# Food-500 Cap: A Fine-Grained Food Caption Benchmark for Evaluating Vision-Language Models

Zheng Ma*
National Key Laboratory for Novel
Software Technology,
Nanjing University,
Nanjing, China
maz@smail.nju.edu.cn

Mianzhi Pan*
National Key Laboratory for Novel
Software Technology,
Nanjing University,
Nanjing, China
panmz@smail.nju.edu.cn

Wenhan Wu
National Key Laboratory for Novel
Software Technology,
Nanjing University,
Nanjing, China
wuwh@smail.nju.edu.cn

Kanzhi Cheng
National Key Laboratory for Novel
Software Technology,
Nanjing University,
Nanjing, China
chengkz@smail.nju.edu.cn

Jianbing Zhang†
National Key Laboratory for Novel
Software Technology,
Nanjing University,
Nanjing, China
zjb@nju.edu.cn

Shujian Huang†
National Key Laboratory for Novel
Software Technology,
Nanjing University,
Nanjing, China
huangsj@nju.edu.cn

Jiajun Chen
National Key Laboratory for Novel
Software Technology,
Nanjing University,
Nanjing, China
chenjj@nju.edu.cn

## ABSTRACT

Vision-language models (VLMs) have shown impressive performance in substantial downstream multi-modal tasks. However, only comparing the fine-tuned performance on downstream tasks leads to the poor interpretability of VLMs, which is adverse to their future improvement. Several prior works have identified this issue and used various probing methods under a zero-shot setting to detect VLMs' limitations, but they all examine VLMs using general datasets instead of specialized ones. In practical applications, VLMs are usually applied to specific scenarios, such as e-commerce and news fields, so the generalization of VLMs in specific domains should be given more attention. In this paper, we comprehensively investigate the capabilities of popular VLMs in a specific field, the food domain. To this end, we build a food caption dataset, Food-500 Cap, which contains 24,700 food images with 494 categories. Each image is accompanied by a detailed caption, including fine-grained attributes of food, such as the ingredient, shape, and color. We also provide a culinary culture taxonomy that classifies each food category based on its geographic origin in order to better analyze the performance differences of VLM in different regions. Experiments on our proposed datasets demonstrate that popular VLMs underperform in the food domain compared with their performance in the general domain. Furthermore, our research reveals severe bias in VLMs' ability to handle food items from different geographic regions. We adopt diverse probing methods and evaluate nine VLMs belonging to different architectures to verify the aforementioned observations. We hope that our study will bring researchers' attention to VLM's limitations when applying them to the domain of food or culinary cultures, and spur further investigations to address this issue.

## CCS CONCEPTS

• **Information systems** → **Multimedia databases**.

## KEYWORDS

Vision-language Models; Food Benchmark; Evaluation

## 1 INTRODUCTION

Despite the remarkable success of vision-language models (VLMs) [18, 33, 34, 38, 45, 47, 52] in substantial uni-modal and multi-modal downstream tasks, they are still poorly understood as yet. The prevalent approach for evaluating VLMs is comparing their performance on downstream tasks after fine-tuning. However, evaluation solely based on the fine-tuning results renders poor interpretability [59], which hinders the further development of VLMs. Consequently, researchers have proposed a range of probing methods and benchmarks [9, 13, 29, 30, 39, 48] in recent years to assess the capabilities of VLMs from various perspectives, providing a more comprehensive understanding of these models. However, these methodologies are still limited in the general domain. They typically construct evaluation benchmarks by employing images from widely used general-domain datasets and subsequently assigning hand-crafted textual annotations to these images. If VLMs perform well in a specific domain, we can directly employ the models in that domain without any modifications. However, the above situation is

*Equal contribution.
†Corresponding author.

**Label:** agedashi tofu

**Region:** Japanese

**Caption:** Five pieces of golden agedashi tofu with some scallions, carrot strips and wooden fish flowers in a long white plate.

**Figure 1: An example from our Food-500 Cap. The image is equipped with the label, geographic origin, and a detailed description. This description is annotated with a class label (red) and hand-curated various fine-grained visible content of the image such as ingredients (blue), food colors (green), and the food container (orange).**

unclear due to only few works studying the generalization of using VLMs directly in specific domains without fine-tuning.

Motivated by this, we focus on evaluating the generalization capacity of VLMs in a specific domain, namely, the food domain. Since food computing [25] has been gaining widespread attention as it has the potential to support numerous food-related applications, such as healthy diets and food choices. To comprehensively evaluate the VLMs' performance on food-related tasks, we introduce a new benchmark named Food-500 Cap, which comprises 24,700 food images with 494 categories, each accompanied by a detailed caption. The Food-500 Cap dataset is created by selecting images from ISIA Food-500 [26] that covers a wide range of food categories. We select 50 images from each category and engage an annotation company to annotate fine-grained descriptions for all 24,700 images. Each description includes the original food category label and fine-grained attributes of the food, such as the color, shape, and ingredients. Such an in-house labeling process guarantees the high quality of our dataset. Besides, as food is always associated with a specific geographic region, we also provide a taxonomy that classifies food categories based on their original place, enabling a more comprehensive investigation of VLMs' performance across culinary cultures. We provide a sample of Food-500 Cap in Figure 1, which contains a Japanese food image labeled *agedashi tofu* from and a description with some related attributes. In contrast to the prevalent food datasets [3, 23, 54], Food-500 Cap are equipped with high-quality image captions containing richer visual information and geographic origin tags, which is more suitable for exploring the performance of VLM in the food domain.

To comprehensively evaluate VLMs' capacity in the food domain, we seriously pick up nine representative models from three popular architectures, including vision-language representation models (e.g. CLIP [33]), image-to-text generative models (e.g. OFA [52]), and text-to-image generative models (e.g. Stable Diffusion [38]). We probe these VLMs with various food-related tasks in a zero-shot setting. For vision-language representation models, we employ food classification and image-text retrieval to assess VLMs' multi-modal

alignment capabilities. As for image-to-text generative models and text-to-image generative models, we utilize image captioning and image synthesis respectively to test their multi-modal generation capabilities. Both qualitative and quantitative analyses are performed on the experimental results, revealing that these models exhibit poor performance in the food domain, in contrast to their performance in the general domain. Moreover, we find that all the models display a significant bias in culinary culture, with their performance in Asian cuisine falling markedly behind that in European, North American, and Latin American cuisine. In summary, this paper makes the following contributions:

- We equip a subset of the ISIA Food-500 dataset with (1) find-grained image descriptions (2) the geographic origin of each food category. Based on this enhanced dataset, we propose Food-500 Cap, which serves as a benchmark to evaluate the vision-language ability of VLMs in the food domain. To the best of our knowledge, Food-500 Cap is the first image-caption dataset that specifically targets the food domain.
- We evaluate nine representative VLMs from diverse architectures on our benchmark and use four probing tasks to analyze the performance of VLMs in the food domain comprehensively.
- The results of our experiments on Food-500 Cap unveil the limitations of VLMs in the food domain, as well as their bias towards specific culinary cultures.

## 2 RELATED WORKS

### 2.1 Probing VLMs

VLMs have achieved state-of-the-art performance in a large number of downstream multi-modal tasks, but they are still poorly understood. Therefore, evaluating VLMs has attracted much attention. Commonly, VLMs are evaluated by comparing their fine-tuned performance in downstream vision-language tasks. However, fine-tuning VLMs in downstream tasks only provides a black-box score, which renders poor interpretability of VLMs.

To acquire a deeper understanding of VLMs, a number of existing works have probed their capability from various perspectives, including verb understanding [9], spatial relation understanding [8, 39], visual abstract reasoning with tangram shapes [13], generalization ability in out-of-domain datasets [61], compositional reasoning ability [30, 48, 57], visual-linguistic grounding capabilities on specific linguistic phenomena [29], robustness to image and text perturbations [32], attribute recognition capability [50, 59], object hallucination problem [4]. These works have revealed that prevalent VLMs have severe shortcomings in certain aspects.

Nevertheless, current probing works are still limited in the general domain. Specifically, they utilize images from the general domain to construct datasets or benchmarks, such as MSCOCO [20], Visual Genome [15], LAION-400M [42], or social media data without domain specification. Instead of investigating VLMs in the general domain, we focus on VLMs' vision-language capability in the food domain, which is closely linked with people's health and daily life. To this end, we introduce a food image-caption dataset and comprehensively evaluate a range of representative VLMs on it

## 2.2 Food Dataset

In recent years, there have been substantial food datasets available. Most of them are proposed for food image classification, such as ETH Food-101 [1], UPMC Food-101 [53] with western food, UEC Food256 [14], Sushi-50 [31] with Japanese food, VIREO Food-172 [2], ChineseFoodNet [3] with Chinese food, ISIA Food-500 [26] comprising miscellaneous food categories worldwide. Besides category labels, UPMC Food-101 [53] and VIREO Food-172 [2] contain additional metadata such as related web text, ingredients, and cooking instructions.

In addition, Yummly-66k [24] annotates images with ingredients, courses and regions for food topic models. FoodSeg103 [54] implements a food image segmentation dataset that tags each image with multiple ingredients and draws the corresponding pixel-wise masks. Recipe1M [22] and Recipe1M+ [23] construct datasets with numerous images and recipes, which is suitable for image-recipe retrieval task. These datasets mainly facilitate specific tasks. However, the visual correlation between their textual annotations and images is relatively weak. Their texts not only neglect plenty of visual attributes of the image but also contain invisible contents, e.g. cooking instructions [22, 23]. VLMs are hard to align these images and texts, which hinders these datasets from serving as probing datasets. To this end, we introduce Food-500 Cap, the first food image-caption dataset. Food-500 Cap has captions describing fine-grained visual content of the image. It also includes food category labels and their geographic origin tags. Hence, Food-500 Cap can serve as a comprehensive benchmark for probing VLMs' generalization ability in the food domain.

## 3 FOOD-500 CAP

This section outlines the construction process of the Food-500 Cap dataset and provides a detailed description of its statistics.

### 3.1 Dataset Construction

***Collecting food images.*** To ensure the diversity of food categories, we utilize images from the ISIA Food-500 [26], a comprehensive dataset for food recognition containing 399,726 samples covering 500 food categories from various countries and regions. For each category, we randomly select 50 images from ISIA Food-500. Note that actually we only use 494 out of the 500 categories currently and six categories are manually removed.

***Annotating food images.*** To obtain high-quality captions depicting fine-grained visual features, we employ a data annotation company and urge the annotators to follow the next three rules. First, annotators must include category labels in each caption, which contain the food's principal information. Second, we encourage annotators to be as careful as possible, marking all visible content of images including not only the food's color, shape, ingredients, seasonings, accessories, etc. but also the container's color, shape, pattern, etc. To ensure the distinctiveness of the captions, some general words should be avoided to the largest degree, such as `fruit` and `vegetables`. At last, every annotator is instructed to integrate the aforementioned information into fluent sentences using diverse syntactic constructions. Eventually, we obtain food

image captions with fine-grained visual content and one example is shown in Figure 1.

***Marking regions.*** Although covering diverse food categories, one insufficiency of ISIA Food-500 is that it mixes food categories from different regions without marking their original regions, hindering further study of culinary cultures. Therefore, we resort to Wikipedia to mark the original region of each food category by ourselves, and we show the detailed process in Appendix A. Consequently, all food categories are divided into seven regions: *Worldwide*, *Western*, *Latin-American*, *Chinese*, *Japanese*, *Indian*, and *Asian*. Table 6 displays the food category distribution over these regions. Note that 90 food categories have ambiguous original places, so we merge them into *Worldwide*.

### 3.2 Dataset Statistics and Characteristics

Food-500 Cap contains 24,700 images that are uniformly divided into 494 categories. Captions are of average length 18.57, and there are 7.26 nouns, 1.96 verbs, and 2.53 adjectives in each caption on average[1]. As shown in Table 1, Food-500 Cap surpasses current food datasets in the following two aspects: (1) Food-500 Cap annotates each image with a human-crafted, fine-grained, fluent visual description, hence containing richer visual information. (2) All food categories are tagged with their geographic origins, enabling culinary culture studies across regions. Whereas almost all current datasets neglect region annotation except for VIREO Food-172 [2]. Therefore, Food-500 Cap can better serve as a comprehensive vision-language benchmark.

## 4 PROBING VLMS IN FOOD DOMAIN

To comprehensively evaluate prevalent VLMs, we pick up three different types of VLMs including vision-language representation models [18, 33, 44, 55, 58], image-to-text generative models [18, 51, 52], and text-to-image generative models [34, 38]. Then we apply four probing methods to them. For vision-language representation VLMs, we utilize classification and retrieval tasks to probe their cross-modal alignment ability. For generative VLMs, we utilize image captioning and image generation tasks to test whether they can generate satisfactory images or descriptions. All tasks are performed in a zero-shot setting to directly assess the generalization of VLMs.

### 4.1 Vision-language Representation Models

#### 4.1.1 Evaluated Models.

***CLIP [33].*** It employs two independent encoders to encode image and text respectively. It is trained with an image-text contrastive (ITC) objective, which encourages the embeddings of paired images and texts to be closer while pushing away those of mismatched pairs. Benefiting from 400 million image-text pairs during training, CLIP exhibits powerful zero-short transfer ability across a wide range of downstream tasks, e.g. cross-modal retrieval.

***TCL [55].*** It consists of an image encoder, a text encoder, and a multi-modal encoder to fuse image and text features from unimodal encoders. Apart from the original cross-modal contrastive objective like CLIP, TCL proposes intra-modal contrastive target

---

[1]https://spacy.io

| Dataset | Image Number | Category | | Annotation | |
|---|---|---|---|---|---|
| | | Number | Coverage | type | source |
| Recipe1M+ [23] | 13M | - | - | Ingredients & Cooking instructions | Web |
| FoodSeg103 [54] | 7,118 | 103 | Worldwide | Ingredients | Manual |
| UPMC Food-101 [53] | 90,840 | 101 | Western | Related web text | Web |
| UEC Food256 [14] | 25,088 | 256 | Japanese | - | - |
| VIREO Food-172 [2] | 110,241 | 172 | Chinese | Ingredients & Cooking instructions | Web |
| Sushi-50 [31] | 3,963 | 50 | Japanese | - | - |
| ChineseFoodNet [3] | 185,628 | 208 | Chinese | - | - |
| Yummly-66k [24] | 66,615 | - | - | Course & ingredients & region | Web |
| ISIA Food-500 [26] | 399,726 | 500 | Worldwide | - | - |
| Food-500 Cap | 24,700 | 494 | Worldwide | Image Captions & region | Manual |

Table 1: Summary of popular food-domain datasets and Food-500 Cap. Our proposed Food-500 Cap has 24,700 images, covering 494 food categories around the world. Compared to existing food datasets, each image from Food-500 Cap has a hand-curated fine-grained image caption and the geographic origin of the food. Captions are annotated by a data annotation company and food origins are tagged by ourselves resorting to Wikipedia.

and local mutual information maximization to robust the uni-modal representations.

**X_VLM** [58]. It shares the same model framework as TCL while better utilizing the region annotations in some datasets. It optimizes the model by predicting the location of bounding boxes in the image given the corresponding caption and meanwhile conducts vision-language alignment in multi-granularity.

**FLAVA** [44]. It inherits the architecture of TCL and X_VLM. Different from VLMs only focus on cross-modal tasks, FLAVA is trained with regard to both cross-modal and uni-modal objectives, including global image-text contrastive learning, masked image modeling, masked language modeling, etc. And FLAVA achieves comparable results on vision-only, language-only, and cross-modal tasks.

**BLIP** [18]. It introduces a novel multi-modal mixture of Encoder-Decoder framework. It can operate as a uni-modal encoder, an image-grounded text encoder, and an image-grounded text decoder, sharing parameters with each other. They are optimized with contrastive loss, image-text matching (ITM) loss, and language modeling loss respectively. To leverage noisy web data effectively, BLIP augments the datasets utilizing captions synthesized by itself.

*4.1.2 Evaluation Task.*

**Food Classification**. Previous works [33, 44] have revealed that VLMs have competitive zero-shot power in general-domain classification benchmarks, such as ImageNet [5], PASCAL VOC [7], CIFAR [16]. Furthermore, within the domain of food, CLIP [33] and FLAVA [44] report their overall accuracy on Food-101 [1], but it has a relatively limited number of 101 food classes. To this end, we employ representation VLMs to undertake zero-shot food classification using our benchmark and elaborate the results.

Following [33, 43, 60], we undertake zero-shot food classification using prompt engineering. Specifically, we utilize a prompt template "a food photo of a {label}" and populate it with related category
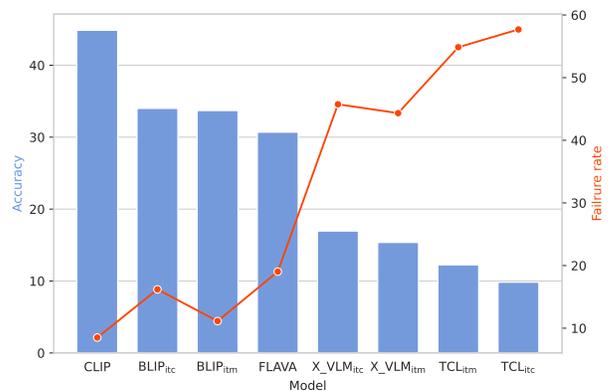


Figure 2: Results of zero-shot classification, including (1) accuracy and (2) failure rate, which represents the percentage of categories where none of the images is correctly classified. Model names with subscript "itm" are of ITM configuration and others are of ITC configuration.

labels[2]. In the evaluation phase, we task the VLMs with identifying the correct prompt for each image among all constructed prompts. To be specific, we evaluate two configurations of vision-language representation models. The first is the ITC configuration, where only the image and text uni-modal encoders are employed. Images and prompts are individually embedded by VLMs' uni-modal encoders, and models select the prompt with the maximum cosine similarity for each image. The second is the ITM configuration, where BLIP, TCL, and X_VLM further use the ITM score from their multi-modal fusion modules to re-rank the top-$k$ nearest prompts in the ITC configuration. Note that we fix $k$ to 128 and re-rank by adding the ITM score to cosine scores.

---

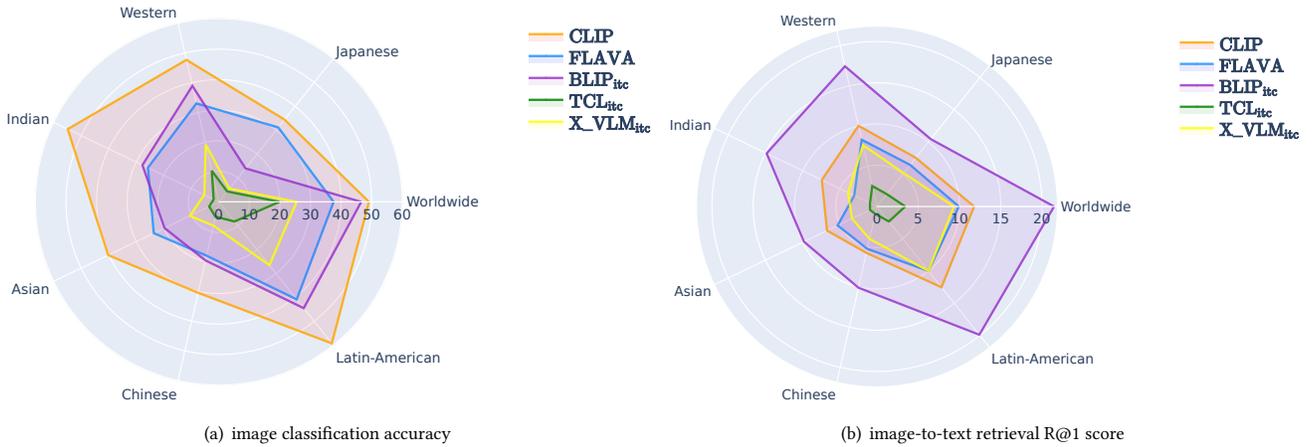[2]We try several prompts and this one is of the overall highest accuracy.

(a) image classification accuracy



(b) image-to-text retrieval R@1 score

**Figure 3: Radar chart of (a) Zero-shot image classification accuracy and (b) image-to-text retrieval R@1 score across regions.**

| Model | IR | | | TR | | |
|---|---|---|---|---|---|---|
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| CLIP | 7.51 | 20.77 | 30.32 | 9.28 | 23.98 | 33.63 |
| FLAVA | 9.82 | 25.85 | 36.05 | 7.60 | 21.28 | 31.04 |
| $BLIP_{itc}$ | 15.28 | 33.12 | 42.37 | 15.90 | 34.15 | 44.02 |
| $TCL_{itc}$ | 3.32 | 9.25 | 13.55 | 2.19 | 6.65 | 10.23 |
| $X\_VLM_{itc}$ | 6.83 | 17.58 | 25.18 | 6.75 | 17.68 | 24.72 |
| $BLIP_{itm}$ | **24.43** | **46.01** | **55.22** | **22.09** | **43.21** | **53.57** |
| $TCL_{itm}$ | 8.63 | 18.29 | 23.91 | 7.11 | 14.98 | 19.40 |
| $X\_VLM_{itm}$ | 20.10 | 37.51 | 44.94 | 15.06 | 31.41 | 40.36 |

**Table 2: Zero-shot cross-modal retrieval results on Food-500 Cap. The overall best result is bold-face.**
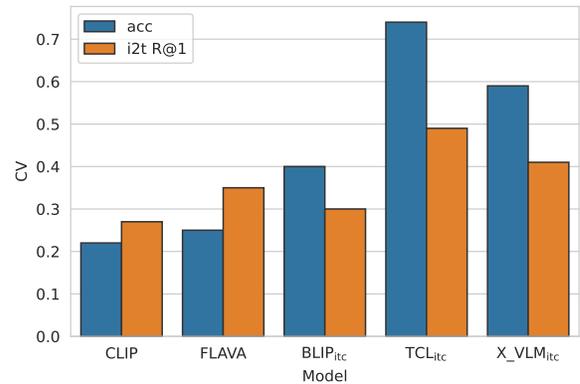


**Figure 4: Coefficient of variation (CV) of classification accuracies (blue) and image-to-text retrieval R@1 scores (orange) across different regions of several VLMs. CV is the ratio of the standard deviation to the mean, which measures the dispersion of a probability distribution. In this figure, the higher the CV value, the more unbalanced performance across regions.**

***Image-text Bidirectional Retrieval.*** Image-text bidirectional retrieval aims to retrieve images using textual queries (text-to-image retrieval) and converse (image-to-text retrieval). This task reveals models' ability to align the semantic space of vision and language. Though current VLMs [18, 33, 44, 55, 58] have achieved superior zero-shot performance on general domain datasets such as Flickr [56], MSCOCO [20], it is unknown whether they also perform well in a specific domain. Consequently, we conduct zero-shot image-text bidirectional retrieval using the food image-caption pair in Food-500 Cap. Like food classification. Similar to food classification, this task is also performed using ITC and ITM configurations. Finally, we report the top-1, 5, and 10 retrieval scores on our benchmark.

### 4.1.3 Results.

***Classification.*** Figure 2 displays the overall food classification results. None of the models achieve an accuracy above 50%. CLIP achieves over 40% accuracy on the zero-shot food classification task, which has the highest accuracy. BLIP and FLAVA also show competitive performance, while TCL and X_VLM exhibit a sizable gap. This phenomenon should be attributed to their relatively small pre-training datasets, which only contain 4M[3] and 16M image-text pairs, while that for FLAVA, BLIP, and CLIP are 70M, 129M, and 400M, respectively. We also obtain the following findings:

**VLMs fail to recognize certain food categories.** The accuracy varies greatly in different food categories. On one hand, VLMs perform nearly perfectly in some categories. For example, CLIP correctly classifies all images in *bandeja paisa* from *Latin-American*. However, On the other hand, VLMs recognize no images from some categories, such as *aburaage* and *doufunao*. The percentage of such categories for all VLMs is displayed in Figure 2 using the failure rate metric. We notice that even the best-performing CLIP fails in

---

[3]TCL does not release its checkpoint pre-trained on 16M image-text pairs.

| Setting | Model | Overall | *Chinese* | *Worldwide* | *Japanese* | *Western* | *Latin-American* | *Asian* | *Indian* |
|---------|-------|---------|-----------|-------------|------------|-----------|------------------|---------|----------|
| Accurate | $BLIP_{Dec}$ | 5.29 | 1.37 | 8.38 | 0 | 7.43 | 8.11 | 0.29 | 0 |
| | GIT | 2.98 | 1.33 | 4.49 | 0 | 3.98 | 6.00 | 0.20 | 0 |
| | OFA | **6.85** | **4.87** | **9.99** | **0.82** | **7.98** | **11.47** | **3.26** | **0.24** |
| Relaxed | $BLIP_{Dec}$ | **30.76** | 29.83 | **45.19** | 0.45 | **37.32** | 20.53 | 11.55 | 4.12 |
| | GIT | 26.27 | **31.70** | 36.43 | 0.73 | 30.56 | 16.74 | 11.06 | 5.30 |
| | OFA | 29.51 | 30.70 | 39.32 | **3.82** | 35.74 | **20.74** | **12.58** | **7.07** |

**Table 3: Semantic Label Accuracy (%) in the entire dataset (Overall) and different regions according to our taxonomy of food categories. *Accurate*: the generated caption exactly contains the whole food category label. *Relaxed*: the generated caption contains some word from the food label.**

| Model | Avg. Len | B@4 | M | R | C | CLIP-S |
|-------|----------|-----|---|---|---|--------|
| $BLIP_{Dec}$ | 13.80 | 2.61 | 8.71 | **20.33** | 13.62 | 0.70 |
| GIT | 20.89 | 2.00 | 8.81 | 16.78 | 9.92 | 0.70 |
| OFA | 20.45 | **2.64** | **9.14** | 17.89 | **14.01** | 0.71 |
| GT | 18.57 | - | - | - | - | **0.78** |

**Table 4: Results of image captioning on various metric. B@4, M, R C and CLIP-S represent BLEU@4, METEOR, ROUGE, CIDEr and CLIPScore respectively.**

nearly 10% of categories. And both X_VLM and TCL fail to identify over 40% categories.

**VLMs exhibit culinary culture bias in zero-shot food classification.** As illustrated in Figure 3 (a), all models exhibit consistency in their performance across different regions. These models can better identify food images from *Western* and *Latin-American* than others except for *Worldwide*. Besides the qualitative results, we furthermore report the coefficient of variation (CV) of scores in different regions in Figure 4, which reveals strong culinary culture bias in TCL and X_VLM. Such bias is probably inherited from the pre-training dataset, where food from some countries or regions, e.g. Japan, appears much less frequently than European and American food.

*Image-text Bidirectional Retrieval*. Table 2 shows the overall performance of the compared models on image-text bidirectional retrieval. Similar to our findings in food classification, the results of retrieval also demonstrate VLMs underperform in the food domain compared to the general domain and suffer from the region bias. To be specific:

**The overall performance is not satisfactory.** For the ITC configuration, BLIP gets the highest score on R@1, but it only reaches 15%. Other models get scores below 10%. In contrast to the classification task, where CLIP achieved the highest accuracy, it does not perform well in retrieval tasks. This implies that while CLIP is better at recognizing the general type of food, it has a weaker ability to distinguish food at a fine-grained level. Using the ITM configuration, this problem can be alleviated a bit. Under this setting, BLIP, TCL, and X_VLM obtain much higher scores on R@1, R@5, and R@10.

**All VLMs also suffer in certain categories in image-text retrieval.** For example, $BLIP_{itm}$'s image-to-text recall@1 score reaches 66.0 for *christmas cake*, a food from *Western*, but it hardly retrieves correct descriptions for *bon bon chicken* which is from *Chinese*. We further investigate the retrieval results in different regions. As shown in Figure 3 (b), we find that all VLMs perform relatively poorly in *Asian*, *Chinese*, *Indian* and *Japanese*. According to the quantitative results in Figure 4, CLIP also maintains the lowest CV, which suggests a relatively weak culinary culture bias. We speculate the reason to be its tremendous amount of pre-training data.

## 4.2 Image-to-text Generative Models

### 4.2.1 Evaluated Models.

*GIT* [51]. It is composed of one swin-like [21] vision transformer and one text decoder. During training, it uses the language modeling task to predict the associated caption given an image. When applied to downstream tasks, GIT first transforms them into text generation and then produces the answer word by word.

*OFA* [52]. It proposes a more generic encoder-decoder framework compared with GIT. It develops a unified multi-modal vocabulary, and both its encoder and decoder can process inputs from different modalities. Therefore, OFA can serve as a task-agnostic and modality-agnostic model. Both multi-modal and uni-modal tasks are combined during pre-training, which renders OFA superior performance on a wide range of tasks, such as image captioning, image generation, image classification, language understanding, etc.

*BLIP* [18]. Owing to its special architecture, BLIP can be regarded as an image-grounded text decoder as well. Hence we use BLIP in this part and denote it as $BLIP_{Dec}$. Detailed introduction can be referred to Section 4.1.1.

### 4.2.2 Evaluation Task.

*Image Captioning*. Image-to-text generative VLMs aggregate multiple tasks into a unified text generation task [18, 51, 52]. Thus we opt to leverage zero-shot image captioning as a means of probing these models. The objective of image captioning is to generate descriptive sentences for given images. We utilize multi-view metrics
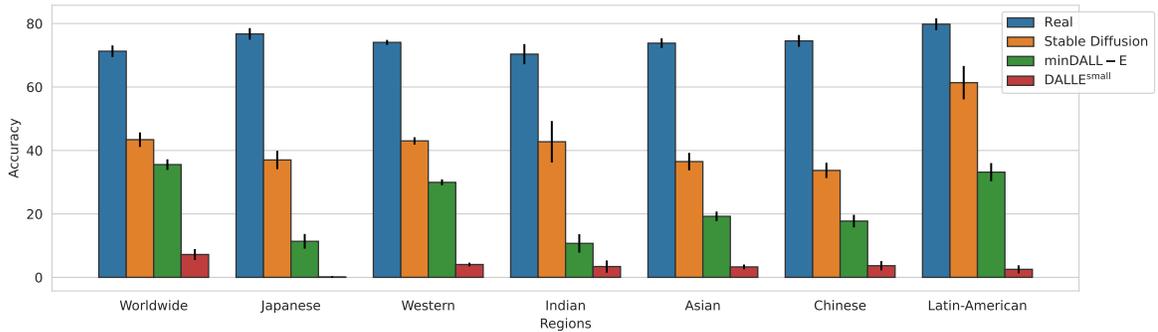
**Figure 5: Per-region accuracy of the classifier trained on real vs. synthetic images. We randomly split the dataset 10 times and reported the test accuracy on real images.**

to better exhibit the generative ability of VLMs. First, we calculate common-used image captioning metrics, including n-gram-based metrics, such as BLEU [19], METEOR [28], ROUGE [6] and CIDEr [49], and a semantic-based metric, CLIPScore [10].

Then, we assess the recognition ability of VLMs. Inspired by Semantic Object Accuracy (SOA) [12] in text-to-image generation evaluation, which evaluates whether the generated image contains the objects mentioned in the text, we check whether the generated captions contain the food labels of the images, and similarly define Semantic Label Accuracy (SLA):

$$ \text{SLA} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}(l^{(i)} \in c^{(i)}) $$

where $N$ is the number of images, $l^{(i)}$ and $c^{(i)}$ are the category label and generated caption of the $i$th image respectively.

### 4.2.3 Results.

We provide the captioning results in Table 4 including the n-gram-based metrics score and CLIPScore for all three models. All models have low scores on n-gram-based metrics.

**All image-to-text generative models hardly generate the fine-grained attributes**. From a perspective of metrics, the low scores on n-gram-based metrics indicate a clear literal mismatch between the generated and reference captions that contain many fine-grained attributes, including the shape, color, and ingredient. We provide some generated descriptions in Table 7. From the Table, we can observe that all VLMs neglect or misidentify a lot of food attributes. In the top example, models do not generate *eggplants*, *potatoes*, *green peppers* compared to the reference in our proposed dataset. Moreover, we find that OFA tends to generate hallucinations [4, 36], and GIT prefers to append meaningless words to the end of the sentence. As shown in the bottom example of Table 7, OFA generates "*ready to be served to the guests at the wedding reception*", and GIT adds "*yum*" and invalid punctuation at the end of the caption. These phenomena lead to the longer captions of GIT and OFA compared with BLIP$_{\text{Dec}}$ (Table 4), but there is no significant advantage in caption metrics. They also have lower CLIP-S scores than the ground truth, suggesting the weaker alignment of the generated captions with the images.

| Model | FID ($\downarrow$) | FID$_{\text{CLIP}}$ ($\downarrow$) | CAS ($\uparrow$) |
|---|---|---|---|
| Stable Diffusion | 25.74 | 10.92 | 41.45 ± 0.89 |
| minDALL-E | 28.64 | 15.46 | 26.60 ± 0.92 |
| DALLE$^{\text{small}}$ | 54.49 | 29.91 | 4.21 ± 0.54 |
| Real | 0 | 0 | 73.79 ± 0.69 |

**Table 5: FID, FID$_{\text{CLIP}}$ and CAS for different text-to-image generation models. For CAS, we randomly split the test data 10 times and report the mean and standard deviation.**

**These models hardly generate correct labels in descriptions** The low scores on metrics may reflect the generated sentence including incorrect labels. To verify this, we display the overall SLA and that in different regions in **??**, revealing that less than 10% captions exactly include whole category labels for all three models. If we relax the requirement and regard it as true if the generated caption contains some word in the label (*Relaxed* setting in **??**) rather than the whole label, SLA obtains significant improvement for all models, especially in *Chinese*, *Worldwide* and *Western*. This is because many category labels (e.g. *lentil soup*) from these regions contain common words like soup, which is easier for VLMs to generate.

**Food labels from different regions pose different levels of difficulty for VLMs to generate.** As shown in **??**, VLMs can hardly generate food category labels from specific regions, indicating a possible bias in culinary culture. In particular, for *Japanese* and *Indian*, BLIP$_{\text{Dec}}$ and GIT fail to generate food labels from these countries, and even the highest performing OFA only achieves 0.82 and 0.24 SLA, respectively.

## 4.3 Text-to-Image Generative Models

### 4.3.1 Evaluated Models.

***DALL-E*** *[34]*. It employs a decoder-only transformer that receives texts and images as a single stream. Given a text prompt, it first predicts the image tokens autoregressively, which have been

pre-defined in the codebook of a pre-trained discrete variational autoencoder (dVAE) [37]. Then the generated image tokens are fed to the decoder of the dVAE to synthesize images. Training on 250M image-text pairs from the internet, DALL-E can create plausible images for various sentences even in the zero-shot setting. Note that the checkpoint is unavailable. We choose two publicly released implements: minDALL-E [40], DALLE$^{small}$[4].

***Stable Diffusion** [38].* It leverages the prevalent diffusion model, which learns to reverse the process of adding noise to images. Unlike diffusion-based AI painters such as GLIDE [27], Imagen [41], Stable Diffusion uses latent diffusion model. The latent diffusion model operates in a compressed image space rather than the high-dimensional pixel space. Consequently, Stable Diffusion can generate high-resolution images from text descriptions with less computation consumption.

*4.3.2 Evaluation Task.*

***Image Synthesis**.* Following the default implementation, we adopt Stable Diffusion, minDALL-E, and DALLE$^{small}$ to synthesize images given food descriptions. To evaluate the overall image quality, we use Fréchet Inception Distance (FID) [11] and FID$_{CLIP}$ [17] scores, which measures the distance between the distributions of the real and synthetic images in the feature space of an ImageNet pre-trained Inception-v3 [46] and CLIP [33], respectively. Then we employ Classification Accuracy Score (CAS) [35] to assess to what extent the generated images manifest the categorical condition, where generated images are used to train a classifier which is then used to predict the of real images. To compute CAS, a classifier is first trained on the generated images, then used to predict labels of real images.

*4.3.3 Results.*

We report the FID and FID$_{CLIP}$ score and the CAS score in Table 5, which shows that all three models exhibit significant differences in their performance on image synthesis tasks. Through quantitative and qualitative analysis, we find the following issues:

**There is a significant gap between synthesis and real images.** From Table 5, all models have a significant gap compared to real images on FID, FID$_{CLIP}$, which measure the similarity between the synthesis and real images. Especially, the performance of DALLE$^{small}$ is worse than the other two models. Figure 8 shows some generated synthetic images. Through Figure 8, we find that images generated by DALLE$^{small}$ are unrealistic. In contrast, those generated by Stable Diffusion appear relatively more realistic and contain more caption content. To investigate whether the models can capture the main features of the category mentioned in the text, we further provide a quantitative evaluation of the synthetic images. Unlike metrics such as FID and FID$_{CLIP}$, CAS ignores some fringe features and is concerned with whether the generated images contain the necessary features to represent the class. We find that all models suffer a performance drop compared to real images, which indicates that text-to-image generative models might have difficulty capturing representative category features.

**Text-to-image Generative models also suffer from region imbalance issues, similar to the previous models.** For all regions, we observe a sizeable accuracy gap between the synthetic and real images. For example, Stable Diffusion's accuracy drops from nearly 20 to 40 across all regions. Furthermore, as shown in Figure 5, all models have some culinary culture bias. Specifically, the classifier trained with images generated by Stable Diffusion achieves particularly higher accuracy in *Latin-American* than other regions. And minDALL-E scores higher in *Worldwide*, *Western*, *Latin-American* than Asian countries. As for DALLE$^{small}$, it fails to recognize almost all food images from *Japanese*. In contrast, using real images to train the classifier results in relatively balanced accuracy across regions, which might be because those text-to-image generative models are trained on a biased dataset, which contains fewer traditional food types from certain regions such as *Chinese* and *Japanese*.

## 5 CONCLUSION

In this work, we introduce Food-500 Cap, a new vision-language benchmark in the food domain. By in-house labeling, Food-500 Cap not only provides each image with a fine-grained visual content description but also labels a novel taxonomy that divides food categories into their geographic origins, which aids in studying different culinary cultures. We adopt four vision-language tasks in the zero-shot setting, including food classification, image-text bidirectional retrieval, image captioning, and image synthesis, and evaluate nine VLMs of three different architectures on our proposed benchmark. Experiments reveal VLMs' limitations in the food domain and their bias against culinary culture. We hope that our proposed benchmark will promote the study of multi-modal food computing and our findings will provide insights into the deployment and application of VLMs in the food domain.

## ACKNOWLEDGMENTS

---

[4]https://github.com/lucidrains/DALLE-pytorch

# REFERENCES

[1] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 2014. Food-101 – Mining Discriminative Components with Random Forests. In *European Conference on Computer Vision*.

[2] Jingjing Chen and Chong-Wah Ngo. 2016. Deep-based Ingredient Recognition for Cooking Recipe Retrieval. In *Proceedings of the 2016 ACM Conference on Multimedia Conference, MM 2016, Amsterdam, The Netherlands, October 15-19, 2016*, Alan Hanjalic, Cees Snoek, Marcel Worring, Dick C. A. Bulterman, Benoit Huet, Aisling Kelliher, Yiannis Kompatsiaris, and Jin Li (Eds.). ACM, 32–41. https://doi.org/10.1145/2964284.2964315

[3] Xin Chen, Hua Zhou, and Liang Diao. 2017. ChineseFoodNet: A large-scale Image Dataset for Chinese Food Recognition. *CoRR* abs/1705.02743 (2017). arXiv:1705.02743 http://arxiv.org/abs/1705.02743

[4] Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. 2022. Plausible May Not Be Faithful: Probing Object Hallucination in Vision-Language Pretraining. *CoRR* abs/2210.07688 (2022). https://doi.org/10.48550/arXiv.2210.07688 arXiv:2210.07688

[5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*.

[6] Michael J. Denkowski and Alon Lavie. 2014. Meteor Universal: Language Specific Translation Evaluation for Any Target Language. In *Proceedings of the Ninth Workshop on Statistical Machine Translation, WMT@ACL 2014, June 26-27, 2014, Baltimore, Maryland, USA*. The Association for Computer Linguistics, 376–380. https://doi.org/10.3115/v1/w14-3348

[7] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. 2010. The Pascal Visual Object Classes (VOC) Challenge. *Int. J. Comput. Vis.* 88, 2 (2010), 303–338. https://doi.org/10.1007/s11263-009-0275-4

[8] Tejas Gokhale, Hamid Palangi, Besmira Nushi, Vibhav Vineet, Eric Horvitz, Ece Kamar, Chitta Baral, and Yezhou Yang. 2022. Benchmarking Spatial Relationships in Text-to-Image Generation. *CoRR* abs/2212.10015 (2022). https://doi.org/10.48550/arXiv.2212.10015 arXiv:2212.10015

[9] Lisa Anne Hendricks and Aida Nematzadeh. 2021. Probing Image-Language Transformers for Verb Understanding. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021 (Findings of ACL, Vol. ACL/IJCNLP 2021)*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, 3635–3644. https://doi.org/10.18653/v1/2021.findings-acl.318

[10] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. CLIPScore: A Reference-free Evaluation Metric for Image Captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 7514–7528. https://doi.org/10.18653/v1/2021.emnlp-main.595

[11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 6626–6637. https://proceedings.neurips.cc/paper/2017/hash/8a1d694707eb0fefe65871369074926d-Abstract.html

[12] Tobias Hinz, Stefan Heinrich, and Stefan Wermter. 2022. Semantic Object Accuracy for Generative Text-to-Image Synthesis. *IEEE Trans. Pattern Anal. Mach. Intell.* 44, 3 (2022), 1552–1565. https://doi.org/10.1109/TPAMI.2020.3021209

[13] Anya Ji, Noriyuki Kojima, Noah Rush, Alane Suhr, Wai Keen Vong, Robert D. Hawkins, and Yoav Artzi. 2022. Abstract Visual Reasoning with Tangram Shapes. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, 582–601. https://aclanthology.org/2022.emnlp-main.38

[14] Yoshiyuki Kawano and Keiji Yanai. 2014. Automatic Expansion of a Food Image Dataset Leveraging Existing Categories with Domain Adaptation. In *Computer Vision - ECCV 2014 Workshops - Zurich, Switzerland, September 6-7 and 12, 2014, Proceedings, Part III (Lecture Notes in Computer Science, Vol. 8927)*, Lourdes Agapito, Michael M. Bronstein, and Carsten Rother (Eds.). Springer, 3–17. https://doi.org/10.1007/978-3-319-16199-0_1

[15] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *Int. J. Comput. Vis.* 123, 1 (2017), 32–73. https://doi.org/10.1007/s11263-016-0981-7

[16] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).

[17] Tuomas Kynkäänniemi, Tero Karras, Miika Aittala, Timo Aila, and Jaakko Lehtinen. 2022. The Role of ImageNet Classes in Fr\`echet Inception Distance. *arXiv preprint arXiv:2203.06026* (2022).

[18] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. 2022. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (Eds.). PMLR, 12888–12900. https://proceedings.mlr.press/v162/li22n.html

[19] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.

[20] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common Objects in Context. In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V (Lecture Notes in Computer Science, Vol. 8693)*, David J. Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (Eds.). Springer, 740–755. https://doi.org/10.1007/978-3-319-10602-1_48

[21] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 9992–10002. https://doi.org/10.1109/ICCV48922.2021.00986

[22] Javier Marín, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2018. Recipe1M: A Dataset for Learning Cross-Modal Embeddings for Cooking Recipes and Food Images. *CoRR* abs/1810.06553 (2018). arXiv:1810.06553 http://arxiv.org/abs/1810.06553

[23] Javier Marín, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2021. Recipe1M+: A Dataset for Learning Cross-Modal Embeddings for Cooking Recipes and Food Images. *IEEE Trans. Pattern Anal. Mach. Intell.* 43, 1 (2021), 187–203. https://doi.org/10.1109/TPAMI.2019.2927476

[24] Weiqing Min, Bing-Kun Bao, Shuhuan Mei, Yaohui Zhu, Yong Rui, and Shuqiang Jiang. 2018. You Are What You Eat: Exploring Rich Recipe Information for Cross-Region Food Analysis. *IEEE Trans. Multim.* 20, 4 (2018), 950–964. https://doi.org/10.1109/TMM.2017.2759499

[25] Weiqing Min, Shuqiang Jiang, Linhu Liu, Yong Rui, and Ramesh C. Jain. 2019. A Survey on Food Computing. *ACM Comput. Surv.* 52, 5 (2019), 92:1–92:36. https://doi.org/10.1145/3329168

[26] Weiqing Min, Linhu Liu, Zhiling Wang, Zhengdong Luo, Xiaoming Wei, Xiaolin Wei, and Shuqiang Jiang. 2020. ISIA Food-500: A Dataset for Large-Scale Food Recognition via Stacked Global-Local Attention Network. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*, Chang Wen Chen, Rita Cucchiara, Xian-Sheng Hua, Guo-Jun Qi, Elisa Ricci, Zhengyou Zhang, and Roger Zimmermann (Eds.). ACM, 393–401. https://doi.org/10.1145/3394171.3414031

[27] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2022. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (Eds.). PMLR, 16784–16804. https://proceedings.mlr.press/v162/nichol22a.html

[28] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*. ACL, 311–318. https://doi.org/10.3115/1073083.1073135

[29] Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. 2022. VALSE: A Task-Independent Benchmark for Vision and Language Models Centered on Linguistic Phenomena. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, 8253–8280. https://doi.org/10.18653/v1/2022.acl-long.567

[30] Dong Huk Park, Samaneh Azadi, Xihui Liu, Trevor Darrell, and Anna Rohrbach. 2021. Benchmark for Compositional Text-to-Image Synthesis. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, Joaquin Vanschoren and Sai-Kit Yeung (Eds.). https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/0a09c8844ba8f0936c20bd791130d6b6-Abstract-round1.html

[31] Jianing Qiu, Frank Po Wen Lo, Yingnan Sun, Siyao Wang, and Benny Lo. 2019. Mining Discriminative Food Regions for Accurate Food Recognition. In *30th British Machine Vision Conference 2019, BMVC 2019, Cardiff, UK, September 9-12, 2019*. BMVA Press, 158. https://bmvc2019.org/wp-content/uploads/papers/0839-paper.pdf

[32] Jielin Qiu, Yi Zhu, Xingjian Shi, Florian Wenzel, Zhiqiang Tang, Ding Zhao, Bo Li, and Mu Li. 2022. Are Multimodal Models Robust to Image and Text Perturbations? *CoRR* abs/2212.08044 (2022). https://doi.org/10.48550/arXiv.2212.

08044 arXiv:2212.08044

[33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. http://proceedings.mlr.press/v139/radford21a.html

[34] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-Shot Text-to-Image Generation. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8821–8831. http://proceedings.mlr.press/v139/ramesh21a.html

[35] Suman V. Ravuri and Oriol Vinyals. 2019. Classification Accuracy Score for Conditional Generative Models. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 12247–12258. https://proceedings.neurips.cc/paper/2019/hash/fcf55a303b71b84d326fb1d06e332a26-Abstract.html

[36] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object Hallucination in Image Captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (Eds.). Association for Computational Linguistics, 4035–4045. https://doi.org/10.18653/v1/d18-1437

[37] Jason Tyler Rolfe. 2017. Discrete Variational Autoencoders. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net. https://openreview.net/forum?id=ryMxXPFex

[38] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 10674–10685. https://doi.org/10.1109/CVPR52688.2022.01042

[39] Philipp J. Rösch and Jindrich Libovický. 2022. Probing the Role of Positional Information in Vision-Language Models. In *Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, Marine Carpuat, Marie-Catherine de Marneffe, and Iván Vladimir Meza Ruíz (Eds.). Association for Computational Linguistics, 1031–1041. https://doi.org/10.18653/v1/2022.findings-naacl.77

[40] Chiheon Kim Doyup Lee Saehoon Kim, Sanghun Cho and Woonhyuk Baek. 2021. minDALL-E on Conceptual Captions. https://github.com/kakaobrain/minDALL-E.

[41] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. *CoRR* abs/2205.11487 (2022). https://doi.org/10.48550/arXiv.2205.11487 arXiv:2205.11487

[42] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. LAION-400M: Open Dataset of CLIP-Filtered 400 Million Image-Text Pairs. *CoRR* abs/2111.02114 (2021). arXiv:2111.02114 https://arxiv.org/abs/2111.02114

[43] Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, 4222–4235. https://doi.org/10.18653/v1/2020.emnlp-main.346

[44] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. FLAVA: A Foundational Language And Vision Alignment Model. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 15617–15629. https://doi.org/10.1109/CVPR52688.2022.01519

[45] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. https://openreview.net/forum?id=SygXPaEYvH

[46] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 2818–2826. https://doi.org/10.1109/CVPR.2016.308

[47] Hao Tan and Mohit Bansal. 2019. LXMERT: Learning Cross-Modality Encoder Representations from Transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (Eds.). Association for Computational Linguistics, 5099–5110. https://doi.org/10.18653/v1/D19-1514

[48] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. 2022. Winoground: Probing Vision and Language Models for Visio-Linguistic Compositionality. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 5228–5238. https://doi.org/10.1109/CVPR52688.2022.00517

[49] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. CIDEr: Consensus-based image description evaluation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. IEEE Computer Society, 4566–4575. https://doi.org/10.1109/CVPR.2015.7299087

[50] Felix Vogel, Nina Shvetsova, Leonid Karlinsky, and Hilde Kuehne. 2022. VL-Taboo: An Analysis of Attribute-based Zero-shot Capabilities of Vision-Language Models. *CoRR* abs/2209.06103 (2022). https://doi.org/10.48550/arXiv.2209.06103 arXiv:2209.06103

[51] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. 2022. GIT: A Generative Image-to-text Transformer for Vision and Language. *CoRR* abs/2205.14100 (2022). https://doi.org/10.48550/arXiv.2205.14100 arXiv:2205.14100

[52] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. OFA: Unifying Architectures, Tasks, and Modalities Through a Simple Sequence-to-Sequence Learning Framework. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (Eds.). PMLR, 23318–23340. https://proceedings.mlr.press/v162/wang22al.html

[53] Xin Wang, Devinder Kumar, Nicolas Thome, Matthieu Cord, and Frédéric Precioso. 2015. Recipe recognition with large multimodal food dataset. In *2015 IEEE International Conference on Multimedia & Expo Workshops, ICME Workshops 2015, Turin, Italy, June 29 - July 3, 2015*. IEEE Computer Society, 1–6. https://doi.org/10.1109/ICMEW.2015.7169757

[54] Xiongwei Wu, Xin Fu, Ying Liu, Ee-Peng Lim, Steven C. H. Hoi, and Qianru Sun. 2021. A Large-Scale Benchmark for Food Image Segmentation. In *MM '21: ACM Multimedia Conference, Virtual Event, China, October 20 - 24, 2021*, Heng Tao Shen, Yueting Zhuang, John R. Smith, Yang Yang, Pablo César, Florian Metze, and Balakrishnan Prabhakaran (Eds.). ACM, 506–515. https://doi.org/10.1145/3474085.3475201

[55] Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath Chanda, Liqun Chen, Belinda Zeng, Trishul Chilimbi, and Junzhou Huang. 2022. Vision-Language Pre-Training with Triple Contrastive Learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 15650–15659. https://doi.org/10.1109/CVPR52688.2022.01522

[56] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Trans. Assoc. Comput. Linguistics* 2 (2014), 67–78. https://doi.org/10.1162/tacl_a_00166

[57] Mert Yüksekgönül, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2022. When and why vision-language models behave like bags-of-words, and what to do about it? *CoRR* abs/2210.01936 (2022). https://doi.org/10.48550/arXiv.2210.01936 arXiv:2210.01936

[58] Yan Zeng, Xinsong Zhang, and Hang Li. 2022. Multi-Grained Vision Language Pre-Training: Aligning Texts with Visual Concepts. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (Eds.). PMLR, 25994–26009. https://proceedings.mlr.press/v162/zeng22c.html

[59] Tiancheng Zhao, Tianqi Zhang, Mingwei Zhu, Haozhan Shen, Kyusong Lee, Xiaopeng Lu, and Jianwei Yin. 2022. VL-CheckList: Evaluating Pre-trained Vision-Language Models with Objects, Attributes and Relations. *CoRR* abs/2207.00221 (2022). https://doi.org/10.48550/arXiv.2207.00221 arXiv:2207.00221

[60] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to Prompt for Vision-Language Models. *Int. J. Comput. Vis.* 130, 9 (2022), 2337–2348. https://doi.org/10.1007/s11263-022-01653-1

[61] Wangchunshu Zhou, Yan Zeng, Shizhe Diao, and Xinsong Zhang. 2022. VLUE: A Multi-Task Benchmark for Evaluating Vision-Language Models. *CoRR* abs/2205.15237 (2022). https://doi.org/10.48550/arXiv.2205.15237 arXiv:2205.15237

# A    PROCESS OF ANNOTATING GEOGRAPHIC ORIGINS

For each food category, we resort to its Wikipedia entry. We can find their places of origin of most food categories. We display two examples in Figure 6, according to which we assign *Wonton noodles* to region *Chinese* and *Tsukemono* to region *Japanese*. However, some food categories have unknown origins, such as the *Carrot salad* shown in Figure 7. We assign these food categories to *Worldwide* in our taxonomy.



(a) Wonton noodles          (b) Tsukemono

**Figure 6: Annotating geographic origins in the case that the clear origin of the food are given. (a) *Wonton noodles*, where the place of origin is directly provided. (b) *Tsukemono*, where the entry lacks additional citations for verification and the place of origin is not provided, but we can still find the origin in the article.**
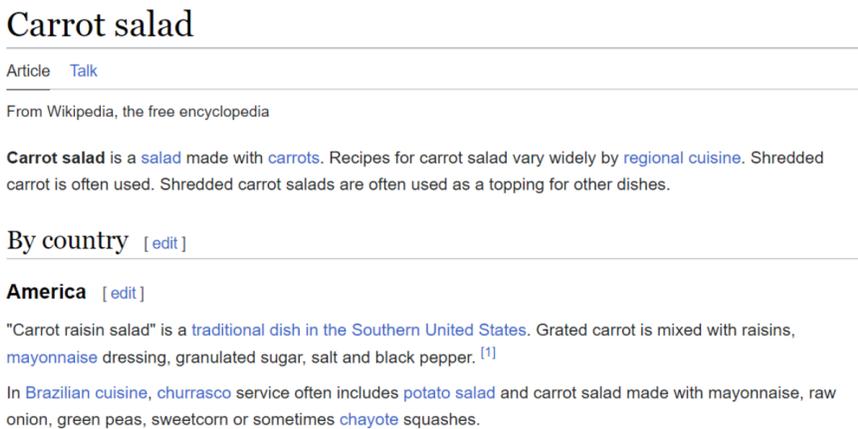


**Figure 7: Annotating geographic origins when some food categories have unknown places of origin.**

| Region | #Categories | Covered Countries |
|---|---|---|
| *Worldwide* | 90 | unknown original regions |
| *Western* | 216 | Europe & America & Canada |
| *Latin-American* | 19 | Latin American countries |
| *Chinese* | 60 | China |
| *Japanese* | 22 | Japan |
| *Indian* | 17 | India |
| *Asian* | 70 | Asia except for China, Japan, India |

**Table 6: Region distribution of food categories in Food-500 Cap.**

| Image | Captions |
|---|---|
|  | **GT** di san xian sauteed with soft eggplants, potatoes and slices of green peppers in a bowl which placed on a red napkin. |
|  | **BLIP$_{Dec}$** a bowl of food with chops and chops on a blue and white floral tablecloth. **OFA** a bowl of food with chopsticks on a blue and white tablecloth with white daisies in the background. **GIT** chicken in a bowl with chopsticks on a red and blue placemat. yum!!! yum yum... |
|  | **GT** a grilled piece of bone-in pork knuckle served with yellow sauerkraut, and decorated with rosemary. |
|  | **BLIP$_{Dec}$** a plate of carrots on a table. **OFA** a plate of sweet potato fries with a drizzle of olive oil on top, ready to be served to the guests at the wedding reception. **GIT** cooked carrots in a white plate with a brown sauce... yum!!! : - ) ( : - - - ) |

**Table 7: Examples of generated captions from three image-to-text generative models compared to ground truth (GT).**



**Figure 8: Some examples of real images and synthesis images from text-to-image generative models**