

---

# Transfer Learning for Deep Convolutional Neural Network on Food Image Dataset

---

## Abstract

Deep Convolutional Neural Network (CNN) has recently achieved great successes in object recognition competitions, many of which have been accomplished using AlexNet-like architecture. Meanwhile, the explorations of the recent GoogLeNet remain insufficient. Moreover, transfer learning with deep CNN has shown impressive results on benchmark datasets and can be an efficient way for specific recognition tasks. In this paper, we mainly focus on transfer learning on the food images recognition task. The results of our experiments show that incorporating the ImageNet dataset model can lead to improved performance over using only the target labeled dataset. We also conduct empirical studies comparing the performance of AlexNet and GoogLeNet on this task and provide the state-of-the-art performance on our task. We reveal that the small receptive field used in GoogLeNet can improve both training and computational efficiency for transfer learning. We also show in our experiments that deep CNN has strong generalization ability in transfer learning across two food datasets. Our findings may help people better understand the architectures of deep CNN and explore new architectures.

## 1. Introduction

Recognizing the objects is the most fundamental function for human to understand the world, the whole procedure of which only takes a few tens of milliseconds for human brain. Recently, Convolutional Neural Network (CNN) shows its potential to replace the human engineered features, such as SIFT (Lowe, 1999), SURF (Bay et al., 2006) and HOG (Dalal & Triggs, 2005) etc, in the large object recognition tasks. The success of CNNs on large scale image set started from Krizhevsky et al and their 8 layer model AlexNet in 2012 ImageNet Large-Scale Visual

Recognition Challenge (ILSVRC2012), reaching a top-5 accuracy at 83% (Krizhevsky et al., 2012). Since then, many attempts have been made to improve the architecture of AlexNet. By reducing the size of the receptive field and stride, Zeiler and Fergus improve AlexNet by 1.7% on top 5 accuracy (Zeiler & Fergus, 2014). With the help of high performance computing system, such as GPU and large scale distributed cluster, researchers can explore larger and more complex architecture. By both adding extra convolutional layers between two pooling layers and reducing the receptive field size, Simonyan and Zisserman built a 19 layer very deep CNN and achieved 92.5% top-5 accuracy (Simonyan & Zisserman, 2014).

After the AlexNet-like deep CNNs won ILSVRC2012 and ILSVRC2013, Szegedy et al built a 22 layers deep network, called GoogLeNet and won the 1st prize on ILSVRC2014 for 93.33% top-5 accuracy, almost as good as human annotation (Szegedy et al., 2014). Different from AlexNet-like architecture, GoogLeNet shows another trend of design, utilizing many  $1 \times 1$  receptive field. However, its unique architecture has not been widely studied and few results can be found for applying GoogLeNet on other benchmark datasets.

Unlike the local features such as SIFT or SURF, which present an intuitive interpretation of spatial property, we still don't have enough knowledge to fully understand the visual features of CNN learned in each layer. However, the truth that deep CNN outperforms other shallow models by a large margin in some real image recognition tasks encourages researchers to build deeper architecture with powerful high performance hardware and larger datasets. With the help of high performance GPU clusters and data argumentation, people are more enthusiastic to explore bigger networks on complex recognition problems without much interpretation which makes other researchers difficult to understand the advantages of these architectures and explore their own ones.

Since these CNN models are trained on a very large image dataset and can learn hierarchical features, they have strong generalization ability and can be applied in many other scenarios. Applying the pre-trained model from ImageNet dataset on other object recognition benchmark datasets shows some impressive results. Zeiler et al. ap-

plied their pre-trained model on Caltech-256 with just 15 instances per class and improved the previous state-of-the-art in which about 60 instances were used, by almost 10% (Zeiler & Fergus, 2014). Chatfield et al used their pre-trained model on VOC2007 dataset and outperformed the previous state-of-the-art by 0.9% (Chatfield et al., 2014). Few results have been shown related to the transfer learning performance of deep CNN on a more challenging and specific real world recognition problem.

In this paper, we apply two kinds typical of deep CNN architecture, AlexNet and GoogLeNet, on a real specific recognition problem, food recognition, and discuss some observations in fine-tuning the existing CNN architectures on this problem. To our best knowledge, little work has done to discuss the architecture of GoogLeNet while the architecture of AlexNet has been widely studied and improved. Also, little work has been found to compare two deep CNNs with different architectures. By comparing some statistics of the weights and feature maps of these two architectures, we find that the  $1 \times 1$  small receptive field used in GoogLeNet can improve both computational and training efficiency which leads to the success of the whole architecture. Then, we conduct several experiments to stimulate a real world scenario when the training labeled data is rare. The results reveal that deep CNN can work well while transferring knowledge from general recognition task to specific one in this scenario. We achieve 95% of the accuracy trained on full dataset while just utilizing half of the dataset.

The rest of this paper is organized as follow: in Section 2, the two deep CNN architectures and food image datasets are introduced. In Section 3, some experimental results are shown and we also compare the performance between the deep CNNs as well as some traditional methods on these two datasets. And some discussion of the Inception's architecture and statistics are shown in Section 3. We also show some further transfer learning results when the target training examples are rare for each class.

## 2. Experimental Setup

In this section, we introduce some details about the two architectures and the food datasets used in our experiments.

### 2.1. Models

In this paper, AlexNet and GoogLeNet are their Caffe (Jia et al., 2014) implementations and all the results for a specific CNN architecture are obtained from single model.

**AlexNet** contains 5 layers followed by the auxiliary classifier which contains 2 fully connected layers (FC) and 1 softmax layer. Each of the first two layers can be subdivided into 3 components: convolutional layer with rec-

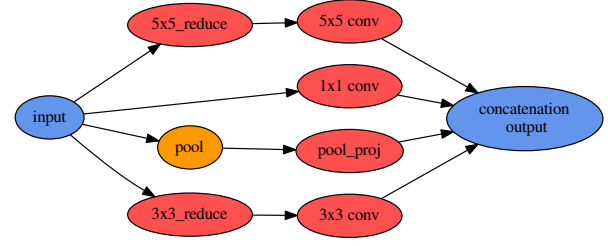


Figure 1. Inception Module.  $n \times n$  stands for size  $n$  receptive field,  $n \times n\_reduce$  stands for the  $1 \times 1$  convolutional layer before the  $n \times n$  convolution layer and  $pool\_proj$  is another  $1 \times 1$  convolutional layer after the MAX pooling layer. The output layer concatenates all its input layers.

tified linear units (ReLUs), local response normalization layer (LRN) and max pooling layer. Layer 3 and layer 4 contain just convolutional layer with ReLUs while layer 5 is similar to the first two layers except for the LRN. For each of the fully connected layer, 1 ReLUs and 1 dropout (Srivastava et al., 2014) layer are followed.

**GoogLeNet** shows another trend of deep CNN architecture with lots of small receptive fields. There are 9 Inception modules in GoogLeNet and Figure 1 shows the architecture of a single inception module. Inspired by (Lin et al., 2013), lots of  $1 \times 1$  convolutional layers are used for computational efficiency. Another interesting feature of GoogLeNet is that there are two extra auxiliary classifiers in intermediate layers. During the training procedure, the loss of these two classifiers are counted into the total loss with a discount weight 0.3, in addition with the loss of the classifier on top. More architecture details can be found from (Szegedy et al., 2014).

### 2.2. Food Datasets

Besides ImageNet dataset, there are many popular benchmark datasets for image recognition such as Caltech dataset and CIFAR dataset, both of which contain hundreds of classes. However, in this paper, we try to focus on a more specific area, food classification. Compared to other recognition tasks, there are some properties of the food (dishes) which make the tasks become a real challenge:

- Food doesn't have any distinctive spatial layout: for other tasks like scene recognition, we can always find some discriminative features such as buildings or trees, etc;
- Food class is a small sub-category among all the categories in daily life, so the inter-class variation is relatively small; on the other hand, the contour of a food

varies depending on many aspects such as the point of the view or even its components.

These properties make food recognition catastrophic for some recognition algorithms. Therefore, the training these two architectures on the food recognition task can reveal some important aspects of themselves and help people better understand them. In this paper, we use two image datasets Food-256 (Kawano & Yanai, 2014a)<sup>1</sup> and Food-101 (Bossard et al., 2014)<sup>2</sup>. It is worthy to mention that PFID dataset is also a big public image database for classification, but their images are collected in a laboratory condition which is considerably not applicable for real recognition task.

**Food-256 Dataset.** This is a relatively small dataset containing 256 kinds of foods and 31644 images from various countries such as French, Italian, US, Chinese, Thai, Vietnamese, Japanese and Indonesia. The distribution among classes is not even and the biggest class (vegetable tempura) contains 731 images while the smallest one contains just 100 images. For this small dataset, we randomly split the data into training and testing set, using around 80% (25361 images) and 20% (6303 images) of the original data respectively and keep the class distribution in these two sets uniform. The collector of this dataset also provides boundary box for each image to separate different foods and our dataset is cropped according to these boundary boxes.

**Food-101 Dataset.** This dataset contains 101-class real-world food (dish) images which were taken and labeled manually. The total number of images is 101,000 and there are exactly 1000 images for each class. Also, each class has been divided into training and testing set containing 750 images and 250 images respectively by its collector. The testing set is well cleaned manually while the training set is not well cleaned on purpose. This noisy training set is more similar to our real recognition situation and it is also a good way to see the effect of the noise on these two architectures.

Data argumentation is an efficient way to enrich the data. There are also some techniques that can applied to enlarge the dataset such as subsampling and mirroring. The original images are firstly resized to  $256 \times 256$  pixels. We crop the 4 corners and center for each image according to the input size of each model and flap the 5 cropped images to obtain 10 crops. For the testing set, the prediction of an image is the average prediction of the 10 crops.

<sup>1</sup>Dataset can be found <http://foodcam.mobi/dataset.html>

<sup>2</sup>Dataset can be found [http://www.vision.ee.ethz.ch/datasets\\_extra/food-](http://www.vision.ee.ethz.ch/datasets_extra/food-101)

### 3. Experimental Discuss

Training a CNN with millions of parameters on a small dataset could easily lead to horrible overfitting. But the idea of supervised pre-training on some huge image datasets could preventing this problem in certain degree. Compared to other randomly initialized strategies with certain distribution, supervised pre-training is to initialize the weights according to the model trained from a specific task. Indeed, initialization with pre-trained model has certain bias as there is no single dataset including all the invariance for natural images (Agrawal et al., 2014), but this bias can be reduced as the pre-trained image dataset increases and the fine-tuning should be benefit from it.

#### 3.1. Pre-training and Fine-tuning

We conduct several experiments on both architectures and use different training initialization strategies for both Food-256 and Food-101 datasets. The scratch models are initialized with Gaussian distribution for AlexNet and Xavier algorithm for GoogLeNet (Glorot & Bengio, 2010). These two initializations are used for training the original models for the ImageNet task. The ft-last and fine-tuned models are initialized with the weights pre-trained from ImageNet dataset. For the ft-last model, we just re-train the fully connected layers while the whole network is fine-tuned for the fine-tune model.

Table 1. Top-5 Accuracy in percent on fine-tuned, ft-last and scratch model for two architectures

	AlexNet		GoogLeNet	
	Food-101	Food-256	Food-101	Food-256
Fine-tune	<b>88.12</b>	<b>85.59</b>	<b>93.51</b>	<b>90.66</b>
Ft-last	76.49	79.26	82.84	83.77
Scratch	78.18	75.35	90.45	81.20

Table 2. Accuracy compared to other method on Food-256 dataset in percent

	fv+linear (Kawano & Yanai, 2014b)	GoogLeNet	AlexNet
Top1	50.1	<b>70.13</b>	63.82
Top5	74.4	<b>90.66</b>	85.59

From Table 1 we can see that fine-tuning the whole network can improve the performance of the CNN for our task. Compared to other traditional computer vision methods (see Table 2 and 3), GoogLeNet outperforms the other methods with large margins and we provide the state-of-the-art performance of these two food image datasets.

In Figure 2 we visualize the feature maps of the pre-trained

Table 3. Top-1 accuracy compared to other methods on Food-101 dataset in percent

	RFDC(Bossard et al., 2014)	MLDS( $\approx$ (Singh et al., 2012))	GoogLeNet	AlexNet
Top1 accuracy	50.76	42.63	<b>78.11</b>	66.40

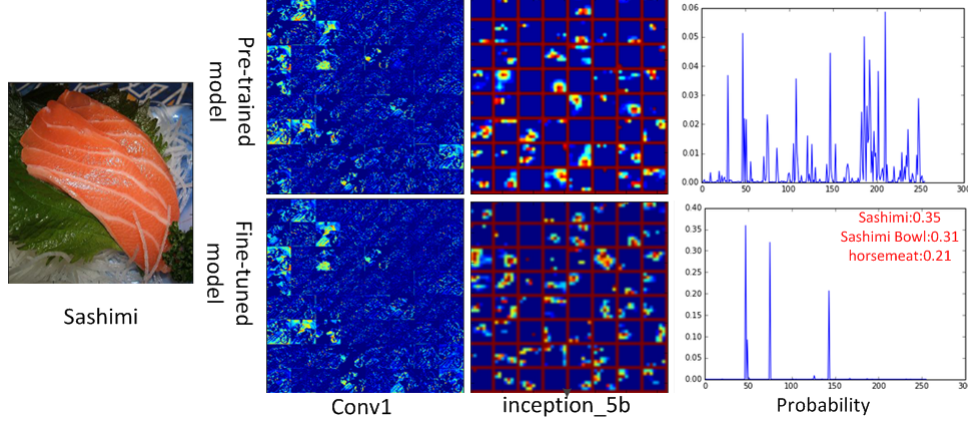


Figure 2. Visualization of some feature maps of different GoogLeNet models in different layers for the same input image. 64 feature maps of each layer are shown. Conv1 is the first convolutional layer and Inception\_5b is the last convolutional layer.

GoogLeNet model and fine-tuned GoogLeNet model with the same input image for some layers. We can see that the feature maps of the lower layer are similar as the lower level features are similar for most recognition tasks. Then we can see that the feature maps in the high-level are different which leads to totally different recognition results. Since only the last layer (auxiliary classifier) of the ft-last model is optimized, we can infer that the higher level features are more important which is consistent with our intuition. Also from Table 1, it is interesting to see that for the Food-101 task, the accuracy of the scratch models outperforms the pre-trained models. Since Food-101 is a relatively large dataset with 750 images per class while Food-256 dataset is an imbalanced small one, this indicates that it is difficult to obtain a good deep CNN model while the data is insufficient.

From Table 1 we can see that GoogLeNet always performs better than AlexNet on both datasets. This implies that the higher level features of GoogLeNet are more discriminative compared to AlexNet and this is due to the special architecture of its basic unit, Inception module. Table 4 and 5 show the weights' cosine similarity of each layer between the fine-tuned models and their pre-trained models. From the results we can see that the weights in the low layer are more similar which implies that these two architectures can learn the hierarchical features. As the low level features are similar for most of the tasks, the difference of the objects is determined by high-level ones which are the combination of these low level features. Also from Table

5, we can observe that, the weights of the pre-trained and fine-tuned models are extremely similar in AlexNet. This can be caused by the size of receptive field. Since ReLUs are used in both architectures, vanishing gradients do not exist. Rectified activation function is mathematically given by:

$$h = \max(w^T x, 0) = \begin{cases} w^T x & w^T x > 0 \\ 0 & \text{else} \end{cases} \quad (1)$$

The ReLU is inactivated when its input is below 0 and its partial derivative is 0 as well. Sparsity can improve the performance of the linear classifier on top, but on the other hand, sparse representations make the network more difficult to train as well as fine-tune. The derivative of the filter is  $\frac{\partial J}{\partial w} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w} = \frac{\partial J}{\partial y} * x$  where  $\frac{\partial J}{\partial y}$  denotes the partial derivative of the activation function,  $y = w^T x$  and  $x$  denotes the inputs of the layer. The sparse input could lead to sparse filter derivative for back propagation which would eventually prevent the errors passing down effectively. Therefore, the filters of the fine-tuned AlexNet is extremely similar. Compared to large receptive field used in AlexNet, the inception module in GoogLeNet employs 2 additional  $n \times n_{reduced}$  convolutional layers before the  $3 \times 3$  and  $5 \times 5$  convolutional layers (see Figure 1). Even though the original purpose of these two  $1 \times 1$  convolutional layer is for computational efficiency, these 2 convolutional layers tend to squeeze their sparse inputs and generate a dense outputs for the following layer. We can see from Table 6 that the sparsity of the  $n \times n_{reduce}$  layers are denser



than other layers within the inception module. This makes the filters in the following layer more easily to be trained for transfer learning and generate efficient sparse representations.

The unique structure of the Inception module guarantees that the sparse outputs from previous layer can be squeezed with the  $1 \times 1$  convolutional layers and feed to convolutional layers with bigger receptive field to generate sparser representation. The squeeze action promises the back propagation error can be transferred more efficiently and makes the whole network more flexible to fit different recognition tasks.

### 3.2. Learning across the datasets

From the previous experiments we can see that pre-training on the ImageNet dataset can improve the performance of the deep convolutional neural network on our specific area. In this part, we will discuss the generalization ability within the food recognition problem. Zhou et al. trained AlexNet for Scene Recognition across two datasets with identical categories (Zhou et al., 2014). But for more complex situation, such as two similar datasets with a little overlapped categories, we are very interested in exploring whether deep CNN can still successfully handle. Therefore, we conduct the following experiment to stimulate a more challenging real world problem: transferring the knowledge from the fine-tuned Food-101 model to a target set, Food-256 dataset. To make the experiment more practical, we limit the number of samples per category from Food-256 for training, because if we want to build a our own model using deep CNN for a specific task, the resource is always limited and it is exhausted to collect hundreds of labeled images for each category.

The Food-101 and Food-256 datasets share about 46 categories of food even though the images in the same category may vary across these two datasets. The types of food in Food-101 are mainly western style while most types of food in Food-256 are typical Asian foods. We compare the top-5 accuracy trained from different size of subset for Food-256 on different pre-trained model and the results are shown in Table 7. The ImageNet columns denote using the model pre-trained from ImageNet dataset as the pre-trained model and Food101\_ft columns denote using the fine-tuned Food-101 model (the same one in Table 1) as the pre-trained model.

From the result of Table 7 we can see that, with this further transfer learning, both CNNs can achieve around 95% of the accuracy trained on full dataset while just utilizing about half of them (50 per class, 12800 of 25361 images). This indicates that when there is not enough labeled data, with its strong generalization ability, deep CNN trained from general task can still achieve satisfying result and per-

form even better when an additional relevant dataset is involved. This encouraging result may attract more people to use deep CNN for their specific task and continue to explore the potential of the existing architecture as well as designing new ones.

## 4. Conclusion

In this paper, we compare two different deep convolutional neural network architectures and their transferring ability on food datasets. Both architectures show their potential on generalization ability and we provide the state-of-the-art on both Food-101 dataset and Food-256 dataset using fine-tuned GoogLeNet. GoogLeNet shows its strong ability on transferring the knowledge between different tasks with the help of the specially designed unit, Inception. We find that not only does the intensively used  $1 \times 1$  convolutional layer in Inception reduce the computational cost, but it also helps to train the filters in the following layer. The filters after the  $1 \times 1$  convolutional layer can be trained more efficiently while the  $1 \times 1$  convolutional layer provides squeezed input and this helps GoogLeNet learn more complex high-level features. Moreover, in a more practical situation such as transfer learning within a specific area with a few labeled instances for the target set, both of these two architectures show some encouraging results in transferring knowledge across the dataset, reaching around 95% of the accuracy trained on full dataset with just half data.

## Acknowledgments

## References

- Agrawal, Pulkit, Girshick, Ross, and Malik, Jitendra. Analyzing the performance of multilayer neural networks for object recognition. In *Computer Vision–ECCV 2014*, pp. 329–344. Springer, 2014.
- Bay, Herbert, Tuytelaars, Tinne, and Van Gool, Luc. Surf: Speeded up robust features. In *Computer vision–ECCV 2006*, pp. 404–417. Springer, 2006.
- Bossard, Lukas, Guillaumin, Matthieu, and Van Gool, Luc. Food-101 – mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.
- Chatfield, K., Simonyan, K., Vedaldi, A., and Zisserman, A. Return of the devil in the details: Delving deep into convolutional nets. In *British Machine Vision Conference*, 2014.
- Dalal, Navneet and Triggs, Bill. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer*

Table 4. Cosine similarity of the layers in inception modules between fine-tuned models and pre-trained model for GoogLeNet

food256						
	1x1	3x3_reduce	3x3	5x5_reduce	5x5	pool_proj
inception_3a	0.72	0.72	0.64	0.67	0.73	0.69
inception_3b	0.59	0.64	0.53	0.70	0.60	0.56
inception_4a	0.46	0.53	0.54	0.50	0.67	0.38
inception_4b	0.55	0.58	0.63	0.52	0.69	0.41
inception_4c	0.63	0.64	0.63	0.57	0.68	0.52
inception_4d	0.60	0.62	0.60	0.58	0.68	0.50
inception_4e	0.60	0.61	0.67	0.61	0.68	0.50
inception_5a	0.51	0.53	0.58	0.48	0.60	0.39
inception_5b	0.40	0.44	0.50	0.41	0.59	0.40

food101						
	1x1	3x3_reduce	3x3	5x5_reduce	5x5	pool_proj
inception_3a	0.71	0.72	0.63	0.67	0.73	0.68
inception_3b	0.56	0.63	0.50	0.71	0.60	0.53
inception_4a	0.43	0.50	0.50	0.47	0.62	0.36
inception_4b	0.48	0.52	0.57	0.50	0.67	0.35
inception_4c	0.57	0.61	0.59	0.53	0.63	0.47
inception_4d	0.54	0.58	0.53	0.54	0.64	0.44
inception_4e	0.53	0.54	0.61	0.55	0.62	0.42
inception_5a	0.43	0.47	0.53	0.45	0.57	0.34
inception_5b	0.36	0.39	0.46	0.38	0.52	0.37

Table 5. Cosine similarity of the layers between fine-tuned models and pre-trained model for AlexNet

	conv1	conv2	conv3	conv4	conv5	fc6	fc7
food256	0.997	0.987	0.976	0.976	0.978	0.936	0.923
food101	0.996	0.984	0.963	0.960	0.963	0.925	0.933

Table 6. Sparsity of the output for each unit in GoogLeNet inception module for training data from Food101 in percent

	1x1	3x3_reduce	3x3	5x5_reduce	5x5	pool_proj
inception_3a	$69.3 \pm 1.3$	$69.6 \pm 1.1$	$80.0 \pm 1.0$	$64.1 \pm 2.2$	$75.8 \pm 1.6$	$76.2 \pm 5.4$
inception_3b	$92.8 \pm 0.9$	$76.5 \pm 0.9$	$94.7 \pm 0.9$	$71.6 \pm 2.3$	$94.4 \pm 0.5$	$94.7 \pm 1.6$
inception_4a	$90.9 \pm 0.9$	$70.0 \pm 1.2$	$93.8 \pm 1.1$	$63.3 \pm 4.0$	$91.9 \pm 1.8$	$95.1 \pm 2.0$
inception_4b	$71.9 \pm 1.6$	$67.5 \pm 1.2$	$75.4 \pm 1.0$	$58.5 \pm 2.6$	$78.9 \pm 1.6$	$85.6 \pm 3.6$
inception_4c	$75.1 \pm 2.4$	$72.6 \pm 1.3$	$81.0 \pm 2.0$	$66.3 \pm 6.1$	$79.7 \pm 3.6$	$88.1 \pm 3.3$
inception_4d	$87.3 \pm 2.7$	$78.0 \pm 2.2$	$88.0 \pm 1.6$	$67.9 \pm 3.1$	$88.9 \pm 2.8$	$93.0 \pm 2.2$
inception_4e	$91.8 \pm 1.1$	$62.3 \pm 2.2$	$91.0 \pm 2.5$	$49.5 \pm 3.7$	$94.0 \pm 1.0$	$92.3 \pm 1.5$
inception_5a	$78.7 \pm 1.6$	$66.5 \pm 1.7$	$82.3 \pm 2.6$	$59.9 \pm 3.2$	$86.4 \pm 2.3$	$87.1 \pm 2.6$
inception_5b	$88.2 \pm 2.3$	$86.8 \pm 1.6$	$83.3 \pm 4.4$	$84.0 \pm 3.1$	$81.4 \pm 5.3$	$94.7 \pm 1.5$

Table 7. Top5 Accuracy for transferring from Food101 to subset of Food256 in percent

instances per class	AlexNet		GoogLeNet	
	ImageNet	Food101_ft	ImageNet	Food101_ft
20	68.80	75.12	74.54	77.77
30	73.15	77.02	79.21	81.06
40	76.04	80.23	81.76	83.52
50	78.90	81.66	84.22	85.84
all	85.59	87.21	90.66	90.65

*Society Conference on*, volume 1, pp. 886–893. IEEE, 2005.

Glorot, Xavier and Bengio, Yoshua. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*, pp. 249–256, 2010.

Jia, Yangqing, Shelhamer, Evan, Donahue, Jeff, Karayev, Sergey, Long, Jonathan, Girshick, Ross, Guadarrama, Sergio, and Darrell, Trevor. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the ACM International Conference on Multimedia*, pp. 675–678. ACM, 2014.

Kawano, Y. and Yanai, K. 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)*, 2014a.

Kawano, Yoshiyuki and Yanai, Keiji. Foodcam-256: A large-scale real-time mobile food recognitionsystem employing high-dimensional features and compression of classifier weights. In *Proceedings of the ACM International Conference on Multimedia*, MM '14, pp. 761–762, New York, NY, USA, 2014b. ACM. ISBN 978-1-4503-3063-3.

Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.

Lin, Min, Chen, Qiang, and Yan, Shuicheng. Network in network. *CoRR*, abs/1312.4400, 2013.

Lowe, David G. Object recognition from local scale-invariant features. In *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, volume 2, pp. 1150–1157. Ieee, 1999.

Simonyan, Karen and Zisserman, Andrew. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.

Singh, Saurabh, Gupta, Abhinav, and Efros, Alexei A. Unsupervised discovery of mid-level discriminative patches. In *Computer Vision–ECCV 2012*, pp. 73–86. Springer, 2012.

Srivastava, Nitish, Hinton, Geoffrey, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1): 1929–1958, 2014.

Szegedy, Christian, Liu, Wei, Jia, Yangqing, Sermanet, Pierre, Reed, Scott, Anguelov, Dragomir, Erhan, Dumitru, Vanhoucke, Vincent, and Rabinovich, Andrew. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014.

Zeiler, Matthew D and Fergus, Rob. Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014*, pp. 818–833. Springer, 2014.

Zhou, Bolei, Lapedriza, Agata, Xiao, Jianxiong, Torralba, Antonio, and Oliva, Aude. Learning deep features for scene recognition using places database. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N.D., and Weinberger, K.Q. (eds.), *Advances in Neural Information Processing Systems 27*, pp. 487–495. Curran Associates, Inc., 2014.