

# Leveraging Selective Prediction for Reliable Image Geolocation

Apostolos Panagiotopoulos, Giorgos Kordopatis-Zilos, and Symeon Papadopoulos

Information Technologies Institute, CERTH, Thessaloniki 60361, Greece  
{apanag, georgekordopatis, papadop}@iti.gr

**Abstract.** Reliable image geolocation is crucial for several applications, ranging from social media geo-tagging to media verification. State-of-the-art geolocation methods surpass human performance on the task of geolocation estimation from images. However, no method assesses the suitability of an image for this task, which results in unreliable and erroneous estimations for images containing ambiguous or no geolocation clues. In this paper, we define the task of image localizability, i.e. suitability of an image for geolocation, and propose a selective prediction methodology to address the task. In particular, we propose two novel selection functions that leverage the output probability distributions of geolocation models to infer localizability at different scales. Our selection functions are benchmarked against the most widely used selective prediction baselines, outperforming them in all cases. By abstaining from predicting non-localizable images, we improve geolocation accuracy from 27.8% to 70.5% at the city-scale, and thus make current geolocation models reliable for real-world applications.

**Keywords:** image localizability · selective prediction · geolocation estimation · spatial entropy · prediction density

## 1 Introduction

A great portion of the images daily uploaded on the Internet are from smartphones and therefore contain geotags, providing information for their geographic location. However, there are numerous cases, such as social media photos, where this information is missing. Therefore, the ability to estimate the geographic location of these images, known as *location estimation* or *geolocation*, is crucial for a number of applications ranging from social media mining to media verification. More formally, image geolocation is the process of inferring the GPS coordinates of the depicted picture based solely on its visual elements.

State-of-the-art approaches for geolocation employ the latest advances in Computer Vision, such as Convolutional Neural Networks (CNNs), to extract representations of the depicted scenes and utilize huge databases of images taken worldwide either for training a classifier [5–7] or for retrieval [2–4, 11]. Classification solutions partition the earth into a set of geographic cells, and the images

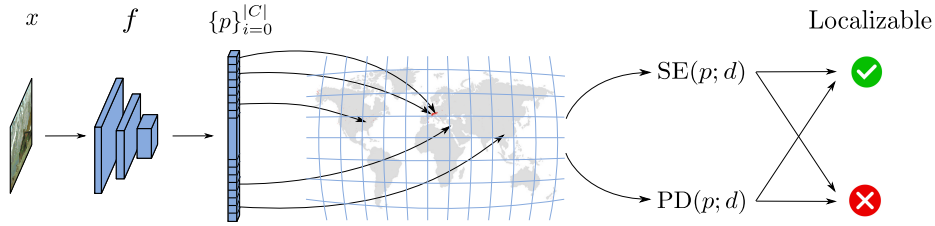


Fig. 1: Each image  $x$  is passed through a base geolocation model  $f$ , yielding a cell probability distribution  $\{p\}_{i=1}^{|C|}$ . Considering this distribution in the context of a world map, we measure whether our model’s confidence is concentrated in a particular region or dispersed over globe through our selection functions  $SE(p; d)$  and  $PD(p; d)$ , to assess the localizability of the input image.

are passed through a CNN to be classified into a single cell. Retrieval solutions compare the test images against the ones from a large-scale database in order to retrieve the most similar images and derive a single estimation by aggregating their locations. In both cases, performance is measured by the percentage of images localized within a certain distance  $d$  from their ground-truth location, denoted as *geolocation accuracy @  $d$  km*.

Ongoing research focuses mainly on improving geolocation accuracy at different granularities (e.g.  $d = 1, 25, 200, 750$  and  $2500$  km). However, contrary to most image classification datasets where images usually contain enough visual cues to be classified in a unique class, a good portion of images in geolocation datasets depict scenes with no apparent visual cues mapping to the image’s location (e.g. indoor spaces, portraits). Such images could have been captured anywhere on the globe, and therefore attempting to localize them would most likely result in erroneous or unreliable estimations. Hence, we deem crucial for the reliability of geolocation systems to estimate not only the geolocation of input images but also their *localizability*.

In a general sense, we consider localizable the images that contain enough visual cues for their accurate geolocation. However, localizability is better approached considering a granularity scale (i.e. the range within which an image can be correctly placed from its true location) and a geolocation model (i.e. a mechanism that derives the image’s location from its visual content). To this end, we introduce the problem of *image localizability detection*, building upon the foundation of *selective prediction* [12]. We propose a methodology that utilizes state-of-the-art geolocation systems to infer localizability at different scales. More precisely, we re-implement an image geolocation model [5] and, instead of interpreting the output probability distribution as pure categorical data and predicting the most probable location, we visualize the whole distribution on the world map. We then devise two novel selection functions, i.e. spatial entropy and prediction density, which measure cell probability dispersion and concentration over the globe, exploiting intrinsic properties of geolocation systems – unlike current state-of-the-art selection functions. We extensively evaluate our

methodology on the two most widely used evaluation datasets, i.e. Im2GPS [2] and Im2GPS3k [11], and highlight the effectiveness of our proposed selection functions compared to state-of-the-art approaches in selective prediction. By discarding images considered non-localizable at city-scale, we boost the accuracy of our base geolocation model at city-scale from 27.8% to 70.5%, making it reliable for real-world applications. To the best of our knowledge, we are the first to propose the task of image localizability detection and leverage selective prediction to address it. Therefore, our work makes the following contributions:

- We introduce the problem of image localizability detection and frame it under a selective prediction framework. This formulation allows current classification models to infer localizability, and hence abstain from predicting non-localizable images.
- We propose two novel selection functions, specifically designed for geolocation, that outperform current state-of-the-art selection functions which do not consider spatial information.
- We extensively evaluate our methodology on the two most widely used datasets, achieving good separation between localizable and non-localizable images, and making current geolocation systems more reliable.

## 2 Related Work

This section gives an overview of some of the fundamental works that have contributed to geolocation estimation and selective prediction.

**Geolocation Estimation:** Hays and Efros [2] introduced the problem of planet-scale image location estimation. They used handcrafted features to retrieve images similar to a query image and infer its location based on theirs. Weyand et al. [5] took advantage of the deep learning advances and trained a Convolutional Neural Network (CNN) to extract features from images. Additionally, they formulated the geolocation problem as a classification task and divided the earth’s surface, using Google’s s2 geometry library<sup>1</sup>, to create a set of classes for the training and test images. More recent works modify the classification pipeline using a hierarchical partitioning of the earth [6], or novel loss functions [7]. Recently, a hybrid scheme called Search within Cell [8] was proposed, which combines a classification and retrieval approach for the final location estimation.

**Selective prediction:** Works in this area focus on machine learning systems that are not only able to make predictions but also to know when to abstain from predicting. Although the field exists for several decades, it was not until recently that a unified formulation was introduced [12] and approaches regarding deep architectures were proposed by El-Yaniv and Wiener [15, 16]. In [15], Softmax Response (the maximum output after the softmax layer) and MC-Dropout [13] were used as confidence functions, and an algorithm that finds the appropriate threshold given a desired risk was proposed. In [16], an algorithm for jointly training the classification network and the selection function was proposed.

<sup>1</sup> <https://s2geometry.io>

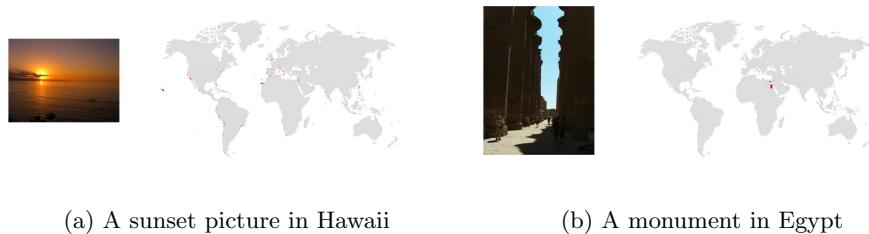


Fig. 2: Placing our model’s cell probability distribution on a map illustrates its ability to associate visual *concepts* to locations.

### 3 Methodology

In this section, we present the proposed methodology for the selection of localizable images; this is illustrated in Fig. 1.

#### 3.1 Geolocation Estimation

A geolocation model  $f$  takes as input an image  $x \in \mathbb{R}^{H \times W \times 3}$  and returns the estimated GPS coordinates  $\hat{y} \in \mathbb{R}^2$  of the location it was captured. Following the classification approach for geolocation, we first divide the earth’s surface into a grid of geographic cells  $C$ , and then we employ a CNN, with  $|C|$  outputs in the final layer, corresponding to the cells of the grid. For each input image  $x$ , the CNN creates a probability distribution over the grid of cells. The predicted location  $\hat{y}$  is the mean coordinates of the cell with the highest probability.

Most geolocation approaches [5–7] consider only the cell with the highest probability, ignoring all the information provided by the cell probability distribution over the grid. We found that this probability distribution can provide valuable insights for the localizability of images. More precisely, a trained network has learned both to estimate the location of images and also to associate *concepts* with locations. For example, when the network is presented with the image of Fig. 2a, the probability of several cells nearby the sea is high. However, when presented with the image from Fig. 2b, all cells around Egypt are activated, since the image contains many visual cues that map to that area. Thus, even though it is challenging to predict the exact location of those images, the model generates *reasonable* estimates to candidate locations.

By inspecting the map, it is evident that the network is more confident for the estimation of Fig. 2b’s location than Fig. 2a’s since the probability distribution is more concentrated on a specific area. Thus, the estimation of the former can be considered more reliable than the latter. The spatial distribution of cells is essential information for the geolocation estimation problem, differentiating it from the general image classification task. Hence, our goal in this paper is to exploit this information to improve the reliability of the model’s predictions.

### 3.2 Localizability

To develop and evaluate a methodology for image localizability, we have to associate all images in a dataset with ground-truth labels that indicate localizability. Moreover, labeling images as localizable or not is highly subjective and depends on the collection of images recognized by the prospective annotator. To address the former issue, we define localizability at a certain scale  $d$  (distance tolerance from the ground truth location). To address the latter issue, we approximate localizability in terms of our model’s ability to infer location from the input image, i.e. we assess which images our employed model is able to predict correctly. Therefore, all images that our model is able to predict within a certain distance from their ground-truth location are labeled as localizable, and all other images are labeled as non-localizable.

More formally, given a geolocation estimation model  $f$ , the localizability of an image  $x \in \mathbb{R}^{M \times N \times 3}$  at distance  $d$  is defined as:

$$\mathcal{L}_f(x; d) = \begin{cases} 1, & \text{if } \text{GCD}(f(x), y) < d \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

where  $\text{GCD}(\cdot, \cdot)$  is the Great Circle Distance between two locations, and  $y$  is the ground-truth location of  $x$ .

### 3.3 Selective Prediction

Predicting which images are localizable according to our model’s geolocation capability can be formulated as a selective prediction scheme following the formulation in [12]. Here, we adapt their definitions to fit the needs of the geolocation estimation task. Our aim is to build a selective geolocation system  $(f, g)$  such that  $f$  is the base geolocation module, as described in section 3.1, and  $g$  is the selection function. Then our selective geolocation system is defined as:

$$(f, g)(x; d) = \begin{cases} f(x; d), & \text{if } g(x; d) = 1 \\ \text{abstain}, & \text{if } g(x; d) = 0 \end{cases} \quad (2)$$

The selection function  $g$  is usually modeled based on a confidence function  $\kappa_f$  (which measures our model’s confidence or uncertainty), a scale  $d$  and a tunable threshold  $\theta$ . For  $\kappa_f$  measuring confidence,  $g(x; d)$  is defined as follows:

$$g(x; d) = \begin{cases} 1, & \text{if } \kappa_f(x; d) \geq \theta \\ 0, & \text{if } \kappa_f(x; d) < \theta \end{cases} \quad (3)$$

Let  $P(X, Y)$  be the distribution over  $\mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{X}$  is the image space and  $\mathcal{Y}$  the coordinate space, characterizing the probability of image  $X$  being captured at geographical coordinates  $Y$ . Given an underlying distribution  $P$ , a confidence function  $\kappa_f$  and a scale  $d$ , varying the parameter  $\theta$  determines the performance of our selective geolocation system, which can be expressed using coverage and risk, as follows:

**Coverage** is the mass probability of the non-rejected region in  $\mathcal{X}$ , and can be approximated given enough i.i.d. samples  $(x_i, y_i)$  from  $P$  as follows:

$$\phi_d(f, g) \triangleq \mathbb{E}_{P(X)} [g(x; d)] \approx \frac{1}{N} \sum_{i=1}^N g(x_i; d) \quad (4)$$

**Risk** is the expected percentage of the kept images that will be predicted outside a radius  $d$ , and can be approximated given enough i.i.d. samples  $(x_i, y_i)$  from  $P$  as follows:

$$R_d(f, g) \triangleq \mathbb{E}_{P(X, Y)} [l_d(f(x), y)] \approx \frac{\frac{1}{N} \sum_{i=1}^N l_d(f(x_i), y_i) g(x_i; d)}{\phi_d(f, g)} \quad (5)$$

where  $l_d$  is a loss function defined as:

$$l_d(l_1, l_2) = \begin{cases} 1, & \text{if } \text{GCD}(l_1, l_2) > d \\ 0, & \text{if } \text{GCD}(l_1, l_2) \leq d \end{cases} \quad (6)$$

### 3.4 Estimating Image Localizability

Our main goal in this work is to find good confidence functions and thresholds for the function  $g$ . For  $f$ , we employ a geolocation system similar to [5]; however, any geolocation system that tackles geolocation as a classification problem can be used.

**Spatial Entropy (SE)** measures the dispersion of cell probabilities around the globe. To calculate the  $\text{SE}(p; d)$  of the cell distribution  $\{p\}_{i=0}^{|C|}$  at scale  $d$ , we initially select the most probable cell and merge all cells within distance  $d$  from it to form a super-cell. The probability of the super-cell derives from the sum of the cell probabilities of the individual cells. Then, we ignore all cells merged to the super-cell and find the next most probable cell from the remaining ones. Similarly, we merge it with its neighboring cells that are inside a radius  $d$ . This process is repeated until the cumulative probability of the super-cells accounts for the 90% of the total confidence<sup>2</sup>. We denote as  $\{\bar{p}\}_{i=0}^{|C'|}$  the new probability distribution of the super-cells; hence, SE is defined as:

$$\text{SE}(p; d) = - \sum_{i=0}^{|C'|} \bar{p}_i \log_2 \bar{p}_i \quad (7)$$

Higher Spatial Entropy indicates lower confidence; therefore, we devise the selection function  $g_{SE}$  as:

$$g_{SE}(x) = \begin{cases} 1, & \text{if } \text{SE}(p; d) \leq \theta_{SE} \\ 0, & \text{if } \text{SE}(p; d) > \theta_{SE} \end{cases} \quad (8)$$

where  $\theta_{SE}$  is a tunable threshold.

<sup>2</sup> We empirically found this to remove noise from cell distributions compared to considering cells accounting for 100% of the total confidence.

**Prediction Density (PD)** measures the concentration of cell probability in a particular region instead of its dispersion around the globe. To calculate the  $\text{PD}(p; d)$  of the cell distribution  $\{p\}_{i=0}^{|C|}$  at scale  $d$ , we accumulate the model’s cell probabilities in a radius  $d$  around the most probable cell, which can be considered as the model’s confidence that an input image can be localized at scale  $d$ . This is formulated as follows:

$$\text{PD}(p; d) = \sum_{c \in C} p_c \cdot \mathbb{1} \left[ \text{GCD}(c, \arg \max_{c' \in C} p_{c'}) \leq d \right] \quad (9)$$

where  $\mathbb{1}$  is the indicator function. Higher Prediction Density denotes higher confidence; therefore, we devise the selection function  $g_{PD}$  as:

$$g_{PD}(x) = \begin{cases} 1, & \text{if } \text{PD}(p; d) \geq \theta_{PD} \\ 0, & \text{if } \text{PD}(p; d) < \theta_{PD} \end{cases} \quad (10)$$

where  $\theta_{PD}$  is a tunable threshold.

## 4 Evaluation Set-up

### 4.1 Datasets

To train our geolocation model, we use the training split of the MediaEval Placing Task 2016 dataset (*MP-16 train*) [9], which is a subset of the Yahoo Flickr Creative Commons 100 Million (*YFCC100M*) [10]. It consists of 4,723,695 images posted on Flickr with their metadata, among which geographical coordinates. We also use the YFCC25k dataset from [6], composed of 25,600 randomly selected images from YFCC100M (excluding images from MP-16 train), for validation. Due to the unavailability of several images, we end up with a total of 23,007 images. Finally, for evaluation, we use the Im2GPS [2] and Im2GPS3k [11] datasets, provided by the original authors, consisting of 237 and 3,000 images, respectively.

### 4.2 Implementation Details

For the cell partitioning described, we adopt the fine partitioning from [6] and terminate cell splitting when each cell contains between 50 and 1,000 images from the MP-16 train. We discard all cells that end up with less than 50 photos. This results in 13,662 cells and 4,071,346 images for training. Although the particular partitioning implementation could affect the geolocation estimation performance, we are primarily interested in the performance of our localizability methods, and hence we do not consider alternate partitionings.

For the geolocation model  $f$ , we use EfficientNet-B4 [1] as our backbone CNN and replace its last layer with a linear layer consisting of 13,662 neurons corresponding to our total number of cells. We replicate the pre-processing and training pipeline of Kordopatis et al. [8], and we do not use any further additions such as hierarchical partitioning [6] or the MvMF loss [7].

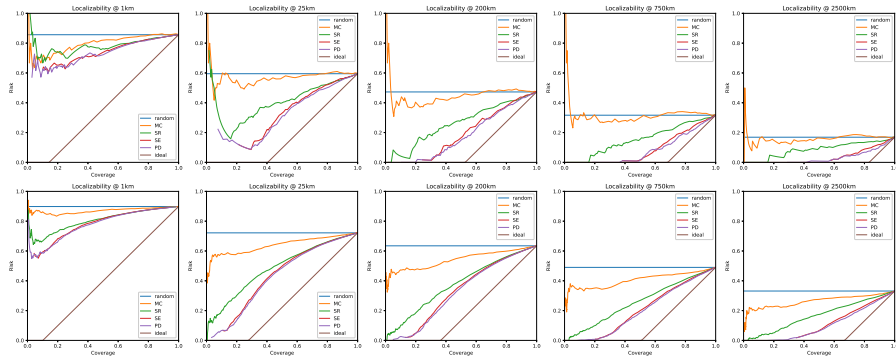


Fig. 3: Risk-Coverage curves for six selection functions on Im2GPS (top row) and Im2GPS3k (bottom row). Lower is better.

### 4.3 Competing Approaches

In Section 5, we compare our selection functions against two baseline runs (which serve to visualize selective performance limits) and two state-of-the-art methods in selective prediction, briefly described below:

**Random** selection function: randomly selects whether to predict  $f(x)$  or abstain from providing a prediction, with an equal probability of 50%.

**Ideal** selection function: selects images based on their ground-truth localizability values, prioritizing the images considered localizable.

**Softmax Response (SR)** [14]: uses the maximum probability after the final softmax layer, i.e. the maximum cell probability, as confidence function. It has been employed for selective prediction in [15]. Note that this selection function cannot be intrinsically adapted for the different scales  $d$ .

**Monte-Carlo Dropout (MC)** [13]: uses as confidence function the variance of the softmax response of multiple forward passes of an input image with dropout applied in the final layer. This is shown to be a good approximation of a Bayesian Neural Network with Gaussian parameter priors [13] and a state-of-the-art method in selective prediction [15]. Following [15] we use a dropout of 0.5. Note that again this selection function cannot be intrinsically adapted for the different scales  $d$ .

## 5 Experiments

### 5.1 Selective Geolocation Performance

We benchmark the selective prediction performance of the proposed selection functions  $g_{SE}$  and  $g_{PD}$  against the competing approaches. We evaluate them in both Im2GPS and Im2GPS3k at scales  $d = 1, 25, 200, 750$  and  $2500$ km, which are the most widely reported scales and correspond to street, city, region, country and continent level granularity scales.



	Acc ↑	F1 ↑	OR ↓	OC ↑		Acc ↑	F1 ↑	OR ↓	OC ↑		Acc ↑	F1 ↑	OR ↓	OC ↑
<b>SR</b>	74%	38%	71%	27%	<b>SR</b>	70%	66%	40%	49%	<b>SR</b>	70%	73%	31%	57%
<b>MC</b>	<b>83%</b>	23%	67%	5%	<b>MC</b>	62%	31%	53%	19%	<b>MC</b>	53%	36%	38%	25%
<b>SE</b>	76%	44%	67%	<b>28%</b>	<b>SE</b>	73%	72%	37%	<b>55%</b>	<b>SE</b>	76%	80%	29%	<b>67%</b>
<b>PD</b>	78%	<b>47%</b>	<b>64%</b>	27%	<b>PD</b>	<b>79%</b>	<b>77%</b>	<b>31%</b>	50%	<b>PD</b>	<b>78%</b>	<b>81%</b>	<b>26%</b>	64%

(a)  $d = 1$  km                      (b)  $d = 25$  km                      (c)  $d = 200$  km

	Acc ↑	F1 ↑	OR ↓	OC ↑		Acc ↑	F1 ↑	OR ↓	OC ↑
<b>SR</b>	70%	79%	24%	74%	<b>SR</b>	82%	90%	14%	<b>88%</b>
<b>MC</b>	49%	53%	28%	41%	<b>MC</b>	61%	74%	16%	62%
<b>SE</b>	79%	86%	20%	<b>78%</b>	<b>SE</b>	85%	91%	10%	85%
<b>PD</b>	<b>84%</b>	<b>88%</b>	<b>12%</b>	70%	<b>PD</b>	<b>87%</b>	<b>92%</b>	<b>8%</b>	85%

(d)  $d = 750$  km                      (e)  $d = 2500$  km

Table 1: Selective prediction performance of our two selection functions against state-of-the-art on Im2GPS. We report the localizability accuracy and the F1-score as well as the optimal risk (OR) and coverage (OC) (for more details read Section 5.1). ↑ indicates that higher is better, and ↓ that lower is better.

	Acc ↑	F1 ↑	OR ↓	OC ↑		Acc ↑	F1 ↑	OR ↓	OC ↑		Acc ↑	F1 ↑	OR ↓	OC ↑
<b>SR</b>	84%	34%	70%	14%	<b>SR</b>	78%	62%	39%	29%	<b>SR</b>	76%	69%	33%	38%
<b>MC</b>	<b>87%</b>	7%	84%	3%	<b>MC</b>	70%	23%	58%	11%	<b>MC</b>	65%	35%	48%	18%
<b>SE</b>	85%	44%	64%	<b>16%</b>	<b>SE</b>	83%	72%	35%	<b>34%</b>	<b>SE</b>	80%	75%	31%	<b>43%</b>
<b>PD</b>	86%	<b>45%</b>	<b>62%</b>	15%	<b>PD</b>	<b>85%</b>	<b>74%</b>	<b>29%</b>	31%	<b>PD</b>	<b>82%</b>	<b>77%</b>	<b>26%</b>	40%

(a)  $d = 1$  km                      (b)  $d = 25$  km                      (c)  $d = 200$  km

	Acc ↑	F1 ↑	OR ↓	OC ↑		Acc ↑	F1 ↑	OR ↓	OC ↑
<b>SR</b>	72%	74%	28%	53%	<b>SR</b>	70%	79%	25%	<b>73%</b>
<b>MC</b>	57%	50%	37%	34%	<b>MC</b>	58%	65%	28%	53%
<b>SE</b>	78%	80%	24%	<b>56%</b>	<b>SE</b>	78%	84%	18%	71%
<b>PD</b>	<b>80%</b>	<b>81%</b>	<b>22%</b>	55%	<b>PD</b>	<b>80%</b>	<b>86%</b>	<b>16%</b>	70%

(d)  $d = 750$  km                      (e)  $d = 2500$  km

Table 2: Selective prediction performance of our two selection functions against state-of-the-art on Im2GPS3k. We report the localizability accuracy and the F1-score as well as the optimal risk (OR) and coverage (OC) (for more details read Section 5.1). ↑ indicates that higher is better, and ↓ that lower is better.

First, we present the Risk-Coverage (RC) curves for each dataset, illustrated in Fig. 3. These curves are obtained by computing the risk and coverage of each method for different values of the threshold  $\theta$ . Both our selection functions achieve state-of-the-art performance, yielding lower risk at every coverage level in all datasets and scales. SE and PD perform similarly, with PD consistently outperforming SE by a small margin. Moreover, in coarser granularity scales, the performance gap between our selection functions and the competing SR and MC widens considerably, with our selection functions reaching close to the ideal. This can probably be attributed to their intrinsic adaptation to different scales. Finally, it is evident that RC curves on Im2GPS3k are smoother and more monotonous, which is expected due to its greater size and variety of images.

Although risk-coverage curves give a comprehensive insight of the selective prediction performance, we need to determine a specific threshold  $\theta$  that sepa-

	1km	25km	200km	750km	2500km		1km	25km	200km	750km	2500km
$f$	14.3	40.5	52.7	68.3	83.1	$f$	10.1	27.8	36.5	51.0	66.8
$(f, g_{SE})^L$	24.4	62.6	74.0	86.2	91.6	$(f, g_{SE})^L$	24.6	65.3	77.4	86.9	92.6
$(f, g_{PD})^L$	27.1	69.5	79.6	89.8	94.0	$(f, g_{PD})^L$	26.6	70.5	80.8	89.1	94.1
$(f, g_{SE})^N$	1.9	13.2	26.4	46.2	72.6	$(f, g_{SE})^N$	2.7	8.4	17.7	35.0	55.6
$(f, g_{PD})^N$	1.6	11.7	26.0	47.0	72.6	$(f, g_{PD})^N$	2.6	8.2	15.3	32.3	54.5

(a) Im2GPS

(b) Im2GPS3k

Table 3: Geolocation accuracies (%) when evaluating the whole dataset and the Localizable (L) and Non-Localizable (N) subsets. The percentage of images on (L) corresponds to the Optimal Coverage (OC) column of Tables 1 and 2.

rates localizable and non-localizable images given a selection function. To do so, we chose the  $\theta$  value that corresponds to the coverage that equals the percentage of images  $f$  can successfully localize. We learn this value on the validation YFCC25K dataset for each selection function and each granularity scale. We call the risk and coverage at this threshold *Optimal Risk* (OR) and *Optimal Coverage* (OC) respectively. We also report the classification accuracy and the F1-score of the positive class.

Tables 1 and 2 display the results of the selection functions on the two evaluation datasets. We note that high accuracy in finer scales  $d$  is not indicative of good separation between localizable and non-localizable images due to the class imbalance; however, combined with F1-score, they provide useful insights. In particular, it is evident that the proposed SE and PD achieve better class separation than the SR and MC, with PD slightly surpassing SE. Moreover, in most cases, the selected threshold  $\theta$  for our methods leads to lower risk and wider coverage compared to the competition.

## 5.2 Selective Geolocation Reliability

We present a quantitative and qualitative assessment of the performance of our selective models  $(f, g_{SE})$  and  $(f, g_{PD})$  compared to  $f$ .

We split both Im2GPS and Im2GPS3k into a localizable and a non-localizable subset using our selective models at city-scale. Table 3 displays the geolocation accuracies on these splits at all granularity scales, compared to the performance of the base model  $f$  without a selection scheme. For fine and medium granularity scales, our selective models achieve more than double the geolocation accuracy of the base model  $f$ , with only a tiny portion of localizable images rejected by our functions. In particular, prediction density increased the geolocation accuracy on Im2GPS3k from 27.8% to 70.5% by discarding non-localizable images, from which only 8.2% could have been successfully localized. This highlights the reliability current image geolocation models can achieve using the selective prediction mechanisms presented.

Fig. 4 depicts image samples randomly selected from Im2GPS3k for qualitative evaluation of our methodology. Images are grouped by their predicted and

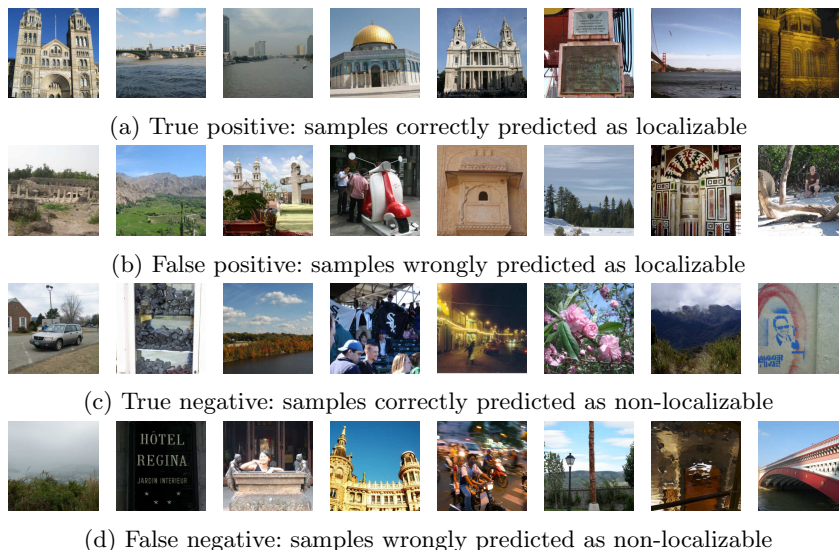


Fig. 4: Sample predictions of Prediction Density (PD) on randomly selected images of the Im2GPS3k dataset.

ground-truth localizability in (a) true positive, (b) false positive, (c) true negative and (d) false negative, using PD at city-scale as the selection function. True positive samples either depict landmarks or characteristic elements that hint at very specific locations (e.g. the Golden Gate Bridge). True negative samples contain mostly generic scenes that should not even be attempted to be localized. The two images with the car and the lights could have been localized if they were in our training dataset, but even in that case, a similar scene can easily exist in multiple cities. False positive samples contain enough visual cues to be worthy of geolocation, however not enough for the required granularity. Finally, all false negative samples besides the Edificio Meneses picture are not localizable and their correct geolocation by our geolocation model could be attributed to presence of very similar images in the train dataset.

## 6 Conclusions

In this paper, we introduced the problem of image localizability detection and used it as a foundation for reliable image geolocation. We adapted a selective prediction methodology to the context of geolocation and presented two novel selection functions, Spatial Entropy and Prediction Density, tailored to the needs of the geolocation task. Our functions achieved superior selective performance compared to state-of-the-art on the two widely-used evaluation datasets. We also demonstrated how they can be exploited to abstain from geolocating non-localizable images, significantly boosting the geolocation performance in all granularity scales, and thus making current geolocation models more reliable. In the

future, we plan to explore the design and evaluation of more sophisticated and trainable selection functions.

**Acknowledgments:** This work has been supported by the projects WeVerify and MediaVerse, partially funded by the European Commission under contract number 825297 and 957252, respectively.

## References

1. Tan M, Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, 2019.
2. Hays J, Efros AA. Im2gps: estimating geographic information from a single image. In *IEEE Computer Vision and Pattern Recognition*, 2008.
3. Hays J, Efros AA. Large-scale image geolocation. In *Multimodal location estimation of videos and images*, 2015.
4. Kordopatis-Zilos G, Popescu A, Papadopoulos S, Kompatsiaris Y. Placing Images with Refined Language Models and Similarity Search with PCA-reduced VGG Features. In *MediaEval*, 2016.
5. Weyand T, Kostrikov I, Philbin J. Planet-photo geolocation with convolutional neural networks. In *European Conference on Computer Vision*, 2016.
6. Muller-Budack E, Pustu-Iren K, Ewerth R. Geolocation estimation of photos using a hierarchical model and scene classification. In *European Conference on Computer Vision*, 2018.
7. Izbicki M, Papalexakis EE, Tsotras VJ. Exploiting the earth’s spherical geometry to geolocate images. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2019.
8. Kordopatis-Zilos G, Galopoulos P, Papadopoulos S, Kompatsiaris I. Leveraging EfficientNet and Contrastive Learning for Accurate Global-scale Location Estimation. In *ACM International Conference on Multimedia Retrieval*, 2021.
9. Larson M, Soleymani M, Gravier G, Ionescu B, Jones GJ. The benchmarking initiative for multimedia evaluation: MediaEval 2016. *IEEE MultiMedia*. 2017 Feb 9;24(1):93-6.
10. Thomee B, Shamma DA, Friedland G, Elizalde B, Ni K, Poland D, Borth D, Li LJ. YFCC100M: The new data in multimedia research. *Communications of the ACM*. 2016 Jan 25;59(2):64-73.
11. Vo N, Jacobs N, Hays J. Revisiting im2gps in the deep learning era. In *IEEE International Conference on Computer Vision*, 2017.
12. El-Yaniv R. On the Foundations of Noise-free Selective Classification. *Journal of Machine Learning Research*. 2010 May 1;11(5).
13. Gal Y, Ghahramani Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning*, 2016.
14. Hendrycks D, Gimpel K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*. 2016.
15. Geifman Y, El-Yaniv R. Selective Classification for Deep Neural Networks. In *Advances in Neural Information Processing Systems*, 2017.
16. Geifman Y, El-Yaniv R. Selectivenet: A deep neural network with an integrated reject option. In *International Conference on Machine Learning*, 2019.