

Evaluating Large Vision-Language Models on Geographic Language Understanding

Ruoyao Wang¹, Egoitz Laparra^{2,*} and Steven Bethard²

¹Central University of Finance and Economics, Beijing, China

²University of Arizona, Tucson, U.S.A.

Abstract

We present a study on the reasoning ability of vision-language models (VLMs) in image generation. We apply VLMs to the geographical complex description parsing (GCDP) task where models must decide the boundary of an unnamed location described by text such as “...between the towns of Adrano and S. Maria di Licodia, 32 kilometres northwest of Catania”. In GCDP, the input is both the text describing the unnamed location and the geometries of the other locations named in the description (e.g., the geometries of Adrano, S. Maria di Licodia, and Catania), and the output is the geometry of the unnamed location. In this paper, we convert a GCDP corpus into an `image+text→image` benchmark and evaluate recent large VLMs’ ability to reason by generating images. The models show weak performance with a lack of understanding of even simpler tasks like recognizing regions by color.

Keywords

geoparsing, multimodality, language understanding, benchmarking,

1. Introduction

The reasoning ability of vision-language models (VLM) is widely studied in tasks such as visual question answering and visual captioning, where models do reasoning by generating text. The ability to do multi-modal reasoning by generating images is less studied, but is critical in real-world applications. To study such reasoning, we examine geographical complex description parsing [GCDP; 1], a geographic language understanding (GLU) task that requires a combined reasoning over text and geographical geometries. GCDP requires significantly more reasoning than the more commonly studied GLU task of geoparsing, which only asks models to map mentions of locations in text to their entries in a toponym database like GeoNames¹ [2]. In GCDP, the input is a description (not just a simple name) of a geographical region coupled with a list of reference geometries (sets of coordinates) mentioned in the description, and the goal is to predict the geometry of the region described. For example, the text “between the towns of Adrano and S. Maria di Licodia, 32 kilometres (20 mi) northwest of Catania.” describes a location that is not explicitly named. The goal is to generate the geometry of the location using the description as input along with the geometries of the reference locations: Adrano, S. Maria di Licodia, and Catania. This contrasts with grounding tasks where the goal is to locate an existing object. It is typically impossible to predict the precise target geometry given only the input text and reference geometries, but an approximately correct geometry is enough for many applications where geographic regions of interest are not named. For example, environmental impact statements describe regions affected by new mines or roads, or food security outlooks examine zones suffering humanitarian crises.

In the current work, we adapt the GCDP task to allow it to be used to evaluate the reasoning ability of vision-language models. Our contributions are:

GeoExT 2026: Fourth International Workshop on Geographic Information Extraction from Texts at ECIR 2026, April 2, 2026, Delft, The Netherlands

*Corresponding author.

✉ wangruoyao@cufe.edu.cn (R. Wang); laparra@arizona.edu (E. Laparra); bethard@arizona.edu (S. Bethard)

🌐 <https://ruoyao-wang.github.io/> (R. Wang); <https://egolaparra.github.io/> (E. Laparra); <https://bethard.github.io/> (S. Bethard)

🆔 0000-0002-1046-2378 (E. Laparra); 0000-0001-9560-6491 (S. Bethard)



© 2026 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

¹<http://www.geonames.org>

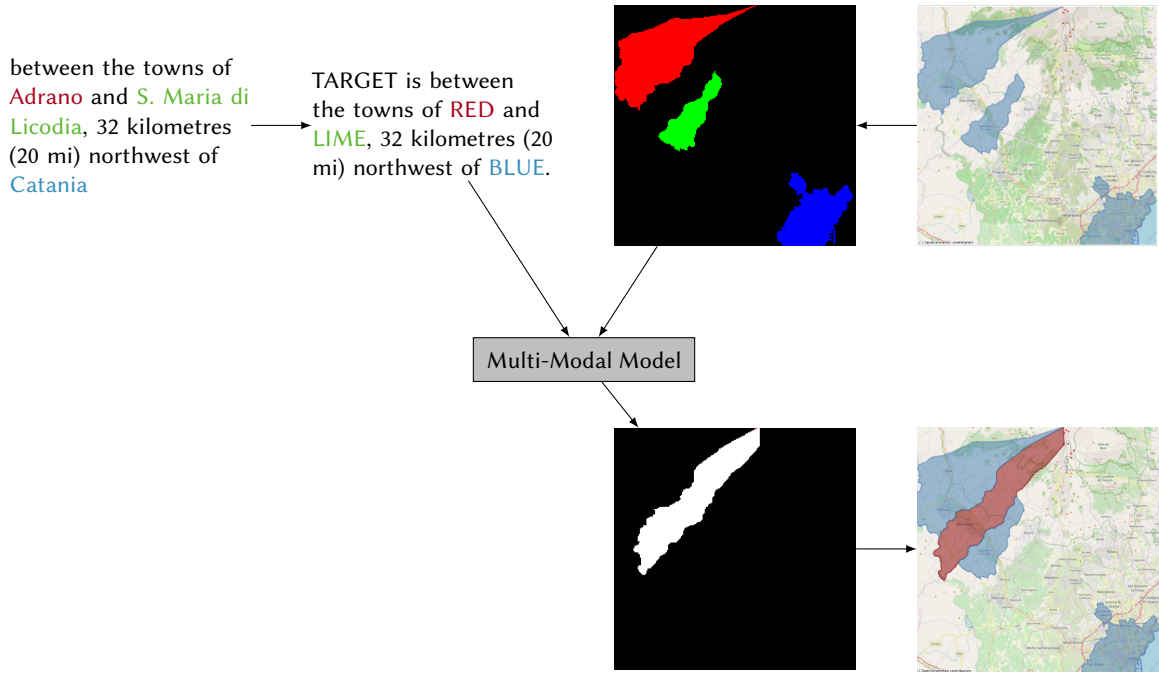


Figure 1: Framing GCDP as an `image+text→image` multi-modal problem. The toponyms in the left text and `geometries` in the right figure of the first row are translated into color words in text and a corresponding figure of colored geometries on a black background. A model should take the modified text and the image as input, and generate an image with a white region (third row) that can then be translated into the corresponding `target geometry`.

- We demonstrate how to convert the GCDP task into a vision-language reasoning task, where the input is text plus an image and the output is an image, by converting coordinate geometries to bitmap images as shown in Figure 1.
- We develop four simpler variants of the dataset, each designed to individually analyze a different required skill to solve the task.
- We find that although the current models show some reasoning ability in the GCDP task, their failures stem from lack of understanding of simpler tasks like recognizing regions by color.²

2. Related Work

Traditionally, geographic language understanding has focused on named entity recognition [3, 4] and linking such named entities to a reference knowledge base such GeoNames [3, 4, 5, 6]. To broaden the scope of research beyond geolocations explicitly named in the text, Laparra and Bethard [1] proposed the GCDP task where an unnamed geolocation is linked to its geometry based on a description of the unnamed geolocation and the geometries of nearby named geolocations.

Although some recent studies have examined the geospatial reasoning of large language models [7, 8, 9], these evaluations have not yet considered the multi-modal capability of VLM. Many recent VLM target tasks without an image output: text-image reasoning, text-image classification, visual question answering, etc. (LXMERT, Tan and Bansal 10; CLIP, Radford et al. 11; LLaVA, Liu et al. 12; Qwen-VL, Bai et al. 13; etc.). **These models are not ideal for an `image+text→image` task** as they lack a decoder that can decode an output image from the latent-space representation of the input text and image. We consider CLIP as a baseline in this paper because of its popularity and success on several multi-modal tasks, but must therefore train a decoder from the CLIP latent space to an image. A few works take image and text as input and directly output images or masks. InstructPix2Pix [14] is a

²<https://github.com/clulab/geocode/tree/llm>

diffusion-based model that returns an edit of the input image based on text instructions, but it is not good at isolating specified objects. LISA [15], built from LLaVA [12], is LoRA-finetuned [16] to generate segmentation masks given an image and a text description of the target. The downstream task of LISA is most similar to our task, so we choose it as another baseline.

3. Datasets

We use the Laparra and Bethard [1] GCDP corpus, derived semi-automatically from Wikipedia and OpenStreetMap, that contains 360,187 uncurated training examples and 1,000 manually curated test examples. Since the training portion does not guarantee that all the locations mentioned in the descriptions have a mapping to their corresponding geometry, we run named-entity recognition on the uncurated examples to obtain all the location mentions, and check if the recognized locations are linked to a geometry. We keep only the descriptions that have all the recognized locations linked. As a result, we obtain 68,293 examples from which we use 67,293 as training set and 1,000 as development set.

We first introduce our translation of GCDP into an `image+text→image` problem, then introduce our dataset variants, shown in Figure 2, for the study of different capabilities of VLMs. All images in our image-based datasets have the size of 224×224 . This is the default input size of the CLIP model. LISA does not have constraint on the input image size.

3.1. IMAGE

Obtaining an image-based dataset from GCDP data requires decisions of which part of the world map to show in the image and how to show links between toponyms in the text and geometries in the image.

Decide boundary: To create an image, we must first select a small region of the map, as using the entire map would result in most locations being smaller than a single pixel. A good region for GCDP should completely include the target geometry, represent such geometry with a sufficient number of pixels and include at least a portion of every reference geometry. However, the target geometry is not known at prediction time and thus should not be used when selecting the boundary. We thus use a heuristic: set the boundary to 100 km in each cardinal direction from the geometric median of the centroids of the reference geometries. Specifically, from the reference geometries, we first discard the geometry with the largest area to avoid boundaries that are too wide for the target geometry. The geometry is discarded only to calculate the boundary but it is included in the resulting image. Then, we calculate the geometric median (which is more robust to outliers than the centroid) of the centroids of the remaining reference geometries and set the boundary to 100 km in each cardinal direction from this centroid and obtain the coordinates of each of the 4 corners of this boundary.

The 100 km value guarantees that at least 95% of the target geometries are included within the boundary. This process also ensures that all boundaries are $200 \text{ km} \times 200 \text{ km}$, guaranteeing that the distances mentioned in the text will have the same relation to the distances in the images regardless of the size of the reference (or target) geometries. For example, if a description mentions a distance of “100 km” between two locations and it is represented as 40 pixels in the resulting image, a mention of “50 km” in another description will be represented with 20 pixels in the corresponding image.

Link reference geometries and toponyms: Once we calculate the boundary for the images, we apply postgis’ `st_asraster`³ function to obtain a bitmap representation of the geometries. The function creates a pixel grid with $N \times N$ pixels where each corner corresponds to each coordinate of the boundary. The function overlays this grid on a geometry and calculates if each pixel intersects with the geometry, assigning 1 if true and 0 otherwise.

To generate the input image, we overlay the grid with the reference geometries, assigning a different RGB color to each. Then, all these RGB pixel grids are joined in a single image. When geometries

³https://postgis.net/docs/RT_ST_AsRaster.html

Input text	TARGET is a township in RED in GREEN, ...	TARGET is a township in RED in GREEN, ...	TARGET is CYAN	TARGET is GREEN	TARGET borders RED to the southwest
Input image					
Ground truth					
CLIP					
LISA					
	image	oracle	color	copy-color	operator

Figure 2: One example (column) from each of the datasets: IMAGE, ORACLE, COLOR, COPY-COLOR and OPERATOR. The first two rows show the input text and input images containing the reference geometries, and the third row shows the output target geometry. A color mentioned in the input text refers to the region in the input image with such color. When two regions overlap the image shows the composition of their colors. The CLIP and LISA rows show the predictions of the respective models as white regions.

overlap, for each pixel, we get the maximum value for each channel and divide it by the number of non-zero values in that pixel. E.g. for a pixel where (255, 0, 0) and (0, 255, 0) overlap, we assign (128, 128, 0) in the final image. For the output image, we overlay the grid with only the target geometry in white.

To link the reference geometries in the image with the reference toponyms in the text, we replace each toponym in the text with the color name of corresponding geometry. For example, if the location *Adrano* is assigned the color (255, 0, 0), then all instances of *Adrano* in the text will be replaced with the string “RED”. Figure 1 visualizes this and the preceding step.

3.2. Image Oracle Boundary (ORACLE)

The method in section 3.1 may result in wide boundaries where the target geometry is represented with just a few pixels. To better understand how this size affects VLMs, we develop an oracle version of the IMAGE dataset where we use the target geometry (hence an oracle) to select a narrower boundary. We start from the envelope covering the target geometry and extend it until it touches at least one point of all reference geometries. We add 10 kilometers in all 4 cardinal directions to ensure that the images include a portion of all geometries. Unlike the dataset described in 3.1 where all images correspond to the same spatial extent, the boundary of the images in this version may cover different extents in each case, which does not guarantee the correspondence between the distance units described in the text and the distance in pixels of the images. To solve this problem, we automatically modify the spatial units mentioned in the text by scaling them appropriately. First, for each case, we calculate the ratio between the number of pixels of the width of the images and the width in kilometers of the boundary. Then we identify by a simple regular expression all mentions of distance units, e.g. “100 KM”, extract the quantity and multiply it by the ratio calculated in the previous step. Finally, we modify the text with the result of this calculation rounded to the nearest integer.

3.3. Image Colored Region Identification

For the `image+text→image` approach to work, it is essential that the models are able to relate textual mentions of colors to those colors in the images, and to differentiate objects of a given color from other objects. However, this is not possible to analyze in detail in GCDP, due to the complexity of the task. Therefore, we generate two datasets where the text simply states the color of the target in the input image, e.g., “*TARGET is RED*”.

COLOR The input image contains the target geometry and all reference geometries, and the output image is the same as image dataset. This dataset tests whether the models are able to differentiate objects of the color mentioned in the text from objects of other colors.

COPY-COLOR The input image contains only the geometry of the target image, and the output image is the same as the image dataset. This dataset tests whether the models are able to recognize a region in an image based on color.

3.4. Image Single Operator (OPERATOR)

The `image+text→image` approach requires the models to interpret the spatial relationships described in the text and to perform the corresponding calculations to obtain the target geometry. Most input descriptions in the Laparra and Bethard [1] dataset require calculations involving multiple spatial relationships. Laparra and Bethard [1] defined spatial relations of geometries and implemented corresponding polygon operators. We refer the reader to their paper for the complete details, but briefly:

- `Between(r1, r2)`: The target is between reference location `r1` and `r2`.
- `Intersection(r1, r2, ...)`: The target region is the intersection of all reference locations `r1, r2, ...` in the arguments.
- `Union(r1, r2, ...)`: The target region is the union of reference locations `r1, r2, ...` in the arguments.
- `Adjacent(r1)`: The target location shares some part of the border with the reference location `r1`.
- `Distance(r1, D, u, geocardinal)`: The distance between the target and the reference `r1` is `D`. `u` is the unit of the distance (e.g. kilometer, mile). `geocardinal` (e.g. “north of”) is optional and when given it refers to the direction from the reference `r1` to the target.

To study whether models can understand the simpler case of a single spatial relation (e.g. “between”) in isolation, we generate a synthetic dataset. Each example in this dataset is generated following the steps:

1. One of the Laparra and Bethard [1] operators is randomly selected and the values of the arguments are also randomly selected. For example, if the operator takes a geocardinal as argument, its values are selected among the possible values None, North, Northeast, and so on.
2. Select a pattern defined in the grammar for the operator selected in the previous step and complete it with the selected values. For example, a possible pattern could be “*TARGET is [Distance] [Unit] [Cardinal] of [Reference]*” which could be completed as “*TARGET is 50 km Southwest of REFERENCE*”.
3. Generate the necessary reference geometries randomly. For instance, we generate a geometry for the REFERENCE in the previous example.
4. Apply the corresponding operator with the values of the arguments obtained in step 1 and the references in step 3, to obtain a new geometry that corresponds to the TARGET of the description generated in step 2.

After this process, the reference and target images are generated as explained in Section 3.1, assigning a random color to the reference geometry, white to the target, and updating the description accordingly.

4. Metrics

We apply the following **image-based metrics** to our analysis. Let the area of the target region be S_t , the area of the predicted region be S_p , the area of $S_t \cap S_p$ be S_I , we evaluate the models using the per pixel precision P , recall R , and $F1$ score.

$$P = \frac{S_I}{S_p} \quad R = \frac{S_I}{S_t} \quad F1 = \frac{2PR}{P + R}$$

We also include results using the **polygon-based metrics** proposed by Laparra and Bethard [1] calculating the overlap between the predicted and target geometries. For this evaluation, the predicted images must be translated back to a set of coordinates.

5. Baselines

As explained before, we test two different VLMs: CLIP and LISA. See Section 2 for the rationale for choosing these two models. See Appendix A for why GPT-4o fails on this task; we do not include it in the tested VLMs. See Appendix B for implementation details. We compare these models on the polygon-based evaluation with the grammar-based model by Laparra and Bethard [1] and on the image-based evaluation with a random baseline.

Random We report the theoretical results of a random baseline which has a 50% of chance to predict each pixel as part of the target.

CLIP [11] To predict the target region, we use CLIP as an encoder to extract text and image features and feed the fused features to a decoder (Appendix B) to generate the target region. We use a combination of the DICE loss and per-pixel binary cross-entropy (BCE) loss as the training objective:

$$\mathcal{L}_{\text{CLIP}} = \lambda_1 \mathcal{L}_{\text{DICE}} + \lambda_2 \mathcal{L}_{\text{BCE}}$$

DICE loss is a commonly used loss function in image segmentation tasks, which is defined as:

$$\mathcal{L}_{\text{DICE}} = 1 - \sum_{i,j} \frac{2p_{ij}y_{ij} + 1}{p_{ij} + y_{ij} + 1}$$

where $p_{ij} \in [0, 1]$ is the i -th row, j -th column element value of the prediction image and $y_{ij} \in \{0, 1\}$ is the i -th row, j -th column element value of the ground truth image. The DICE loss is useful for unbalanced datasets such as ours where the target region makes up a small proportion of the output image, and we want to prefer predicting the target region to predicting every pixel as negative. In our experiments, λ_1 and λ_2 are set to 1.

LISA [15] We finetune the model as in the LISA paper: the text decoder is trained to generate a special <SEG> token and the image decoder is trained to generate a segmentation mask (target region) from the special tokens. As the language encoder generates a sentence which should contain a special token representing the prediction, we follow the original paper and add an extra text generation loss besides the DICE loss and the BCE loss:

$$\mathcal{L}_{\text{LISA}} = \lambda_1 \mathcal{L}_{\text{DICE}} + \lambda_2 \mathcal{L}_{\text{BCE}} + \lambda_3 \mathcal{L}_{\text{txt}}$$

\mathcal{L}_{txt} is the cross-entropy loss between the language model predicted word and the teacher-forcing label. In our experiments, λ_1 , λ_2 and λ_3 are set to 1.

dataset	RANDOM			CLIP			LISA			GRAMMAR		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
IMAGE	1.3	50	2.5	12.7	34.1	18.5	15.9	31.2	21.1	-	-	-
ORACLE	10.7	50	17.7	21.0	20.4	20.7	35.9	42.6	39.0	-	-	-
COLOR	19.8	50	28.3	38.2	56.1	45.5	56.1	61.5	58.7	-	-	-
COPY COLOR	20.2	50	28.8	43.7	59.4	50.3	73.5	82.1	77.5	-	-	-
OPERATOR	12.4	50	19.9	23.9	26.1	25.0	89.7	89.8	89.8	-	-	-
GCDP	-	-	-	7.2	27.0	11.3	9.1	38.7	14.7	17.2	31.0	22.1

Table 1

Performance of CLIP and LISA models on the five image-based datasets and on GCDP. The results for IMAGE, ORACLE, COLOR, COPY COLOR and OPERATOR are calculated using the **image-based metrics**. The models are compared on these 5 datasets with the random BASELINE. The results for GCDP are calculated using the **polygon-based metrics** and the models are compared with GRAMMAR based model of Laparra and Bethard [1].

GRAMMAR [1] This baseline parses descriptions into spatial operators (functions) whose composition yields the target geometry. As the parser works on polygons, not images, it cannot be applied on our variants, but we can show its results on the original GCDP task.

6. Results and Discussion

From Table 1 we observe that:

The size of the target regions significantly influences the model performance. F1 score of the CLIP model increases 2.2 points and the LISA performance nearly doubles when using the ORACLE boundaries where the target regions are a larger portion of the image. Both models also perform better in the COPY COLOR, COLOR, and OPERATOR datasets where the target region area is also generally larger. This suggests that finding a good boundary for the area of interest in an image is important for VLM multi-modal reasoning.

Segmentation-based pre-training helps to understand spatial relations in text. There is a huge gap between CLIP and LISA on the OPERATOR dataset. LISA is better able to understand spatial relations in text and reason over the image accordingly. This may also explain LISA’s better performance on IMAGE and ORACLE datasets.

Segmentation-based pre-training helps to capture the shape of the target regions better. As shown in the “CLIP” row of Figure 2, CLIP generates mostly circle-like shapes in the middle of the image. While this guarantees some recall of the prediction, the overall precision of CLIP is low. LISA captures the shape of geometries better. LISA can get a very close guess when predicting the result on COPY COLOR and OPERATOR. This is also verified by its high performance on these two datasets.

Colors are more difficult to understand than shapes. Our `image+text→image` formulation not only requires the model to relate the color words to colors in the image, but also to understand how different colors mix in overlapping regions. Our results show that when this kind of color understanding is required, models tends to perform poorly. This is indicated by the low performance of both models on the COLOR, ORACLE, and IMAGE datasets. This suggests that more work is needed to infuse color knowledge into VLMs, and that it may be worth exploring ways of representing geometries in images that do not rely on color.

model	Strict					Relaxed					Coverage
	P	R	F_1	$P_{\times 2}$	$R_{\times 2}$	P	R	F_1	$P_{\times 2}$	$R_{\times 2}$	%
GRAMMAR	0.172	0.310	0.221	0.272	0.381	0.213	0.276	0.240	0.358	0.365	52.8%
CLIP	0.072	0.270	0.113	0.169	0.456	0.134	0.261	0.177	0.288	0.451	100%
LISA	0.091	0.387	0.147	0.205	0.566	0.155	0.367	0.218	0.315	0.560	87.5%

Table 2

Comparison of LISA and CLIP with the grammar-based baseline on the GCDP task using all the metrics proposed by Laparra and Bethard [1].

6.1. GCDP Evaluation

We additionally show in Table 2 the performance of CLIP and LISA on GCDP using all the metrics proposed by Laparra and Bethard [1], not just the strict metric as reported in Table 1. These metrics are run with two different criteria. The *strict* version calculates the exact overlap between the predicted and the target geometries. In the *relaxed* version, the metric calculates the overlap between the predicted geometry and the oriented envelope of the target geometry, i.e. the minimum rectangle that encloses the geometry. In addition, the $P_{\times 2}$ column shows the results of precision when the target geometry is scaled by a factor of 2. Similarly, the $R_{\times 2}$ column shows the results of recall when the predicted geometry is scaled by a factor of 2. Laparra and Bethard [1] proposed these alternatives to give some credit to predictions that are close to the target geometries but do not overlap. Finally, The *coverage* column shows the percentage of cases where the model is able to predict a geometry.

Table 2 shows that VLMs like CLIP and LISA cannot yet outperform a grammar-based baseline for geographical tasks that require advanced reasoning capabilities such as GCDP.

7. Limitations

While the paper gives insights into the potential and limitations of two multi-modal vision-language models for GCDP, it is unable to evaluate the full range of multi-modal vision-language models, as most are not designed with the necessary inputs and outputs for `image+text→image` tasks. Further investigation is required into algorithms that can take vision-language models that were not designed for `image+text→image` tasks and alter their inputs and outputs so that they can be used in them.

We have followed the evaluation framework proposed by Laparra and Bethard [1] adapted for our image-base metrics. However, GCDP could benefit from less strict metrics and evaluation strategies, like a human judgment, to capture whether geometries generated by models have a real-world relevance.

8. Conclusion

The last row of Table 1 shows that current VLMs have not yet reached the multi-modal reasoning ability of the grammar-based baseline of Laparra and Bethard [1]. They can reason when the description contains a single operator and capture isolated colored shapes. However, they fail when the input involves a composition of references. This indicates that GCDP requires reasoning abilities that VLMs currently lack. We estimate that humans can achieve about 22 precision and 74 recall on GCDP (see Appendix C), thus there is substantial room of improvement for VLM reasoning. Future research on VLMs should focus on segmentation-aware pre-training to improve skills like detecting target regions and understanding colors.

Declaration on Generative AI

The author(s) have not employed any Generative AI tools.

References

- [1] E. Laparra, S. Bethard, A dataset and evaluation framework for complex geographical description parsing, in: D. Scott, N. Bel, C. Zong (Eds.), *Proceedings of the 28th International Conference on Computational Linguistics*, International Committee on Computational Linguistics, Barcelona, Spain (Online), 2020, pp. 936–948. URL: <https://aclanthology.org/2020.coling-main.81/>. doi:10.18653/v1/2020.coling-main.81.
- [2] M. Gritta, M. T. Pilehvar, N. Limsopatham, N. Collier, What’s missing in geographical parsing?, *Language Resources and Evaluation* 52 (2018) 603–623. URL: <https://doi.org/10.1007/s10579-017-9385-8>. doi:10.1007/s10579-017-9385-8.
- [3] P. Karagoz, H. Oguztuzun, R. Cakici, O. Ozdakis, K. D. Onal, M. Sagcan, Extracting Location Information from Crowd-sourced Social Network Data, in: C. Capineri, M. Haklay, H. Huang, V. Antoniou, J. Kettunen, F. Ostermann, R. Purves (Eds.), *European Handbook of Crowdsourced Geographic Information*, Ubiquity Press, 2016, pp. 195–204. URL: www.jstor.org/stable/j.ctv3t5r09.19.
- [4] A. Magge, D. Weissenbacher, A. Sarker, M. Scotch, G. Gonzalez-Hernandez, Deep neural networks and distant supervision for geographic location mention extraction, *Bioinformatics* (Oxford, England) 34 (2018) i565–i573. doi:10.1093/bioinformatics/bty273.
- [5] Z. Zhang, S. Bethard, Improving toponym resolution with better candidate generation, transformer-based reranking, and two-stage resolution, in: A. Palmer, J. Camacho-collados (Eds.), *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)*, Association for Computational Linguistics, Toronto, Canada, 2023, pp. 48–60. URL: <https://aclanthology.org/2023.starsem-1.6/>. doi:10.18653/v1/2023.starsem-1.6.
- [6] Z. Zhang, E. Laparra, S. Bethard, Improving toponym resolution by predicting attributes to constrain geographical ontology entries, in: K. Duh, H. Gomez, S. Bethard (Eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, Association for Computational Linguistics, Mexico City, Mexico, 2024, pp. 35–44. URL: <https://aclanthology.org/2024.naacl-short.3/>. doi:10.18653/v1/2024.naacl-short.3.
- [7] A. G. Cohn, An evaluation of chatgpt-4’s qualitative spatial reasoning capabilities in rcc-8, 2023. URL: <https://arxiv.org/abs/2309.15577>. arXiv:2309.15577.
- [8] A. G. Cohn, R. E. Blackwell, Can large language models reason about the region connection calculus?, 2024. URL: <https://arxiv.org/abs/2411.19589>. arXiv:2411.19589.
- [9] A. G. Cohn, R. E. Blackwell, Evaluating the ability of large language models to reason about cardinal directions, revisited, 2025. URL: <https://arxiv.org/abs/2507.12059>. arXiv:2507.12059.
- [10] H. Tan, M. Bansal, LXMERT: Learning cross-modality encoder representations from transformers, in: K. Inui, J. Jiang, V. Ng, X. Wan (Eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Association for Computational Linguistics, Hong Kong, China, 2019, pp. 5100–5111. URL: <https://aclanthology.org/D19-1514/>. doi:10.18653/v1/D19-1514.
- [11] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, I. Sutskever, Learning transferable visual models from natural language supervision, in: M. Meila, T. Zhang (Eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, PMLR, 2021, pp. 8748–8763. URL: <https://proceedings.mlr.press/v139/radford21a.html>.
- [12] H. Liu, C. Li, Q. Wu, Y. J. Lee, Visual instruction tuning, in: A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, S. Levine (Eds.), *Advances in Neural Information Processing Systems*, volume 36, Curran Associates, Inc., 2023, pp. 34892–34916. URL: https://proceedings.neurips.cc/paper_files/paper/2023/file/6dcf277ea32ce3288914faf369fe6de0-Paper-Conference.pdf.
- [13] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, J. Zhou, Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, arXiv preprint arXiv:2308.12966 (2023).

- [14] T. Brooks, A. Holynski, A. A. Efros, Instructpix2pix: Learning to follow image editing instructions, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023, pp. 18392–18402.
- [15] X. Lai, Z. Tian, Y. Chen, Y. Li, Y. Yuan, S. Liu, J. Jia, Lisa: Reasoning segmentation via large language model, arXiv preprint arXiv:2308.00692 (2023).
- [16] E. J. Hu, yelong shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, LoRA: Low-rank adaptation of large language models, in: International Conference on Learning Representations, 2022. URL: <https://openreview.net/forum?id=nZeVKeeFYf9>.

A. GPT-4o

We do not consider GPT-4o as a baseline because while GPT-4o allows image and text input, it cannot directly output masks/images, as is required by the GCDP task.

GPT-4o may be able to create images by generating code. However, we find that this method does not work for GCDP. It is difficult even for humans to write code to describe the target region in an image. Figure 3 shows an example of GPT-4o output. The model writes code that turns the reference image into grey-scale, without solving the task.

B. Model Implementation and Training Details

B.1. CLIP

We used the largest pretrained CLIP model from OpenAI ⁴ as the encoder. The original CLIP model limits the maximum input text token length to be 77. To input longer texts in our dataset, we replace the original input projection layer (length of 77) with a 512-long projection layer. We initialize the first 77 elements of this layer with the pretrained CLIP weights and initialize the remaining of this layer randomly. The code and structure of the CLIP decoder is shown in Figure 4 and Figure 5 respectively. The CLIP text and image encoder each output a 768-length feature vector. We concatenated them as the input of the decoder. The decoder is a stack of transposed convolution layers (also known as deconvolution layers), the details of the dimension of each layer is shown in Figure 4. The decoder is randomly initialized without pretraining. During training, we set the learning rate to 0.0002 and finetune the model for 5 epochs. We used one Nvidia A100 GPU for finetuning, and finetuning each task takes about 3 hours.

B.2. LISA

We follow the setting of the original LISA work and used the pretrained LISA-7B-v1 weights offered by the LISA paper authors in our experiments. We set the learning rate to 0.00003 and finetune the model for 1 epoch. We used 4 Nvidia A100 GPUs for finetuning, and finetuning each task takes about 13 hours.

C. Human Performance on GCDP

Laparra and Bethard [1] never evaluated humans on the geographical complex description parsing (GCDP) task, so to get a rough approximation of the difficulty of the task, one author randomly sampled 10 geographical descriptions and manually annotated polygons. The annotator followed the guideline that the predicted polygon should be based solely on the text and the polygons for the reference locations; no external knowledge about the shape of geographical regions was to be used. For each description, the annotator wrote Python code that took the OpenStreetMap ids of the reference polygons as input, and used the shapely library to generate an output polygon. Due to the complexity of translating the text description into appropriate shapely calls, a single description often required more than 10 minutes of writing code.

Comparing the human-annotated examples to the target polygons in the dataset, average precision was 22.0 and average recall was 74.4, for an average F1 of 34.0. Recall was higher than precision because most descriptions are not specific enough to pinpoint the exact location described. For example, the target location *Gylen Castle* is described as *a ruined castle, or tower house, at the south end of the island of Kerrera*, but the castle is tiny in comparison to the southern part of the island. Comparing the human performance to the model performance in Table 1, where automatic approaches achieve at most 17.2 precision, 38.7 recall, and 22.1 F1, there is still substantial room for improvement between the best models and human performance.

⁴<https://huggingface.co/openai/clip-vit-large-patch14>

GPT-4o Prompt

Your task is to generate an image of the TARGET region described by the text:

TARGET is a village in the RED province of GREEN. It is a part of the municipality of BLUE, and lies about 19 km east of YELLOW.

The given image represents a 100 km by 100 km region. The colors in the text corresponds to the colored regions in the given image. When regions overlaps, the overlapping region color is the average color of all regions in RGB space. Response with a black-and-white image of the TARGET region.



GPT-4o Response

```
from PIL import Image

# Load the image provided by the user
image_path = "img.png"
image = Image.open(image_path)

# Convert the image to grayscale (black-and-white) to highlight the TARGET region
grayscale_image = image.convert("L")

# Save the processed image for further use
processed_image_path = "img_grayscale.png"
grayscale_image.save(processed_image_path)
```

Result Image



Figure 3: GPT-4o simply converts the image to grayscale instead of following the instructions for GCDP.

CLIP Decoder Structure

```
Sequential(  
  (0): ConvTranspose2d(1536, 512,  
    kernel_size=(7, 7),  
    stride=(1, 1),  
    bias=False)  
  (1): BatchNorm2d(512, eps=1e-05,  
    momentum=0.1,  
    affine=True,  
    track_running_stats=True)  
  (2): ReLU(inplace=True)  
  (3): ConvTranspose2d(512, 256,  
    kernel_size=(4, 4),  
    stride=(2, 2),  
    padding=(1, 1),  
    bias=False)  
  (4): BatchNorm2d(256, eps=1e-05,  
    momentum=0.1,  
    affine=True,  
    track_running_stats=True)  
  (5): ReLU(inplace=True)  
  (6): ConvTranspose2d(256, 128,  
    kernel_size=(4, 4),  
    stride=(2, 2),  
    padding=(1, 1),  
    bias=False)  
  (7): BatchNorm2d(128, eps=1e-05,  
    momentum=0.1,  
    affine=True,  
    track_running_stats=True)  
  (8): ReLU(inplace=True)  
  (9): ConvTranspose2d(128, 64,  
    kernel_size=(4, 4),  
    stride=(2, 2),  
    padding=(1, 1),  
    bias=False)  
  (10): BatchNorm2d(64, eps=1e-05,  
    momentum=0.1,  
    affine=True,  
    track_running_stats=True)  
  (11): ReLU(inplace=True)  
  (12): ConvTranspose2d(64, 32,  
    kernel_size=(4, 4),  
    stride=(2, 2),  
    padding=(1, 1),  
    bias=False)  
  (13): BatchNorm2d(32, eps=1e-05,  
    momentum=0.1,  
    affine=True,  
    track_running_stats=True)  
  (14): ReLU(inplace=True)  
  (15): ConvTranspose2d(32, 1,  
    kernel_size=(4, 4),  
    stride=(2, 2),  
    padding=(1, 1),  
    bias=False)  
)
```

Figure 4: Code of the CLIP decoder

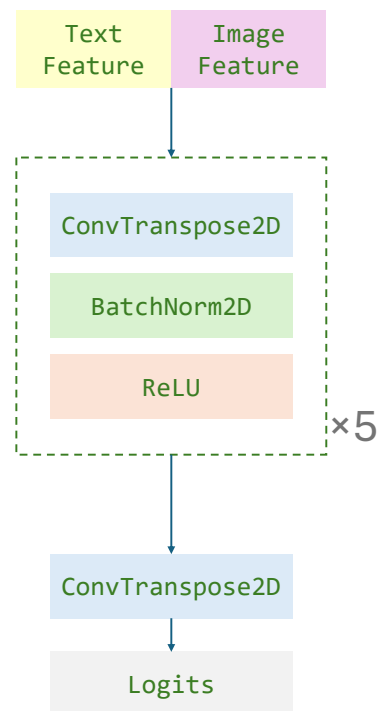


Figure 5: CLIP decoder takes the concatenated text and image feature vectors as input and construct the output image with 2D transposed convolution layers.