

Long-Tailed Continual Learning For Visual Food Recognition

Jiangpeng He, *Member, IEEE*, Luotao Lin, Jack Ma,
Heather A. Eicher-Miller, and Fengqing Zhu, *Senior Member, IEEE*

Abstract—Deep learning based food recognition has achieved remarkable progress in predicting food types given an eating occasion image. However, there are two major obstacles that hinder deployment in real world scenario. First, as new foods appear sequentially overtime, a trained model needs to learn the new classes continuously without causing catastrophic forgetting for already learned knowledge of existing food types. Second, the distribution of food images in real life is usually long-tailed as a small number of popular food types are consumed more frequently than others, which can vary in different populations. This requires the food recognition method to learn from class-imbalanced data by improving the generalization ability on instance-rare food classes. In this work, we focus on long-tailed continual learning and aim to address both aforementioned challenges. As existing long-tailed food image datasets only consider healthy people population, we introduce two new benchmark food image datasets, VFN-INSULIN and VFN-T2D, which exhibits on the real world food consumption for insulin takers and individuals with type 2 diabetes without taking insulin, respectively. We propose a novel end-to-end framework for long-tailed continual learning, which effectively addresses the catastrophic forgetting by applying an additional predictor for knowledge distillation to avoid misalignment of representation during continual learning. We also introduce a novel data augmentation technique by integrating class-activation-map (CAM) and CutMix, which significantly improves the generalization ability for instance-rare food classes to address the class-imbalance issue. The proposed method is evaluated on Food101-LT, VFN-LT, VFN-INSULIN and VFN-T2D and show promising performance with large margin improvements compared with existing methods. We conduct an ablation study and discuss potential techniques that can further improve the performance, demonstrating great potential to deploy our method in real world food recognition applications.

Index Terms—Continual learning, long-tailed distribution, food recognition, knowledge distillation, data augmentation

I. INTRODUCTION

The emergence of modern deep learning technologies have enabled the automatic food nutrition content analysis, including image-based dietary assessment [1]–[4], to help monitor and inform interventions to improve dietary intake and to prevent chronic diseases such as diabetes. As the first and fundamental step of image-based dietary assessment, food recognition aims to identify food types given an input image and the overall dietary assessment performance greatly relies on the precise food recognition results. Though existing deep learning based food recognition methods [5]–[10] have achieved remarkable performance by training off-the-shelf Convolutional Neural Networks (*e.g.* ResNet [11]) using static datasets (*e.g.* Food-101 [12], Food2K [13]), there are still two major challenges when applying them in a real world

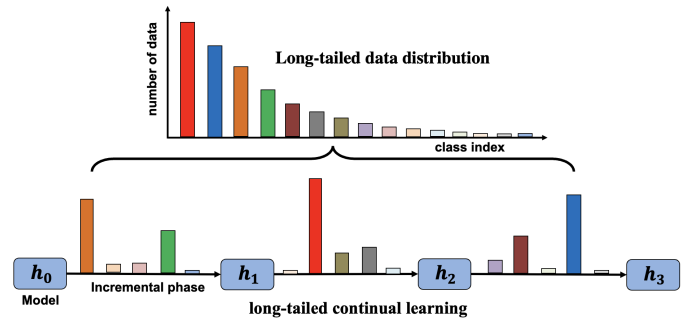


Fig. 1. The overview of long-tailed continual learning. New classes with imbalanced distribution arrive sequentially over time at each incremental learning phase. The updated model should be able to learn new classes continuously and perform classification on all the classes seen so far.

scenario including (i) how to update the model when new food classes appear sequentially overtime, and (ii) how to address the severe class-imbalance issue since the real life food images are usually in long-tailed distribution as shown in [14] where a minority of foods classes (*i.e.* instance-rich or head classes) are consumed more frequently than the remaining majority food classes (*i.e.* instance-rare or tail classes). The food recognition performance could drop dramatically without considering both two challenges. As shown in Figure 1, an ideal food recognition system in the real world should be able to learn new foods incrementally in long-tailed distribution with class-imbalance training samples at each incremental phase.

Continual learning, also known as incremental learning or lifelong learning, has been studied with the objective of learning new classes continuously without catastrophic forgetting [15] of learned knowledge. Compared with the naive solution of retraining the model from scratch whenever encountering a new class, continual learning is more practical as it only requires the data from new classes during the learning process, leading to the improvements in time, computation and memory efficiency [16]. However, the problem becomes more challenging if the data exhibits long-tailed distribution where the model needs to address both catastrophic forgetting and class-imbalance issues. Though the most recent work [17] aimed to solve this problem by introducing a 2-stage framework, the detached training process also presents challenges when implemented in real-world applications due to the inefficiency caused by manual fine-tuning of each individual stage. In addition, none of the existing work target on food images, which can be more challenging due to the

intra-class variation and inter-class similarity.

Existing continual learning methods show the effectiveness of applying knowledge distillation and storing a small fixed number of seen images as exemplars to mitigate catastrophic forgetting. However, both techniques become less effective in long-tailed distribution. Specifically, the knowledge distillation [18] may even harm the performance when teacher’s model is not trained on balanced data due to the bias in output logits as shown in a recent study [19]. On the other hand, distilling knowledge through learned representations impose a new challenge of feature space misalignment [20] as the learned representation needs to evolve during continual learning to accommodate new classes. Regarding using an exemplar set, most classes in long-tailed distribution may contain only a few training samples. Consequently, the overall performance may still be hindered even when all available samples are stored for instance-rare classes due to the poor generalization ability.

In this work, we focus on designing an end-to-end long-tailed continual learning framework for visual food recognition. We leverage feature-based knowledge distillation while incorporating an additional prediction head that projects the current representation space to the past, which addresses the misalignment issue by providing more freedom to the student model and encourages the retention of the learned knowledge. In addition, inspired by the most recent work [21] that uses the context-rich information in head classes to help the tail classes, we introduce a new data augmentation technique by integrating class-activation-map (CAM) and CutMix [22], which cuts the most important region calculated by CAM in instance-rare classes data as foreground and pastes into the instance-rich classes images. With minimal computational overhead, this method significantly enhances the generalization capabilities of tail classes. We evaluate our method on existing long-tailed food image datasets including Food101-LT and VFN-LT. We further expand the VFN-LT to include two additional population groups, namely Insulin Takers and those with Type 2 diabetes without taking insulin, denoted as VFN-INSULIN and VFN-T2D, respectively. Our proposed framework achieves the best performance with a large improvement margin compared to existing methods while not requiring detached training stages. Finally, we conduct an ablation study to evaluate the effectiveness of each component in our proposed method and discuss potential techniques that can boost the accuracy for implementing in real world applications. The main contributions of this work are summarized as in the following:

- We study the long-tailed continual learning in this work, which is closely related to food recognition in various real world scenarios. Despite its significance, this area of research is currently under-explored.
- Two new benchmark long-tailed food image datasets including VFN-INSULIN and VFN-T2D are introduced, which exhibit the real world food consumption for the population of insulin takers and type 2 diabetes without taking insulin, respectively.
- We propose a novel framework to address both catastrophic forgetting and class-imbalance effectively by using feature-based knowledge distillation with a prediction

head and a novel CAM-based CutMix for data augmentation. In addition, an integrated loss is introduced to strike the balance between learning new classes and maintaining the learned knowledge.

- We conduct a set of comprehensive experiments on all long-tailed continual learning benchmarks for food recognition and discuss potential techniques that could boost the accuracy for facilitating the deployment in real world food recognition.

II. RELATED WORK

In this section, we summarize existing methods that are most related to our work including the food recognition, long-tailed recognition and continual learning.

A. Food Recognition

Food image recognition is a challenging yet practical task that is related to various real world applications such as the image-based dietary assessment [23]–[25]. The performance of nutritional content analysis such as energy [26] or macronutrients are heavily reliant on the accuracy of recognizing and predicting the correct food types from food images. Most existing deep learning based work leverage off-the-shelf models [11], [27]–[29] and train on static food image datasets [9], [12], [13], [30]–[32]. In order to address the issue of inter-class similarity and intra-class variability, a manually built hierarchy based food recognition is proposed in [10]. Later, Mao *et al.* [9] propose to construct a food hierarchy based on visual similarity and nutrition content [33] without requiring human efforts. In addition, food recognition has been studied under different scenarios such as ingredient recognition [34], [35], fine-grained recognition [36]–[38], few-shot learning [39], [40], long-tailed recognition [14], [41] and continual learning [42], [43]. However, none of the existing method are capable of learning new classes continuously in long-tailed distributions, which closely relates to real-world food recognition as new foods appear sequentially overtime with a minority of foods consumed more frequently than the others [14]. Though the most recent work [17] aims to integrate continual learning and long-tailed recognition, they decouple the training process to multi-stages and do not focus on food images. In this work, we target on long-tailed continual learning for visual food recognition to fill this gap and introduce a novel end-to-end framework to address both class-imbalance issue and catastrophic forgetting simultaneously.

B. Long-tailed Recognition

Image recognition in long-tailed distribution has been widely studied in recent years [44]. Existing work can be mainly categorized into two main groups including *re-weighting* and *re-sampling* based methods. The major challenge of long-tailed recognition is the imbalance training data between instance-rich (head) and instance-rare (tail) classes. The **re-weighting** based methods aim to balance the loss or gradients during the training process. The class-level re-weighting loss assigns weights to each class such as the

Balanced Softmax loss [45], Label-Distribution-Aware Margin loss [46] where the inverse of class frequency [47] is widely used as the class weights. In addition, **re-sampling** based techniques aim to construct a balanced training set by over-sampling the tail classes or under-sampling the head classes. However, the naive over-sampling [48] by simply repeating the tail classes data for balanced training further intensifies the over-fitting issue and the naive under-sampling [49] by randomly removing the training data of head classes results in performance degradation. Therefore, most existing work performs efficient data augmentation techniques to improve the generalization ability of tail classes and achieve better overall recognition performance. Gao *et al.* [41] propose Dynamic Mixup for multi-label long-tailed recognition problem, which dynamically adjusts the selection of images based on the previous training performance and set the label of synthetic image as the union of two images. The most recent work proposed CMO [21], which applies CutMix [22] for data augmentation by cutting the foreground region in tail classes images and paste in head classes background. The center idea of CMO is to leverage the context rich information from the head classes to help the generalization of tail classes. Later, He *et al.* [14] improves the CMO to use visually similar image pairs and allows for multi-image CutMix to achieve improved performance. In spite of the efficiency of CutMix for data augmentation, one of the limitations is that it suffers from loss of semantic information of the original image since the cut region is generated randomly. Inspired by [50], we introduce a novel CAM based CutMix, which is able to seamlessly combine the images without losing semantic information, which will be illustrated in Section IV.

C. Continual Learning

Continual learning, also known as incremental learning or lifelong learning, has been studied under different scenarios such as class-incremental, task-incremental and domain-incremental learning as categorized in [51]. From the perspective of real world application, we primarily focus on class-incremental learning in this work, which aims to learn new classes continuously and perform classification on all classes seen so far during the inference phase, which does not require the task index or use multi-head classifier [52] compared to task-incremental learning. The domain-incremental learning instead does not involve new classes but only the domain shift for learned knowledge. The major challenge of class-incremental learning is catastrophic forgetting [15] where the updated model quickly forgets the learned knowledge due to the unavailability of learned classes data. There are two main groups of existing class-incremental learning methods including *regularization based* and *replay based* methods.

Regularization based methods address forgetting by restricting the change of learned parameters during the learning of new classes. The initial work freeze the learned parameters in fully connected layer [53] or constraining the weights update [54], which also confines the model's ability to learn from new data. Later, Li *et al.* [55] proposed to use knowledge distillation [18] to maintain the learned knowledge by

mimicking the output logits distribution from a teacher model, which is then improved in [56] by adding stronger constraints. In addition, feature-based knowledge distillation for class-incremental learning are applied in [57] by minimizing the discrepancy of learned representations between student and teacher models, which is further developed in [58] integrating the logits and feature-based distillation. However, existing knowledge distillation methods are developed on balanced data and becomes less efficient in long-tailed distribution. As revealed in [19], the knowledge distillation using output logits may even harm the overall performance if the teacher model is not trained on balanced data due to the bias towards instance-rich classes. Furthermore, directly applying feature based distillation may also impose new challenge such as feature space misalignment [20] due to the evolving of feature space when learning new classes especially in a long-tailed scenario where the data distribution may vary a lot for each incremental learning step. We address this problem by adding a prediction head to map the current representation space to the past, enabling more efficient knowledge transfer. **Replay based** methods assume the availability of a memory budget to store a small of number of learned classes data as exemplars to perform knowledge replay during the class-incremental learning. Herding algorithm [59] for exemplar selection based on the class mean vector is firstly proposed in [60], which becomes one of the most popular strategies in existing work [56], [57], [61]–[63]. Besides, He *et al.* [42] further improves herding by applying an additional clustering step to select exemplars from each cluster based on herding, which mitigates the issue of intra-class variability. Despite the effectiveness of exemplar replay to maintain the learned knowledge, existing methods assume the training data is balanced and each class has more training data than the memory budget (*e.g.* 20 exemplars per class). However, the majority of classes in the long-tailed scenario only contain a few training samples, resulting in class-imbalance issues in the exemplar set and may even harm the overall performance when performing the knowledge replay. In this work, we address this issue by constructing a balanced exemplar set by augmenting the tail classes data with the proposed CAM based CutMix, which not only improves the efficiency of knowledge replay to maintain the knowledge but increases the generalization ability on tail classes to achieve better overall performance.

III. LONG-TAILED FOOD DATASETS FOR DIABETES

The most recent work [14] introduced VFN-LT, which is a new long-tailed version of VFN [9] dataset for food recognition where the data distribution exhibits the real world food consumption frequencies [64], [65] for healthy people aged 18 to 65 in U.S.. However, around 34.2 million U.S. individuals (10.5 % of the U.S. population) have diabetes [66], which can cause various health problem such as heart disease, vision loss, and kidney disease. The management of diabetes is significantly influenced by diet, therefore, food recognition for predicting appropriate dietary choices to maintain control of diabetes is essential. Nevertheless, no existing studies have specifically focused on this special population so far. In

TABLE I
EXAMPLE OF MATCHED FOOD CODES IN VFN. (NFS DENOTES NOT FURTHER SPECIFIED.)

Food Type	Food Code	Main food description from FNDDS	Additional food description
Apple	63101000	Apple, raw	apple, NFS
Bagel	51180010	Bagel; flavored bagel; egg bagel; Bagel Thins	
Baked potato	71100100	Potato, baked, NFS	
Boiled egg	31103010	Egg, whole, boiled or poached	soft boiled or hard boiled egg; egg, cooked, no fat or milk added
Chicken wings	24160110	Chicken wing, NFS as to cooking method	drummette without sauce or seasoning
Cookies	53201000	Cookie, NFS	
Fried rice	58150310	Rice, fried, NFS	Chinese rice
Ice cream	13110000	Ice cream, NFS	NFS as to flavor
Pancake	55100005	Pancakes, NFS	hot cakes; flapjacks
Yogurt	11400000	Yogurt, NFS	

this work, we aim to fill this gap by introducing two new benchmark long-tailed datasets including VFN-INSULIN and VFN-T2DM, which are constructed based on Viper FoodNet (VFN) datasets [9] and target on food image recognition for dietary assessment among insulin takers and those with type 2 diabetes without taking insulin, respectively.

The original VFN datasets consists of 74 most frequently consumed foods (exclude drink and beverages) selected based on What We Eat In America (WWEIA)¹. Following the similar process in [14], we first manually match each of the 74 food types in VFN [9] with one 8-digit USDA food codes from the Food and Nutrient Database for Dietary Studies (FNDDS)². Each 8-digit USDA food code represents a specific food item reported and consumed in the food supply . Table I shows examples of matched food codes from FNDDS in the VFN dataset. Next, food codes were labeled with their corresponding consumption frequency as determined by Lin *et al* [64] using the nationally representative dietary data collected through the National Health and Nutrition Examination Survey (NHANES)³ from 2009–2016 among U.S adults aged from 20 to 65. The consumption frequency represents how many times a particular food item was reported on one day in the particular age group among the U.S. population. As such, the frequency shows the prominence of the food in comparison with other foods reported as consumed in the population with 774 and 2,758 participants for insulin takers and type 2 diabetes without taking insulin. Finally, we construct VFN-INSULIN and VFN-T2D by reducing the number of training samples for food classes in original VFN based on the matched consumption frequency. Specifically, given the consumption frequency f_i for the i -th food class, we calculate the number of training samples by (1)

$$s_i = n_i \times \frac{f_i}{f_{max}} \quad (1)$$

where n_i is the original number of training samples in VFN for food class i and s_i denotes how many training images are kept for this class, which are randomly selected among original training data. f_{max} refers to the maximum matched consumption frequency among the 74 food types in VFN datasets.

¹<https://data.nal.usda.gov/dataset/what-we-eat-america-wweia-database>
²<https://www.ars.usda.gov/northeast-area/beltsville-md-bhnrc/beltsville-human-nutrition-research-center/food-surveys-research-group/docs/fndds-download-databases/>
³<https://www.cdc.gov/nchs/nhanes/index.html>

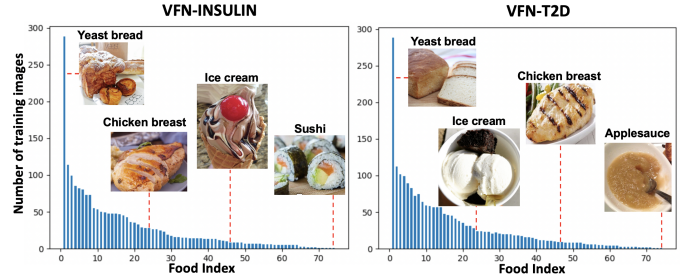


Fig. 2. The distribution of VFN-INSULIN and VFN-T2D shown in descending order based on number of training samples.

Overall, there are 2,056 training images from 74 food classes in VFN-INSULIN with a maximum of 288 training images per class and a minimum of 1 image per class. In addition, there are 2,145 training images in VFN-T2D with a maximum of 288 and a minimum 1 image per class. We observe that the food type *Yeast bread*, has a very high consumption among the represented adults and dominates the consumption frequency for both insulin takers and type 2 diabetes without taking insulin, resulting in the same imbalance ratio $\rho = \frac{\max_i\{s_i\}}{\min_i\{s_i\}}$ calculated as 288 for both datasets. However, the consumption frequencies of other food types may vary between the two groups. Figure 2 shows the distribution for each food type in VFN-INSULIN and VFN-T2D in descending order based on the number of training samples per class.

IV. METHOD

In this work, we introduce a novel end-to-end long-tailed continual learning framework for visual food recognition. The overview of our method is shown in Figure 3. To address catastrophic forgetting, we leverage the teacher model learned from the last incremental step and perform feature-based knowledge distillation with an additional prediction head to enable the efficient knowledge transfer. The exemplar set selected based a novel CAM-based data augmentation for tail classes. Finally, we replace the cross-entropy with the balanced softmax loss [45] based on the current training data distribution to learn class-balance visual representation. In this section, we first introduce the preliminaries for continual learning in long-tailed distribution in Section IV-A and then illustrate the detail of each proposed components in Section IV-B, IV-C and IV-D, respectively.

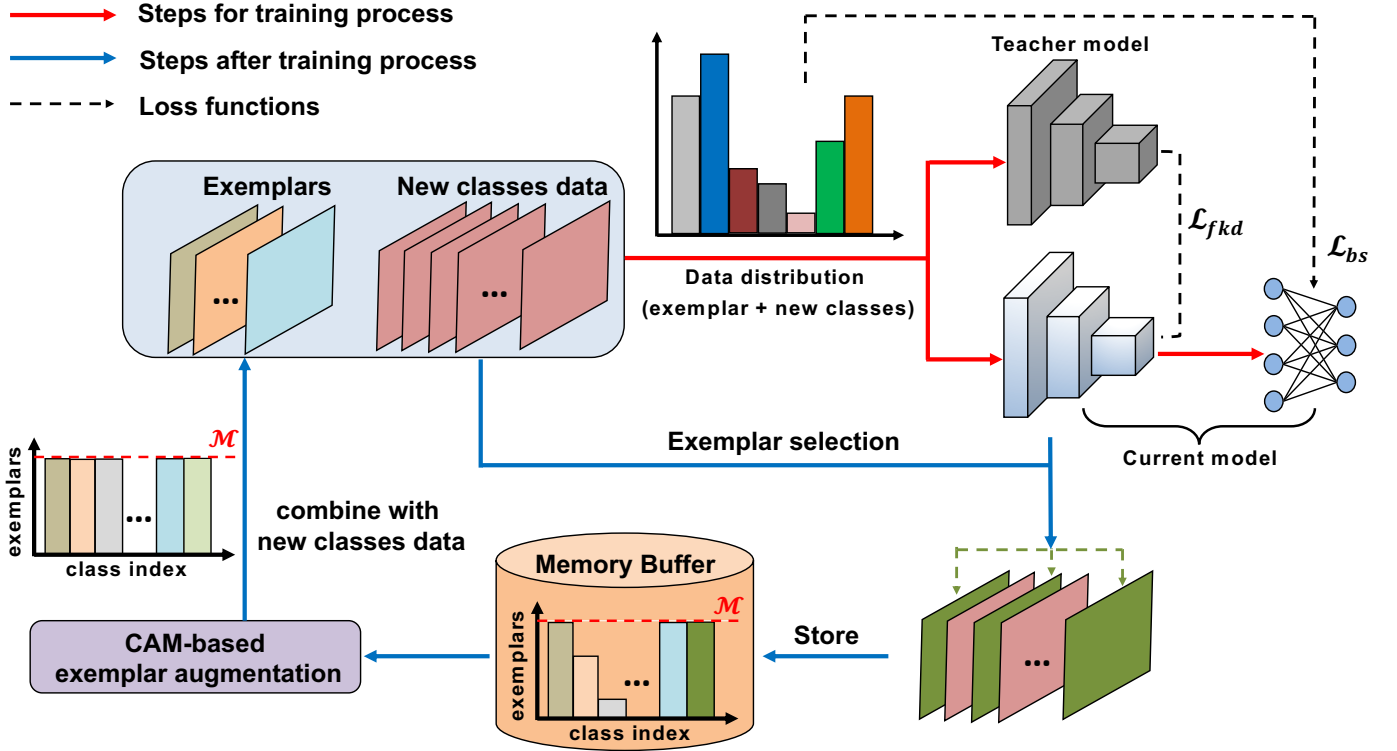


Fig. 3. The overview of our proposed framework. The red arrow shows the training procedure where at each incremental learning phase, the input training data contains the new classes images and the exemplars for all seen classes. We leverage the balanced softmax loss [45] \mathcal{L}_{bs} based on the distribution of input training data to learn class-balance representation. We also employed a fixed teacher model obtained from the last incremental step to perform feature-based knowledge distillation \mathcal{L}_{fkd} to address catastrophic forgetting. The blue arrows denote the steps after the training process where we select a small fixed number of \mathcal{M} exemplars per class and store in the memory buffer. Finally, before the next learning phase, we perform the CAM-based exemplar augmentation to construct a balanced exemplar set and combine with the new classes data for the next training phase.

A. Preliminaries

We focus on continual learning in class-incremental settings where the objective is to learn new classes incrementally and perform classification on all classes seen so far during the inference phase. Specifically, the continual learning in the class-incremental scenario can be formulated as applying an initial model h_0 to learn a sequence of N tasks denoted as $\mathcal{T} = \{\mathcal{T}^1, \dots, \mathcal{T}^N\}$ where each task \mathcal{T}^i contains C_i non-overlapped new classes, which is also known as the incremental step size. During the learning phase of each new task, only the training data $D_i = \{\mathbf{x}_i^j, y_i^j\}$ of the current task is available where \mathbf{x}_i^j and y_i^j denote the j -th input image and label, respectively. After each incremental learning step, the updated model h_i needs to classify $C_{1:i}$ classes encountered so far. The major challenge of continual learning is catastrophic forgetting [15] where the updated model h_i after learning the task \mathcal{T}^i forgets the knowledge of previous tasks $\{\mathcal{T}^1, \dots, \mathcal{T}^{i-1}\}$, resulting in significant performance degradation to classify $C_{1:i-1}$. In the conventional setup, the training data D_i for each task \mathcal{T}^i is balanced distributed, containing $|D_i|/C_i$ samples per class. However, this assumption simplifies the real world complexities especially for food recognition where data is usually long-tailed distributed and exhibits imbalance among food classes. Formally, the training data D_i for each task in long-tailed continual learning is a class-imbalanced distributed with each class containing $(0, |D_i|)$ training samples. The entire training

data D for all the N tasks \mathcal{T} exhibits the long-tail distribution.

1) *Knowledge distillation*: Most existing work [55], [56], [60]–[62] applies knowledge distillation [18] on output logits to maintain the performance on previously learned classes. Specifically, during the learning step of the task \mathcal{T}^i , a teacher model $h_t = h_{i-1}$ learned from the last task with fixed parameters is employed. The knowledge distillation aims to minimize the difference between the output logits of the current model $L = [o^1, o^2, \dots, o^{C_{1:i}}] \in \mathbb{R}^{C_{1:i} \times 1}$ and the outputs of the teacher model $\hat{L} = [\hat{o}^1, \hat{o}^2, \dots, \hat{o}^{C_{1:i-1}}] \in \mathbb{R}^{C_{1:i-1} \times 1}$ by

$$\mathcal{L}_{kd} = - \sum_{j=1}^{C_{1:i-1}} \hat{L}_T^{(j)} \log(L_T^{(j)}) \quad (2)$$

where T is the temperature scalar to learn the hidden knowledge by softening the output distribution as

$$\hat{L}_T^{(j)} = \frac{\exp(\hat{o}^{(j)}/T)}{\sum_{k=1}^{C_{1:i-1}} \exp(\hat{o}^{(k)}/T)} \quad (3)$$

Finally, the knowledge distillation is integrated with cross-entropy during the training process by using a hyper-parameter α to learn new classes as well as maintaining the learned knowledge.

$$\mathcal{L} = \alpha \mathcal{L}_{kd} + (1 - \alpha) \mathcal{L}_{cn} \quad (4)$$

2) *Exemplar replay*: As one of the most commonly used strategies to address catastrophic forgetting, the exemplar replay based methods [57], [60], [61] assume the availability of a reasonable memory budget to select a small fixed number of data as exemplars for each seen class and store them in memory buffer (also known as exemplar set). Specifically, after learning each task \mathcal{T}^i , the lower layers of updated model h_i is used to extract feature embeddings for the new classes training data $D_i = \{\mathbf{x}_i^j, y_i^j\}$. The Herding algorithm [59] is widely applied to select the most representative data for each new class based on euclidean distance between feature embedding and the class mean vector. Therefore, given a memory budget of \mathcal{M} data per class (also known as memory capacity), a subset of $E_i \subseteq D_i$ is selected with $|E_i| = \mathcal{M} \times C_i$ and stored in the memory buffer. Finally, at the beginning of the next new task \mathcal{T}^{i+1} , all the exemplars in memory buffer are combined with the new classes training data to construct $E_i + D_{i+1}$ for continual learning. In this work, we use Herding as the exemplar selection algorithm while others latest work [42], [67] could also be applied.

B. Feature-based Knowledge Distillation

Despite the effectiveness of knowledge distillation in conventional continual learning setup as described in Section IV-A1, it is difficult to apply when new classes exhibit long-tailed distribution since the output logits of the teacher model could be heavily biased towards instance-rich classes due to the severe class-imbalance problem [62], [63]. Directly applying knowledge distillation as in (2) on biased output logits may even harm the overall performance [19]. Therefore, we explore the feature-based knowledge distillation in this work for more efficient knowledge transfer in the long-tailed continual learning scenario. As discussed in [20], one of the challenges when applying feature-based distillation is called the feature space misalignment where the representation of student and teacher models could mismatch in terms of both magnitude and direction. This problem is also relevant in the long-tailed continual learning scenario as the new knowledge is acquired sequentially, which requires the evolution of the continual learning model in the feature space to accommodate the addition of new classes. In this work, we address this challenge by introducing a simple yet effective method as shown in Figure 4. Specifically, instead of directly mimicking the feature from teacher model, we apply an additional predictor g on the head of the continual learning model to map the current representation space to the past in the teacher model. Given an image \mathbf{x} , the predictor g takes the feature representation from the current model (*i.e.* student model) $h_i(\mathbf{x})$ as input and outputs the mapped feature $g(h_i(\mathbf{x}))$. Then we distill the knowledge from the teacher model h_{i-1} by

$$\mathcal{L}_{fkd}(\mathbf{x}) = 1 - \langle g(h_i(\mathbf{x})), h_{i-1}(\mathbf{x}) \rangle \quad (5)$$

where \langle, \rangle measures the cosine similarity. By applying the predictor, we provide the student model with more freedom to accommodate the previous learned representation into the current feature space, enabling more efficient knowledge distillation in long-tailed continual learning. The predictor g

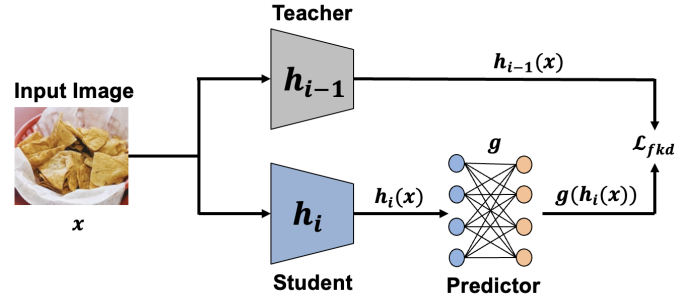


Fig. 4. The overview of proposed feature-based knowledge distillation by applying an additional predictor g .

is be removed after each incremental learning phase. Note that although we apply cosine embedding loss for knowledge distillation, our method can be integrated with other loss functions such as the Mean Squared Error (MSE) loss.

C. CAM-based Exemplar Augmentation

Existing exemplar replay based methods assume each class should contain at least \mathcal{M} images given \mathcal{M} as the memory budget in Section IV-A2. However, most classes in long-tailed distribution may contain only a few training samples $n < \mathcal{M}$, which imposes two new challenges including (i) inefficiency of knowledge replay due to the insufficient of training samples and (ii) intensification of the class-imbalance issue if we directly combine the stored exemplars with training data from new class due to the imbalanced nature of memory buffer. Therefore, we propose a novel data augmentation method in this work to construct a balanced exemplar set by augmenting the tail classes images to address both aforementioned issues. The overview of proposed data augmentation technique is illustrated in Figure 5. To address the issue of losing semantic information when performing data augmentation [14], [21] as described in Section II-B, we propose to use a class activation map (CAM) [68] to identify the most important region from instance-rare classes images and then preserve the semantic information by performing CutMix [22] to cut and paste the identified region into the images with rich context that are selected based on visual similarity. Specifically, we construct class-balanced memory buffer before each new task \mathcal{T}^{i+1} by augmenting stored images for food classes C_t with less than \mathcal{M} exemplars through CutMix in conjunction with images selected from food classes C_h containing \mathcal{M} exemplars. Given an input image $\mathbf{x}_t \in C_t$, we first select the most visually similar candidate $\mathbf{x}_h \in C_h$ by comparing the cosine similarity with h_i as feature extractor where $\mathbf{x}_h = \operatorname{argmax}_{\mathbf{x}_k \in C_h} \langle h_i(\mathbf{x}_t), h_i(\mathbf{x}_k) \rangle$. The lower half of Figure 5 illustrates the procedure to identify the the region to cut and paste into \mathbf{x}_h . Formally, given $\mathbf{x}_t \in \mathbb{R}^{c \times h \times w}$, the class-activation map $M(\mathbf{x}_t) \in \mathbb{R}^{h \times w}$ is calculated by

$$M(\mathbf{x}_t) = \sum_k^d v_{y_{\mathbf{x}_t}}^k h_i(\mathbf{x}_t) \quad (6)$$

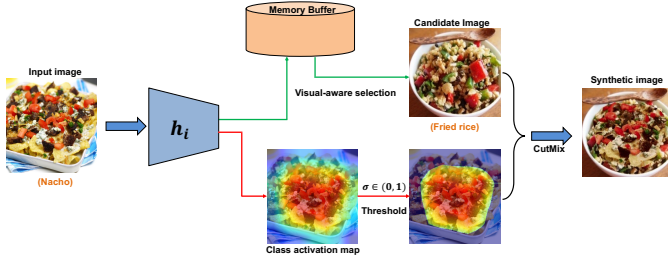


Fig. 5. The overview of proposed CAM-based data augmentation technique. The green arrow describes the selection of the most visually similar candidate image and the red arrow illustrates the steps to obtain the most important region of input image to perform CutMix [22].

where $v_{y_{x_t}} \in \mathbb{R}^d$ refers to the weight vector in classifier of the current model corresponding to the seen class $y_{x_t} \in C_{1:i}$. The value of CAM ranges from $[0, 1]$ and a higher value indicates the more discriminative class-specific region. Therefore, we apply a random threshold $\sigma \in (0, 1)$ to select the region $M(\mathbf{x}_t)^T \in \mathbb{R}^{h \times w}$ where

$$M(\mathbf{x}_t)^T = \begin{cases} M(\mathbf{x}_t) & M(\mathbf{x}_t) \geq \sigma \\ 0 & M(\mathbf{x}_t) < \sigma \end{cases} \quad (7)$$

without losing semantic information of the input image. Finally, we perform CutMix to generate a synthetic image $\tilde{\mathbf{x}}_t$ by

$$\tilde{\mathbf{x}}_t = (1 - S(\mathbf{x}_t)^T) \odot \mathbf{x}_h + S(\mathbf{x}_t)^T \odot \mathbf{x}_t \quad (8)$$

where \odot refers to element-wise multiplication and $S(\mathbf{x}_t)^T$ denotes the binary mask obtained from $M(\mathbf{x}_t)^T$ that $\mathbf{1}$ indicates the region with $M(\mathbf{x}_t)^T > 0$. The class label \tilde{y}_t of the synthetic image is calculated by the area of the replaced region in \mathbf{x}_h as

$$\tilde{y}_t = \frac{1 - A_r}{A} y_h + \frac{A_r}{A} y_t \quad (9)$$

where A_r and A denote the area of the replaced region and the total area of \mathbf{x}_h . y_h and y_t refer to the original class label for \mathbf{x}_h and \mathbf{x}_t , respectively. Note that the exemplar augmentation is performed at the beginning of each new task and the augmented images are not stored in the memory buffer. Note that the Grad-CAM [69], which can be regarded as the generalization of CAM [68], could also be applied in our method.

D. Integrated Loss

While the exemplar augmentation mitigates the class-imbalance issue by constructing a balanced memory buffer, the number of available training data between new classes and the stored classes may still vary a lot during the training phase due to the limited memory budget. Existing work [17] addresses this problem by decoupling the training process into two stages to first learn a feature extractor and then fine-tune the classifier using class-balanced sampler. In this work, we propose to use Balanced Softmax (BS) [45] by extending it into a long-tailed continual learning scenario without requiring a decoupled training process. Specifically, during the training phase of the new task \mathcal{T}^i , a distribution vector $v_d \in \mathbb{R}^{C_{1:i}}$ is

generated by counting the number of training data of each food class for input images in the current task. Recall $L \in \mathbb{R}^{C_{1:i}}$ is the output logits from current model h_i , the distribution vector is then used as the prior information when calculating the loss as shown in (10)

$$\mathcal{L}_{bs} = \sum_{k=1}^{C_{1:i}} -y^k \log[\Phi(\bar{v}_d^k + L^k)] \quad (10)$$

where $\bar{v}_d = v_d / \text{sum}(v_d)$ is the normalized distribution vector and $\Phi(\cdot)$ denotes the *Softmax* function. Therefore, the larger value in the distribution vector achieves smaller gradients when we compute the cross-entropy using the adjusted logits $v_t + L$ and vice versa, which address the class-imbalance issue and enables the end-to-end training pipeline.

The overall training loss function is the weighted sum of feature-based knowledge distillation as described in (5) and the balanced softmax \mathcal{L}_{bs} , which can be expressed as

$$\mathcal{L} = \mathcal{L}_{bs} + \lambda \mathcal{L}_{fkd} \quad (11)$$

where λ is the adaptive ratio to tune the two losses. In this work, as the number of training data D_i may vary a lot for each task \mathcal{T}^i , we propose to calculate $\lambda = \sqrt{|D_i|/|D_{1:i}|}$ as the ratio of training data for the current task and the learned tasks. Therefore, the ratio λ increases when there are more training data from new classes.

V. EXPERIMENT

In this section, we evaluate our proposed long-tailed continual learning framework for visual food recognition as illustrated in Section IV. Specifically, we first introduce the experimental setup including the split of datasets and implementation detail described in Section V-A and V-B, respectively. Then we compare our method with existing work in Section V-C and conduct an ablation study to show the effectiveness of each individual component in Section V-D. Finally, we discuss potential techniques that can boost the performance for real world food related applications in Section V-E.

A. Datasets

We conduct experiments on four benchmark long-tailed food image datasets including: (1) Food101-LT [14], (2) VFN-LT [14], (3) VFN-INSULIN and (4) VFN-T2D where the latter two datasets are introduced in Section III.

Food101-LT is the long-tailed version of Food-101 [12] constructed by applying *Pareto distribution* [71] with the power ratio $\alpha = 6$. We randomly partition the 101 food classes into 5, 10, and 20 tasks for continual learning. Therefore, each task comprises 20, 10, and 5 new classes, respectively, with the exception of the first task, which includes one additional class. The test set is kept as balanced with 125 images per class.

VFN-LT, **VFN-INSULIN** and **VFN-T2D** are all the long-tailed version of VFN [9] but constructed based on the food consumption frequency of different populations including (i) healthy people (ii) insulin takers and (iii) individuals with type 2 diabetes without taking insulin, respectively. We randomly

TABLE II

RESULTS ON FOOD101-LT, VFN-LT, VFN-INSULIN AND VFN-T2D BY COMPARING WITH EXISTING CONTINUAL LEARNING METHODS IN TERMS OF LAST STEP ACCURACY (A_L) AND AVERAGE ACCURACY A_M . BEST RESULTS ARE MARKED IN BOLD.

Datasets	Food101-LT						VFN-LT		VFN-INSULIN		VFN-T2D	
	$N = 5$		$N = 10$		$N = 20$		$N = 7$		$N = 7$		$N = 7$	
Number of tasks	A_L	A_M	A_L	A_M	A_L	A_M	A_L	A_M	A_L	A_M	A_L	A_M
Accuracy (%)	A_L	A_M	A_L	A_M	A_L	A_M	A_L	A_M	A_L	A_M	A_L	A_M
LwF [55]	8.62	10.02	5.86	10.10	0.83	3.58	1.02	4.85	4.69	12.35	3.80	7.10
EWC [54]	4.29	5.05	3.70	7.23	0.83	3.96	1.58	6.31	2.97	12.05	1.50	4.47
iCaRL [60]	11.48	12.42	12.46	13.42	11.04	12.77	11.84	18.76	1.91	4.49	10.64	12.36
EEIL [61]	12.88	12.68	10.57	14.63	6.98	12.66	16.87	20.77	14.77	17.36	13.02	12.95
LwM [70]	10.62	10.82	7.22	12.27	2.45	6.82	7.41	12.32	6.24	15.47	4.24	8.13
IL2M [19]	12.55	11.45	10.97	14.21	6.81	12.16	14.87	18.68	14.19	15.24	10.71	11.93
BiC [62]	14.94	16.72	12.39	16.22	10.38	14.36	15.40	20.89	13.53	19.37	15.32	20.70
LUCIR [57]	13.87	16.94	10.17	15.75	6.74	9.32	16.31	23.37	15.35	20.68	13.42	16.38
PODNet [58]	12.04	12.22	10.46	12.07	8.99	12.41	14.75	18.21	14.18	15.06	12.78	15.55
EEIL-2stage [17]	14.16	14.96	13.29	16.38	9.76	14.71	17.61	22.98	13.64	18.60	12.86	14.47
LUCIR-2stage [17]	16.34	18.90	13.03	18.75	10.85	16.53	17.33	24.26	14.49	22.49	12.82	17.85
PODNet-2stage [17]	13.94	17.89	11.12	17.21	10.28	16.28	17.14	25.58	16.04	22.28	14.24	19.71
Ours	17.54	21.83	13.57	19.25	11.78	17.43	20.43	29.33	18.49	27.51	18.54	28.17

divide the 74 food classes into $N = 7$ tasks with the first task contain 14 new classes and the remaining tasks contain 10 new classes for continual learning. The test set is balanced with 25 images per class. To facilitate an equatable analysis, we use the same testing data in VFN-LT, VFN-INSULIN and VFN-T2D, which is balanced as 25 samples per class with total 1,850 images.

B. Implementation Detail

Our implementation of neural networks are based on Pytorch framework [72] and we apply the ResNet-18 from scratch as the backbone network for experiments for all datasets. The ResNet implementation follows the setting suggested in [11]. We train 90 epochs for each new task with the learning rate starts from 0.1 and decreased with ratio 1/10 for every 30 epochs. The batch size is set as 128 and we apply the stochastic gradient descent (SGD) optimizer with weight decay 0.0001. The memory budget is set as $\mathcal{M} = 20$ to store at most 20 images per food class in the memory buffer.

Evaluation protocol: We use Top-1 classification accuracy as the evaluation metric and evaluate the updated model after learning each new task \mathcal{T}^i on test data belonging to all classes seen so far $\mathcal{C}_{1:i}$. Besides, we report the last step accuracy A_L defined as the classification accuracy on the entire test set after learning all N tasks and the average accuracy A_M calculated by averaging the accuracy obtained after learning each new task, which shows the overall performance for the entire continual learning procedure. Each experiment is run five times and the average performance is presented.

Existing methods: We compare our proposed framework with existing continual learning methods including **LwF** [55], **EWC** [54], **iCaRL** [60], **EEIL** [61], **LwM** [70], **IL2M** [19], **BiC** [62], **LUCIR** [57] and **PODNet** [58]. Also, we apply the 2-stage framework for long-tailed continual learning in [17] and incorporate it into three existing methods including **EEIL-2stage**, **LUCIR-2stage** and **PODNet-2stage** for comparisons.

C. Comparisons With Existing Methods

Table II summarizes the results on Food101-LT, VFN-LT, VFN-INSULIN and VFN-T2D in terms of the last step accuracy A_L and the average accuracy A_M . We observe noticeable improvements for both A_L and A_M of our proposed method by comparing with existing work especially on VFN-LT, VFN-INSULIN and VFN-T2D, which contains a heavier tail with a larger imbalance ratio as introduced in Section III. For example, we achieve around 5% increase of average accuracy A_M and 2% increase of the last step accuracy A_L on VFN-INSULIN compared with the 2-stage framework even without requiring the decoupled training process. Furthermore, we obtain the best performance on Food101-LT with different number of tasks $N \in \{5, 10, 20\}$. We notice both A_L and A_M accuracy will drop as the total number of tasks N increase since we need to address the catastrophic forgetting to maintain the learned knowledge at each learning phase of new tasks except for the first task.

Figure 6 shows the results of top-1 classification accuracy evaluated on all the classes seen so far after learning each new task. Our method achieve promising performance at each learning phase of new task. In addition, different from conventional continual learning where the accuracy usually drops monotonously overtime when new tasks arrive due to the forgetting of previous knowledge, we notice that the performance may even increase after learning a new task in the long-tailed scenario. For example on Food101-LT with $N = 5$, we observe the increase of accuracy for most methods after learning the third task. This occurs because the number of training samples varies significantly among different tasks in long-tailed continual learning where the model gains better knowledge for tasks with a larger number of training images. However, this imposes new challenges of learning from class-imbalanced data for each new task and also hyper-parameter tuning (e.g. the factor of knowledge distillation term as in Equation 4). In this work, we not only address the catastrophic forgetting by introducing a feature-based knowledge distillation and an effective exemplar augmentation technique but also alleviate the class-imbalance by integrating the balanced

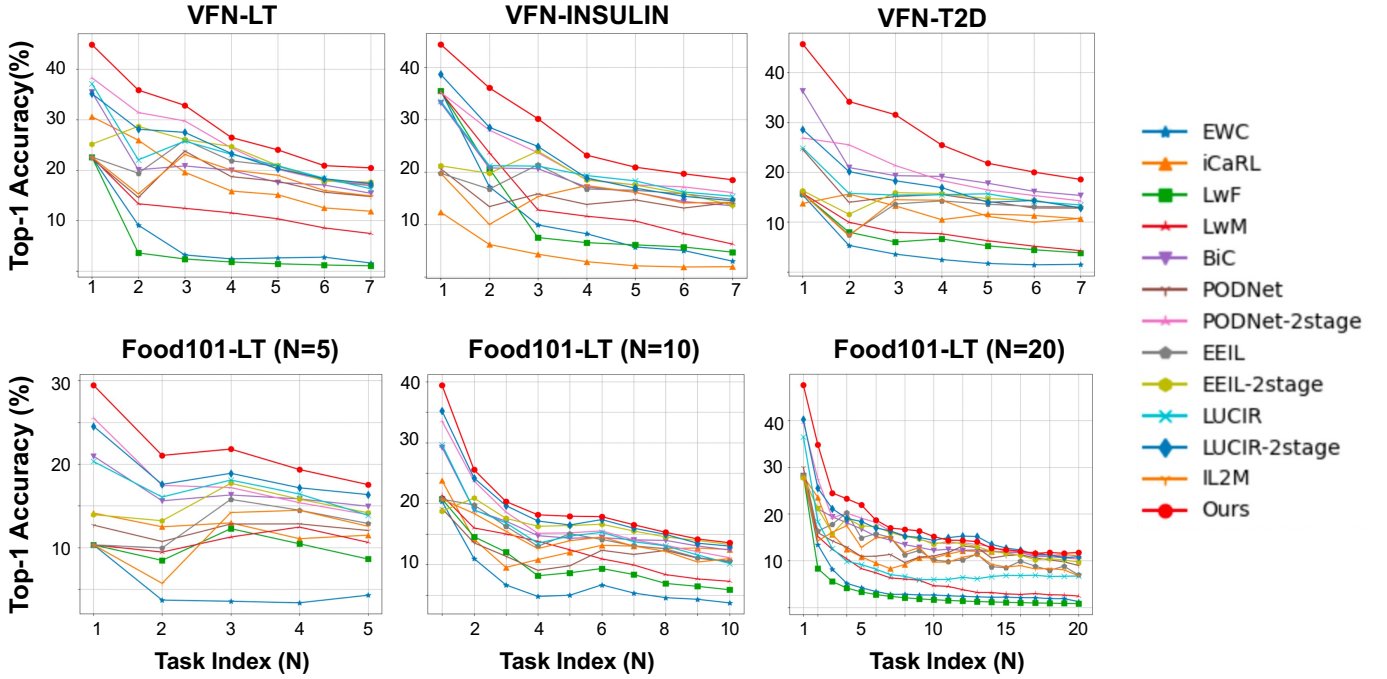


Fig. 6. Results on VFN-LT, VFN-INSULIN, VFN-T2D and Food101-LT with different number of tasks N . Each marker represents the Top-1 classification accuracy evaluated on all classes seen so far after learning each task.

TABLE III
ABLATION STUDY ON FOOD101-LT, VFN-LT, VFN-INSULIN AND VFN-T2D IN TERMS OF AVERAGE ACCURACY A_M .

\mathcal{L}_{fkd}	CAM-CutMix	\mathcal{L}_{bs}	Food101-LT			VFN-LT	VFN-INSULIN	VFN-T2D
			$N = 5$	$N = 10$	$N = 20$	$N = 7$	$N = 7$	$N = 7$
			5.90	8.79	10.55	12.21	11.47	11.93
✓			17.42	15.83	13.96	22.53	21.19	20.73
	✓		13.27	12.99	11.64	16.73	15.41	14.93
		✓	16.52	14.20	12.02	22.96	21.33	21.05
✓	✓		19.31	17.26	15.49	24.18	22.31	22.19
✓	✓	✓	21.83	19.25	17.43	29.33	27.51	28.17

softmax with an adaptive ratio to adjust the impact of each loss term as illustrated in Equation 11 and strike a balance of stability-plasticity.

D. Ablation Study

In this section, we evaluate the effectiveness of each individual component in our proposed framework including (i) the feature-based knowledge distillation (\mathcal{L}_{fkd}), (ii) the cam-based data augmentation (CAM-CutMix) and (iii) the integration of balanced softmax with adaptive ratio (\mathcal{L}_{bs}). Formally, we consider the *baseline* method as using imbalanced memory buffer ($M = 20$) with cross-entropy loss and integrating each of the aforementioned components to conduct experiments. The results in terms of average accuracy A_M are summarized in Table III. We observe consistent performance improvements compared with *baseline* by adding our proposed techniques. Specifically, the feature-based knowledge distillation \mathcal{L}_{fkd} achieves the largest improvements on Food101-LT dataset, demonstrating that catastrophic forgetting is a crucial issue and the integration with CAM-CutMix is able to achieve higher

accuracy. On the other hand, as VFN-LT, VFN-INSULIN, and VFN-T2D exhibit more severe class-imbalance issue due to higher imbalance ratio, the balanced softmax \mathcal{L}_{bs} term has the most significant impact, resulting in the largest performance improvements. Our proposed framework by integrating all the three components obtains the best classification accuracy on all datasets.

We also evaluate our proposed CAM-CutMix by replacing it with existing data augmentation based methods including the CutMix [22] based approaches: (a) the original CutMix used in CMO [21], (b) Visual-Multi CutMix (VM-CMO) [14], (c) SnapMix [50] and the Mixup [73] based approach: (d) D-Mixup [41]. We conduct experiments on VFN-LT and Food101-LT with $N = 10$ as shown in Table IV. Generally, the CutMix based approaches work better in long-tailed continual learning scenarios than D-Mixup, which is usually applied in multi-label recognition scenarios. In addition, the SnapMix achieve a slightly better performance than CMO and VM-CMO as it also considers the class-activation map (CAM) when generating mixed labels. Our method achieves the best performance as it not only preserves the most important re-

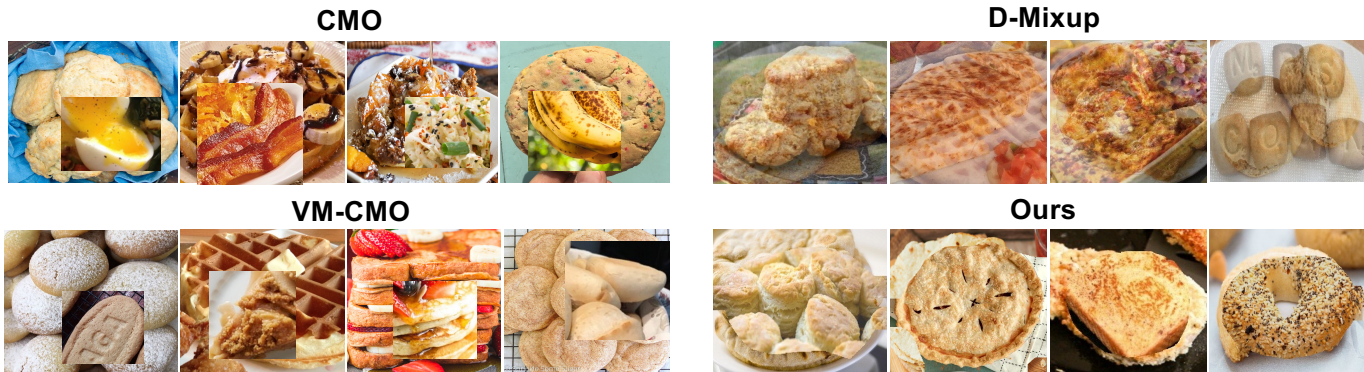


Fig. 7. Examples of augmented food images on VFN-LT using CMO [21], VM-CMO [14], D-Mixup [41] and our proposed CAM-CutMix.

TABLE IV
ABLATION STUDY OF DIFFERENT DATA AUGMENTATION METHODS ON FOOD101-LT ($N = 10$) AND VFN-LT WITH AVERAGE ACCURACY A_M .

	Food101-LT ($N = 10$)	VFN-LT
CMO [21]	17.28	25.93
VM-CMO [14]	16.47	26.41
SnapMix [50]	18.31	27.62
D-Mixup [41]	15.93	25.14
CAM-CutMix (Ours)	19.25	29.33

gions based on CAM but also enables seamless CutMix, rather than relying on a randomly generated bounding box. The example augmented food images using VFN-LT are shown in Figure 7. Note that we do not visualize SnapMix [50] as it has the same synthetic image as in CMO [21] but with a different mixed label.

E. Discussions

Despite the performance improvements our proposed framework demonstrates in comparison to existing methods as shown in Table II, the deployment in real-world applications still remains challenging based on the current classification accuracy. Therefore, in this section we discuss potential techniques that could be applied to boost the performance including (1) increasing the memory buffer capacity to store more exemplar images for knowledge replay and (2) performing transfer learning by pre-training the backbone on large scale image datasets.

1) *Memory buffer capacity*: As one of the most efficient techniques to address catastrophic forgetting, the performance of knowledge replay greatly relies on the capacity of memory buffer (*i.e.* and how many exemplar images can be stored). In this part, we evaluate the long-tailed continual learning performance by varying the memory buffer capacity $\mathcal{M} \in \{10, 20, 30, 40, 50, 100\}$. Figure 8 shows the results in terms of average accuracy A_M on Food101-LT ($N = 10$) and VFN-LT where we observe consistent performance improvements by increasing the memory capacity. However, the memory buffer capacity is a significant constraint for continual learning in real world application as it requires larger memory storage and also poses challenges related to privacy concerns when storing original images as exemplars. Additionally, we observe a performance bottleneck where increasing memory capacity

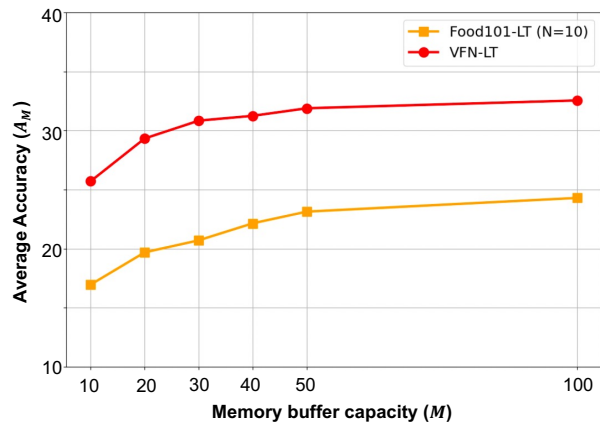


Fig. 8. Average accuracy (A_M) on Food101-LT ($N = 10$) and VFN-LT by varying the memory buffer capacity $\mathcal{M} \in \{10, 20, 30, 40, 50, 100\}$.

does not substantially boost performance. This is predominantly due to dual challenges of catastrophic forgetting and class-imbalance problems that arise in the long-tailed continual learning scenario.

2) *Transfer learning*: Applying the deep models pre-trained on large scale image datasets as backbone a network has been a common strategy to enhance performance across a multitude of vision tasks [13], [74], [75]. In this part, we evaluate the efficacy by leveraging a variety of pre-trained models within the context of long-tailed continual learning for visual food recognition. Formally, we consider different network structures with various depth such as *ResNet-50* [11], *MobileNet* [76], [77], *EfficientNet* [78], [79] *Vision Transformers (ViT)* [80]–[82] and its variants including *DeiT* [81] and *Swin* [82] transformers. In addition, we leverage both ImageNet-1K [83] and ImageNet-21K [84] as the pre-training datasets. ImageNet-1K contains 1,000 classes of general objects, which is the subset of full ImageNet-21K that contains 21,841 classes with over 14,197,122 training images. The results on VFN-LT in terms of average accuracy A_M is shown in Figure 9. We observe significant performance improvements for over 20% by using pre-trained models on large-scale datasets compared to our results in Table II with a model from scratch. This is mainly attributed to the fact that pre-training enhances the feature extraction capabilities of the backbone network,

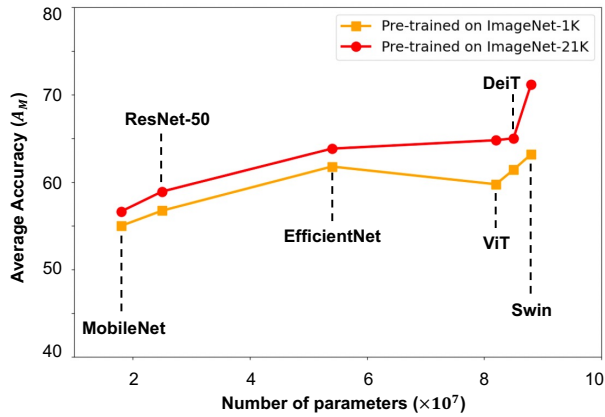


Fig. 9. Average accuracy (A_M) on VFN-LT by leveraging pre-trained models.

thereby enabling it to procure the most discriminative features that are essential for downstream tasks. In addition, pre-training on larger scale datasets that contain more images and classes is able to achieve higher accuracy. However, there is trade-off between the computation complexity and the performance where the increase of model parameters would require longer training time and higher computation capability, which may not be practical for specific real world applications with limited resources available. Note that we intentionally refrain from utilizing food datasets for pre-training in this part to prevent potential overlap with any food class in VFN [9]. However, it's reasonable to anticipate a more substantial performance enhancement if the pre-training were conducted on large-scale food datasets such as Food2K [13].

VI. CONCLUSION

In this work, we focus on visual food recognition within the context of long-tailed continual learning, which is closely related to various real world food based applications. We first introduce two benchmark long-tailed food image datasets including VFN-INSULIN and VFN-T2D, which are generated based on VFN to exhibit the real life food consumption frequency for groups with diabetes. In addition, we propose a novel end-to-end framework that is capable of learning new food classes incrementally in long-tailed data distribution without forgetting the learned knowledge. The proposed framework consists of an effective feature-based knowledge distillation structure by leveraging additional prediction head and a novel data augmentation method based on class-activation-map to address both catastrophic forgetting and class-imbalance issues simultaneously. Our method shows improved performance on all datasets, Food101-LT, VFN-LT, VFN-INSULIN and VFN-T2D, comparing with existing work. Finally, we conduct an ablation study to evaluate the effectiveness of each individual component and discuss the trade-off of potential techniques to achieve higher accuracy including the memory buffer capacity and transfer learning. Our future work includes designing an exemplar-free long-tailed continual learning framework. If successful, this could address the current challenges associated with large memory buffer and potential privacy concerns associated with storing food images.

REFERENCES

- [1] C. J. Boushey, M. Spoden, F. M. Zhu, E. J. Delp, and D. A. Kerr, "New mobile methods for dietary assessment: review of image-assisted and image-based dietary assessment methods," *Proceedings of the Nutrition Society*, vol. 76, no. 3, pp. 283–294, August 2017.
- [2] Fengqing Zhu, Marc Bosch, Insoo Woo, SungYe Kim, Carol J Boushey, David S Ebert, and Edward J Delp, "The use of mobile devices in aiding dietary assessment and evaluation," *IEEE journal of selected topics in signal processing*, vol. 4, no. 4, pp. 756–766, 2010.
- [3] Jiangpeng He, Runyu Mao, Zeman Shao, Janine L Wright, Deborah A Kerr, Carol J Boushey, and Fengqing Zhu, "An end-to-end food image analysis system," *Electronic Imaging*, vol. 2021, no. 8, pp. 285–1, 2021.
- [4] Jiangpeng He, Zeman Shao, Janine Wright, Deborah Kerr, Carol Boushey, and Fengqing Zhu, "Multi-task image-based dietary assessment for food recognition and portion size estimation," *2020 IEEE Conference on Multimedia Information Processing and Retrieval*, pp. 49–54, 2020.
- [5] Weiqing Min, Shuqiang Jiang, Linhu Liu, Yong Rui, and Ramesh Jain, "A survey on food computing," *ACM Computing Surveys (CSUR)*, vol. 52, no. 5, pp. 1–36, 2019.
- [6] Shuqiang Jiang, Weiqing Min, Linhu Liu, and Zhengdong Luo, "Multi-scale multi-view deep feature aggregation for food recognition," *IEEE Transactions on Image Processing*, vol. 29, pp. 265–276, 2020.
- [7] Marc Bolaños and Petia Radeva, "Simultaneous food localization and recognition," *2016 23rd International Conference on Pattern Recognition*, pp. 3140–3145, 2016.
- [8] Jianing Qiu, Frank P.-W. Lo, Yingnan Sun, Siyao Wang, and Benny P. L. Lo, "Mining discriminative food regions for accurate food recognition," *British Machine Vision Conference*, 2019.
- [9] Runyu Mao, Jiangpeng He, Zeman Shao, Sri Kalyan Yarlagadda, and Fengqing Zhu, "Visual aware hierarchy based food recognition," *Proceedings of the International Conference on Pattern Recognition Workshop*, pp. 571–598, February 2021.
- [10] Hui Wu, Michele Merler, Rosario Uceda-Sosa, and John R Smith, "Learning to make better mistakes: Semantics-aware visual food recognition," *Proceedings of the 24th ACM international conference on Multimedia*, pp. 172–176, 2016.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [12] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool, "Food-101 – mining discriminative components with random forests," *Proceedings of the European Conference on Computer Vision*, 2014.
- [13] Weiqing Min, Zhiling Wang, Yuxin Liu, Mengjiang Luo, Liping Kang, Xiaoming Wei, Xiaolin Wei, and Shuqiang Jiang, "Large scale visual food recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [14] Jiangpeng He, Luotao Lin, Heather Eicher-Miller, and Fengqing Zhu, "Long-tailed food classification," *arXiv preprint arXiv:2210.14748*, 2022.
- [15] Michael McCloskey and Neal J Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," vol. 24, pp. 109–165. Elsevier, 1989.
- [16] Jiangpeng He, *Continual Learning: Towards Image Classification From Sequential Data*, Ph.D. thesis, Purdue University, West Lafayette, IN, August 2022.
- [17] Xialei Liu, Yu-Song Hu, Xu-Sheng Cao, Andrew D Bagdanov, Ke Li, and Ming-Ming Cheng, "Long-tailed class incremental learning," *European Conference on Computer Vision*, pp. 495–512, 2022.
- [18] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean, "Distilling the knowledge in a neural network," *Proceedings of the NIPS Deep Learning and Representation Learning Workshop*, 2015.
- [19] Eden Belouadah and Adrian Popescu, "IL2m: Class incremental learning with dual memory," *Proceedings of the IEEE International Conference on Computer Vision*, pp. 583–592, 2019.
- [20] Guo-Hua Wang, Yifan Ge, and Jianxin Wu, "Distilling knowledge by mimicking features," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 11, pp. 8183–8195, 2021.
- [21] Seulki Park, Youngkyu Hong, Byeongho Heo, Sangdoon Yun, and Jin Young Choi, "The majority can help the minority: Context-rich minority oversampling for long-tailed classification," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6887–6896, 2022.

- [22] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Jun-suk Choe, and Youngjoon Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6023–6032, 2019.
- [23] Keigo Kitamura, Toshihiko Yamasaki, and Kiyoharu Aizawa, "Food log by analyzing food images," *Proceedings of the 16th ACM international conference on Multimedia*, pp. 999–1000, 2008.
- [24] Zeman Shao, Yue Han, Jiangpeng He, Runyu Mao, Janine Wright, Deborah Kerr, Carol Jo Boushey, and Fengqing Zhu, "An integrated system for mobile image-based dietary assessment," *Proceedings of the 3rd Workshop on AIFood*, p. 19–23, 2021.
- [25] Kei Nakamoto, Kohei Kumazawa, Hiroaki Karasawa, Sosuke Amano, Yoko Yamakata, and Kiyoharu Aizawa, "Foodlog athl: Multimedia food recording platform for dietary guidance and food monitoring," pp. 1–2, 2022.
- [26] S. Fang, Z. Shao, D. A. Kerr, C. J. Boushey, and F. Zhu, "An end-to-end image-based automatic food energy estimation technique based on learned energy distribution images: Protocol and methodology," *Nutrients*, vol. 11, no. 4, pp. 877, 2019.
- [27] Christian Szegegy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich, "Going deeper with convolutions," *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- [28] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [29] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger, "Densely connected convolutional networks," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, July 2017, Honolulu, HI.
- [30] Y. Kawano and K. Yanai, "Automatic expansion of a food image dataset leveraging existing categories with domain adaptation," *Proc. of ECCV Workshop on Transferring and Adapting Source Knowledge in Computer Vision (TASK-CV)*, 2014.
- [31] Xin Wang, D. Kumar, N. Thome, M. Cord, and F. Precioso, "Recipe recognition with large multimodal food dataset," *2015 IEEE International Conference on Multimedia Expo Workshops (ICMEW)*, pp. 1–6, June 2015.
- [32] Jingjing Chen and Chong-Wah Ngo, "Deep-based ingredient recognition for cooking recipe retrieval," *Proceedings of the 24th ACM international conference on Multimedia*, pp. 32–41, 2016.
- [33] Runyu Mao, Jiangpeng He, Luotao Lin, Zeman Shao, Heather A. Eicher-Miller, and Fengqing Zhu, "Improving dietary assessment via integrated hierarchy food classification," *2021 IEEE 23rd International Workshop on Multimedia Signal Processing (MMSP)*, pp. 1–6, 2021.
- [34] Jingjing Chen and Chong-Wah Ngo, "Deep-based ingredient recognition for cooking recipe retrieval," *Proceedings of the 24th ACM international conference on Multimedia*, pp. 32–41, 2016.
- [35] Jingjing Chen, Bin Zhu, Chong-Wah Ngo, Tat-Seng Chua, and Yu-Gang Jiang, "A study of multi-task and region-wise deep learning for food ingredient recognition," *IEEE Transactions on Image Processing*, vol. 30, pp. 1514–1526, 2020.
- [36] Chengxu Liu, Yuanzhi Liang, Yao Xue, Xueming Qian, and Jianlong Fu, "Food and ingredient joint learning for fine-grained recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 6, pp. 2480–2493, 2020.
- [37] Berker Arslan, Sefer Memiş, Elena Battini Sönmez, and Okan Zafer Batur, "Fine-grained food classification methods on the ucf food-100 database," *IEEE Transactions on Artificial Intelligence*, vol. 3, no. 2, pp. 238–243, 2021.
- [38] Javier Ródenas, Bhalaji Nagarajan, Marc Bolaños, and Petia Radeva, "Learning multi-subset of classes for fine-grained food recognition," *Proceedings of the 7th International Workshop on Multimedia Assisted Dietary Management*, pp. 17–26, 2022.
- [39] Shuqiang Jiang, Weiqing Min, Yongqiang Lyu, and Linhu Liu, "Few-shot food recognition via multi-view representation learning," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 16, no. 3, pp. 1–20, 2020.
- [40] Yanqi Wu, Xue Song, and Jingjing Chen, "Few-shot food recognition with pre-trained model," *Proceedings of the 1st International Workshop on Multimedia for Cooking, Eating, and related Applications*, pp. 45–48, 2022.
- [41] Jixiang Gao, Jingjing Chen, Huazhu Fu, and Yu-Gang Jiang, "Dynamic mixup for multi-label long-tailed food ingredient recognition," *IEEE Transactions on Multimedia*, 2022.
- [42] Jiangpeng He and Fengqing Zhu, "Online continual learning for visual food classification," *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pp. 2337–2346, October 2021.
- [43] Ghalib Ahmed Tahir and Chu Kiong Loo, "An open-ended continual learning for food recognition using class recognition extreme learning machines," *IEEE Access*, vol. 8, pp. 82328–82346, 2020.
- [44] Yifan Zhang, Bingyi Kang, Bryan Hooi, Shuicheng Yan, and Jiashi Feng, "Deep long-tailed learning: A survey," *arXiv preprint arXiv:2110.04596*, 2021.
- [45] Jiawei Ren, Cunjun Yu, Xiao Ma, Haiyu Zhao, Shuai Yi, et al., "Balanced meta-softmax for long-tailed visual recognition," *Advances in neural information processing systems*, vol. 33, pp. 4175–4186, 2020.
- [46] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arachiga, and Tengyu Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," *Advances in neural information processing systems*, vol. 32, 2019.
- [47] Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang, "Learning deep representation for imbalanced classification," *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5375–5384, 2016.
- [48] Jason Van Hulse, Taghi M Khoshgoftaar, and Amri Napolitano, "Experimental perspectives on learning from imbalanced data," *Proceedings of the 24th international conference on Machine learning*, pp. 935–942, 2007.
- [49] Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural networks*, vol. 106, pp. 249–259, 2018.
- [50] Shaoli Huang, Xinchao Wang, and Dacheng Tao, "Snapmix: Semantically proportional mixing for augmenting fine-grained data," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, pp. 1628–1636, 2021.
- [51] Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira, "Re-evaluating continual learning scenarios: A categorization and case for strong baselines," *arXiv preprint arXiv:1810.12488*, 2018.
- [52] Davide Maltoni and Vincenzo Lomonaco, "Continuous learning in single-incremental-task scenarios," *Neural Networks*, vol. 116, pp. 56–73, 2019.
- [53] Heechul Jung, Jeongwoo Ju, Minju Jung, and Junmo Kim, "Less-forgetting learning in deep neural networks," *arXiv preprint arXiv:1607.00122*, 2016.
- [54] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al., "Overcoming catastrophic forgetting in neural networks," *The National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [55] Zhizhong Li and Derek Hoiem, "Learning without forgetting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2935–2947, 2017.
- [56] Jiangpeng He, Runyu Mao, Zeman Shao, and Fengqing Zhu, "Incremental learning in online scenario," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 13926–13935, 2020.
- [57] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin, "Learning a unified classifier incrementally via rebalancing," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 831–839, 2019.
- [58] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle, "Podnet: Pooled outputs distillation for small-tasks incremental learning," *Proceedings of the European Conference on Computer Vision*, pp. 86–102, 2020.
- [59] Max Welling, "Herding dynamical weights to learn," *Proceedings of the International Conference on Machine Learning*, pp. 1121–1128, 2009.
- [60] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert, "iCaRL: Incremental classifier and representation learning," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, July 2017.
- [61] Francisco M. Castro, Manuel J. Marin-Jimenez, Nicolas Guil, Cordelia Schmid, and Karteek Alahari, "End-to-end incremental learning," *Proceedings of the European Conference on Computer Vision*, September 2018.
- [62] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu, "Large scale incremental learning," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2019.
- [63] Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia, "Maintaining discrimination and fairness in class incremental learning," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 13208–13217, 2020.

- [64] Luotao Lin, Fengqing Zhu, Edward J Delp, and Heather A Eicher-Miller, "Differences in dietary intake exist among us adults by diabetic status using nhanes 2009–2016," *Nutrients*, vol. 14, no. 16, pp. 3284, 2022.
- [65] Luotao Lin, Fengqing Zhu, Edward Delp, and Heather Eicher-Miller, "The most frequently consumed and the largest energy contributing foods of us insulin takers using nhanes 2009–2016," *Current Developments in Nutrition*, vol. 5, pp. 426–426, 2021.
- [66] Centers for Disease Control and Prevention, "National diabetes statistics report 2020," <https://www.cdc.gov/diabetes/pdfs/data/statistics/national-diabetes-statistics-report.pdf>.
- [67] Yaoyao Liu, Yuting Su, An-An Liu, Bernt Schiele, and Qianru Sun, "Mnemonics training: Multi-class incremental learning without forgetting," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 12245–12254, 2020.
- [68] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba, "Learning deep features for discriminative localization," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2016.
- [69] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 618–626, 2017.
- [70] Prithviraj Dhar, Rajat Vikram Singh, Kuan-Chuan Peng, Ziyang Wu, and Rama Chellappa, "Learning without memorizing," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5138–5146, 2019.
- [71] Stuart A. Klugman, Harry H. Panjer, and Gordon E. Willmot, *Loss Models: From Data to Decisions*, John Wiley & Sons, 4th edition, 2012.
- [72] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer, "Automatic differentiation in PyTorch," *Proceedings of the Advances Neural Information Processing Systems Workshop*, 2017.
- [73] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz, "Mixup: Beyond empirical risk minimization," *International Conference on Learning Representations*, 2018.
- [74] Dan Hendrycks, Kimin Lee, and Mantas Mazeika, "Using pre-training can improve model robustness and uncertainty," *Proceedings of the 36th International Conference on Machine Learning*, vol. 97, pp. 2712–2721, 09–15 Jun 2019.
- [75] Tz-Ying Wu, Gurumurthy Swaminathan, Zhizhong Li, Avinash Ravichandran, Nuno Vasconcelos, Rahul Bhotika, and Stefano Soatto, "Class-incremental learning with strong pre-trained models," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9601–9610, 2022.
- [76] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [77] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam, "Searching for mobilenetv3," *Proceedings of the IEEE/CVF International Conference on Computer Vision*, October 2019.
- [78] Mingxing Tan and Quoc Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," *International conference on machine learning*, pp. 6105–6114, 2019.
- [79] Mingxing Tan and Quoc Le, "Efficientnetv2: Smaller models and faster training," *International conference on machine learning*, pp. 10096–10106, 2021.
- [80] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *International Conference on Learning Representations*, 2021.
- [81] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou, "Training data-efficient image transformers & distillation through attention," *International conference on machine learning*, pp. 10347–10357, 2021.
- [82] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- [83] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [84] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, 2009.