

Journal of Visualized Experiments

Deep Neural Networks for Image-Based Dietary Assessment

--Manuscript Draft--

Article Type:	Invited Methods Collection - Author Produced Video
Manuscript Number:	JoVE61906R1
Full Title:	Deep Neural Networks for Image-Based Dietary Assessment
Corresponding Author:	Simon Mezgec, M.S. Jožef Stefan International Postgraduate School Ljubljana, Ljubljana SLOVENIA
Corresponding Author's Institution:	Jožef Stefan International Postgraduate School
Corresponding Author E-Mail:	simon.mezgec@gmail.com
Order of Authors:	Simon Mezgec, M.S. Barbara Koroušič Seljak
Additional Information:	
Question	Response
Please indicate whether this article will be Standard Access or Open Access.	Open Access (US\$3000)
Please confirm that you have read and agree to the terms and conditions of the author license agreement that applies below:	I agree to the Author License Agreement
Please specify the section of the submitted manuscript.	Engineering
Please provide any comments to the journal here.	

1 TITLE:

2 Deep Neural Networks for Image-Based Dietary Assessment

3

4 AUTHORS AND AFFILIATIONS:

5 Simon Mezgec¹, Barbara Koroušić Seljak²

6

7 ¹Jožef Stefan International Postgraduate School, Ljubljana, Slovenia

8 ²Computer Systems Department, Jožef Stefan Institute, Ljubljana, Slovenia

9

10 Email addresses of co-authors:

11 Barbara Koroušić Seljak (barbara.korousic@ijs.si)

12

13 Corresponding author:

14 Simon Mezgec (simon.mezgec@gmail.com)

15

16 KEYWORDS:

17 food image recognition; food image segmentation; deep neural networks; deep learning;
18 convolutional neural networks; dietary assessment

19

20 SUMMARY:

21 The goal of the work presented in this article is to develop technology for automated recognition
22 of food and beverage items from images taken by mobile devices. The technology comprises of
23 two different approaches – the first one performs food image recognition while the second one
24 performs food image segmentation.

25

26 ABSTRACT:

27 Due to the issues and costs associated with manual dietary assessment approaches, automated
28 solutions are required to ease and speed up the work and increase its quality. Today, automated
29 solutions are able to record a person's dietary intake in a much simpler way, such as by taking an
30 image with a smartphone camera. In this article, we will focus on such image-based approaches
31 to dietary assessment. For the food image recognition problem, deep neural networks have
32 achieved the state of the art in recent years, and we present our work in this field. In particular,
33 we first describe the method for food and beverage image recognition using a deep neural
34 network architecture, called NutriNet. This method, like most research done in the early days of
35 deep learning-based food image recognition, is limited to one output per image, and therefore
36 unsuitable for images with multiple food or beverage items. That is why approaches that perform
37 food image segmentation are considerably more robust, as they are able to identify any number
38 of food or beverage items in the image. We therefore also present two methods for food image
39 segmentation – one is based on fully convolutional networks (FCNs), and the other on deep
40 residual networks (ResNet).

41

42 INTRODUCTION:

43 Dietary assessment is a crucial step in determining actionable areas of an individual's diet.
44 However, performing dietary assessment using traditionally manual approaches is associated

45 with considerable costs. These approaches are also prone to errors as they often rely on self-
46 reporting by the individual. Automated dietary assessment addresses these issues by providing a
47 simpler way to quantify and qualify food intake. Such an approach can also alleviate some of the
48 errors present in manual approaches, such as missed meals, inability to accurately assess food
49 volume, etc. Therefore, there are clear benefits to automating dietary assessment by developing
50 solutions that identify different foods and beverages and quantify food intake¹. These solutions
51 can also be used to enable an estimation of nutritional values of food and beverage items
52 (henceforth 'food items'). Consequently, automated dietary assessment is useful for multiple
53 applications – from strictly medical uses, such as allowing dietitians to more easily and accurately
54 track and analyze their patients' diets, to the usage inside well-being apps targeted at the general
55 population.

56
57 Automatically recognizing food items from images is a challenging computer vision problem. This
58 is due to foods being typically deformable objects, and due to the fact that a large amount of the
59 food item's visual information can be lost during its preparation. Additionally, different foods can
60 appear to be very similar to each other, and the same food can appear to be substantially
61 different on multiple images². Furthermore, the recognition accuracy depends on many more
62 factors, such as image quality, whether the food item is obstructed by another item, distance
63 from which the image was taken, etc. Recognizing beverage items presents its own set of
64 challenges, the main one being the limited amount of visual information that is available in an
65 image. This information could be the beverage color, beverage container color and structure,
66 and, under optimal image conditions, the beverage density².

67
68 To successfully recognize food items from images, it is necessary to learn features of each food
69 and beverage class. This was traditionally done using manually-defined feature extractors³⁻⁶ that
70 perform recognition based on specific item features like color, texture, size, etc., or a
71 combination of these features. Examples of these feature extractors include multiple kernel
72 learning⁴, pairwise local features⁵ and the bag-of-features model⁶. Due to the complexity of food
73 images, these approaches mostly achieved a low classification accuracy – between 10% and
74 40%³⁻⁵. The reason for this is that the manual approach is not robust enough to be sufficiently
75 accurate. Because a food item can vary significantly in appearance, it is not feasible to encompass
76 all these variances manually. Higher classification accuracy can be achieved with manually-
77 defined feature extractors when either the number of food classes is reduced⁵, or different image
78 features are combined⁶, thus indicating that there is a need for more complex solutions to this
79 problem.

80
81 This is why deep learning proved to be so effective for the food image recognition problem. Deep
82 learning, or deep neural networks, was inspired by biological brains, and allows computational
83 models composed of multiple processing layers to automatically learn features through training
84 on a set of input images^{7,8}. Because of this, deep learning has substantially improved the state of
85 the art in a variety of research fields⁷, with computer vision, and subsequently food image
86 recognition, being one of them².

87
88 In particular, deep convolutional neural networks (DCNNs) are most popular for food image

89 recognition – these networks are inspired by the visual system of animals, where individual
90 neurons try to gain an understanding of the visual input by reacting to overlapping regions in the
91 visual field⁹. A convolutional neural network takes the input image and performs a series of
92 operations in each of the network layers, the most common of which are convolutional, fully-
93 connected and pooling layers. Convolutional layers contain learnable filters that respond to
94 certain features in the input data, whereas fully-connected layers compose output data from
95 other layers to gain higher-level knowledge from it. The goal of pooling layers is to down-sample
96 the input data². There are two approaches to using deep learning models that proved popular:
97 taking an existing deep neural network definition^{10,11}, referred to as a deep learning architecture
98 in this article, or defining a new deep learning architecture^{12,13}, and training either one of these
99 on a food image dataset. There are strengths and weaknesses to both approaches – when using
100 an existing deep learning architecture, an architecture that performed well for other problems
101 can be chosen and fine-tuned for the desired problem, thus saving time and ensuring that a
102 validated architecture has been chosen. Defining a new deep learning architecture, on the other
103 hand, is more time-intensive, but allows the development of architectures that are specifically
104 made to take into account the specifics of a problem and thus theoretically perform better for
105 that problem.

106
107 In this article, we present both approaches. For the food image recognition problem, we
108 developed a novel DCNN architecture called NutriNet², which is a modification of the well-known
109 AlexNet architecture¹⁴. There are two main differences compared to AlexNet: NutriNet accepts
110 512x512-pixel images as input (as opposed to 256x256-pixel images for AlexNet), and NutriNet
111 has an additional convolutional layer at the beginning of the neural network. These two changes
112 were introduced in order to extract as much information from the recognition dataset images as
113 possible. Having higher-resolution images meant that there is more information present on
114 images and having more convolutional layers meant that additional knowledge could be
115 extracted from the images. Compared to AlexNet’s around 60 million parameters, NutriNet
116 contains less parameters: approximately 33 million. This is because of the difference in
117 dimensionality at the first fully-connected layer caused by the additional convolutional layer².
118 **Figure 1** contains a diagram of the NutriNet architecture. The food images that were used to train
119 the NutriNet model were gathered from the Internet – the procedure is described in the protocol
120 text.

121
122 For the food image segmentation problem, we used two different existing architectures: fully
123 convolutional networks (FCNs)¹⁵ and deep residual networks (ResNet)¹⁶, both of which
124 represented the state of the art for image segmentation when we used them to develop their
125 respective food image segmentation solutions. There are multiple FCN variants that were
126 introduced by Long et al.: FCN-32s, FCN-16s and FCN-8s¹⁵. FCN-32s outputs a pixel map based on
127 the predictions by the FCN’s final layer, whereas the FCN-16s variant combines these predictions
128 with those by an earlier layer. FCN-8s considers yet another layer’s predictions and is therefore
129 able to make predictions at the finest grain, which is why it is suitable for food image recognition.
130 The FCN-8s that we used was pre-trained on the PASCAL Visual Object Classes (PASCAL VOC)
131 dataset¹⁷ and trained and tested on images of food replicas (henceforth ‘fake food’)¹⁸ due to their
132 visual resemblance to real food and due to a lack of annotated images of real food on a pixel

133 level. Fake food is used in different behavioral studies and images are taken for all dishes from
134 all study participants. Because the food contents of these images are known, it makes the image
135 dataset useful for deep learning model training. Dataset processing steps are described in the
136 protocol text.

137
138 The ResNet-based solution was developed in the scope of the Food Recognition Challenge
139 (FRC)¹⁹. It uses the Hybrid Task Cascade (HTC)²⁰ method with a ResNet-101¹⁶ backbone. This is a
140 state-of-the-art approach for the image segmentation problem that can use different feature
141 extractors, or backbones. We considered other backbone networks as well, particularly other
142 ResNet variants such as ResNet-50¹⁶, but ResNet-101 was the most suitable due to its depth and
143 ability to represent input images in a complex enough manner. The dataset used for training the
144 HTC ResNet-101 model was the FRC dataset with added augmented images. These
145 augmentations are presented in the protocol text.

146
147 This article is intended as a resource for machine learning experts looking for information about
148 which deep learning architectures and data augmentation steps perform well for the problems
149 of food image recognition and segmentation, as well as for nutrition researchers looking to use
150 our approach to automate food image recognition for use in dietary assessment. In the following
151 section, deep learning solutions and datasets from the food image recognition field are
152 presented. In the protocol text, we detail how each of the three approaches was used to train
153 deep neural network models that can be used for automated dietary assessment. Additionally,
154 each protocol section contains a description of how the food image datasets used for training
155 and testing were acquired and processed.

156
157 DCNNs generally achieved substantially better results than other methods for food image
158 recognition and segmentation, which is why the vast majority of recent research in the field is
159 based on these networks. Kawano et al. used DCNNs to complement manual approaches²¹ and
160 achieved a classification accuracy of 72.26% on the UEC-FOOD100 dataset²². Christodoulidis et
161 al. used them exclusively to achieve a higher accuracy of 84.90% on a self-acquired dataset²³.
162 Tanno et al. developed DeepFoodCam – a smartphone app for food image recognition that uses
163 DCNNs²⁴. Liu et al. presented a system that performs an Internet of Things-based dietary
164 assessment using DCNNs²⁵. Martinel et al. introduced a DCNN-based approach that exploits the
165 specifics of food images²⁶ and reported an accuracy of 90.27% on the Food-101 dataset²⁷. Zhou
166 et al. authored a review of deep learning solutions in the food domain²⁸.

167
168 Recently, Zhao et al. proposed a network specifically for food image recognition in mobile
169 applications²⁹. This approach uses a smaller ‘student’ network that learns from a larger ‘teacher’
170 network. With it, they managed to achieve an accuracy of 84% on the UEC-FOOD256³⁰ and an
171 accuracy of 91.2% on the Food-101 dataset²⁷. Hafiz et al. used DCNNs to develop a beverage-only
172 image recognition solution and reported a very high accuracy of 98.51%³¹. Shimoda et al.
173 described a novel method for detecting plate regions in food images without the usage of pixel-
174 wise annotation³². Ciocca et al. introduced a new dataset containing food items from 20 different
175 food classes in 11 different states (solid, sliced, creamy paste, etc.) and presented their approach
176 for training recognition models that are able to recognize the food state, in addition to the food

177 class³³. Knez et al. evaluated food image recognition solutions for mobile devices³⁴. Finally,
178 Furtado et al. conducted a study on how the human visual system compares to the performance
179 of DCNNs and found out that human recognition still outperforms DCNNs with an accuracy of
180 80% versus 74.5%³⁵. The authors noted that with a small number of food classes, the DCNNs
181 perform well, but on a dataset with hundreds of classes, human recognition accuracy is higher³⁵,
182 highlighting the complexity of the problem.

183
184 Despite its state-of-the-art results, deep learning has a major drawback – it requires a large input
185 dataset to train the model on. In the case of food image recognition, a large food image dataset
186 is required, and this dataset needs to encompass as many different real-world scenarios as
187 possible. In practice this means that for each individual food or beverage item, a large collection
188 of images is required, and as many different items need to be present in the dataset. If there are
189 not enough images for a specific item in the dataset, that item is unlikely to be recognized
190 successfully. On the other hand, if only a small number of items is covered by the dataset, the
191 solution will be limited in scope, and only able to recognize a handful of different foods and
192 beverages.

193
194 Multiple datasets were made available in the past. The Pittsburgh Fast-Food Image Dataset
195 (PFID)³ was introduced to encourage more research in the field of food image recognition. The
196 University of Electro-Communications Food 100 (UEC-FOOD100)²² and University of Electro-
197 Communications Food 256 (UEC-FOOD256)³⁰ datasets contain Japanese dishes, expanded with
198 some international dishes in the case of the UEC-FOOD256 dataset. The Food-101 dataset
199 contains popular dishes acquired from a website²⁷. The Food-50³⁶ and Video Retrieval Group
200 Food 172 (VireoFood-172)³⁷ datasets are Chinese-based collections of food images. The
201 University of Milano-Bicocca 2016 (UNIMIB2016) dataset is composed of images of food trays
202 from an Italian canteen³⁸. Recipe1M is a large-scale dataset of cooking recipes and food images³⁹.
203 The Food-475 dataset⁴⁰ collects four previously published food image datasets^{27,30,36,37} into one.
204 The Beijing Technology and Business University Food 60 (BTBUFood-60) is a dataset of images
205 meant for food detection⁴¹. Recently, the ISIA Food-500 dataset⁴² of miscellaneous food images
206 was made available. In comparison to other publicly available food image datasets, it contains a
207 large number of images, divided into 500 food classes, and is meant to advance the development
208 of multimedia food recognition solutions⁴².

209
210 **PROTOCOL:**

211
212 **1. Food image recognition with NutriNet**

213
214 **1.1. Obtaining the food image dataset**

215
216 1.1.1. Gather a list of different foods and beverages that will be the outputs of the food image
217 recognition model. A varied list of popular foods and beverages is preferred, as that will allow
218 the training of a robust food image recognition model.

219
220 1.1.2. Save the food and beverage list in a text file (e.g., 'txt' or 'csv').

221

222 NOTE: The text file used by the authors of this article can be found in the supplemental files
223 ('food_items.txt') and includes a list of 520 Slovenian food items.

224

225 1.1.3. Write or download a Python⁴³ script that uses the Google Custom Search API⁴⁴ to download
226 images of each food item from the list and saves them into a separate folder for each food item.

227

228 NOTE: The Python script used by the authors of this article can be found in the supplemental files
229 ('download_images.py'). If this script is used, the Developer Key (variable 'developerKey', line 8
230 in the Python script code) and Custom Search Engine ID (variable 'cx', line 28 in the Python script
231 code) need to be replaced with values specific to the Google account being used.

232

233 1.1.4. Run the Python script from step 1.1.3 (e.g., with the command: 'python
234 download_images.py').

235

236 **1.2. (Optional) Cleaning the food image dataset**

237

238 1.2.1. Train a food image detection model in the same way as in section 1.4, except use only two
239 outputs (food, non-food) as opposed to the list of outputs from step 1.1.1.

240

241 NOTE: The authors of this article used images combined from recipe websites and the ImageNet
242 dataset⁴⁵ to train the food image detection model. Since the focus here is on food image
243 recognition and this is an optional step for cleaning the recognition dataset, further details are
244 omitted. Instead, more details about this approach can be found in Mezgec et al.².

245

246 1.2.2. Run the detection model from step 1.2.1 on the food image dataset that is the result of
247 step 1.1.4.

248

249 1.2.3. Delete every image that was tagged as non-food by the detection model from step 1.2.1.

250

251 1.2.4. Manually check the food image dataset for other erroneous or low-quality images, and for
252 image duplicates.

253

254 1.2.5. Delete images found in step 1.2.4.

255

256 **1.3. Augmenting the food image dataset**

257

258 1.3.1. Create a new version of each image from the food image dataset by rotating it by 90° using
259 the CLODSA library⁴⁶ (lines 19 to 21 in the included Python script).

260

261 NOTE: The Python script containing all the CLODSA commands used by the authors of this article
262 can be found in a file included in the supplemental files ('nutrinet_augmentation.py'). If this script
263 is used, the Input Path (variable 'INPUT_PATH', line 8 in the Python script code) and Output Path
264 (variable 'OUTPUT_PATH', line 11 in the Python script code) need to be replaced with paths to

265 the desired folders.

266

267 1.3.2. Create a new version of each image from the food image dataset by rotating it by 180°
268 using the CLoDSA library (lines 19 to 21 in the included Python script).

269

270 1.3.3. Create a new version of each image from the food image dataset by rotating it by 270°
271 using the CLoDSA library (lines 19 to 21 in the included Python script).

272

273 1.3.4. Create a new version of each image from the food image dataset by flipping it horizontally
274 using the CLoDSA library (lines 23 and 24 in the included Python script).

275

276 1.3.5. Create a new version of each image from the food image dataset by adding random color
277 noise to it using the CLoDSA library (lines 26 and 27 in the included Python script).

278

279 1.3.6. Create a new version of each image from the food image dataset by zooming into it by 25%
280 using the CLoDSA library (lines 29 and 30 in the included Python script).

281

282 1.3.7. Save images from steps 1.3.1–1.3.6, along with the original images (lines 16 and 17 in the
283 included Python script), into a new food image dataset (in total, 7 variants per food image). This
284 is done by executing the command in line 32 of the included Python script.

285

286 **1.4. Performing food image recognition**

287

288 1.4.1. Import the food image dataset from step 1.3.7 into the NVIDIA DIGITS environment⁴⁷,
289 dividing the dataset into training, validation and testing subsets in the NVIDIA DIGITS user
290 interface.

291

292 1.4.2. Copy and paste the definition text of the NutriNet architecture² into NVIDIA DIGITS as a
293 custom network.

294

295 NOTE: The NutriNet architecture definition text can be found in the supplemental files
296 ('nutrinet.prototxt').

297

298 1.4.3. (Optional) Define training hyperparameters in the NVIDIA DIGITS user interface.

299

300 NOTE: Hyperparameters are parameters that are used to define the training process prior to its
301 start. The hyperparameters used by the authors of this article can be found in a file included in
302 the supplemental files ('nutrinet_hyperparameters.prototxt'). While experimentation is needed
303 for each dataset to find the optimal hyperparameters, the file contains a hyperparameter
304 configuration which can be copied into the NVIDIA DIGITS user interface. Furthermore, NVIDIA
305 DIGITS populates the hyperparameters with default values which can be used as a baseline. This
306 step is therefore optional.

307

308 1.4.4. Run the training of the NutriNet model.

309
310 1.4.5. After training is complete, take the best performing NutriNet model iteration. This model
311 is then used for testing the performance of this approach.

312
313 NOTE: There are multiple ways to determine the best-performing model iteration. A
314 straightforward way to do this is as follows. NVIDIA DIGITS outputs a graph of accuracy measures
315 for each training epoch. Check which epoch achieved the lowest loss value for the validation
316 subset of the food image dataset – that model iteration can be considered best-performing. An
317 optional step in determining the best-performing model iteration is to observe how the loss value
318 for the training subset changes from epoch to epoch and if it starts to drop continuously while
319 the loss value for the validation subset remains the same or rises continuously, take the epoch
320 prior to this drop in training loss value, as that can signal when the model started overfitting on
321 the training images.

322
323 **2. Food image segmentation with FCNs**

324
325 **2.1. Obtaining the fake-food image dataset**

326
327 2.1.1. Obtain a dataset of fake-food images. Fake-food images are gathered by researchers
328 conducting behavioral studies using food replicas.

329
330 NOTE: The authors of this article received images of fake food that were collected in a lab
331 environment¹⁸.

332
333 2.1.2. Manually annotate every food image on a pixel level – each pixel in the image must contain
334 information about which food class it belongs to. The result of this step is one annotation image
335 for each image from the food image dataset, where each pixel represents one of the food classes.

336
337 NOTE: There are many tools to achieve this – the authors of this article used JavaScript Segment
338 Annotator⁴⁸.

339
340 **2.2. Augmenting the fake-food image dataset**

341
342 2.2.1. Perform the same steps as in section 1.3, but only on images from the training subset of
343 the food image dataset.

344
345 NOTE: With the exception of step 1.3.5, all data augmentation steps need to be performed on
346 corresponding annotation images as well. If the script from section 1.3 is used, the Input Path
347 (variable 'INPUT_PATH', line 8 in the Python⁴³ script code) and Output Path (variable
348 'OUTPUT_PATH', line 11 in the Python script code) need to be replaced with paths to the desired
349 folders. In addition, set the Problem (variable 'PROBLEM', line 6 in the Python script code) to
350 'instance_segmentation' and the Annotation Mode (variable 'ANNOTATION_MODE', line 7 in the
351 Python script code) and Output Mode (variable 'OUTPUT_MODE', line 10 in the Python script
352 code) to 'coco'.

353

354 **2.3. Performing fake-food image segmentation**

355

356 2.3.1. Perform the same steps as in section 1.4, with the exception of step 1.4.2 In place of that
357 step, perform steps 2.3.2 and 2.3.3.

358

359 NOTE: Hyperparameters are parameters that are used to define the training process prior to its
360 start. The training hyperparameters used by the authors of this article for the optional step 1.4.3
361 can be found in a file included in the supplemental files ('fcn-8s_hyperparameters.prototxt').
362 While experimentation is needed for each dataset to find the optimal set of hyperparameters,
363 the file contains a hyperparameter configuration which can be copied into the NVIDIA DIGITS⁴⁷
364 user interface. Furthermore, NVIDIA DIGITS populates the hyperparameters with default values
365 which can be used as a baseline.

366

367 2.3.2. Copy and paste the definition text of the FCN-8s architecture¹⁵ into the NVIDIA DIGITS
368 environment as a custom network.

369

370 NOTE: The FCN-8s architecture definition text is publicly available on GitHub⁴⁹.

371

372 2.3.3. Enter the path to the pre-trained FCN-8s model weights into the NVIDIA DIGITS user
373 interface.

374

375 NOTE: These model weights were pre-trained on the PASCAL VOC dataset¹⁷ and can be found on
376 the Internet⁴⁹.

377

378 **3. Food image segmentation with HTC ResNet**

379

380 **3.1. Obtaining the food image dataset**

381

382 3.1.1. Download the food image dataset from the FRC website¹⁹.

383

384 **3.2. Augmenting the food image dataset**

385

386 3.2.1. Perform steps 1.3.1–1.3.4.

387

388 NOTE: The Python⁴³ script containing all the CLODSA⁴⁶ commands used by the authors of this
389 article can be found in a file included in the supplemental files ('frc_augmentation.py'). If this
390 script is used, the Input Path (variable 'INPUT_PATH', line 8 in the Python script code) and Output
391 Path (variable 'OUTPUT_PATH', line 11 in the Python script code) need to be replaced with paths
392 to the desired folders.

393

394 3.2.2. Create a new version of each image from the food image dataset by adding Gaussian blur
395 to it using the CLODSA library (lines 26 and 27 in the included Python script).

396

397 3.2.3. Create a new version of each image from the food image dataset by sharpening it using
398 the CLoDSA library (lines 29 and 30 in the included Python script).

399
400 3.2.4. Create a new version of each image from the food image dataset by applying gamma
401 correction to it using the CLoDSA library (lines 32 and 33 in the included Python script).

402
403 3.2.5. Save images from steps 3.2.1–3.2.4, along with the original images (lines 16 and 17 in the
404 included Python script), into a new food image dataset (in total, 8 variants per food image). This
405 is done by executing the command in line 35 of the included Python script.

406
407 3.2.6. Save images from steps 3.2.2–3.2.4, along with the original images (lines 16 and 17 in the
408 included Python script), into a new food image dataset (in total, 4 variants per food image). This
409 is done by deleting lines 19 to 24 of the included Python script and executing the command in
410 line 35.

411
412 **3.3. Performing food image segmentation**

413
414 3.3.1. Modify the existing HTC²⁰ ResNet-101 architecture¹⁶ definition from the MMDetection
415 library⁵⁰ in sections ‘model settings’ and ‘dataset settings’ of the architecture definition file so
416 that it accepts the food image datasets from steps 3.1.1, 3.2.5 and 3.2.6.

417
418 3.3.2. (Optional) Modify the HTC ResNet-101 architecture definition from step 3.3.1 to define
419 training hyperparameters: batch size in section ‘dataset settings’, solver type and learning rate
420 in section ‘optimizer’, learning policy in section ‘learning policy’ and number of training epochs
421 in section ‘runtime settings’ of the architecture definition file.

422
423 NOTE: The modified HTC ResNet-101 architecture definition file can be found in the supplemental
424 files (‘htc_resnet-101.py’). Hyperparameters are parameters that are used to define the training
425 process prior to its start. While experimentation is needed for each dataset to find the optimal
426 set of hyperparameters, the file already contains a hyperparameter configuration which can be
427 used without modification. This step is therefore optional.

428
429 3.3.3. Run the training of the HTC ResNet-101 model on the food image dataset from step 3.1.1
430 using the MMDetection library (e.g., with the command: ‘python mmdetection/tools/train.py
431 htc_resnet-101.py’).

432
433 3.3.4. After the training from step 3.3.3 is complete, take the best-performing HTC ResNet-101
434 model iteration and fine-tune it by running the next phase of training on the food image dataset
435 from step 3.2.5.

436
437 NOTE: There are multiple ways to determine the best-performing model iteration. A
438 straightforward way to do this is as follows. The MMDetection library outputs values of accuracy
439 measures for each training epoch in the command line interface. Check which epoch achieved
440 the lowest loss value for the validation subset of the food image dataset – that model iteration

441 can be considered best-performing. An optional step in determining the best-performing model
442 iteration is to observe how the loss value for the training subset changes from epoch to epoch
443 and if it starts to drop continuously while the loss value for the validation subset remains the
444 same or rises continuously, take the epoch prior to this drop in training loss value, as that can
445 signal when the model started overfitting on the training images.

446
447 3.3.5. After the training from step 3.3.4 is complete, take the best-performing HTC ResNet-101
448 model iteration and fine-tune it by running the next phase of training on the food image dataset
449 from step 3.2.6.

450
451 NOTE: See note for step 3.3.4.

452
453 3.3.6. After the training from step 3.3.5 is complete, take the best-performing HTC ResNet-101
454 model iteration and fine-tune it by again running the next phase of training on the food image
455 dataset from step 3.2.5.

456
457 NOTE: See note for step 3.3.4.

458
459 3.3.7. After the training from step 3.3.6 is complete, take the best-performing HTC ResNet-101
460 model iteration. This model is then used for testing the performance of this approach.

461
462 NOTE: See note for step 3.3.4. Steps 3.3.3–3.3.7 yielded the best results for the purposes defined
463 by the authors of this article. Experimentation is needed for each dataset to find the optimal
464 sequence of training and data augmentation steps.

465
466 **REPRESENTATIVE RESULTS:**

467 NutriNet was tested against three popular deep learning architectures of the time: AlexNet¹⁴,
468 GoogLeNet⁵¹ and ResNet¹⁶. Multiple training parameters were also tested for all architectures to
469 define the optimal values². Among these is the choice of solver type, which determines how the
470 loss function is minimized. This function is the primary quality measure for training neural
471 networks as it is better suited for optimization during training than classification accuracy. We
472 tested three solvers: Stochastic Gradient Descent (SGD)⁵², Nesterov’s Accelerated Gradient
473 (NAG)⁵³ and the Adaptive Gradient algorithm (AdaGrad)⁵⁴. The second parameter is batch size,
474 which defines the number of images that are processed at the same time. The depth of the deep
475 learning architecture determined the value of this parameter, as deeper architectures require
476 more space on the GPU memory – the consequence of this approach was that the memory was
477 completely filled with images for all architectures, regardless of depth. The third parameter is
478 learning rate, which defines the speed with which the neural network parameters are being
479 changed during training. This parameter was set in unison with the batch size, as the number of
480 concurrently processed images dictates the convergence rate. AlexNet models were trained
481 using a batch size of 256 images and a base learning rate of 0.02; NutriNet used a batch size of
482 128 images and a rate of 0.01; GoogLeNet 64 images and a rate of 0.005; and ResNet 16 images
483 and a rate of 0.00125. Three other parameters were fixed for all architectures: learning rate
484 policy (step-down), step size (30%) and gamma (0.1). These parameters jointly describe how the

485 learning rate is changing in every epoch. The idea behind this approach is that the learning rate
486 is being gradually lowered to fine-tune the model the closer it gets to the optimal loss value.
487 Finally, the number of training epochs was also fixed to 150 for all deep learning architectures².

488
489 The best result among all the parameters tested that NutriNet achieved was a classification
490 accuracy of 86.72% on the recognition dataset, which was around 2% higher than the best result
491 for AlexNet and slightly higher than GoogLeNet’s best result. The best-performing architecture
492 overall was ResNet (by around 1%), however the training time for ResNet is substantially higher
493 compared to NutriNet (by a factor of approximately five), which is important if models are
494 continuously re-trained to improve accuracy and the number of recognizable food items.
495 NutriNet, AlexNet and GoogLeNet achieved their best results using the AdaGrad solver, whereas
496 ResNet’s best model used the NAG solver. NutriNet was also tested on the publicly available
497 UNIMIB2016 food image dataset³⁸. This dataset contains 3,616 images of 73 different food items.
498 NutriNet achieved a recognition accuracy of 86.39% on this dataset, slightly outperforming the
499 baseline recognition result of the authors of the dataset, which was 85.80%. Additionally,
500 NutriNet was tested on a small dataset of 200 real-world images of 115 different food and
501 beverage items, where NutriNet achieved a top-five accuracy of 55%.

502
503 To train the FCN-8s fake-food image segmentation model, we used Adam⁵⁵ as the solver type, as
504 we found that it performed optimally for this task. The base learning rate was set very low – to
505 0.0001. The reason for the low number is the fact that only one image could be processed at a
506 time, which is a consequence of the pixel-level classification process. The GPU memory
507 requirements for this approach are significantly greater than image-level classification. The
508 learning rate thus had to be set low so that the parameters were not being changed too fast and
509 converge to less optimal values. The number of training epochs was set to 100, while the learning
510 rate policy, step size and gamma were set to step-down, 34% and 0.1, respectively, as these
511 parameters produced the most accurate models.

512
513 Accuracy measurements of the FCN-8s model were performed using the pixel accuracy
514 measure¹⁵, which is analogous to the classification accuracy of traditional deep learning
515 networks, the main difference being that the accuracy is computed on the pixel level instead of
516 on the image level:

517

$$518 \quad PA = \frac{\sum_i n_{ii}}{\sum_i t_i}$$

519
520 where PA is the pixel accuracy measure, n_{ij} is the number of pixels of class i predicted to belong
521 to class j and $t_i = \sum_j n_{ij}$ is the total number of pixels from class i in the ground-truth labels¹. In
522 other words, the pixel accuracy measure is computed by dividing correctly predicted pixels with
523 the total number of pixels. The final accuracy of the trained FCN-8s model was 92.18%. Figure 2
524 shows three example images from the fake-food image dataset (one from each of the training,
525 validation and testing subsets), along with the corresponding ground-truth and model prediction
526 labels.

527

528 The parameters to train the HTC²⁰ ResNet-101 model for food image segmentation were set as
529 follows: the solver type used was SGD because it outperformed other solver types. The base
530 learning rate was set to 0.00125 and the batch size to 2 images. The number of training epochs
531 was set to 40 per training phase, and multiple training phases were performed – first on the
532 original FRC dataset without augmented images, then on the 8x-augmented and 4x-augmented
533 FRC dataset multiple times in an alternating fashion, each time taking the best-performing model
534 and fine-tuning it in the next training phase. More details on the training phases are found in
535 section 3.3. of the protocol text. Finally, the step-down learning policy was used, with fixed
536 epochs for when the learning rate decreased (epochs 28 and 35 for the first training phase). An
537 important thing to note is that while this sequence of training phases produced the best results
538 in our testing in the scope of the FRC, using another dataset might require a different sequence
539 to produce optimal results.

540

541 This ResNet-based solution for food image segmentation was evaluated using the following
542 precision measure¹⁹:

543

$$544 \quad P_{IoU \geq 0.5} = \frac{TP_{IoU \geq 0.5}}{TP_{IoU \geq 0.5} + FP_{IoU \geq 0.5}}$$

545

546 where P is precision, TP is the number of true positive predictions by the food image
547 segmentation model, FP is the number of false positive predictions and IoU is Intersection over
548 Union, which is computed with this equation:

549

$$550 \quad IoU = \frac{Area\ of\ Overlap}{Area\ of\ Union}$$

551

552 where *Area of Overlap* represents the number predictions by the model that overlap with the
553 ground truth, and *Area of Union* represents the total number of predictions by the model
554 together with the ground truth, both on a pixel level and for each individual food class. Recall is
555 used as a secondary measure and is calculated in a similar way, using the following formula¹⁹:

556

$$557 \quad R_{IoU \geq 0.5} = \frac{TP_{IoU \geq 0.5}}{TP_{IoU \geq 0.5} + FN_{IoU \geq 0.5}}$$

558

559 where R is recall and FN is the number of false negative predictions by the food image
560 segmentation model. The precision and recall measures are then averaged across all classes in
561 the ground truth. Using these measures, our model achieved an average precision of 59.2% and
562 an average recall of 82.1%, which ranked second in the second round of the Food Recognition
563 Challenge¹⁹. This result was 4.2% behind the first place and 5.3% ahead of the third place in terms
564 of the average precision measure. Table 1 contains the results for the top-4 participants in the
565 competition.

566

567 **FIGURE AND TABLE LEGENDS:**

568 **Figure 1: Diagram of the NutriNet deep neural network architecture.** This figure has been

569 published in Mezgec et al.².

570

571 **Figure 2: Images from the fake-food image dataset.** Original images (left), manually-labelled
572 ground-truth labels (middle) and predictions from the FCN-8s model (right). This figure has been
573 published in Mezgec et al.¹.

574

575 **Table 1: Top-4 results from the second round of the Food Recognition Challenge.** Average
576 precision is taken as the primary performance measure and average recall as a secondary
577 measure. Results are taken from the official competition leaderboard¹⁹.

578

579 **DISCUSSION:**

580 In recent years, deep neural networks have been validated multiple times as a suitable solution
581 for recognizing food images^{10–12,21,23,25,26,29,31,33}. Our work presented in this article serves to
582 further prove this^{1,2}. The single-output food image recognition approach is straightforward and
583 can be used for simple applications where images with only one food or beverage item are
584 expected².

585

586 The food image segmentation approach seems particularly suitable for recognizing food images
587 in general, without any restriction on the number of food items¹. Because it works by classifying
588 each individual pixel of the image, it is able to not only recognize any number of food items in
589 the image, but also specify where a food item is located, as well as how large it is. The latter can
590 then be used to perform food weight estimation, particularly if used with either a reference
591 object or a fixed-distance camera.

592

593 There has been some work done regarding the availability of food image datasets^{3,22,27,30,36–42},
594 and we hope more will be done in the future, particularly when it comes to aggregating food
595 image datasets from different regions across the world, which would enable more robust
596 solutions to be developed. Currently, the accuracy of automatic food image recognition solutions
597 has not yet reached human-level accuracy³⁵, and this is likely in large part due to food image
598 datasets of insufficient size and quality.

599

600 In the future, our goal will be to further evaluate the developed procedures on real-world images.
601 In general, datasets in this field often contain images taken in controlled environments or images
602 that were manually optimized for recognition. This is why it is important to gather a large and
603 diverse real-world food image dataset to encompass all the different food and beverage items
604 that individuals might want to recognize. The first step towards this was provided by the Food
605 Recognition Challenge, which included a dataset of real-world food images¹⁹, but further work
606 needs to be done to validate this approach on food images from all around the world and in
607 cooperation with dietitians.

608

609 **ACKNOWLEDGMENTS:**

610 The authors would like to thank Tamara Bucher from the University of Newcastle, Australia, for
611 providing the fake-food image dataset. This work was supported by the European Union's
612 Horizon 2020 research and innovation programs (grant numbers 863059 – FNS-Clouds, 769661 –

613 SAAM); and the Slovenian Research Agency (grant number P2-0098). The European Union and
614 Slovenian Research Agency had no role in the design, analysis or writing of this article.

615

616 **DISCLOSURES:**

617 The authors have nothing to disclose.

618

619 **REFERENCES:**

620 1. Mezgec, S., Eftimov, T., Bucher, T., Koroušič Seljak, B. Mixed Deep Learning and Natural
621 Language Processing Method for Fake-Food Image Recognition and Standardization to Help
622 Automated Dietary Assessment. *Public Health Nutrition*. **22** (7), 1193–1202 (2019).

623 2. Mezgec, S., Koroušič Seljak, B. NutriNet: A Deep Learning Food and Drink Image Recognition
624 System for Dietary Assessment. *Nutrients*. **9** (7), 657 (2017).

625 3. Chen, M. et al. PFID: Pittsburgh Fast-Food Image Dataset. *Proceedings of the ICIP 2009*. 289–
626 292 (2009).

627 4. Joutou, T., Yanai, K. A Food Image Recognition System with Multiple Kernel Learning.
628 *Proceedings of the ICIP 2009*. 285–288 (2009).

629 5. Yang, S., Chen, M., Pomerlau, D., Sukthankar, R. Food Recognition using Statistics of Pairwise
630 Local Features. *Proceedings of the CVPR 2010*. 2249–2256 (2010).

631 6. Anthimopoulos, M. M., Gianola, L., Scarnato, L., Diem, P., Mougiakakou, S. G. A Food
632 Recognition System for Diabetic Patients Based on an Optimized Bag-of-Features Model. *IEEE*
633 *Journal of Biomedical and Health Informatics*. **18** (4), 1261–1271 (2014).

634 7. LeCun, Y., Bengio, Y., Hinton, G. Deep Learning. *Nature*. **521**, 436–444 (2015).

635 8. Deng, L., Yu, D. Deep Learning: Methods and Applications. *Foundations and Trends in Signal*
636 *Processing*. **7** (3–4), 197–387 (2014).

637 9. Hubel, D. H., Wiesel, T. N. Receptive Fields, Binocular Interaction and Functional Architecture
638 in the Cat's Visual Cortex. *The Journal of Physiology*. **160** (1), 106–154 (1962).

639 10. Singla, A., Yuan, L., Ebrahimi, T. Food/Non-Food Image Classification and Food Categorization
640 using Pre-Trained GoogLeNet Model. *Proceedings of the MADiMa'16*. 3–11 (2016).

641 11. Yanai, K., Kawano, Y. Food Image Recognition using Deep Convolutional Network with Pre-
642 Training and Fine-Tuning. *Proceedings of the ICMEW 2015*. 1–6 (2015).

643 12. Liu, C. et al. DeepFood: Deep Learning-Based Food Image Recognition for Computer-Aided
644 Dietary Assessment. *Proceedings of the ICOST 2016*. 37–48 (2016).

645 13. De Sousa Ribeiro, F. et al. An End-to-End Deep Neural Architecture for Optical Character
646 Verification and Recognition in Retail Food Packaging. *Proceedings of the ICIP 2018*. 2376–2380
647 (2018).

648 14. Krizhevsky, A., Sutskever, I., Hinton, G. ImageNet Classification with Deep Convolutional
649 Neural Networks. *Proceedings of the NIPS'12*. 1097–1105 (2012).

650 15. Long, J., Shelhamer, E., Darrell, T. Fully Convolutional Networks for Semantic Segmentation.
651 *Proceedings of the CVPR 2015*. 3431–3440 (2015).

652 16. He, K., Zhang, X., Ren, S., Sun, J. Deep Residual Learning for Image Recognition. *Proceedings*
653 *of the CVPR 2016*. 770–778 (2016).

654 17. PASCAL VOC Project. PASCAL Visual Object Classes. <http://host.robots.ox.ac.uk/pascal/VOC>
655 (2020).

656 18. Bucher, T., van der Horst, K., Siegrist, M. Fruit for Dessert. How People Compose Healthier

657 Meals. *Appetite*. **60** (1), 74–80 (2013).

658 19. AICrowd. Food Recognition Challenge. [https://www.aicrowd.com/challenges/food-](https://www.aicrowd.com/challenges/food-recognition-challenge)

659 recognition-challenge (2020).

660 20. Chen, K. et al. Hybrid Task Cascade for Instance Segmentation. *Proceedings of the CVPR 2019*.

661 4974–4983 (2019).

662 21. Kawano, Y., Yanai, K. Food Image Recognition with Deep Convolutional Features. *Proceedings*

663 *of the UbiComp 2014*. 589–593 (2014).

664 22. Matsuda, Y., Hoashi, H., Yanai, K. Recognition of Multiple-Food Images by Detecting

665 Candidate Regions. *Proceedings of the ICME 2012*. 25–30 (2012).

666 23. Christodoulidis, S., Anthimopoulos, M. M., Mougiakakou, S. G. Food Recognition for Dietary

667 Assessment using Deep Convolutional Neural Networks. *Proceedings of the ICIAP 2015*. 458–465

668 (2015).

669 24. Tanno, R., Okamoto, K., Yanai, K. DeepFoodCam: A DCNN-Based Real-Time Mobile Food

670 Recognition System. *Proceedings of the MADiM α '16*. 89–89 (2016).

671 25. Liu, C. et al. A New Deep Learning-Based Food Recognition System for Dietary Assessment on

672 An Edge Computing Service Infrastructure. *IEEE Transactions on Services Computing*. **11** (2), 249–

673 261 (2017).

674 26. Martinel, N., Foresti, G. L., Micheloni, C. Wide-Slice Residual Networks for Food Recognition.

675 *Proceedings of the IEEE WACV 2018*. 567–576 (2018).

676 27. Bossard, L., Guillaumin, M., Van Gool, L. Food-101—Mining Discriminative Components with

677 Random Forests. *Proceedings of the ECCV'14*. 446–461 (2014).

678 28. Zhou, L., Zhang, C., Liu, F., Qiu, Z., He, Y. Application of Deep Learning in Food: A Review.

679 *Comprehensive Reviews in Food Science and Food Safety*. **18**, 1793–1811 (2019).

680 29. Zhao, H., Yap, K.-H., Kot, A. C., Duan, L. JDNet: A Joint-Learning Distilled Network for Mobile

681 Visual Food Recognition. *IEEE Journal of Selected Topics in Signal Processing*. **14** (4), 665–675

682 (2020).

683 30. Kawano, Y., Yanai, K. Automatic Expansion of a Food Image Dataset Leveraging Existing

684 Categories with Domain Adaptation. *Proceedings of the ECCV'14*. 3–17 (2014).

685 31. Hafiz, R., Haque, M. R., Rakshit, A., Uddin, M. S. Image-Based Soft Drink Type Classification

686 and Dietary Assessment System using Deep Convolutional Neural Network with Transfer

687 Learning. *Journal of King Saud University – Computer and Information Sciences*. In Press (2020).

688 32. Shimoda, W., Yanai, K. Weakly-Supervised Plate and Food Region Segmentation. *Proceedings*

689 *of the ICME 2020*. 1–6 (2020).

690 33. Ciocca, G., Micali, G., Napoletano, P. State Recognition of Food Images using Deep Features.

691 *IEEE Access*. **8**, 32003–32017 (2020).

692 34. Knez, S., Šajn, L. Food Object Recognition using a Mobile Device: Evaluation of Currently

693 Implemented Systems. *Trends in Food Science & Technology*. **99**, 460–471 (2020).

694 35. Furtado, P., Caldeira, M., Martins, P. Human Visual System vs Convolution Neural Networks

695 in Food Recognition Task: An Empirical Comparison. *Computer Vision and Image Understanding*.

696 **191**, 102878 (2020).

697 36. Chen, M.-Y. et al. Automatic Chinese Food Identification and Quantity Estimation. *SA'12*

698 *Technical Briefs*. 1–4 (2012).

699 37. Chen, J., Ngo, C.-W. Deep-Based Ingredient Recognition for Cooking Recipe Retrieval.

700 *Proceedings of the MM'16*. 32–41 (2016).

701 38. Ciocca, G., Napoletano, P., Schettini, R. Food Recognition: A New Dataset, Experiments, and
702 Results. *IEEE Journal of Biomedical and Health Informatics*. **21** (3), 588–598 (2017).

703 39. Salvador, A. et al. Learning Cross-Modal Embeddings for Cooking Recipes and Food Images.
704 *Proceedings of the IEEE CVPR 2017*. 3020–3028 (2017).

705 40. Ciocca, G., Napoletano, P., Schettini, R. CNN-Based Features for Retrieval and Classification
706 of Food Images. *Computer Vision and Image Understanding*. **176–177**, 70–77 (2018).

707 41. Cai, Q., Li, J., Li, H., Weng, Y. BTBUFood-60: Dataset for Object Detection in Food Field.
708 *Proceedings of the IEEE BigComp 2019*. 1–4 (2019).

709 42. Min, W. et al. ISIA Food-500: A Dataset for Large-Scale Food Recognition via Stacked Global-
710 Local Attention Network. *Proceedings of the MM'20*. 393–401 (2020).

711 43. Python Software Foundation. Python. <https://www.python.org> (2020).

712 44. Google. Google Custom Search API. [https://developers.google.com/resources/api-](https://developers.google.com/resources/api-libraries/documentation/customsearch/v1/python/latest/customsearch_v1.cse.html)
713 [libraries/documentation/customsearch/v1/python/latest/customsearch_v1.cse.html](https://developers.google.com/resources/api-libraries/documentation/customsearch/v1/python/latest/customsearch_v1.cse.html) (2020).

714 45. Stanford Vision Lab. ImageNet. <http://www.image-net.org> (2020).

715 46. Heras, J. CLoDSA. <https://github.com/joheras/CLoDSA> (2020).

716 47. NVIDIA. NVIDIA DIGITS. <https://developer.nvidia.com/digits> (2020).

717 48. Yamaguchi, K. JavaScript Segment Annotator. [https://github.com/kyamagu/js-segment-](https://github.com/kyamagu/js-segment-annotator)
718 [annotator](https://github.com/kyamagu/js-segment-annotator) (2020).

719 49. Shelhamer, E. Fully Convolutional Networks for Semantic Segmentation.
720 <https://github.com/shelhamer/fcn.berkeleyvision.org> (2020).

721 50. Multimedia Laboratory, CUHK. MMDetection. [https://github.com/open-](https://github.com/open-mmlab/mmdetection)
722 [mmlab/mmdetection](https://github.com/open-mmlab/mmdetection) (2020).

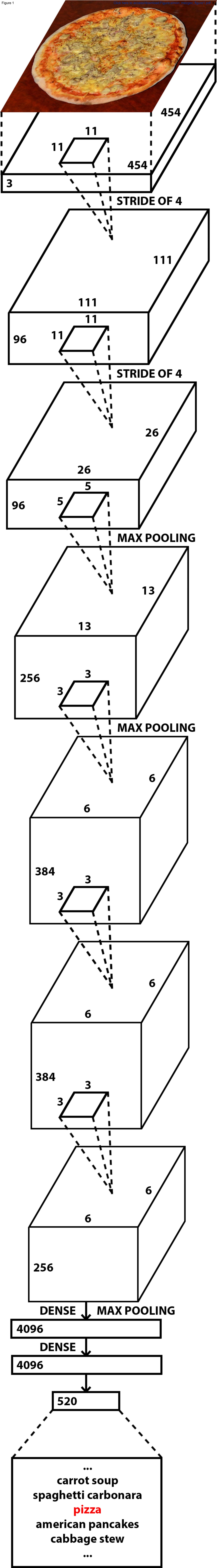
723 51. Szegedy, C. et al. Going Deeper with Convolutions. *Proceedings of the CVPR 2015*. 1–9 (2015).

724 52. Bottou, L. Large-Scale Machine Learning with Stochastic Gradient Descent. *Proceedings of the*
725 *COMPSTAT'2010*. 177–186 (2010).

726 53. Nesterov, Y. A Method of Solving a Convex Programming Problem with Convergence Rate
727 $O(1/k^2)$. *Doklady Akademii Nauk SSSR*. **27**, 372–376 (1983).

728 54. Duchi, J., Hazan, E., Singer, Y. Adaptive Subgradient Methods for Online Learning and
729 Stochastic Optimization. *Journal of Machine Learning Research*. **12**, 2121–2159 (2011).

730 55. Kingma, D. P., Ba, J. Adam: A Method for Stochastic Optimization. *arXiv Preprint*.
731 [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) (2017).



Original image

Ground truth

Predictions

Training



Validation



Testing



- apple
- apple juice
- apple tart
- banana
- banana slice
- beans
- beetroot
- biscuit
- broccoli
- burger
- butter
- cappuccino
- carrot
- cheesecake
- chicken
- coffee
- cola
- cream
- cucumber
- dark bread
- fish
- french dressing
- fries
- fruit flan
- grapes
- herring
- italian dressing
- ketchup
- kiwi
- kiwi slice
- mayonnaise
- meat loaf
- mohrenkopf
- muffin
- natural schnitzel
- onion
- orange
- orange juice
- pasta
- pear
- platzi
- potatoes
- praline
- rice
- sacher cake
- salad leaf
- salmon
- sausage
- schnitzel
- steak
- sugar
- tea
- tomato
- water
- white bread

Team Name	Placement	Average Precision	Average Recall
rssfete	1	63.4%	88.6%
simon_mezgec	2	59.2%	82.1%
arimboux	3	53.9%	73.5%
latentvec	4	48.7%	71.1%

Name of Material / Equipment	Company	Catalog Number
HARDWARE		
NVIDIA GPU	NVIDIA	N/A
SOFTWARE		
Caffe	Berkeley AI Research	N/A
CLoDSA	Jónathan Heras	N/A
Google API Client	Google	N/A
JavaScript Segment Annotator	Kota Yamaguchi	N/A
MMDetection	Multimedia Laboratory, CUHK	N/A
NVIDIA DIGITS	NVIDIA	N/A
OpenCV	Intel	N/A
Python	Python Software Foundation	N/A
PyTorch	Facebook AI Research	N/A
Ubuntu OS	Canonical	N/A

Comments/Description

An NVIDIA GPU is needed as some of the software frameworks below will not work otherwise. <https://www.nvidia.com>

Caffe is a deep learning framework. <https://caffe.berkeleyvision.org>

CLoDSA is a Python image augmentation library. <https://github.com/joheras/CLoDSA>

Google API Client is a Python client library for Google's discovery based APIs. <https://github.com/googleapis/google-api-python-client>

JavaScript Segment Annotator is a JavaScript image annotation tool. <https://github.com/kyamagu/js-segment-annotator>

MMDetection is an object detection toolbox based on PyTorch. <https://github.com/open-mmlab/mmdetection>

NVIDIA DIGITS is a wrapper for Caffe that provides a graphical web interface. <https://developer.nvidia.com/digits>

OpenCV is a library for computer vision. <https://opencv.org>

Python is a programming language. <https://www.python.org>

PyTorch is a machine learning framework. <https://pytorch.org>

Ubuntu 14.04 is the OS used by the authors and offers compatibility with all of the software frameworks and tools above. <https://ubuntu.com>

n

Dear All,

All the comments have been addressed in this revision of the article. Changes have been made to the manuscript using the Track Changes option in Microsoft Word, and changes can be viewed in the Simple Markup (for easier readability) and All Markup (all changes visible) modes in Word. The video has also been updated according to the comments.

Please note that in order to address all the reviewers' comments, some parts of the manuscript have been significantly updated. Consequently, the Introduction section turned out to be almost 2100 words long, which is over the 1500-word limit described in the manuscript template and instructions. In order to address this, and still include the necessary changes as outlined by the reviewers, the Introduction section has been split into Introduction and Related Work sections, where the latter includes the description of deep-learning-based food image recognition solutions and datasets. This way, the new Introduction section is around 1400 words long and the Related Work section has around 700 words. If necessary, this can be changed.

Included below are all the editor and reviewer comments (in an italic font), as well as our responses to them (in a blue font for better legibility). In them, we detail how the comments have been addressed and describe the changes that have been made to the manuscript and video.

In the hopes of a positive response I wish you all the best,

Simon Mezgec

Editorial and Production Comments:

Changes to be made by the Author(s) regarding the written manuscript:

1. Please take this opportunity to thoroughly proofread the manuscript to ensure that there are no spelling or grammar issues

The manuscript has been checked thoroughly and some errors have been corrected throughout the manuscript.

2. Please add more details to your protocol steps. Please ensure you answer the "how" question, i.e., how is the step performed?

Significant detail has been added to several protocol steps that should clarify how they are to be performed. For example, details have been added to augmentation steps (1.3.1 to 1.3.7. and 3.2.1. to 3.2.6.), to hyperparameter-defining steps (1.4.3., 2.3.1. and 3.3.2.), to training steps (3.3.4. to 3.3.7.), and others. More on this can be found in our replies to Reviewer #1 below.

3. 1.4.4/2.3.3: What does it mean to enter the definition? What are the user input commands here?

It is a copy-and-paste process from the definition text file into the NVIDIA DIGITS user interface. The instructions for both steps have been reworded to clarify this, and information on where to get the FCN-8s architecture definition text has been added as a note in step 2.3.3 (now step 2.3.2.).

4. 3.2.2: How is CLoDSA being used?

CLoDSA is being used to perform image segmentation steps (steps 1.3.1. to 1.3.7. and new steps 3.2.1. to 3.2.6.). The Python scripts containing all the CLoDSA commands are now attached to the manuscript (files 'nutrinet_augmentation.py' and 'frc_augmentation.py'). To further detail what exactly needs to be done in each of these steps, the CLoDSA library requirement and the specific CLoDSA command location in the included Python scripts for each step have been added. This allows the reader to simply copy-paste the commands, and only change the path of the image dataset.

Changes to be made by the Author(s) regarding the video:

1. Hold on all of the section titles for 2 or 3 more seconds before moving on. For example, the section 1 title at 0:43 is only on screen for a second and a half before it is gone.

The display duration of section titles has been prolonged to address this issue. New section titles are now held for 5 seconds (which is between 2 and 3 seconds more than before), with the exception of the title for section 1.2, which is held a little over 6 seconds due to the extra text present on the title card – this should allow the viewer to read everything without hurrying. Furthermore, the steps in the video have been updated to align with the new version of the manuscript. Lastly, the design of the section titles has been updated.

Please upload a revised high-resolution video here:

<https://www.dropbox.com/request/chF0mlCkyvhVtOprB52d?oref=e>

The revised high-resolution video has been uploaded to the above link.

Thank you for your feedback!

Reviewers' Comments:

Reviewer #1:

Manuscript Summary:

This paper describes a method for recognizing food-related images (based on a

deep neural network architecture proposed by the authors), and also two methods for segmentation of food images (down to food items). The paper presents methods for gathering the required training sets, creating the ground truth, installing and using the necessary software, and training the models. Experimental evaluation on a publicly available dataset and an in-house dataset shows promising results.

Major Concerns:

The main concern with the paper is the different levels of detail that the various steps are described. This makes it hard to identify the core audience of this work: is it machine-learning experts who want to adapt their tools and methods to food image recognition and segmentation? is it non-experts who want to apply a relatively automated method for food image recognition and segmentation? The following examples make this issue clearer.

The protocol steps have been updated significantly to equalize the level of detail across all steps. To this end, the installation steps have been removed and additional details have been added to a large number of steps. Apart from the steps mentioned below, details have been added to other steps as well (e.g., step 1.1.1., step 2.1.1., etc.). The target audience are both of the groups mentioned – machine learning experts can gain insight into what architectures and data augmentation steps perform well for the problems of food image recognition and segmentation, while non-experts can just take the provided instructions and files and run the deep learning model training as-is. To better suit the latter, significant detail has been added to the protocol steps. The target audience description has also been added to the manuscript itself (lines 147-150).

1. Software requirements: some steps in the protocol are related to installing software (python, caffe, google API client, pytorch, and more), while another step is to "write a python script that uses the google custom search API". There are some issues with this: (a) the level of detail is not consistent, (b) there is no clear definition of the underlying system (e.g. type/version of operating system, hardware requirements), (c) software requirements should be better mentioned separately, instead of listing the installation steps as part of a protocol.

The software requirements are listed in the Materials table (Excel spreadsheet file) that is attached to the manuscript. Because of this, the installation steps have been removed from the Protocol text in the manuscript to improve the consistency and help equalize the level of details in each protocol step. Furthermore, links have been added to each of the requirements in the Materials table, and three more requirements have been added to it: the OS requirement and GPU requirement (listed in a separate Hardware table section), as well as an additional software requirement (OpenCV) that was missed before.

As for step 1.1.5. (now step 1.1.3.), which contains the instruction to write a Python script – the text has been slightly altered to make it clearer that it is not mandatory to write a new Python script from scratch, as we provide our own Python script in the supplemental files of the manuscript (this is mentioned in the step 1.1.3. note).

That said, the reader who follows these steps still has the option to write their own script, and the following steps should still be applicable. This is why step 1.1.3. now reads "Write or download a Python script...". Due to the fact that our script is included, we did not go into further details in this step – if the reader is interested in a quick way to get started with deep learning model training and maybe does not have a programming background, they can just run the script without going into details. For others, they can view the code, analyze it and either adapt it for their purpose, or write another script entirely.

2. Some steps are not clear: for example what does step 1.1.7 mean?

Upon analyzing the Python script mentioned in step 1.1.6. (now step 1.1.4.), we found that step 1.1.7. is superfluous as the script already saves the resulting images in the correct format, so there is no need for an additional step after that. Thanks for pointing out this error – it has been corrected by removing step 1.1.7.

3. Machine-learning related steps: steps of 1.3 describe a process of dataset augmentation. Additional dataset augmentations are described within 3.2. On the other hand, other parts of the machine-learning process are not adequately described. For example, the note of step 3.3.4 mentions "experimentation is needed for each dataset to find the optimal hyperparameters": this is a broader challenge and a focus of relevant research. While it might be impossible to provide an "optimal" solution, authors could suggest some approach (or specific configuration), which would be more in line with a protocol description. Step 1.4.7 mentions "take the best-performing ... model": what is the best-performing? This might be clear to someone familiar with machine-learning (who understands that this is referring to validation accuracy or validation loss), but it might not be clear to someone who is not familiar and is using the paper and described protocol to implement the proposed method/system.

Step 3.3.4. (now step 3.3.2.) has been updated in multiple ways. More details have been added about what parameters exactly should be changed in the architecture definition file (if the reader wants to change them). More importantly, the note that mentions experimentation with the hyperparameters has been changed to reflect the fact that the included architecture definition file already contains a possible hyperparameter configuration that the reader can use without modifying the hyperparameters themselves. It is also now stressed that this step is optional. Since the possible hyperparameter values are already defined in the file and adding more details about the process of trying to find optimal hyperparameters would make this step very long, we did not go into further details on this topic. Steps 1.4.3. and 2.3.1. have been similarly updated.

As for step 1.4.7. (now step 1.4.5.) – this is true, and the previous step description could have confused some readers. This is why we have added a long note to this step that details how the best-performing model iteration can be determined in the

NVIDIA DIGITS software, along with an optional step that might help with overfitting. Steps 3.3.4. to 3.3.7. have been similarly updated.

Overall, the protocol described in this work follows these steps: (a) dataset and ground-truth collection, (b) dataset curation and augmentation, and (c) training of the recognition/segmentation models. Most of the steps are related to dataset augmentation and installation of various software components. As a result, the target audience as well as the desired outcome of the protocol is not clear.

As mentioned in our first reply, the installation steps were removed and the target audience are both machine learning experts and non-experts (more on this above). As for the desired outcome of the protocol – the purpose is to publicly document the methods that have led to promising results in the food image recognition field. With these instructions, other researchers (both those with machine learning expertise as well as nutrition experts who are looking to apply this approach to their own research) in the field can either take the exact approach and apply it to their own food image dataset or gather ideas from specific steps of the protocol and apply them to their own existing deep learning approach.

Thank you for your feedback!

Reviewer #2:

Manuscript Summary:

This article deals with an interesting topic image processing by the usage of deep neural networks, called 'NutriNet' with a specific application on food image classification for dietary assessment. Today, Deep neural network is considered as one of the most important topics in computer science. It is used in many studies by researchers and authors in the field of machine learning because it has achieved good results in classification and recognition.

Major Concerns:

None

Minor Concerns:

Suggestions which would improve the quality of the article but are not essential for publication:

1. Starting from the introduction authors may provide the reader with a clear picture of why automatic dietary assessment is necessary nowadays. They can discuss from health/ medical or any other point of view.

A more in-depth explanation of the usefulness of automating dietary assessment has been added to the introduction, including some applications (lines 46-49 and 52-55).

2. I suggest to add the following reference to the paper.

Hafiz, Rubaiya, et al. "Image-based soft drink type classification and dietary assessment system using deep convolutional neural network with transfer learning." Journal of King Saud University-Computer and Information Sciences (2020).

The suggested reference has been added to the paper (reference 31, cited in line 173).

Thank you for your feedback!

Reviewer #3:

Manuscript Summary:

This paper presents some deep neural networks techniques for image-based dietary assessment. This paper is well written. However, the data are newly collected but the presented techniques are not new. This paper tends to focus more on the applications instead of new technique. The novelty of this paper is fairly acceptable. I will suggest to accept this paper for publication in JoVE after the authors has successfully addressed the listed comments as below.

Major Concerns:

1. *The authors should include some comparative tables of prior works considering different parameters.*

NutriNet was mainly tested on a self-collected/self-acquired dataset (which includes the smaller real-world image dataset), and the fake-food image recognition solution was tested on the fake-food buffet (FFB) dataset, provided by Tamara Bucher from the University of Newcastle, Australia. Due to legal reasons, we cannot share either of the datasets, which means there are no other results on these two datasets that we can compare our results to. We could offer comparisons to other works on other datasets, but that comparison would not be fair, since the accuracy results are highly dependent on the size, quality and class distribution of the dataset, as well as the number of food classes present. NutriNet was also tested on a publicly available dataset (University of Milano-Bicocca 2016 – UNIMIB2016) and the comparison to the authors' baseline result is given in lines 497-500.

As for the ResNet-based solution for food image segmentation, which was tested on the Food Recognition Challenge benchmark dataset – we have included a table (Table 1) comparing our accuracy result to the results of the other top-4 participants.

2. *Experimental results are not clear. More scientific reasoning should be added in the experimental results' explanation. What are the parameters used in the proposed system and how their values are set?*

The Representative Results section of the paper has been significantly expanded to address this. For each of the three solutions (NutriNet food image recognition, FCN-8s food image segmentation and HTC ResNet-101 food image segmentation), the parameters have been described and their values given, as well as the reasoning behind their values. Additionally, the results and measures used to determine them have also been elaborated on. Lines 468-502 now contain a more in-depth description of the NutriNet food image recognition solution's parameters and results, lines 504-527 of the FCN-8s food image segmentation solution, and lines 529-566 of the HTC ResNet-101 food image segmentation solution.

3. The authors should provide some details (e.g. formula, algorithm, etc.) on the presented techniques. The presented protocols are impressive but this makes this paper like a handbook if excluded the techniques' explanations and details.

Apart from the Representative Results and the Protocol sections, the Introduction section has also been significantly expanded to elaborate further on the presented techniques and approaches used. In particular, further details about the NutriNet architecture have been added, additional description of fully convolutional networks and the fake-food dataset, and the Hybrid Task Cascade has also been expanded upon. These additions are located in lines 107-145. Some formulas for accuracy measures have been added to the Representative Results section (lines 519 and 558).

Apart from that, Figure 1 already contains a complete visual representation of the NutriNet architecture, and the entire definition is included as a supplemental file to the paper ('nutrinet.prototxt'). For the FCN-8s solution, the FCN architecture was taken from the FCN paper's GitHub page, and the definition is linked in the paper (line 371). For the HTC ResNet-101 solution, the definition is also included as a supplemental file ('htc_resnet-101.py'). All the architectures' definitions are therefore either attached to or linked in the paper, so we do not think it would be beneficial to repeat them in the manuscript text as algorithms, as that would make the paper very long. Other algorithms used, such as for data augmentation, are also now included as supplemental files ('frc_augmentation.py' and 'nutrinet_augmentation.py') and described in the protocol text.

Minor Concerns:

1. Some important recent references are missing in the fields. The authors should consider more recent research done in the field of their study (especially in the year 2020 onwards).

7 recent references from the food image recognition research field, published in the last year, have been added to the paper (references 29, 31-35 and 42, described in lines 169-183 and 206-209).

Thank you for your feedback!



[Click here to access/download](#)

Supplemental Coding Files

[simon_mezgec_supplemental_files_revision1.zip](#)

