
Swarm Intelligence in Geo-Localization: A Multi-Agent Large Vision-Language Model Collaborative Framework

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Visual geo-localization demands in-depth knowledge and advanced reasoning skills
2 to associate images with real-world geographic locations precisely. In general,
3 traditional methods based on data-matching are hindered by the impracticality
4 of storing adequate visual records of global landmarks. Recently, Large Vision-
5 Language Models (LVLMs) have demonstrated the capability of geo-localization
6 through Visual Question Answering (VQA), enabling a solution that does not
7 require external geo-tagged image records. However, the performance of a single
8 LVLM is still limited by its intrinsic knowledge and reasoning capabilities. Along
9 this line, in this paper, we introduce a novel visual geo-localization framework
10 called smileGeo that integrates the inherent knowledge of multiple LVLM agents
11 via inter-agent communication to achieve effective geo-localization of images.
12 Furthermore, our framework employs a dynamic learning strategy to optimize the
13 communication patterns among agents, reducing unnecessary discussions among
14 agents and improving the efficiency of the framework. To validate the effectiveness
15 of the proposed framework, we construct GeoGlobe, a novel dataset for visual geo-
16 localization tasks. Extensive testing on the dataset demonstrates that our approach
17 significantly outperforms state-of-the-art methods. The source code is available at
18 <https://anonymous.4open.science/r/VisualGeoLocalization-F8F5/> and the dataset
19 will also be released after the paper is accepted.

20 1 Introduction

21 Visual geo-localization, referred to the task of estimating geographical identification for a given
22 image, is vital in various fields such as human mobility analysis [1, 2, 3, 4, 5] and robotic navigation
23 [6, 7, 8, 9, 10, 11]. In general, accurate visual geo-localization without the help of any localization
24 equipment (*e.g.*, GPS sensors) is a complex task that requires abundant geospatial knowledge and
25 strong reasoning capabilities. Traditional methods [12, 13, 14, 15] typically formulate it as an image
26 retrieval problem where to geo-localize the given image by retrieving similar images with known
27 geographical locations. Thus, their effectiveness is limited by the scope and quality of the geo-tagged
28 image records.

29 Recently, the success of Large Vision-Language Models (LVLMs) has enabled Visual Question
30 Answering (VQA) to become a unified paradigm for multi-modal problems [16, 17], providing a
31 novel solution for visual geo-localization without the need for external geo-tagged image records.
32 However, the performance of a single LVLM on the geo-localization task is still limited by its
33 inherent geospatial knowledge and reasoning capabilities. Along this line, in this paper, we introduce
34 a novel multi-agent framework, named **swarm intelligence Geo-localization (smileGeo)**, which
35 aims to adaptively integrate the inherent knowledge and reasoning capabilities of multiple LVLMs

36 to effectively and efficiently geo-localize images. Specifically, for a given image, the framework
37 initially elects K suitable LVLM agents as answer agents for initial location analysis. Then, each
38 answer agent chooses several review agents via an adaptive social network, which imitates the
39 collaborative relationships between agents with a target on the visual geo-localization task, to
40 discuss and share their knowledge for refining its location analysis. Finally, our framework conducts
41 free discussion among all of the answer agents to reach a consensus. Besides, we also design
42 a novel dynamic learning strategy to optimize the election mechanism along with the adaptive
43 collaboration social network of agents. We hope that by the effectiveness of the election mechanism
44 and the review mechanism, our framework can discover the mode of communication among agents,
45 thereby enhancing geo-localization performance through multi-agent collaboration while minimizing
46 unnecessary discussions. In summary, our contributions are demonstrated as follows:

- 47 • A novel swarm intelligence geo-localization framework, smileGeo, is proposed to adaptively
48 integrate the inherent knowledge and reasoning capability of multiple LVLMs through
49 discussion for visual geo-localization tasks.
- 50 • A dynamic learning strategy is introduced to discover the most appropriate discussion mode
51 among LVLM agents for enhancing the effectiveness and efficiency of the framework.
- 52 • A new visual geo-localization dataset named GeoGlobe¹ is collected, containing a wide
53 variety of images globally. The diversity and richness of GeoGlobe allow us to evaluate
54 the performance of different models more accurately. Moreover, extensive experiments
55 demonstrate our competitive performance compared to state-of-the-art methods.

56 The remainder of this paper is organized as follows: Section 2 discusses the related literature. In
57 Section 3, the proposed framework is introduced. Section 4 provides the performance evaluation, and
58 Section 5 concludes the paper.

59 2 Related Work

60 **Visual Geo-localization.** Recent research in visual geo-localization, commonly referred to as
61 geo-tagging, primarily focuses on developing image retrieval systems to address this challenge
62 [3, 18, 19, 20, 21, 22]. These systems utilize learned embeddings generated by a feature extraction
63 backbone, which includes an aggregation or pooling mechanism [23, 24, 25, 26]. However, the
64 applicability of these retrieval systems to globally geo-localize landmarks or natural attractions is
65 often limited by the constraints of the available database knowledge and the restrictions imposed by
66 national or regional geo-data protection laws. Alternatively, some studies treat visual geo-localization
67 as a classification problem [27, 28, 29, 30]. These approaches posit that two images from the same
68 geographical region, despite depicting different scenes, typically share common semantic features.
69 Practically, these methods organize the geographical area into discrete cells and categorize the
70 image database accordingly. This cell-based categorization facilitates scaling the problem globally,
71 provided the number of categories remains manageable. However, while the number of countries
72 globally remains relatively constant, accurately enumerating cities in real-time at a global scale is
73 challenging due to frequent administrative changes, such as city reorganizations or mergers, which
74 reflect shifts in national policies. Additionally, in the context of globalization, this strategy has
75 inherent limitations. The recent advent of LVLMs offers promising compensatory mechanisms for
76 the deficiencies observed in traditional geo-localization methodologies, making the exploration of
77 LVLM-based approaches significantly relevant in current research.

78 **Multi-agent Framework for LLM/LVLMs.** LLM/LVLM agents have demonstrated the potential
79 to act like human [31, 32, 33], and a large number of studies have focused on developing robust
80 architectures for collaborative LLM/LVLM agents [34, 35, 36, 37, 38]. These architectures enable
81 each LLM/LVLM agent that endows with unique capabilities to engage in debates or discussions.
82 For instance, [34] proposes an approach to aggregate multiple LLM/LVLM responses by generating
83 candidate responses from various LLM/LVLM in a single round and employing pairwise ranking to
84 synthesize the most effective response. While some studies [34] utilize a static architecture potentially
85 limiting the performance and generalization of LLM/LVLM, others like [38] have implemented
86 dynamic interaction architectures that adjust according to the query and incorporate user feedback.

¹Because GeoGlobe is relatively large (about 32GB), we are unable to provide it as an attachment during the double-blind review stage. We will publish it once the paper is accepted.

87 Recent advancements also demonstrate the augmentation of LLM/LVLM as autonomous agents
88 capable of utilizing external tools to address challenges in interactive settings. These techniques
89 include retrieval augmentation [39, 40, 41], mathematical tools [40, 42, 43], and code interpreters
90 [44, 45]. With these capabilities, LLM/LVLMs are well-suited for various tasks, especially for
91 geo-localization. However, most LLM/LVLM agent frameworks mandate participation from all
92 agents in at least one interaction round, leading to significant computational overhead. To address
93 this issue, our framework introduces a dynamic learning strategy electing only a small number of
94 agents to geo-localize different images, which significantly enhances the efficiency of LLM/LVLM
95 agents by reducing unnecessary interactions.

96 **3 Methodology**

97 In this section, we first present the overall framework and then introduce each part of smileGeo in
98 detail for geo-localization tasks.

99 **3.1 Model Overview**

100 In this paper, we denote the social network of LVLM agents by \mathcal{G} , where $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$. \mathcal{V} stands for the
101 agent set and \mathcal{E} presents the edge set. Each agent $v_i \in \mathcal{V}, i \in [N]$ is an LVLM, which is pre-trained
102 by massive vision-language data and can infer the possible location Y of a given image X . Besides,
103 each edge $e_{ij} \in \mathcal{E}, i, j \in [N]$ is the connection weighted by the improvement effect of agent v_i to
104 agent v_j via discussion regarding the geo-localization performance.

105 As illustrated in Figure 1, smileGeo contains the process of the review mechanism in agent discussions
106 along with a dynamic learning strategy of agent social networks:

107 The review mechanism in agent discussions is a 3-stage anonymous collaboration approach to allow
108 LVLM agents to reach a consensus via discussion. In the first stage, for a given image X , our
109 framework elects the most suitable K agents as answer agents by agent election probability Lst . In
110 the second stage, these answer agents respectively select R review agents by the adaptive collaboration
111 social network A to refine their answer via discussion. Finally, our framework facilitates consensus
112 among all agents through open discussion to reach a final answer. Both Lst and A are analyzed
113 from the given image X , allowing our framework to minimize unnecessary discussions, thereby
114 significantly enhancing its efficiency while maintaining its accuracy. Moreover, the multi-stage
115 discussion facilitates communication among agents, maximizing the integration of their knowledge
116 and reasoning abilities to generate an accurate response Y .

117 To get Lst and A , we specifically design a dynamic learning module, which initially deploys the
118 encoder component of a pre-trained image variational autoencoder (VAE) to extract features from the
119 given image X . The extracted features, combined with agent embeddings Emb , are employed to
120 determine the suitability of agents *w.r.t.* Lst for agent discussions and predict agent collaboration
121 connections A in the geo-localization task.

122 **3.2 Review Mechanism in Agent Discussions**

123 LLM/LVLM have demonstrated remarkable capabilities in complicated tasks and some pioneering
124 works have further proven that the performances can be further enhanced by ensembling multiple
125 LLM/LVLM agents. Thus, to improve the geo-localization capability of LVLMs, we propose a
126 cooperation framework to effectively integrate the diverse knowledge and reasoning abilities of
127 multiple LVLMs. Inspired by the fact that community review mechanisms can improve the quality of
128 manuscripts, an iterative 3-stage anonymous reviewing mechanism is proposed for helping agents
129 share knowledge and reasoning capability with each other through their collaboration social network:
130 i) answer agent election & answering, ii) review agent selection & reviewing, and iii) final answer
131 conclusion.

132 **Stage 1: Answer Agent Election & Answering**

133 Initially, we select K agents with the highest agent election probabilities Lst as answer agents and
134 let them geo-localize independently as the preliminary step for further discussion. By initiating
135 the discussion with a limited number of agents, we aim to reduce potential chaos and maintain the
136 efficiency of our framework as the number of participating agents increases.

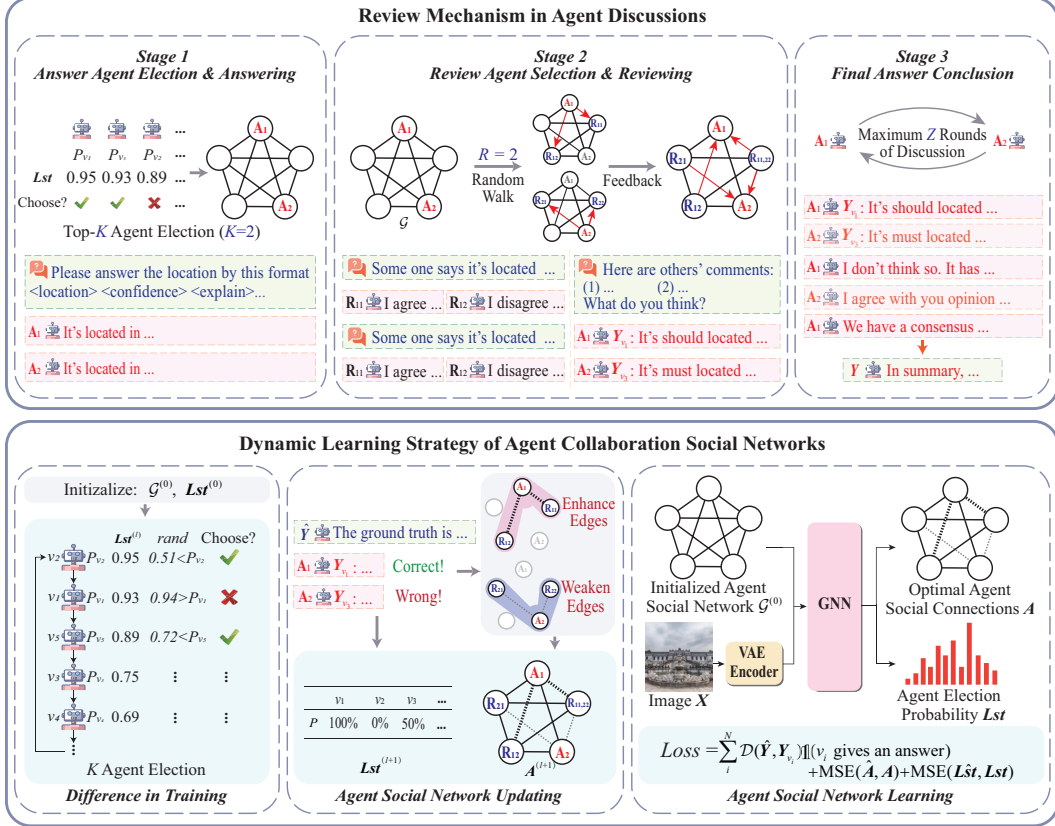


Figure 1: The framework overview of smileGeo. It contains the process of review mechanism in agent discussions along with a dynamic learning strategy of agent collaboration social networks. The first part deploys a review mechanism for LVLMs to discuss and share their knowledge anonymously, which could enhance the overall performance of geo-localization tasks. The second one mainly utilizes the GNN-based learning module to improve efficiency by reducing unnecessary discussions among agents while showing the process of updating the agent collaboration social network during the training process.

137 After the answer agents are elected, we send the image X to all answer agents and let them give the
 138 primary analysis. Each answer must contain three parts: one location (city, country, and so on), one
 139 confidence (a percentage number), and a detailed explanation.

140 Stage 2: Review Agent Selection & Reviewing

141 In this stage, for each answer agent, we choose R review agents by performing a transfer-probability-
 142 based random walk on the agent collaboration social network \mathcal{G} for answer reviewing. The transfer
 143 probability $p(v_i, v_j)$ from node v_i to node v_j can be calculated as follows:

$$p(v_i, v_j) = \begin{cases} \frac{A_{ij}}{\sum_{k \in \mathcal{N}(v_i)} A_{ik}}, & \text{if } e_{ij} \in \mathcal{E} \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

144 where $\mathcal{N}(v_i)$ is the 1-hop neighbor node set of node v_i .

145 For each selected review agent, it reviews the results as well as the explanations generated by the
 146 corresponding answer agent and gives its own comments. After that, each answer agent would
 147 summarize their preliminary analysis and the feedback from all of its review agents to get the final
 148 answer, which must include three parts as well: one location, one confidence, and an explain.

149 Stage 3: Final Answer Conclusion

150 In the previous stage, each answer agent produces a refined result based on feedback. When $K > 1$ in
 151 Stage 1, the proposed framework generates multiple independent results, which may not be consistent.

152 However, we aim to provide a definitive answer rather than multiple options for people to choose
 153 from. To address this, we allow up to Z rounds of free discussion among those answer agents to
 154 reach a unified answer:

155 First, we maintain a global dialog history list, $diag$, recording all replies agents respond. In addition,
 156 discussions are executed asynchronously, which means that any answer agent can always reply based
 157 on the latest $diag$, and replies would be added to the end of $diag$ as soon as they are posted. Each
 158 answer agent is allowed to speak only once in each discussion round, and after Z rounds of free
 159 discussion, we determine the final result using a minority-majority approach, *i.e.*, we choose the reply
 160 with the most agreement as the final conclusion. If all agents reach a consensus, we early stop this
 161 stage and adopt the consensus answer as the final answer. If none of any consensus is reached, we
 162 only select the reply of the first answer agent elected from Stage 1 as the final result.

163 3.3 Dynamic Learning Strategy of Agent Collaboration Social Networks

164 In our framework, choosing the appropriate answer agents and review agents for knowledge sharing
 165 and discussion is vital to its effectiveness and efficiency. Therefore, we propose a dynamic learning
 166 strategy to optimize them. Specifically, for each training sample, *i.e.*, a geo-tagged image, we would
 167 first estimate the optimal answer agent election probability \hat{Lst} and the optimal collaboration social
 168 network of agent $\hat{\mathcal{G}}$ by its actual location. Then we train an attention-based graph neural network,
 169 which aims to predict Lst and \mathcal{G} , by such estimated ground truth.

170 To estimate the optimal \hat{Lst} and $\hat{\mathcal{A}}$ for agents to geo-localize image \mathbf{X} , we first initialize the agent
 171 social network $\mathcal{G}^{(0)}$ by a fully connected graph with the agent set \mathcal{V} . Besides, we initialize the agent
 172 election probability $Lst^{(0)} = [0.5, 0.5, \dots]$, with all agents having 50% probability of being chose
 173 as answer agents.

174 Then, we iteratively conduct our 3-stage discussion framework to get the prediction answer. $Lst^{(l)}$
 175 and $\mathcal{G}^{(l)}$ is updated at the end of each round $l \in L$ by comparing the answers $\mathbf{Y}_{v_i}^{(l)}$ from each answer
 176 agent with the ground truth $\hat{\mathbf{Y}}$.

177 After L rounds of agent discussions, the updated agent election probability for an image \mathbf{X} , $\hat{Lst} :=$
 178 $Lst^{(L)}(\mathbf{X}) = [P_{v_1}^{(L)}, P_{v_2}^{(L)}, \dots, P_{v_N}^{(L)}]$, determines whether an agent v_i gives the correct/wrong
 179 answers $\mathbf{Y}_{v_i}^{(L)}$ by comparing it with the ground truth $\hat{\mathbf{Y}}$. Here, the definition of $P_{v_i}^{(l)}$ of agent v_i at
 180 round l is as follows:

$$P_{v_i}^{(l)} := \begin{cases} 0, & \text{if } \mathcal{D}(\hat{\mathbf{Y}}, \mathbf{Y}_{v_i}^{(l)}) > th \\ 1, & \text{if } \mathcal{D}(\hat{\mathbf{Y}}, \mathbf{Y}_{v_i}^{(l)}) \leq th \\ \frac{1}{2}, & \text{if } v_i \text{ did not participate in the discussion} \end{cases} \quad (2)$$

181 where th is a pre-defined threshold for determining whether the predicted location is close enough
 182 to the actual location. In the distance function $\mathcal{D}(\cdot)$, we first deploy geocoding to convert natural
 183 language into location intervals in a Web Mercator coordinate system (WGS84) by utilizing OSM
 184 APIs, and then compute the shortest distance between two two location intervals.

185 Please note that, rather than electing the top- K answer agents in each round, we choose each agent
 186 with probability P_{v_i} during the training period to ensure that every agent has the opportunity to
 187 participate in the discussion for more accurate estimation, as shown at the left part of the dynamic
 188 learning strategy module of agent collaboration social networks in Figure 1.

189 In addition, the agent collaboration social network would also be updated by comparing the actual
 190 location with the generated answer of each answer agent at the same time. For l -th round, we
 191 strengthen the link between the correctly answered agent and the corresponding review agents while
 192 weakening the link between the incorrectly answered agent and the corresponding review agents:

$$\hat{\mathbf{A}}_{ij} := \mathbf{A}_{ij}^{(l)}(\mathbf{X}) = \begin{cases} \frac{tt+1}{2tt} \mathbf{A}_{ij}^{(l-1)}(\mathbf{X}), & \text{if agent } v_i \text{ answers correctly} \\ \frac{2tt-1}{2tt} \mathbf{A}_{ij}^{(l-1)}(\mathbf{X}), & \text{if agent } v_i \text{ answers incorrectly} \end{cases} \quad (3)$$

193 where $\mathbf{A}_{ij}^{(l-1)}(\mathbf{X})$ is the weight of the connection between answer agent v_i and review agent v_j
 194 at round $l - 1$ when geo-locating image \mathbf{X} , $\mathbf{A}_{ij}^{(0)}(\mathbf{X}) = 1, i \neq j, \mathbf{A}_{ii}^{(0)}(\mathbf{X}) = 0, i, j \in [N]$, tt
 195 is the number of consecutive times an agent has answered correctly, which is used to attenuate
 196 the connection weights when updating them, preventing the performance of an agent on a certain
 197 portion of the continuous dataset from interfering with the model’s evaluation of the current agent’s
 198 performance on the entire dataset.

199 Then, we try to learn an attention-based graph neural network to predict the corresponding optimal
 200 agent election probability $\mathbf{Lst} = h(\mathbf{X}, \mathcal{G}|\Theta)$ and the optimal agent collaboration connections
 201 $\mathbf{A} = f(\mathbf{X}, \mathcal{V}|\Theta)$:

$$\begin{aligned} \mathbf{A} &= \text{Att}_{\text{GNN}}(\mathbf{Fea}, \mathbf{Fea}, \mathbf{1}) \\ &= \text{softmax}\left(\frac{\mathbf{Fea} \cdot \mathbf{Fea}^\top}{\sqrt{d_k}}\right) \mathbf{1}, \\ \mathbf{Lst} &= \sigma'(\text{Linear}(\text{Flatten}(\sigma(\mathbf{A} \cdot \mathbf{Fea} \cdot \mathbf{W}))))), \\ \mathbf{Fea} &= \text{Linear}(\mathbf{Emb} + \text{VAE}_{\text{Enc}}(\mathbf{X})), \end{aligned} \quad (4)$$

202 where $\mathbf{W}, \mathbf{Emb} \in \Theta$ are two learnable parameters, $\mathbf{Emb} := [\mathbf{Emb}_{v_1}, \mathbf{Emb}_{v_2}, \dots]^\top$ is the agent
 203 embedding and \mathbf{W} is the weight matrix, $\sigma(\cdot)$ is the LeakyReLU function, $\sigma'(\cdot)$ is the Sigmoid
 204 function, $\text{VAE}_{\text{Enc}}(\cdot)$ is the encoder of the image VAE that compresses and maps the image data
 205 into the latent space. It is used to align the image features with the agent embedding, and d_k is the
 206 dimension of the \mathbf{Fea} . Our learning target can be formalized as:

$$\arg \min_{\Theta} \sum_i^N \mathcal{D}(\hat{\mathbf{Y}}, \mathbf{Y}_{v_i}) \mathbb{1}(v_i \text{ gives an answer}) + \text{MSE}(\hat{\mathbf{Lst}}, \mathbf{Lst}) + \text{MSE}(\hat{\mathbf{A}}, \mathbf{A}), \quad (5)$$

207 where $\mathcal{D}(\cdot)$ denotes the distance between the places an LVLM agent answered and the ground truth,
 208 $\mathbb{1}(\cdot)$ is the indicator function, $\mathbf{Y}_{v_i} := \mathbf{Y}_{v_i}^{(L)} = g_{v_i}(\mathbf{X}, \mathbf{Y}_{v_j}^{(L-1)})$, $g_{v_i}(\cdot)$ represent the LVLM agent v_i
 209 with fixed parameters and $\mathbf{Y}_{v_i}^{(0)} = g_{v_i}(\mathbf{X})$ is the answer that LVLM agent v_i generates at the initial
 210 stage of discussion.

211 4 Experiments

212 To evaluate the performance of our framework, we conducted experiments on the real-world dataset
 213 that was gathered from the Internet to answer the following research questions:

- 214 • **RQ1:** Can smileGeo outperform state-of-the-art methods in open-ended geo-localization tasks?
- 215 • **RQ2:** Are LVLM agents with diverse knowledge and reasoning abilities more suitable for building
 216 a collaboration social network of agents?
- 217 • **RQ3:** How does the setting of hyperparameters affect the performance of smileGeo?

218 4.1 Experiment Setup

219 **Datasets.** In this paper, we newly construct a geo-localization dataset named GeoGlobe. It contains a
 220 variety of man-made landmarks or natural attractions from nearly 150 countries with different cultural
 221 and regional styles. The diversity and richness of GeoGlobe allow us to evaluate the performance of
 222 different models more accurately. More details can be found in Appendix B.

223 **Implementation Details.** We select both open-source and close-source LVLMs with different scales
 224 trained by different datasets as agents in the proposed framework. As for the open-source LVLMs,
 225 we utilize several open-source fine-tuned LVLMs: Infi-MM², Qwen-VL³, vip-llava-7b&13b⁴, llava-

²<https://huggingface.co/Infi-MM/infimm-zephyr>

³<https://huggingface.co/Qwen/Qwen-VL>

⁴<https://huggingface.co/llava-hf/vip-llava-xxx>

226 1.5–7b–base&mistral&vicuna⁵, llava–1.6–7b&13b&34b–mistral&vicuna⁶, CogVLM⁷. As for the
 227 closed-source LVLMs, we chose the models provided by three of the most famous companies in the
 228 world: Claude–3–opus⁸, GPT–4V⁹, and Gemini–1.5–pro¹⁰. Besides, 99% of images (about 290,000
 229 samples) from the original dataset are randomly chosen as training samples. For the open-world
 230 geolocation problem, we construct the test dataset using approximately 4,000 samples, of which
 231 nearly 66.67% samples reflected different locations not present in the training dataset. More details
 232 about the deployment of smileGeo and the related parameter settings can be found in Appendix C.

233 **Baselines.** In this work, we compare the proposed framework with three kinds of baselines: single
 234 LVLMs, LLM/LVLM-based multi-agent frameworks, and image retrieval approaches. Firstly, we use
 235 each LVLM alone as an agent directly for the geo-localization task and compute the performance of
 236 these single LVLMs under the same dataset. In addition, we experiment with multi-agent collaborative
 237 frameworks, including LLM-Blender [34], PHP [35], Reflexion [36], LLM Debate [37], and DyLAN
 238 [38]. Finally, several state-of-the-art image retrieval approaches, including NetVLAD [3], GeM
 239 [26], and CosPlace [46], are also used to be part of the baselines. We set the training dataset as the
 240 geo-tagged image database of each image retrieval system and use images in the test dataset for the
 241 retrieval system to generate answers.

242 **Evaluation Metrics.** We use *Accuracy* (Acc) to evaluate the performance: $Accuracy = \frac{N_{correct}}{N_{total}}$,
 243 where $N_{correct}$ is the number of samples that the proposed framework correctly geo-localizes, and
 244 N_{total} refers to the total number of testing samples.

245 In this paper, we first geo-encode the answers with the ground truth, *i.e.*, we transform the addresses
 246 described through natural language into latitude-longitude coordinates. Then, we calculate the
 247 distance between the two coordinates. When the distance between the two coordinates is less than
 248 $th = 50km$ (city-level), we consider the answer of the framework to be correct.

249 4.2 Performance Comparison

250 We divide the baseline comparison experiment into three parts: i) comparison with single LVLMs,
 251 ii) comparison with LLM/LVLM-based agent frameworks, and iii) comparison with image retrieval
 252 systems.

Table 1: Results of different single LVLM baselines.

	Without Web Searching			With Web Searching		
	Natural	ManMade	Overall	Natural	ManMade	Overall
Infi-MM	19.2547	21.4133	20.9883	0.9938	0.3351	0.4648
Qwen-VL	42.4845	37.4657	38.4540	4.9689	11.2093	9.9804
vip-llava-13b	20.6211	15.4127	16.4384	8.323	4.3558	5.137
vip-llava-7b	21.9876	18.4892	19.1781	31.9255	56.5032	51.6634
llava-1.5-7b	17.3913	16.3265	16.5362	27.205	47.2129	43.273
llava-1.6-7b-mistral	0.3727	0.0914	0.1468	0.8696	2.1627	1.908
llava-1.6-7b-vicuna	2.2360	2.0713	2.1037	6.9565	15.8696	14.1145
llava-1.6-13b	10.4348	8.8943	9.1977	12.1739	28.2668	25.0978
llava-1.6-34b	10.3106	9.1379	9.3689	52.795	77.1855	72.3826
CogVLM	7.7019	7.5845	7.6076	6.8323	10.3564	9.6624
claude-3-opus	22.06	37.38	16.5468	33.0435	40.7125	39.2027
GPT-4V	27.5776	35.3443	33.8145	61.9876	87.6028	82.5587
Gemini-1.5-pro	55.6522	60.3107	59.3933	62.2360	82.8206	78.7671
smileGeo	58.6111	64.3968	63.2730	78.0448	87.0069	85.2630

Bold indicates the statistically significant improvements
(i.e., two-sided t-test with $p < 0.05$) over the best baseline.

⁵<https://huggingface.co/llava-hf/llava-1.5-xxx>

⁶<https://huggingface.co/liuhaotian/llava-v1.6-xxx>

⁷<https://github.com/THUDM/CogVLM>


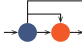
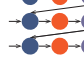
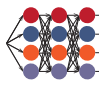
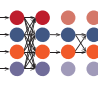
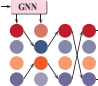
⁸<https://anthropic.com/>

⁹<https://openai.com/>

¹⁰<https://gemini.google.com/>

253 Firstly, the performance of all single LVLM baselines is shown in Table 1, in terms of the metric
 254 Acc. The data in Table 1 indicate that open-source LVLMs with diverse knowledge and reasoning
 255 capabilities exhibit significant variations, particularly in geo-localization tasks. This may be due
 256 to the difference in the overlap between the pre-training datasets used by different LVLMs and
 257 the dataset we constructed. Therefore, in addition to querying the LVLM locations about images,
 258 we also incorporated real-time image search results from Google to provide the model with more
 259 comprehensive information. These results from Internet retrievals are incorporated into the chain-of-
 260 thoughts (CoT) [47] of LVLMs as external knowledge. At this time, models with larger parameters,
 261 such as llava-1.6-34b, demonstrate superior reasoning abilities compared to smaller models (7b or
 262 13b). In addition, closed-source large models also show more consistent performance than their open-
 263 source counterparts and are more adept at analyzing and utilizing external knowledge for accurate
 264 inferences. Compared to all single LVLMs, our proposed LVLM agent framework surpasses all
 265 single LVLM baselines in accuracy. This improvement confirms the effectiveness of different LVLMs
 266 collaborating by engaging in discussions and analyzing various types of images, thus producing more
 267 precise results.

Table 2: Results of different agent frameworks without web searching.

Framework	LLM-Blender	PHP	Reflexion	LLM Debate	DyLAN	smileGeo
Structure						
Acc \uparrow	55.7802%	60.9809%	62.3412%	57.0119%	62.8187%	63.2730%
Tks \downarrow	23,662	154,520	109,524	260,756	159,320	18,876

'Acc' stands for the accuracy of the framework;

'Tks' means the average tokens a framework costs per query (including image tokens).

268 Secondly, the comparative results across various LLM/LVLM agent frameworks are presented in
 269 Table 2. It is evident that the majority of LLM/LVLM agent frameworks surpass individual LVLMs
 270 in terms of geo-localization accuracy. This improvement can primarily be attributed to the ability to
 271 integrate knowledge from multiple LVLM agents, thereby enhancing the overall precision of these
 272 frameworks. However, LLM-Blender and LLM Debate exhibit lower accuracy due to statements of
 273 some agents misleading others during discussions, which impedes the generation of correct outcomes.
 274 Our framework, smileGeo, guarantees the highest accuracy while being able to accomplish the
 275 geo-localization task with the lowest token costs. The average number of tokens our framework
 276 spent per query is 18,876, and it is less than the computational overhead of LLM-Blender (23,662),
 277 which has the simplest agent framework structure but the lowest accuracy among all baselines. This
 278 is mainly due to a 'small' GNN-based dynamic learning model being deployed for agent selection
 279 stages and significantly reducing unnecessary discussions among agents.

280 Finally, Table 3 presents the comparison between the proposed framework and existing
 281 image retrieval systems. Our framework, smileGeo, consistently outperforms all other
 282 retrieval-based approaches. This superior performance can be attributed to the fact that other image
 283 retrieval methods rely on a rich geo-tagged image database. In our test
 284 dataset, however, two-thirds of the images
 285 are new and localized in completely different areas from those in the training dataset. This highlights
 286 the shortages of conventional database-based retrieval systems due to the limitations of the geo-tagged
 287 image databases and demonstrates the effectiveness of our proposed framework in solving open-world
 288 geo-localization tasks.

Table 3: Comparison with image retrieval systems.

	Natural	ManMade	Overall
NetVLAD	26.5134	28.9955	28.6047
GeM	23.1022	25.4175	25.0749
CosPlace	28.1688	30.2782	29.8701
smileGeo	58.6111	64.3968	63.2730

Bold indicates the statistically significant improvements (i.e., two-sided t-test with $p < 0.05$) over the best baseline.

293 4.3 Ablation Study

294 **Number of Agents.** We further demonstrate the relationships between the number of agents and the
 295 framework performance. We conduct experiments in two ways: i) by calling the same closed-source
 296 LVLM API (Here, we use Gemini-1.5-pro because it performs best without the help of the Internet)
 297 under different prompts (e.g., You are good at recognizing natural attractions; You're a traveler around

298 Europe) to simulate different agents, and ii) by using different LVLM backbones to represent distinct
 299 agents. The results are shown in Figure 2.

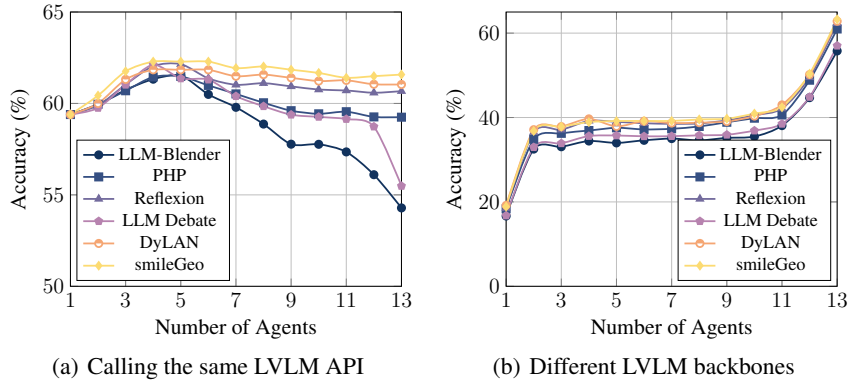


Figure 2: Results of model performance in relation to the number of agents.

300 As illustrated in Figure 2(a), the framework achieves optimal accuracy with 4 or 5 agents. Beyond
 301 this number, the framework’s performance begins to deteriorate. This shows that using models
 302 with the same knowledge and reasoning capabilities as different agents has limited improvement
 303 in the accuracy of the framework. Despite this decline, the performance of frameworks other than
 304 LLM-Blender and LLM Debate remains superior to that of a single agent. LLM-Blender and LLM
 305 Debate, however, have a significant decrease in model accuracy when the number of agents exceeds
 306 11. This is mainly because both of them involve all LVLMs in every discussion, which suffers from
 307 excessive repetitive and redundant discussions. Figure 2(b) reveals that the accuracy of the framework
 308 improves with the incorporation of more LVLM backbones, indicating that the diversity of LVLMs
 309 can enhance the quality of discussions.

310 **Hyperparameter K & R .** There are two hyperpa-
 311 rameters, K and R , that need to be pre-defined in the
 312 proposed framework: K is the number of agents (an-
 313 swer agents) that respond in each round of discussion,
 314 and R is the number of agents (review agents) used
 315 to review answers from answer agents. Therefore, we
 316 conduct experiments under different combinations of
 317 $K \in [1, 8]$ and $R \in [1, 8]$, as shown in Figure 3. The re-
 318 sults indicate that optimal performance can be achieved
 319 with relatively small values of K or R . However, the
 320 computational cost, measured in tokens, increases expo-
 321 nentially with higher values of K and R . To balance
 322 both the efficiency and the accuracy of smileGeo, for
 323 the experiments presented in this paper, we set both K and R equal to 2.

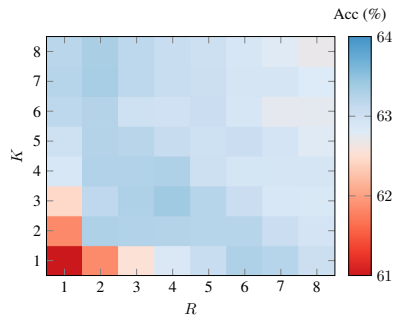


Figure 3: Results under different K and R .

324 5 Conclusion

325 This work introduces a novel LVLM agent framework, smileGeo, specifically designed for geo-
 326 localization tasks. Inspired by the review mechanism, it integrates various LVLMs to discuss
 327 anonymously and geo-localize images worldwide. Additionally, we have developed a dynamic
 328 learning strategy for agent collaboration social networks, electing appropriate agents to geo-localize
 329 each image with different characteristics. This enhancement reduces the computational burden
 330 associated with collaborative discussions among LVLM agents. Moreover, we have constructed a
 331 geo-localization dataset called GeoGlobe and will open-source it. Overall, smileGeo demonstrates
 332 significant improvements in geo-localization tasks, achieving superior performance with lower
 333 computational demands compared to contemporary state-of-the-art LLM/LVLM agent frameworks.

334 Looking ahead, we aim to expand the capabilities of smileGeo to incorporate more powerful external
 335 tools beyond just web searching. Additionally, we plan to explore extending its application to complex
 336 scenarios, such as high-precision global positioning and navigation for robots, laying the cornerstone
 337 for exploring LVLM agent collaboration to handle different complex open-world tasks efficiently.

References

- 338
- 339 [1] B. Huang and K. M. Carley, "A large-scale empirical study of geotagging behavior on twitter," in
340 *ASONAM '19: International Conference on Advances in Social Networks Analysis and Mining,*
341 *Vancouver, British Columbia, Canada, 27-30 August, 2019*, F. Spezzano, W. Chen, and X. Xiao,
342 Eds. ACM, 2019, pp. 365–373. [Online]. Available: <https://doi.org/10.1145/3341161.3342870>
- 343 [2] J. Luo, D. Joshi, J. Yu, and A. C. Gallagher, "Geotagging in multimedia and computer vision
344 - a survey," *Multim. Tools Appl.*, vol. 51, no. 1, pp. 187–211, 2011. [Online]. Available:
345 <https://doi.org/10.1007/s11042-010-0623-y>
- 346 [3] R. Arandjelovic, P. Gronát, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: CNN architecture for
347 weakly supervised place recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 6,
348 pp. 1437–1451, 2018. [Online]. Available: <https://doi.org/10.1109/TPAMI.2017.2711011>
- 349 [4] M. Zaffar, S. Garg, M. Milford, J. F. P. Kooij, D. Flynn, K. D. McDonald-Maier, and S. Ehsan,
350 "Vpr-bench: An open-source visual place recognition evaluation framework with quantifiable
351 viewpoint and appearance change," *Int. J. Comput. Vis.*, vol. 129, no. 7, pp. 2136–2174, 2021.
352 [Online]. Available: <https://doi.org/10.1007/s11263-021-01469-5>
- 353 [5] A. Torii, R. Arandjelovic, J. Sivic, M. Okutomi, and T. Pajdla, "24/7 place recognition by view
354 synthesis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 2, pp. 257–271, 2018. [Online].
355 Available: <https://doi.org/10.1109/TPAMI.2017.2667665>
- 356 [6] Z. Chen, A. Jacobson, N. Sünderhauf, B. Upcroft, L. Liu, C. Shen, I. D. Reid, and
357 M. Milford, "Deep learning features at scale for visual place recognition," in *2017*
358 *IEEE International Conference on Robotics and Automation, ICRA 2017, Singapore,*
359 *Singapore, May 29 - June 3, 2017*. IEEE, 2017, pp. 3223–3230. [Online]. Available:
360 <https://doi.org/10.1109/ICRA.2017.7989366>
- 361 [7] Z. Chen, L. Liu, I. Sa, Z. Ge, and M. Chli, "Learning context flexible attention model for
362 long-term visual place recognition," *IEEE Robotics Autom. Lett.*, vol. 3, no. 4, pp. 4015–4022,
363 2018. [Online]. Available: <https://doi.org/10.1109/LRA.2018.2859916>
- 364 [8] Z. Chen, F. Maffra, I. Sa, and M. Chli, "Only look once, mining distinctive landmarks from
365 convnet for visual place recognition," in *2017 IEEE/RSJ International Conference on Intelligent*
366 *Robots and Systems, IROS 2017, Vancouver, BC, Canada, September 24-28, 2017*. IEEE,
367 2017, pp. 9–16. [Online]. Available: <https://doi.org/10.1109/IROS.2017.8202131>
- 368 [9] S. Garg, N. Sünderhauf, and M. Milford, "Semantic-geometric visual place recognition: a new
369 perspective for reconciling opposing views," *Int. J. Robotics Res.*, vol. 41, no. 6, pp. 573–598,
370 2022. [Online]. Available: <https://doi.org/10.1177/0278364919839761>
- 371 [10] S. Hausler, A. Jacobson, and M. Milford, "Multi-process fusion: Visual place recognition using
372 multiple image processing methods," *IEEE Robotics Autom. Lett.*, vol. 4, no. 2, pp. 1924–1931,
373 2019. [Online]. Available: <https://doi.org/10.1109/LRA.2019.2898427>
- 374 [11] A. Khaliq, S. Ehsan, Z. Chen, M. Milford, and K. D. McDonald-Maier, "A holistic visual
375 place recognition approach using lightweight cnns for significant viewpoint and appearance
376 changes," *IEEE Trans. Robotics*, vol. 36, no. 2, pp. 561–569, 2020. [Online]. Available:
377 <https://doi.org/10.1109/TRO.2019.2956352>
- 378 [12] M. M. ElQadi, M. Lesiv, A. G. Dyer, and A. Dorin, "Computer vision-enhanced selection of
379 geo-tagged photos on social network sites for land cover classification," *Environ. Model. Softw.*,
380 vol. 128, p. 104696, 2020. [Online]. Available: <https://doi.org/10.1016/j.envsoft.2020.104696>
- 381 [13] M. Campbell and M. Wheeler, "A vision based geolocation tracking system for uav's," in *AIAA*
382 *Guidance, Navigation, and Control Conference and Exhibit*, 2006, p. 6246.
- 383 [14] F. Deng, L. Zhang, F. Gao, H. Qiu, X. Gao, and J. Chen, "Long-range binocular
384 vision target geolocation using handheld electronic devices in outdoor environment,"
385 *IEEE Trans. Image Process.*, vol. 29, pp. 5531–5541, 2020. [Online]. Available:
386 <https://doi.org/10.1109/TIP.2020.2984898>

- 387 [15] L. Zhang, F. Deng, J. Chen, Y. Bi, S. K. Phang, X. Chen, and B. M. Chen,
388 “Vision-based target three-dimensional geolocation using unmanned aerial vehicles,” *IEEE*
389 *Trans. Ind. Electron.*, vol. 65, no. 10, pp. 8052–8061, 2018. [Online]. Available:
390 <https://doi.org/10.1109/TIE.2018.2807401>
- 391 [16] X. Feng, Z.-Y. Chen, Y. Qin, Y. Lin, X. Chen, Z. Liu, and J.-R. Wen, “Large language model-
392 based human-agent collaboration for complex task solving,” *arXiv preprint arXiv:2402.12914*,
393 2024.
- 394 [17] W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song *et al.*,
395 “Cogvlm: Visual expert for pretrained language models,” *arXiv preprint arXiv:2311.03079*,
396 2023.
- 397 [18] V. Paolicelli, G. M. Berton, F. Montagna, C. Masone, and B. Caputo, “Adaptive-attentive
398 geolocation from few queries: A hybrid approach,” *Frontiers Comput. Sci.*, vol. 4, p.
399 841817, 2022. [Online]. Available: <https://doi.org/10.3389/fcomp.2022.841817>
- 400 [19] Y. Ge, H. Wang, F. Zhu, R. Zhao, and H. Li, “Self-supervising fine-grained region similarities
401 for large-scale image localization,” in *Computer Vision - ECCV 2020 - 16th European*
402 *Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part IV*, ser. Lecture Notes in
403 Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12349. Springer,
404 2020, pp. 369–386. [Online]. Available: https://doi.org/10.1007/978-3-030-58548-8_22
- 405 [20] H. Jin Kim, E. Dunn, and J.-M. Frahm, “Learned contextual feature reweighting for image
406 geo-localization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern*
407 *Recognition*, 2017, pp. 2136–2145.
- 408 [21] L. Liu, H. Li, and Y. Dai, “Stochastic attraction-repulsion embedding for large scale image
409 localization,” in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019,*
410 *Seoul, Korea (South), October 27 - November 2, 2019.* IEEE, 2019, pp. 2570–2579. [Online].
411 Available: <https://doi.org/10.1109/ICCV.2019.00266>
- 412 [22] F. Warburg, S. Hauberg, M. Lopez-Antequera, P. Gargallo, Y. Kuang, and J. Civera, “Mapillary
413 street-level sequences: A dataset for lifelong place recognition,” in *Proceedings of the IEEE/CVF*
414 *conference on computer vision and pattern recognition*, 2020, pp. 2626–2635.
- 415 [23] G. Peng, Y. Yue, J. Zhang, Z. Wu, X. Tang, and D. Wang, “Semantic reinforced attention
416 learning for visual place recognition,” in *IEEE International Conference on Robotics and*
417 *Automation, ICRA 2021, Xi’an, China, May 30 - June 5, 2021.* IEEE, 2021, pp. 13 415–13 422.
418 [Online]. Available: <https://doi.org/10.1109/ICRA48506.2021.9561812>
- 419 [24] S. Ibrahim, N. van Noord, T. Alpherts, and M. Worring, “Inside out visual
420 place recognition,” in *32nd British Machine Vision Conference 2021, BMVC 2021,*
421 *Online, November 22-25, 2021.* BMVA Press, 2021, p. 362. [Online]. Available:
422 <https://www.bmvc2021-virtualconference.com/assets/papers/0467.pdf>
- 423 [25] S. Hausler, S. Garg, M. Xu, M. Milford, and T. Fischer, “Patch-netvlad: Multi-scale
424 fusion of locally-global descriptors for place recognition,” in *IEEE Conference on*
425 *Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021.*
426 Computer Vision Foundation / IEEE, 2021, pp. 14 141–14 152. [Online]. Available: https://openaccess.thecvf.com/content/CVPR2021/html/Hausler_Patch-NetVLAD_Multi-Scale_Fusion_of_Locally-Global_Descriptors_for_Place_Recognition_CVPR_2021_paper.html
- 429 [26] F. Radenovic, G. Toliás, and O. Chum, “Fine-tuning CNN image retrieval with no human
430 annotation,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 7, pp. 1655–1668, 2019.
431 [Online]. Available: <https://doi.org/10.1109/TPAMI.2018.2846566>
- 432 [27] M. Izbicki, E. E. Papalexakis, and V. J. Tsotras, “Exploiting the earth’s spherical geometry to
433 geolocate images,” in *Machine Learning and Knowledge Discovery in Databases - European*
434 *Conference, ECML PKDD 2019, Würzburg, Germany, September 16-20, 2019, Proceedings,*
435 *Part II*, ser. Lecture Notes in Computer Science, U. Brefeld, É. Fromont, A. Hotho, A. J.
436 Knobbe, M. H. Maathuis, and C. Robardet, Eds., vol. 11907. Springer, 2019, pp. 3–19.
437 [Online]. Available: https://doi.org/10.1007/978-3-030-46147-8_1

- 438 [28] G. Kordopatis-Zilos, P. Galopoulos, S. Papadopoulos, and I. Kompatsiaris, “Leveraging
439 efficientnet and contrastive learning for accurate global-scale location estimation,” in *ICMR*
440 *'21: International Conference on Multimedia Retrieval, Taipei, Taiwan, August 21-24, 2021*,
441 W. Cheng, M. S. Kankanhalli, M. Wang, W. Chu, J. Liu, and M. Worring, Eds. ACM, 2021,
442 pp. 155–163. [Online]. Available: <https://doi.org/10.1145/3460426.3463644>
- 443 [29] E. Müller-Budack, K. Pustu-Iren, and R. Ewerth, “Geolocation estimation of photos
444 using a hierarchical model and scene classification,” in *Computer Vision - ECCV 2018 -*
445 *15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part*
446 *XII*, ser. Lecture Notes in Computer Science, V. Ferrari, M. Hebert, C. Sminchisescu,
447 and Y. Weiss, Eds., vol. 11216. Springer, 2018, pp. 575–592. [Online]. Available:
448 https://doi.org/10.1007/978-3-030-01258-8_35
- 449 [30] P. H. Seo, T. Weyand, J. Sim, and B. Han, “Cplanet: Enhancing image geolocalization
450 by combinatorial partitioning of maps,” in *Computer Vision - ECCV 2018 - 15th*
451 *European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part*
452 *X*, ser. Lecture Notes in Computer Science, V. Ferrari, M. Hebert, C. Sminchisescu,
453 and Y. Weiss, Eds., vol. 11214. Springer, 2018, pp. 544–560. [Online]. Available:
454 https://doi.org/10.1007/978-3-030-01249-6_33
- 455 [31] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
456 K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell,
457 P. Welinder, P. F. Christiano, J. Leike, and R. Lowe, “Training language models to follow
458 instructions with human feedback,” in *Advances in Neural Information Processing Systems*
459 *35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New*
460 *Orleans, LA, USA, November 28 - December 9, 2022*, S. Koyejo, S. Mohamed, A. Agarwal,
461 D. Belgrave, K. Cho, and A. Oh, Eds., 2022. [Online]. Available: [http://papers.nips.cc/paper_](http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html)
462 [files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html)
- 463 [32] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T.
464 Lee, Y. Li, S. M. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro, and Y. Zhang, “Sparks of
465 artificial general intelligence: Early experiments with GPT-4,” *CoRR*, vol. abs/2303.12712,
466 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2303.12712>
- 467 [33] R. Schaeffer, B. Miranda, and S. Koyejo, “Are emergent abilities of large language models
468 a mirage?” in *Advances in Neural Information Processing Systems 36: Annual Conference*
469 *on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA,*
470 *December 10 - 16, 2023*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and
471 S. Levine, Eds., 2023. [Online]. Available: [http://papers.nips.cc/paper_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/adc98a266f45005c403b8311ca7e8bd7-Abstract-Conference.html)
472 [adc98a266f45005c403b8311ca7e8bd7-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/adc98a266f45005c403b8311ca7e8bd7-Abstract-Conference.html)
- 473 [34] D. Jiang, X. Ren, and B. Y. Lin, “Llm-blender: Ensembling large language models
474 with pairwise ranking and generative fusion,” in *Proceedings of the 61st Annual Meeting*
475 *of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023,*
476 *Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds.
477 Association for Computational Linguistics, 2023, pp. 14 165–14 178. [Online]. Available:
478 <https://doi.org/10.18653/v1/2023.acl-long.792>
- 479 [35] C. Zheng, Z. Liu, E. Xie, Z. Li, and Y. Li, “Progressive-hint prompting improves
480 reasoning in large language models,” *CoRR*, vol. abs/2304.09797, 2023. [Online]. Available:
481 <https://doi.org/10.48550/arXiv.2304.09797>
- 482 [36] N. Shinn, F. Cassano, A. Gopinath, K. Narasimhan, and S. Yao, “Reflexion: language agents
483 with verbal reinforcement learning,” in *Advances in Neural Information Processing Systems*
484 *36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New*
485 *Orleans, LA, USA, December 10 - 16, 2023*, A. Oh, T. Naumann, A. Globerson, K. Saenko,
486 M. Hardt, and S. Levine, Eds., 2023. [Online]. Available: [http://papers.nips.cc/paper_files/](http://papers.nips.cc/paper_files/paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html)
487 [paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html)
- 488 [37] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch, “Improving factuality and
489 reasoning in language models through multiagent debate,” *CoRR*, vol. abs/2305.14325, 2023.
490 [Online]. Available: <https://doi.org/10.48550/arXiv.2305.14325>

- 491 [38] Z. Liu, Y. Zhang, P. Li, Y. Liu, and D. Yang, “Dynamic llm-agent network: An llm-agent
492 collaboration framework with agent team optimization,” *CoRR*, vol. abs/2310.02170, 2023.
493 [Online]. Available: <https://doi.org/10.48550/arXiv.2310.02170>
- 494 [39] W. Shi, S. Min, M. Yasunaga, M. Seo, R. James, M. Lewis, L. Zettlemoyer, and W. Yih,
495 “REPLUG: retrieval-augmented black-box language models,” *CoRR*, vol. abs/2301.12652, 2023.
496 [Online]. Available: <https://doi.org/10.48550/arXiv.2301.12652>
- 497 [40] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. R. Narasimhan, and Y. Cao, “React: Synergizing
498 reasoning and acting in language models,” in *The Eleventh International Conference on
499 Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net,
500 2023. [Online]. Available: https://openreview.net/pdf?id=WE_vluYUL-X
- 501 [41] G. Izacard, P. S. H. Lewis, M. Lomeli, L. Hosseini, F. Petroni, T. Schick, J. Dwivedi-Yu,
502 A. Joulin, S. Riedel, and E. Grave, “Atlas: Few-shot learning with retrieval augmented
503 language models,” *J. Mach. Learn. Res.*, vol. 24, pp. 251:1–251:43, 2023. [Online]. Available:
504 <http://jmlr.org/papers/v24/23-0037.html>
- 505 [42] T. Schick, J. Dwivedi-Yu, R. Dessi, R. Raileanu, M. Lomeli, E. Hambro, L. Zettlemoyer,
506 N. Cancedda, and T. Scialom, “Toolformer: Language models can teach themselves to
507 use tools,” in *Advances in Neural Information Processing Systems 36: Annual Conference
508 on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA,
509 December 10 - 16, 2023*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and
510 S. Levine, Eds., 2023. [Online]. Available: [http://papers.nips.cc/paper_files/paper/2023/hash/
511 d842425e4bf79ba039352da0f658a906-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/d842425e4bf79ba039352da0f658a906-Abstract-Conference.html)
- 512 [43] P. Lu, B. Peng, H. Cheng, M. Galley, K. Chang, Y. N. Wu, S. Zhu, and J. Gao,
513 “Chameleon: Plug-and-play compositional reasoning with large language models,” in
514 *Advances in Neural Information Processing Systems 36: Annual Conference on Neural
515 Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December
516 10 - 16, 2023*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and
517 S. Levine, Eds., 2023. [Online]. Available: [http://papers.nips.cc/paper_files/paper/2023/hash/
518 871ed095b734818cfba48db6aeb25a62-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/871ed095b734818cfba48db6aeb25a62-Abstract-Conference.html)
- 519 [44] L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig,
520 “PAL: program-aided language models,” in *International Conference on Machine Learning,
521 ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, ser. Proceedings of Machine
522 Learning Research, A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and
523 J. Scarlett, Eds., vol. 202. PMLR, 2023, pp. 10764–10799. [Online]. Available:
524 <https://proceedings.mlr.press/v202/gao23f.html>
- 525 [45] X. Wang, S. Li, and H. Ji, “Code4struct: Code generation for few-shot structured
526 prediction from natural language,” *CoRR*, vol. abs/2210.12810, 2022. [Online]. Available:
527 <https://doi.org/10.48550/arXiv.2210.12810>
- 528 [46] G. M. Berton, C. Masone, and B. Caputo, “Rethinking visual geo-localization for large-scale
529 applications,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR
530 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 4868–4878. [Online].
531 Available: <https://doi.org/10.1109/CVPR52688.2022.00483>
- 532 [47] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le,
533 and D. Zhou, “Chain-of-thought prompting elicits reasoning in large language models,” in
534 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural
535 Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November
536 28 - December 9, 2022*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and
537 A. Oh, Eds., 2022. [Online]. Available: [http://papers.nips.cc/paper_files/paper/2022/hash/
538 9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html)

539 **A Notations**

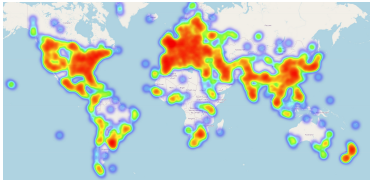
540 We summarize all notations in this paper and list them in Table 4.

Table 4: Notations in this paper.

Notation	Description
X	The image to be recognized.
$Y (\hat{Y})$	The predicted (ground truth of) geospatial location in the natural language form.
$\mathcal{G} (\hat{\mathcal{G}})$	The predicted (ground truth of) LVLM-based agent collaboration social network.
$A (\hat{A})$	The predicted (ground truth of) adjacency matrix of the agent social network.
$Lst (\hat{Lst})$	The predicted (ground truth of) scalar of agent election probability.
\mathcal{V}	The set of LLM agents.
\mathcal{E}	The set of connections between LLM agents.
N	The number of agents.
K	The number of agents to be elected as answer agent(s).
R	The number of agents to be selected as review agent(s).
L	The number of agent discussion rounds.
Z	The maximum number of rounds in which answer agents harmonize opinions.
Θ	The learnable parameters of the agent social network learning model.

541 **B Dataset Details**

542 The images in this dataset are copyright-free images obtained from the Internet via a crawler. We
 543 divide the images into two main categories: man-made landmarks as well as natural attractions. Then,
 544 we filter out the data samples that could clearly identify the locations of the landmarks or attractions
 545 in the images. As a result, we filter out nearly three hundred thousand data samples, and please
 546 refer to Table 5 and Figure 4 for details. Due to the fact that a large number of natural attractions in
 547 different geographical regions with high similarity are cleaned, the magnitude of the data related to
 548 natural attractions in this dataset is smaller than that of man-made attractions.



549 Figure 4: The data distribution around the world.

Table 5: Statistics of the dataset GeoGlobe.

	Images	Cities	Countries	Attractions
Man-made	253,118	2,313	143	10,492
Natural	40,087	1,044	97	1,849

550 For an open-world geo-localization task, the relationship between the training and test samples in
 551 the experiment could greatly affect the results. We label the training samples as $\mathcal{Z}_{\text{train}}$, and the test
 552 sample set as $\mathcal{Z}_{\text{test}}$, and use two metrics, *coverage* as well as *consistency*, to portray this relationship:

$$\begin{aligned}
 \text{coverage} &= \frac{\mathcal{Z}_{\text{train}} \cap \mathcal{Z}_{\text{test}}}{\mathcal{Z}_{\text{train}}} \times 100\% \\
 \text{consistency} &= \frac{\mathcal{Z}_{\text{train}} \cap \mathcal{Z}_{\text{test}}}{\mathcal{Z}_{\text{test}}} \times 100\%
 \end{aligned} \tag{6}$$

553 As for the samples in this paper, *coverage* \approx 4.6564%, and *consistency* \approx 33.2957%.

554 **C Implementation Details**

555 In all experiments, we employ a variety of LVLMs, encompassing both open-source and closed-source
 556 models, to be agents in the proposed framework. Unless specified otherwise, zero-shot prompting is
 557 applied. Each open-source LVLM is deployed on a dedicated A800 (80G) GPU server with 200GB
 558 memory. As for each closed-source LVLM, we cost amounting to billions of tokens by calling APIs
 559 as specified by the official website. To avoid the context length issue that occurs in some LVLMs, we
 560 truncate the context before submitting it to the agent for questions based on the maximum number of

Algorithm 1 The smileGeo framework

Input: A set of pre-trained LLMs $\mathcal{V} = \{v_1, v_2, \dots\}$, the input image \mathbf{X} , and the ground truth $\hat{\mathbf{Y}}$ (if has);

Output: The geospatial location \mathbf{Y} .

Initialization Stage:

- 1: Initialize (Load) the parameter of the agent selection model: Θ
- 2: Calculate: $\mathbf{A} \leftarrow f(\mathbf{X}, \mathcal{V}|\Theta)$
- 3: Initialize the agent collaboration social network: \mathcal{G}
- 4: Calculate: $\mathbf{Lst} \leftarrow f(\mathbf{X}, \mathcal{G}|\Theta)$

Stage 1:

- 5: Elect answer agents: $\mathcal{V}^1 = \{v_a^1, v_b^1, \dots\} \leftarrow \mathbf{Lst}$, where $|\mathcal{V}^1| = K$
- 6: **for** each answer agent v^1 **do**
- 7: Obtain the location: $\mathbf{Y}_{v^1}^1 \leftarrow \text{Ask}_{v^1}(\mathbf{X})$
- 8: Get the confidence percentage: $C_{v^1}^1 \leftarrow \text{Ask}_{v^1}(\mathbf{X}, \mathbf{Y}_{v^1}^1)$
- 9: Store the further explanation: $T_{v^1}^1 \leftarrow \text{Ask}_{v^1}(\mathbf{X}, \mathbf{Y}_{v^1}^1)$
- 10: **end for**

Stage 2:

- 11: **for** each selected answer agent v^1 **do**
- 12: Select the review agents: $\mathcal{V}^2 = \{v_a^2, v_b^2, \dots\} \leftarrow \text{RandomWalk}_{v^1}(\mathcal{G})$, where $|\mathcal{V}^2| = R$
- 13: **for** each review agent v^2 **do**
- 14: Obtain the comment $T_{v^2}^2 \leftarrow \text{Review}_{v^2}(\mathbf{X}, \mathbf{Y}_{v^1}^1, C_{v^1}^1)$
- 15: Get the confidence percentage: $C_{v^2}^2 \leftarrow \text{Ask}_{v^2}(\mathbf{X}, T_{v^2}^2)$
- 16: **end for**
- 17: **end for**

Stage 3:

- 18: **for** each selected answer agent v^1 **do**
- 19: Summary the final answer: $\mathbf{Y}_{v^1}^3 \leftarrow \text{Summary}_{v^1}(\mathbf{Y}_{v^1}^1, C_{v^1}^1, T_{v^1}^2, C_{v^1}^2, T_{v^1}^2, C_{v^1}^2, \dots)$
- 20: Get the final confidence percentage: $C_{v^1}^3 \leftarrow \text{Ask}_{v^1}(\mathbf{Y}_{v^1}^1, C_{v^1}^1, T_{v^1}^2, C_{v^1}^2, T_{v^1}^2, C_{v^1}^2, \dots)$
- 21: **end for**
- 22: Generate the final answer: $\mathbf{Y} \leftarrow \text{Discussion}_Z(\mathbf{Y}_{v^1}^3, C_{v^1}^3, \mathbf{Y}_{v^2}^3, C_{v^2}^3, \dots)$

The dynamic learning strategy module:

- 23: Initialize $\mathbf{Lst}^{(0)}, \mathcal{G}^{(0)}$
- 24: **for** round l in total L rounds **do**
- 25: **for** each selected answer agent v^1 **do**
- 26: Obtain coordinates: $\text{Coors} \leftarrow \text{GeoEmb}(\mathbf{Y}_{v^1}^3), \text{Coors}_{\text{Truth}} \leftarrow \text{GeoEmb}(\mathbf{Y}_{\text{Truth}})$
- 27: **if** $\text{Dis}(\text{Coors}, \text{Coors}_{\text{Truth}}) \leq th$ **then**
- 28: $\mathbf{A}^{(l)} \leftarrow \text{Enhance}(e|e \text{ contains } v^1, e \in \mathcal{E})$
- 29: Update $\mathbf{Lst}^{(l)}[v^1] = 1$
- 30: **else**
- 31: $\mathbf{A}^{(l)} \leftarrow \text{Weaken}(e|e \text{ contains } v^1, e \in \mathcal{E})$
- 32: Update $\mathbf{Lst}^{(l)}[v^1] = 0$
- 33: **end if**
- 34: **end for**
- 35: **end for**
- 36: $\hat{\mathbf{A}} \approx \mathbf{A}^{(L)}, \hat{\mathbf{Lst}} \approx \mathbf{Lst}^{(L)}$
- 37: Update: $\Theta \leftarrow \text{Loss}(\hat{\mathbf{Y}}, \mathbf{Y}, \hat{\mathbf{A}}, \mathbf{A}, \hat{\mathbf{Lst}}, \mathbf{Lst})$

561 tokens that each agent supports. Besides, noting that images are token consuming, we only keep the
562 freshest response for agent discussions.

563 The detailed algorithm of smileGeo is illustrated in Algorithm 1. In the initialization stage, we
564 initialize or load the parameters of the agent social network learning model, as delineated in line 1.
565 Next, we treat each LVL agent as a node, establishing the LVL agent collaboration social network
566 and computing the adjacency relationships among LVL agents as well as the probability that each
567 agent is suited for responding to image \mathbf{X} , as shown in line 2. Then, line 3 initializes the agent

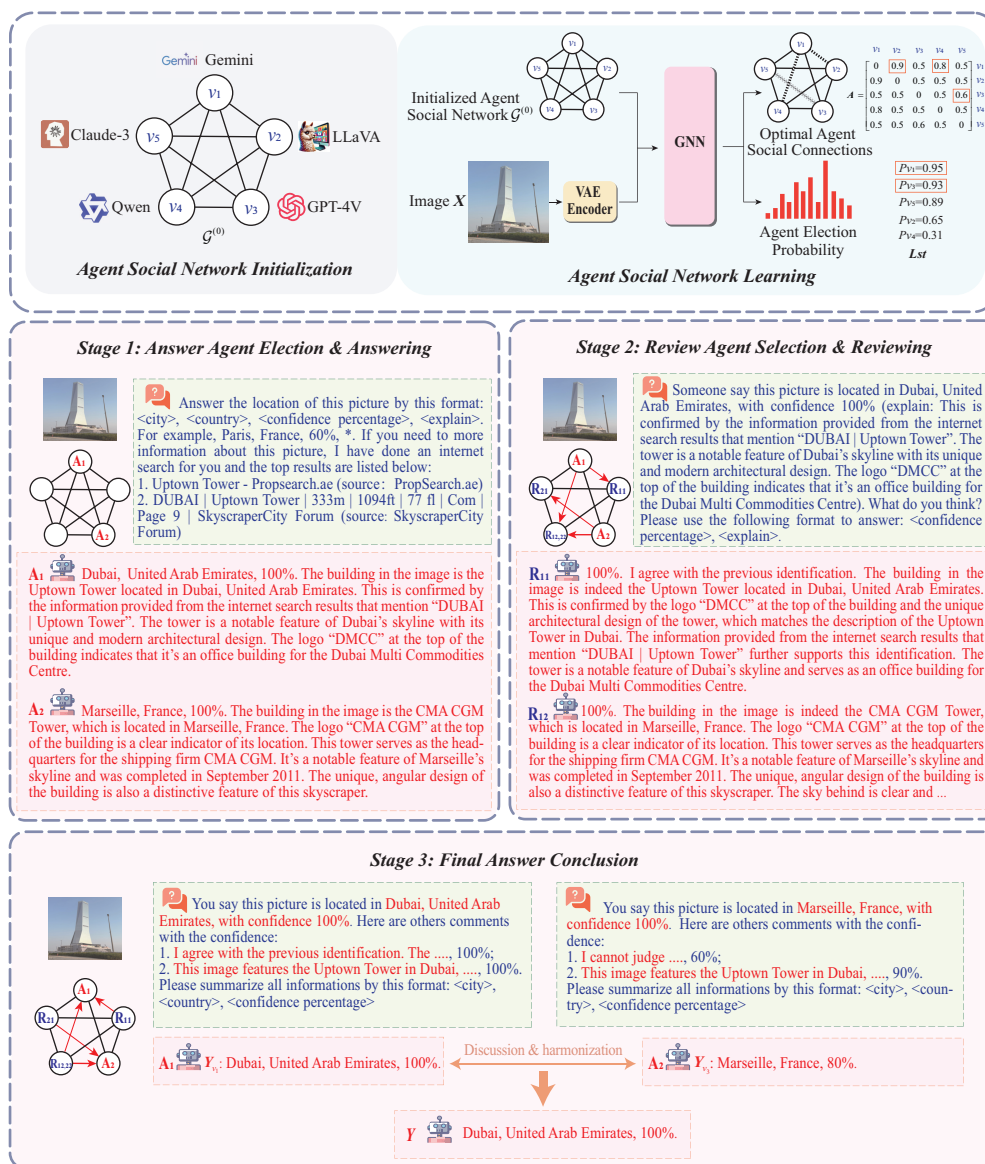


Figure 5: A case study on the geo-localization process via a given image.

568 collaboration social network and line 4 computes the agent election probability. In Stage 1, line 5
 569 involves electing appropriate answer agents based on the calculated probabilities. Subsequently, lines
 570 6-10 detail the process through which each chosen answer agent formulates their response. Stage 2
 571 begins by employing the random walk algorithm to assign review agents to each answer agent, as
 572 depicted in lines 11-12. Lines 13-16 then describe how these review agents generate feedback based
 573 on the answers provided. In Stage 3, each answer agent consolidates feedback from their assigned
 574 review agents to finalize their response, as illustrated in lines 18-21. Line 22 concludes the final
 575 answer with up to Z rounds (we set $Z = 10$ in experiments) of intra-discussion among all answer
 576 agents only. The dynamic learning strategy module involves L -round (we set $L = 20$ in experiments)
 577 comparing the generated answers against the ground truth and updating the connections between the
 578 answer and review agents accordingly, as shown in lines 23-36. In line 37, the process concludes
 579 with the updating of the learning parameters of the dynamic agent social network learning model.

580 Here, for the agent social network learning model, we first deflate each image to be recognized to
 581 512x512 pixels and then use the pre-trained VAE model¹¹ to compress the image again (compression

¹¹<https://huggingface.co/stabilityai/sd-vae-ft-mse>

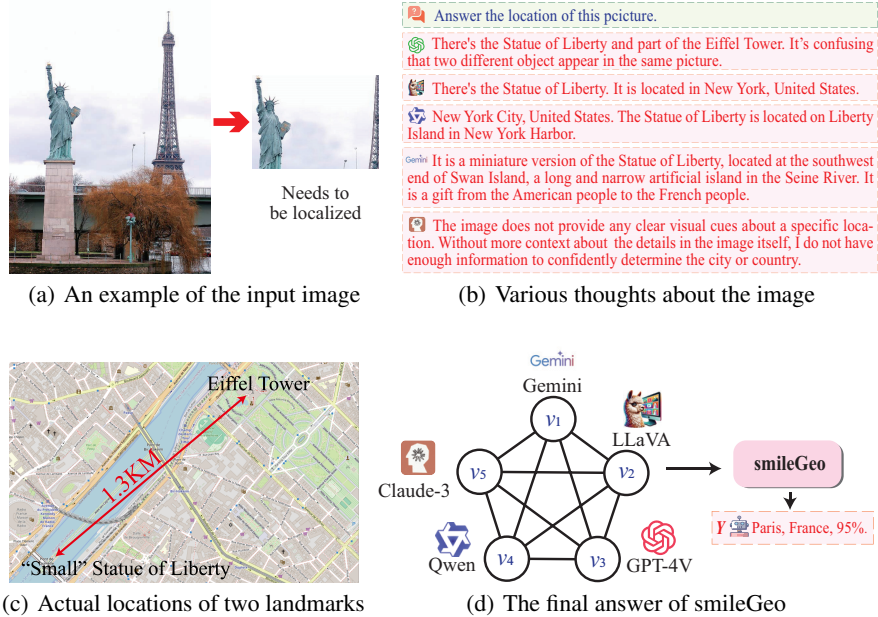


Figure 6: A case study illustrating the reasoning capabilities of smileGeo.

582 ratio 1:8) and extract its representations. We define the embedding dimension of the nodes to be 1024
 583 and the hidden layer dimension of the network layer to be 1024. we use Adam as an optimizer for
 584 gradient descent with a learning rate of $1e^{-5}$. For each stage of the LVL agent discussion, we use a
 585 uniform template to ask questions to different LVL agents and ask them to make a response in the
 586 specified format. In addition, the performance of our proposed framework is the average of the last
 587 100 epochs in a total training of 2500 epochs.

588 D Additional Experiments

589 D.1 Case Study

590 **Case 1:** In Figure 5, we illustrate the application of smileGeo in a visual geo-localization task.
 591 For this demonstration, we randomly select an image from the test dataset and employ five distinct
 592 LVLs: LLaVA, GPT-4V, Claude-3, Gemini, and Qwen. The agent selection model selects two
 593 answer agents, as depicted in the top part of the figure. Subsequently, stages 1 through 3 detail the
 594 process of generating the accurate geo-location. Initially, only one answer agent provided the correct
 595 response. However, after several rounds of discussion, the agent that initially responded incorrectly
 596 revised the confidence level of its answer. During the final internal discussion, this agent aligned its
 597 response with the correct answer. This outcome validates the efficacy of our proposed framework,
 598 demonstrating its ability to integrate the knowledge and reasoning capabilities of different agents to
 599 enhance the overall performance of the proposed LVL agent framework.

600 **Case 2:** This case study illustrates the need to pinpoint the geographical location of a complete
 601 image based on only a portion of it, as demonstrated in 6(a). As illustrated in Figure 6(b), all agents
 602 recognized the Statue of Liberty in Figure 6(a), and some identified the presence of part of the Eiffel
 603 Tower at the edge of the picture. For instance, GPT-4V concluded that the buildings in these two
 604 locations appeared in the same image. However, as is known through the knowledge of other agents
 605 (Gemini), a scaled-down version of the Statue of Liberty has been erected on Swan Island, an artificial
 606 island in the Seine River in France. By marking both the Eiffel Tower and the island on the Open
 607 Street Map (OSM) manually, as shown in Figure 6(c), it is evident that they are merely 1.3 kilometers
 608 apart in a straight line. By utilizing the proposed framework, agents discuss and summarize the
 609 location depicted in the picture to be Paris, France, as shown in Figure 6(d). Thus, without human
 610 intervention, this framework demonstrates the effectiveness of doing geo-localization tasks.

611 **NeurIPS Paper Checklist**

612 **1. Claims**

613 Question: Do the main claims made in the abstract and introduction accurately reflect the
614 paper's contributions and scope?

615 Answer: [\[Yes\]](#)

616 Justification: Our work proposes a swarm intelligence geo-localization framework, smileGeo,
617 which contains the process of the review mechanism in agent discussions along with a
618 dynamic learning strategy of agent collaboration social network, to achieve open-world
619 geo-localization tasks. In addition, we construct a novel geo-localization dataset, GeoGlobe
620 for evaluation and it will be public. All of the contributions we claimed in both abstract and
621 introduction.

622 Guidelines:

- 623 • The answer NA means that the abstract and introduction do not include the claims
624 made in the paper.
- 625 • The abstract and/or introduction should clearly state the claims made, including the
626 contributions made in the paper and important assumptions and limitations. A No or
627 NA answer to this question will not be perceived well by the reviewers.
- 628 • The claims made should match theoretical and experimental results, and reflect how
629 much the results can be expected to generalize to other settings.
- 630 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
631 are not attained by the paper.

632 **2. Limitations**

633 Question: Does the paper discuss the limitations of the work performed by the authors?

634 Answer: [\[Yes\]](#)

635 Justification: At present, the LVLM agent framework we proposed can only search the
636 Internet autonomously. Our agent still has shortcomings in the use of other multiple tools.
637 We stated in our future outlook that our follow-up work will solve this problem.

638 Guidelines:

- 639 • The answer NA means that the paper has no limitation while the answer No means that
640 the paper has limitations, but those are not discussed in the paper.
- 641 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 642 • The paper should point out any strong assumptions and how robust the results are to
643 violations of these assumptions (e.g., independence assumptions, noiseless settings,
644 model well-specification, asymptotic approximations only holding locally). The authors
645 should reflect on how these assumptions might be violated in practice and what the
646 implications would be.
- 647 • The authors should reflect on the scope of the claims made, e.g., if the approach was
648 only tested on a few datasets or with a few runs. In general, empirical results often
649 depend on implicit assumptions, which should be articulated.
- 650 • The authors should reflect on the factors that influence the performance of the approach.
651 For example, a facial recognition algorithm may perform poorly when image resolution
652 is low or images are taken in low lighting. Or a speech-to-text system might not be
653 used reliably to provide closed captions for online lectures because it fails to handle
654 technical jargon.
- 655 • The authors should discuss the computational efficiency of the proposed algorithms
656 and how they scale with dataset size.
- 657 • If applicable, the authors should discuss possible limitations of their approach to
658 address problems of privacy and fairness.
- 659 • While the authors might fear that complete honesty about limitations might be used by
660 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
661 limitations that aren't acknowledged in the paper. The authors should use their best
662 judgment and recognize that individual actions in favor of transparency play an impor-
663 tant role in developing norms that preserve the integrity of the community. Reviewers
664 will be specifically instructed to not penalize honesty concerning limitations.

665 **3. Theory Assumptions and Proofs**

666 Question: For each theoretical result, does the paper provide the full set of assumptions and
667 a complete (and correct) proof?

668 Answer: [NA]

669 Justification: This work is a solution to the problem of geo-localization in application
670 scenarios. We have provided the source code and will release the related dataset, as the
671 dataset is relatively large (about 32 GB) and cannot be uploaded as an attachment.

672 Guidelines:

- 673 • The answer NA means that the paper does not include theoretical results.
- 674 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
675 referenced.
- 676 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 677 • The proofs can either appear in the main paper or the supplemental material, but if
678 they appear in the supplemental material, the authors are encouraged to provide a short
679 proof sketch to provide intuition.
- 680 • Inversely, any informal proof provided in the core of the paper should be complemented
681 by formal proofs provided in appendix or supplemental material.
- 682 • Theorems and Lemmas that the proof relies upon should be properly referenced.

683 **4. Experimental Result Reproducibility**

684 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
685 perimental results of the paper to the extent that it affects the main claims and/or conclusions
686 of the paper (regardless of whether the code and data are provided or not)?

687 Answer: [Yes]

688 Justification: We provided the source code and will release the related dataset once the paper
689 is accepted, as the dataset is relatively large (about 32 GB) and cannot be uploaded as an
690 attachment.

691 Guidelines:

- 692 • The answer NA means that the paper does not include experiments.
- 693 • If the paper includes experiments, a No answer to this question will not be perceived
694 well by the reviewers: Making the paper reproducible is important, regardless of
695 whether the code and data are provided or not.
- 696 • If the contribution is a dataset and/or model, the authors should describe the steps taken
697 to make their results reproducible or verifiable.
- 698 • Depending on the contribution, reproducibility can be accomplished in various ways.
699 For example, if the contribution is a novel architecture, describing the architecture fully
700 might suffice, or if the contribution is a specific model and empirical evaluation, it may
701 be necessary to either make it possible for others to replicate the model with the same
702 dataset, or provide access to the model. In general, releasing code and data is often
703 one good way to accomplish this, but reproducibility can also be provided via detailed
704 instructions for how to replicate the results, access to a hosted model (e.g., in the case
705 of a large language model), releasing of a model checkpoint, or other means that are
706 appropriate to the research performed.
- 707 • While NeurIPS does not require releasing code, the conference does require all submis-
708 sions to provide some reasonable avenue for reproducibility, which may depend on the
709 nature of the contribution. For example
 - 710 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
711 to reproduce that algorithm.
 - 712 (b) If the contribution is primarily a new model architecture, the paper should describe
713 the architecture clearly and fully.
 - 714 (c) If the contribution is a new model (e.g., a large language model), then there should
715 either be a way to access this model for reproducing the results or a way to reproduce
716 the model (e.g., with an open-source dataset or instructions for how to construct
717 the dataset).

718 (d) We recognize that reproducibility may be tricky in some cases, in which case
719 authors are welcome to describe the particular way they provide for reproducibility.
720 In the case of closed-source models, it may be that access to the model is limited in
721 some way (e.g., to registered users), but it should be possible for other researchers
722 to have some path to reproducing or verifying the results.

723 5. Open access to data and code

724 Question: Does the paper provide open access to the data and code, with sufficient instruc-
725 tions to faithfully reproduce the main experimental results, as described in supplemental
726 material?

727 Answer: [Yes]

728 Justification: We provide the anonymous code link: [https://anonymous.4open.science/
729 r/ViusalGeoLocalization-F8F5/](https://anonymous.4open.science/r/ViusalGeoLocalization-F8F5/). In this link, we also provide a small-scale dataset we
730 collected for people to reproduce the results.

731 Guidelines:

- 732 • The answer NA means that paper does not include experiments requiring code.
- 733 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
734 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 735 • While we encourage the release of code and data, we understand that this might not be
736 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
737 including code, unless this is central to the contribution (e.g., for a new open-source
738 benchmark).
- 739 • The instructions should contain the exact command and environment needed to run to
740 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.
741 cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 742 • The authors should provide instructions on data access and preparation, including how
743 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 744 • The authors should provide scripts to reproduce all experimental results for the new
745 proposed method and baselines. If only a subset of experiments are reproducible, they
746 should state which ones are omitted from the script and why.
- 747 • At submission time, to preserve anonymity, the authors should release anonymized
748 versions (if applicable).
- 749 • Providing as much information as possible in supplemental material (appended to the
750 paper) is recommended, but including URLs to data and code is permitted.

751 6. Experimental Setting/Details

752 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
753 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
754 results?

755 Answer: [Yes]

756 Justification: We explain all the settings in both the main paper (Experiments) and the
757 appendix (Implementation Details).

758 Guidelines:

- 759 • The answer NA means that the paper does not include experiments.
- 760 • The experimental setting should be presented in the core of the paper to a level of detail
761 that is necessary to appreciate the results and make sense of them.
- 762 • The full details can be provided either with the code, in appendix, or as supplemental
763 material.

764 7. Experiment Statistical Significance

765 Question: Does the paper report error bars suitably and correctly defined or other appropriate
766 information about the statistical significance of the experiments?

767 Answer: [Yes]

768 Justification: We deploy a two-sided t-test with $p < 0.05$ for our baseline experiments.

769 Guidelines:

- 770 • The answer NA means that the paper does not include experiments.
- 771 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 772 dence intervals, or statistical significance tests, at least for the experiments that support
- 773 the main claims of the paper.
- 774 • The factors of variability that the error bars are capturing should be clearly stated (for
- 775 example, train/test split, initialization, random drawing of some parameter, or overall
- 776 run with given experimental conditions).
- 777 • The method for calculating the error bars should be explained (closed form formula,
- 778 call to a library function, bootstrap, etc.)
- 779 • The assumptions made should be given (e.g., Normally distributed errors).
- 780 • It should be clear whether the error bar is the standard deviation or the standard error
- 781 of the mean.
- 782 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 783 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 784 of Normality of errors is not verified.
- 785 • For asymmetric distributions, the authors should be careful not to show in tables or
- 786 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 787 error rates).
- 788 • If error bars are reported in tables or plots, The authors should explain in the text how
- 789 they were calculated and reference the corresponding figures or tables in the text.

790 8. Experiments Compute Resources

791 Question: For each experiment, does the paper provide sufficient information on the com-
 792 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 793 the experiments?

794 Answer: [Yes]

795 Justification: We announce the compute resources in the appendix (Implementation Details).

796 Guidelines:

- 797 • The answer NA means that the paper does not include experiments.
- 798 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 799 or cloud provider, including relevant memory and storage.
- 800 • The paper should provide the amount of compute required for each of the individual
- 801 experimental runs as well as estimate the total compute.
- 802 • The paper should disclose whether the full research project required more compute
- 803 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 804 didn't make it into the paper).

805 9. Code Of Ethics

806 Question: Does the research conducted in the paper conform, in every respect, with the
 807 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

808 Answer: [Yes]

809 Justification: The codes used in our paper are all open source, and the data used in the paper
 810 come from copyright-free images on the Internet.

811 Guidelines:

- 812 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 813 • If the authors answer No, they should explain the special circumstances that require a
- 814 deviation from the Code of Ethics.
- 815 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 816 eration due to laws or regulations in their jurisdiction).

817 10. Broader Impacts

818 Question: Does the paper discuss both potential positive societal impacts and negative
 819 societal impacts of the work performed?

820 Answer: [Yes]

821 Justification: We have an outlook on our research in the section Conclusion, which can be
822 widely used in robot positioning and navigation in the future.

823 Guidelines:

- 824 • The answer NA means that there is no societal impact of the work performed.
- 825 • If the authors answer NA or No, they should explain why their work has no societal
826 impact or why the paper does not address societal impact.
- 827 • Examples of negative societal impacts include potential malicious or unintended uses
828 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
829 (e.g., deployment of technologies that could make decisions that unfairly impact specific
830 groups), privacy considerations, and security considerations.
- 831 • The conference expects that many papers will be foundational research and not tied
832 to particular applications, let alone deployments. However, if there is a direct path to
833 any negative applications, the authors should point it out. For example, it is legitimate
834 to point out that an improvement in the quality of generative models could be used to
835 generate deepfakes for disinformation. On the other hand, it is not needed to point out
836 that a generic algorithm for optimizing neural networks could enable people to train
837 models that generate Deepfakes faster.
- 838 • The authors should consider possible harms that could arise when the technology is
839 being used as intended and functioning correctly, harms that could arise when the
840 technology is being used as intended but gives incorrect results, and harms following
841 from (intentional or unintentional) misuse of the technology.
- 842 • If there are negative societal impacts, the authors could also discuss possible mitigation
843 strategies (e.g., gated release of models, providing defenses in addition to attacks,
844 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
845 feedback over time, improving the efficiency and accessibility of ML).

846 11. Safeguards

847 Question: Does the paper describe safeguards that have been put in place for responsible
848 release of data or models that have a high risk for misuse (e.g., pretrained language models,
849 image generators, or scraped datasets)?

850 Answer: [Yes]

851 Justification: The data sets we collect have been manually reviewed twice, and all data
852 containing various types of sensitive information or copyright risks have been filtered out.

853 Guidelines:

- 854 • The answer NA means that the paper poses no such risks.
- 855 • Released models that have a high risk for misuse or dual-use should be released with
856 necessary safeguards to allow for controlled use of the model, for example by requiring
857 that users adhere to usage guidelines or restrictions to access the model or implementing
858 safety filters.
- 859 • Datasets that have been scraped from the Internet could pose safety risks. The authors
860 should describe how they avoided releasing unsafe images.
- 861 • We recognize that providing effective safeguards is challenging, and many papers do
862 not require this, but we encourage authors to take this into account and make a best
863 faith effort.

864 12. Licenses for existing assets

865 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
866 the paper, properly credited and are the license and terms of use explicitly mentioned and
867 properly respected?

868 Answer: [Yes]

869 Justification: We list and acknowledge all other open-source codes we used in the file
870 'README.md' and we follow the license for existing assets.

871 Guidelines:

- 872 • The answer NA means that the paper does not use existing assets.
- 873 • The authors should cite the original paper that produced the code package or dataset.

- 874 • The authors should state which version of the asset is used and, if possible, include a
875 URL.
876 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
877 • For scraped data from a particular source (e.g., website), the copyright and terms of
878 service of that source should be provided.
879 • If assets are released, the license, copyright information, and terms of use in the
880 package should be provided. For popular datasets, `paperswithcode.com/datasets`
881 has curated licenses for some datasets. Their licensing guide can help determine the
882 license of a dataset.
883 • For existing datasets that are re-packaged, both the original license and the license of
884 the derived asset (if it has changed) should be provided.
885 • If this information is not available online, the authors are encouraged to reach out to
886 the asset's creators.

887 13. **New Assets**

888 Question: Are new assets introduced in the paper well documented and is the documentation
889 provided alongside the assets?

890 Answer: [Yes]

891 Justification: In this paper, we provide the algorithm of the code and introduce the dataset in
892 detail (in the appendix).

893 Guidelines:

- 894 • The answer NA means that the paper does not release new assets.
895 • Researchers should communicate the details of the dataset/code/model as part of their
896 submissions via structured templates. This includes details about training, license,
897 limitations, etc.
898 • The paper should discuss whether and how consent was obtained from people whose
899 asset is used.
900 • At submission time, remember to anonymize your assets (if applicable). You can either
901 create an anonymized URL or include an anonymized zip file.

902 14. **Crowdsourcing and Research with Human Subjects**

903 Question: For crowdsourcing experiments and research with human subjects, does the paper
904 include the full text of instructions given to participants and screenshots, if applicable, as
905 well as details about compensation (if any)?

906 Answer: [NA]

907 Justification: This paper aims to address visual geo-localization tasks and does not contain
908 any experiments with human subjects.

909 Guidelines:

- 910 • The answer NA means that the paper does not involve crowdsourcing nor research with
911 human subjects.
912 • Including this information in the supplemental material is fine, but if the main contribu-
913 tion of the paper involves human subjects, then as much detail as possible should be
914 included in the main paper.
915 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
916 or other labor should be paid at least the minimum wage in the country of the data
917 collector.

918 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 919 Subjects**

920 Question: Does the paper describe potential risks incurred by study participants, whether
921 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
922 approvals (or an equivalent approval/review based on the requirements of your country or
923 institution) were obtained?

924 Answer: [NA]

925 Justification: This paper does not contain any experiments with human subjects.

926
927
928
929
930
931
932
933
934
935
936

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.