

Aplicaciones de los embeddings satelitales de Google AlphaEarth Foundations en el análisis geoespacial: clasificación de uso del suelo, modelado de distribución de especies y diseño de muestreo

José Ramón Martínez Batlle

Universidad Autónoma de Santo Domingo (UASD)

Actualizado: 2025-11-19

¿Qué son los embeddings de Google AlphaEarth Foundations (AEF)?

- Los embeddings de **Google AlphaEarth Foundations (AEF)** condensan información espacio–temporal derivada de series multiespectrales completas (p.ej., Sentinel, Landsat, mapas de cobertura, LiDAR, bioclima).
- Cada píxel es representado por un **vector numérico de alta dimensión** que captura características espectrales, temporales y espaciales.
- Su principal fortaleza es que **permiten realizar tareas de clasificación con muy pocos datos de entrenamiento**, evitando el entrenamiento de modelos de deep learning desde cero.

AlphaEarth Foundations: An embedding field model for accurate and efficient global mapping from sparse label data

Christopher F. Brown^{*,1}, Michal R. Kazmierski^{*,1}, Valerie J. Pasquarella^{*,2}, William J. Rucklidge², Masha Samsikova¹, Chenhui Zhang¹, Evan Shelhamer¹, Estefania Lahera², Olivia Wiles¹, Simon Ilyushchenko², Noel Gorelick², Lihui Lydia Zhang¹, Sophia Alj¹, Emily Schechter², Sean Askay², Oliver Guinan², Rebecca Moore², Alexis Boukouvalas¹ and Pushmeet Kohli¹

^{*}Equal contributions, ¹Google DeepMind, ²Google

Unprecedented volumes of Earth observation data are continually collected around the world, but high-quality labels remain scarce given the effort required to make physical measurements and observations. This has led to considerable investment in bespoke modeling efforts translating sparse labels into maps. Here we introduce AlphaEarth Foundations, an embedding field model yielding a highly general, geospatial representation that assimilates spatial, temporal, and measurement contexts across multiple sources, enabling accurate and efficient production of maps and monitoring systems from local to global scales. The embeddings generated by AlphaEarth Foundations are the only to consistently outperform all previous featurization approaches tested on a diverse set of mapping evaluations without re-training. We will release a dataset of global, annual, analysis-ready embedding field layers from 2017 through 2024.

Introduction

Management of global food supplies, public health, and disaster response all start from maps that geographically anchor questions like "which forests pose an unacceptable wildfire risk?" or "where are soybeans grown?". The launch of the first Landsat satellite in 1972 marked the dawn of an era where spaceborne monitoring could serve the interests of global environmental policy-making and provide critical insights into our changing planet (Cohen and Goward, 2004). Over the following decades Earth observation (EO) data became widely available, and streams from both historic and modern EO instruments are now routinely used to create maps that answer questions about the past, present, and future of Earth's ecosystems and climate (Wulder et al., 2022). Nonetheless, advancements in deriving planetary-scale insights from petabytes of satellite imagery and other environmental datasets remain hamstrung by the relative scarcity of ground-based measurements and annotations, and a new problem: the overwhelming volume of geospatial data (Tuia et al., 2024). In this work, we introduce a foundational geospatial

embedding model that solves fundamental challenges in the institution of mapping through the generation of a universal feature space. The features produced by our model consistently achieve top performance in all application domains tested when compared to other general and even domain specific approaches (Figure 1A). This marks a shift from the previous state-of-the-art for which no single approach was dominant.

From sparse labels to maps

High-quality maps depend on high-quality labeled data, yet when working at global scales, a balance must be struck between measurement precision and spatial coverage. Many global mapping efforts focus on individual ecosystems like forests (Hansen et al., 2013), water (Pekel et al., 2016), tidal wetlands (Murray et al., 2022a) or other broad legends, e.g., (Brown et al., 2022; Zanaga et al., 2022). This simplifies the label collection process, allowing trained interpreters to collect larger volumes at scale at the expense of descriptive power for certain use cases. In the cases where high-quality annotations and/or field

AlphaEarth Foundations: An embedding field model for accurate and efficient global mapping from sparse label data

Christopher F. Brown^{*†}, Michal R. Kazmierski^{*†}, Valerie J. Pasquarella^{*‡}, William J. Ruckridge², Masha Samsikova¹, Chenhui Zhang¹, Evan Shelhamer¹, Estefania Lahera², Olivia Wiles¹, Simon Ilyushchenko², Noel Gorelick², Lihui Lydia Zhang¹, Sophia Alj¹, Emily Schechter², Sean Askay², Oliver Guinan², Rebecca Moore², Alexis Boukouvalas¹ and Pushmeet Kohli¹

^{*}Equal contributions, [†]Google DeepMind, [‡]Google

Unprecedented volumes of Earth observation data are continually collected around the world, but high-quality labels remain scarce given the effort required to make physical measurements and observations. This has led to considerable investment in bespoke modeling efforts translating sparse labels into maps. Here we introduce AlphaEarth Foundations, an embedding field model yielding a highly general, geospatial representation that assimilates spatial, temporal, and measurement contexts across multiple sources, enabling accurate and efficient production of maps and monitoring systems from local to global scales. The embeddings generated by AlphaEarth Foundations are the only to consistently outperform all previous featurization approaches tested on a diverse set of mapping evaluations without re-training. We will release a dataset of global, annual, analysis-ready embedding field layers from 2017 through 2024.

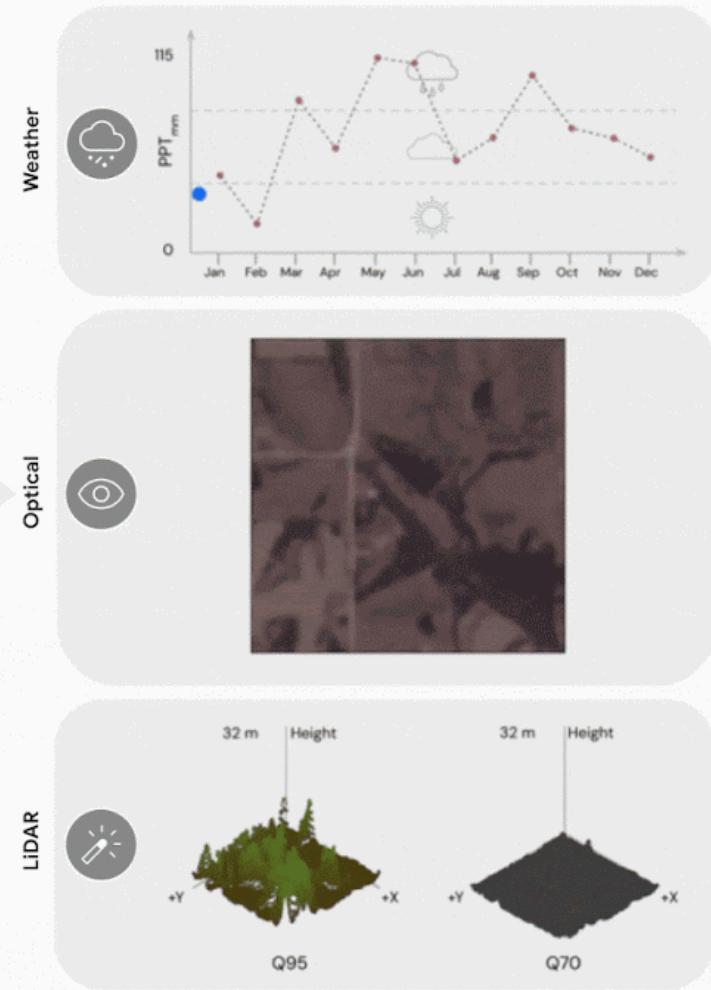
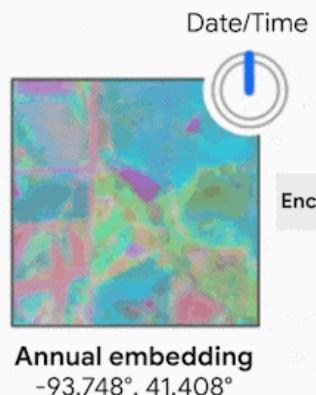
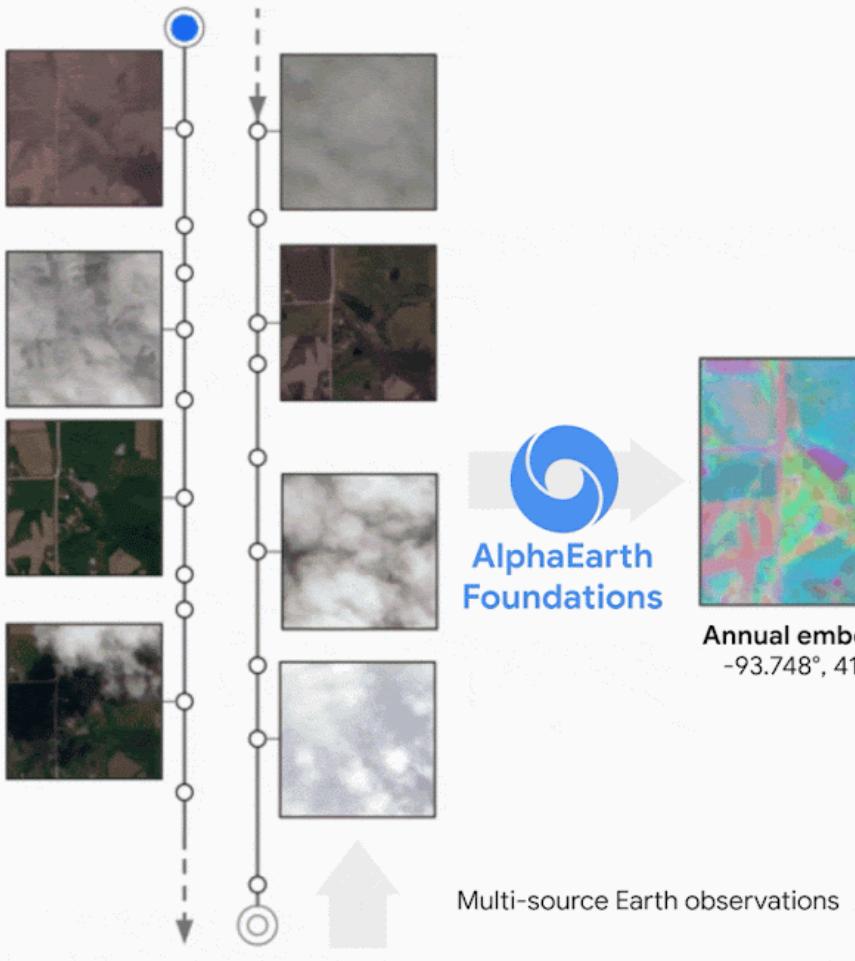
Introduction

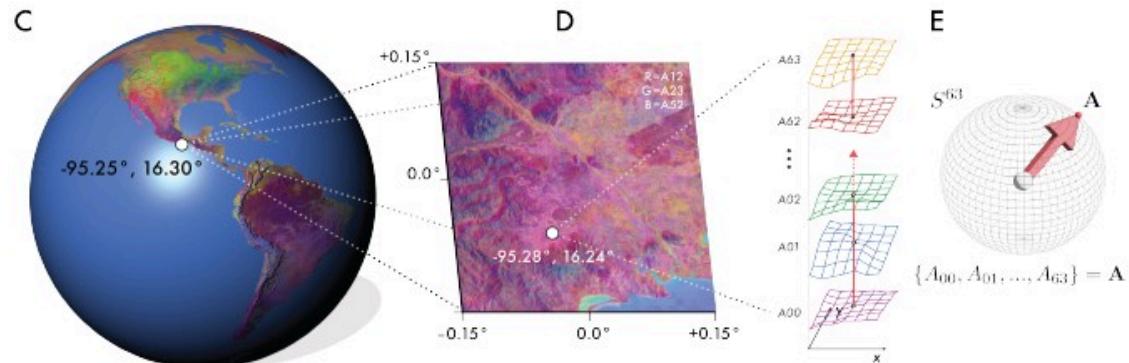
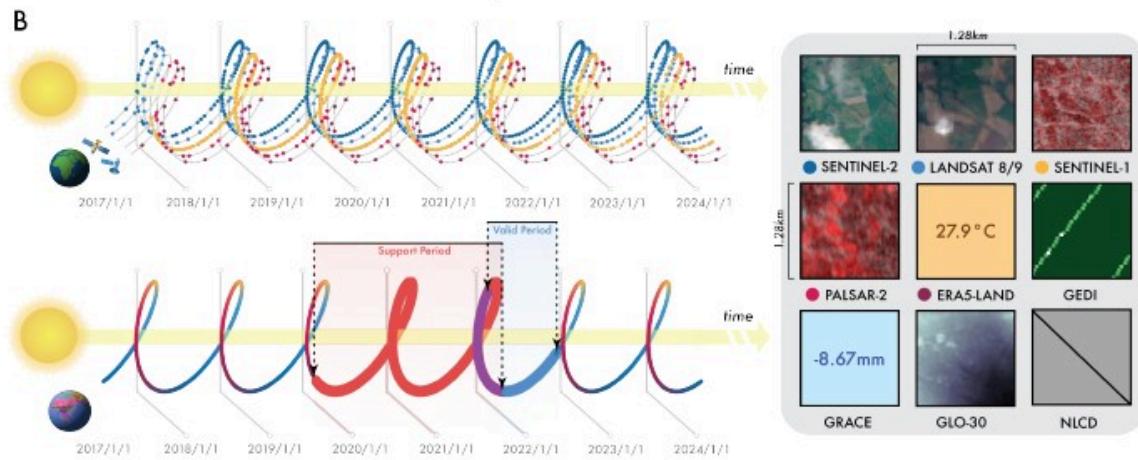
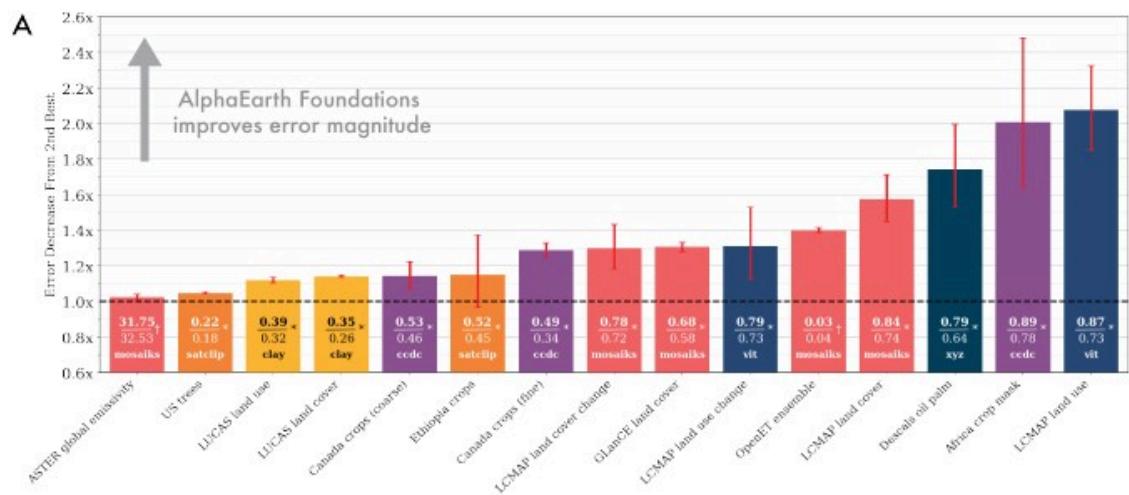
Management of global food supplies, public health, and disaster response all start from maps that geographically anchor questions like "which forests pose an unacceptable wildfire risk?" or "where are soybeans grown?". The launch of the first Landsat satellite in 1972 marked the dawn of an era where spaceborne monitoring could serve the interests of global environmental policy-making and provide critical insights into our changing planet (Cohen and Goward, 2004). Over the following decades Earth observation (EO) data became widely available, and streams from both historic and modern EO instruments are now routinely used to create maps that answer questions about the past, present, and future of Earth's ecosystems and climate (Wulder et al., 2022). Nonetheless, advancements in deriving planetary-scale insights from petabytes of satellite imagery and other environmental datasets remain hamstrung by the relative scarcity of ground-based measurements and annotations, and a new problem: the overwhelming volume of geospatial data (Tuia et al., 2024). In this work, we introduce a foundational geospatial

embedding model that solves fundamental challenges in the institution of mapping through the generation of a universal feature space. The features produced by our model consistently achieve top performance in all application domains tested when compared to other general and even domain specific approaches (Figure 1A). This marks a shift from the previous state-of-the-art for which no single approach was dominant.

From sparse labels to maps

High-quality maps depend on high-quality labeled data, yet when working at global scales, a balance must be struck between measurement precision and spatial coverage. Many global mapping efforts focus on individual ecosystems like forests (Hansen et al., 2013), water (Pekel et al., 2016), tidal wetlands (Murray et al., 2022a) or other broad legends, e.g., (Brown et al., 2022; Zanaga et al., 2022). This simplifies the label collection process, allowing trained interpreters to collect larger volumes at scale at the expense of descriptive power for certain use cases. In the cases where high-quality annotations and/or field

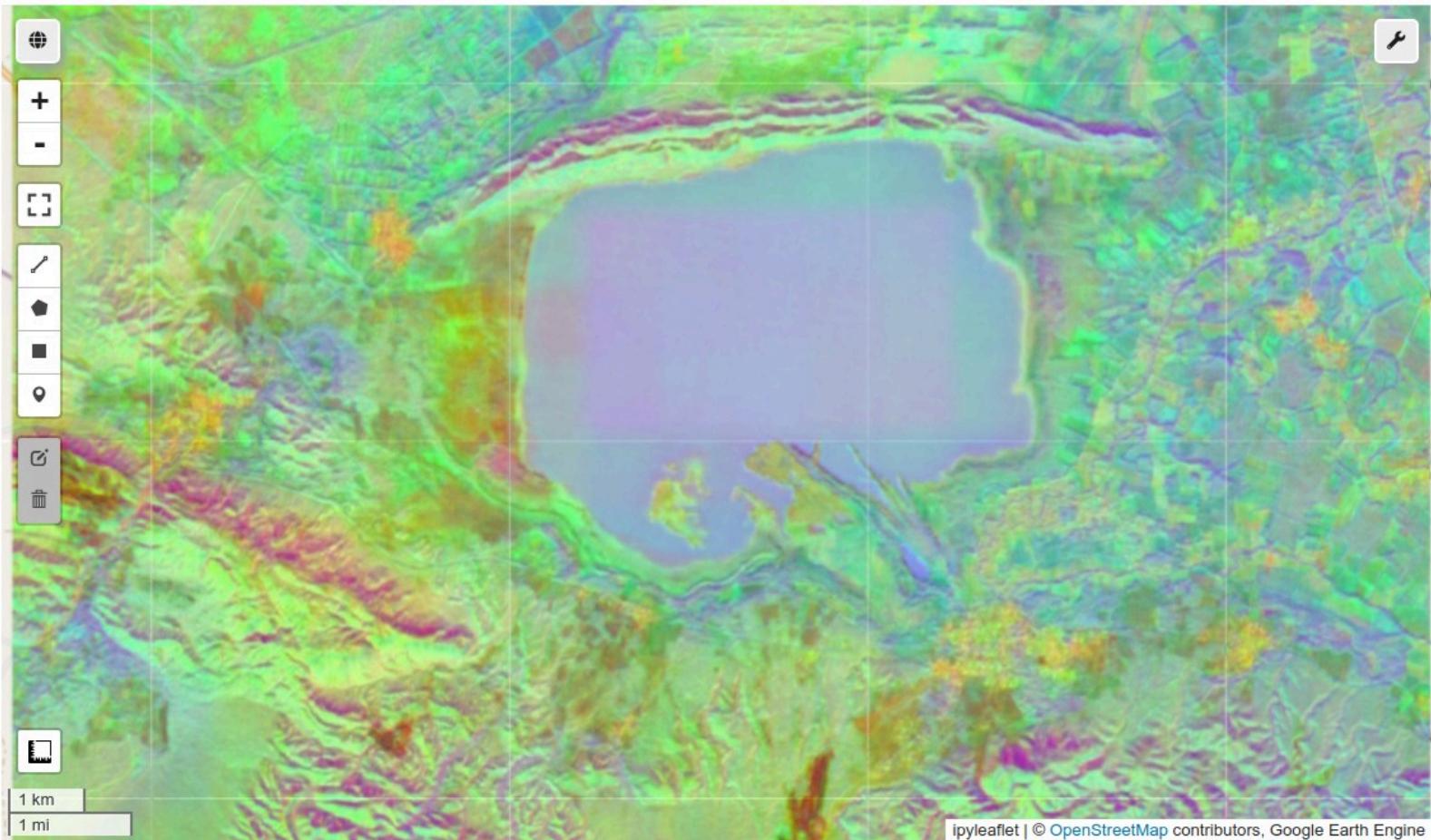




Visualización "tal cual" de los embeddings (composición RGB)

↑ ↓ ← → ⌂

```
embedding_rgb_bands = ["A00", "A21", "A42"]
emb_vis_params = {
    "bands": embedding_rgb_bands,
    "min": -0.25,
    "max": 0.25,
}
m2 = geemap.Map(center=[(bbox[1] + bbox[3]) / 2, (bbox[0] + bbox[2]) / 2], zoom=11)
# Embeddings en RGB pseudo-color
m2.addLayer(emb_image, emb_vis_params, "Embeddings AEF (RGB)")
# AOI
m2.addLayer(aoi, {"color": "black"}, "AOI", False)
m2
```



Objetivos

- Objetivo general: evaluar la utilidad práctica de los **embeddings de AlphaEarth Foundations** para apoyar diferentes procesos de análisis espacial en contextos ecológicos y territoriales de la República Dominicana.
- Objetivos específicos:
 - Probar su desempeño en la modelación de la **distribución espacial de especies** empleando registros de presencia y técnicas de aprendizaje supervisado.
 - Evaluar la capacidad para **discriminar coberturas** en zonas agrícolas, forestales y urbanas.
 - Analizar su aplicación en el **diseño de muestreo estratificado**, utilizando la similitud entre embeddings para optimizar la representatividad espacial y ambiental

Objetivo 1

Probar desempeño de embeddings de AEF en la modelación de la **distribución espacial de especies** empleando registros de presencia y técnicas de aprendizaje supervisado.

Colaboración para mis estudiantes del semestre 2023-02,
asignatura biogeografía: Adrián Montás, Ángel González,
Arisleydi De la Cruz, Bryan Ramos, Claribel Ramírez, Manuel
Reyes, Ramona Muñoz, Saderis Carmona, Yenny Santana

[biogeografia-202302 / manuscrito](#) Public template[Code](#) [Issues](#) [Pull requests](#) [Actions](#) [Projects](#) [Security](#) [Insights](#)[main](#) Branch[1 Branch](#)[0 Tags](#) Go to file[Code](#) **geofis** ADD manuscript assignment Saderis data/report

59ec7a6 · 2 years ago 28 Commits

R	UPDATE advancing results template	2 years ago
data	ADD assignment 1 sample collection	2 years ago
fuentes/manuscrito	ADD manuscript assignment Saderis data/report	2 years ago
odk	UPDATE symlinks deleted	2 years ago
.gitignore	UPDATE .gitignore	2 years ago
LICENSE	Initial commit	2 years ago
README.md	UPDATE README	2 years ago
README	GPL-3.0 license	

Asignaciones de manuscrito

[Asignación de manuscrito 1. Diseño de muestreo y colecta de datos de campo](#)

[Asignación de manuscrito 2. Técnicas de procesamiento y analíticas \(subsección de la Metodología\)](#)

[Asignación de manuscrito 3. Introducción](#)

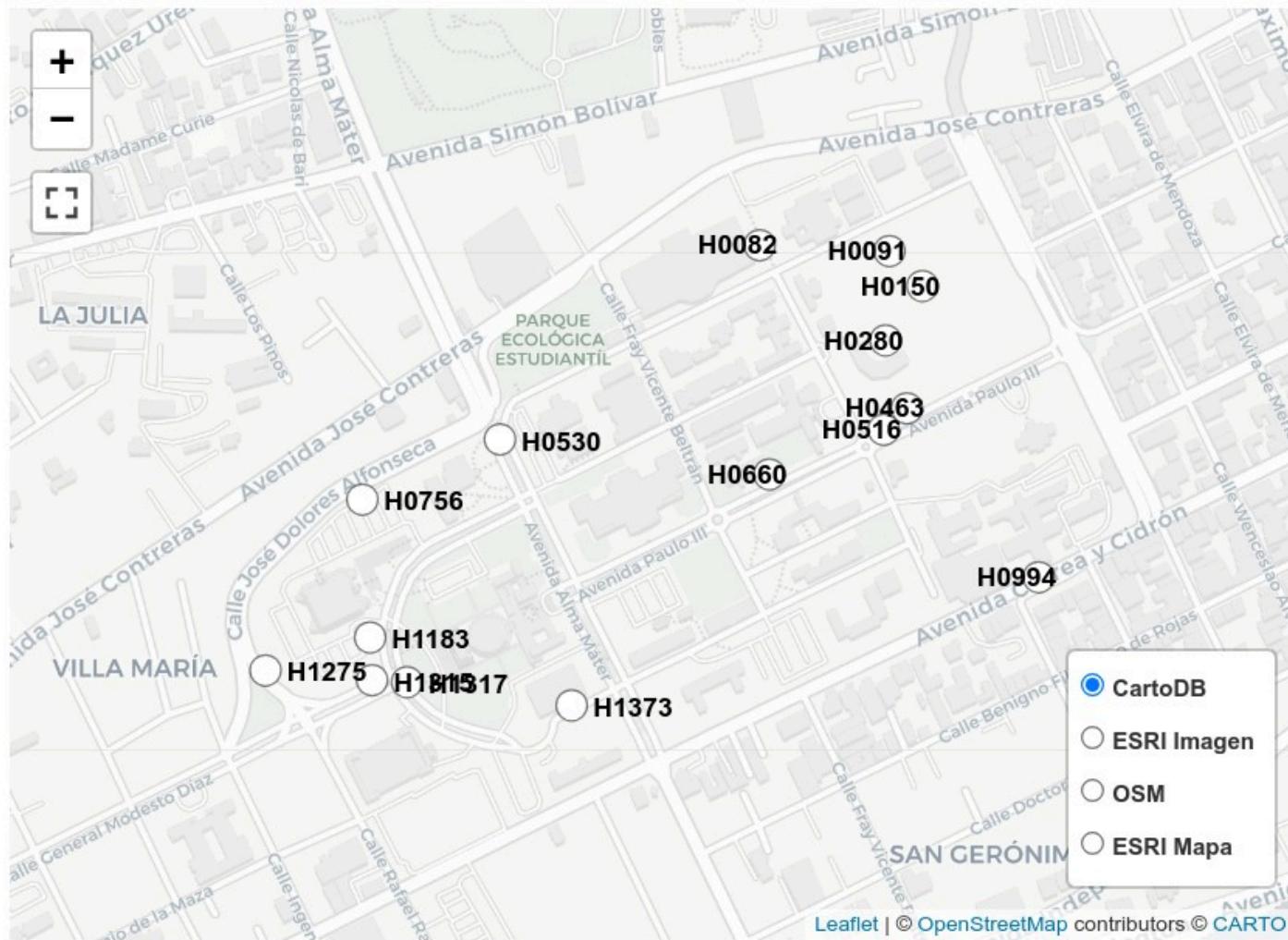
[Asignación de manuscrito 4. Resultados](#)

Estudiante	Ruta de informe
Adrian-Rafael-Diplan-Montas	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Adrian-Rafael-Diplan-Montas.html
Angel-Rolando-Gomez-Gonzalez	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Angel-Rolando-Gomez-Gonzalez.html
Arisleydi-Mejia-De-La-Cruz	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Arisleydi-Mejia-De-La-Cruz.html
Bryan-Josue-Funez-Ramos	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Bryan-Josue-Funez-Ramos.html
Claribel-Reyes-Ramirez	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Claribel-Reyes-Ramirez.html
Manuel-Enrique-Urena-Reyes	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Manuel-Enrique-Urena-Reyes.html
Ramona-Geraldo-Munoz	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Ramona-Geraldo-Munoz.html
Saderis-Carmona-Marte	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Saderis-Carmona-Marte.html
Yenny-Mabel-Santana	https://biogeografia-202302.github.io/manuscrito/fuentes/manuscrito/informe-resultados-Yenny-Mabel-Santana.html

2.2 Mapa de mis puntos colectados

Nota: sólo figuran los puntos con coordenadas. El total de formularios podría ser mayor al número de puntos mostrado en el mapa. Para un mapa comprensivo, ver más adelante el mapa bajo el texto “Mapa de hexágonos visitados”.

Code



2.3 Fotos



1696783022562.jpg



1696772480069.jpg

Code ▾

Datos y análisis para resultados de Adrian-Rafael-Diplan-Montas

Elaborado por: José-Ramón Martínez-Batlle (jmartinez19@uasd.edu.do)
Facultad de Ciencias, Universidad Autónoma de Santo Domingo (UASD)
Santo Domingo, República Dominicana

1 Carga de paquetes

Son muchos los paquetes empleados en estos análisis. Puedes consultar en el ChatGPT qué hace cada uno. Considera un aspecto también importante: algunas funciones escritas por mí se cargan con `source_url` y `source`; dentro de algunas de dichas funciones, también se cargan paquetes adicionales.

Code

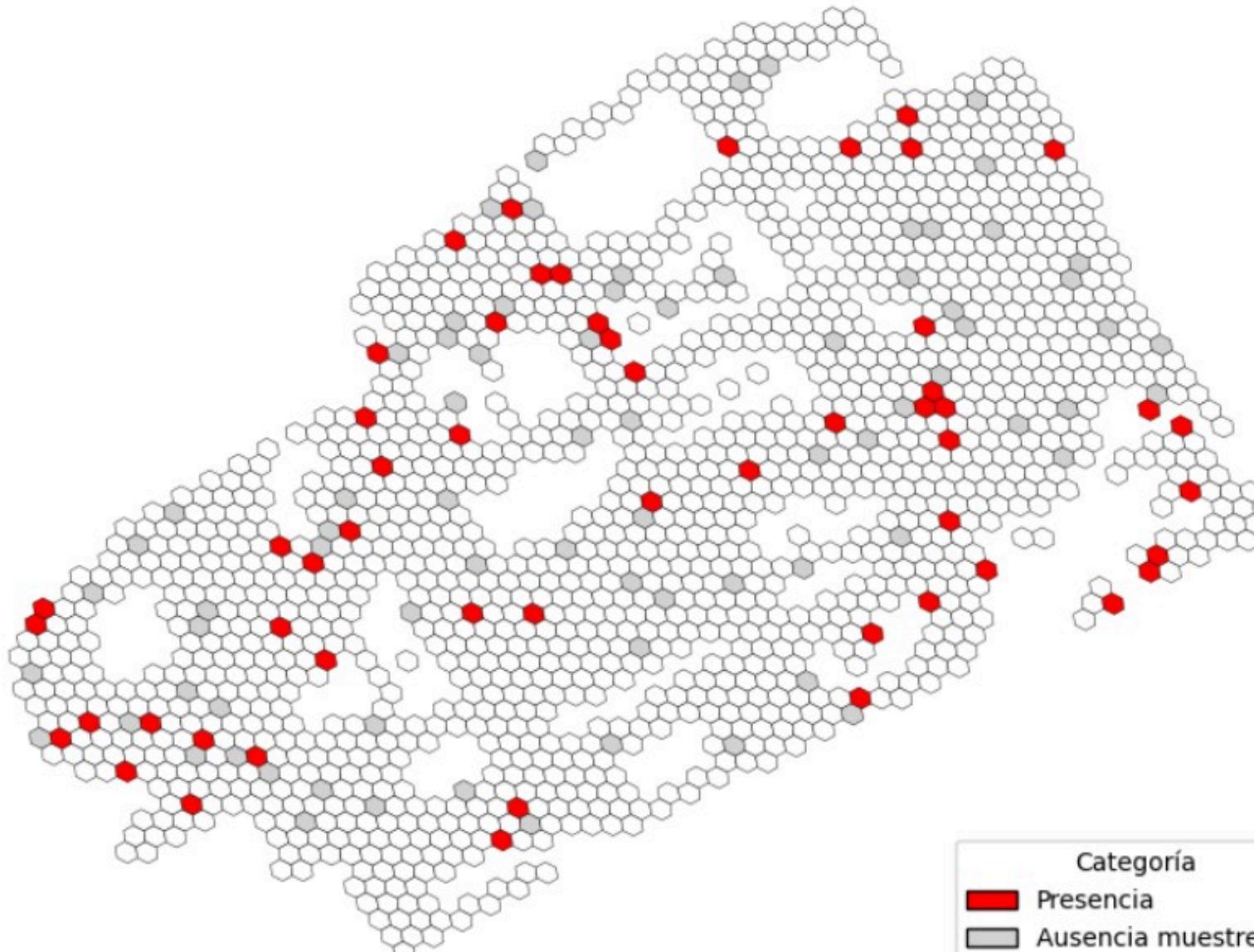
2 Leer y preparar datos

2.1 Formularios de Adrian Rafael Diplan Montas

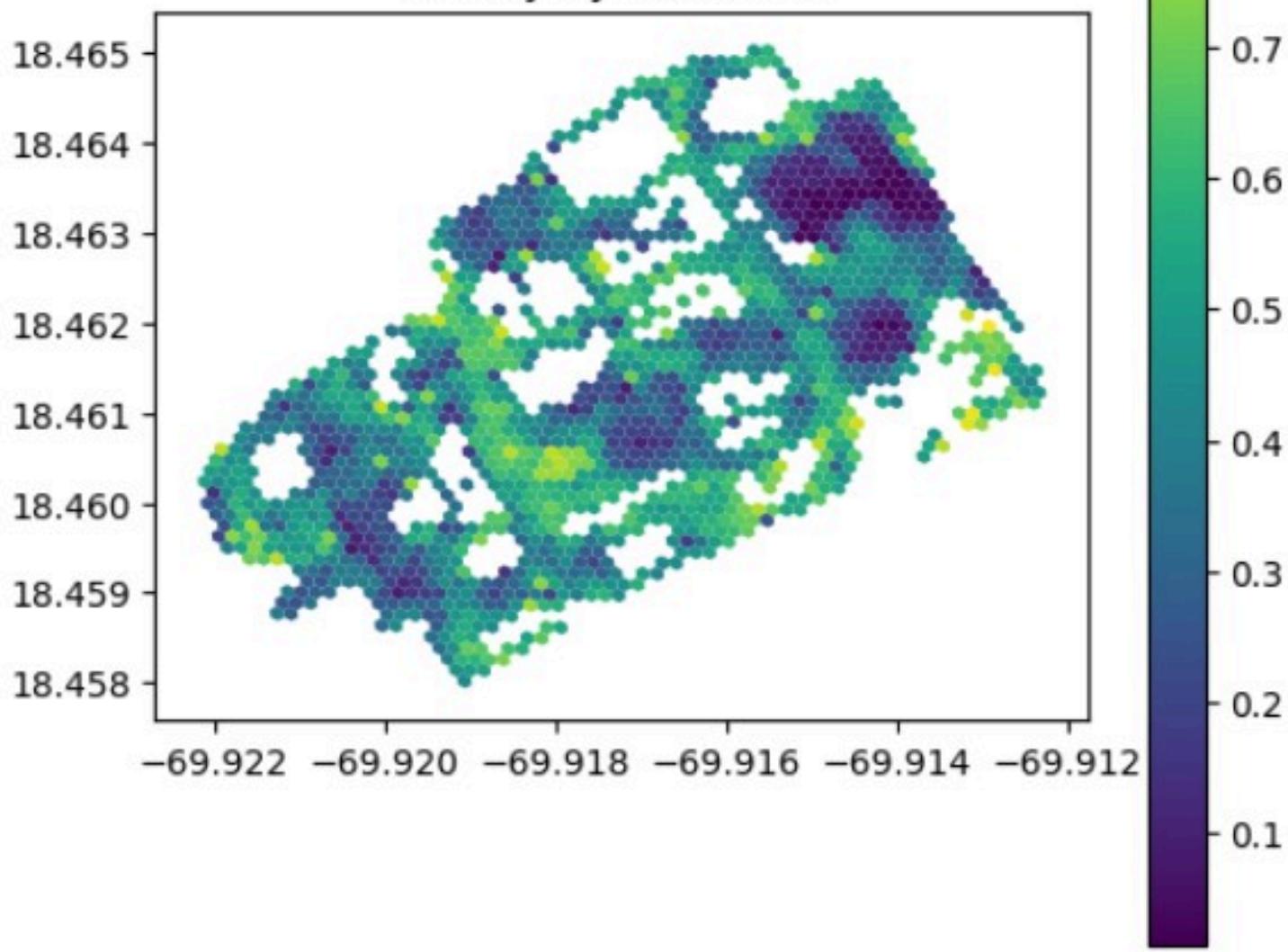
Code



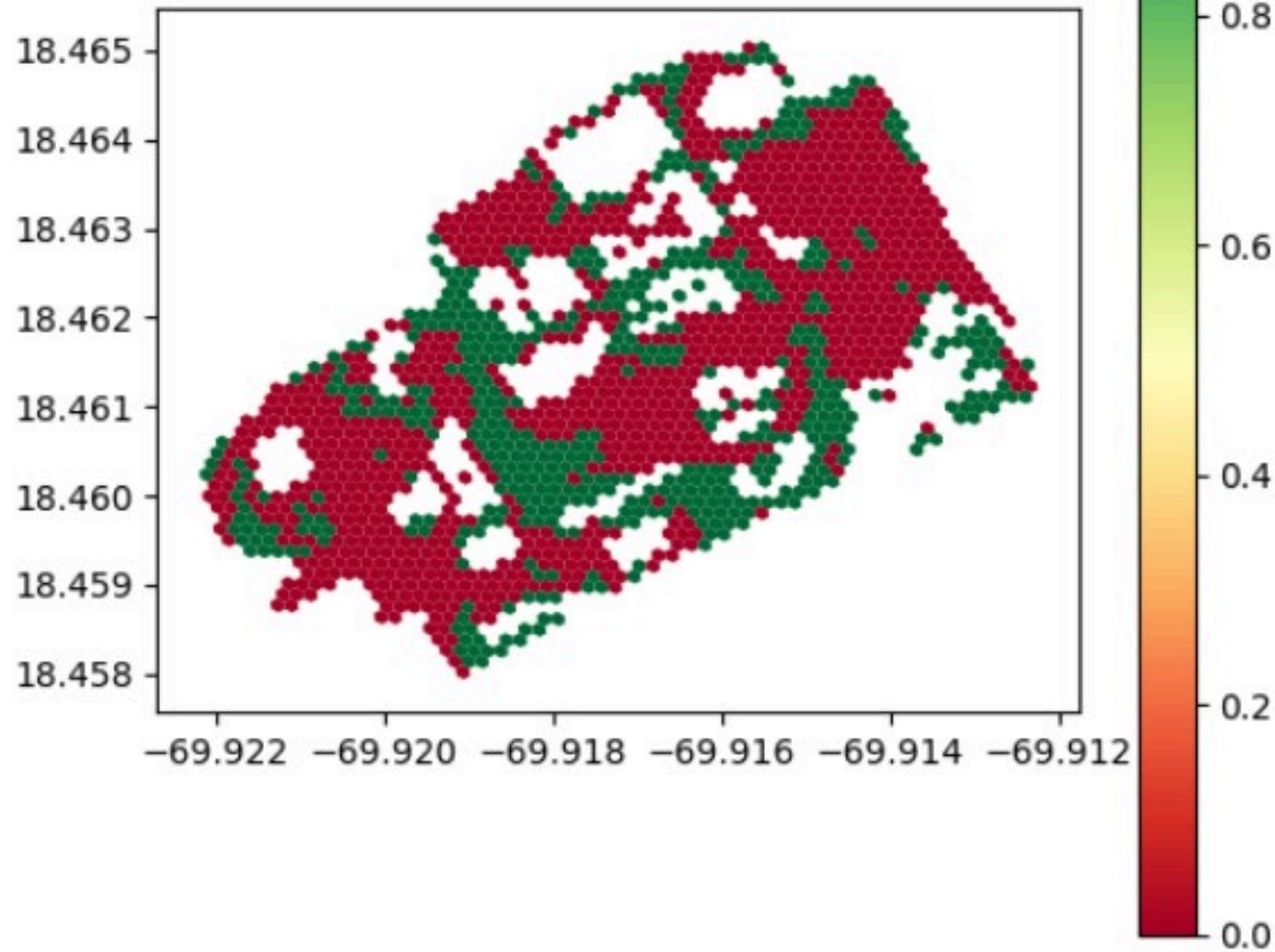
Registros reales de *Brachymyrmex heeri*



Probabilidad de presencia
Brachymyrmex heeri



Presencia/ausencia (umbral 0.5)



Modelo

```
[29]: results = model_species_distribution(  
        gdf_full,  
        especie_elegida,  
        classifier_name="et",  
        threshold=0.5  
    )
```

```
Presencias de Brachymyrmex heeri: 54  
Filas finales para entrenar: 130  
Nº de bandas de embeddings usadas: 64  
Matriz de confusión (test):  
[[16  7]  
 [ 8  8]]  
ROC AUC (test): 0.569
```

```
Especie: Brachymyrmex heeri
Hexágonos muestreados: 130 (presencias=54, ausencias=76)
Fold 1: Boyce = 0.517 | presencias val = 14
Fold 2: Boyce = 0.429 | presencias val = 9
Fold 3: Boyce = 0.810 | presencias val = 13
Fold 4: Boyce = 0.548 | presencias val = 12
Fold 5: Boyce = 1.000 | presencias val = 6

==== Resumen Boyce (validación cruzada) ====
Especie: Brachymyrmex heeri
Boyce medio = 0.660
Desviación estándar = 0.237

{'species': 'Brachymyrmex heeri',
 'boyce_mean': 0.6604761904761905,
 'boyce_sd': 0.23700703702824313,
 'folds': [{ 'fold': 1, 'boyce': 0.5166666666666667, 'n_pres_val': 14},
            { 'fold': 2, 'boyce': 0.4285714285714286, 'n_pres_val': 9},
            { 'fold': 3, 'boyce': 0.8095238095238096, 'n_pres_val': 13},
            { 'fold': 4, 'boyce': 0.5476190476190477, 'n_pres_val': 12},
            { 'fold': 5, 'boyce': 0.9999999999999999, 'n_pres_val': 6}]} 
```

Conclusiones del Objetivo 1

- **Los embeddings AEF reflejan gradientes ambientales finos** dentro del campus, capturando variación espacial relevante.
- **La extracción de valores ambientales por hexágono H3 fue exitosa**, generando una base homogénea para integrar las presencias.
- **Los modelos exploratorios muestran potencial predictivo** de los embeddings como variables para la distribución

Objetivo 2

Evaluar la capacidad de los embeddings satelitales de Google AEF para discriminar coberturas agrícolas, forestales y urbanas

Colaboración para Kénnida Polanco.

Contexto

Para este objetivo se seleccionó un área con presencia simultánea de:

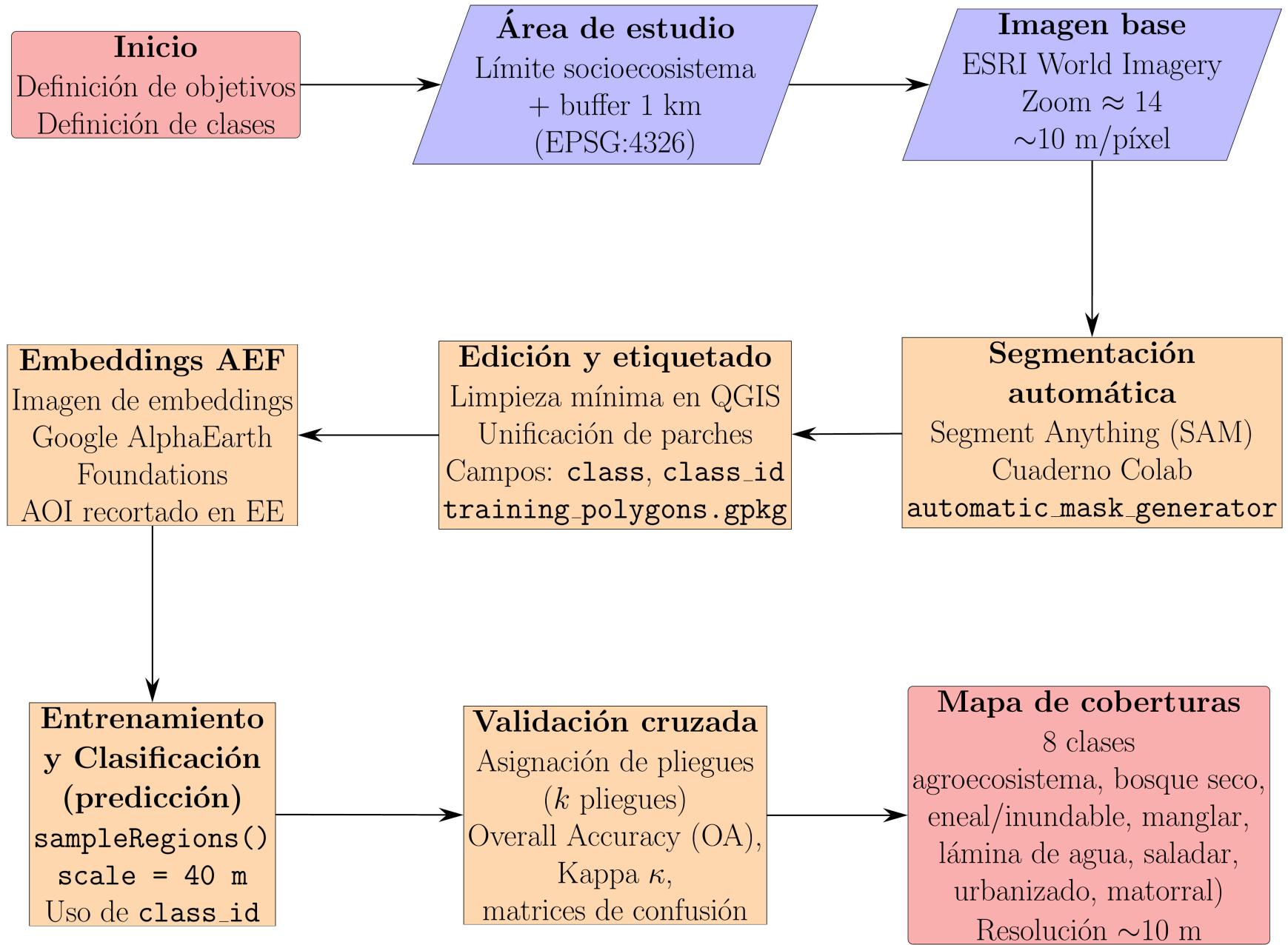
- **Agroecosistemas**
- **Bosque seco / matorral** (clases forestales/arbustivas)
- **Humedales** (varias)
- **Áreas urbanizadas**

Flujo metodológico general

A continuación se muestra el resumen metodológico seguido para este objetivo, basado en tres componentes:

1. **Obtención de las máscaras de entrenamiento** mediante SAM (Segment Anything) usando el cuaderno adaptado de *Qiusheng Wu*, paquetes samgeo, geemap, en Google Colab.
2. **Curación y unificación de clases** en QGIS.
3. **Clasificación supervisada con embeddings AEF** mediante Random Forest y validación cruzada *k-folds* en Google Earth Engine.

Metodología





Qiusheng Wu

giswqs

Follow

♥ Sponsor

Associate Professor at the
University of Tennessee, Knoxville
| Amazon Scholar

8 7.2k followers · 178 following

University of Tennessee

Knoxville, TN

01:18 - 1h behind

✉ qwu18@utk.edu

giswqs / README.md

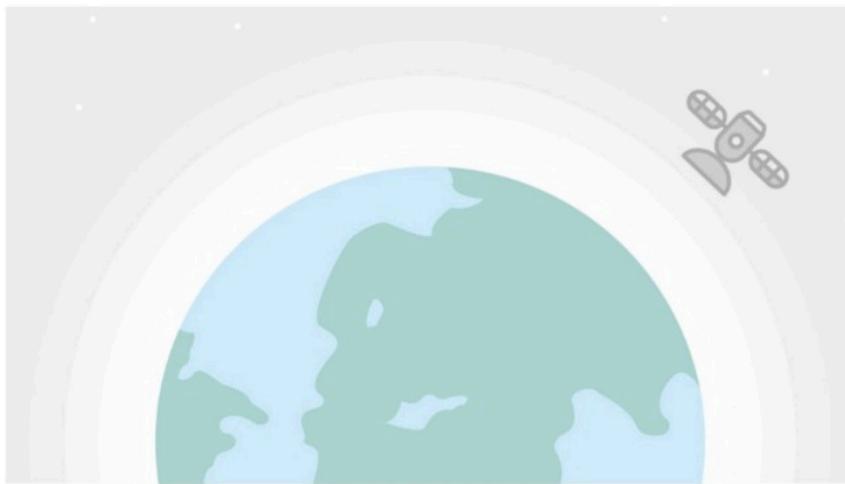
Qiusheng Wu

Followers 7.2k | Follow @giswqs | Google Scholar | UTK Faculty
My YouTube | My LinkedIn | My CV | Donate | Buy me a coffee | visitors 171262

Dr. Qiusheng Wu is an Associate Professor and the Director of Graduate Studies in the Department of Geography & Sustainability at the University of Tennessee, Knoxville. He also serves as an Amazon Scholar. Dr. Wu's research focuses on geospatial data science and open-source software development, with an emphasis on leveraging big geospatial data and cloud computing to study environmental change, particularly surface water and wetland inundation dynamics. He is the creator of several widely used open-source Python packages, including geemap, leafmap, segment-geospatial, and geoui, which support advanced geospatial analysis and interactive visualization. His open-source work is available at <https://github.com/opengeos>.

Open-source Projects

- Linux: [manjaro-linux](#)
- R packages: [whiteboxR](#)
- Python packages: [geemap](#) | [leafmap](#) | [eefolium](#) | [geehydro](#) | [lidar](#) | [whitebox](#) | [whiteboxgui](#) | [geospatial](#) | [pygis](#) | [pypackage](#)
- ArcGIS Toolboxes: [WhiteboxTools-ArcGIS](#) | [Depression Analysis Toolbox](#) | [Wetland Hydrology Analyst](#)
- Google Earth Engine: [Awesome-GEE](#) | [earthengine-py-notebooks](#) | [qgis-earthengine-examples](#) | [earthengine-apps](#)



Welcome to Google Earth Engine

1. Obtención de imágenes de alta resolución (ESRI World Imagery)

Se descargó una imagen de referencia del área seleccionada empleando *ESRI World Imagery* con resolución aproximada de **5–10 m**, de manera que fuera coherente con la escala representada en los embeddings AEF.

2. Segmentación automática con Segment Anything (SAM)

La segmentación se realizó usando cuaderno Jupyter. El procesamiento incluyó:

- Uso del modelo **ViT-H** (SAM), adecuado para imágenes de alta resolución pero útil también para imágenes de mediana resolución. Eficiente por su relación sensibilidad–costo computacional.
- Aplicación de `sam.generate()` con `unique=True` para obtener **instancias diferenciadas**.
- Ajuste de hiperparámetros (`points_per_side`, `pred_iou_thresh`, `crop_n_layers`, etc.) para aumentar la densidad de máscaras.

Automatic mask generation options

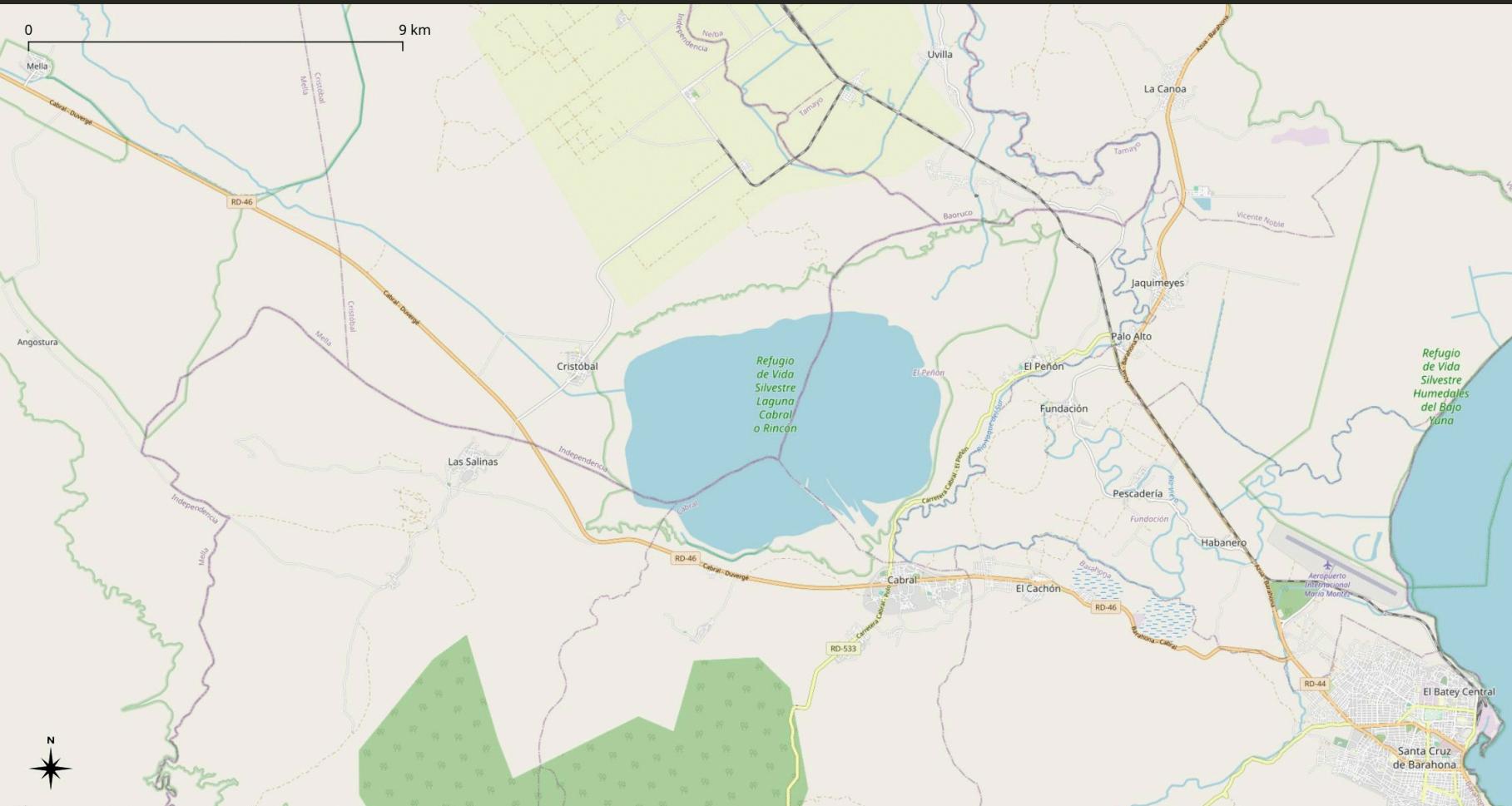
There are several tunable parameters in automatic mask generation that control how densely points are sampled and what the thresholds are for removing low quality or duplicate masks. Additionally, generation can be automatically run on crops of the image to get improved performance on smaller objects, and post-processing can remove stray pixels and holes. Here is an example configuration that samples more masks:

```
# sam_kwargs = {
#     "points_per_side": 32,
#     "pred_iou_thresh": 0.86,
#     "stability_score_thresh": 0.92,
#     "crop_n_layers": 1,
#     "crop_n_points_downscale_factor": 2,
#     "min_mask_region_area": 100,
# }

# Para zoom 14
sam_kwargs = {
    "points_per_side": 24,           # más que 16, menos que 32 (menos VRAM que 32)
    "pred_iou_thresh": 0.5,         # bajar de 0.86 -> acepta más máscaras
    "stability_score_thresh": 0.8,  # bajar de 0.92 -> idem
    "crop_n_layers": 1,            # mantiene una capa extra de crops
    "crop_n_points_downscale_factor": 2,
    "min_mask_region_area": 25,    # baja de 100 -> deja pasar más regiones pequeñas
}

sam = SamGeo(
    model_type="vit_h",
    sam_kwargs=sam_kwargs,
)

sam.generate(image, output="masks2.tif", foreground=False)
```



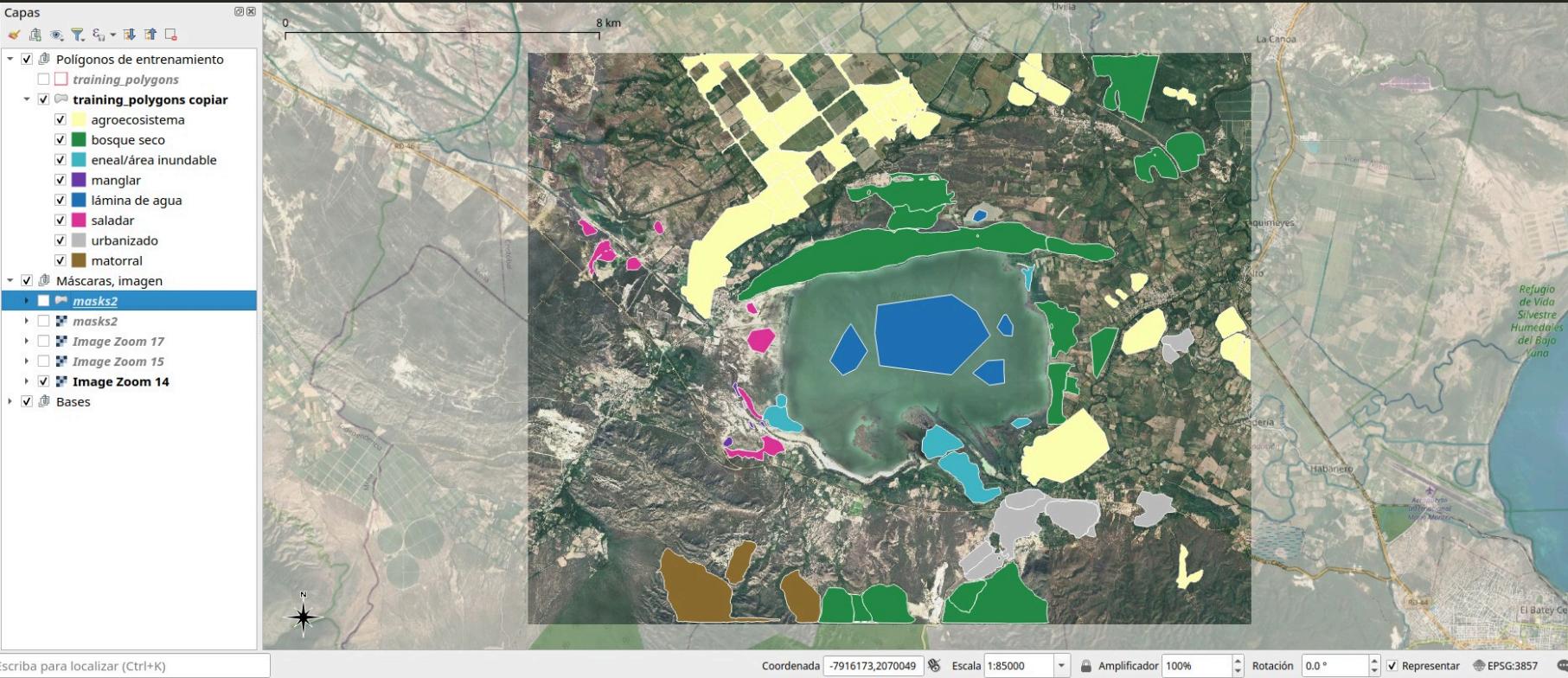
3. Limpieza y edición mínima en QGIS

Las máscaras obtenidas se exportaron a GeoPackage, donde se realizó una edición mínima:

- Corrección de pequeñas islas y polígonos erróneos.
- Eliminación de duplicados.
- Fusión de parches contiguos cuando correspondían a una misma cobertura.

Se añadieron dos campos:

- `class_id` (entero 1–8)
- `class_label` (nombre de la clase)



4. Clasificación supervisada con embeddings AEF en Google Earth Engine

Pasos principales

- Se definió un **AOI** a partir del *buffer de 1 km* del socioecosistema.
- Se extrajeron los embeddings mediante:

```
emb_image = ee.Image("GOOGLE/AE/IMG_V1").clip(aoi)
```

- Se muestraron los píxeles según los polígonos curados:

```
sample = emb_image.sampleRegions(  
    collection=fc,  
    properties=["class_id"],  
    scale=40,    # ajustado para evitar EEEexception: memory exceed  
    geometries=False  
)
```

- Se añadieron **5 folds** para validación cruzada:

```
sample = sample.randomColumn("rand").divide(k).int()
```

- Se entrenó un clasificador:

```
classifier = ee.Classifier.smileRandomForest(  
    number_of_trees=100,  
    min_leaf_population=5  
)
```

5. Resultados de la validación cruzada (5-fold CV)

Los resultados fueron evaluados en términos de:

- **Kappa de Cohen**
- **Exactitud global (OA)**
- **Matriz de confusión por pliegue**

Ejecutar k-fold CV

```
# Alternativa 1: usando list comprehension
results = [train_and_evaluate_fold(i) for i in range(k)]

for r in results:
    print(f"Fold {r['fold']}: "
          f"Train OA={r['train_accuracy']:.3f}, κ={r['train_kappa']:.3f} | "
          f"Val OA={r['val_accuracy']:.3f}, κ={r['val_kappa']:.3f}")
```

```
Fold 0: Train OA=0.986, κ=0.981 | Val OA=0.971, κ=0.961
Fold 1: Train OA=0.986, κ=0.981 | Val OA=0.972, κ=0.962
Fold 2: Train OA=0.986, κ=0.981 | Val OA=0.970, κ=0.960
Fold 3: Train OA=0.985, κ=0.980 | Val OA=0.972, κ=0.963
Fold 4: Train OA=0.986, κ=0.981 | Val OA=0.975, κ=0.967
```

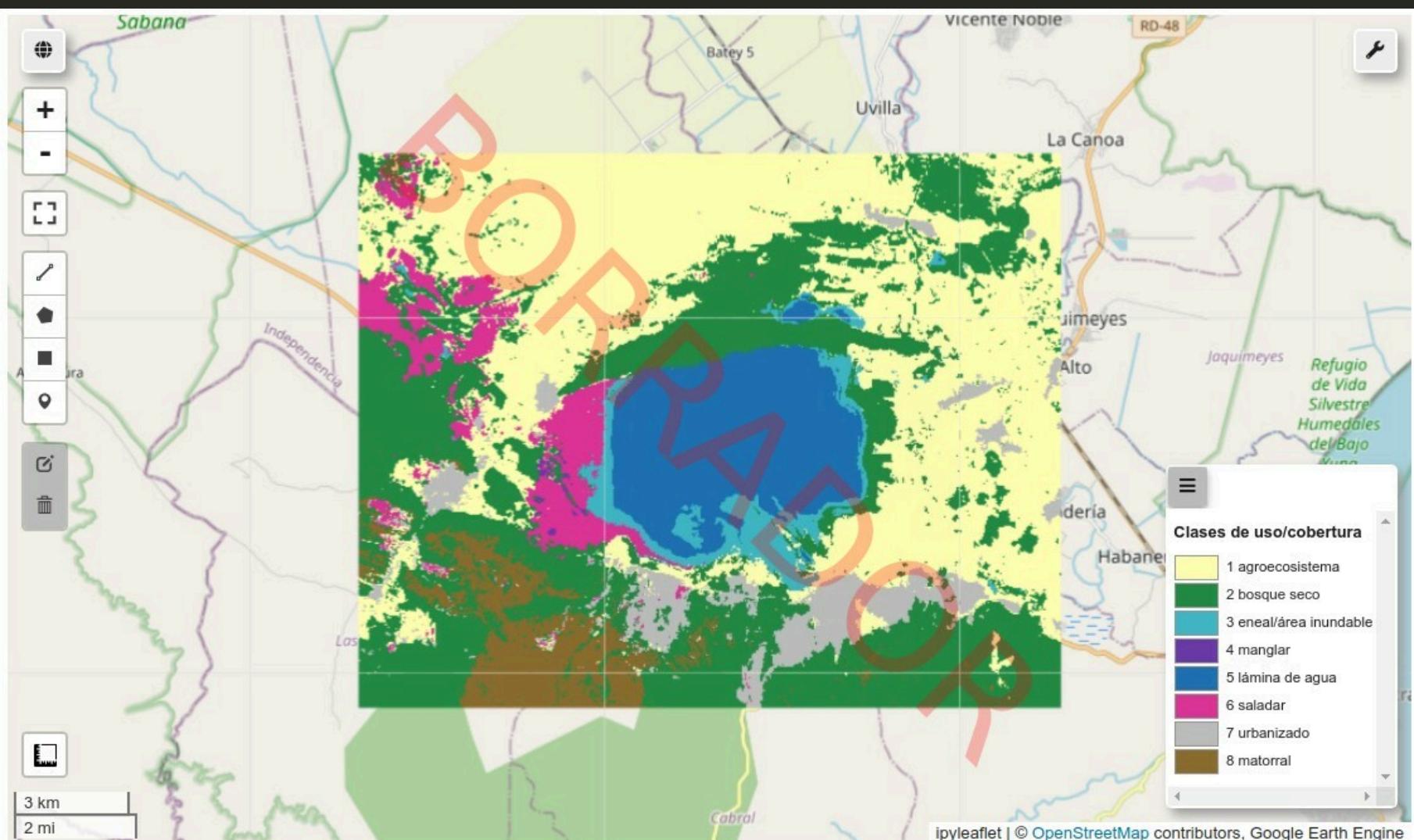
Métricas promedio de validación

```
val_accs = [r["val_accuracy"] for r in results]
val_kappas = [r["val_kappa"] for r in results]

print("Promedio OA validación:", np.mean(val_accs))
print("Desv. estándar OA:", np.std(val_accs))
print("Promedio κ validación:", np.mean(val_kappas))
print("Desv. estándar κ:", np.std(val_kappas))
```

```
Promedio OA validación: 0.9718087715864383
Desv. estándar OA: 0.001769670358446398
Promedio κ validación: 0.9623802751491193
Desv. estándar κ: 0.002317866532848917
```

Mapa de clasificación



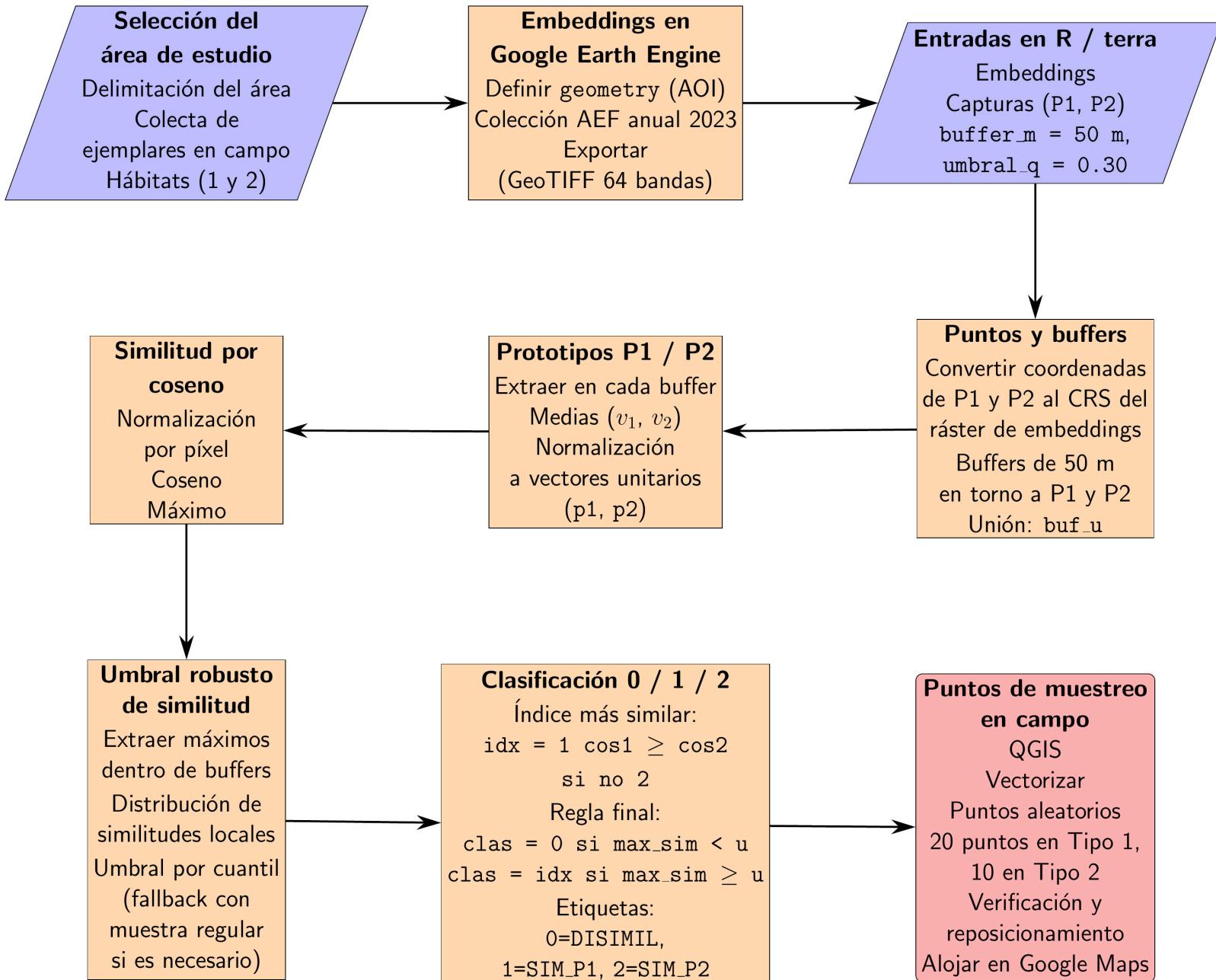
Conclusiones del Objetivo 2

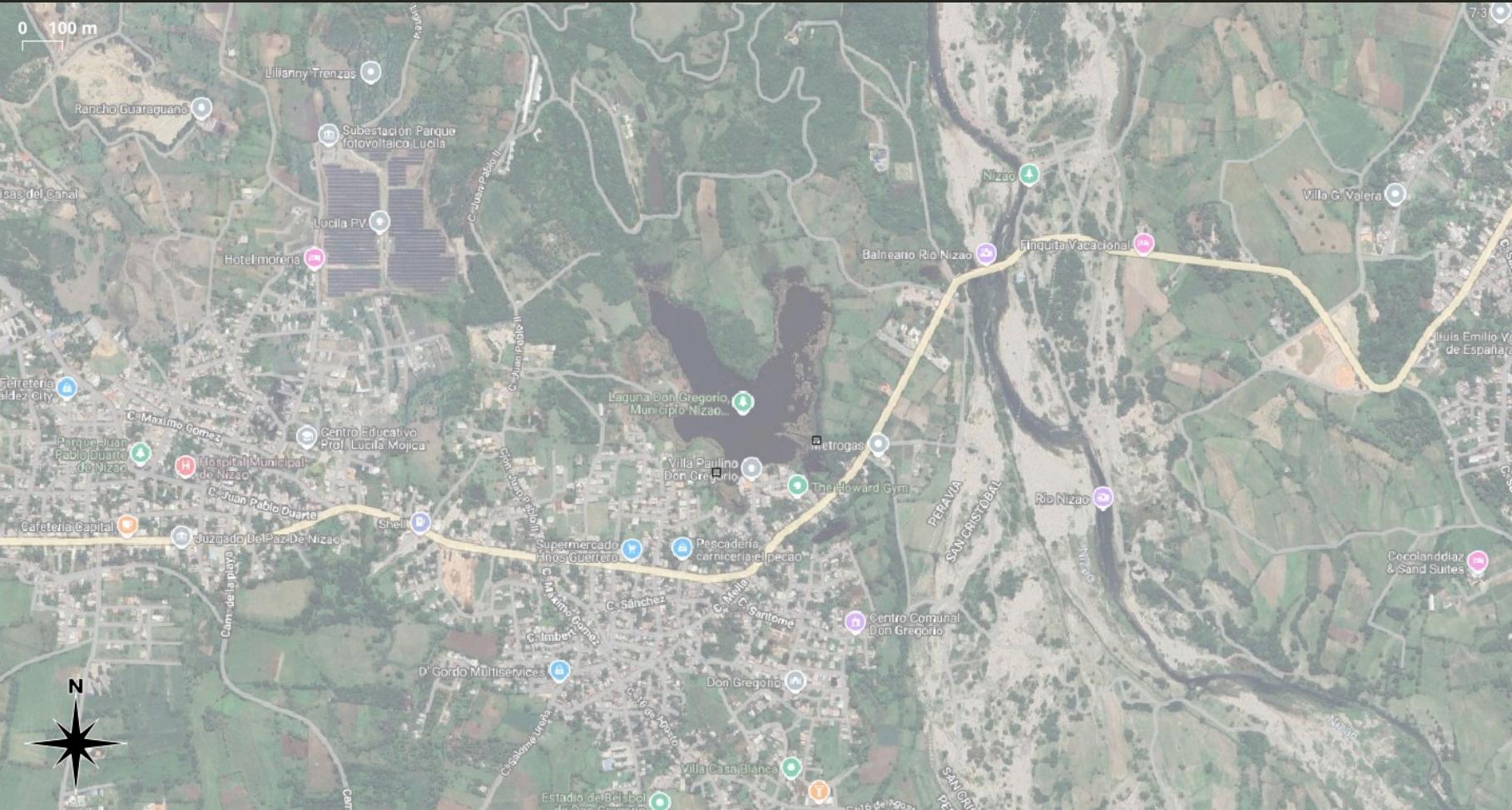
- Los embeddings de AEF **sí capturan diferencias estructurales** entre agroecosistemas, áreas forestales/arbustivas y zonas urbanas.
- Con solo unos pocos polígonos de entrenamiento (segmentados vía SAM), se logró una clasificación robusta.
- El uso combinado de **SAM + QGIS + AEF embeddings + Random Forest** constituye un flujo de trabajo reproducible y eficiente para análisis de cobertura sin necesidad de deep learning complejo.
- La validación cruzada mostró niveles consistentes de desempeño entre pliegues, indicando estabilidad del clasificador.

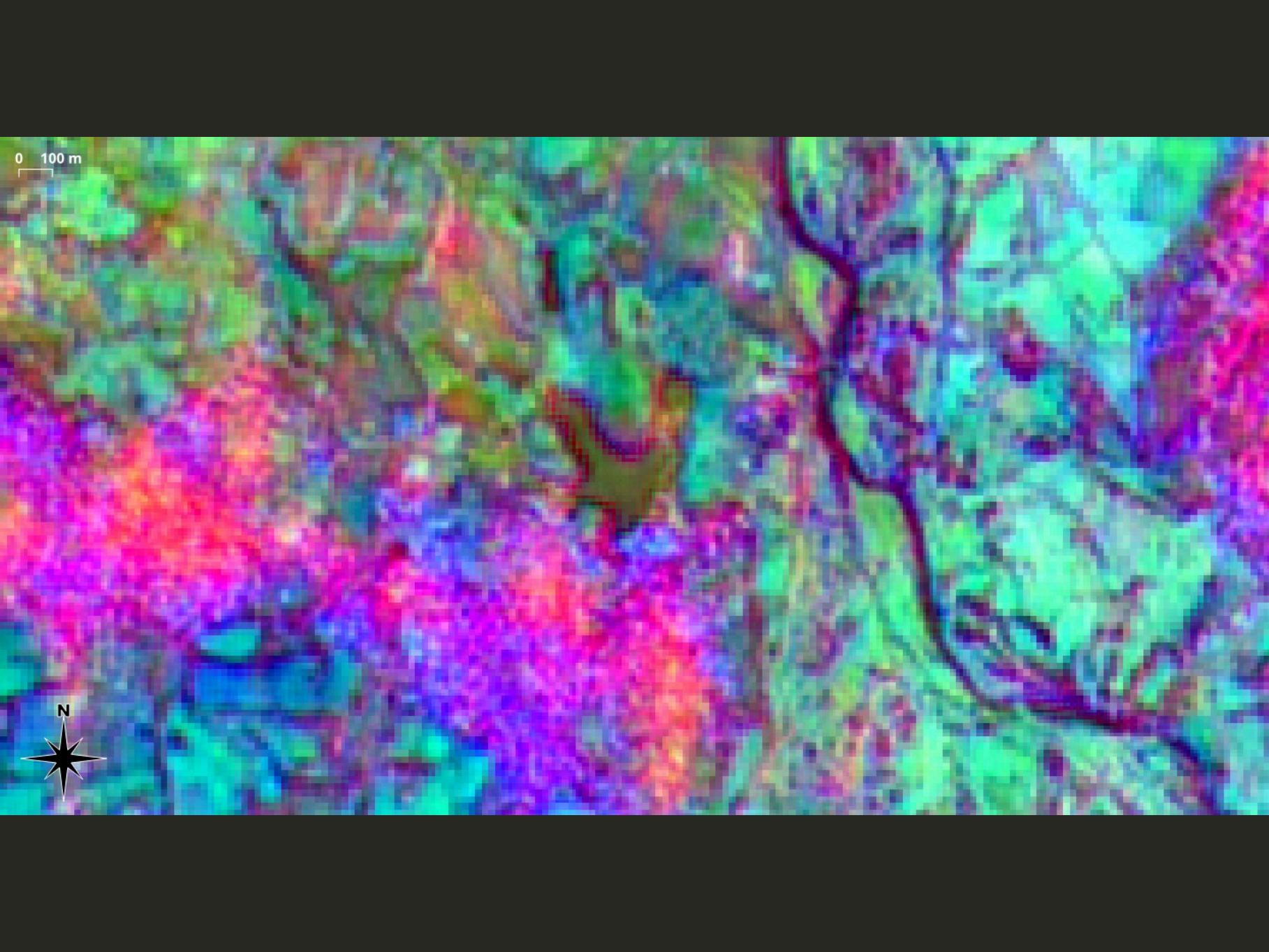
Objetivo 3

Analizar la aplicación de embeddings de AEF en el **diseño de muestreo estratificado**, utilizando la similitud entre embeddings para optimizar la representatividad espacial y ambiental

Colaboración para María Fernanda Rodríguez y Wellin Brito.

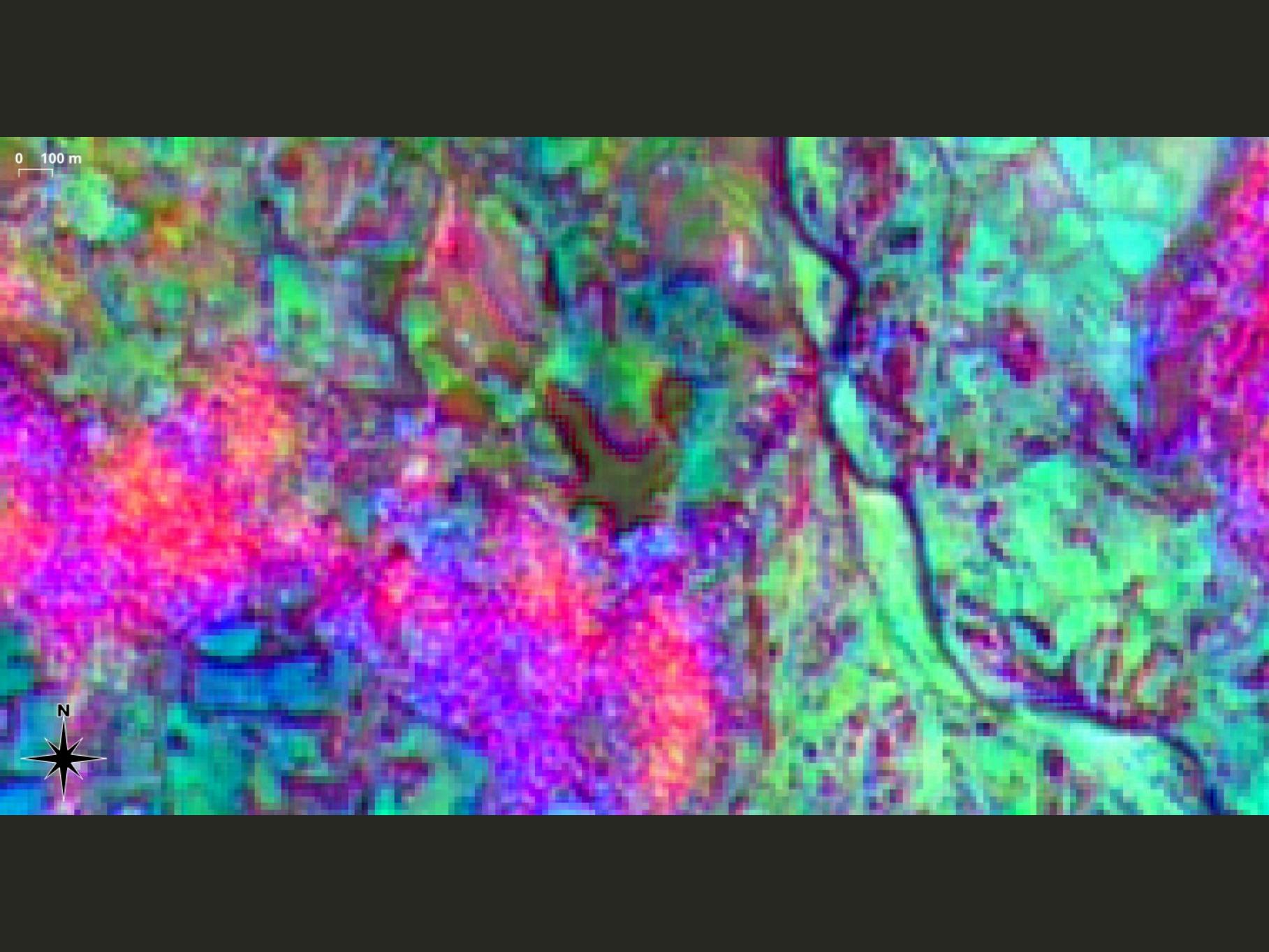


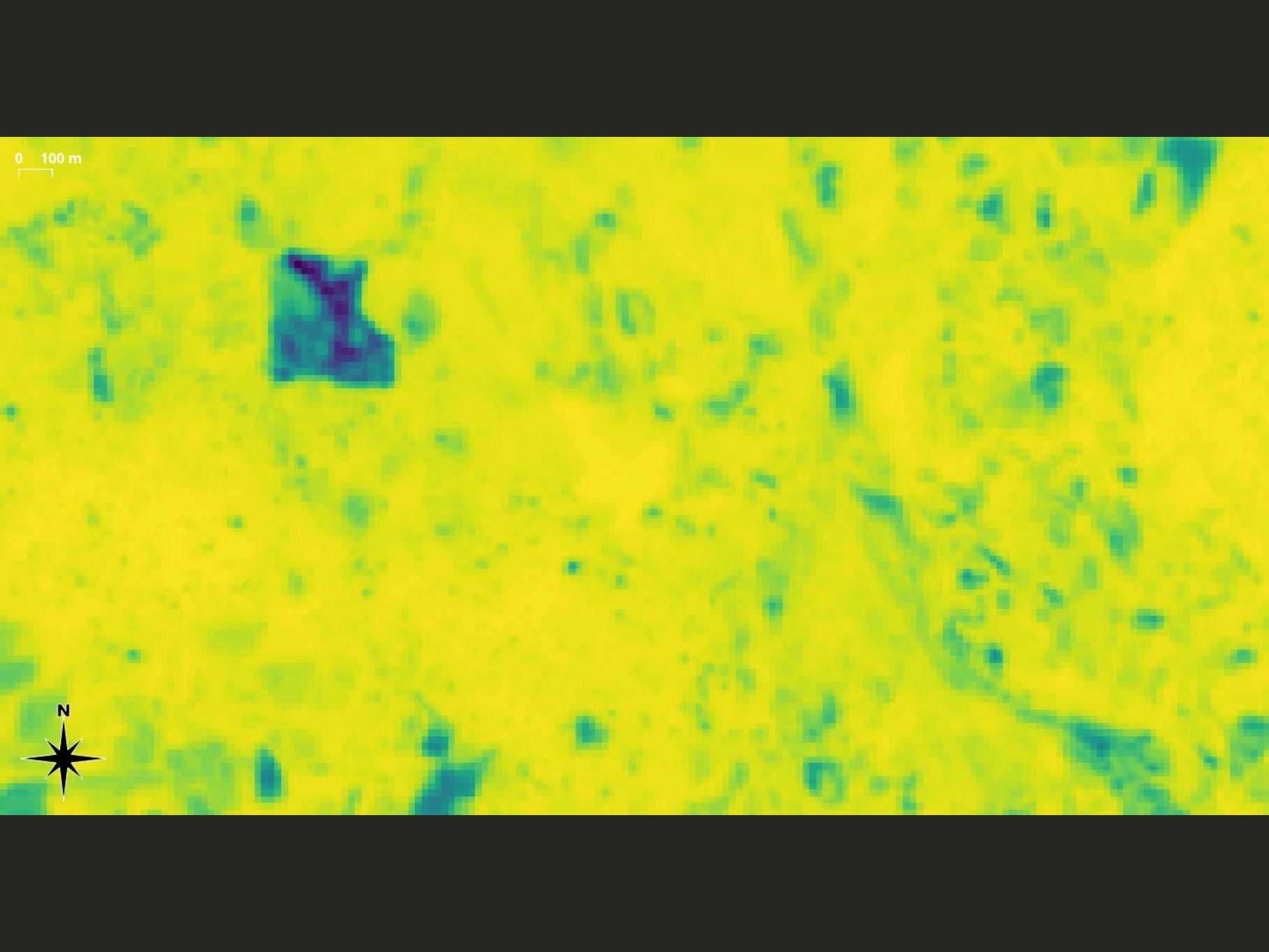


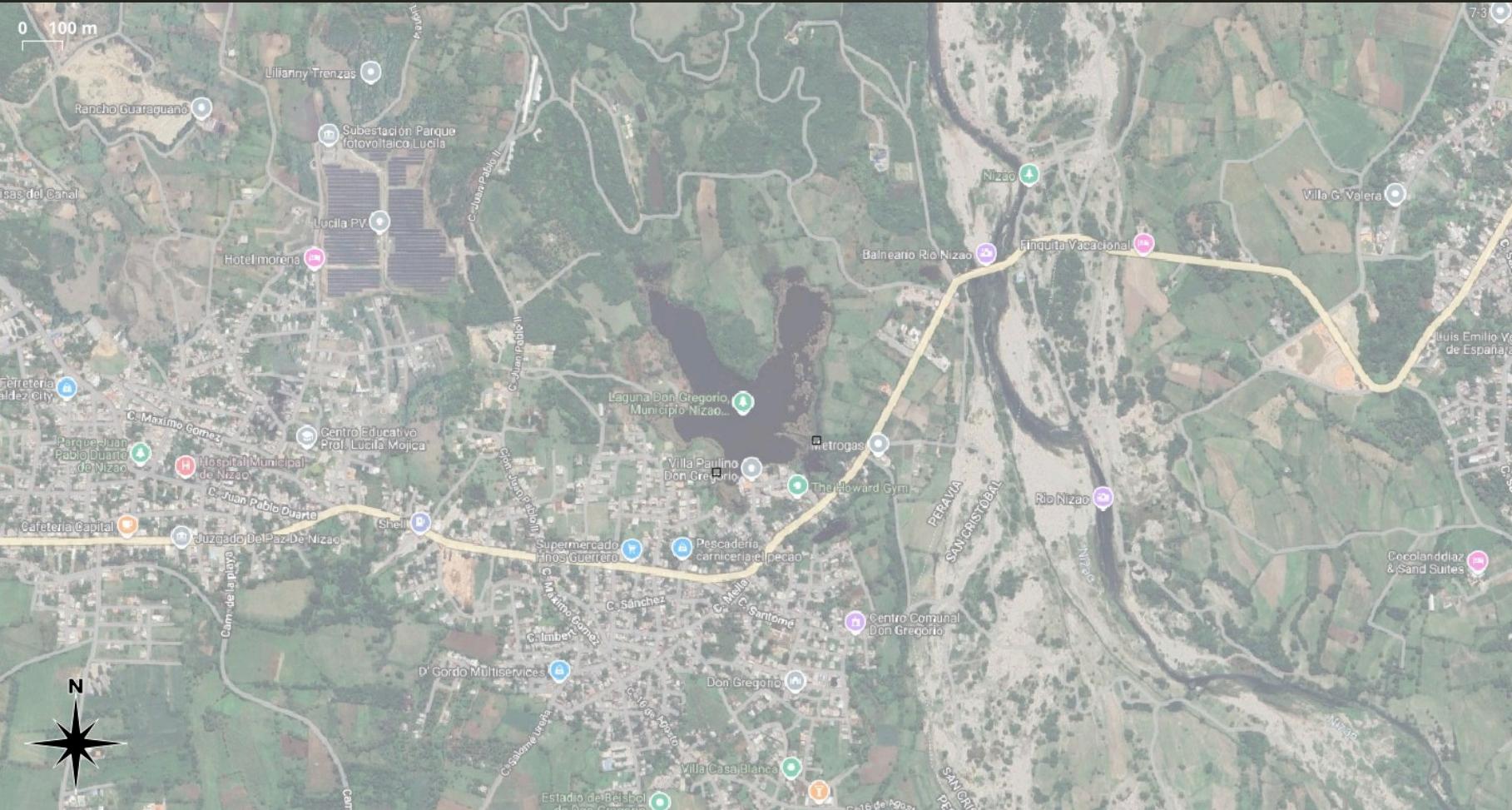


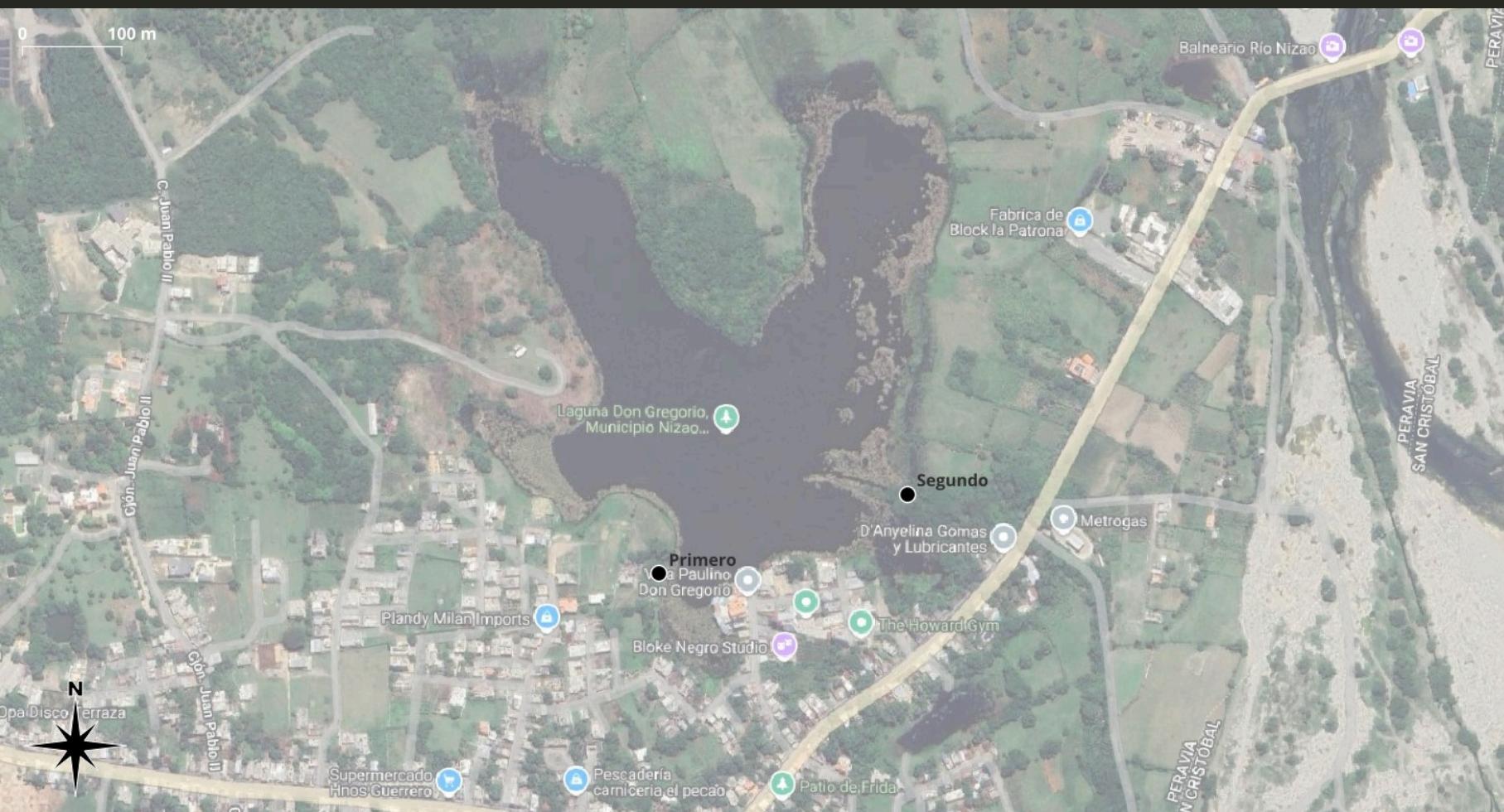
0 100 m



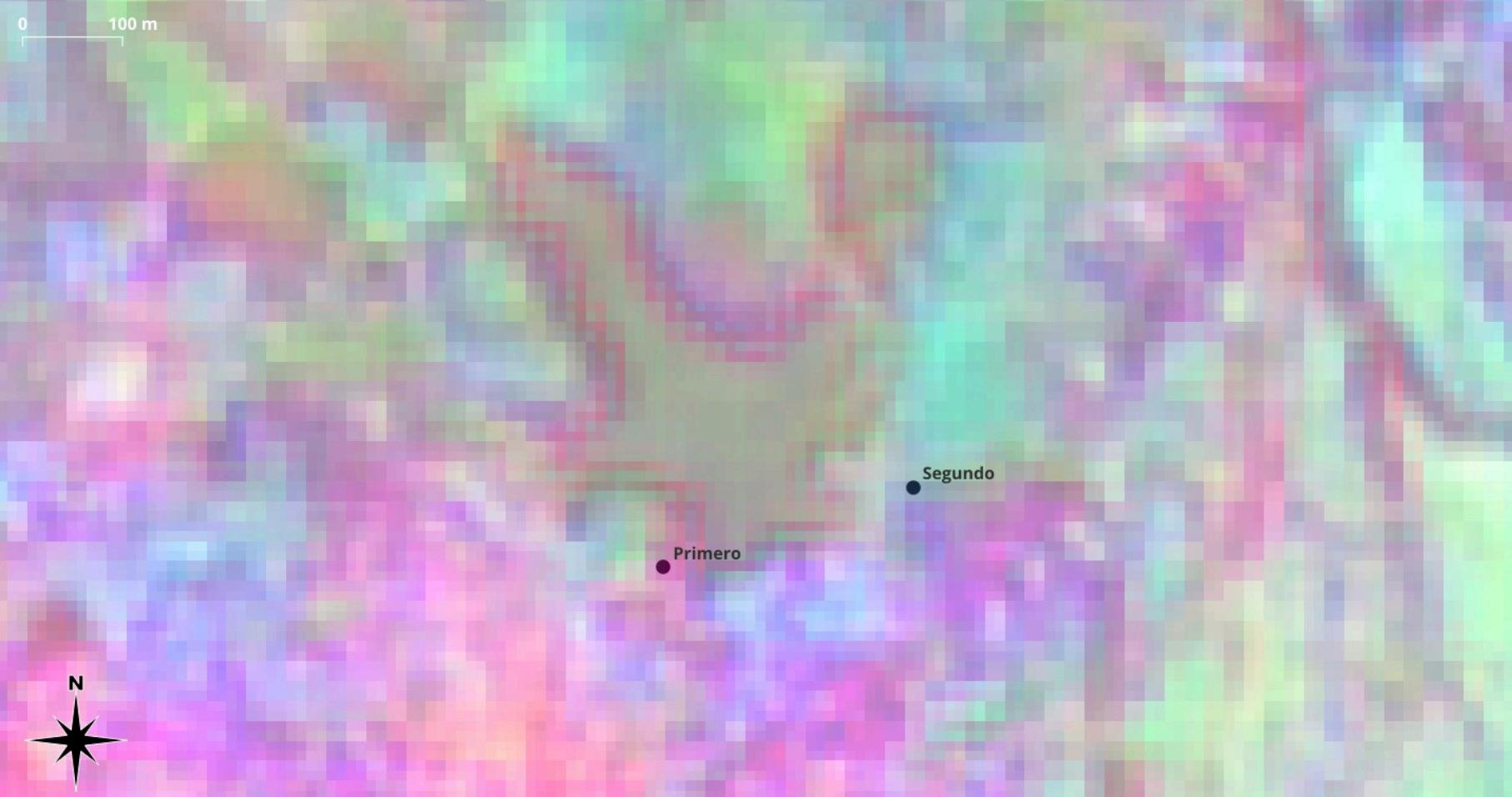


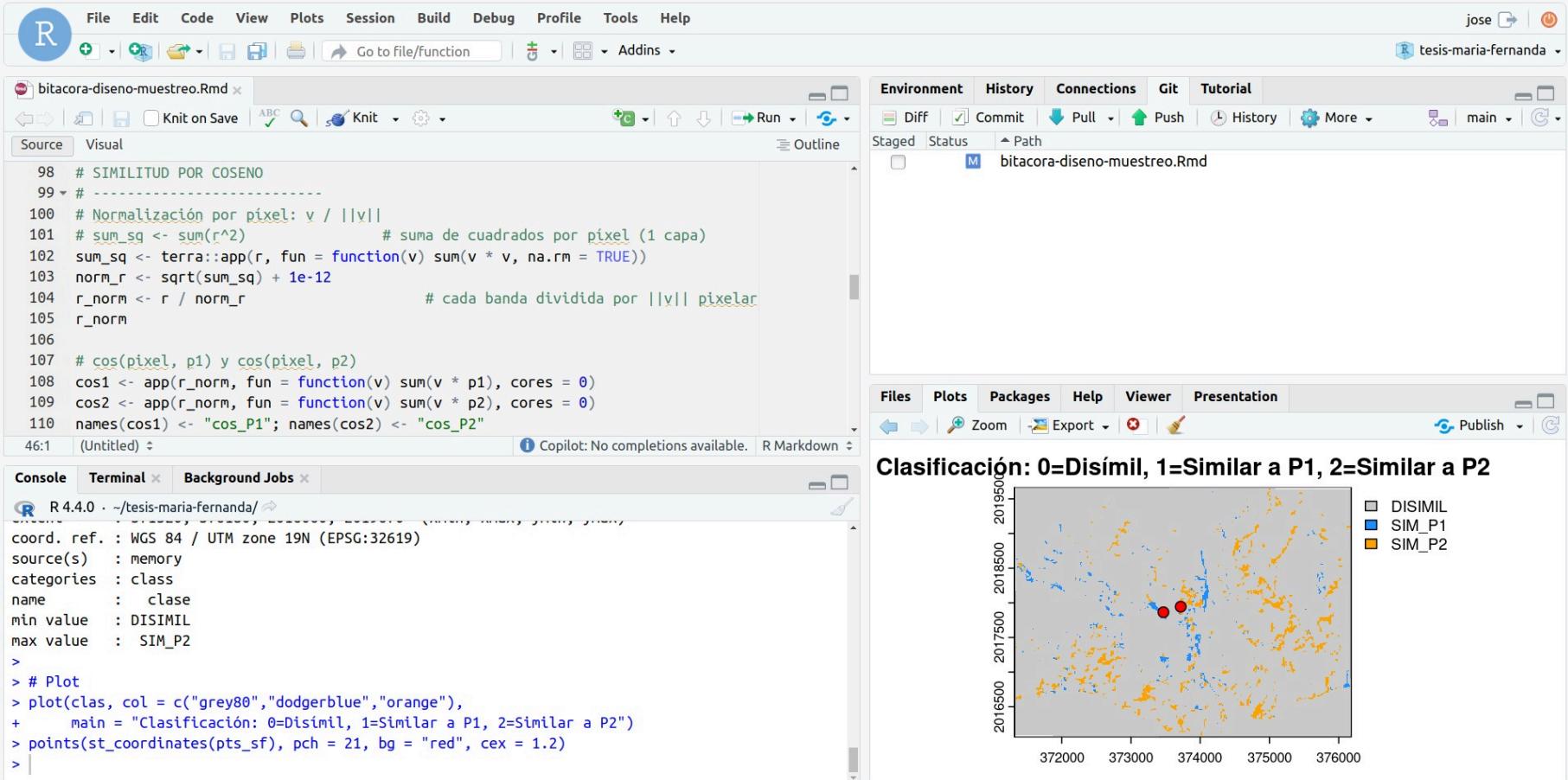


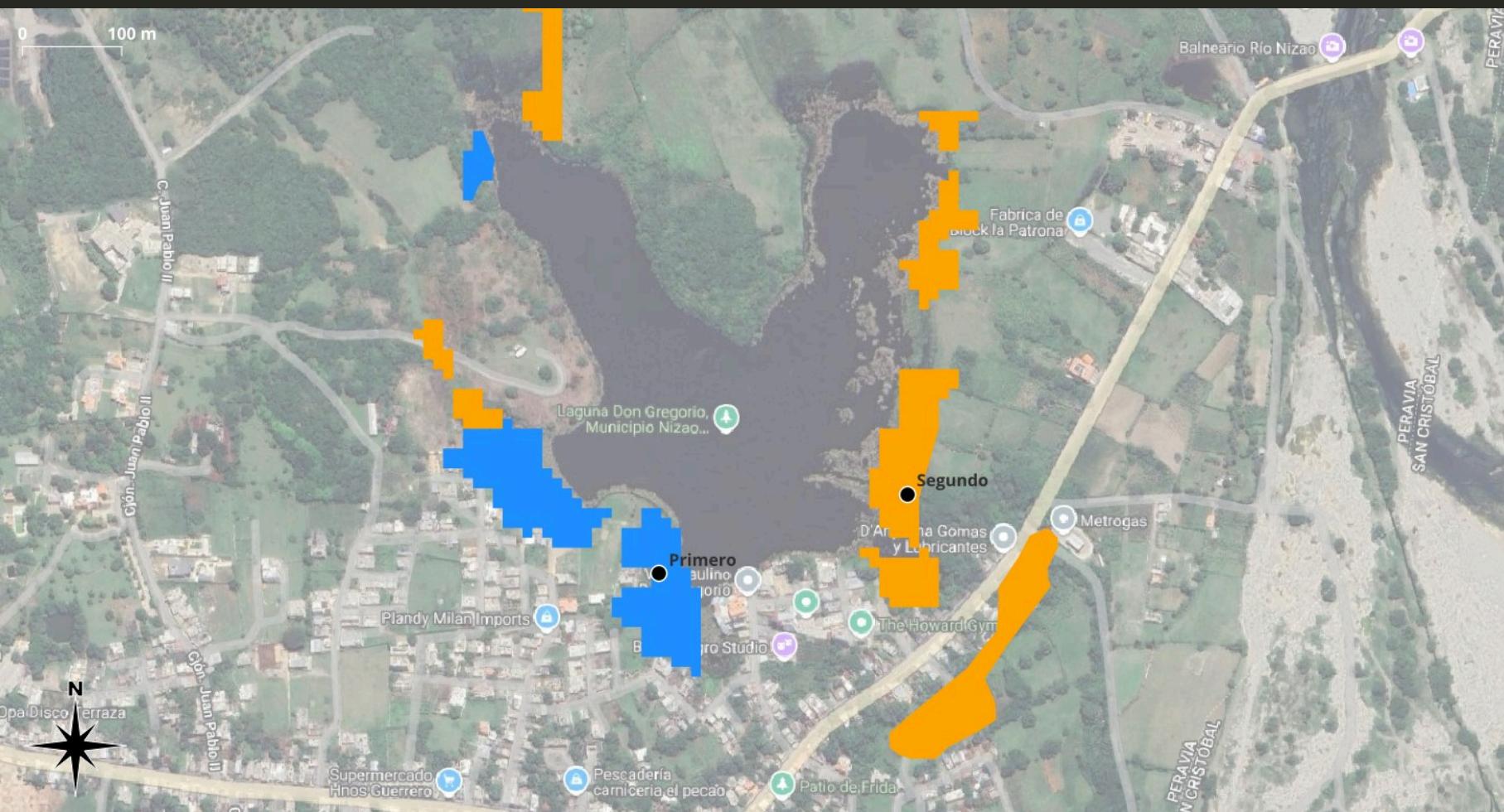




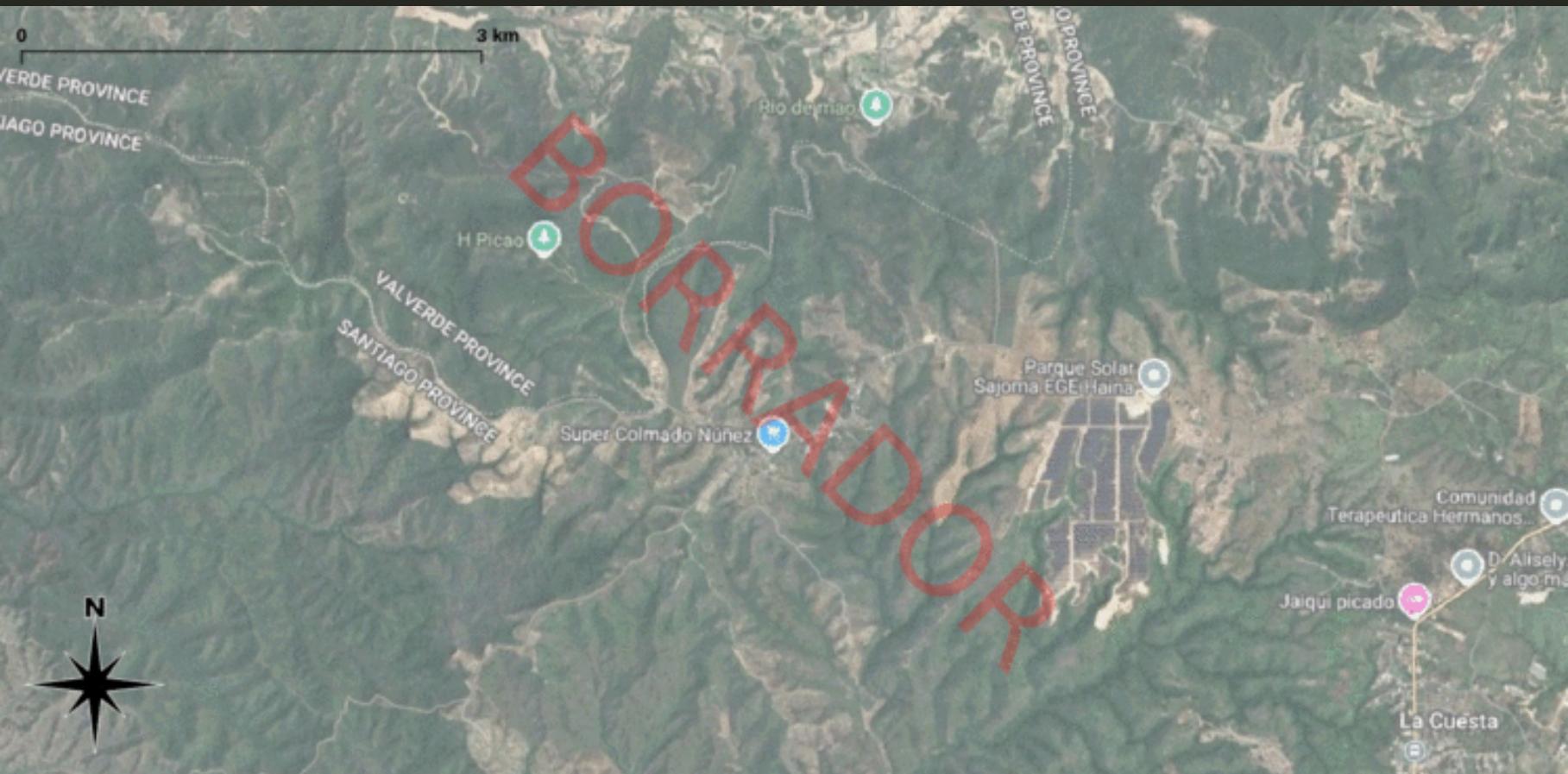












Conclusiones del Objetivo 3

- **Los embeddings AEF permiten definir estratos ambientales** sin mapas previos, usando solo similitud por coseno.
- **El umbral aprendido localmente genera estratos estables**, coherentes con los hábitats asociados a los puntos de terreno.
- **El muestreo estratificado resultante es más representativo**, asignando puntos según la estructura ambiental real del área.

Gracias por su atención



jmartinez19@uasd.edu.do



geofis