

# Climate TRACE: Estimating global greenhouse gas emissions from buildings

MIDS Capstone 2025  
Duke University



# Team



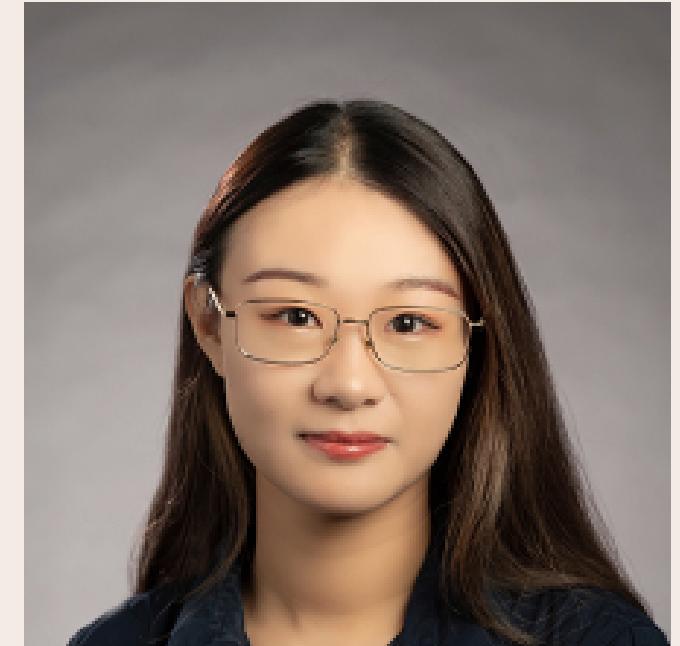
Bárbara Flores



Meixiang Du



Jiechen Li



Yulei (Alicia) Xia

Mentor & Client: Kyle Bradbury, Ph.D.

Duke University





# Problem Statements

In 2024, global carbon dioxide (CO<sub>2</sub>) emissions reached a record 41.6 billion metric tons.<sup>[1]</sup>

Emissions inventories remain limited—52 countries lack data after 2012, with even larger gaps for subnational and local governments.<sup>[2] [3]</sup>

Existing data is often outdated, spatially limited, or not free.

Available emission data is at a coarse spatial resolution.

[1]. Source: <https://www.livescience.com/planet-earth/climate-change/global-carbon-emissions-reach-new-record-high-in-2024-with-no-end-in-sight-scientists-say>

[2]. Source: <https://climatetrace.org/news/more-than-70000-of-the-highest-emitting-greenhouse-gas>

[3]. Source: Luers, A., Yona, L., Field, C. B., Jackson, R. B., Mach, K. J., Cashore, B. W., Elliott, C., Gifford, L., Honigsberg, C., Klaassen, L., & Matthews, H. D. (2022).

Make greenhouse-gas accounting reliable—Build interoperable systems. *Nature*, 607(7920), 653–656.

# Problem Statements

In 2024, global carbon dioxide (CO<sub>2</sub>) emissions reached a record 41.6 billion metric tons.<sup>[1]</sup>

Emissions inventories remain limited—52 countries lack data after 2012, with even larger gaps for subnational and local governments.<sup>[2] [3]</sup>

Existing data is often outdated, spatially limited, or not free.

Available emission data is at a coarse spatial resolution.

# Climate TRACE

Climate Trace is a global non-profit coalition.

They estimate and track global GHG emissions since 2015.

They provide fast, easy, and free data access to the public.

They provide comprehensive granular data across the sectors.

[1]. Source: <https://www.livescience.com/planet-earth/climate-change/global-carbon-emissions-reach-new-record-high-in-2024-with-no-end-in-sight-scientists-say>

[2]. Source: <https://climatetrace.org/news/more-than-70000-of-the-highest-emitting-greenhouse-gas>

[3]. Source: Luers, A., Yona, L., Field, C. B., Jackson, R. B., Mach, K. J., Cashore, B. W., Elliott, C., Gifford, L., Honigsberg, C., Klaassen, L., & Matthews, H. D. (2022).

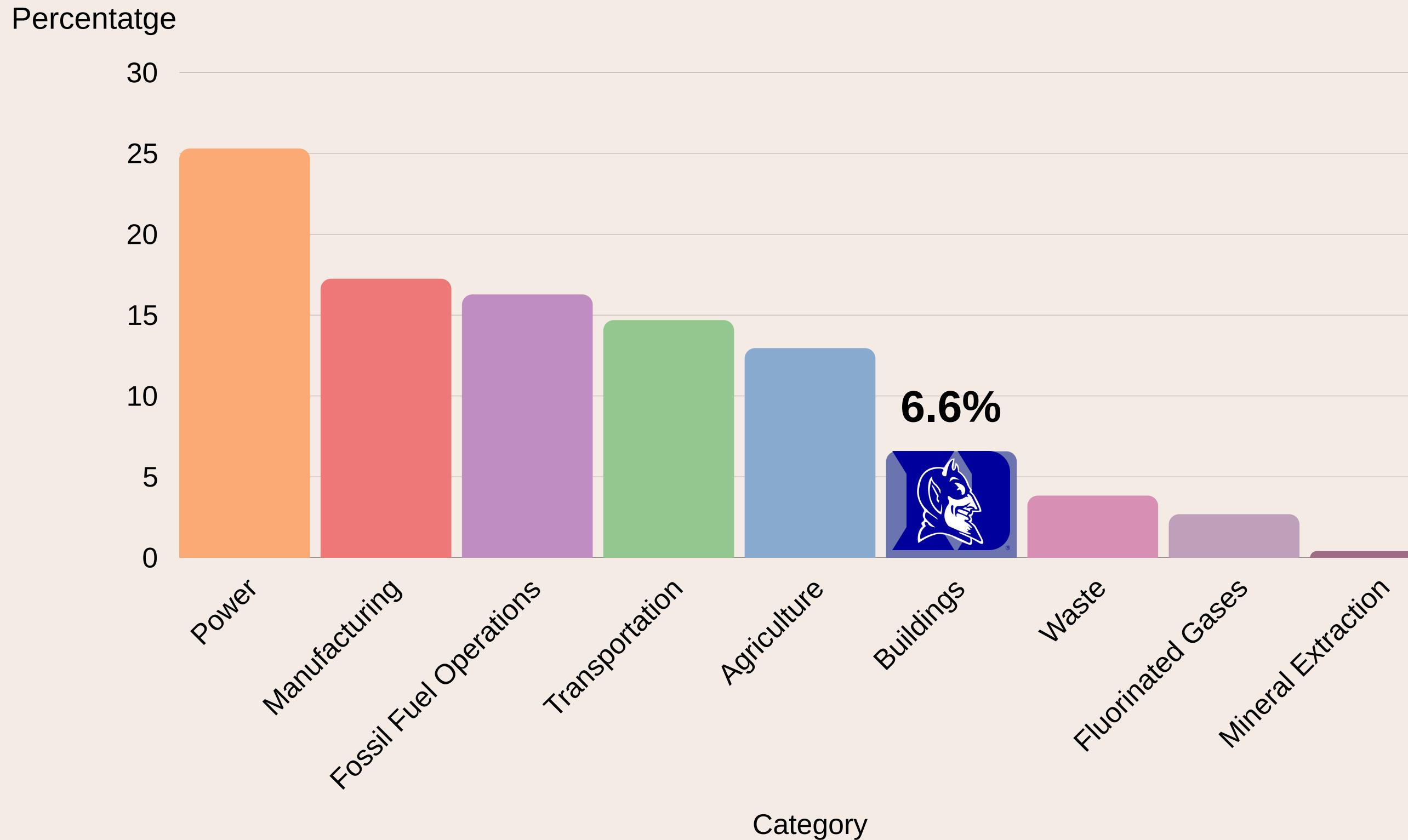
Make greenhouse-gas accounting reliable—Build interoperable systems. *Nature*, 607(7920), 653–656.

# Climate TRACE



## Global Greenhouse Gas Emissions by Source

Source: <https://climatetrace.org/sectors>



We focus on the **building** sector, which accounts for **6.6%** of total emissions.

# How to Calculate GHG

$$A \times EUI \times EF = GHG$$

Estimate from direct observation through satellite imagery

**Building Floor Area**

Gathered from statistical assessments and adjusted using machine learning

**Energy Use Intensity**  
measures the energy consumption per square meter of building space

Gathered from statistical assessments

**Equivalent Emissions Factor**

**Residential**

**Non-Residential**

# Dataset



## Data Source

- World Bank's CURB Tool provides comprehensive and reliable data on urban buildings and energy use.



## 482 Entries

- These **city-level** entries represent 482 distinct locations across the globe, serving as the dependent variable in our analysis.



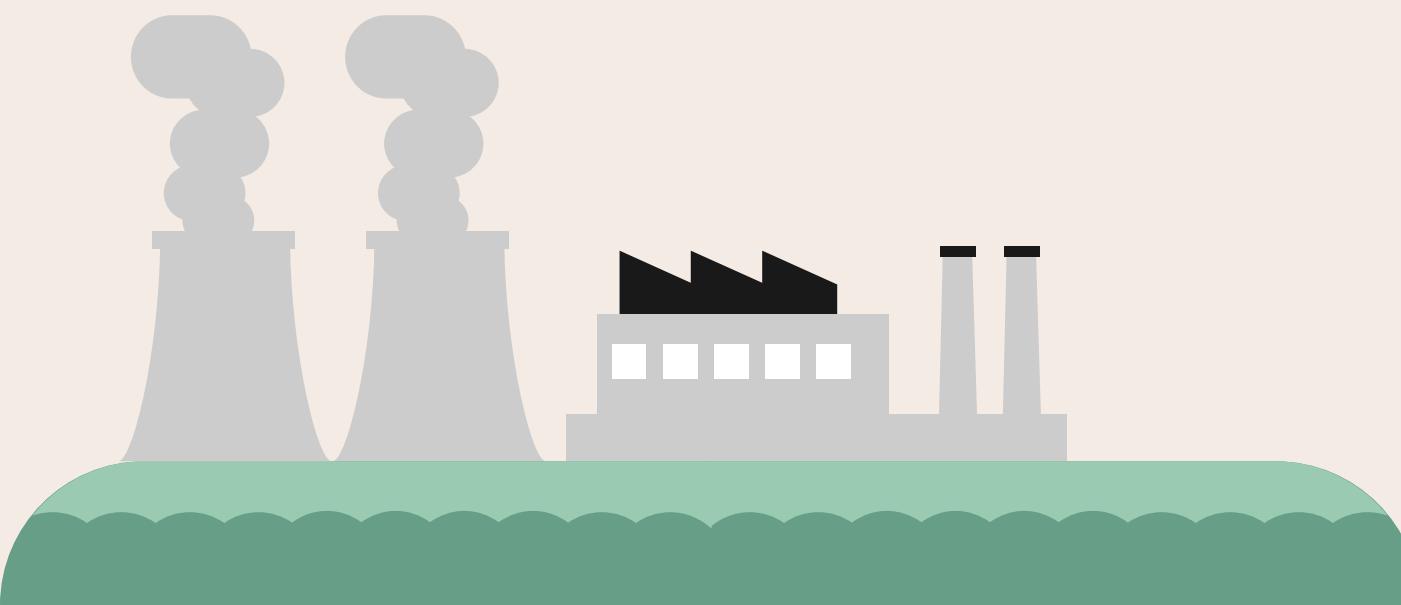
## 1 Entry = 1 Geographical Location

- Each entry represents a geographical location, defined by a specific latitude and longitude, which is treated as a data point in our calculations.

# Goal of the Project

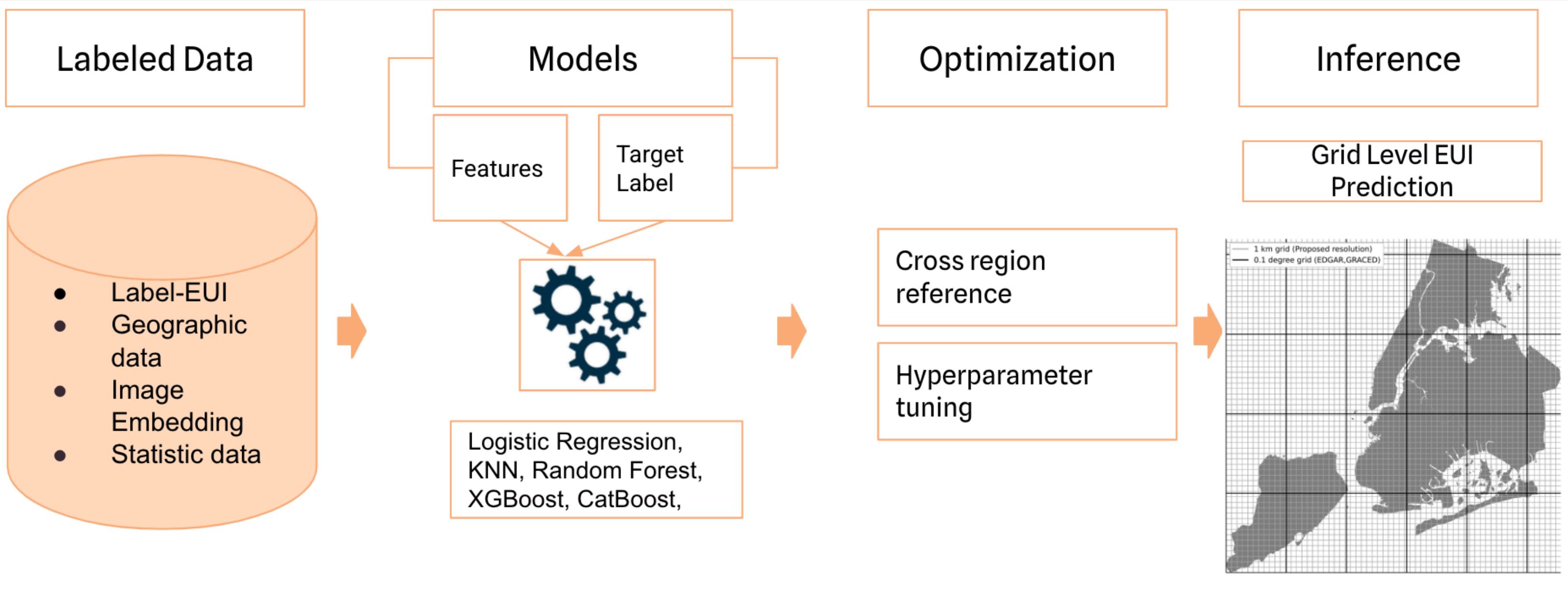
Develop a machine learning model to predict the EUI of buildings using globally available features, enabling the estimation of global greenhouse gas emissions.

Predict EUI



# Methodology

- Employing geographic information techniques, image embedding retrieval methods to process features.
- Use supervised learning models for predicting target variables and conduct cross-region evaluation.



\*\*grid image from paper: High-resolution Global Building Emissions Estimation using Satellite Imagery.

# Identify Variables–Feature Map

From previous study, literature review and further interpretation of the project.

## **Building Geometry Data**

Embeddings from  
Satellite Image  
data(Sentinel-2  
images)

## **Weather**

Temperature,  
DewPoint  
Temperature,  
Heating Degree  
Days(HDD),  
Cooling Degree  
Days(CDD)

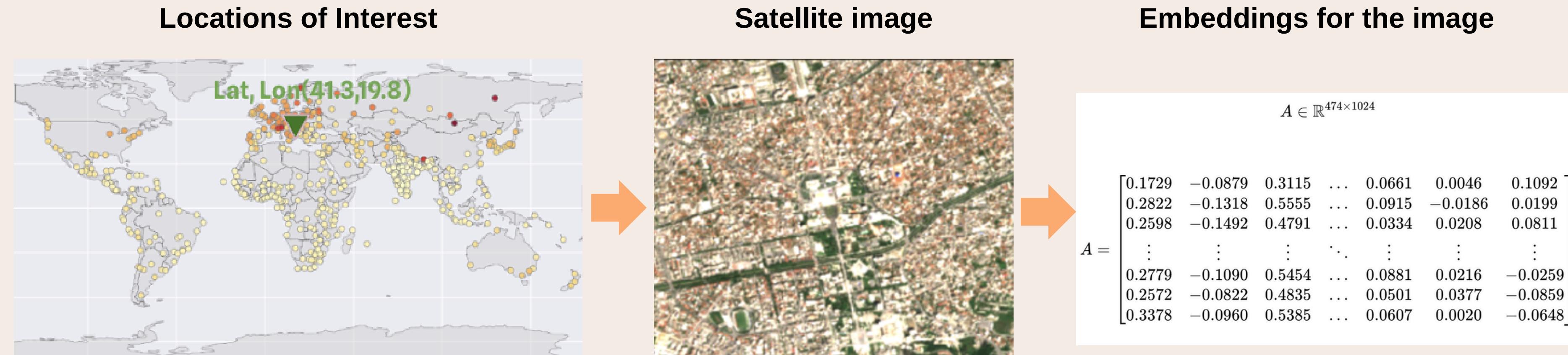
## **Socioeconomics**

Human Development  
Index(HDI),  
GDP,  
Population,  
Educational Index,  
Income Index,  
Urbanization

## **Policy/Law**

Paris Agreement

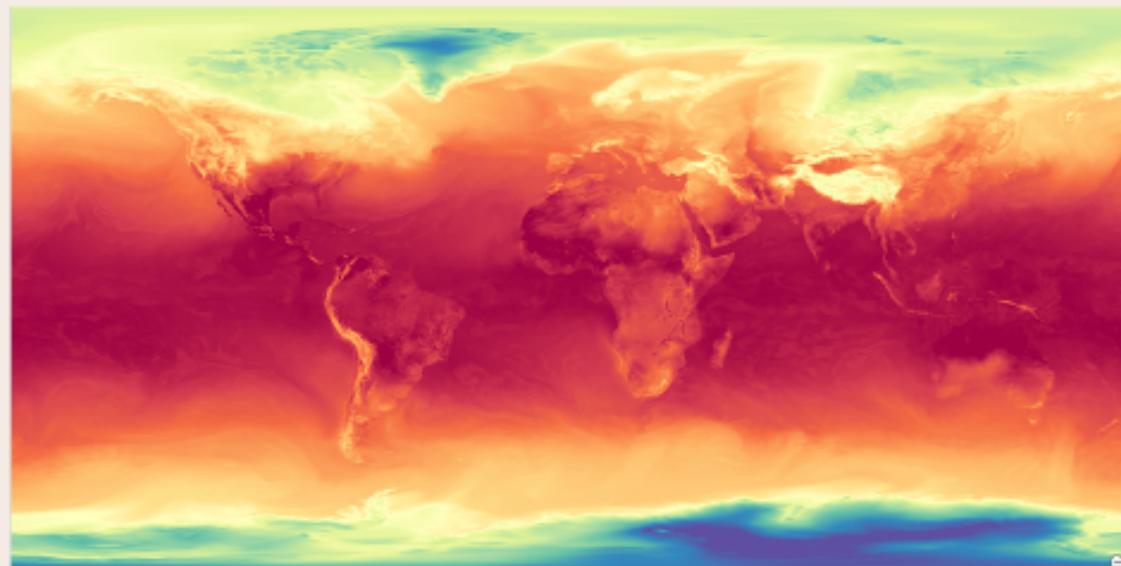
# Represent Variables–Image Embedding



1. Retrieve satellite images from Sentinel-2 for the locations of interest.
2. Employ Clay (a pretrained Vision Transformer MAE model specialized in satellite images) to generate embeddings.
3. Apply dimension reduction and/or classification methods to optimize the representational ability of embeddings.

# Represent Variables–Geographic information

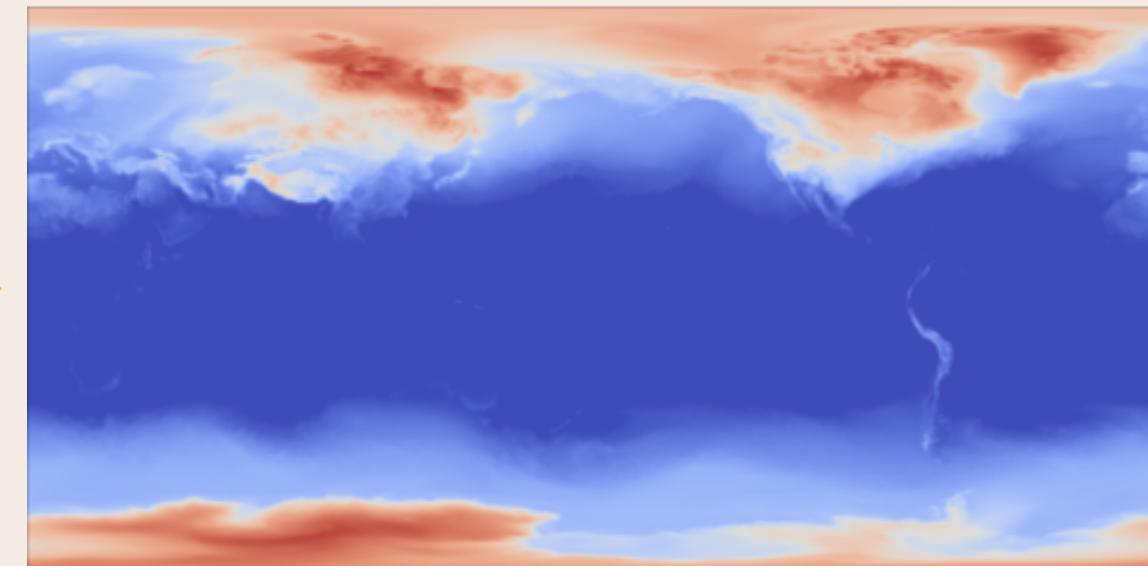
**Retrieve Temperature**



Low High

Retrieve air temperature at 2m above the surface, data source from European Commission and the Group on Earth Observations.

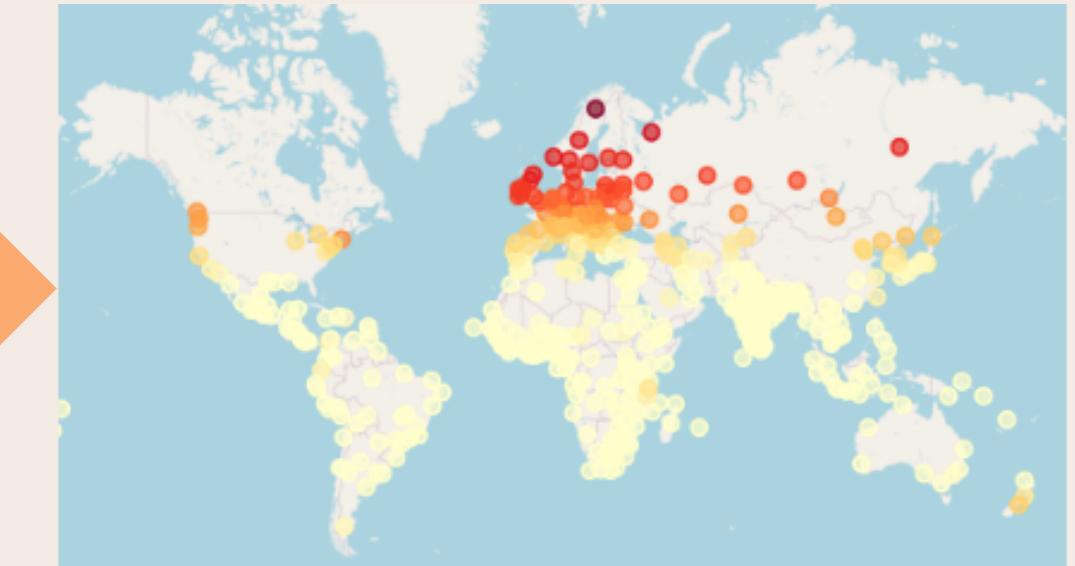
**Heating Degree Days**



Low High

Calculate HDD to measure the demand for heating energy based on the difference between outdoor temperature and a baseline "comfort" temperature, typically 65°F (18°C).

**Heating Degree Days for City**



Low High

Assign HDD to the nearest data point(city level) of interest

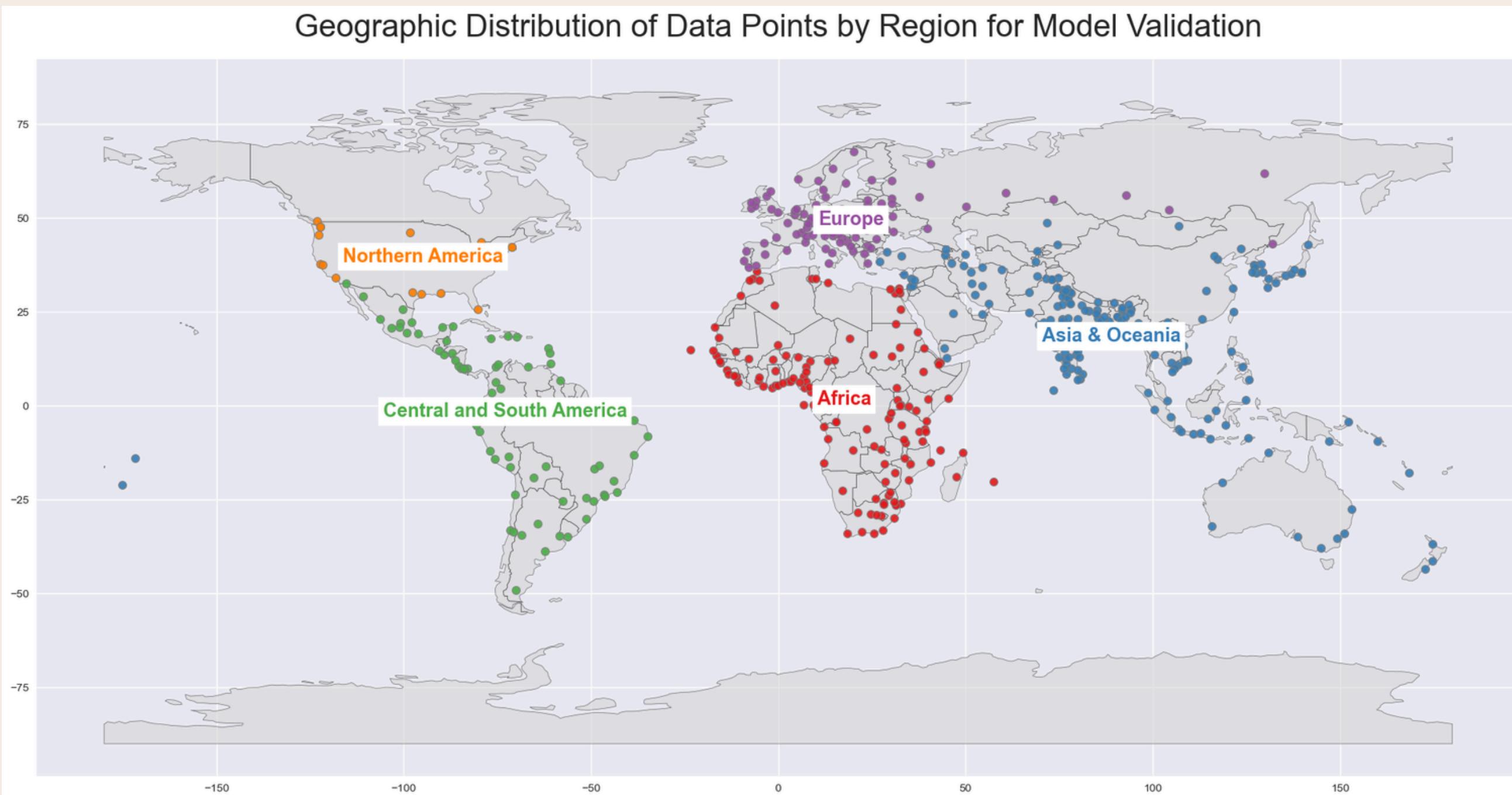


# Experimental Design - Modeling Energy Use Intensity

- **Objective:** Predict Energy Use Intensity
- **Models Used:**
  - K-NN
  - Linear Regression (Lasso and Ridge)
  - Random Forest
  - XGBoost
  - CatBoost
- **Feature Selection:** Removed less important features based on feature importance analysis.
- **Hyperparameter Tuning:** Grid Search for MAPE optimization.
- **Validation Strategy:** Region-based split for validation

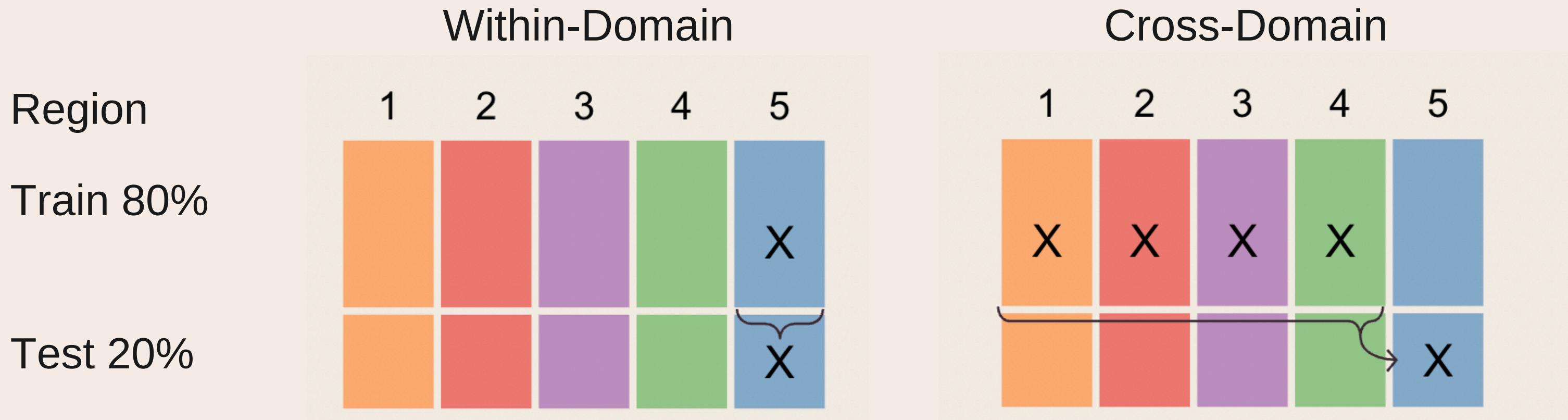


# Experimental Design - Regional Validation Approach



- One challenge with global data is the difficulty of extrapolating results due to regional variations.
- To address this, we will **validate** our predictions at the regional level and estimate Energy Use Intensity while considering these effects. This approach helps us identify biases and enhances the model's robustness for better **extrapolation**.
- We de are considering 5 regions

# Experimental Design - Regional Validation Approach



We train our model on 80% of the data from each region and test it on the remaining 20% of the same region.

To obtain a global result, we calculate the average of our evaluation metrics across the 5 regions.

We train our model using 80% of the data from 4 regions and test it on the 20% of the 5th region. We repeat this to evaluate the model's extrapolation.

The global result is the average of all outcomes.

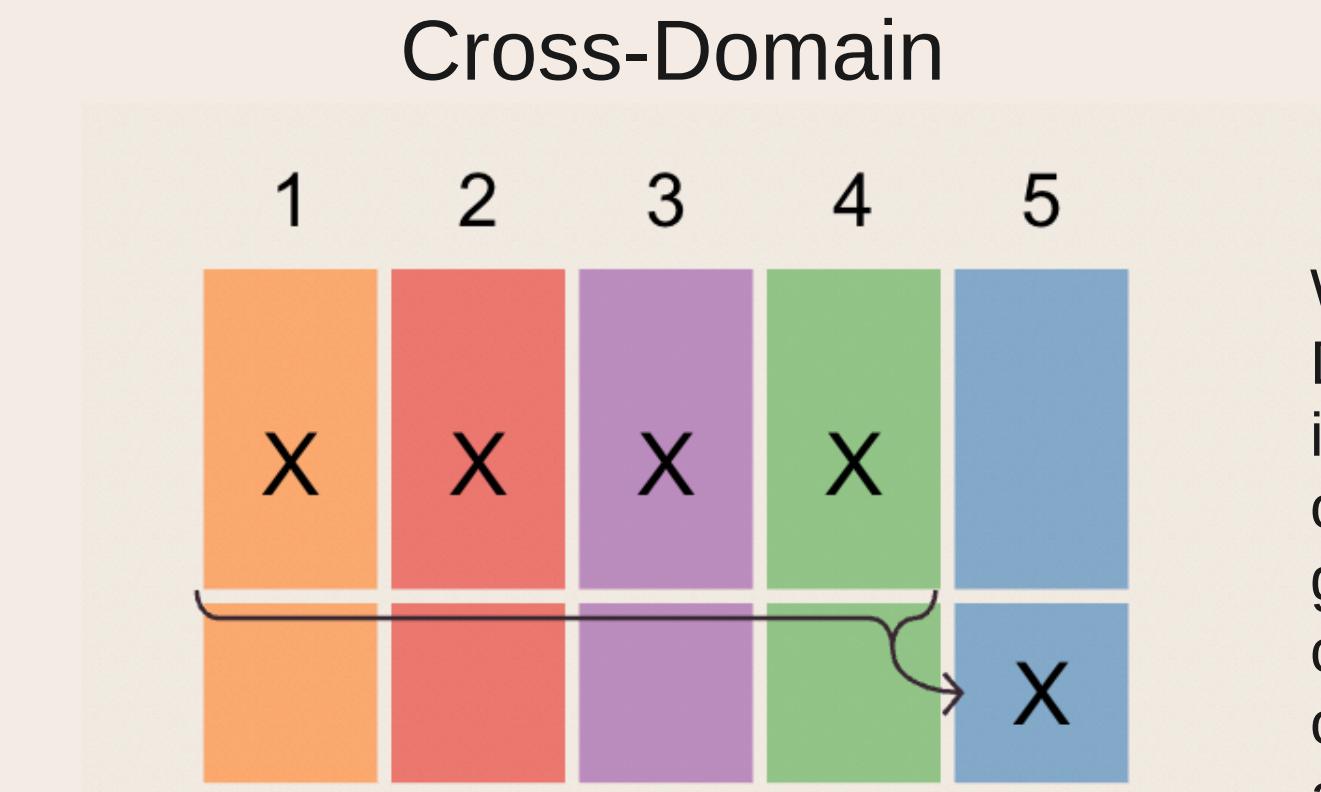
# Experimental Design - Regional Validation Approach



We train our model on 80% of the data from each region and test it on the remaining 20% of the same region.

To obtain a global result, we calculate the average of our evaluation metrics across the 5 regions.

Optimistic Validation



We train our model using 80% of the data from 4 regions and test it on the 20% of the 5th region. We repeat this to evaluate the model's extrapolation.

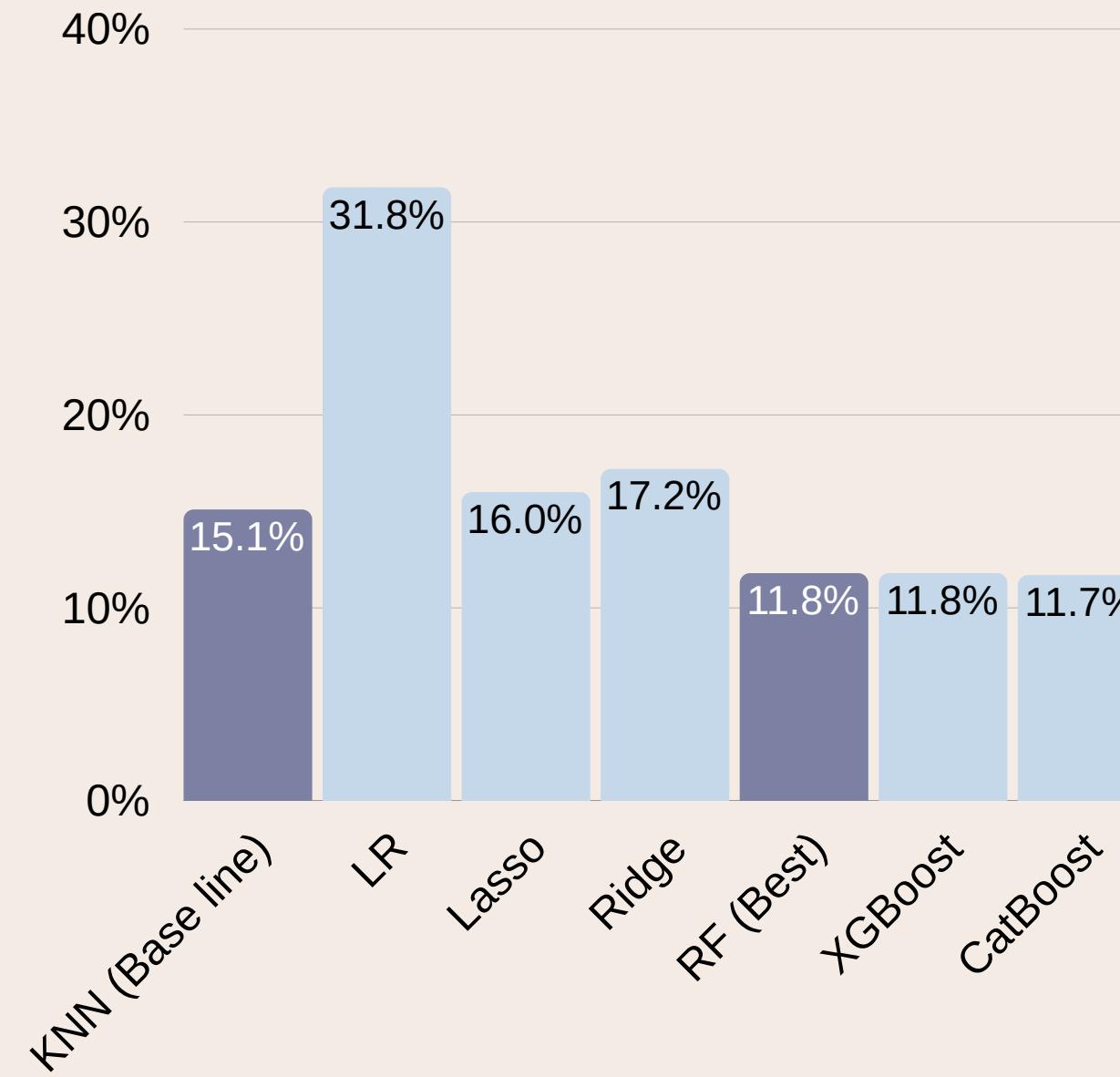
The global result is the average of all outcomes.

Conservative Validation

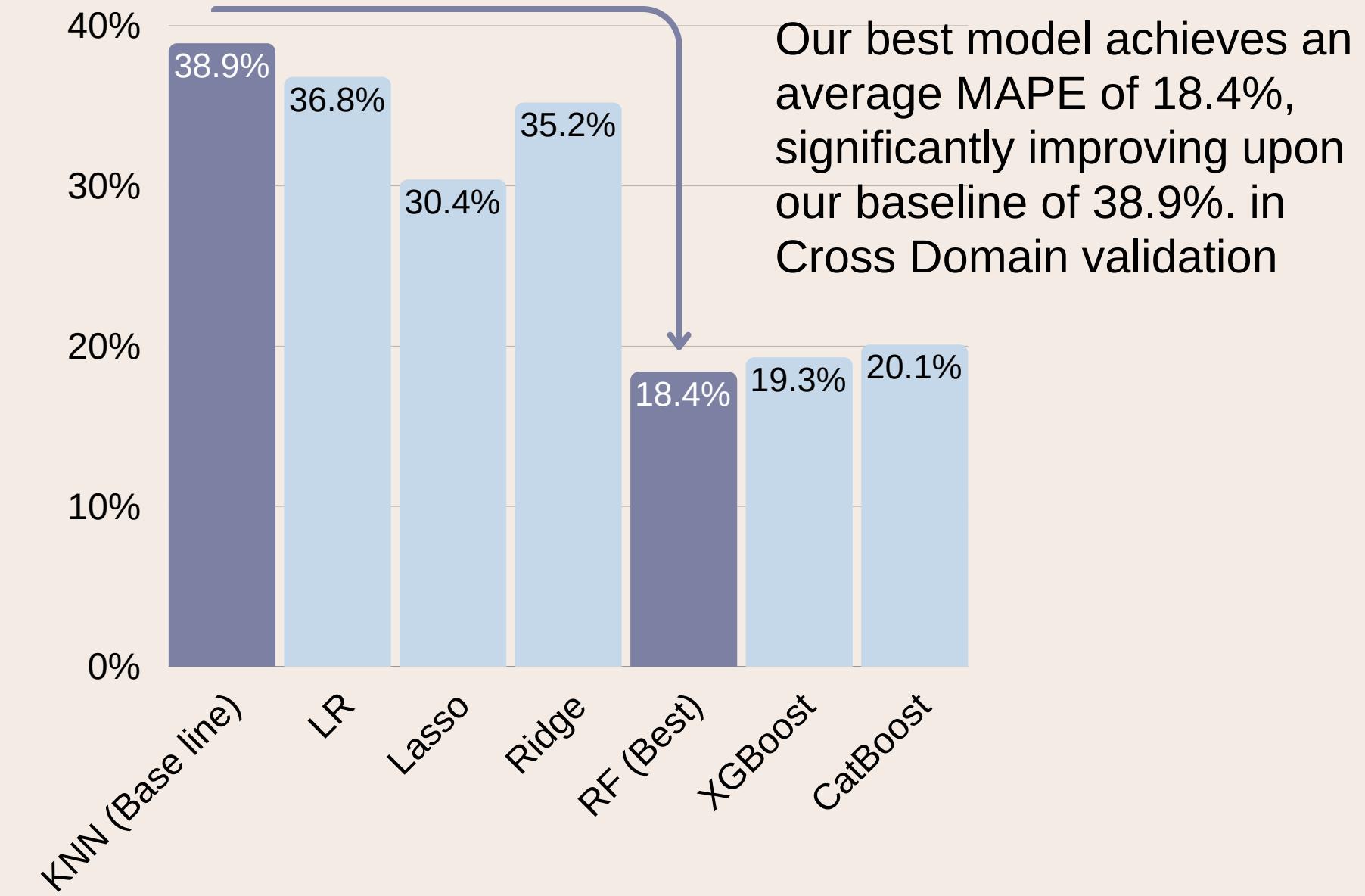
We prioritize Cross-Domain validation as it better represents our model's real-world generalization capabilities and optimize our models accordingly

# Results - Average MAPE across regions and building types

Within-Domain



Cross-Domain

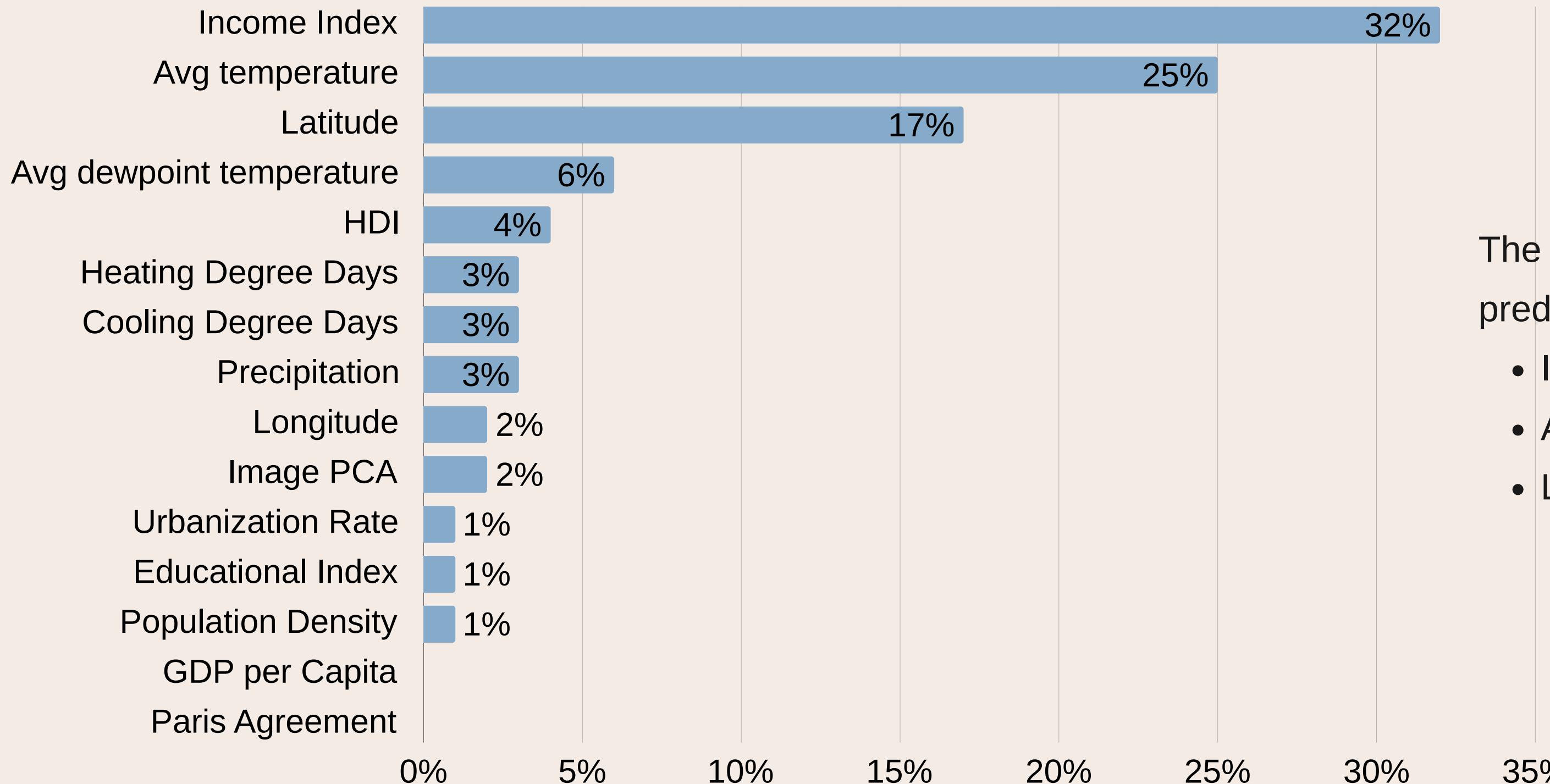


Optimistic Validation

Conservative Validation

**MAPE:** Mean Absolute Percentage Error. Is the average of the absolute percentage errors between predicted and actual values.  
A lower **MAPE** indicates better model performance.

# Results - Feature Importance Using Random Forest



The most important variables in predicting EUI in buildings are:

- Income Index
- Average Temperature
- Latitude

**Heating Degree Days:** quantify heating demand based on temperatures below a set threshold

**Cooling Degree Days:** quantify cooling demand based on temperatures exceeding a set threshold

**HDI Human Development Index:** measures a country's development based on life expectancy, education, and income per capita.

**GDP Gross Domestic Product:** is the total value of goods and services produced in a country.

# Results

- Ensemble models outperformed individual algorithms, with **Random Forest** achieving the **best results** overall
- Our best model achieved an **average MAPE of 18.4%**, representing a **53% improvement from the baseline** (38.9%) in cross-domain validation
- Socioeconomic factors proved crucial - **Income Index** is the most influential predictor (**32% importance**), followed by **Temperature (25%)**



# Implications

- **Filling the Knowledge Gap:** This project fills the knowledge gap of local and municipal emissions inventories by having accurate building emissions data, crucial for estimating and reducing global GHG emissions.
- **Informing Climate Policy:** Access to accurate emissions data enables policymakers to develop targeted, data-driven strategies for reducing GHG emissions.
- **Expanding Data Sources:** Given data limitations, we explored additional public EUI datasets. However, U.S. city-level data offered no new insights beyond our original dataset, underscoring the challenge of obtaining comprehensive global building energy data.

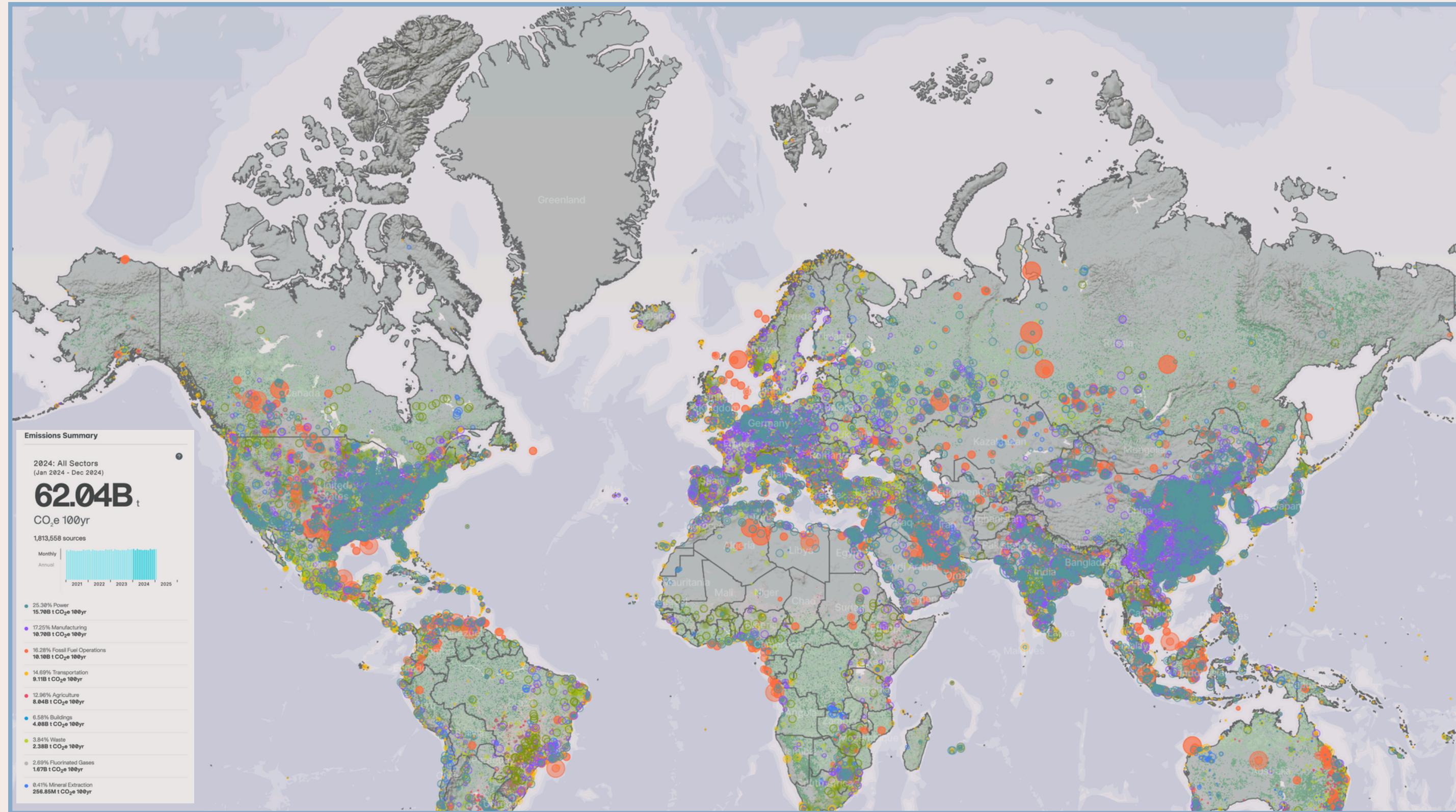


# Limitations

- **Data Sparsity:** With only 482 data entries, the dataset is insufficient to accurately represent global trends and predict the EUI for the entire world.
- **Lack of Building-Specific Features:** Key factors such as age, materials, and energy codes are missing due to the limited granularity of the dataset.
- **EUI Data Constraints:** Derived from a single, possibly imprecise time point and centered on major cities (overlooking suburban/rural areas). Aggregation data that applies US-based building-use assumptions globally and converting units adds further uncertainty, limiting overall accuracy and usefulness.



# Future Work



Global Scale Extrapolation

Estimating GHG Emissions

Exploring New Features and Model Refinements

# Thanks!

We are ready for your questions!

