

# **generative geolocation**

## Interpreting PLONK: Geographic Knowledge in Vision Encoders

**BOUKHARI Imed Eddine**

19/01/2026

# Table of Contents

- ① Introduction & Motivation
- ② PLONK Overview
- ③ Linear Probing Experiments
- ④ Attention Analysis
- ⑤ Discussion
- ⑥ Conclusions

# Introduction & Motivation

# Outline

**Research Focus:** Understanding geographic knowledge in PLONK's encoder

## Two Experiments:

### ① Linear Probing

- How much does the frozen encoder already know?
- Train simple classifier on StreetCLIP features
- Test: Country, Region, City classification

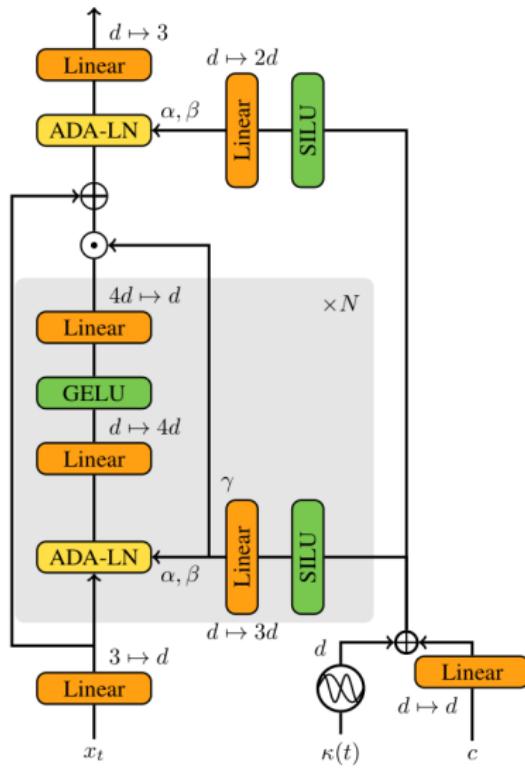
### ② Attention Analysis

- What visual features does the model use?
- Extract attention maps from 50,000 images
- Analyze: Entropy, concentration, spatial patterns

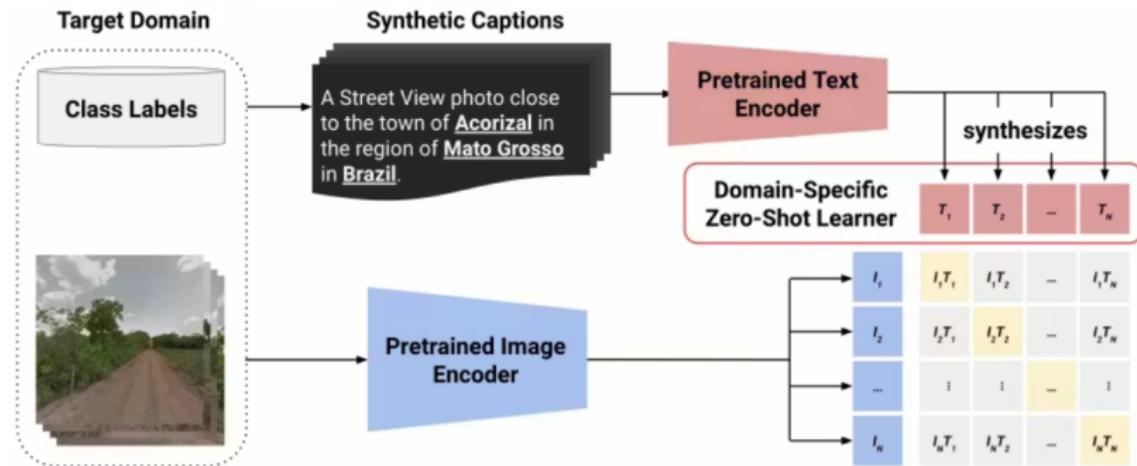
**Goal:** Decompose PLONK to understand where geographic knowledge comes from

# PLONK Overview

# PLONK: Generative Geolocation on Earth's Sphere



# StreetCLIP: Vision Transformer Encoder



# Linear Probing Experiments

# Linear Probing: Where Does Geographic Knowledge Come From?

**Question:** Does StreetCLIP already know geography?

**Method:**

- Extract features from *frozen* StreetCLIP encoder
- Train simple linear classifier (logistic regression)
- Evaluate on three levels: Country, Region, City

If accuracy is high → encoder already learned geography

If accuracy is low → flow matching learns from scratch

# Linear Probing: Dataset & Setup

## Dataset: OSV-5M Test Set

- 50,000 images sampled from test set
- Geographic labels: Country, Region, City
- 217 countries, ~2,050 regions, 7,292 cities

## Method:

- ① Extract 1024-dim features from frozen StreetCLIP encoder
- ② Filter classes: Keep only labels with  $\geq 2$  samples
- ③ Train/test split: 80/20 (40k train, 10k test)
- ④ Train: Logistic Regression
- ⑤ Evaluate: Accuracy on held-out test set

Goal: Measure what encoder learned *before* flow matching

# Linear Probing Results

Level	Linear Probe	PLONK (10000)	PLONK Paper	Classes
Country	<b>85.0%</b>	78%	76.2%	217
Region	<b>62%</b>	39%	44.2%	2,050
City	<b>8.2%</b>	6%	5.4%	7,292

## Observations:

- Linear probe achieves strong accuracy on different levels
- Encoder captures substantial geographic knowledge from pre-training
- Different tasks: classification vs. coordinate regression

# Understanding the Comparison

## What We're Measuring:

- **Linear Probe:** Direct classification on frozen features
  - Shows: What encoder learned during pre-training
  - Baseline: How much knowledge before flow matching
- **PLONK Full Model:** GPS prediction → label lookup
  - Shows: End-to-end system performance
  - Reference: Published results from paper

Comparison helps quantify: How much geographic knowledge comes from the encoder vs. the flow matching process?

# Attention analysis

# Attention Analysis: Motivation

**Question:** What parts of images does StreetCLIP look at for geolocation?

## Why This Matters:

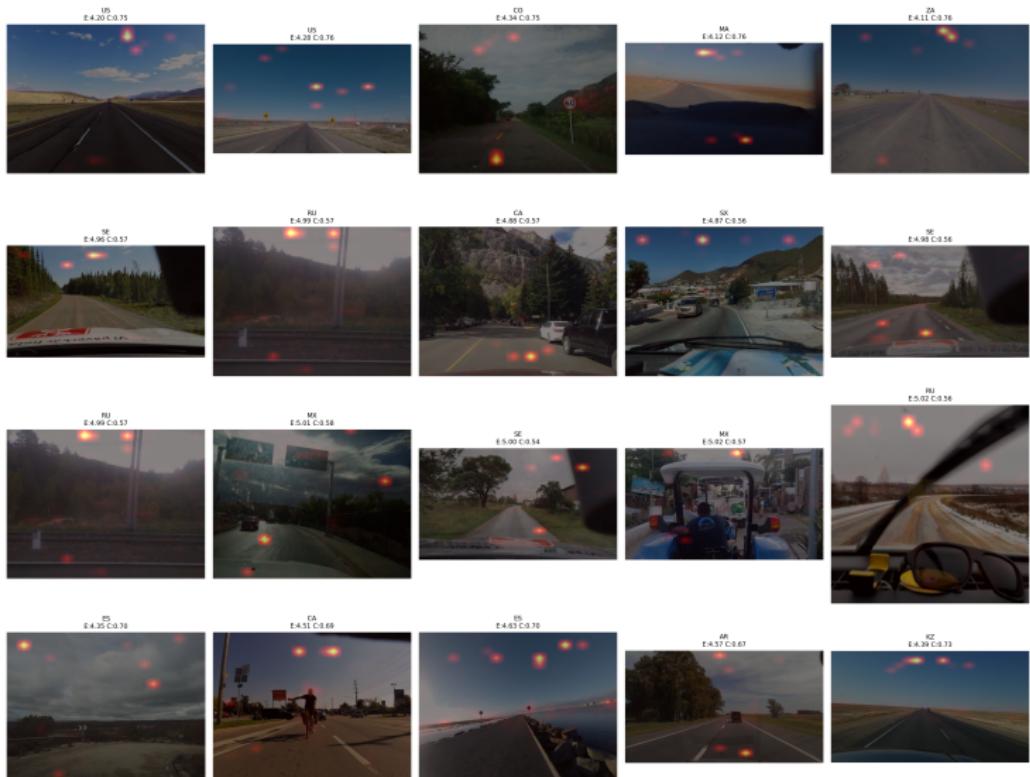
- Does it only focus on obvious landmarks?
- Does it use contextual cues (sky, vegetation)?
- Is attention consistent across images?

## Approach:

- ① Extract attention weights from last layer (50,000 images)
- ② Visualize where model "looks" on  $16 \times 16$  image patches
- ③ Compute statistics: spread, focus, patterns

Understanding attention helps explain how the encoder learns geography

# Diverse examples



# Attention Analysis Results

Metric	Value
Images Analyzed	<b>50,000</b>
Mean Entropy	<b>4.53</b>
Top-10% Concentration	<b>66.6%</b>
High Attention Patches	<b>95.2 (avg)</b>

## Interpretation:

- **Moderate entropy:** Model doesn't hyper-focus on single features
- **Selective attention:** 67% weight on top 10% of patches
- **Contextual coverage:** Uses 95 patches (37% of image)
- **Finding:** Model combines landmarks AND environmental cues (sky, vegetation, architectural patterns)

# Discussion

# PLONK's Two-Stage Architecture

## Design Overview:

### Stage 1: Encoder

- Pre-trained StreetCLIP
- Extracts visual features
- Captures geographic patterns
- Frozen during training

### Stage 2: Flow Matching

- Maps features to coordinates
- Riemannian geometry (sphere)
- Generates probability distributions
- Trained end-to-end

Combining strong pre-trained encoder with geometric flow matching enables both high accuracy and uncertainty estimation

# What Our Analysis Reveals

- Encoder already captures substantial geographic knowledge
- Strong performance
- Model attends to diverse visual features (not just landmarks)
- Uses environmental context: sky patterns, vegetation, architecture
- Attention distributed across 37% of image patches

**Implication:** The pre-trained encoder provides a strong foundation that flow matching can build upon for coordinate prediction

# Conclusions

# Summary

## What We Did:

- Linear probing on frozen StreetCLIP features (country, region, city)
- Attention analysis on 50,000 images (1/4 of test set)

## Key Findings:

- Pre-trained encoder contains substantial geographic knowledge
- Attention distributed across contextual features (sky, vegetation, environment)
- Model uses diverse visual cues, not just landmarks
- Attention patterns consistent across different countries
- Linear probing validates PLOANK's frozen encoder design choice

**Main Insight:** Strong pre-trained encoder is critical for PLOANK's performance

# Future Directions

## Potential Extensions:

- Replace frozen encoder with fine-tuned version and retrain flow matching
- Analyze attention patterns at different encoder layers
- Study encoder behavior on failure cases

**Thank you for your attention!**

Questions?