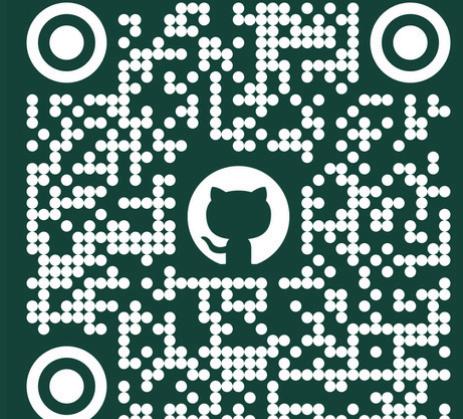


# Skyfall: A Hierarchical Methodology for Probing the Limits of VLMs in Earth Observation

Devansh Lodha<sup>1</sup>, Rishabh Mondal<sup>1</sup>, Ujjwal Gupta<sup>1,2</sup>, Nipun Batra<sup>1</sup><sup>1</sup>Indian Institute of Technology Gandhinagar, <sup>2</sup>Space Applications Centre, ISRO

Scan for the open-source code.

## Problem Statement: Opaque Models

The deployment of Vision-Language Models (VLMs) in scientific domains such as Earth Observation is critically impeded by their black-box nature. Standard evaluation offers insights into what a model predicts, but fails to provide a formal, reproducible methodology to diagnose how and why it fails in its geospatial reasoning.

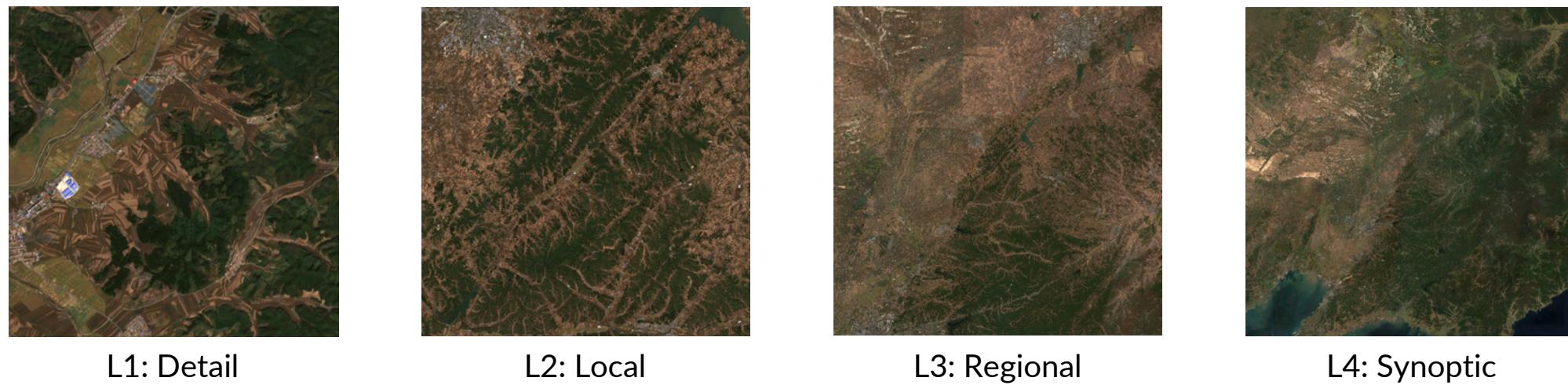
## Our Contribution

We introduce **Skyfall**, a comprehensive, open-source framework designed to address this gap.

### Skyfall Global Scene Packets

View Name	Assigned Tasks	Data Source	GSD (m)	Area (km <sup>2</sup> )
L1: Detail	Land Cover Classification	ESA WorldCover	10	5.12 × 5.12
L2: Local	State/Province Classification	FAO GAUL	100	51.2 × 51.2
L3: Regional	Biome, Climate Zone Classification	RESOLVE / Köppen-Geiger	500	256 × 256
L4: Synoptic	Country, Continent Classification	USDOS LSIB	1500	768 × 768

Views obtained from Sentinel - 2 (COPERNICUS/S2\_SR\_HARMONIZED)



L1: Detail

L2: Local

L3: Regional

L4: Synoptic

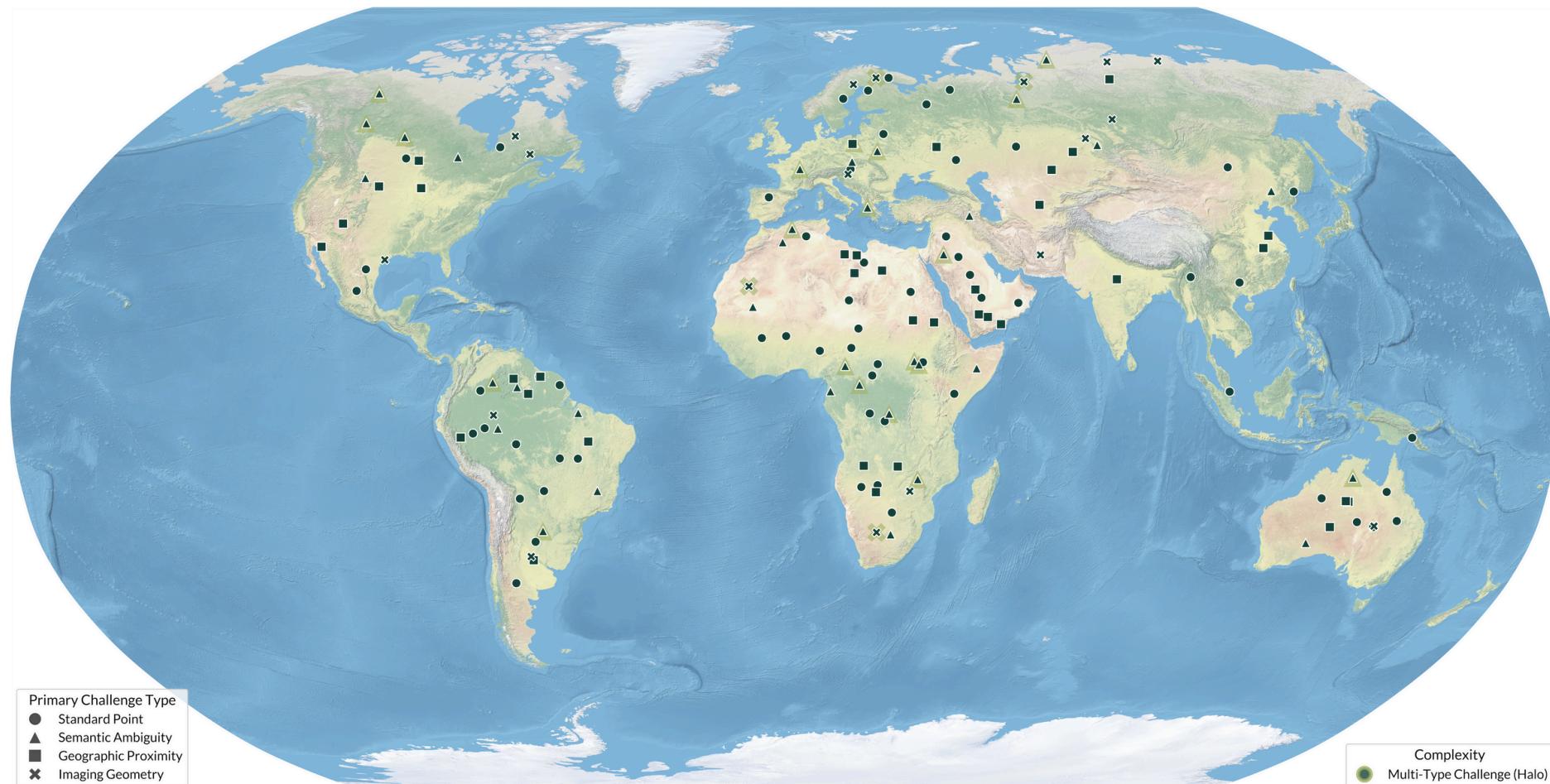
### Skyfall Atlas

A knowledge graph containing the geographic, semantic, and probabilistic data including entity-level environmental distributions, a climate distance matrix, and real-world conditional probabilities.

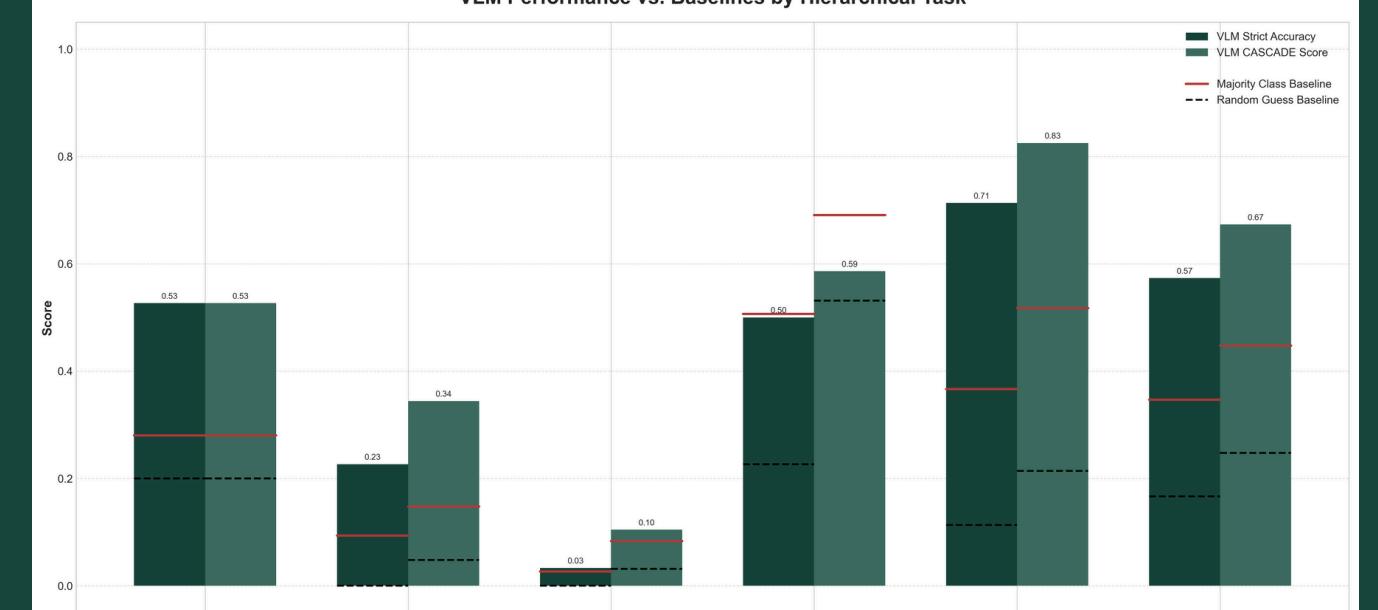
### Skyfall - X

A framework to quantify performance degradation against explicit, real-world challenges (e.g., long shadows, semantic ambiguity, mosaic artifacts).

#### Global Profile of the Skyfall Dataset: Semantics & Challenges



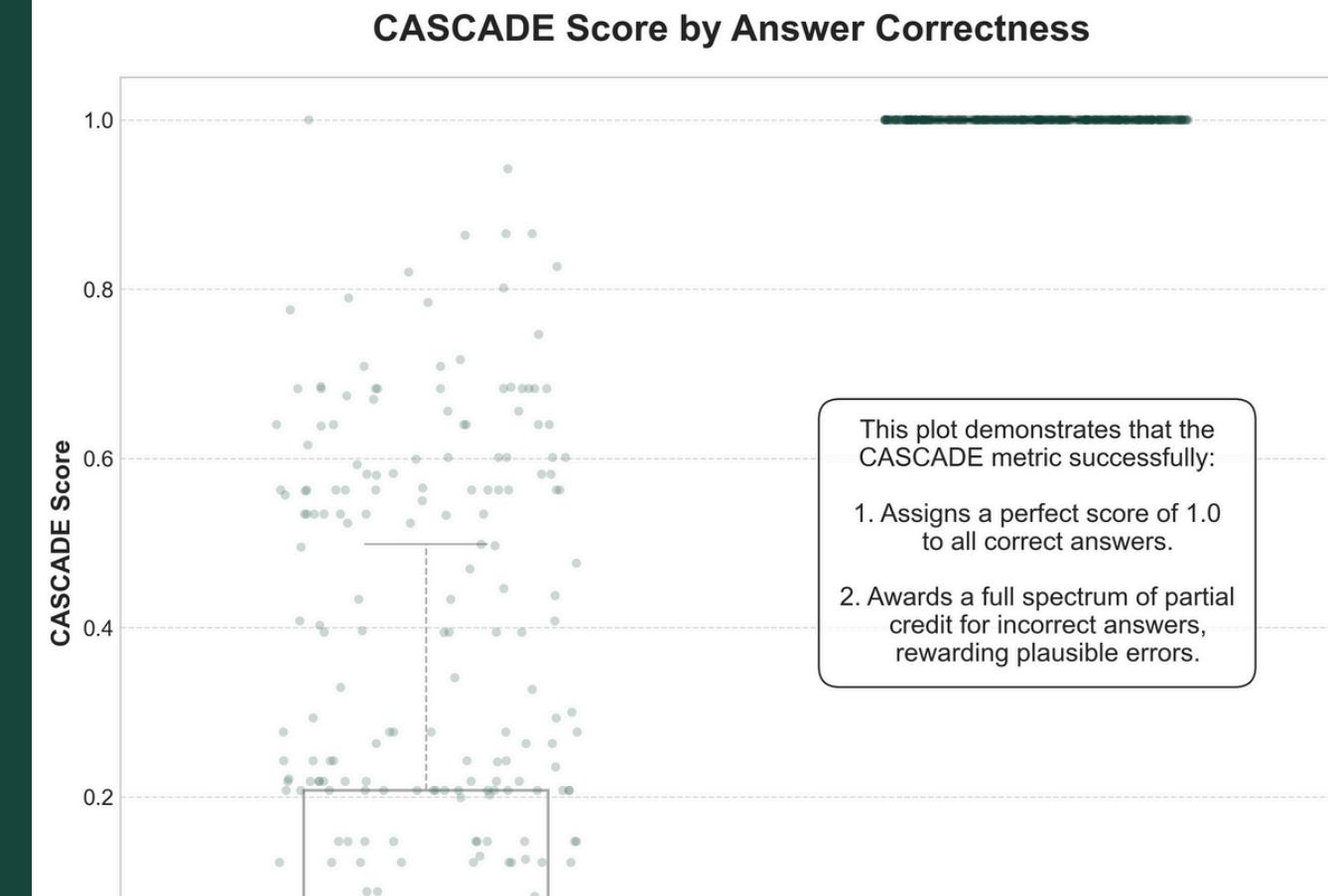
### VLM Performance vs. Baselines by Hierarchical Task



The VLM significantly outperforms baselines across all tasks, but performance degrades as geographic specificity increases from Continent to Country and State/Province.

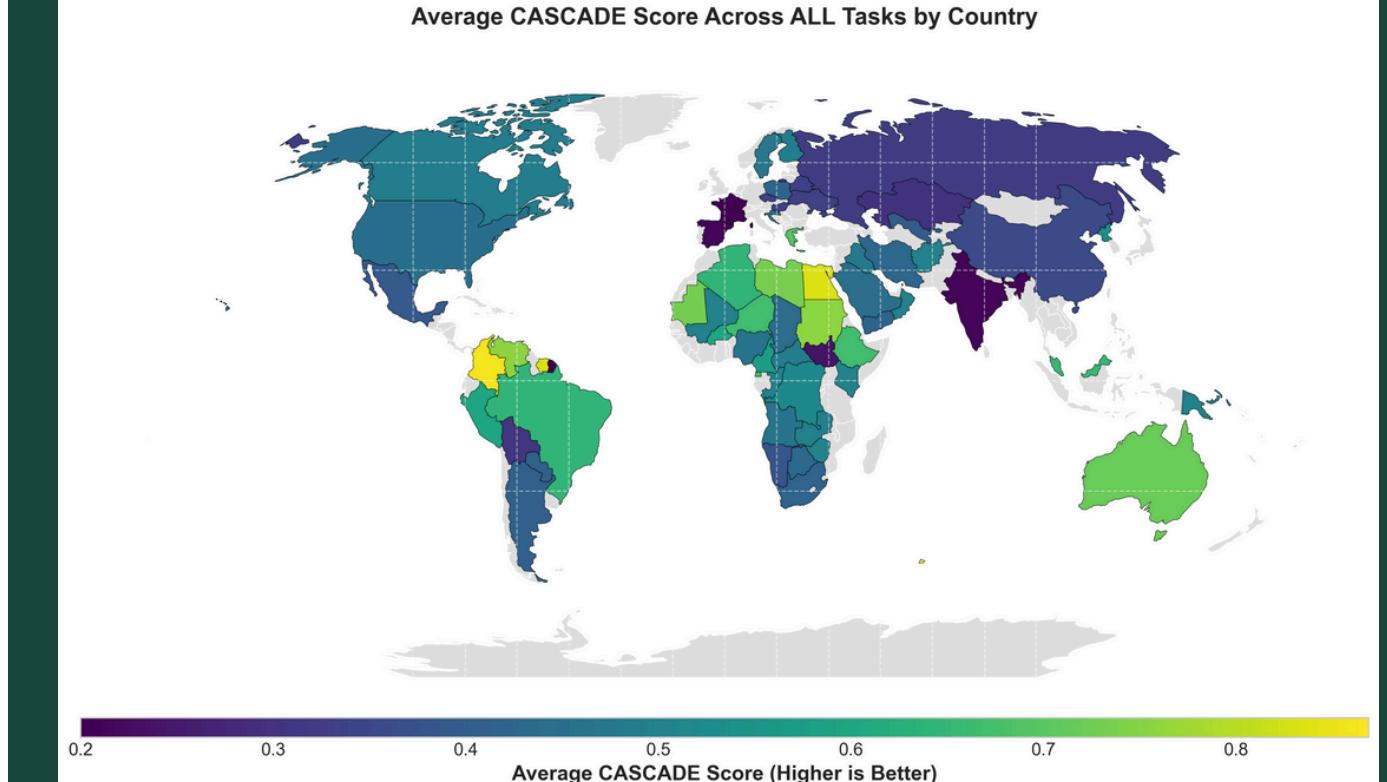
## Results and Discoveries

### CASCADE Score by Answer Correctness



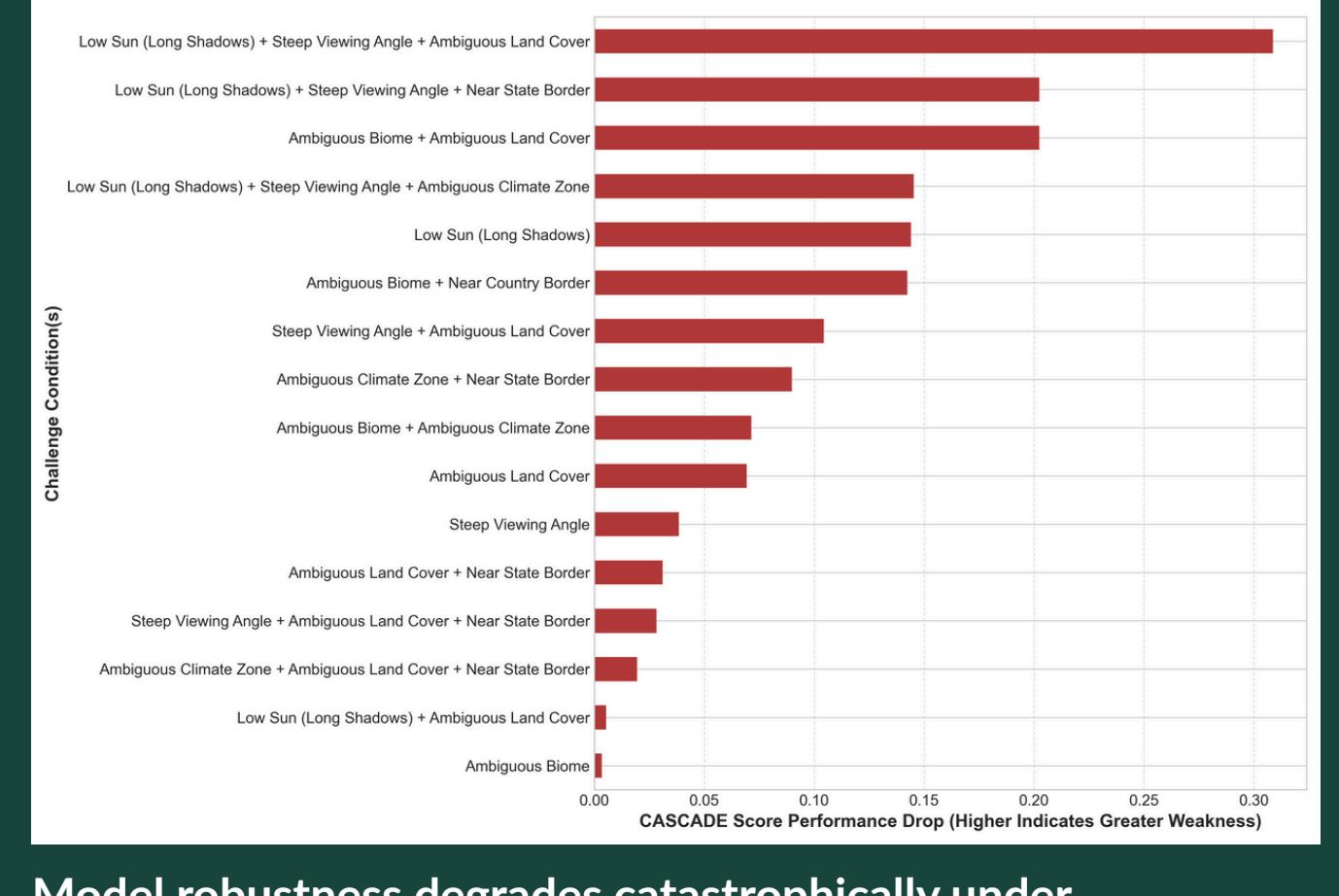
This plot demonstrates that the CASCADING metric successfully:  
1. Assigns a perfect score of 1.0 to all correct answers.  
2. Awards a full spectrum of partial credit for incorrect answers, rewarding plausible errors.

### Average CASCADING Score Across ALL Tasks by Country



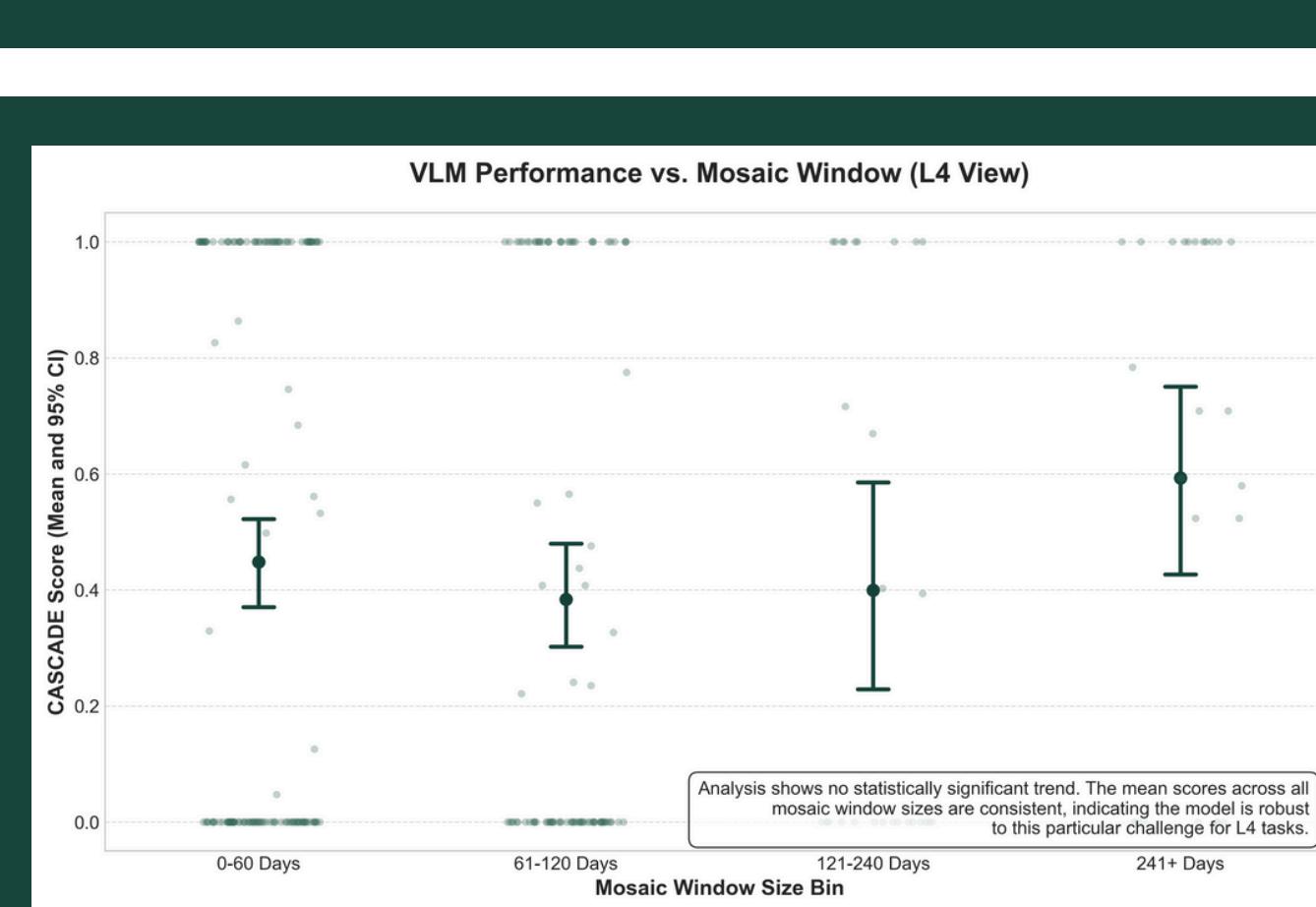
The model's performance reveals a strong training data bias, favoring Western nations while underperforming in diverse, non-Latin script regions like India, China, and Russia.

### VLM Robustness to Compounding Challenges



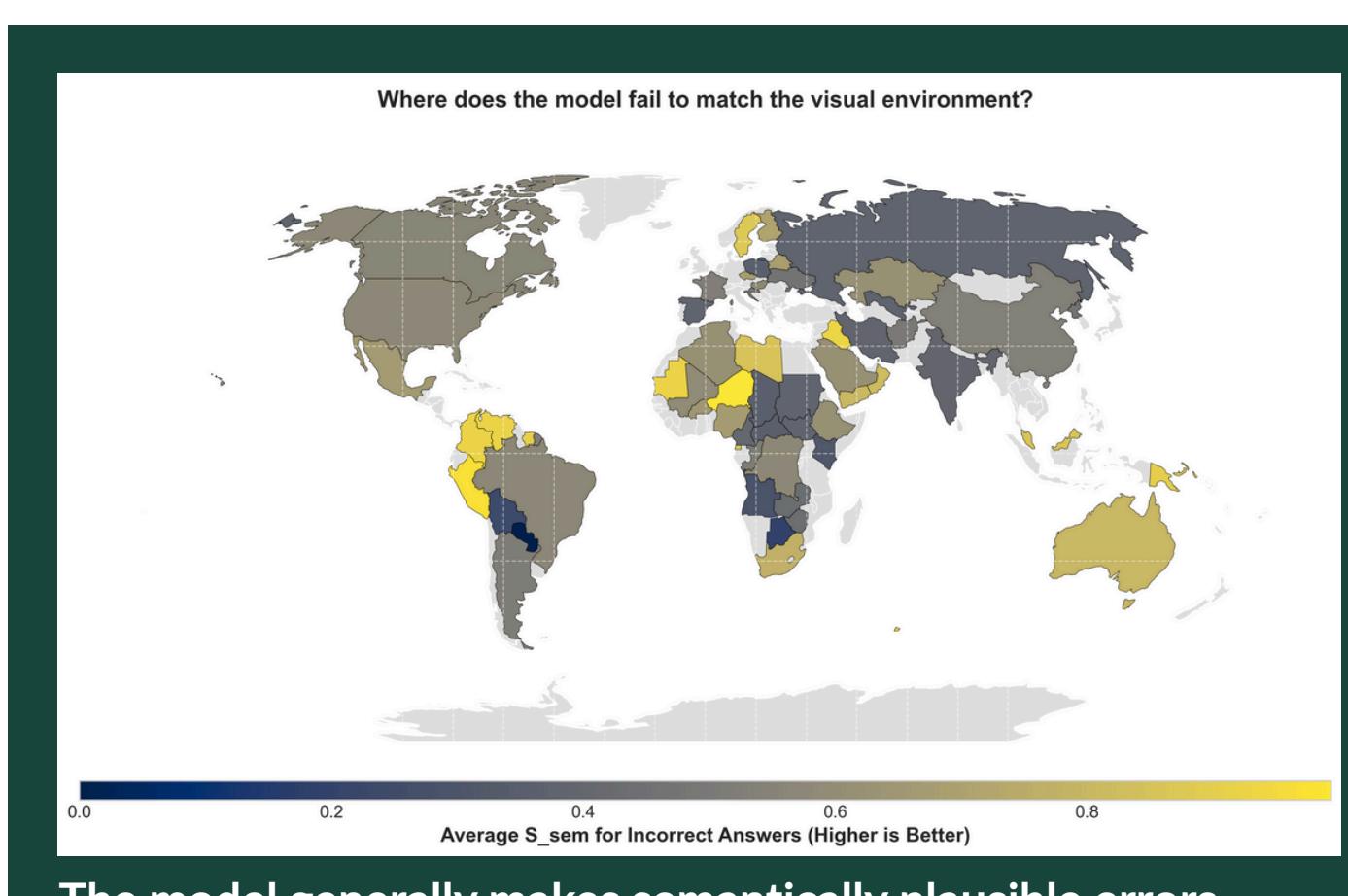
Model robustness degrades catastrophically under compounding challenges.

### VLM Performance vs. Mosaic Window (L4 View)



The model is robust to temporal incoherence in mosaics.

### Where does the model fail to match the visual environment?



The model generally makes semantically plausible errors indicating strong visual but poor geographic reasoning.