# ChineseFoodNet: A Large-scale Image Dataset for Chinese Food Recognition

Xin Chen[†], Yu Zhu[†], Hua Zhou, Liang Diao, and Dongyan Wang*

*Abstract*—In this paper, we introduce a new and challenging large-scale food image dataset called "ChineseFoodNet", which aims to automatically recognizing pictured Chinese dishes. Most of the existing food image datasets collected food images either from recipe pictures or selfie. In our dataset, images of each food category of our dataset consists of not only web recipe and menu pictures but photos taken from real dishes, recipe and menu as well. ChineseFoodNet contains over 180,000 food photos of 208 categories, with each category covering a large variations in presentations of same Chinese food. We present our efforts to build this large-scale image dataset, including food category selection, data collection, and data clean and label, in particular how to use machine learning methods to reduce manual labeling work that is an expensive process. We share a detailed benchmark of several state-of-the-art deep convolutional neural networks (CNNs) on ChineseFoodNet. We further propose a novel two-step data fusion approach referred as "TastyNet", which combines prediction results from different CNNs with voting method. Our proposed approach achieves top-1 accuracies of 81.43% on the validation set and 81.55% on the test set, respectively. The latest dataset is public available for research and can be achieved at https://sites.google.com/view/chinesefoodnet/.

*Index Terms*—dish recognition, deep learning, ChineseFood-Net, TastyNet
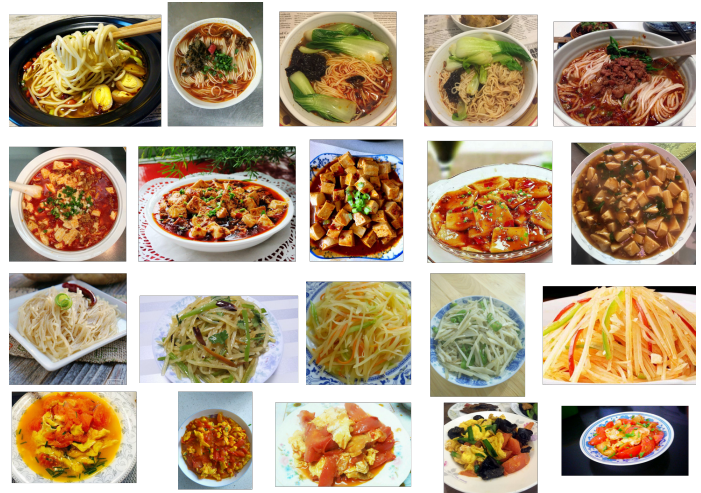


Fig. 1. Example images from our dataset. Each row shows five images from one category of Chinese food. From top to bottom, the food names are Sichuan noodles with peppery sauce, Mapo tofu, potato silk, and scrambled egg with tomato, respectively. Variations in visual appearance of images of Chinese food caused by complex background, various illumination, different angle of view, different ingredients of the same category, etc. show challenges of visual food recognition. All of these image keep their original size.

## I. INTRODUCTION

**F**OOD plays an essential role in everyone's lives, and the behaviour of diet and eating impacts everyone's health [1]. Underestimating food intake directly relates to diverse psychological implications [2]. In recent years, photographing foods and sharing them on social networks have become a part of daily life. Consequently, several applications have been developed to record daily meal activities in personal food log system [3] [4] [5], which are employed to computer-aided dietary assessment [6], further usage preference experiments [7] [8], calorie measurement [9] and nutrition balance estimation [10] [11]. As one of user-friendly ways to input of the food log, automatic recognition of dish pictures gives rise of a research field of interest.

Deep convolutional neural networks (CNNs) have achieved state-of-the-art in a variety of computer vision tasks [12] [13]. The visual dish recognition task is the same situation [14]. The quality of training datasets always plays an important role for

[†] These authors contributed equally to this work.

* means corresponding author.

Xin Chen, Hua Zhou, Yu Zhu and Dongyan Wang are with Midea Emerging Technology Center, San Jose, 95134, USA. Xin's email: chen1.xin@midea.com, Yu's email: zhu.yu@midea.com, Hua's email: hua.zhou@midea.com, and Dongyan's email: dongyan.wang@midea.com.

Liang Diao is with Midea Artificial Intelligence Research Institute, Shenzhen, Guangdong, 528311, P. R. China. email: liang.diao@midea.com.

training a deep neural network, where the high performance of the deep model is still data-driven to some extent [15] [16].

However, to the best of our knowledge, there still exist no effective Chinese food recognition system matured enough to be used in real-world. The major reason is absence of large-scale and high quality image datasets. In [17], the Chinese food dataset includes 50 categories, each of which has only 100 images. Obviously, the size of this dataset is not sufficient to satisfy deep learning training requirements.

The visual dish recognition problem has widely been considered as one of challenging computer vision and pattern recognition tasks [14] [18]. Compared to other types of food such as Italian food and Japanese food, it is more difficult to recognize the images of Chinese dish as the following reasons:

1) The images of same category appear differently. Since most of the same Chinese dish have different ingredients and different cooking methods, the images are greatly visual different, even for human vision;
2) The noise of images of Chinese dishes is hard to model because of complex noise and a variety of backgrounds.The images of Chinese food are taken in various environment and complex background, for example dim light, vapour environment, strong reflection, various utensils of Chinese dishes such as color, shape,

ornament, etc.

In order to give impetus to the progress of visual food classification and related computer vision tasks, we build a large-scale image dataset of Chinese dish, named by **ChineseFoodNet**. This dataset contains 185,628 images of 208 food categories covering most of popular Chinese food, and these images include web images and photos taken in real world under unconstrained conditions. To the best of our knowledge, ChineseFoodNet is the largest and most comprehensive dataset for visual Chinese food recognition. Some of images of ChineseFoodNet are shown in Figure. 2.

We benchmark nine CNNs models of four state-of-the-art deep CNNs, SqueezeNet [19], VGG [20], ResNet [21], and DenseNet [22], on our dataset. Experimental results reveal that ChineseFoodNet is capable of learning complex models.

In this paper, we also propose a novel two-step data fusion approach with voting. Although simple, voting is an effective way to fuse results [23] [24]. Guided by our benchmarks, we try some combination of different CNNs models Based on results on ChineseFoodNet, we take ResNet152, DenseNet121, DeneseNet169, DenseNet201 and VGG19-batch normalization (BN) [25] as our predictive models. [1] Then we fusing these results with voting as a final result. This method is designated as" **TastyNet**". Our proposed method has achieved top-1 accuracy 81.43% in validation set and 81.55% in test set, respectively. Compared to best results of the approaches with a single network structure, the improvements of 2.38% in validation set and 2.33% in these sets have been achieved, respectively.

This paper takes three major contributions as following:

1) We present a large-scale image dataset, ChineseFoodNet, for Chinese food recognition tasks. ChineseFoodNet is made up with 185,628 images of 208 categories, and most of the food image are from users' daily life. It is public available for research in related topics. [2]

2) We provide a benchmark on our dataset. Totally nine different models of four state-of-the-art CNNs architectures are evaluated. We presents the details of the methodology used in the evaluation and the pre-trained models will be public available for further research.

3) We propose a novel two-step data fusion approach for visual food recognition, which combines predictive results of different CNNs with voting. Experimental results on ChineseFoodNet have shown that approach improves the performance compared to one deep CNNs model. It has shown that data fusion should be an alternative way to improve accuracy instead of only increasing numbers of layers in CNNs.

The paper is organized as follows. Section II briefly reviews some public food datasets and the state-of-the-art visual food recognition methods. Section III describes the procedure of building and tagging the ChineseFoodNet dataset. In section IV, several state-of-the-art CNNs methods are benchmarked on ChineseFoodNet. Section V details our proposed data fusion approach and present our results on Chinese-FoofNet. This paper closes with a conclusion of our work and some future directions in section VI.

## II. RELATED WORK

### A. Food Dataset

The scholars have developed some public food datasets[3] for food-related applications such as dietary assessment, computational cooking, food recipe retrieval and so on. Pittsburgh Food Image Dataset (PFID) collects 4,556 fast food images [26]. The UNICT-FD889 dataset of 3,583 images related to 889 distinct dishes are used for Near Duplicate Image retrieval (NDIR) [18]. UEC-Food100 [27] and UEC-Food256 [28] are both Japanese food datasets and contain 100 and 256 categories, respectively. The UPMC-FOOD-101 [29] and ETHZ-FOOD-101 [30] datasets are twin datasets and have same 101 food categories but different images. The images of UPMC-FOOD-101 are recipe images, in which each has the additional textual information, and the images of ETHZ-FOOD-101 are selfies. VIREO-172 [31] is a Chinese Food dataset containing a total of 353 ingredient labels and 110,241 images. However, it aims at cooking recipe retrieval with ingredient recognition.

### B. Visual Dish Recognition

Before introducing deep learning techniques to classification, traditional approaches with hand-crafted features have been applied to visual food recognition, including the pairwise feature distribution (PED) [32], Gabor filters [33], SIHT-based Bag of Visual Words (BoW) [34] [4], optimized bag-of-features model [35], co-occurrence [36], textons [37], Random Forests (RF) [30], and Fisher Vector [38]. Like deep learning applied to other computer vision tasks, CNNs models have outperformed all of traditional methods and achieve higher and higher accuracy with deeper and deeper CNNs [4] [6] [14] [39] [40] [41].

However, all of these approaches of both traditional methods and deep learning haven't been tested on a large-scale image dataset of Chinese food.

## III. CHINESEFOODNET: A LARGE-SCALE CHINESE FOOD IMAGE DATASET

To the best of our knowledge, there is no such large-scale image datasets for Chinese dish recognition which is mature enough to provided necessary resources for the data-driven techniques, e.g. deep learning, to train complex food recognition models. In this section, we present our procedures to build ChineseFoodNet. Labelling image is an expensive step in building large-scale dataset. In this paper, we design and develop a semi-supervised method to accelerate the whole process.

---

[1]The name of CNNs networks consists of letters+numbers. Letters are type of CNNs, and following numbers are the number of layers.

[2]Our dataset can be accessed from https://sites.google.com/view/chinesefoodnet/.

[3]In order to review fairly, we only discuss the data that are available for download in this paper. The last access date is June 1, 2017

Fig. 2. Fifty sample images of ChineseFoodNet dataset. The dataset contains 185,628 Chinese food images organized into 208 categories. All images in the dataset are color. Images are resized for better presentation

### A. Category Selection

Various cooking styles exist in Chinese food culture, such as Sichuan cuisine, Canton cuisine, etc. Our Chinese food dataset must cover the most popular of Chinese cuisines from different styles of cooking. In this subsection, we present our efforts to meet this goal.

First, 250 food categories are gathered from the internet.[4] However, some dishes are missed in search engine yet because they are too popular to be searched such as Tomato omelette. In order to cover them, we conduct a survey of favorite Chinese dishes within our group. Combining with results of the survey, we select about 300 categories. Since Chinese cruise categories is complex and some dishes are very similar visually, such as Braised Chicken Wings and Cola Chicken Wings, we manually merge related categories. After this process, 208 categories of Chinese dish are taken.[5]

### B. Data Collection

There are two resources of our images, web images and taken photos. The web images in our dataset are coming from social network of the Chinese food and drink/cooking,[6] where users uploaded their Chinese food pictures and also provided the tags (labels) of the image. Also some partial of the images in this dataset are collected by our group in daily life.

After these steps, the number of images we brought together achieves more than 500,000. However, those images may contain missing labels, incorrect labels or unclear labels.

### C. Data Clean and Label

After collecting large number of food images, the next step is to clean these data and generate proper labels for each image. In this step, we first remove the images with irregular height or width (too large or too small) which usually are irrelevant images. Then we use entropy to clean the images without content. Entropy is a quantitative metric of image content [42]. We calculate the value of entropy of each channel. If the value of any channel is small, we remove it because the image doesn't have enough useful information. The following step is to remove duplicate and/or very similar images with two steps. First, we calculate 1,024 deep features with the last full connection layer of AlexNet [43]. Second, we calculate the Euclidean distance to measure the similarity. If the distance is below a threshold, we consider the images are very similar and remove one.

Some of these images are clearly categorized with specific Chinese food name, such as most of recipe and menu images. The ground truth of this type of image can be directly extracted. However, the number of such images is very limited and the quality of those images are usually very high, e.g., the images are shot with sufficient light condition, good presentation of the food, and good angels, etc. Thus this type of images shows very different distributions comparing to the images captured in daily life, and brings a potential impact for the food recognition tasks in real life.

The other images are usually not well-labeled, and the food photos are taken in real world conditions. Those images are mainly from the users' daily uploads which show very preferred data distributions in food images in the wild. Besides, this type of images is usually associated with metadata. The metadata can be viewed as an description of each image in text format, which often describes the name, cooking recipe

---

[4] www.top.baidu.com

[5] The names of Chinese dish in ChineseFoodNet are also listed at https://sites.google.com/view/chinesefoodnet/
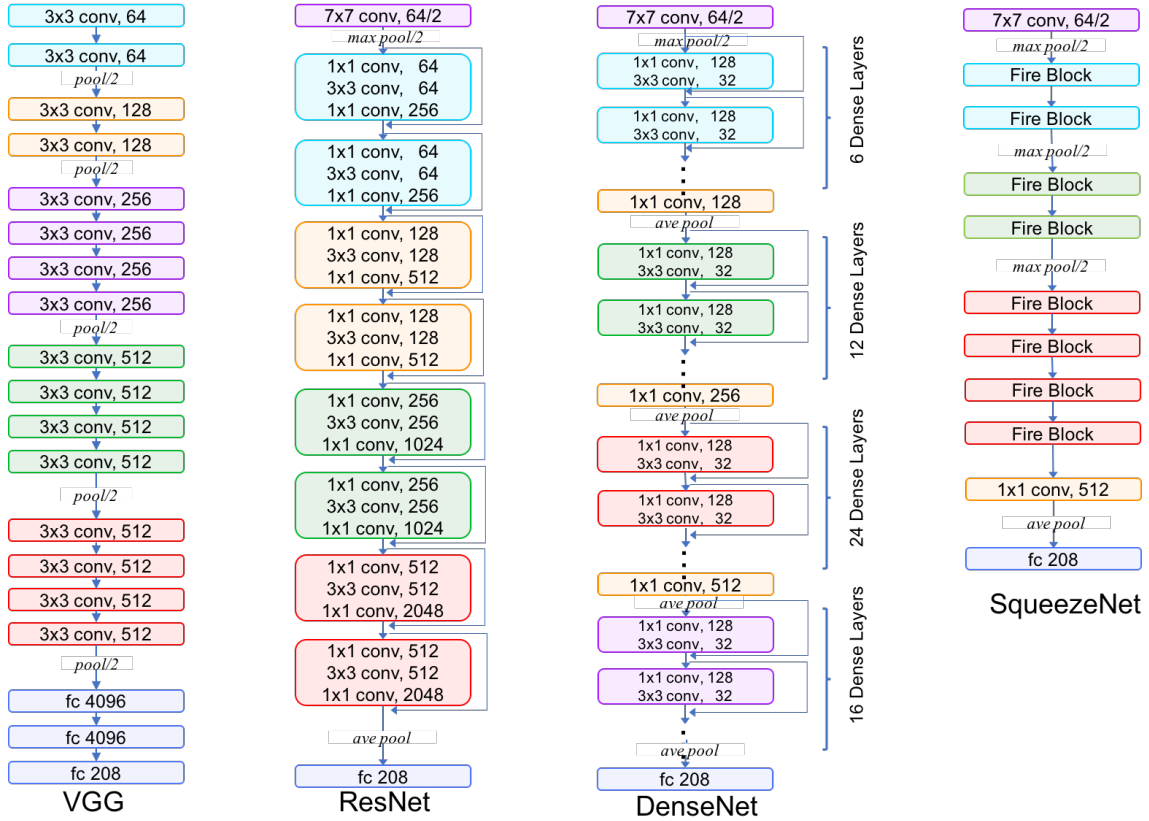
[6] www.douguo.com

Fig. 3. We show basic architectures of four well-known CNNs in our evaluation. From left to right, the architectures are VGG, Resnet, Densenet, and Squeezenet, respectively.

and other information about the food in that image. In our procedure, this metadata is utilized to filter the useful images with correct labels. Particularly, we manually generate a set of keywords for each food class in our database, and use each set of the keyword to match the image metadata. Images with metadata which contains the keywords of certain class are selected and labeled with that class.

It should be noted that, after the aforementioned step, there are still a number of incorrect labels, which are either caused by unclear descriptions in metadata or irrelevant images. Label validation by human labor on this large number of images is an expensive task in terms of both time and costs. Here we accelerate to label these image by some already labeled samples in advance. We first collect a small database of food images using the crowd-sourcing platform (no overlap between our current dataset) with same class labels. Then a shallow CNN model is trained for the food recognition task on this small database. Given this CNNs model, we classify our collected images into different classes representing candidate labels. Specifically, top $n$ (e.g., 5) predictions from the shallow network are selected as the candidate labels for one image. Finally, we perform manually label validation to finalize the dataset by eliminating the wrong labelled images.

### D. Dataset Description

After work of category selection, data collection, the data collection and cleaning mentioned in previous subsections,

finally the ChineseFoodNet dataset contains 185,628 images, with total size of 19.4 Gigabyte (GB). Images in the dataset are kept their original size without any processing and color. The total number of categories is 208 for the current version of dataset, and each image is labelled with only one label from 0 to 207.

We split the whole dataset into training, testing and validation sets, approximately in the ratio 80%, 10% and 10%, respectively. Specifically, there are 145,066, 20,254 and 20,310 images for training, validation and testing set, respectively. Figure 1 and Figure 2 show some example images in our dataset.

## IV. BENCHMARK ON CHINESEFOODNET DATASET

In this section, we conducted benchmark experiments for the ChineseFoodNet dataset. First the experimental settings are described, then we introduce the experimental protocol and finally we provide the experimental results and analysis.

### A. Experimental settings

Our experiments were all conducted using PyTorch [44] deep learning framework. In the training phase, the initial learning rate is set to 0.01, momentum is set 0.9, and weight_decay is set to 1e-4. We set the learning rate to the initial learning rate decayed by 10 every 30 epoch. The number of epoch for the training is set to 90. Training optimization

TABLE I
RECOGNITION RATES OF DIFFERENT DEEP NETWORKS ON OUR FOOD DATASET. BOTH TOP-1 AND TOP-5 ACCURACY ARE SHOWN ON VALIDATION SET AND TEST SET.

| Method | Validation | | Test | |
|--------|------------|------------|------|------|
| | Top-1 Accuracy | Top-5 Accuracy | Top-1 Accuracy | Top-5 Accuracy |
| Squeezenet1_1 | 58.42% | 85.02% | 58.24% | 85.43% |
| VGG19-BN | 78.96% | 95.73% | **79.22%** | **95.99%** |
| ResNet18 | 73.64% | 93.53% | 73.67% | 93.62% |
| ResNet34 | 75.51% | 94.29% | 75.82% | 94.56% |
| ResNet50 | 77.31% | 95.20% | 77.84% | 95.44% |
| ResNet152 | 78.34% | 95.51% | 79.00% | 95.79% |
| DenseNet121 | 78.07% | 95.42% | 78.25% | 95.53% |
| DenseNet169 | 78.87% | **95.80%** | 78.72% | 95.83% |
| DenseNet201 | **79.05%** | 95.79% | 78.78% | 95.72% |

method is selected to stochastic gradient descent (SGD) with momentum. No augmentation process is applied except the resizing and mirror. Training images are firstly resized to 256x256, then a random crop of size 224x224 with hoizontal flip (probability 0.5) is applied. We have used the pretrained models from imagenet dataset [16] and fine-tuned the network with our food data. During the testing, images are resized to 256x256 and then we use center crop of size 224x224 to feed into the network. All the experiments were conducted on CentOS 7 operation system, with Intel Xeon E5 CPU (2.2G), 128GB RAM and Nvidia P100 Tesla GPUs hardware with 16G memory.

### B. Experimental Protocol

The dataset is split into training, validation, and test sets by random selection. There are 145,065 images in the training set. There are 20254 images in the validation set and the rest 20310 images are used for testing. Comprehensive experiments have been conducted using various popular deep learning network architectures with different structures and different number of layers. Specifically, we have benchmarked the performance of: Squeezenet (version 1.1) [19], VGG19 (with BN layer) [20], Resnet (18, 34, and 50) [21], DenseNet (121, 169, and 201) [22]. In order to have a fair comparison, all the experiments are using same input image size and same pre-processing/postprocessing procedures. Some implementation details of ResNet and Squeezenet are illuminated in Figure. 4.
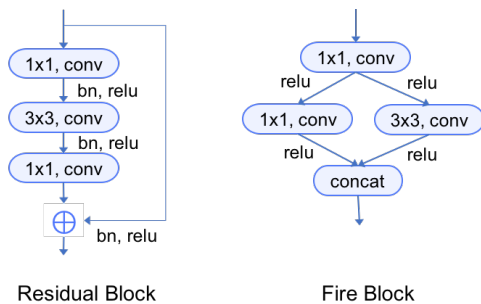


Fig. 4. Illustration of residual block in ResNet and fire block in Squeezenet.

### C. Experimental Results

The recognition performance of different deep networks are shown in Table I, both top 1 accuracy and top 5 accuracy

are presented. Table I has shown that the best top-1 performance on validation set is 79.05%, which is achieved by DenseNet201. The accuracies of VGG19 and DenseNet169 are also very close to the best results. On the test set, the best recognition rate is 79.22%, obtained by VGG19, the second best results is obtained by Resnet152, which is 0.22% lower than then VGG19.

Deeper CNNs models generally achieve better performance [45] [46]. From the results, we can see that, CNN models obtains significant improvements in performance when number of layers in same network architecture are increased. E.g., ResNet with 18 layers has recognition rate 73.64%, while the deeper mode ResNet with 152 layers achieves about 5% improvement in both validation and test sets. Similar results can be observed in DenseNet architecture. On the other hand, deep models with wider structure also shows promising performance, e.g., VGG19-BN obtains the best results in test set, and the worst result (58.42 % and 58.24% on validation and test sets, respectively) is achieved by Squeezenet v1.1, which is a shallow and narrow network structure designed for fast and efficient inference.

## V. TASTYNET: A TWO-STEP DATA FUSION APPROACH

### A. Methodology

As shown in Table I, the accuracy has higher and higher with deeper and deeper model. If we would improve furthermore, a possible way to use much deeper CNNs models. However, it needs much computation and memory resources. What is more, deeper models easily lead to overfitting problem. The alternative way is the data confusion approach. Its idea is to fuse the inference results of different models. As shown in Figure. 5, predictions from different networks are gathered and a voting approach is utilized to obtain the final fused prediction.

Based on some results of different combinations, as shown in Table. II, we select the combination of models that achieves the best top 1 result, ResNet152, DenseNet121, DenseNet169, DenseNet201 and VGG19-BN. The voting method is to average the results of all models. The algorithm is details in Algorithm.!1.

### B. Results and Analysis

Different combinations of network architectures are applied and the experimental results are shown in Table II. From
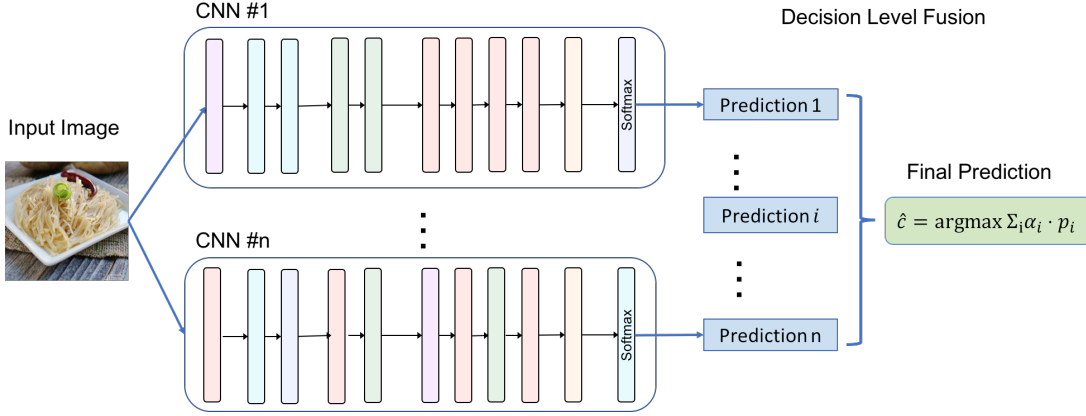
Fig. 5. Basic scheme of the two-step data fusion approach. The first one is to obtain some predictive results from different models. In TastyNet, we use Resnet152, DenseNet121, DenseNet169, DenseNet201 and VGG19-BN. The second one is to combine these result to one final result with voting policy. It his paper, we use weighted coefficient of the results of the first step.

TABLE II
EXPERIMENTAL RESULTS (RECOGNITION ACCURACIES) OF DIFFERENT FUSION SCHEMES

| Fusion Method | Top 1 Accuracy | | Top 5 Accuracy | |
|---|---|---|---|---|
| | Validation | Test | Validation | Test |
| ResNet (18 + 34 + 50 + 152) | 79.19% | 79.46% | 96.03% | 96.16% |
| DenseNet (121 + 169 + 201) | 80.47% | 80.17% | 96.26% | 96.30% |
| ResNet (18 + 34 + 50 + 152) + Densenet (121 + 169 + 201) | 80.89% | 81.08% | 96.60% | 96.67% |
| ResNet (18 + 34 + 50 + 152) + Densenet (121 + 169 + 201) + VGG19-BN | 81.23% | 81.12% | **96.79%** | **96.76%** |
| ResNet152 + DenseNet (121 + 169 + 201) + VGG19-BN | **81.43%** | **81.55%** | 96.73% | **96.76%** |

---

**Algorithm 1** Algorithm of TastyNet.

1: Input:
2: Image
3: Output:
4: Number ▷ Range from 0-207
5: Predictive result from Resnet152, $p(i)$, $i$ from 0 to 207
6: Predictive result from DenseNet121, $p(i)$, $i$ from 0 to 207
7: Predictive result from DenseNet169, $p(i)$, $i$ from 0 to 207
8: Predictive result from DenseNet201, $p(i)$, $i$ from 0 to 207
9: Predictive result from Resnet152, $p(i)$, $i$ from 0 to 207
10: Get average result $\overline{p}(i)$ of all $p(i)$, $i$ from 0 to 207
11: Find maximum $\overline{p}(i)$ and get $i$
12: The output is number $i$

---

this table, we can conclude that the overall performance is generally increasing for different combinations with ensemble more deep networks. The fusion results of ResNet with different number of layers, obtained higher performance (79.46% top 1 accuracy on test set) than single ResNet (ResNet 152, 77.84% top 1 accuracy on test set). Also the fusion results on DenseNet achieved a 1.12% improvement on test set than the best results achieved for single DenseNet architecture. Furthermore, combination of different types of CNNs networks (e.g., ResNets, DenseNets and VGG shown in Row 3 and 4 in Table II) further improves the overall recognition performance. The best result is obtained by fusing Resnet152 and Densenet 121, Densenet 169, Densenet 201, and VGG19-BN, the recognition accuracy is 81.43% on the validation set and 81.55% for the test set. This results is 2.38% and 2.33% higher than the single

network on validation and test set, respectively. Based on the experimental results, we select five CNNs models ,Resnet152, DenseNet121, DenseNet169, DenseNet201 and VGG19-BN, as components of TastyNet.

From our proposed approach, we get two conclusions as followings:

1) By applying data fusing approach on different deep networks, the overall performance can be further boosted than using the single deep network;
2) Combination of different network architectures show more benefits in improving the performance than the combinations with same network architectures, and combination of deeper and wider networks obtains the best results in our evaluation;

## VI. CONCLUSION AND FUTURE WORK

In this paper, we have successfully created a very large-scale image dataset for Chinese dish recognition, ChineseFoodNet. It contains 185,628 images of 208 food categories, in which the images are from not only web images but also real world. As a consequence, the models trained on our dataset should have covered most of food recognition applications. Also, we present the benchmarks of nine state-of-the-art CNNs models of four well-known CNNs architectures on ChineseFoodNet. Finally, we propose a novel two-step data fusion approach, "TastyNet". Based on experimental results, we select Resnet 152, Densenet 121, Densenet 169, Densenet 201 and VGG19+BN models. After voting the results of these model, we obtain final inference result. It has shown the state-of-the-art results on ChineseFoodNet. What is more, our

proposed approach has shown that data fusion is an effective way to obtain a better result instead of only working on one type CNNs model.

For our future work, we are extending the number of food category to over 500 that should be applied in much applications. Also, we will investigate new fusion methods to fuse the different results with different models to obtain the better performance.

### REFERENCES

[1] A. Mesas, M. Muñoz-Pareja, E. López-García, and F. Rodríguez-Artalejo, "Selected eating behaviours and excess body weight: a systematic review," *Obesity Reviews*, vol. 13, no. 2, pp. 106–135, 2012.

[2] M. B. E. Livingstone and A. E. Black, "Markers of the validity of reported energy intake," *The Journal of nutrition*, vol. 133, no. 3, pp. 895S–920S, 2003.

[3] K. Aizawa, "Multimedia foodlog: Diverse applications from self-monitoring to social contributions," *ITE Transactions on Media Technology and Applications*, vol. 1, no. 3, pp. 214–219, 2013.

[4] Y. Kawano and K. Yanai, "Real-time mobile food recognition system," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2013, pp. 1–7.

[5] O. Beijbom, N. Joshi, D. Morris, S. Saponas, and S. Khullar, "Menu-match: restaurant-specific food logging from images," in *Applications of Computer Vision (WACV), 2015 IEEE Winter Conference on*. IEEE, 2015, pp. 844–851.

[6] C. Liu, Y. Cao, Y. Luo, G. Chen, V. Vokkarane, and Y. Ma, "Deepfood: Deep learning-based food image recognition for computer-aided dietary assessment," in *International Conference on Smart Homes and Health Telematics*. Springer, 2016, pp. 37–48.

[7] K. Aizawa, K. Maeda, M. Ogawa, Y. Sato, M. Kasamatsu, K. Waki, and H. Takimoto, "Comparative study of the routine daily usability of foodlog a smartphone-based food recording tool assisted by image retrieval," *Journal of diabetes science and technology*, vol. 8, no. 2, pp. 203–208, 2014.

[8] K. Takahashi, K. Doman, Y. Kawanishi, T. Hirayama, I. Ide, D. Deguchi, and H. Murase, "Estimation of the attractiveness of food photography focusing on main ingredients," in *Proceedings of the 9th Workshop on Multimedia for Cooking and Eating Activities in conjunction with The 2017 International Joint Conference on Artificial Intelligence*. ACM, 2017, pp. 1–6.

[9] P. Pouladzadeh, S. Shirmohammadi, and R. Al-Maghrabi, "Measuring calorie and nutrition from food image," *IEEE Transactions on Instrumentation and Measurement*, vol. 63, no. 8, pp. 1947–1956, 2014.

[10] K. Aizawa, Y. Maruyama, H. Li, and C. Morikawa, "Food balance estimation by using personal dietary tendencies in a multimedia food log," *IEEE Transactions on multimedia*, vol. 15, no. 8, pp. 2176–2185, 2013.

[11] S. Mezgec and B. Koroušić Seljak, "Nutrinet: A deep learning food and drink image recognition system for dietary assessment," *Nutrients*, vol. 9, no. 7, p. 657, 2017.

[12] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[13] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural networks*, vol. 61, pp. 85–117, 2015.

[14] H. Hassannejad, G. Matrella, P. Ciampolini, I. De Munari, M. Mordonini, and S. Cagnoni, "Food image recognition using very deep convolutional networks," in *Proceedings of the 2nd International Workshop on Multimedia Assisted Dietary Management*. ACM, 2016, pp. 41–49.

[15] Y. Bengio *et al.*, "Learning deep architectures for AI," *Foundations and trends® in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.

[16] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009, pp. 248–255.

[17] M.-Y. Chen, Y.-H. Yang, C.-J. Ho, S.-H. Wang, S.-M. Liu, E. Chang, C.-H. Yeh, and M. Ouhyoung, "Automatic Chinese food identification and quantity estimation," in *SIGGRAPH Asia 2012 Technical Briefs*. ACM, 2012.

[18] G. M. Farinella, D. Allegra, and F. Stanco, "A benchmark dataset to study the representation of food images," in *European Conference on Computer Vision*. Springer, 2014, pp. 584–599.

[19] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and $< 0.5$ mb model size," *arXiv preprint arXiv:1602.07360*, 2016.

[20] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

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

[22] F. Iandola, M. Moskewicz, S. Karayev, R. Girshick, T. Darrell, and K. Keutzer, "Densenet: Implementing efficient convnet descriptor pyramids," *arXiv preprint arXiv:1404.1869*, 2014.

[23] X. Chen, L. Lin, and Y. Gao, "Parallel nonparametric binarization for degraded document images," *Neurocomputing*, vol. 189, pp. 43–52, 2016.

[24] C. Macdonald and I. Ounis, "Voting for candidates: adapting data fusion techniques for an expert search task," in *Proceedings of the 15th ACM international conference on Information and knowledge management*. ACM, 2006, pp. 387–396.

[25] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning*, 2015, pp. 448–456.

[26] M. Chen, K. Dhingra, W. Wu, L. Yang, R. Sukthankar, and J. Yang, "PFID: Pittsburgh fast-food image dataset," in *Image Processing (ICIP), 2009 16th IEEE International Conference on*. IEEE, 2009, pp. 289–292.

[27] Y. Matsuda, H. Hoashi, and K. Yanai, "Recognition of multiple-food images by detecting candidate regions," in *Proc. of IEEE International Conference on Multimedia and Expo (ICME)*, 2012.

[28] Y. Kawano and K. Yanai, "Automatic expansion of a food image dataset leveraging existing categories with domain adaptation," in *Proc. of ECCV Workshop on Transferring and Adapting Source Knowledge in Computer Vision (TASK-CV)*, 2014.

[29] X. Wang, D. Kumar, N. Thome, M. Cord, and F. Precioso, "Recipe recognition with large multimodal food dataset," in *Multimedia & Expo Workshops (ICMEW), 2015 IEEE International Conference on*. IEEE, 2015, pp. 1–6.

[30] L. Bossard, M. Guillaumin, and L. Van Gool, "Food-101–mining discriminative components with random forests," in *European Conference on Computer Vision*. Springer, 2014, pp. 446–461.

[31] J. Chen and C.-W. Ngo, "Deep-based ingredient recognition for cooking recipe retrieval," in *Proceedings of the 2016 ACM on Multimedia Conference*. ACM, 2016, pp. 32–41.

[32] S. Yang, M. Chen, D. Pomerleau, and R. Sukthankar, "Food recognition using statistics of pairwise local features," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE, 2010, pp. 2249–2256.

[33] F. Zhu, A. Mariappan, C. J. Boushey, D. Kerr, K. D. Lutes, D. S. Ebert, and E. J. Delp, "Technology-assisted dietary assessment," in *Proceedings of SPIE*, vol. 6814. NIH Public Access, 2008, p. 681411.

[34] F. Kong and J. Tan, "Dietcam: Automatic dietary assessment with mobile camera phones," *Pervasive and Mobile Computing*, vol. 8, no. 1, pp. 147–163, 2012.

[35] M. M. Anthimopoulos, L. Gianola, L. Scarnato, P. Diem, and S. G. Mougiakakou, "A food recognition system for diabetic patients based on an optimized bag-of-features model," *IEEE journal of biomedical and health informatics*, vol. 18, no. 4, pp. 1261–1271, 2014.

[36] Y. Matsuda and K. Yanai, "Multiple-food recognition considering co-occurrence employing manifold ranking," in *Pattern Recognition (ICPR), 2012 21st International Conference on*. IEEE, 2012, pp. 2017–2020.

[37] G. M. Farinella, M. Moltisanti, and S. Battiato, "Classifying food images represented as bag of textons," in *Image Processing (ICIP), 2014 IEEE International Conference on*. IEEE, 2014, pp. 5212–5216.

[38] Y. Kawano and K. Yanai, "Foodcam: A real-time mobile food recognition system employing fisher vector," in *International Conference on Multimedia Modeling*. Springer, 2014, pp. 369–373.

[39] S. Christodoulidis, M. Anthimopoulos, and S. Mougiakakou, "Food recognition for dietary assessment using deep convolutional neural net-

works," in *International Conference on Image Analysis and Processing*. Springer, 2015, pp. 458–465.

[40] N. Martinel, G. L. Foresti, and C. Micheloni, "Wide-slice residual networks for food recognition," *arXiv preprint arXiv:1612.06543*, 2016.

[41] Y. Kawano and K. Yanai, "Food image recognition with deep convolutional features," in *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing: Adjunct Publication*. ACM, 2014, pp. 589–593.

[42] M. Sonka, V. Hlavac, and R. Boyle, *Image processing, analysis, and machine vision*. Cengage Learning, 2014.

[43] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[44] http://pytorch.org/, Last access on August 16, 2017.

[45] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

[46] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826.