

FoodX-251: A Dataset for Fine-grained Food Classification

Parneet Kaur* Karan Sikka^{∇†} Weijun Wang[‡] Serge Belongie^{II} Ajay Divakaran[∇]

[∇]SRI International, Princeton, NJ

[‡]Google, Los Angeles, CA

^{II}Cornell Tech, New York, NY

Abstract

Food classification is a challenging problem due to the large number of categories, high visual similarity between different foods, as well as the lack of datasets for training state-of-the-art deep models. Solving this problem will require advances in both computer vision models as well as datasets for evaluating these models. In this paper we focus on the second aspect and introduce FoodX-251, a dataset of 251 fine-grained food categories with 158k images collected from the web. We use 118k images as a training set and provide human verified labels for 40k images that can be used for validation and testing. In this work, we outline the procedure of creating this dataset and provide relevant baselines with deep learning models. The FoodX-251 dataset has been used for organizing iFood-2019 challenge¹ in the Fine-Grained Visual Categorization workshop (FGVC6 at CVPR 2019) and is available for download.²

1. Introduction

A massive increase in the use of smartphones has generated interest in developing tools for monitoring food intake and trends [24, 28, 21]. Being able to estimate calorie intake can aid users to modify their food habits and maintain a healthy diet. Current food journaling applications such as Fitbit App [1], MyFitnessPal [3] and My Diet Coach [2] require users to enter their meal information manually. A study of 141 participants in [11] reports that 25% of the participants stopped food journaling because of the effort involved while 16% stopped because they found it to be time consuming. On the other hand, designing a computer vision based solution to measure calories from clicked images would make the process very convenient. Such an algorithm would generally be required to solve several sub-problems



Figure 1. FoodX-251 Dataset. We introduce a new dataset of 251 fine-grained classes with 118k training, 12k validation and 28k test images. Human verified labels are made available for the training and test images. The classes are fine-grained and visually similar, for example, different types of cakes, sandwiches, puddings, soups, and pastas.

– classify, segment and estimate 3D volume of the given food items. Our focus in this work is to provide a dataset to facilitate the first task of classifying food items in still images.

Food classification is a challenging task due to several reasons: large number of food categories that are fine-grained in nature, resulting in high intra-class variability and low inter-class variability (*e.g.*, different varieties of pasta), prevalence of non-rigid objects, and high overlap in food item composition across multiple food dishes. Further, in comparison to standard computer vision problems such as object detection [20] and scene classification [29], the datasets for food classification are limited in both quantity and quality to train and evaluate deep neural networks. In this work we push the current research in food classification by introducing a new dataset of 251 fine-grained classes with 158k images that supersedes prior datasets in number of classes and data samples.

2. Related Work

Earlier works have tried to tackle the issue of limited datasets for food classification by collecting train-

*Part of the work done while an intern at SRI International.

[†]Corresponding author, karan.sikka@sri.com

¹<https://www.kaggle.com/c/ifood-2019-fgvc6>

²https://github.com/karansikka1/iFood_2019



Figure 2. Noise in web data. *Cross-domain noise*: Along with the images of specific food class, web image search also includes images of processed and packaged food items and their ingredients. *Cross-category noise*: An image may have multiple food items but only one label as its ground truth.

Dataset	Classes	Total Images	Source	Food-type
ETHZ Food-101 [7]	101	101,000	foodspotting.com	Misc.
UPMC Food-101 [26]	101	90,840	Web	Misc.
Food50 [16]	50	5000	Web	Misc.
Food85 [15]	85	8500	Web	Misc.
CHO-Diabetes [4]	6	5000	Web	Misc.
Bettadapura et al. [5]	75	4350	Web, smartphone	Misc.
UEC256 [18]	256	at least 100 per class	Web	Japanese
ChineseFoodNet [10]	208	185,628	Web	Chinese
NutriNet dataset [22]	520	225,953	Web	Central European
Food-251	251	158,846	Web	Misc.

Table 1. Datasets for food recognition. In comparison to prior work, the FoodX-251 dataset (1) provides more classes and images than existing datasets and (2) features miscellaneous classes as opposed to a specific cuisine/food type.

ing data using human annotators or crowd-sourcing platforms [13, 8, 18, 28, 21]. Such data curation is expensive and limits the scalability in terms of number of training categories as well as number of training samples per category. Moreover, it is challenging to label images for food classification tasks as they often have co-occurring food items, partially occluded food items, and large variability in scale and viewpoints. Accurate annotation of these images would require bounding boxes, making data curation even more time and cost prohibitive. Thus, it is important to build food datasets with minimal data curation so that they can be scaled to novel categories based on the final application. Our solution is motivated by recent advances in exploiting the knowledge available in web-search engines and using it to collect a large-scale dataset with minimal supervision [17].

Unlike data obtained by human supervision, web data is freely available in abundance but contains different types of noise [9, 27, 25]. Web images collected via search engines may include images of processed and packaged food items as well as ingredients required to prepare the food items as shown in Figure 2. We refer to this noise as cross-domain noise as it is introduced by the bias due to specific search engine and user tags. In addition, the web data may also include images with multiple food items while being labeled for a single food category (cross-category noise). For example, in images labeled as Guacamole, Nachos can be predominant (Figure 2). Further, the web results may also include images not belonging to any particular class.

Table 1 lists prior datasets for food classification. ETHZ Food-101 [7] consists of 101,000 images of 101 categories. The images are downloaded from a photo sharing website for food items (foodspotting.com). The test data was manually cleaned by the authors whereas the training data consists of cross-category noise, *i.e.*, images with multiple food items labeled with a single class. UPMC Food-101 [26] consists of 90,840 images for the same 101 categories as ETHZ Food-101 but the images are downloaded using web search engine. Some other food recognition datasets with fewer food categories [16, 15, 4, 5] are also listed in Table 1. In comparison to these datasets, our dataset consists of more classes (251) and images (158k).

UEC256 [18] consists of 256 categories with bounding box indicating the location of its category label. However, it mostly contains Japanese food items. ChineseFoodNet [10] consists of 185,628 images from 208 categories but is restricted to Chinese food items only. NutriNet dataset [22] contains 225,953 images from 520 food and drink classes but is limited to Central European food items. In comparison to these datasets, our dataset consists of miscellaneous food items from various cuisines.

3. FoodX-251 Dataset

We introduce a new dataset of 251 fine-grained (prepared) food categories with 158k images collected from the web. We provide a training set of 118k images and human verified labels for both the validation set of 12k images and the test set of 28k images. The classes are fine-grained and

visually similar, for example, different types of cakes, sandwiches, puddings, soups and pastas.

3.1. Data Collection

We start with the 101 food categories in Food-101 dataset [7] and extract their sibling categories from WordNet [23, 6]. We first manually filter and remove all non-food or ambiguous classes.³ Since our primary aim is fine-grained food classification task, we also remove general food classes. For example, different types of pastas and cakes are included but “pasta” and “cake” are removed from the list. This gives us 251 food classes.

For each class, we use web image search to download the corresponding images. Due to the nature of images on these search engines, these images often include images of processed and packaged food items and their ingredients resulting in cross-domain noise. We also observe cross-category noise when for a image search with a single food-item, some images that have multiple food items are downloaded (see Figure 2).

We further filter exact as well as near-exact duplicate images from the dataset. We then randomly selected 200 images from each class and have human raters (3 replications) do verification on this set. From the verified set, we randomly select 70% images for testing and 30% for validation. We use all the remaining images as the training set. The human verification step ensures that the validation and the test set are clean of any cross-domain or cross-category noise. Example of categories with large numbers of samples are generally popular food items such as “churro” or “meatball,” while examples of categories with lower numbers of samples are less popular items such “marble cake,” “lobster bisque,” and “steak-tartare” (Figure 3).

3.2. Evaluation Metric

We follow a similar metric to the classification tasks of the ILSVRC [12]. For each image i , an algorithm will produce 3 labels $l_{ij}, j = 1, 2, 3$, and has one ground truth label g_i . The error for that image is:

$$e_i = \min_j d(l_{ij}, g_i), \quad (1)$$

where,

$$d(x, y) = \begin{cases} 0, & \text{if } x = y. \\ 1, & \text{otherwise.} \end{cases} \quad (2)$$

The overall error score for an algorithm is the average error over all N test images:

$$score = \frac{1}{N} \sum_i e_i. \quad (3)$$

³By ambiguous we refer to those food classes where people do not seem to have a visual consensus.

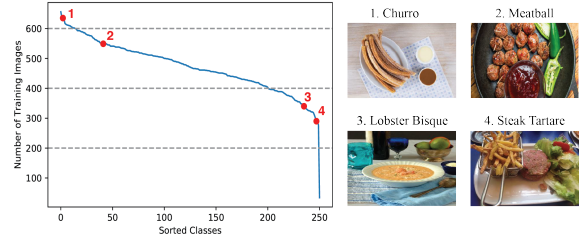


Figure 3. [Left] Training images distribution per class. [Right] Representative images for 4 sampled classes. Food items such as “churro” and “meatball” have large numbers of training images while food items such as “lobster bisque” and “steak-tartare” have relatively fewer training images.

Method	Top-3 Error %		
	Val.	Test	
		Public	Private
ResNet-101 (<i>finetune</i> last-layer)	0.36	0.37	0.37
ResNet-101 (<i>finetune</i> all-layers)	0.16	0.17	0.17

Table 2. Table reports the baseline performance on the FoodX-251 dataset on the validation and the test set.

3.3. Baseline Performance

We implement a naive baseline using a pre-trained ResNet-101 network [14]. We train the model using ADAM optimizer [19] with a learning rate of $5e^{-5}$, which is dropped by a factor of 10 after every 10 epochs. The model is trained for a maximum of 50 epochs with early stopping criteria based on the performance on the validation set. We use random horizontal flips and crops for data augmentation. We use the model checkpoint with the best performance on the validation set for computing test set performance. We have shown results for the validation splits and test splits (as per the Kaggle challenge page) in Table 2.

We observe that ResNet-101 model fine-tuning only the last layer shows a significantly lower performance as compared to the model with fine-tuning all the layers (0.37 vs. 0.17 respectively). We believe that this occurs since the original pre-trained filters are not well suited to the food classification task. As a result, fine-tuning the entire network helps in improving the performance on the fine-grained classification task by a noticeable margin.

4. iFood Challenge at FGVC workshop

The FoodX-211 dataset was used in the iFood-2019 challenge⁴ in Fine-Grained Visual Categorization workshop at CVPR 2019 (FGVC6).⁵ The dataset is also available for

⁴<https://www.kaggle.com/c/ifood-2019-fgvc6>

⁵<https://sites.google.com/view/fgvc6>

download.⁶

This dataset is an extension of FoodX-211 dataset which was used to host iFood-2018 challenge⁷ at FGCV5 (CVPR 2018). FoodX-211 had 211 classes with 101k training images, 10k validation images and 24k test images.

5. Conclusions

In this work, we compiled a new dataset of food images with 251 classes and 158k images. We also provide human-verified labels for 40k images. The baseline results using state-of-the-art ResNet-101 classifier shows 17% top-3 error rate. There is an opportunity for the research community to use more sophisticated approaches on this dataset to further improve the classifier performance. We hope that this dataset will provide an opportunity to develop methods for automated food classification as well as serve as a unique dataset for the computer vision research community to explore fine-grained visual categorization.

6. Acknowledgements

We are thankful to the FGVC workshop organizers for the opportunity to host the iFood competition. We gratefully acknowledge SRI International for providing resources for data collection and Google for providing resources for labeling the data. We are also thankful to Tsung-Yi Lin and CVDF for helping with uploading the data, and also Maggie Demkin, Elizabeth Park, and Wendy Kan from Kaggle for helping us set up the challenge.

References

- [1] Fitbit app. <https://www.fitbit.com/app>. Accessed: 2017-11-14. **1**
- [2] My diet coach. <https://play.google.com/store/apps/details?id=com.dietcoacher.sos>. Accessed: 2017-11-14. **1**
- [3] Myfitnesspal. <https://www.myfitnesspal.com>. Accessed: 2017-11-14. **1**
- [4] Marios Anthimopoulos, Joachim Dehais, Peter Diem, and Stavroula Mougiakakou. Segmentation and recognition of multi-food meal images for carbohydrate counting. In *BIBE*, pages 1–4. IEEE, 2013. **2**
- [5] Vinay Bettadapura, Edison Thomaz, Aman Parnami, Gregory D Abowd, and Irfan Essa. Leveraging context to support automated food recognition in restaurants. In *WACV*, pages 580–587. IEEE, 2015. **2**
- [6] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. ” O’Reilly Media, Inc.”, 2009. **3**
- [7] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *ECCV*, pages 446–461. Springer, 2014. **2, 3**
- [8] Mei-Yun Chen, Yung-Hsiang Yang, Chia-Ju Ho, Shih-Han Wang, Shane-Ming Liu, Eugene Chang, Che-Hua Yeh, and Ming Ouhyoung. Automatic chinese food identification and quantity estimation. In *SIGGRAPH Asia 2012 Technical Briefs*, page 29. ACM, 2012. **2**
- [9] Xinlei Chen and Abhinav Gupta. Webly supervised learning of convolutional networks. In *ICCV*, pages 1431–1439, 2015. **2**
- [10] Xin Chen, Yu Zhu, Hua Zhou, Liang Diao, and Dongyan Wang. ChineseFoodNet: A large-scale image dataset for chinese food recognition. *arXiv preprint arXiv:1705.02743*, 2017. **2**
- [11] Felicia Cordeiro, Elizabeth Bales, Erin Cherry, and James Fogarty. Rethinking the mobile food journal: Exploring opportunities for lightweight photo-based capture. In *HFCS*, pages 3207–3216. ACM, 2015. **1**
- [12] J Deng, A Berg, S Satheesh, H Su, A Khosla, and L Fei-Fei. ILSVRC-2012, 2012. URL <http://www.image-net.org/challenges/LSVRC>, 2012. **3**
- [13] Giovanni Maria Farinella, Dario Allegra, Marco Moltisanti, Filippo Stanco, and Sebastiano Battiato. Retrieval and classification of food images. *Computers in biology and medicine*, 77:23–39, 2016. **2**
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. **3**
- [15] Hajime Hoashi, Taichi Joutou, and Keiji Yanai. Image recognition of 85 food categories by feature fusion. In *ISM*, pages 296–301. IEEE, 2010. **2**
- [16] Taichi Joutou and Keiji Yanai. A food image recognition system with multiple kernel learning. In *ICIP*, pages 285–288. IEEE, 2009. **2**
- [17] Parneet Kaur, Karan Sikka, and Ajay Divakaran. Combining weakly and webly supervised learning for classifying food images. *arXiv preprint arXiv:1712.08730*, 2017. **2**
- [18] Yoshiyuki Kawano and Keiji Yanai. Automatic expansion of a food image dataset leveraging existing categories with domain adaptation. In *ECCV*, pages 3–17, 2014. **2**
- [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. **3**
- [20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014. **1**
- [21] Austin Meyers, Nick Johnston, Vivek Rathod, Anoop Korattikara, Alex Gorban, Nathan Silberman, Sergio Guadarrama, George Papandreou, Jonathan Huang, and Kevin P Murphy. Im2calories: towards an automated mobile vision food diary. In *ICCV*, pages 1233–1241, 2015. **1, 2**
- [22] Simon Mezgec and Barbara Koroušić Seljak. Nutrinet: a deep learning food and drink image recognition system for dietary assessment. *Nutrients*, 9(7):657, 2017. **2**
- [23] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995. **3**
- [24] Manika Puri, Zhiwei Zhu, Qian Yu, Ajay Divakaran, and Harpreet Sawhney. Recognition and volume estimation of

⁶https://github.com/karansikka1/iFood_2019

⁷<https://github.com/karansikka1/Foodx>

- food intake using a mobile device. In *WACV*, pages 1–8. IEEE, 2009. [1](#)
- [25] Sainbayar Sukhbaatar and Rob Fergus. Learning from noisy labels with deep neural networks. *arXiv preprint arXiv:1406.2080*, 2(3):4, 2014. [2](#)
- [26] Xin Wang, Devinder Kumar, Nicolas Thome, Matthieu Cord, and Frederic Precioso. Recipe recognition with large multi-modal food dataset. In *ICMEW*, pages 1–6. IEEE, 2015. [2](#)
- [27] Xin-Jing Wang, Lei Zhang, Xirong Li, and Wei-Ying Ma. Annotating images by mining image search results. *TPAMI*, 30(11):1919–1932, 2008. [2](#)
- [28] Weiyu Zhang, Qian Yu, Behjat Siddiquie, Ajay Divakaran, and Harpreet Sawhney. snap-n-eat food recognition and nutrition estimation on a smartphone. *JDST*, 9(3):525–533, 2015. [1](#), [2](#)
- [29] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *TPAMI*, 2017. [1](#)