

# Development of a Food Image Recognition Algorithm Using Machine Learning

Serge Assaad (✉ [serge.assaad@duke.edu](mailto:serge.assaad@duke.edu))

Duke University <https://orcid.org/0000-0001-9310-0972>

Lawrence Carin

Duke University

Anthony Joseph Viera

Duke Medicine

---

## Research

**Keywords:** machine learning, diet tracking, nutrition apps

**Posted Date:** March 3rd, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-15859/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

RESEARCH

# Development of a Food Image Recognition Algorithm Using Machine Learning

Serge Assaad<sup>1\*</sup>

, Lawrence Carin<sup>1</sup>

and Anthony J Viera<sup>2</sup>

\*Correspondence:

[serge.assaad@duke.edu](mailto:serge.assaad@duke.edu)

<sup>1</sup>Department of Electrical & Computer Engineering, Duke University, Science Drive, 27708 Durham, NC, USA

Full list of author information is available at the end of the article

## Abstract

**Background:** Researchers and consumers have limited options for objectively collecting or tracking data related to food choices.

**Objective:** To develop and pilot test an algorithm that could accurately categorize food items from a meal photograph.

**Methods:** We used a dataset of 7721 meal photographs taken by patrons in a cafeteria setting. We designed 22 broad categories recognizable by image that are parents of the original 1239 types of items in the photographs. We split the dataset into 3 mutually exclusive subsets: a training set (5250 images), a validation set (1312 images), and a test set (1159 images). Using a convolutional neural network and standard machine learning techniques, we tested the operating characteristics of the algorithm.

**Results:** Salad recognition had the lowest specificity (0.74), while multiple categories had specificities close to 1.0 (e.g. cereals, pastries, sushi, yogurt). Areas under the ROC curve (AUCs), reflecting trade-offs between sensitivity and specificity, ranged from 0.73 (for yogurt) to 0.97 (for sushi).

**Conclusions:** This work provides proof-of-concept for an algorithm that can categorize food items from a meal photograph.

**Keywords:** machine learning; diet tracking; nutrition apps

1

2

### 3 Introduction

4 Obesity and physical inactivity are leading causes of premature death in the United  
5 States and across the globe. The costs of obesity and its health related problems are  
6 also a substantial societal burden. Obese adults have more visits to physicians, and  
7 a higher number of inpatient hospital days per year [1],[2]. Given the public health  
8 and medical consequences of obesity and unhealthy eating, nutrition research is of  
9 critical importance.

10 Methods used to measure nutritional content of foods that people eat include  
11 various forms of food surveys. A food diary (or food/diet record) has people record  
12 all foods and drinks consumed over a specified time period (e.g. consecutive days).  
13 Measuring the items is critical to obtaining reliable estimates. The method does not  
14 rely on recall, but it is burdensome and the measurement process itself might alter  
15 eating behaviors. The Food Frequency Questionnaire [3] asks people to report how  
16 frequently they consumed certain foods and drinks over a specified time period.  
17 It is less burdensome and does not alter eating behaviors, but it relies on recall.  
18 Additionally, it is not as precise in terms of quantifying intake. The Automated Self-  
19 Administered 24-Hour (ASA24®) [4] Dietary Assessment Tool provides a more  
20 quantifiable estimate of intake while keeping burden fairly low. The ASA24 may  
21 capture food intake with less bias than food frequency questionnaires, but it still  
22 relies on recall.

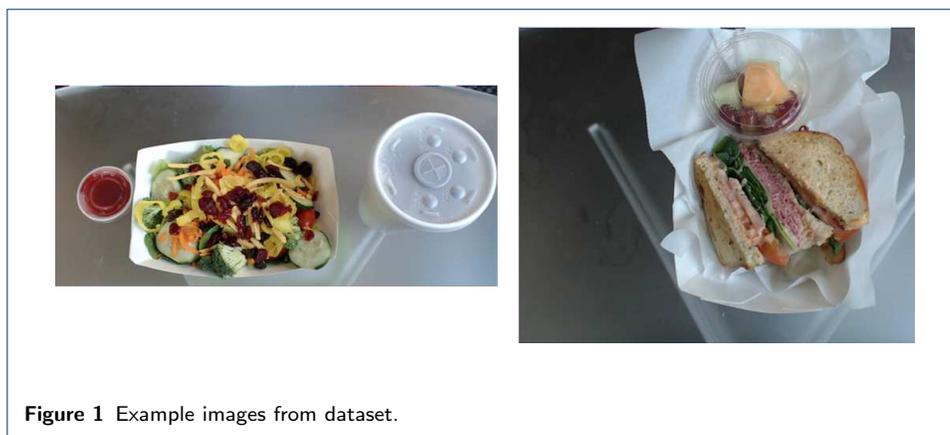
23 A method that facilitates accurate capture of food information with little or no  
24 reliance on recall would be a valuable research tool. Such a method could also  
25 revolutionize our approach to helping people with their nutrition needs. An image-  
26 based dietary assessment tool could provide quantifiably accurate assessments with  
27 no reliance on recall and extremely low burden. We developed a prototype machine  
28 learning algorithm to demonstrate that foods can be recognized via a smartphone  
29 meal photograph, which can then be used to estimate calories and other nutritional

30 content. This paper reports the methods we used to employ machine learning to  
31 develop a prototype algorithm along with results of our initial testing.

## 32 **Methods**

### 33 Description of Image Dataset

34 Our image dataset was derived from a study [5] that used photographs of lunchtime  
35 meals to estimate calories purchased. The dataset contains 7721 images of meals,  
36 with each image containing several different food items. The labels for each image  
37 come from a list of food items from  $F = 1239$  categories. Otherwise stated, we have a  
38 dataset  $D = \{x_i, \tilde{y}_i\}_{i=1}^N$ , where  $x_i$  is the  $i$ -th image in the dataset, and  $\tilde{y}_i$  is a binary  
39 vector of length 1239 with ones at the indices corresponding to food items present  
40 in image  $x_i \in \mathbb{R}^{224 \times 224 \times 3}$ . Note that this could be framed as an object detection  
41 problem, but we do not have bounding boxes around the objects of interest (i.e.  
42 the food items), as that would require significant manual work to obtain. Instead,  
43 we were limited to binary labels indicating presence or absence of the food items.  
44 A few example pictures are shown in Figure 1.



### 45 Development of Broader Categories Recognizable by Image

46 Since there are  $F = 1239$  categories for only  $N = 7721$  images, we needed to  
47 make the categories coarser to have more examples per category to properly train a  
48 machine learning system. We designed  $B = 22$  broader categories, which are parents

of the original  $F = 1239$  categories. We achieved this task via a manually crafted assignment matrix  $A \in \{0, 1\}^{22 \times 1239}$ , which is essentially a binary table with a 1 at location  $(i, j)$  if original class  $j$  is a subcategory of coarse class  $i$ .

For example, the entry (Soup, 12oz Chicken Noodle) is a 1, but the entry (Desserts, Meat Lovers Pizza) is a 0 since “12oz Chicken Noodle” is a subcategory of “Soup”, but “Meat Lovers Pizza” is not a subcategory of “Desserts”.

In order to convert a label  $\tilde{y}_i$  for an image  $x_i$  into a coarse label vector  $y_i$ , we defined the following:

$$y_i = \mathbb{1}(A\tilde{y}_i > 0) \quad (1)$$

where  $\mathbb{1}(\cdot)$  is the elementwise indicator function. This process established a dataset of  $N$  images with corresponding coarse label vectors – i.e.  $D = \{x_i, y_i\}_{i=1}^N$ , with  $x_i \in \mathbb{R}^{224 \times 224 \times 3}$  and  $y_i \in \{0, 1\}^{22}$ .

Having reduced the number of categories from 1239 to 22, the number of instances of each label in the dataset is sufficient to train a deep learning system to label food items according to the broad categories (the number of instances for each coarse label are shown in Table 1).

**Table 1** Sample counts for all 22 broad categories. Note the counts do not add up to  $N = 7721$  since there are multiple food items per image.

Bread	Burger	Cereal	Cheese	Desserts	Dressing	Drink	Fried sides	Fruit	Meats	Pasta
926	318	56	901	551	2956	2715	1164	621	1825	400
Pastry	Pizza	Rice	Salad	Salty snacks	Sandwich	Soup	Sushi	Veggies	Wrap	Yogurt
204	224	856	4147	705	1671	899	163	2241	515	174

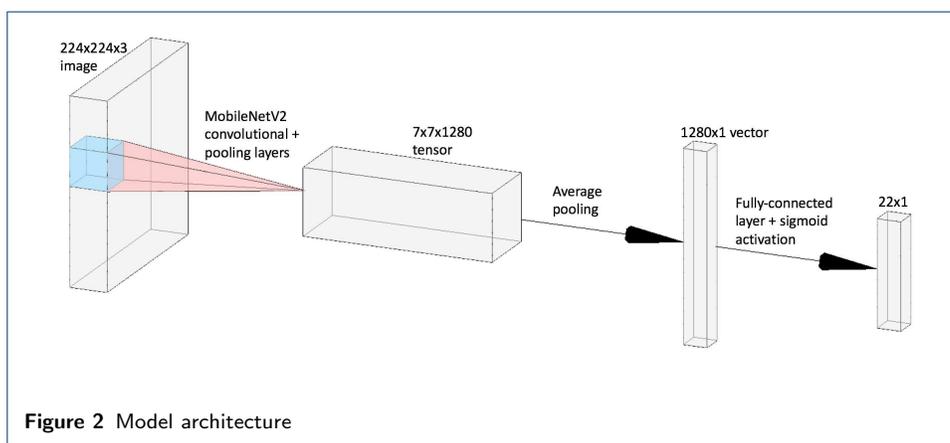
## Data split & Augmentation

From our images and corresponding coarse labels, we split the dataset into 3 mutually exclusive subsets: a training set  $D_{train}$  (5250 images), a validation set  $D_{val}$  (1312 images), and a test set  $D_{test}$  (1159 images). Our system is trained on  $D_{train}$ . We use  $D_{val}$  to validate the performance of the system and tune hyperparameters. Finally, our results are reported on the test set  $D_{test}$ . Furthermore, we augment

68 our training set using random rotations of the images between 0 and 20 degrees,  
 69 random horizontal & vertical flips (each with probability 50%), random horizontal  
 70 shifts between  $-0.2W$  and  $0.2W$  pixels (where  $W$  is the image width), random  
 71 vertical shifts between  $-0.2H$  and  $0.2H$  pixels (where  $H$  is the image height), and  
 72 random central crops (between 1.1x and 1.6x zoom since most of the items are cen-  
 73 trally located). These augmentations help the model gain a better understanding  
 74 of the data distributions for different food items.

75 Model architecture

The convolutional neural network [6] architecture consists of the convolutional & pooling layers of MobileNetV2 [7] (since the intended use case would be on a mobile device), followed by an average pooling operation and a single fully connected layer – a diagram detailing the model architecture can be found in Figure 2. For



simplicity of notation, denote the convolutional & pooling layers of MobileNetV2 to be a function  $MobileNet(\cdot) : \mathbb{R}^{224 \times 224 \times 3} \rightarrow \mathbb{R}^{7 \times 7 \times 1280}$ . Specifically, the model is:

$$a = \text{MobileNet}(x) \quad (2)$$

$$a_{avg}^{(k)} = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W a^{(i,j,k)} \quad \forall k \in \{1, \dots, 1280\} \quad (3)$$

$$\hat{y} = \sigma(a_{avg}V + b) \quad (4)$$

76 where  $\sigma(\cdot)$  is the elementwise sigmoid function,  $V \in \mathbb{R}^{1280 \times 22}$  is a learned weight  
 77 matrix,  $b \in \mathbb{R}^{22}$  is a learned bias vector,  $x \in \mathbb{R}^{224 \times 224 \times 3}$  is the input image, and  
 78  $\hat{y} \in \mathbb{R}^{22}$  is the prediction vector.  $a^{(i,j,k)}$  denotes the  $(i, j, k)$ -th element of  $a$ , and  $a_{avg}^{(k)}$   
 79 denotes the  $k$ -th element of  $a_{avg}$ . Henceforth, we will use the notation  $\hat{y} = \text{CNN}(x)$ ,  
 80 where  $\text{CNN}(\cdot) : \mathbb{R}^{224 \times 224 \times 3} \rightarrow \mathbb{R}^{22}$  is the function described in equations (2)-(4).

### 81 Loss & Training

We use an elementwise binary cross-entropy loss to train the model:

$$\hat{y}_i = \text{CNN}(x_i) \quad (5)$$

$$\mathcal{L}_i = -\frac{1}{22} \sum_{k=1}^{22} y_i^{(k)} \cdot \log(\hat{y}_i^{(k)}) + (1 - y_i^{(k)}) \cdot \log(1 - \hat{y}_i^{(k)}) \quad (6)$$

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i \quad (7)$$

82 Here,  $\mathcal{L}$  represents the loss on a batch of size  $n$ , and  $i$  is used to denote the  $i$ -th  
 83 image in the batch. We train the entire network for 200 epochs using the Adam  
 84 optimizer [8] with a learning rate of 0.001 and a batch size of 32.

### 85 Threshold selection

86 To convert the vector  $\hat{y}_i \in \mathbb{R}^{22}$  into a binary prediction vector  $\hat{y}_i^B \in \{0, 1\}^{22}$ , we  
 87 have to select a threshold for each dimension/category of  $\hat{y}_i$ . To do so, we sweep  
 88 over a range of thresholds between 0 and 1, and – for each dimension/category –

89 we select the threshold which maximized balanced accuracy for that category on  
 90 the validation dataset  $D_{val}$ .

91 **Results**

**Table 2** Per-class results. The rows are #Train (the number of training examples), followed by AUC, accuracy, sensitivity, specificity.

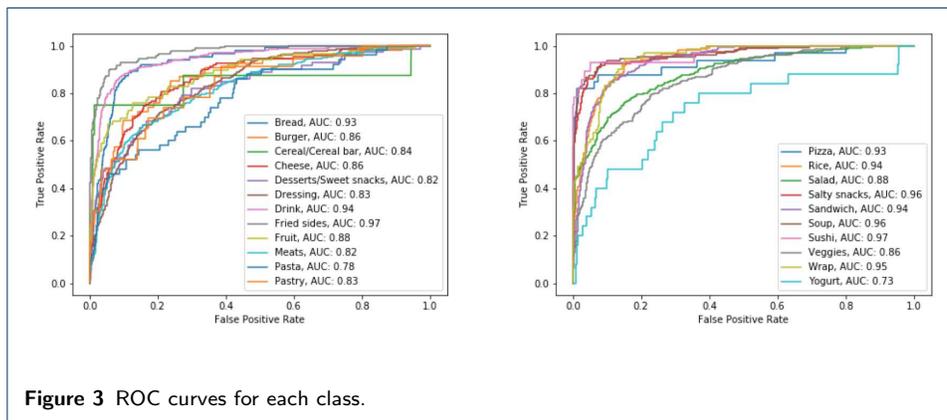
	Bread	Burger	Cereal	Cheese	Desserts	Dressing	Drink	Fried sides	Fruit	Meats	Pasta
#Train	625	222	41	597	373	1989	1851	791	421	1271	272
AUC	0.93	0.86	0.84	0.86	0.82	0.83	0.94	0.97	0.88	0.82	0.78
Acc.	0.89	0.91	0.95	0.85	0.79	0.75	0.89	0.91	0.92	0.78	0.84
Sens.	0.87	0.56	0.75	0.69	0.69	0.77	0.87	0.92	0.67	0.68	0.56
Spec.	0.89	0.93	0.95	0.87	0.80	0.73	0.91	0.91	0.94	0.81	0.85

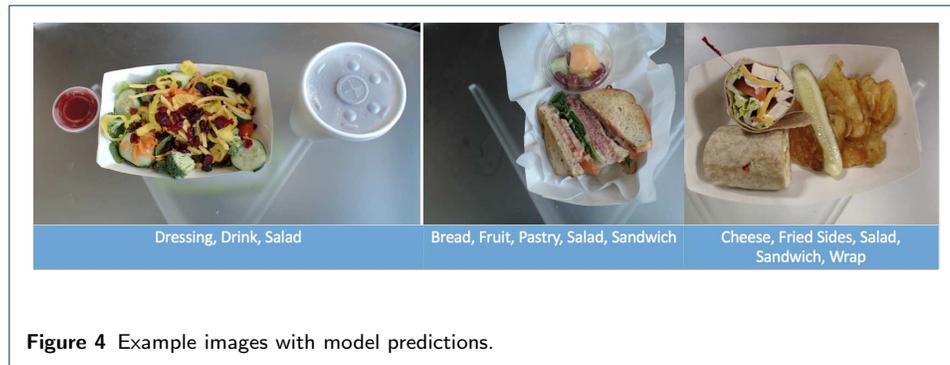
  

	Pastry	Pizza	Rice	Salad	Salty Snacks	Sandwich	Soup	Sushi	Veggies	Wrap	Yogurt
#Train	151	153	587	2807	489	1158	618	100	1523	350	116
AUC	0.83	0.93	0.94	0.88	0.96	0.94	0.96	0.97	0.86	0.95	0.73
Acc.	0.90	0.96	0.86	0.79	0.92	0.88	0.94	0.96	0.76	0.86	0.91
Sens.	0.52	0.82	0.88	0.83	0.91	0.83	0.88	0.89	0.80	0.92	0.40
Spec.	0.90	0.96	0.85	0.74	0.92	0.90	0.95	0.96	0.75	0.85	0.92

92 We report our results on the test set  $D_{test}$ . The per-class Area Under the ROC  
 93 Curve (AUC), accuracies, sensitivities, and specificities are shown in Table 2. Accu-  
 94 racy is a little misleading here since there are so many negative examples for each  
 95 class (i.e. one could build a classifier which predicts 0 items for every image and still  
 96 get good accuracy), so we instead look to the AUC values, as well as sensitivities  
 97 and specificities to get a better sense of our model’s performance.

98 AUC values ranged from 0.73 for yogurt to 0.97 for sushi. Figure 3 shows the  
 99 ROC curves for each food category, which show that the model works very well for  
 100 most items.





101 A few test set images and their corresponding predictions are shown in Figure  
102 4. As can be seen, the algorithm correctly identifies the categories of foods in each  
103 photo. These categories are broad, however, and additional information would need  
104 to be gathered from the user to maximize the usefulness of the program.

## 105 Discussion

106 An algorithm that provides recognition and proper categorization of foods from  
107 a photograph holds great promise for research and practice applications. We de-  
108 veloped such a prototype algorithm and demonstrated its ability to predict food  
109 categories with high accuracy. With the food categories linked to a database of  
110 nutritional content (and a portion-size estimator), the applications are numerous.

111 For the research community, the ability to capture food information without re-  
112 liance on recall or burdensome diaries will be a valuable resource. Such a method  
113 will greatly mitigate measurement bias in such studies. From an end-user perspec-  
114 tive, food recognition can be incorporated into programs and applications to track  
115 calories or other nutritional content, or provide guidance on specific diets to promote  
116 health (e.g., the DASH diet for hypertension) or weight loss.

117 The food categories are not granular enough alone to be useful for the kinds  
118 of applications we envision. However, programs could include a menu for users to  
119 input additional, or if needed, corrective, information. For example, the sandwich in  
120 Figure 4 is recognized and properly categorized as a sandwich, but the user would

121 have to indicate what is on the sandwich and the size of it. A program could be  
122 customizable to a person’s usual preferences for food (e.g., the category “sandwich”  
123 could default to a “half-size ham & cheese on white bread”, and “drink” could  
124 default to a “16 oz sweetened iced tea”). Additionally, crowd-sourcing could be  
125 employed to gather additional training data for the model.

## 126 **Conclusion**

127 We used machine learning to develop an algorithm that recognizes categories of  
128 food from a meal photograph. Our initial tests show a high degree of accuracy for  
129 placing food items in a meal photo into the correct categories. Next steps include in-  
130 corporation of this type of algorithm into an application that research participants,  
131 and—when more refined—consumers and can utilize.

## 132 **Declarations**

133 Ethical Approval and Consent to participate

134 Not applicable.

135 Consent for publication

136 Not applicable.

137 Availability of supporting data and materials

138 Data may be made available by contacting AJV.

139 Below are the details to access the code used for this project:

- 140 • Project name: Food Labeling with Deep Learning
- 141 • Project home page: <https://github.com/sergeassaad/nutrition>
- 142 • DOI: 10.5281/zenodo.3675840
- 143 • Operating system(s): Platform independent
- 144 • Programming language: Python
- 145 • Other requirements: Tensorflow, Keras
- 146 • License: MIT

147 Competing interests

148 The authors declare that they have no competing interests.

149 **Funding**

150 The PACE Study was funded by R01 CA184473-01A1 from the National Cancer Institute, National Institutes of  
151 Health. The funding body played no role in study design, data collection, analysis, interpretation of data, or writing  
152 the manuscript.

153 Author's contributions

154 AJV conceived of the study. LC, SA, and AJV designed the analysis plan. SA performed the analyses. SA and AJV  
155 drafted the manuscript. LC provided critical review of the manuscript.

156 Acknowledgements

157 Not applicable.

158 Author details

159 <sup>1</sup>Department of Electrical & Computer Engineering, Duke University, Science Drive, 27708 Durham, NC, USA.

160 <sup>2</sup>Department of Family Medicine & Community Health, Duke University, 2200 West Main Street, Suite 400, 27705  
161 Durham, NC, USA.

162 References

- 163 1. Quesenberry Jr, C.P., Caan, B., Jacobson, A.: Obesity, Health Services Use, and Health Care Costs Among  
164 Members of a Health Maintenance Organization. *Archives of Internal Medicine* **158**(5), 466–472 (1998).  
165 doi:[10.1001/archinte.158.5.466](https://doi.org/10.1001/archinte.158.5.466)
- 166 2. Thompson, D., Brown, J.B., Nichols, G.A., Elmer, P.J., Oster, G.: Body mass index and future healthcare costs:  
167 A retrospective cohort study. *Obesity Research* **9**(3), 210–218 (2001). doi:[10.1038/oby.2001.23](https://doi.org/10.1038/oby.2001.23)
- 168 3. Gonzalez Carrascosa, R., Garcia Segovia, P., Martinez Monzo, J.: Paper and pencil vs online self-administered  
169 food frequency questionnaire (FFQ) applied to university population: a pilot study. *Nutricion hospitalaria* **26**(6),  
170 1378–1384 (2011). doi:[10.1590/S0212-16112011000600027](https://doi.org/10.1590/S0212-16112011000600027)
- 171 4. Subar, A.F., Kirkpatrick, S.I., Mittl, B., Zimmerman, T.P., Thompson, F.E., Bingley, C., Willis, G., Islam, N.G.,  
172 Baranowski, T., McNutt, S., Potischman, N.: The automated self-administered 24-hour dietary recall (ASA24):  
173 A resource for researchers, clinicians, and educators from the national cancer institute. *Journal of the Academy*  
174 *of Nutrition and Dietetics* **112**(8), 1134–1137 (2012). doi:[10.1016/j.jand.2012.04.016](https://doi.org/10.1016/j.jand.2012.04.016)
- 175 5. Viera, A.J., Tuttle, L., Olsson, E., Gras-Najjar, J., Gizlice, Z., Hales, D., Linnan, L., Lin, F.-C., Noar, S.M.,  
176 Ammerman, A.: Effects of physical activity calorie expenditure (PACE) labeling: study design and baseline  
177 sample characteristics. *BMC Public Health* **17**(1) (2017). doi:[10.1186/s12889-017-4710-0](https://doi.org/10.1186/s12889-017-4710-0)
- 178 6. LeCun, Y., Haffner, P., Bottou, L., Bengio, Y.: Object recognition with gradient-based learning. In: *Shape,*  
179 *Contour and Grouping in Computer Vision*, p. 319. Springer, London, UK, UK (1999).  
180 <http://dl.acm.org/citation.cfm?id=646469.691875>
- 181 7. Sandler, M., Howard, A.G., Zhu, M., Zhmoginov, A., Chen, L.: Inverted residuals and linear bottlenecks: Mobile  
182 networks for classification, detection and segmentation. *CoRR abs/1801.04381* (2018). [1801.04381](https://arxiv.org/abs/1801.04381)
- 183 8. Kingma, D.P., Ba, J.: Adam: A Method for Stochastic Optimization (2014). [1412.6980](https://arxiv.org/abs/1412.6980)

184 Figure legend

- 185 ● **Figure 1:** Example images from the dataset. The corresponding labels are (Dressing, Salad, Drink) for the  
186 left image, and (Sandwich, Fruit) for the right image.
- 187 ● **Figure 2:** Model architecture used for image labeling. The model employs the convolutional and pooling  
188 layers of the MobileNetV2 architecture and an additional fully connected layer for prediction.
- 189 ● **Figure 3:** Per-class ROC curves. Note, a perfect classifier is one for which the ROC curve reaches the top  
190 left corner (i.e. True Positive Rate=1.0, False Positive Rate=0.0), and for which AUC=1.0. The model  
191 shows strong performance for most food categories.
- 192 ● **Figure 4:** Examples of images from the dataset with the corresponding predictions from the model.

# Figures



Figure 1

Example images from the dataset. The corresponding labels are (Dressing, Salad, Drink) for the left image, and (Sandwich, Fruit) for the right image.

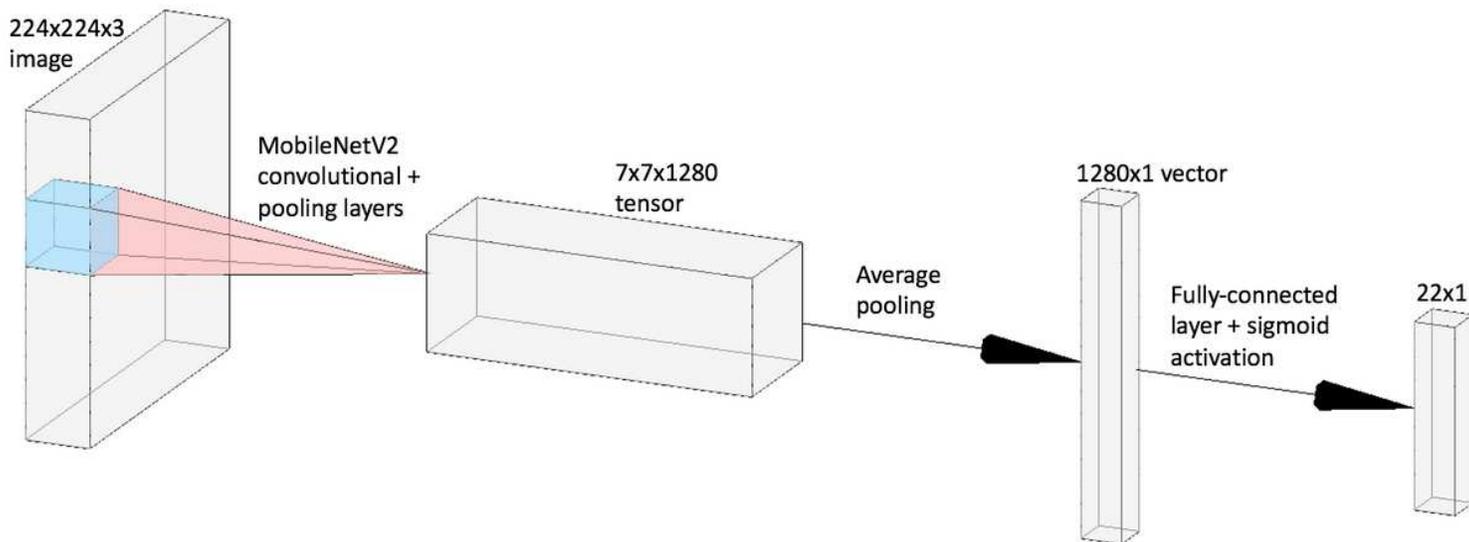
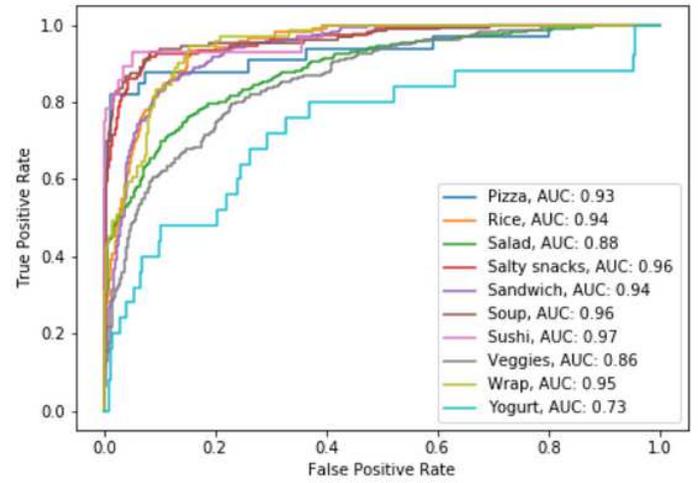
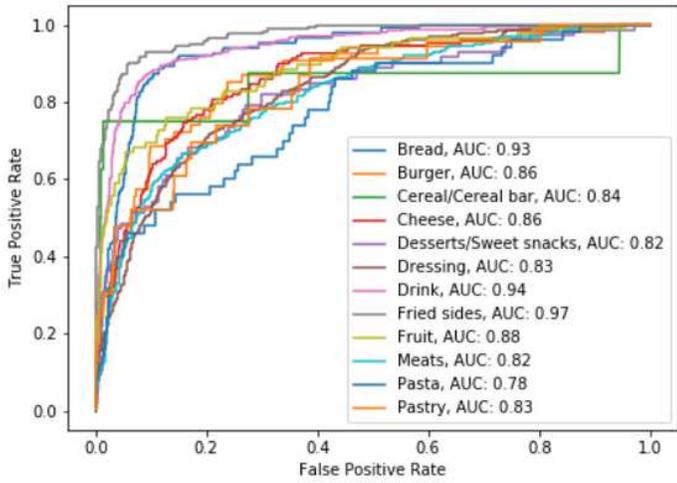


Figure 2

Model architecture used for image labeling. The model employs the convolutional and pooling layers of the MobileNetV2 architecture and an additional fully connected layer for prediction.



**Figure 3**

Per-class ROC curves. Note, a perfect classifier is one for which the ROC curve reaches the top left corner (i.e. True Positive Rate=1.0, False Positive Rate=0.0), and for which AUC=1.0. The model shows strong performance for most food categories.



**Figure 4**

Examples of images from the dataset with the corresponding predictions from the model.