A	UTOMATED OSINT GEOLOCATION
<b>A</b> Pa	nonymous authors aper under double-blind review
	ABSTRACT
	Open source intelligence (OSINT) investigators face the challenge of veri- fying the location of media shared online. Traditional geolocation requires manual effort and cannot scale with the ever-growing volume of images and videos shared on social media. We present GeoFT, a fine-tuned version of GeoCLIP specifically optimized for geolocation in Russia and Ukraine. By focusing on street-level imagery and leveraging community-validated datasets, our model achieves significantly improved accuracy compared to existing solutions. On our test set, GeoFT reduces the average error from 3,520km to 2,150km while maintaining interpretable confidence scores. We demonstrate the model's potential for aiding OSINT investigations and dis- cuss pathways for deployment in real-world applications.
1	INTRODUCTION
effo dat ima tha solu	rt relying on human investigators using resources like Google Street View and I abases. While effective, this manual process cannot keep pace with the thousand ges appearing daily on platforms like Decrypt. We present GeoFT, a fine-tune t specifically targets street-level imagery in regions of active conflict, where tra- tions like Google Street View may be outdated or unavailable.
2	Related Work
Rec GE tras (3,5 imp ima pro	ent approaches to AI geolocation include GeoCLIP (Vivanco Cepeda et al., 20 OTTO (Haas et al., 2024), and GeoDecoder (Qi et al., 2024). GeoCLIP utility tive learning between image and location embeddings but exhibits high avera 20km). PIGEOTTO and GeoDecoder showed promising results but lack open elementations. Commercial solutions like GeoSpy (GeoSpy, 2024) rely heavily of ge databases, limiting their effectiveness in rapidly changing environments. The aches demonstrate the challenge of accurate geolocation at scale.
3	Methodology
3.1	Data Collection and Filtering
We	curated a dataset combining two primary sources:
	1. Eyes on Russia (EoR): 2,887 community-geotagged images from the conflic (Eyes on Russia Project, 2024)
	2. Google Street View (GSV): 16,159 street-level images gathered via API LLC, 2024)
To the imp	ensure data quality, we employed GPT-4-mini as a binary classifier to filter image prompt: "Does this image contain street features?" This preprocessing step sign proved training data relevance by removing indoor scenes and irrelevant image

## GEOFT: FINE-TUNING FOUNDATION MODELS FOR

000

vdsourced building ids of new ied model raditional

2023), PIlizes conage error en-source on static These ap-

- ict region
- I (Google
- ages with gnificantly gery. The

054 final dataset contained approximately 19,000 images after filtering. The test set consisted 055 of only EoR images to reflect realistic images from social media. 056

3.2 Model Architecture 058

057

060

061

062

063

064

065

066 067

068

069

070

071

078

079 080

081

- 059 We build upon GeoCLIP's architecture, consisting of:
  - 1. Location encoder L(\*) that transforms 2D GPS coordinates into high-dimensional features through:
    - (a) Equal Earth Projection to minimize geographic distortion
    - (b) Random Fourier Features (RFF) for positional encoding
    - (c) Hierarchical representation using multiple RFF frequencies
    - (d) MLP layers for processing
    - 2. Image encoder V(\*) based on a pre-trained CLIP ViT-L/14 model with two additional trainable linear layers (h1: 768 dimensions, h2: 512 dimensions)

During training, we keep the CLIP backbone frozen and only train the linear layers of the image encoder and the location encoder components. Training uses contrastive learning with batch and queue-based negatives to minimize the loss:

$$\mathcal{L}_{i} = -\sum_{j=1}^{P} \log \frac{\exp(V_{ij} \cdot L_{ij}/\tau)}{\sum_{i=0}^{|B|} \exp(V_{ij} \cdot L_{ij}/\tau) + \sum_{i=0}^{S} \exp(V_{ij} \cdot \bar{L}_{i}/\tau)}$$
(1)

where  $\tau$  is the temperature parameter, |B| is the batch size, S is the queue size, and  $L_i$ represents the queue of GPS embeddings used as additional negatives.

RESULTS  $\Delta$ 

082 Our evaluation reveals three key findings that demonstrate GeoFT's effectiveness and prac-083 tical applicability.

First, when trained on the combined Eyes on Russia (EoR) and Google Street View (GSV) 085 dataset, GeoFT demonstrates significant improvements over baseline models. Specifically, 086 GeoFT achieves state-of-the-art performance on the Im2GPS3k dataset with improvements 087 of +1.31%, +0.97%, +3.95%, +8.67%, and +3.72% at 1km, 25km, 200km, 750km, and 880 2500km thresholds respectively compared to prior work.

Second, the fine-tuned model exhibits substantially higher confidence in its predictions. 090 While the baseline model shows uncertainty in its geolocation estimates (mean confidence: 091 0.03), GeoFT produces well-calibrated confidence scores reaching up to 0.61 for high-092 confidence predictions, with confidence scores strongly correlating with prediction accuracy. The standard deviation in confidence scores increased from 0% to 5%, indicating a better 094 ability to distinguish between high and low confidence predictions. 095

Third, and perhaps most significantly for practical applications, GeoFT demonstrates 096 strong performance even when trained solely on GSV data, without requiring region-specific human-labeled datasets like EoR. This finding is crucial for extending the approach to new 098 geographic regions where specialized datasets may not be available. Table 1 compares per-099 formance between models trained on the combined dataset versus GSV data only. 100

As shown in Table 1, while performance does decrease without EoR data, particularly at 101 finer granularities, the model still achieves strong results at continental (2500km) and coun-102 try (750km) scales using only GSV data. This is particularly noteworthy since GSV data 103 is widely available globally, making our approach extensible to new regions without requir-104 ing specialized datasets. The model trained only on GSV data still outperforms baseline 105 methods, demonstrating the effectiveness of our approach even with limited training data. 106

The confidence metrics also reflect this pattern - the GSV-only model exhibits lower but 107 still well-calibrated confidence scores (max 2% vs 61% with combined data), appropriately 

109	Table 1: Performance comparison between models trained on combined EoR+GSV data
110	versus GSV data only. While performance decreases without EoR data, the model still
111	achieves strong results at larger scales using only widely-available GSV data.
110	Distance EoR+GSV GSV Only Difference

112	Distance EoR+GSV GSV Only Difference
113	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
114	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
115	$\frac{200 \text{ km}}{25 \text{ km}} = \frac{38\%}{88\%} = \frac{6\%}{6\%} = \frac{32\%}{25 \text{ km}}$
116	1  km $6%$ $0%$ $-6%$
117	
118	
119	indicating its reduced certainty in predictions. This ability to produce reliable confidence
120	estimates, even with limited training data, is crucial for practical applications where under-
121	standing prediction reliability is important.
122	
123	5 Discussion and Future Work
124	
125	While GeoFT shows promising results, several extensions could enhance its utility:
126	
127	1. Integration with production systems for continuous model improvement using new
128	validated data
129	2. Expansion to video and aerial imagery analysis
130	2. Extension to other perions of interest with similar data collection methodology
131	5. Extension to other regions of interest with similar data conection methodology
132	The model can be deployed both as a standalone tool for OSINT investigators and integrated
133	into existing intelligence platforms, providing automated first-pass location estimates for
134	human verification.
135	
136	6 Conclusion
137	0 CONCLUSION
138	CeoFT demonstrates that fine-tuning foundation models on carefully curated regional data
139	can significantly improve geolocation accuracy. The success of our approach using only
140	Google Street View data suggests that this methodology could be extended to other regions
141	of interest, even without access to specialized OSINT datasets. This work represents a step
142	toward scalable, automated support for OSINT investigations.
143	
144	References
145	THE ENERGES
146	Eyes on Russia Project. Eyes on russia. https://eyesonrussia.org, 2024. Accessed:
147	February 2024.
148	GeoSpy Geospy: Ai-powered image geolocation https://geospy.ai 2024 Accessed:
149	February 2024.
150	
151	Google LLC. Google street view. https://www.google.com/streetview, 2024. Accessed:
152	February 2024.
153	Lukas Haas, Michal Skreta, Silas Alberti, and Chelsea Finn. Pigeon: Predicting image
154	geolocations. Preprint. May 2024.
155	
156	Feng Qi, Mian Dai, Zixian Zheng, and Chao Wang. Geodecoder: Empowering multimodal
157	map understanding. arXiv preprint, Feb 2024. arXiv:2401.2401.15118.
158 159	Vicente Vivanco Cepeda, Gaurav Kumar Nayak, and Mubarak Shah. Geoclip: Clip-inspired
	ang michi octacen iocanono and mages for enective autuande geo-iocalization. In Au-

vances in Neural Information Processing Systems, 2023. NeurIPS 2023.