, diabetes as well as hypertension, and hyperlipidemia are all considered in urgent need of attention in Singapore [3].

Motivated by this, we explore and investigate food image analysis aided by dietary nutrition, as an infrastructure to assist in stopping or slowing the progression and complication development in these prevalent diseases. Developed up till now, our food recognition-based diet tracking and management platform `FoodSG` has supported a broad range of healthcare-oriented application scenarios in Singapore as a service. In this section, we present four representative ongoing scenarios, spanning three different levels from clinical medicine to healthcare, and further to public use.

**3H prevention.** 3H refers to (i) hyperglycemia (diabetes), (ii) hypertension (high blood pressure), and (iii) hyperlipidemia (high cholesterol) [3]. If a person develops these 3H problems without proper management, his/her lifetime risk of heart disease, stroke, kidney failure, etc could be increased greatly, leading to long-term medical attention and financial burden, and diminished quality of life. A practical strategy for combating 3H problems
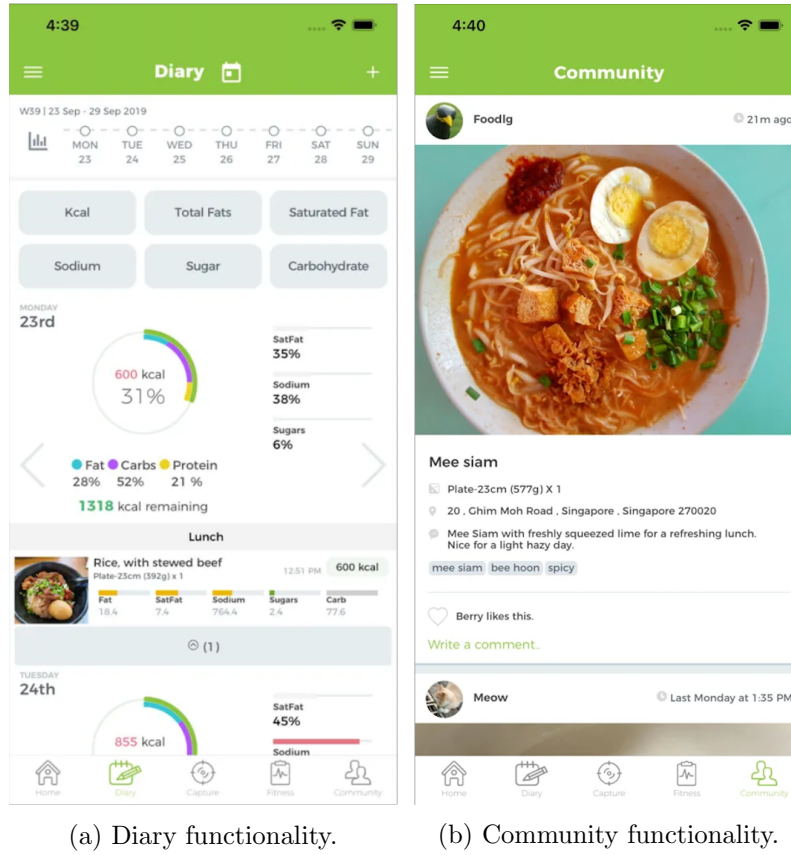
(a) Diary functionality.      (b) Community functionality.

Figure 3: Screenshots of `Food(lg)` powered by `FoodSG`.

lies in advocating healthier lifestyles, e.g., to exercise more, to eat healthier food, so that we can manage or even prevent 3H.

Let's take hyperglycemia, i.e., diabetes as an example, the tackling of which is a top priority in Singapore [2]. Since we believe that prevention is better than cure, we aim to incentivize people to adopt healthy lifestyles. In the Lifestyle Intervention (LIVEN) program in the Ng Teng Fong General Hospital in Singapore, we develop the JurongHealth Food Log (`JHFoodlg`) app powered by `FoodSG` to urge prediabetic patients to make behavioral lifestyle changes that are sustainable in the long term. We enroll the patients who are diagnosed with prediabetes, which is a pre-diagnosis of diabetes, indicating that the patient's blood sugar exceeds a normal range, but not reaching diabetic levels. 14.4% of Singapore's population is affected by prediabetes, among whom one-third will progress to type 2 diabetes in eight years eventually [1]. Fortunately, prediabetes is reversible with long-term lifestyle changes, such as disciplined diets and exercises.

In the program, we require the patients to use `JHFoodlg` in their everyday life to record their diets and exercises. `JHFoodlg` harnesses state-of-the-art food recognition models as the driving force, detects the food type and analyses the corresponding nutrients based on patients' uploaded photos, compiles their diets and exercise amount per day into a diary so that the patients can review their diet summaries against their weight-loss goals and exercise targets. `JHFoodlg` also embeds a health coaching component, which connects patients with the hospital's dietitians and physiotherapists to provide valuable information, from education on nutrition and fitness, to personalized feedback, and further

encouragement and recommendations. In this way, they can monitor the patients' progress and communicate with patients seamlessly via `JHFoodlg`, during the journey to reverse prediabetes.

We obtain promising results from a 6-month study in the LIVEN program. Almost all the patients who participate in the program experience a weight loss of between 4% to 5% of their initial body weight supported by `JHFoodlg` [1]. As an effective and essential app for 3H prevention in clinical medicine, we plan to roll out `JHFoodlg` for a larger cohort of patients in the future.

**Dental care management.** We develop a dental care management system `eDental` [46] powered by `FoodSG` to record, detect and analyze users' daily diets for potential risk factors for dental decay, with dentists and oral surgeons from the Faculty of Dentistry in National University of Singapore and National University Hospital. `eDental`'s key features include: (i) assisting users in logging their diets for comprehensive analysis, (ii) monitoring with attractive user interfaces and informative visual reports, (iii) triggering diet risk alerts for dental decay when necessary, (iv) incentivizing users to continue using the system through goal setting functions, and (v) providing essential educational support to users. Thereby, with `eDental`, we manage to facilitate dentists and patients to co-manage patients' daily dietary risk factors by analyzing their diet diaries.

This scenario demonstrates the applicability of `FoodSG`'s service among clinical medicine, healthcare, and public use in that the end users cover dentists, their patients, and the general public who wishes to practice healthy diets for preventing dental problems. We are currently working on a user study to assess the effectiveness of `eDental` in dental practice.

**Athletes' diet planning.** Our `FoodSG` platform also provides services for Singapore Sport Institute (SSI) [4], under Sport Singapore [5] the core purpose of which is to transform Singapore through sports and encourage individuals to live better through sports. In SSI, we have a more specific goal, i.e., we endeavor to provide the best support to athletes so that they can achieve their full potential and further fulfill their sporting aspirations. To achieve this, we turn to `FoodSG` for its capability of planning the athletes' diets, guaranteeing their food diversity to achieve nutrient balance, and facilitating their exercise accordingly to maintain appropriate body weights, body mass index, etc. On the basis of the athletes' diets and exercise, we can also suggest nutrients for them to supplement. Therefore, supported by `FoodSG`, we are able to unlock the potential in the athletes and help them pursue high-performance sports. In this scenario, we target to deliver high-quality healthcare to athletes via `FoodSG`. We are currently working on a 22-month user study to evaluate `FoodSG`'s positive influence on athletes' sports careers.

**Public use.** In order to serve public users, we release a public version of `FoodSG` as the "`Food(lg)`" app in iOS, Android, and Web browsers, establishing a uniform and consistent user experience. We detect the food types from users' journaled entries in images or text via state-of-the-art food recognition models and calculate users' daily nutrient estimates against the standard nutritional guidelines and food composition data from the Singapore Health Promotion Board (HPB) [9]. We also encourage users to exercise more, share food photos, and comment on them. This social feature lets the users motivate one another for maintaining a healthy lifestyle. In this way, `Food(lg)` helps general users to achieve a well-balanced diet, train their physical fitness, and hence, improve their health conditions and their quality of life in the long run.

In a nutshell, the aforementioned deployed scenarios showcase that our developed `FoodSG` platform manages to cater to diverse applications while relying on a robust data collection and curation pipeline that will be described in the next section.

# 4    Data Collection and Curation

Localized food datasets capture the unique characteristics of each country's dish varieties, cooking styles, and food ingredients, and are therefore highly indispensable and medically meaningful in contributing to people's health management in the specific country. In this section, we elaborate on the data collection and curation process of our released localized Singaporean food dataset `FoodSG-233`, which is crucial for food recognition as in Figure 2. Although our released dataset is specific to Singaporean food dishes, our proposed pipeline for data preparation is general and applicable to other types of food dishes, fostering similar research in other countries.

## 4.1    Data Collection

We incorporate popular ready-to-eat Singaporean food dishes by referring to HPB [9] for both food groups and dishes' names. HPB is a government organization committed to promoting healthy living in Singapore, which provides a large database of local Singaporean food and their corresponding nutrition facts [8]. The database is hierarchically structured with different coarse-grained food groups, each of which contains fine-grained food categories, e.g., the group "Sugars, sweets and confectionery" contains categories such as "Parfait" and "Popcorn". To build a comprehensive Singaporean food dataset, we select 233 popular ready-to-eat Singaporean dishes from 13 main food groups and then crawl candidate images from different search engines. Further, to expand the range of images retrieved, we increase the number of crawled images using a set of synonyms for each food category in our queries. Finally, we collect $226,809$ images for further curation.

## 4.2    Data Curation Pipeline

In the data curation process, we follow three rules of thumb.

- Release as many images as possible so that the released dataset can support different users' practical needs.

- Add necessary preprocessing to curate and clean the dataset in order to guarantee that after cleaning, (i) the images are of high quality, and (ii) the food recognition performance is satisfactory.

- Reduce human participation as much as possible, i.e., since manual efforts are generally expensive, involve human participation only when necessary.

Driven by these principles above, we design a pipeline shown in Figure 4 to curate our collected data and derive the `FoodSG-233` dataset. Next, we introduce each curation step in detail.
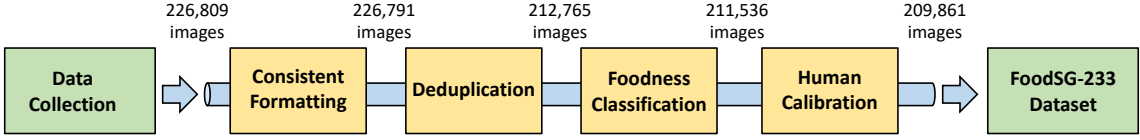
Figure 4: Data curation pipeline to derive the `FoodSG-233` dataset.

### 4.2.1 Consistent Formatting

Given the $226,809$ images collected, we first convert their formats to JPEG, and then remove 10 truncated images in the process. Meanwhile, we constrain that each image has a minimum size of $32 \times 32$ to guarantee the quality of retrieved images, and thus remove 8 images thereafter. After consistent formatting, $226,791$ images are fed to the next step.

### 4.2.2 Deduplication

We then remove exact or near duplicate images within each food category. Specifically, we calculate three different image hashes for each image, namely average hashing, perceptual hashing [28], and difference hashing. Based on the concatenation of these three hashes, we filter out similar images and keep a single image with the largest size when there are exact or near duplicate ones. We have $212,765$ images after deduplication.

### 4.2.3 Foodness Classification

Due to the existence of non-food images in the dataset, we further build a separate classifier to differentiate food and non-food images for filtering out the non-food ones, i.e., a foodness classification model. Specifically, we construct the training dataset via collecting $189,705$ non-food images from the Imagenet dataset [19] as negative images, and $226,037$ food images from prevailing food datasets including VIREO Food-172 [16], UEC-Food256 [23] and Food101 [14] as positive images. We then train a DenseNet169 model [21] for foodness classification and test the model on the current curated dataset. In the meantime, we conduct a first pass of manual checking to determine if each image is food, and such manual labels are then utilized as the ground truth to evaluate the foodness model's performance. Through this foodness classification step, $211,536$ food images are passed to the next step of the curation pipeline, with the foodness model yielding an accuracy of $82.9\%$.

### 4.2.4 Human Calibration

To guarantee the data quality, we conduct a second pass of manual checking as human calibration to correct the images that are assigned to the wrong food categories. We end up with $209,861$ images, which is the `FoodSG-233` dataset.

### 4.2.5 **FoodSG-233**

After curation, `FoodSG-233` has at least 400 food images per food category. The sorted distribution of food image number per food category in `FoodSG-233` is illustrated in
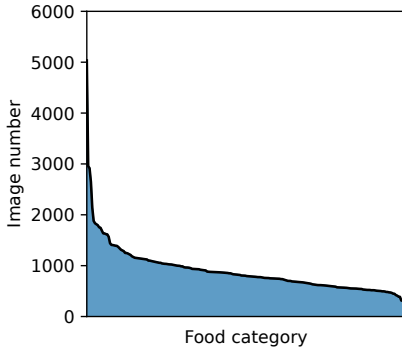
Figure 5: Sorted distribution of food image number per food category in `FoodSG-233`.
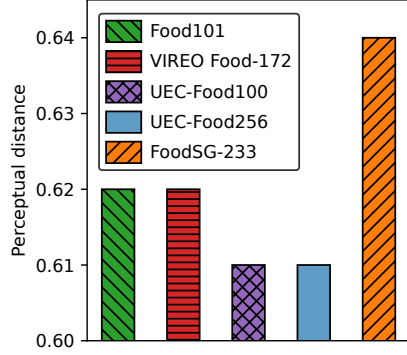


Figure 6: Diversity comparison in the perceptual distance (the larger, the more diverse).

Table 1: `FoodSG-233` vs. Existing Food Datasets.

| Dataset | Category # | Image # | Area | Image #/Category |
|---|---|---|---|---|
| Food101 [14] | 101 | 101,000 | Western | 1000 |
| VIREO Food-172 [16] | 172 | 110,241 | Chinese | 641 |
| UEC-Food100 [29] | 100 | 14,361 | Japanese | 144 |
| UEC-Food256 [23] | 256 | 25,088 | Multiple | 98 |
| `FoodSG-233` | 233 | 209,861 | Singaporean | 901 |

Figure 5, which tends to follow a power-law distribution, exhibiting the popularity of different food categories in the real world.

## 4.3   Comparison with Existing Datasets

We conduct a comparison between `FoodSG-233` and several widely adopted food datasets in Table 1. `FoodSG-233` is superior in terms of the total image number compared with its counterparts. In addition, as a localized Singaporean food dataset, `FoodSG-233` has a competitively large category number and image number per category.

As illustrated in Figure 1, the same food category tends to exhibit varied appearances in the real world. Hence, it is desired to release a representative dataset that covers the variations per food category and thereby, assists in building the food recognition model to tackle all sorts of user-uploaded food images effectively. With this goal, `FoodSG-233` is constructed to contain food images with a high intra-class variation in terms of color, shape, and texture in the collection process. To validate the superiority of `FoodSG-233` in this aspect, we conduct a quantitative diversity comparison from two perspectives below.

To start with, we employ a novel diversity measure based on a perceptual distance calculated via deep visual representations [41]. We use the open-source model to calculate the pairwise distance within each food category, and then average different categories' distances as the diversity metric. This metric has a value range of $[0, 1]$, and the larger the value is, the more diverse the dataset is in appearance. The comparison is shown in Figure 6, where `FoodSG-233` exhibits the highest value among all the datasets in the perceptual distance, confirming its diversity.

(a) Diversity comparison  (b) `FoodSG-233`  (c) UEC-Food256

Figure 7: `FoodSG-233` contains diversified images. (a) Diversity comparison in the lossless JPEG file sizes of five food categories' average images, i.e., the smaller the value, the more diverse. (b) Each compared food category's example images and the average image in `FoodSG-233`. (c) Each compared food category's example images and the average image in UEC-Food256.

We take a step further to calculate the average image of each food category, upon which we can compute the lossless JPEG file size to measure the dataset diversity from another perspective. This metric reflects the information amount in an image. The underlying rationale is that a food category with more diversified images corresponds to an average image that is vaguer, whereas a category with less diversified images corresponds to a clearer average image, e.g., with a more structured shape or a sharper appearance [19]. The average images are $256 \times 256$ and compressed into the lossless JPEG file format. Therefore, a more diversified food image set should exhibit a smaller size of the average image's lossless JPEG.

We compare `FoodSG-233` with UEC-Food256 which has the most shared food categories. The results are illustrated in Figure 7, including the comparison in the lossless JPEG file sizes of the average images for five food categories, the corresponding example food images, and the average images, respectively. As shown, `FoodSG-233` provides more diversified food images than UEC-Food256.

## 5 `SCL`-based Food Recognition

In this section, we introduce our proposed `FoodSG-SCL` model, which is essential to facilitate effective food recognition in `FoodSG` as discussed in Section 2.

### 5.1 Model Overview

The overview of `FoodSG-SCL` based on `SCL` [24] for food recognition is shown in Figure 8. `FoodSG-SCL` consists of four modules for data augmentation, encoding, projection, and prediction, respectively. The design rationale is that we first augment each data sample in the Data Augmentation Module to generate its multiple views, which facilitates `FoodSG-SCL` to learn generalizable and robust representations. The augmented samples
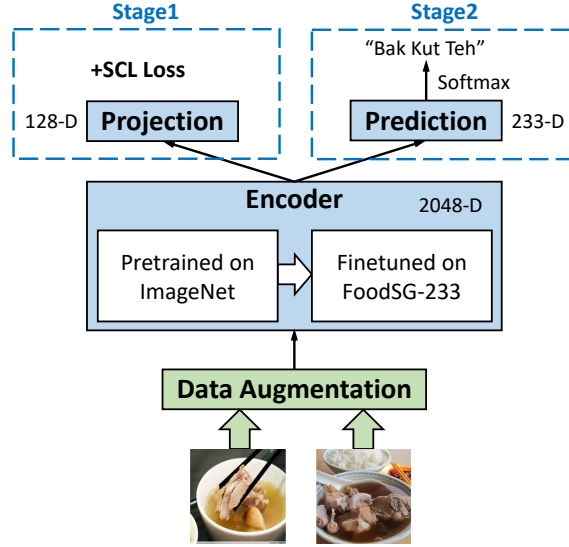
Figure 8: `FoodSG-SCL` model for effective food recognition.

are then input to the Encoder Module to unveil the underlying information in the augmented samples, using the state-of-the-art models pretrained on the ImageNet dataset, and then finetuning them on our `FoodSG-233` dataset. The Encoder Module generates a compact representation that is fed to both the Projection Module for further abstraction and integration of the `SCL` loss and the Prediction Module for providing the final prediction result.

## 5.2 Model Architecture

**Data Augmentation Module Aug$(\cdot)$.** Given a batch of data as input, we first adopt two advanced data augmentation mechanisms to generate the "multiviewed batch" [24] composed of $2N$ samples, where $N$ is the original batch size. We denote this multiviewed batch as $\mathcal{K} = \{1, \ldots, 2N\}$. As discussed above, such augmented data will be fed to the Encoder Module for further modeling.

In contrastive learning (CL) [18], a contrastive loss is proposed to maximize the agreement between each sample's different augmented view representations in the latent space. For instance, given an example $\mathbf{x}_i$, after the two data augmentations $\text{Aug}_1(\cdot)$ and $\text{Aug}_2(\cdot)$, we have $\tilde{\mathbf{x}}_i = \text{Aug}_1(\mathbf{x}_i)$, and $\tilde{\mathbf{x}}_{\rho(i)} = \text{Aug}_2(\mathbf{x}_i)$, where $i$ and $\rho(i)$ denote the two augmented samples from the same source sample. Hence, the index $i$ corresponds to the anchor sample, the positive sample set $\mathcal{P} = \{\rho(i)\}$ contains only one sample, and the remaining samples constitute the negative sample set.

`SCL`, short for supervised contrastive learning [24], extends the CL formulation to include the samples falling into the same class (with the same $y$ label) as $i$ in its positive sample set:

$$\mathcal{P} = \{p | p \in \mathcal{K} \wedge y_p = y_i \wedge p \neq i\} \tag{1}$$

In comparison, `SCL` extends the self-supervised setting in CL to a fully-supervised setting, which effectively leverages the label information and meanwhile, contributes to mining hard positive/negative samples for boosted performance.

11

**Encoder Module Enc($\cdot$).** We then employ a state-of-the-art neural network model as the Encoder Module to extract the compact representations from the augmented data samples. For a higher capacity of this module, we use the model pretrained on ImageNet, and then finetune it on `FoodSG-233`, grounded in the widely-acknowledged ability of convolutional neural networks in learning transferable features generalizable to diverse computer vision tasks [40]. After encoding, each sample $\tilde{\mathbf{x}}_i$ is converted to:

$$\mathbf{e}_i = \text{Enc}(\tilde{\mathbf{x}}_i) \tag{2}$$

where $\mathbf{e}_i \in \mathbb{R}^{d_e}$ and $d_e$ is the dimension size of the generated representation. $d_e$ is set to 2048 in our evaluation. We note that $\mathbf{e}_i$ is further normalized to fall in the unit hypersphere in $\mathbb{R}^{d_e}$ as suggested in [24], i.e., $\|\mathbf{e}_i\| = 1$. Such normalization is considered to bring performance improvement to prediction tasks.

**Projection Module Pro($\cdot$).** Next, we project the derived $\mathbf{e}_i$ into a more compact representation in this Projection Module:

$$\mathbf{s}_i = \text{Pro}(\mathbf{e}_i) = \varphi_{MLP}(\mathbf{e}_i), \varphi_{MLP} : \mathbb{R}^{d_e} \longmapsto \mathbb{R}^{d_p} \tag{3}$$

where a multi-layer perceptron (MLP) model $\varphi_{MLP}$ with one hidden layer is for projecting and abstracting the representation from $d_e$-dimensional to $d_p$-dimensional. We set $d_p$ to 128 in the following experiments. Similarly, $\mathbf{s}_i$ is normalized as $\|\mathbf{s}_i\| = 1$.

The `SCL` loss is then applied to pull positive samples close to each other and push them away from negative samples. As shown in Figure 8, the integration of `SCL` corresponds to **"Stage1"** in `FoodSG-SCL`, which learns a strong backbone for `FoodSG-SCL` in order to support the food recognition task. The learned $\mathbf{s}_i$ is discarded at the end of Stage1, while the output of the Encoder Module $\mathbf{e}_i$ will be used for prediction.

**Prediction Module Pre($\cdot$).** We input $\mathbf{e}_i$ to a linear prediction model $\psi$ for food recognition:

$$\mathbf{q}_i = \text{Pre}(\mathbf{e}_i) = \psi(\mathbf{e}_i), \psi : \mathbb{R}^{d_e} \longmapsto \mathbb{R}^{|\mathcal{C}|} \tag{4}$$

where $\mathcal{C}$ denotes the set of food categories for prediction, $|\mathcal{C}| = 233$ in the `FoodSG-233` dataset. We further parameterize $\hat{\mathbf{y}}_i = \sigma(\mathbf{q}_i)$, where $\sigma(\cdot)$ is the softmax function. The derived $\hat{\mathbf{y}}_i$ is the predicted probability distribution over $|\mathcal{C}|$ classes and will be used for training `FoodSG-SCL` in the optimization process. Besides, this Prediction Module corresponds to **"Stage 2"** as in Figure 8, the focus of which lies in employing the well-trained `FoodSG-SCL` model for food recognition in an effective manner.

## 5.3   Optimization

**Stage1.** The `SCL` loss for each sample $\tilde{\mathbf{x}}_i$ is:

$$l_i^{scl} = -\frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \log \frac{\exp(sim(\mathbf{s}_i, \mathbf{s}_p)/T)}{\sum_{k \in \mathcal{K} \setminus \{i\}} \exp(sim(\mathbf{s}_i, \mathbf{s}_k)/T)} \tag{5}$$

The contrastive prediction task is: given a sample $\tilde{\mathbf{x}}_i$, identify the positive samples in $\mathcal{P}$ out of all remaining samples $\mathcal{K} \setminus \{i\}$. In Equation 5, $sim(\cdot, \cdot)$ is the cosine similarity function, i.e., $sim(\mathbf{s}_i, \mathbf{s}_p) = \mathbf{s}_i^\top \mathbf{s}_p/(\|\mathbf{s}_i\|\|\mathbf{s}_p\|)$. As both representations are normalized, $sim(\mathbf{s}_i, \mathbf{s}_p)$ is equivalent to $\mathbf{s}_i^\top \mathbf{s}_p$. Moreover, $T$ denotes the temperature parameter. The overall `SCL` loss in Stage1 sums all samples' losses:

$$L^{scl} = \sum_{i \in \mathcal{K}} l_i^{scl} \tag{6}$$

With this loss function in Stage1, `FoodSG-SCL` (except the Prediction Module) can be effectively trained via gradient-based optimizers such as stochastic gradient descent (SGD).

**Stage2.** For the food recognition task as a multi-class classification problem, we adopt the cross-entropy loss in Stage2.

$$L^{ce} = -\frac{1}{|\mathcal{K}|} \sum_{i \in \mathcal{K}} \sum_{c \in \mathcal{C}} y_i^c \log(\hat{y}_i^c) \tag{7}$$

With the Encoder Module's output $\mathbf{e}_i$ and the Prediction Module, we can readily train the Predictor Module by optimizing this cross-entropy loss function.

Finally, after both stages, `FoodSG-SCL` is equipped with the contrastive power, contributing to the learning of hard positive/negative samples and further to the prediction performance.

# 6 Experimental Evaluation

In this section, we first introduce the experimental set-up and then describe the experimental results. After that, we analyze the experimental results and discuss the derived insights for practitioners on using the `FoodSG` platform for their specific data-intensive applications, which is imperative to the data management research community in setting its foot in food computing and starting to play an active role.

## 6.1 Experimental Set-up

**Evaluated Models.** We select five state-of-the-art image models for evaluation, and compare their performance against integrating them into `FoodSG-SCL` as the Encoder Module, in order to validate the efficacy of `SCL` in `FoodSG-SCL`.

- **ResNet50** [20] presents a residual learning framework that is easier to train and optimize, bringing considerably improved accuracy on numerous visual recognition tasks.

- **DenseNet169** [21] establishes direct connections between each layer and every other layer in a feed-forward manner, which enables the model to scale to more layers easily without causing optimization difficulties.

- **MobileNetV2** [33] improves the efficiency of mobile models via an inverted residual structure and emphasizes the importance of linear bottlenecks to maintain the representation power.

- **EfficientNetV2** [35] employs training-aware neural architecture search, model scaling, and progressive learning, for accuracy, training speed, and parameter efficiency.

- **ConvNeXt** [27] retains the advantages of standard ConvNets in simplicity and efficiency, but redesigns the components in a ConvNet towards the modern Transformer-based models. Therefore, ConvNeXt achieves superior performance matching up to such Transformer-based models in accuracy and scalability.

**Implementation Details.** The food recognition on `FoodSG-233` is formulated as a 233-class classification. We randomly split the food images per category into 80%, 10%, and 10% for training, validation, and testing, respectively. We adopt both the top-1 and the
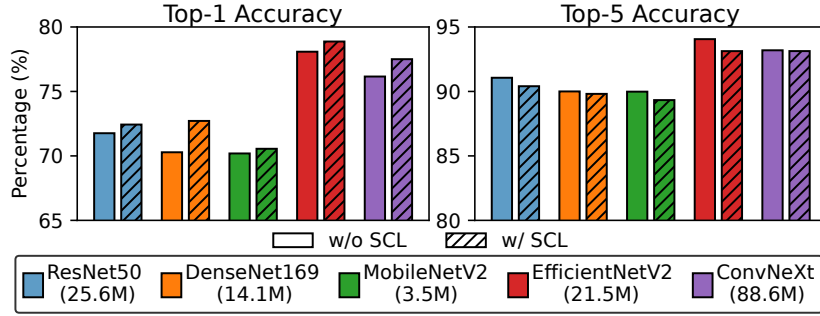
Figure 9: Comparison results of evaluated models w/o `SCL` and their advanced variants w/ `SCL` in `FoodSG-SCL`. Each model's parameter number is shown in brackets.

top-5 accuracy as evaluation metrics, and an accurate recognition model should yield high results in both metrics. We select the hyperparameters that achieve the best performance on validation data and apply such settings to the testing data for reporting the average experimental results of three independent runs.

For each of the aforementioned evaluated models, we adopt its model weights pretrained on ImageNet and then finetune them on `FoodSG-233` for evaluation. We train these models via SGD [34] with a learning rate of 0.05 and a weight decay of 0.0001. We adopt batch size differently across different models, as a larger model corresponds to a smaller batch and we set the batch size as large as possible within the device memory capacity. Specifically, the batch size is 512 in EfficientNetV2, 1024 in DenseNet169 and ConvNeXt, 2048 in ResNet50 and MobileNetV2, respectively.

For validating the effectiveness of `FoodSG-SCL`, we integrate the evaluated models as the Encoder Module of `FoodSG-SCL`, start with the same model parameters pretrained on ImageNet and then finetune on `FoodSG-233` with the other essential modules. SGD is used as the optimizer with the learning rate set to 0.1 and the epoch number set to 200. As `FoodSG-SCL` utilizes two views for each sample, its batch size needs to be halved accordingly.

**Experimental Environment.** We conduct the experiments in a server with Intel(R) Xeon(R) Gold 6248R $\times$ 2, 3.0GHz, 24 cores per chip. The server has 768GB memory, and 8 NVIDIA V100 with 300GB/s NVLINK. We implement the models with PyTorch 1.12.1.

## 6.2 Experimental Results

In Figure 9, we illustrate the experimental results of the evaluated image models without `SCL` and their respective advanced variants with `SCL` as well as other modules in our proposed `FoodSG-SCL`.

As shown, the advanced methods with `SCL` outperform the ones without `SCL` in top-1 accuracy by a large margin consistently, which confirms the effectiveness of the `SCL` mechanism. However, the integration of `SCL` leads to a slight degradation in top-5 accuracy. The underlying reasons for these phenomena will be explained in Section 6.3 in detail. Among the five evaluated models, EfficientNetV2 is the most accurate in terms of both top-1 and top-5 accuracy. In addition, despite the parameter number being one magnitude smaller than the others, MobileNetV2 still achieves moderate performance in both evaluation metrics.
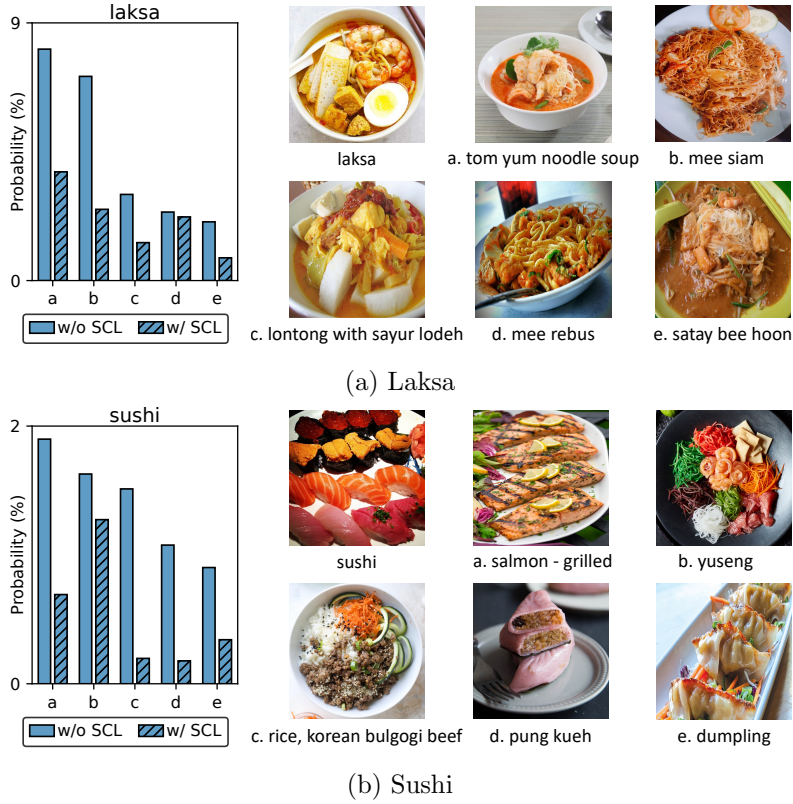
(a) Laksa



(b) Sushi

Figure 10: Representative food categories and their top 5 misclassified categories.

To further investigate the intrinsic property of SCL in FoodSG-SCL, we probe into the cases where food images are misclassified. We choose the following two ground truth food categories to explore: (i) the representative local food category "laksa" which is a spicy noodle dish popular in Singapore, and (ii) "sushi" which is relatively more widely known. For each ground truth category, we select the best-performing checkpoint among all evaluated models without SCL, i.e., EfficientNetV2 without SCL (according to Figure 9), and evaluate it on all the testing samples to derive the predicted probability distribution per testing sample. We then calculate the average of all testing samples' results, based on which we illustrate the top 5 misclassified categories and their respective probabilities predicted. For comparison, we display the results of the same five food categories derived by the most accurate model checkpoint of EfficientNetV2 with SCL as well.

The experimental results of both ground truth categories, including the predicted probabilities of their top 5 misclassified categories and the corresponding example food images are shown in Figure 10. For both laksa and sushi, the predicted top 5 misclassified categories generated by EfficientNetV2 without SCL are highly similar to the ground truth category in terms of color, texture, and shape. On the contrary, EfficientNetV2 with SCL effectively reduces the predicted probabilities of these five misclassified categories in both cases, the reason for which will be explicated in Section 6.3.

## 6.3 Discussions

**Top-1 accuracy vs. top-5 accuracy.** Generally, the top 1 predicted class is used as the final result. However, in practice, the top 1 prediction may not be completely

correct. To alleviate this issue, sometimes users are offered a set of five food category candidates and allowed to select the correct one with "one more touch". If the true food category still cannot be located within the candidate set, then the user experience will be negatively affected. In this sense, both top-1 accuracy and top-5 accuracy are crucial in food recognition, and they depict how accurately we can recognize the food dishes on the released `FoodSG-233` dataset.

In Figure 9, the highest top-1 accuracy is 78.87% yielded by EfficientNetV2 with `SCL`. Besides, EfficientNetV2 without `SCL` performs best in top-5 accuracy, i.e., 94.06%, whereas EfficientNetV2 with `SCL` achieves 93.13% with a minor decrease. This means that in most cases, the returned five candidate food categories successfully contain the ground truth, and users could select the correct one with one more touch conveniently. Both findings validate the effectiveness of `FoodSG-SCL` for food recognition.

**SCL in food recognition.** Food recognition is challenging in that certain food categories highly resemble each other and it could be difficult to distinguish between them even for humans. Furthermore, such subtle differences in the appearance of the two categories may correspond to a large difference in their constituent nutrients, which will have a huge impact on diet assessment. This consequently bolsters the demand for accurate food recognition when handling visually similar food images.

`SCL` is able to meet this demand. Specifically, `SCL`'s underlying principle lies in drawing samples with the same label close while separating them away from samples with other labels. This in nature improves top-1 accuracy, which is shown in Figure 9 where the integration of `SCL` brings performance benefits consistently.

However, such benefits of `SCL` do not come for free. The models without `SCL` adopt the standard cross-entropy loss and generate the top 5 categories that are the most similar to the ground truth category in appearance shown in Figure 10, due to the representation learning capability of convolutional neural networks. Different from them, `SCL`-based models tend to push the negative samples with different labels away from the positive samples, in spite of their resemblance in appearance. Therefore, with `SCL`, the predicted probabilities of the same five similar categories will be reduced to a large extent illustrated in Figure 10, resulting in slightly compromised top-5 accuracy as observed in Figure 9.

According to our real-world practice, we firmly believe that top-1 accuracy should outweigh top-5 accuracy in food recognition regardless of whether we allow users to further select from a candidate category set. Hence, we introduce the `SCL` mechanism into `FoodSG-SCL` to reduce the disturbing influence exerted by similar yet incorrect categories, providing excellent gain in top-1 accuracy.

**Edge Computing.** In healthcare-oriented applications, there are certain scenarios where users' data cannot be processed in a centralized manner, e.g., users are not willing to upload their food photos or share their nutrition intake information for privacy concerns. This requires the computation involved in food recognition to be performed on edge devices, such as mobile phones, which is thus affected by the computing capacity and space limitation of the devices. Under such scenarios, more lightweight models are preferred as the Encoder Module of `FoodSG-SCL`. For instance, MobileNetV2 exhibits a prominent advantage of parameter number being one magnitude smaller than other models, at the expense of accuracy decreasing from 78.87% (by EfficientNetV2 with `SCL`) to 70.55%.

# 7 Related Work

As an interdisciplinary area, food computing targets to apply computational methods to analyze heterogeneous sources of food data for supporting diverse food-related investigations in healthcare, gastronomy, and agronomy [30]. Among the involved tasks, food recognition is a cornerstone, which predicts the food items contained in food images.

Prior studies in the early stage exploit hand-crafted features for food recognition, such as color histograms and SIFT features [17], statistics of pairwise local features [39], and bag of features [12]. Recently, due to the record-breaking performance achieved by deep learning models in numerous areas such as computer vision [25, 31, 36], various convolutional neural network models [22, 38] are employed in food recognition for extracting the features in food images and hence, delivering more satisfactory recognition performance than the traditional approaches based on hand-crafted features.

In the meantime, a large number of food-related datasets are released to further promote the development of food computing. For example, Food101 [14] is constructed to cover 101 food categories (mostly western), containing $101,000$ images in total. Different from Food101, VIREO Food-172 [16] turns the focus to Chinese food and includes 172 food categories, and $110,241$ images, respectively. Moreover, a Japanese food dataset UEC-Food100 [29] is built to incorporate 100 popular categories. This UEC-Food100 dataset is then expanded to UEC-Food256 [23] introducing more categories from other countries, namely French, Italian, American, Chinese, Thai, Vietnamese, and Indonesian. These released multifarious food-related datasets present researchers with ample opportunity for food computing and analysis. However, these datasets have their specific geographic areas and therefore, only include particular types of cuisines. Different from existing work, driven by `FoodSG`, we investigate Singaporean food dishes and release the localized dataset `FoodSG-233`, which is of vital significance to `FoodSG`'s provided healthcare services for Singaporeans.

# 8 Conclusions

Focusing on the cuisines and healthcare-oriented applications in Singapore, we abstract the shared requirements among diverse applications and develop the dietary nutrition-aided platform `FoodSG` as a service to support all sorts of application scenarios in Singapore. We collect, curate, and release a Singaporean food dataset `FoodSG-233` systematically to promote future research directions for the data management community in food computing. Through delving into `FoodSG-233` to analyze its issues of intra-class dissimilarity and inter-class similarity, we propose to integrate `SCL` into food recognition and devise the `FoodSG-SCL` model accordingly to facilitate the learning from hard positive/negative samples for performance gains. By evaluating `FoodSG-SCL` from multiple perspectives, we deliver fresh insights and valuable experience in data-intensive healthcare applications to practitioners.

# References

[1] `https://www.todayonline.com/singapore/get-snap-happy-new-app-pre-diabetics-help-manage-health-condition`, 2019.

[2] `https://www.pmo.gov.sg/Newsroom/PM-Lee-Hsien-Loong-WHO-Global-Diabetes-Compact`, 2021.

[3] `https://aisingapore.org/grand-challenges/health/`, 2022.

[4] `https://www.sportsingapore.gov.sg/athletes-coaches/singapore-sport-institute`, 2022.

[5] `https://www.sportsingapore.gov.sg/`, 2022.

[6] Angular. `https://angular.io/`, 2022.

[7] ejabberd. `https://www.ejabberd.im/`, 2022.

[8] Energy & nutrient composition of food. `https://focos.hpb.gov.sg/eservices/ENCF/`, 2022.

[9] Health promotion board. `https://hpb.gov.sg/`, 2022.

[10] Ionic. `https://ionicframework.com/`, 2022.

[11] P. Achananuparp, E. Lim, and V. Abhishek. Does journaling encourage healthier choices?: Analyzing healthy eating behaviors of food journalers. In DH, pages 35–44. ACM, 2018.

[12] M. Anthimopoulos, L. Gianola, L. Scarnato, P. Diem, and S. G. Mougiakakou. A food recognition system for diabetic patients based on an optimized bag-of-features model. IEEE J. Biomed. Health Informatics, 18(4):1261–1271, 2014.

[13] Y. M. Bee, E. S. Tai, and T. Y. Wong. Singapore's "war on diabetes". The Lancet Diabetes & Endocrinology, 10(6):391–392, 2022.

[14] L. Bossard, M. Guillaumin, and L. V. Gool. Food-101 - mining discriminative components with random forests. In ECCV (6), volume 8694 of Lecture Notes in Computer Science, pages 446–461. Springer, 2014.

[15] Z. Che, D. C. Kale, W. Li, M. T. Bahadori, and Y. Liu. Deep computational phenotyping. In KDD, pages 507–516. ACM, 2015.

[16] J. Chen and C. Ngo. Deep-based ingredient recognition for cooking recipe retrieval. In ACM Multimedia, pages 32–41. ACM, 2016.

[17] M. Chen, K. Dhingra, W. Wu, L. Yang, R. Sukthankar, and J. Yang. PFID: pittsburgh fast-food image dataset. In ICIP, pages 289–292. IEEE, 2009.

[18] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton. A simple framework for contrastive learning of visual representations. In ICML, volume 119 of Proceedings of Machine Learning Research, pages 1597–1607. PMLR, 2020.

[19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In CVPR09, 2009.

[20] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In CVPR, pages 770–778. IEEE Computer Society, 2016.

[21] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In CVPR, pages 2261–2269. IEEE Computer Society, 2017.

[22] H. Kagaya, K. Aizawa, and M. Ogawa. Food detection and recognition using convolutional neural network. In ACM Multimedia, pages 1085–1088. ACM, 2014.

[23] Y. Kawano and K. Yanai. Automatic expansion of a food image dataset leveraging existing categories with domain adaptation. In ECCV Workshops (3), volume 8927 of Lecture Notes in Computer Science, pages 3–17. Springer, 2014.

[24] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan. Supervised contrastive learning. In NeurIPS, 2020.

[25] Y. LeCun, Y. Bengio, and G. E. Hinton. Deep learning. Nat., 521(7553):436–444, 2015.

[26] Z. C. Lipton, D. C. Kale, C. Elkan, and R. C. Wetzel. Learning to diagnose with LSTM recurrent neural networks. In ICLR (Poster), 2016.

[27] Z. Liu, H. Mao, C. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. A convnet for the 2020s. In CVPR, pages 11966–11976. IEEE, 2022.

[28] D. Marr and E. Hildreth. Theory of edge detection. Proceedings of the Royal Society of London. Series B. Biological Sciences, 207(1167):187–217, 1980.

[29] Y. Matsuda and K. Yanai. Multiple-food recognition considering co-occurrence employing manifold ranking. In ICPR, pages 2017–2020. IEEE Computer Society, 2012.

[30] W. Min, S. Jiang, L. Liu, Y. Rui, and R. C. Jain. A survey on food computing. ACM Comput. Surv., 52(5):92:1–92:36, 2019.

[31] B. C. Ooi, K. Tan, S. Wang, W. Wang, Q. Cai, G. Chen, J. Gao, Z. Luo, A. K. H. Tung, Y. Wang, Z. Xie, M. Zhang, and K. Zheng. SINGA: A distributed deep learning platform. In ACM Multimedia, pages 685–688. ACM, 2015.

[32] L. M. Ow Yong and L. W. P. Koe. War on diabetes in singapore: a policy analysis. Health Research Policy and Systems, 19(1):1–10, 2021.

[33] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L. Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In CVPR, pages 4510–4520. Computer Vision Foundation / IEEE Computer Society, 2018.

[34] I. Sutskever, J. Martens, G. E. Dahl, and G. E. Hinton. On the importance of initialization and momentum in deep learning. In ICML (3), volume 28 of JMLR Workshop and Conference Proceedings, pages 1139–1147. JMLR.org, 2013.

[35] M. Tan and Q. V. Le. Efficientnetv2: Smaller models and faster training. In ICML, volume 139 of Proceedings of Machine Learning Research, pages 10096–10106. PMLR, 2021.

[36] W. Wang, M. Zhang, G. Chen, H. V. Jagadish, B. C. Ooi, and K. Tan. Database meets deep learning: Challenges and opportunities. SIGMOD Rec., 45(2):17–22, 2016.

[37] X. Wang, D. A. Sontag, and F. Wang. Unsupervised learning of disease progression models. In KDD, pages 85–94. ACM, 2014.

[38] H. Wu, M. Merler, R. Uceda-Sosa, and J. R. Smith. Learning to make better mistakes: Semantics-aware visual food recognition. In ACM Multimedia, pages 172–176. ACM, 2016.

[39] S. Yang, M. Chen, D. Pomerleau, and R. Sukthankar. Food recognition using statistics of pairwise local features. In CVPR, pages 2249–2256. IEEE Computer Society, 2010.

[40] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In NIPS, pages 3320–3328, 2014.

[41] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In CVPR, pages 586–595. Computer Vision Foundation / IEEE Computer Society, 2018.

[42] K. Zheng, S. Cai, H. R. Chua, M. Herschel, M. Zhang, and B. C. Ooi. Dyhealth: Making neural networks dynamic for effective healthcare analytics. Proc. VLDB Endow., 15(12):3445–3458, 2022.

[43] K. Zheng, S. Cai, H. R. Chua, W. Wang, K. Y. Ngiam, and B. C. Ooi. TRACER: A framework for facilitating accurate and interpretable analytics for high stakes applications. In SIGMOD Conference, pages 1747–1763. ACM, 2020.

[44] K. Zheng, G. Chen, M. Herschel, K. Y. Ngiam, B. C. Ooi, and J. Gao. PACE: learning effective task decomposition for human-in-the-loop healthcare delivery. In SIGMOD Conference, pages 2156–2168. ACM, 2021.

[45] K. Zheng, J. Gao, K. Y. Ngiam, B. C. Ooi, and J. W. L. Yip. Resolving the bias in electronic medical records. In KDD, pages 2171–2180. ACM, 2017.

[46] K. Zheng, T. Nguyen, C. Liu, C. E. Goh, and B. C. Ooi. edental: Managing your dental care in diet diaries. In CIKM, pages 5059–5063. ACM, 2022.

[47] K. Zheng, W. Wang, J. Gao, K. Y. Ngiam, B. C. Ooi, and J. W. L. Yip. Capturing feature-level irregularity in disease progression modeling. In CIKM, pages 1579–1588. ACM, 2017.