

---

# Turath-150K: Image Database of Arab Heritage

---

**Dani Kiyasseh**  
Department of Engineering Science  
University of Oxford  
Oxford, UK  
dani.kiyasseh@eng.ox.ac.uk

**Rasheed El-Bouri**  
Department of Engineering Science  
University of Oxford  
Oxford, UK  
rasheed.el-bouri@eng.ox.ac.uk

## Abstract

1 Large-scale image databases remain largely biased towards objects and activities  
2 encountered in a select few cultures. This absence of culturally-diverse images,  
3 which we refer to as the “hidden tail”, limits the applicability of pre-trained neural  
4 networks and inadvertently excludes researchers from under-represented regions.  
5 To begin remedying this issue, we curate Turath-150K, a database of images of  
6 the Arab world that reflect objects, activities, and scenarios commonly found there.  
7 In the process, we introduce three benchmark databases, Turath Standard, Art,  
8 and UNESCO, specialised subsets of the Turath dataset. After demonstrating  
9 the limitations of existing networks pre-trained on ImageNet when deployed on  
10 such benchmarks, we train and evaluate several networks on the task of image  
11 classification. As a consequence of Turath, we hope to engage machine learning  
12 researchers in under-represented regions, and to inspire the release of additional  
13 culture-focused databases. The database can be accessed here: [danikiyasseh.](https://github.com/danikiyasseh/Turath)  
14 [github.io/Turath](https://github.com/danikiyasseh/Turath).

## 15 1 Introduction

16 Deep neural networks have exhibited great success in performing various computer vision tasks,  
17 such as image classification [1], object detection [2], and segmentation [3]. One of the key factors  
18 and driving forces behind the success of such networks is access to large-scale, annotated datasets  
19 that consist of samples that are mostly representative of the underlying data distribution. To that  
20 end, publicly-available datasets, such as ImageNet [4], SUN [5], and Places [6], attempt to capture  
21 a diverse set of images that are reflective of objects and scenarios encountered “in the wild”. Such  
22 images typically belong to categories guided by the WordNet hierarchy [7] and which are diversified  
23 by incorporating various adjectives into search queries (e.g., night, foggy, etc.)

24 Despite these efforts, existing databases remain largely biased towards objects, activities, and sce-  
25 narios commonly encountered in a small subset of cultures [8], define “diversity” narrowly, and  
26 do not account for the long-tail of image categories that are common in other cultures. For ex-  
27 ample, items and activities common in other parts of the world, such as those in the Arab world,  
28 are under-represented, if at all, in existing image databases [9]. Examples include traditional daily  
29 clothing items, such as the “thobe”, and sporting activities, such as falconry. We refer to these  
30 under-represented categories, in which *no* images are available in existing databases, as the “hidden  
31 tail”. This is analogous to the “long tail” of image categories, in which *few* images are available, that  
32 the machine learning community has dedicated substantial effort to better representing.

33 Such an exclusion of culturally-diverse images has a technical, societal, and ethical impact on  
34 the machine learning community. From a technical perspective, the absence of diverse images in  
35 existing databases violates the assumption that samples are from “the wild” and representative of  
36 the underlying data distribution. By evaluating networks on such narrow samples, their performance  
37 tends to be an over-estimate. Moreover, culturally-diverse image categories are effectively out-of-

38 distribution (OOD) samples notorious for degrading the performance of trained networks [10], a  
39 phenomenon shown to be more prominent when transferring across geographical regions [11]. On  
40 a societal level, pre-trained networks are less likely to be of direct value to researchers residing  
41 in, or operating with, under-represented communities. This is driven by the poor performance of  
42 such networks on OOD samples, which is a direct consequence of the cultural bias inherent in the  
43 datasets used to train such networks. With this imbalance in the applicability of networks across  
44 cultures, under-represented communities are unlikely to capture the benefits of computer vision-based  
45 advancements. Furthermore, the machine learning community’s lack of exposure to data from diverse  
46 cultures suggests that researchers have less of an opportunity to learn about such cultures. Such  
47 dataset-based learning, the acquisition of skills and knowledge via datasets, has been evident with,  
48 for example, the Caltech-UCSD Birds 200 database [12] and ornithology. On an ethical level, the  
49 absence of data to which researchers can relate implicitly excludes these researchers from more  
50 actively engaging with the machine learning community. As such, it is to the advantage of the  
51 community to build the infrastructure that incentivizes the involvement of practitioners from a more  
52 diverse background in machine learning.

53 In this work, we aim to increase the cultural diversity of images that are available for training neural  
54 networks. Hence, we present the Turath-150K<sup>1</sup> database, a large-scale dataset of images depicting  
55 objects, activities, and scenarios that are rooted in the Arab world and culture. We chose this culture  
56 as an example, particularly due to its under-representation in existing publicly-available datasets,  
57 and hope other researchers follow suit with publishing datasets depicting cultures from around the  
58 globe. Specifically, our contributions are the following: (1) we build a large-scale database of images,  
59 entitled Turath-150K, the first of its kind that centres around life in the Arab world. For benchmarking  
60 purposes, we split the database into three distinct subsets; Turath-Standard, Turath-Art (focusing  
61 on art from the Arab world), and Turath-UNESCO (focusing on heritage sites located in the Arab  
62 world). (2) We shed light on the limitations of deep neural networks pre-trained on ImageNet by  
63 showing that they are unable to deal with the out-of-distribution samples of the Turath database.  
64 (3) We evaluate various networks on the Turath benchmark databases and demonstrate their image  
65 classification performance on both high and low-level categories.

## 66 2 Related work

67 There exists a multitude of publicly-available image databases that have been exploited for the training  
68 of deep neural networks. We outline several that we believe are most similar to our work and also  
69 elucidate how our database, Turath, differs significantly in motivation, scope, and content.

70 **Scene recognition databases** The task of scene recognition involves identifying scenes based on  
71 images. To facilitate achieving this task, the SUN397 database [5] was designed to contain 100K  
72 images of 397 scenes. The vast majority of these scene categories are motivated by the WordNet  
73 hierarchy [7]. Similarly, the Places database [6] was designed to contain 2.5 million images of 365  
74 high-level scenes, such as coffee-shop, nursery, and train station. Although extensive in terms of the  
75 number of samples, the scene categories lack the granularity that we offer and do not trivially extend  
76 to the Arab world. Moreover, Turath is not exclusively limited to scenes (see Sec. 3) and goes beyond  
77 the narrow WordNet hierarchy by explicitly accounting for entities in the Arab world.

78 **Object classification databases** The task of object classification focuses on identifying object(s)  
79 in an image. To propel research on this front, the Caltech 256 database [13] was designed to contain  
80 30K images of everyday objects, such as cameras and laptops. The COCO database [14] is much  
81 more extensive with 330K images corresponding to 80 object categories and consisting of multiple  
82 annotations, including segmentation maps at various levels of detail. Nonetheless, such databases  
83 differ in motivation, scope, and content from our database. In order to increase the cultural diversity  
84 of datasets, we turn our attention to objects, activities, and scenarios commonly found in the Arab  
85 world. Moreover, our image annotations are not only absent from existing databases but also offer a  
86 finer resolution of class label. We explain this in further depth in the next section.

87 **Out-of-distribution databases** Researchers have adopted various approaches to handle the gen-  
88 eralization of their models to out-of-distribution samples. These approaches can be split according

---

<sup>1</sup>Turath roughly means heritage in Arabic

89 to whether they are implemented during training or evaluation, with the latter being more relevant  
90 to our work. For example, ImageNet-R [11] is an evaluation database of 30K images, spanning  
91 200 ImageNet categories, rendered in different styles and textures. While their approach augments  
92 existing ImageNet categories, our database includes image samples from categories *beyond* the  
93 ImageNet-1K. ImageNet-O [10] is an evaluation database that claims to reflect label distribution  
94 shift, yet still only comprises images from 200 categories in ImageNet-1K. Whereas ImageNet-O  
95 is focused on evaluating out-of-distribution detectors, the Turath database is primarily focused on  
96 increasing the representation of image categories that are under-represented in ImageNet.

### 97 **3 Design and construction of the Turath database**

98 In light of our emphasis on increasing the cultural diversity of images, we aimed to construct a  
99 database that satisfies the following desiderata:

- 100 1. **Heritage** - Categories of images must be specific to the cultures of the Arab world; we reiterate  
101 that although our particular choice of culture stems from its under-representation in existing  
102 publicly-available databases, it is simply an example. There remains a multitude of rich cultures  
103 that are under-represented and we hope other researchers eventually publish such culture-specific  
104 databases, be they in the form of images, audio, or video.
- 105 2. **Quantity** - Each category must contain a sufficient number of images to facilitate learning;  
106 although the term “sufficient” is nebulous and category-dependent, existing databases have demon-  
107 strated success with at least 50 images per category. We quadruple that amount and aim for at least  
108 200 images per category.
- 109 3. **Real World** - Images in each category must reflect those commonly encountered “in the wild”;  
110 networks trained on image databases have a number of applications but they are, arguably, most  
111 useful when applied in the real world to challenges afflicting stakeholders from patients to farmers.  
112 To that end, we aim to collect natural RGB images.

113 The construction of the Turath database consisted of three main stages. We first defined keywords  
114 to guide the download of images from web-based search engines. We then used these keywords to  
115 assign images an annotation. Lastly, and as a form of noise reduction, we trained several classifiers  
116 to distinguish between categories and removed images that were likely to be associated with the  
117 incorrect annotation. We now describe these stages in more depth.

118 **Stage 1: Defining keywords and downloading the images** Existing image databases such as  
119 ImageNet and Places were created by performing query-based searches using online search engines.  
120 In this setting, the choice of queries determines the type and quality of images that are retrieved. In  
121 our context, and in contrast to the aforementioned work, the WordNet hierarchy [7] did not satisfy  
122 our outlined desiderata. This is primarily because WordNet was not designed for the Arab world  
123 and thus does not contain categories that are directly relevant for our purposes. Although an Arabic  
124 WordNet [15] does exist, it is unable to capture the cultural focus and the *micro* categories (described  
125 next) that we are searching for.

126 Given our emphasis on the Arab world as an example, we conducted query-based searches of entities  
127 engrossed in the diverse cultures of the region. This ranged from categories of images with a low level  
128 of detail, such as cities and architecture, to those with a high level of detail, such as traditional food  
129 and clothing. Each of these *macro* categories are formed by grouping several *micro* categories. For  
130 example, the *macro* category of Cities comprises 25+ *micro* categories of images from specific cities  
131 in the Arab world, e.g., Damascus, Cairo, and Casablanca. To emphasize the under-representation of  
132 images of these cities in existing databases, we note that the largest image database of cities, World  
133 Cities [16], with 2.25M images, covers a single city (Dubai) in the Arab world. In Fig. 1, we present  
134 image samples from three macro categories, Dates, Architecture, and Souq, each containing four  
135 *micro* categories.

136 In addition to retrieving images from the categories mentioned above, we dedicate time and effort  
137 to curating two additional *macro* categories that comprise a large number of *micro* categories.  
138 Specifically, these revolve around Arab Art and United Nations Educational, Scientific and Cultural  
139 Organization (UNESCO) sites. When retrieving images that belong to the Arab Art category, we  
140 followed the same strategy of query-based searches. However, given the breadth of this field and to

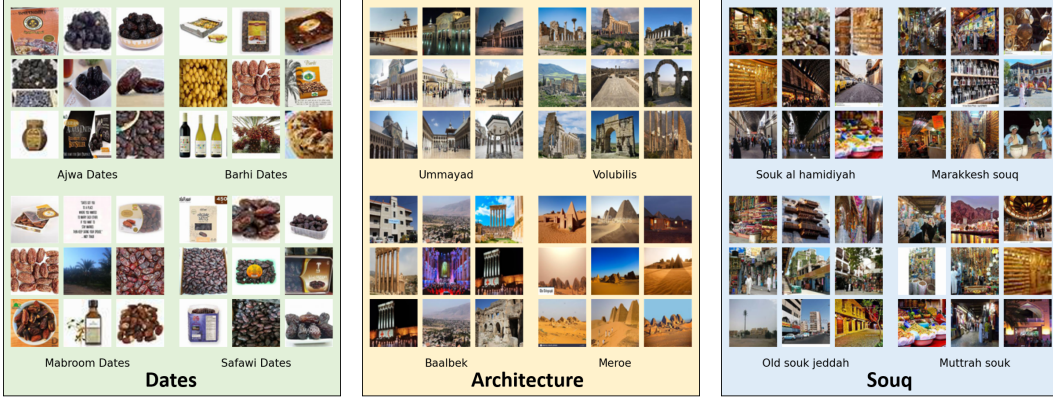


Figure 1: **Images samples from a subset of categories available in Turath.** Four *micro* categories are shown for each of the three *macro* categories, Dates, Architecture, and Souq. The image categories range from objects with low-level details, such as dates, to locations with high-level details, such as architecture.

141 keep the task of downloading images tractable and organized, our search queries were based on artists’  
 142 names. To that end, we identified 425 names available on the Barjeel Art Foundation website<sup>2</sup>. As  
 143 for the UNESCO category, our search queries were based on the names of 88 recognized UNESCO  
 144 sites in the Arab world<sup>3</sup>.

145 **Stage 2: Labelling the images using keywords** Each image in the Turath database has two image-  
 146 level annotations; a *micro* label and a *macro* label. To assign downloaded images to *micro* categories,  
 147 we follow the strategy proposed by Marin *et al.* [17] where each category is defined by the query  
 148 used to search for those images. Similar to their conclusions, we also find that such an approach  
 149 leads to relatively high quality images that are relevant to the search query. We then grouped *micro*  
 150 categories with similar themes into *macro* categories. As an example, we grouped seven types of  
 151 dates (*micro*) into a single Dates category (*macro*).

152 **Stage 3: Filtering the images with classifier-based labelling** Despite our effort to conduct  
 153 searches using queries that are unambiguous and descriptive, upon further inspection, we found that  
 154 certain categories contained images that were irrelevant. This was most prominent amongst images  
 155 that belonged to artists. For example, the query inji efflatoun art returned art pieces associated with  
 156 the artist Inji Efflatoun, as desired, but also images of the artist herself.

157 To remedy this situation, we exploited the prior knowledge that out-of-distribution (OOD) image  
 158 samples are likely to be of artists’ faces. Therefore, given our emphasis on retaining images of  
 159 art pieces, we designed a binary classifier that distinguished between images of art and those of  
 160 faces. To train such a classifier, we needed images with relatively high quality labels. For those in  
 161 the “art” domain, we grouped all the categories in ImageNet-R [11], which comprises images from  
 162 ImageNet rendered artistically, into a single category. For those in the “faces” domain, we exploited  
 163 images from the LFW database [18], which comprises 13K images of faces, and grouped them into a  
 164 single category. After training this classifier, we performed inference on *our* set of artistic images.  
 165 Given that the majority of images are those of art pieces, we would expect the distribution of output  
 166 probabilities to be bi-modal and skewed towards the value zero (i.e., corresponding to art images).  
 167 This is indeed what we find empirically, as shown in Fig. 2. Upon manual inspection of the images,  
 168 we chose a threshold value of 0.1, whereby approximately 26.1% of image samples believed to have  
 169 been of art are instead identified as a face. These 27,302 images are removed from the database.

170 Detecting OOD images of human faces exploited the implicit bias that human faces comprised the  
 171 majority of the OOD images. However, not all OOD images contain human faces. To investigate this,  
 172 we explored more general approaches involving one-class SVMs [19], deep autoencoding GMMs  
 173 [20], adversarial networks [21], geometric transformations [22] and self-supervised classification

<sup>2</sup><https://www.barjeelartfoundation.org/>

<sup>3</sup><https://whc.unesco.org/en/list/???order=region>

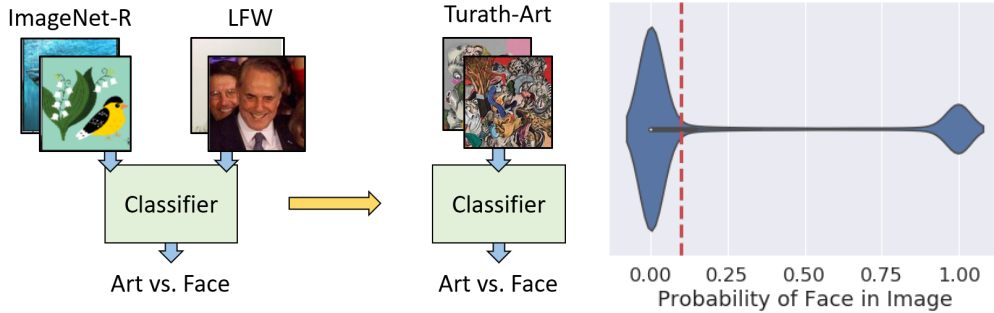


Figure 2: **Pipeline for cleaning data in Turath database.** (Left) Classifier-based cleaning of data. We trained a binary classifier to distinguish between images of art (ImageNet-R) and faces (LFW) and deployed it on Turath-Art. (Right) Distribution of probabilities output by binary classifier deployed on all images of Turath Art. We found that, when a threshold of 0.1 is chosen, approximately 26.1% of images are identified as a face.

174 networks [23]. We empirically found that although this self-supervised approach was preferable to  
 175 the remaining methods, it was still unable to reliably identify OOD samples.

#### 176 4 Turath benchmark databases

177 The Turath database comprises three specialized subsets of data that contain images from mutually-  
 178 exclusive categories. Hereafter, these subsets will be referred to as Turath Standard, Turath Art, and  
 179 Turath UNESCO, respectively, and, in this section, will be described in depth. We chose to separate  
 180 the database along these dimensions to account for the different resolution of the categories, as will  
 181 be shown next.

182 **Turath Standard** The Turath Standard benchmark database comprises images reflecting the diverse  
 183 range of objects, activities, and scenarios commonly encountered in the Arab world. Each image has  
 184 a *macro* and *micro* image-level category annotation. The twelve macro categories are Cities, Food,  
 185 Nature, Architecture, Dessert, Clothing, Instruments, Activities, Drinks, Souq, Dates, and  
 186 Religious Sites. The complete list of the more granular *micro* categories can be found in Appendix A.  
 187 The number of images in each of these micro categories is presented in Fig. 3a. We can see that  
 188 each micro category has anywhere between 50 – 500 images. This is by design since we explicitly  
 189 searched for *up to* 500 images per category and excluded categories with fewer than 50 images. We  
 190 applied this strategy to all benchmark databases to avoid categories with too few images which may  
 191 contain noise and thus hinder a network’s ability to learn.

192 Table 1: **Overview of training, validation, and test splits**  
 193 **for the Turath benchmark databases.** The number of  
 194 macro categories is shown in brackets.

	Turath Database		
	Standard	Art	UNESCO
Training	38,894	46,665	9,540
Valid.	6,418	7,531	1,558
Test	19,472	22,969	4,778
Categories	269 (12)	419	79

201  
 202

203 paintings, sculptures, etc.) created by Arab artists alongside annotations, at the image-level, of such  
 204 artists. We purposefully excluded these categories from the Turath Standard benchmark for the  
 205 following reasons. First, the large number of *micro* categories (419) that would have fallen under  
 206 the *macro* category of Art would have overwhelmed the categories outlined in the Turath Standard  
 207 benchmark. Second, distinguishing between images containing intricate, low-level details reflected  
 208 by paintings, sculptures, etc., poses a difficult task, in and of itself. As a result, this warranted a

For benchmarking, the Turath Standard database contains 38,894 images in the training set, 6,418 images in the validation set, and 19,472 images in the test set (see Table 1). Unless otherwise specified, all data splits are performed uniformly at random with a ratio of 70:10:20 for the training, validation, and test sets, respectively.

**Turath Art** The Turath Art benchmark comprises images of art (e.g.,

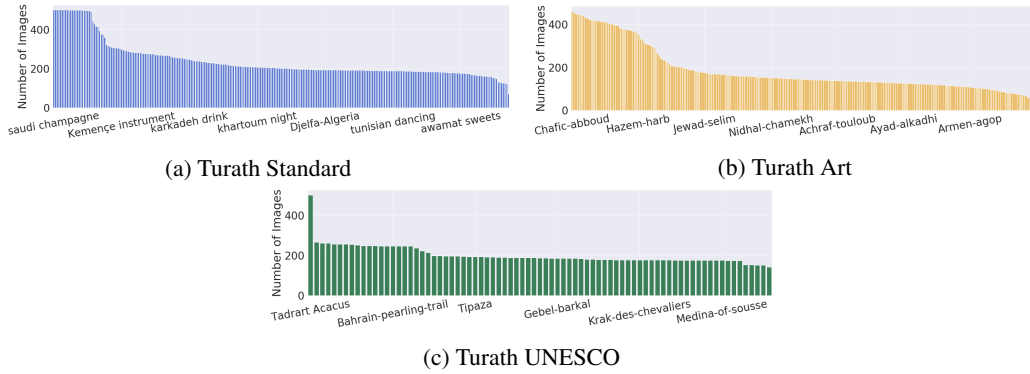


Figure 3: **Number of images per *micro* category in each of the benchmark databases.** Each micro category contains anywhere between 50-500 images. For clarity, we present only a subset of the micro category names. The full list of categories can be found in Appendix A.

209 distinct specialized benchmark, which we refer to as Turath Art. In Fig. 3b, we present the number of  
 210 images in each of the 419 artist categories, and include a subset of the artists’ names for clarity. For  
 211 benchmarking, the Turath Art database contains 38,445 images in the training set, 6,354 images in  
 212 the validation set, and 19,324 images in the test set.

213 **Turath UNESCO** The Turath UNESCO benchmark comprises images of UNESCO world heritage  
 214 sites in the Arab world alongside annotations, at the image-level, of these sites. We present, in Fig. 3c,  
 215 the total number of images in each of the 79 categories. For benchmarking, the Turath UNESCO  
 216 database contains 9,540 images in the training set, 1,558 images in the validation set, and 4,778  
 217 images in the test set.

## 218 5 Experimental results

### 219 5.1 Limitations of networks pre-trained on ImageNet

220 The utility of a pre-trained neural network is contingent upon the similarity of the upstream task, on  
 221 which the network was trained, and the downstream task, on which the network is deployed [24]. To  
 222 qualitatively evaluate this utility in the context of the Turath database, we randomly sample images  
 223 from each of the benchmark databases, perform a forward pass through an EfficientNet [25] pre-  
 224 trained on ImageNet, and compare the Top-5 predictions to the ground-truth label (see Fig. 4). We find  
 225 that, across the benchmarks, EfficientNet assigns a high probability mass to incorrect image categories.  
 226 For example, it classified a sculpture by the artist Maysaloun Faraj as an envelope with a confidence  
 227 score (0.564) and Gebel Barkal, pyramids in Sudan, as a seashore with a confidence score (0.266).  
 228 These results also suggest that confidence-based decisions, such as network classification abstention  
 229 and out-of-distribution detection [26], may be of little value in this context. We show that these  
 230 limitations also extend to other neural architectures (see Appendix C).

### 231 5.2 Image classification on Turath benchmark databases

232 In this section, we adapt networks pre-trained on ImageNet using data from the Turath database  
 233 benchmarks. We do so by introducing, and randomly initializing, a classification head,  $p_\theta : h \rightarrow \hat{y} \in$   
 234  $\mathbb{R}^C$ , that maps the penultimate representation,  $h$ , of the feature extractor network to the predicted  
 235 probability distribution,  $\hat{y}$ , over the set of image categories,  $C \in \{12, 269, 419, 79\}$  depending on the  
 236 benchmark database. In the linear evaluation phase, we freeze the parameters of the feature extractor  
 237 network whereas in the fine-tuning phase, we use those parameters as an initialization and update  
 238 them accordingly. In both phases, we train networks using the Adam optimizer with a categorical  
 239 cross-entropy loss and a learning rate,  $lr \in [1e^{-3}, 1e^{-4}]$ . Further implementation details can be  
 240 found in Appendix B.

241 In Table 2, we present the Top-1 and Top-5 accuracy achieved by networks in these experiments.  
 242 The Top-1 accuracy refers to the percentage of image samples whose ground-truth category matches  
 243 the category most confidently predicted by the network. In contrast, Top-5 accuracy refers to the









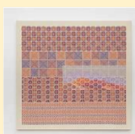



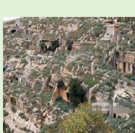

Standard		truth: Tyre, Lebanon top-1: seashore, coast, (0.245) top-2: lakeside, lakeshore (0.123) top-3: promontory, headland, (0.108) top-4: valley, vale (0.081) top-5: cliff, drop, (0.050)		truth: Gouraya National Park top-1: monastery (0.123) top-2: megalith, megalithic (0.094) top-3: triumphal arch (0.054) top-4: cliff dwelling (0.041) top-5: cliff, drop, (0.036)
		truth: Grape Leaves top-1: French loaf (0.142) top-2: tile roof (0.130) top-3: dough (0.097) top-4: bath towel (0.056) top-5: velvet (0.041)		truth: Tabbouleh top-1: pot, flowerpot (0.250) top-2: mortar (0.076) top-3: consomme (0.058) top-4: mixing bowl (0.047) top-5: guacamole (0.039)
Art		truth: Seif Wanly top-1: jigsaw puzzle (0.111) top-2: quilt, comforter, (0.109) top-3: vault (0.041) top-4: mosquito net (0.031) top-5: book jacket, (0.022)		truth: Abdullah AL Qassar top-1: book jacket, (0.131) top-2: shower curtain (0.094) top-3: comic book (0.076) top-4: jack-o'-lantern (0.017) top-5: guillotine (0.015)
		truth: Lamyia Gargash top-1: doormat, welcome (0.115) top-2: window shade (0.070) top-3: Band Aid (0.067) top-4: envelope (0.057) top-5: prayer rug, (0.051)		truth: Maysaloun Faraj top-1: envelope (0.564) top-2: carton (0.053) top-3: paper towel (0.037) top-4: pedestal, plinth, (0.033) top-5: perfume, essence (0.018)
UNESCO		truth: Palmyra, Syria top-1: mosque (0.292) top-2: monastery (0.129) top-3: triumphal arch (0.081) top-4: seashore, coast, (0.051) top-5: palace (0.031)		truth: Gebel Barkal, Sudan top-1: seashore, coast, (0.266) top-2: alp (0.109) top-3: yawl (0.063) top-4: valley, vale (0.056) top-5: Arabian camel, (0.049)
		truth: Cyrene, Libya top-1: cliff, drop, (0.737) top-2: cliff dwelling (0.053) top-3: valley, vale (0.035) top-4: alp (0.025) top-5: promontory, headland, (0.024)		truth: Al-Balad, Saudi Arabia top-1: mosque (0.079) top-2: seashore, coast, (0.061) top-3: lakeside, lakeshore (0.061) top-4: street sign (0.035) top-5: bell cote, (0.034)

Figure 4: **Top-5 predictions (and confidence) made by an EfficientNet, pre-trained on ImageNet and directly deployed on image samples from the Turath benchmark databases.** We also present the ground-truth *micro* category of each of the image samples. Many of the predictions assign a high probability mass to the incorrect category, lack the finer resolution of our *micro* categories, and do not have a cultural emphasis.

244 percentage of images samples whose ground-truth category can be found in the Top-5 most confident  
245 predictions made by the network<sup>4</sup>. On average, we find that EfficientNet outperforms MobileNetV2  
246 and ResNet50 uniformly across the benchmark databases. For example, on the UNESCO database,  
247 EfficientNet, in the linear evaluation phase, achieves Top-1= 39.5 whereas MobileNetV2 and  
248 ResNet50 achieve Top-1= 32.1 and 33.2, respectively. We also show that the *micro* category image  
249 classification tasks across benchmark databases differ in their level of difficulty. This is evident by the  
250 large range of reported accuracy scores. For example, Turath Standard poses the least difficult task  
251 with a best Top-1= 46.1 whereas Turath Art poses the most challenging task with a best Top-1= 16.5.  
252 This is expected given the high similarity of images in the Art database. We believe these accuracy  
253 scores, which remain relatively lower than those achieved on ImageNet (Top-1=90.2), stand to benefit  
254 from further advancements in neural architecture design, transfer learning, and domain adaptation.  
255 We also find that fine-tuning networks, regardless of the architecture, is more advantageous than a  
256 linear evaluation of such networks. This suggests that the fixed features extracted from a network  
257 pre-trained on ImageNet are relatively constraining.

<sup>4</sup>We provide demos of these networks in action at [danikiyasseh.github.io/Turath/](https://danikiyasseh.github.io/Turath/) [benchmark] Demo where benchmark  $\in$  [Standard, Art, UNESCO].

Table 2: **Image classification test accuracy on the Turath Standard, Art, and UNESCO benchmark databases.** Results are averaged across five random seeds and standard deviation is shown in brackets. Bold results reflect the best-performing network architecture in each benchmark.

Architecture	Standard ( <i>macro</i> )		Standard ( <i>micro</i> )		Art		UNESCO	
	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
<i>Linear evaluation</i>								
MobileNetV2	70.1 (0.7)	96.8 (0.1)	39.1 (0.1)	62.6 (0.1)	12.7 (0.2)	22.4 (0.2)	32.1 (0.4)	53.6 (0.2)
EfficientNet	<b>71.2</b> (0.3)	96.6 (0.1)	<b>46.1</b> (0.2)	<b>69.5</b> (0.1)	<b>16.5</b> (0.3)	<b>25.2</b> (0.3)	<b>39.5</b> (0.4)	<b>60.6</b> (0.2)
ResNet50	69.7 (0.2)	96.9 (0.2)	39.6 (0.5)	63.4 (0.3)	13.2 (0.2)	23.2 (0.3)	33.2 (0.3)	54.0 (0.2)
<i>Fine-tuning</i>								
MobileNetV2	65.6 (1.9)	95.6 (0.3)	41.7 (1.2)	65.9 (1.3)	12.9 (0.6)	23.6 (0.6)	34.4 (0.7)	56.1 (0.7)
EfficientNet	<b>77.2</b> (0.6)	<b>97.6</b> (0.0)	<b>49.9</b> (0.3)	<b>73.8</b> (0.3)	<b>19.0</b> (0.3)	<b>31.2</b> (0.4)	<b>43.2</b> (0.4)	<b>64.2</b> (0.7)
ResNet50	71.4 (0.7)	96.8 (0.1)	41.2 (1.3)	65.9 (1.0)	14.2 (0.8)	25.0 (1.1)	35.7 (1.7)	56.7 (1.4)

258 To gain better insight on the type of misclassifications committed on Turath Standard, we present,  
 259 in Fig. 5 (left), the confusion matrix of macro-category predictions made by EfficientNet on image  
 260 samples in the test set of the Turath Standard benchmark. This is complemented by Fig. 5 (right) in  
 261 which we illustrate the UMAP embedding of the penultimate representations ( $\mathbb{R}^{640}$ ) of the same set  
 262 of image samples. We chose the fine-tuned EfficientNet for these visualizations given its superior  
 263 performance (see Table 2). In light of Fig. 5, we find that the network is capable of comfortably  
 264 distinguishing between macro categories. This is evident by the relatively darker diagonal elements in  
 265 the confusion matrix and the high degree of category-specific separability of the UMAP embeddings.  
 266 On the other hand, we find that images in the Food category are occasionally misclassified as Dessert,  
 267 an error which makes sense given the semantic proximity of these categories.

268 Having shown that an EfficientNet can adequately learn to distinguish between the various categories  
 269 in the Turath benchmark databases, we wanted to explore whether its classifications were inferred  
 270 from the appropriate components of the input image. To do so, we exploit an established deep neural  
 271 network interpretability method, Grad-CAM [27], which attempts to identify the salient regions of the  
 272 input image in the form of a heatmap. Even though saliency methods have come under scrutiny [28],  
 273 we find that, in practice, they can be insightful. In Fig. 6, we illustrate the Grad-CAM-derived heatmap  
 274 overlaid on the original input image presented to a trained EfficientNet alongside the ground-truth  
 275 annotation of the image. In the case of Leptis Magna (Fig. 6c), we see that the ancient Carthaginian  
 276 arches are appropriately identified.

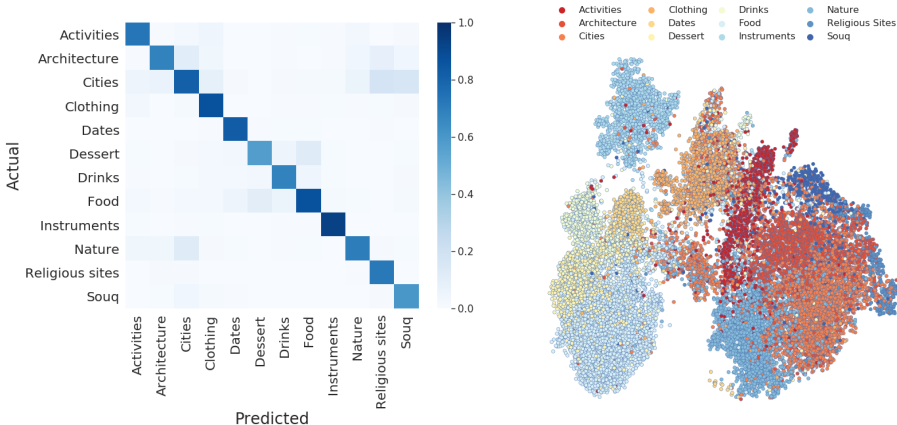


Figure 5: **Performance of EfficientNet fine-tuned on the Turath Standard benchmark database.** (Left) Confusion matrix of predictions made on the test set of the Turath Standard benchmark database. Normalization is performed across columns. (Right) UMAP embedding of the penultimate layer representations ( $\mathbb{R}^{640}$ ) of image samples in the test set. We find that the representations exhibit a high degree of separability amongst the macro categories.





Figure 6: **Heatmap of the most pertinent regions of the image for the category prediction.** We used Grad-CAM with an EfficientNet trained on the Turath (a) Standard, (b) Art, or (c) UNESCO benchmark databases. Red and blue regions are of high and low importance, respectively. We see that the network is able to identify regions in the image appropriate to the image category.

## 277 6 Discussion

278 In this paper, we discussed how existing image databases under-represent objects, activities, and  
 279 scenarios commonly found in certain cultures. To increase the cultural diversity of image databases,  
 280 we introduced Turath, a database of approximately 150K images of Arab heritage. Moreover, we  
 281 proposed three specialized benchmark databases, Turath Standard, Art, and UNESCO, that reflect a  
 282 range of entities within the Arab world and evaluated several deep networks on such benchmarks. Of  
 283 the networks evaluated, we found that EfficientNet performed best achieving Top-1 accuracy of 49.9,  
 284 19.0, and 43.2, on Turath Standard, Art, and UNESCO, respectively. We hope that our benchmark  
 285 databases can spur the research community to further advance neural architecture design, transfer  
 286 learning, and domain adaptation. That being said, it is vital that we consider the limitations and  
 287 broader societal impact of our work.

288 **Limitations** When searching for and cleaning the data, we opted out of a crowd-sourcing approach  
 289 (e.g., Mechanical Turk) in order to scale the database with minimal cost. The machine learning  
 290 community stands to benefit from the challenge of more independent data cleaning. Despite efforts  
 291 to clean the data, they exhibit some label noise and may thus benefit from innovative labelling  
 292 procedures, a challenge we leave to the community. Furthermore, any endeavour dependent on the  
 293 delineation of categories faces potential biases. Categories simplify and freeze nuanced narratives and  
 294 obscure political and moral reasoning [8]. Despite our cultural domain knowledge, niche categories  
 295 that remain undiscovered or unavailable online with sufficient images will not be represented in  
 296 our database. We aim to continue to engage with artists and heritage specialists to improve the  
 297 representativeness of our categories.

298 **Ethics and societal impact** Turath was primarily motivated by the need to increase the cultural  
 299 diversity of image databases, to improve the applicability of neural networks to under-represented  
 300 regions, and to actively engage researchers in such regions in the field of machine learning. However,  
 301 the cultural focus of this database may be prone to abuse by, for example, government and private  
 302 entities looking to delineate and target cultures for nefarious reasons. To mitigate the abuse of  
 303 our database for commercial purposes, we are releasing it under a CC BY-NC license, allowing  
 304 researchers to share and adapt the database in non-commercial settings. More broadly, our belief is  
 305 that by improving the awareness and understanding of cultures from around the globe, we can better  
 306 appreciate what they have to offer. Moving forward, we envision the Turath initiative expanding in  
 307 scope to encompass modalities such as text, audio, and video. Such a path can contribute to research  
 308 on language preservation, speech recognition, and video analysis.

## 309 References

- 310 [1] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and  
 311 Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb  
 312 model size. *arXiv preprint arXiv:1602.07360*, 2016.
- 313 [2] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time  
 314 object detection with region proposal networks. *arXiv preprint arXiv:1506.01497*, 2015.

- 315 [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille.  
316 DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution,  
317 and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
318 40(4):834–848, 2017.
- 319 [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
320 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern*  
321 *Recognition*, pages 248–255. Ieee, 2009.
- 322 [5] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:  
323 Large-scale scene recognition from abbey to zoo. In *2010 IEEE Computer Society Conference*  
324 *on Computer Vision and Pattern Recognition*, pages 3485–3492. IEEE, 2010.
- 325 [6] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A  
326 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and*  
327 *Machine Intelligence*, 40(6):1452–1464, 2017.
- 328 [7] Christiane Fellbaum. Wordnet. In *Theory and applications of ontology: computer applications*,  
329 pages 231–243. Springer, 2010.
- 330 [8] Abeba Birhane and Vinay Uday Prabhu. Large image datasets: A pyrrhic win for computer  
331 vision? In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer*  
332 *Vision*, pages 1537–1547, 2021.
- 333 [9] Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, and Olga Russakovsky. Towards fairer datasets:  
334 Filtering and balancing the distribution of the people subtree in the imagenet hierarchy. In  
335 *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages  
336 547–558, 2020.
- 337 [10] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural  
338 adversarial examples. *arXiv preprint arXiv:1907.07174*, 2019.
- 339 [11] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,  
340 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A  
341 critical analysis of out-of-distribution generalization. *arXiv preprint arXiv:2006.16241*, 2020.
- 342 [12] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The  
343 caltech-ucsd birds-200-2011 dataset. 2011.
- 344 [13] Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007.
- 345 [14] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
346 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*  
347 *Conference on Computer Vision*, pages 740–755. Springer, 2014.
- 348 [15] William Black, Sabri Elkateb, Horacio Rodriguez, Musa Alkhalifa, Piek Vossen, Adam Pease,  
349 and Christiane Fellbaum. Introducing the arabic wordnet project. In *Proceedings of the third*  
350 *international WordNet conference*, pages 295–300. Citeseer, 2006.
- 351 [16] Giorgos Toliás and Yannis Avrithis. Speeded-up, relaxed spatial matching. In *2011 International*  
352 *Conference on Computer Vision*, pages 1653–1660. IEEE, 2011.
- 353 [17] Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar  
354 Weber, and Antonio Torralba. Recipe1m+: A dataset for learning cross-modal embeddings for  
355 cooking recipes and food images. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.
- 356 [18] Gary B. Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the  
357 wild: A database for studying face recognition in unconstrained environments. Technical Report  
358 07-49, University of Massachusetts, Amherst, October 2007.
- 359 [19] Sarah M Erfani, Sutharshan Rajasegarar, Shanika Karunasekera, and Christopher Leckie. High-  
360 dimensional and large-scale anomaly detection using a linear one-class svm with deep learning.  
361 *Pattern Recognition*, 58:121–134, 2016.

- 362 [20] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and  
 363 Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection.  
 364 In *International Conference on Learning Representations*, 2018.
- 365 [21] Dan Li, Dacheng Chen, Jonathan Goh, and See-kiong Ng. Anomaly detection with generative  
 366 adversarial networks for multivariate time series. *arXiv preprint arXiv:1809.04758*, 2018.
- 367 [22] Izhak Golan and Ran El-Yaniv. Deep anomaly detection using geometric transformations. *arXiv*  
 368 *preprint arXiv:1805.10917*, 2018.
- 369 [23] Elad Amrani and Alex Bronstein. Self-supervised classification network. *arXiv preprint*  
 370 *arXiv:2103.10994*, 2021.
- 371 [24] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding  
 372 transfer learning for medical imaging. *arXiv preprint arXiv:1902.07208*, 2019.
- 373 [25] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural  
 374 networks. In *International Conference on Machine Learning*, pages 6105–6114. PMLR, 2019.
- 375 [26] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution  
 376 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 377 [27] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi  
 378 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based  
 379 localization. In *Proceedings of the IEEE international conference on computer vision*, pages  
 380 618–626, 2017.
- 381 [28] Richard Tomsett, Dan Harborne, Supriyo Chakraborty, Prudhvi Gurram, and Alun Preece.  
 382 Sanity checks for saliency metrics. In *Proceedings of the AAAI Conference on Artificial*  
 383 *Intelligence*, volume 34, pages 6021–6029, 2020.
- 384 [29] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu  
 385 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for  
 386 large-scale machine learning. In *12th {USENIX} symposium on operating systems design and*  
 387 *implementation ({OSDI} 16)*, pages 265–283, 2016.

## 388 Checklist

- 389 1. For all authors...
- 390 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 391 contributions and scope? [Yes] We claim and indeed introduce a database (see Sec. 3)  
 392 and evaluate several networks on such a database (see Sec. 5).
- 393 (b) Did you describe the limitations of your work? [Yes] We discuss the limitations of  
 394 category definitions and dataset bias (see Sec.6)
- 395 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We discuss  
 396 potential abuse of the dataset by government and non-government entities (see Sec. 6)
- 397 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 398 them? [Yes]
- 399 2. If you are including theoretical results...
- 400 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 401 (b) Did you include complete proofs of all theoretical results? [N/A]
- 402 3. If you ran experiments...
- 403 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
 404 mental results (either in the supplemental material or as a URL)? [Yes] We include the  
 405 URL to the corresponding website (which contains code and data) in the abstract. We  
 406 also include links to demos in Sec. 5
- 407 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 408 were chosen)? [Yes] We include data splits in Table 1. Implementation details are  
 409 included in Appendix B.

- 410 (c) Did you report error bars (e.g., with respect to the random seed after running exper-  
411 iments multiple times)? [Yes] We report the standard deviation (across five random  
412 seeds) of Top-1 and Top-5 accuracy scores in Table 2.
- 413 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
414 of GPUs, internal cluster, or cloud provider)? [Yes] We used Google Colab’s GPU  
415 resources and outline the duration of each training epoch in Appendix B.
- 416 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 417 (a) If your work uses existing assets, did you cite the creators? [Yes] We reference the  
418 creators of TensorFlow in Appendix B.
- 419 (b) Did you mention the license of the assets? [Yes] We are releasing the database and the  
420 code under a CC BY-NC license (see Sec. 6)
- 421 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
422 We include a link in the abstract to our website which has code, data, and models.
- 423 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
424 using/curating? [N/A]
- 425 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
426 information or offensive content? [N/A]
- 427 5. If you used crowdsourcing or conducted research with human subjects...
- 428 (a) Did you include the full text of instructions given to participants and screenshots, if  
429 applicable? [N/A] We did not crowd-source image annotations.
- 430 (b) Did you describe any potential participant risks, with links to Institutional Review  
431 Board (IRB) approvals, if applicable? [N/A] Since we did not crowd-source image  
432 annotations nor did we involve human subjects, IRB approval was not required.
- 433 (c) Did you include the estimated hourly wage paid to participants and the total amount  
434 spent on participant compensation? [N/A] Since we did not involve human participants,  
435 payment details are not applicable.