Image Geolocalization with Directed Feature Identification

Claire Gong, Madeline Brumley

Motivation

• Image geolocalization is known to be a

challenging task for vision models, but

Dataset construction for task is challenging:

Most high-performing models require lots

identifying characteristic features in an

image such as bollards or license plates

human guessers can do very well

hard to represent every region

Human guessers achieve success by

Can models perform better and extract

better representations from data by

attending to extracted features like

Experimental Setup

We fine-tune pretrained ViT-B-16 models on raw

Mask2Former model for instance segmentation.

feature-to-country dataset into 4 main groups:

We then split each group into test/train splits &

fine-tune a ViT-B-16 model for each group for

We evaluate each model on their respective test

measuring top-1 percent country-label accuracy.

We fine-tune 3 ViT-B-16 models on the entire

feature-to-country dataset from G3, and a

selection of 15k features from G50 equally

We evaluate on respective test splits and

IM2GPS with top-1 percent accuracy.

distributed across 6 countries.

feature-to-country dataset from G50, the entire

splits by top-1 percent labeling accuracy.

We then assemble the model ensemble with

Mask2Former and evaluate it on IM2GPS,

images from G50 and G3 datasets to use as

• We construct feature-to-country datasets

from G50 and G3 using Facebook's

We group feature classes from the G50

Expert Feature Classifiers:

Ground (road, terrain)

Structures (buildings)

Vehicles (cars, trucks, etc.)

feature-to-country classification.

Vegetation (foliage)

of data to train

humans do?

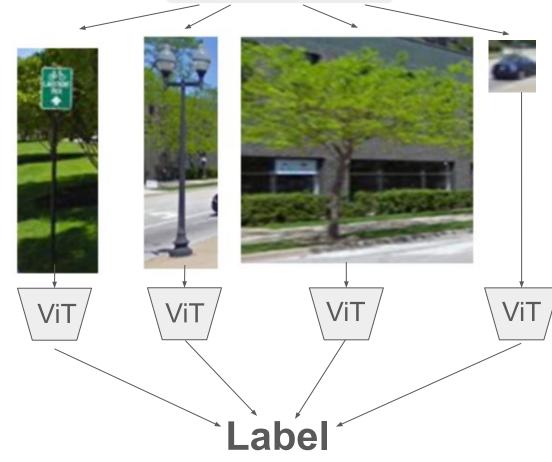
baselines.

Architecture

• 3-step model ensemble for

Input Image

Instance Segmentation Model



- At inference time, input images sent
 - Road markings, vegetation, buildings, cars, ect.
- Each category of feature sent to a respective expert Vision Transformer fine-tuned for image-to-country
- Most frequently appearing country label across features chosen as final label

Results

Model	Accuracy
G50 Baseline	
G3 Baseline	0.3155
G50 Feature Classifier	0.5393
G50 6 Country Feature Classifier	0.8181
G3 Feature Classifier	0.5626
G50 Ground Classifier	0.4946
G50 Structure Classifier	0.3919
G50 Vegetation Classifier	0.3721
G50 Vehicle Classifier	0.1716
GeoCLIP	

Table 1. Results from ViT Test Accuracies

Model	Accuracy
G50 Feature Classifier	
G50 6 Country Feature classifier	0.2911
G3 Feature Classifier	0.0000
G50 Feature Ensemble Classifier	0.0302
G50 Ground Classifier	0.0251
G50 Structure Classifier	0.0101
G50 Vegetation Classifier	0.0754
G50 Vehicle Classifier	0.0000

Datasets

We use a mix of original and publicly available

evaluation:

countries.

image-to-country pair datasets for training and

1. (G50): Geolocation - Geoguessr Images,

(G3): Original dataset of 3k Google

equally distributed across countries

We use **instance segmentation** to extract

features from images in both datasets to

1. 99k features from G50, all countries

3. 15k features from G50, equally distributed

Sample image with label "Australia

fence" from G3 feature-to-country

We evaluate models on original test sets

derived from above feature-to-country

datasets, and on IM2GPS test set of

create feature-to-country datasets.

2. 6k features from G3, all countries

across 6 countries

geo-tagged Flickr images.

Sample image with label "Norway"

from original Street View images

(G3) dataset

Street View images from 55 countries,

50k Google Street View images from 124

- image-to-country classification

- through instance segmentation model to extract important geographical signifiers
- classification

99

Performance varies between different expert feature classifiers, dragging down overall ensemble model performance significantly

Discussion

dataset

- Likely due to differences in extracted feature quality.
- Feature extraction is still promising performance of models trained on feature-to-country data on test splits is markedly improved from baseline model performance
- Distribution of data across countries in G50 extremely unequal, likely causing reductions in accuracy.
- IM2GPS images are very dissimilar to Street View training data, likely contributing to abysmal performance.

Future Work!

- Is directed feature extraction better or worse than random crop and scaling data augmentation?
- Is there a better way to extract features so that all features are the same quality?
- Are accuracy scores on IM2GPS higher on the continent-level?

Model	Accuracy
G50 Baseline	
G3 Baseline	0.3155
G50 Feature Classifier	0.5393
G50 6 Country Feature Classifier	0.8181
G3 Feature Classifier	0.5626
G50 Ground Classifier	0.4946
G50 Structure Classifier	0.3919
G50 Vegetation Classifier	0.3721
G50 Vehicle Classifier	0.1716
GeoCLIP	

References

All-Feature Classifiers:

- [1] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask trans
 - former for universal image segmentation. 2022. 3 J. Hays, A. Efros, et al. Im2gps: estimating geographic information from a single image. In Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, pages

1–8. IEEE, 2008. 2, 3, 4, 5, 7, 8

//www.kaggle.com/datasets/ubitquitin/ geolocation - geoguessr - images - 50k / data, 2021. [Online; accessed 27-Feb-2024]. 3 Table 2. Model Evalulation Accuracies on Im2GPS

Rohan K. Geolocation - geoguessr images (50k). https: