

RegionGPT: Towards Region Understanding Vision Language Model

What is the problem?

Vision language models (VLMs) struggle with detailed regional visual understanding due to limited spatial awareness of the vision encoder and the use of coarse-grained training data, which lacks detailed, region-specific captions. This impedes the models' ability to handle tasks that require fine-grained image region analysis(Guo_RegionGPT_Towards_R...).

What has been done earlier?

Prior works have explored inputting regions of interest in textual form (e.g., bounding boxes) and some methods such as GPT4RoI have introduced spatial boxes with RoI-aligned features. However, these approaches are limited by restricted positional formats (e.g., fixed boxes) and do not fully exploit the potential for fine-grained region-specific visual representation(Guo_RegionGPT_Towards_R...)..

Proposed Solutions and Challenges

What are the remaining challenges?

The main challenges are the need for better spatial awareness of regional representations in vision-language models and the lack of detailed region-level annotations. Existing datasets provide simplistic descriptions of regions, and current models fail to fully capture fine-grained details like color, shape, and spatial relationships.

What novel solution proposed by the authors to solve the problem?

- The authors propose **RegionGPT (RGPT)**, a novel framework designed for region-level captioning and understanding. Key innovations include:
- Enhancing visual encoders with simple modifications to improve spatial awareness.
- Utilizing task-guided instruction prompts for better task-specific outputs.
- Developing an automated region caption data generation pipeline to create detailed region-specific captions, addressing the lack of rich region-level annotations