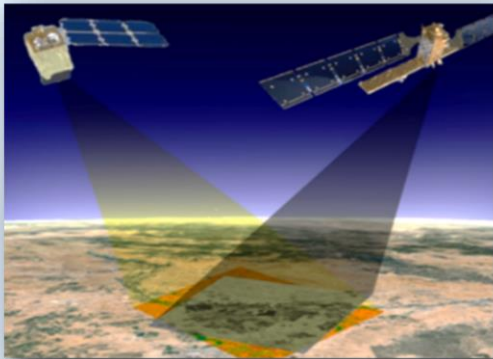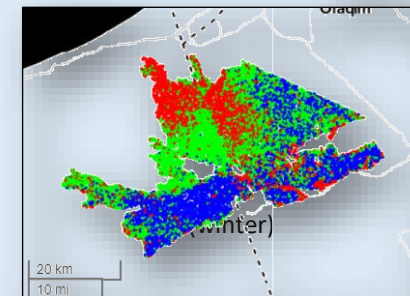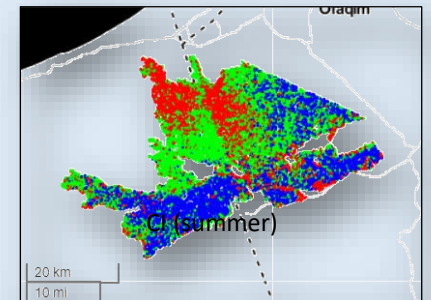# Mining Seasonal Patterns in Earth Observations

*** Thesis Final Exam Presentation ***

By:
**Ricky Shama**

Under the Supervision of:
Prof. Mark Last
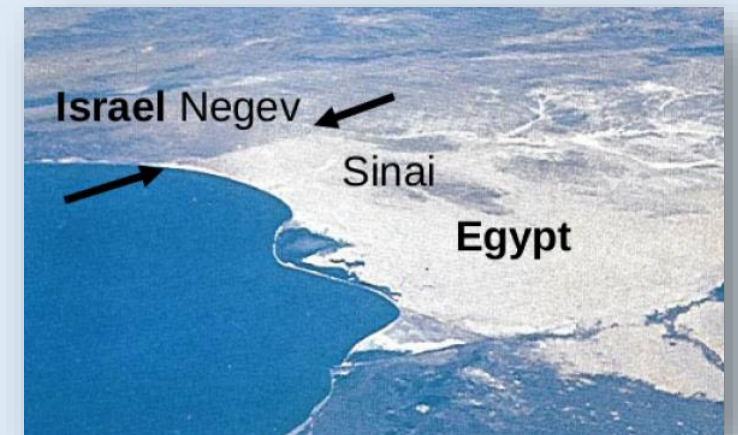Prof. Arnon Karnieli

14.6.22

# Agenda

- Introduction

- Background

- Research Definition

- Methodology

- Implementation
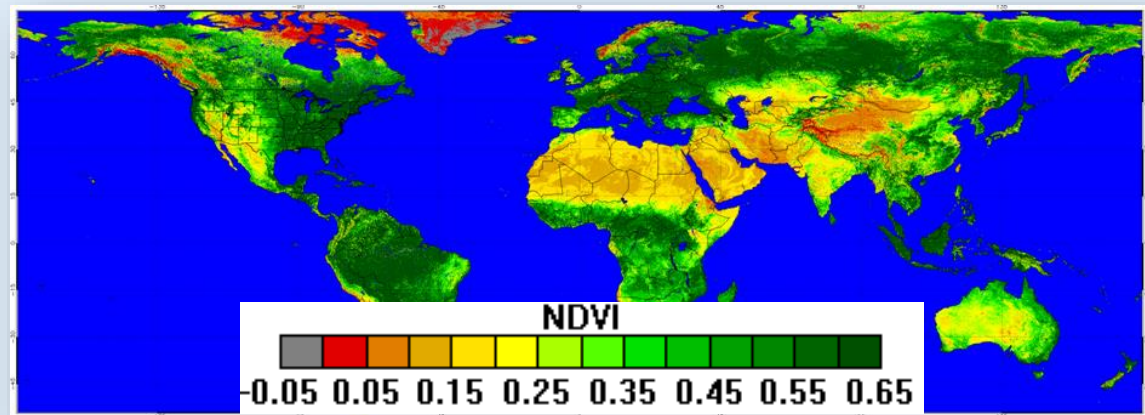
- Results

- Conclusions

- EO (earth observations) are the **main** resource used for monitoring and understanding climate changes.

- They are considered as **Big Data** (Volume, Velocity, Variety).

- **GEE** (Google Earth Engine): analytical platform.

- **Periodicity/ Seasonality Detection**: periodic pattern mining in time-series data. Finding periodic behaviors is useful in other time-series tasks, including: forecasting, clustering, etc.

➔ In this research I developed a seasonality detection model

(using GEE, machine-learning) in earth observations;

**Case Study Area**: *The Israeli-Egyptian Sandfield* (unique seasonal

optical dynamics, human activities vs. nature conservation).

- **Remote Sensing (RS) Index:** a numerical indicator which is achieved using a mathematical formula applied on various **spectral bands** of an image per **pixel**. The **level** of the index indicates a rate of a cover (e.g vegetation), content, etc.

- RS indices used in this research: **CI** (Crust Index), **LST** (Land Surface Temperature), **NDVI /NDWI** (Normalized Difference Vegetation/ Water Index), **precipitation**, **RADAR** backscatter coefficient.
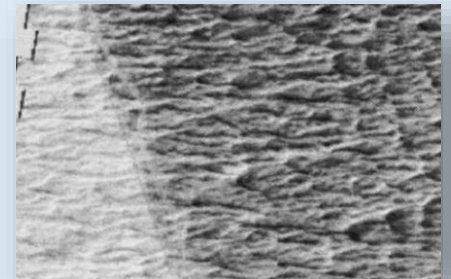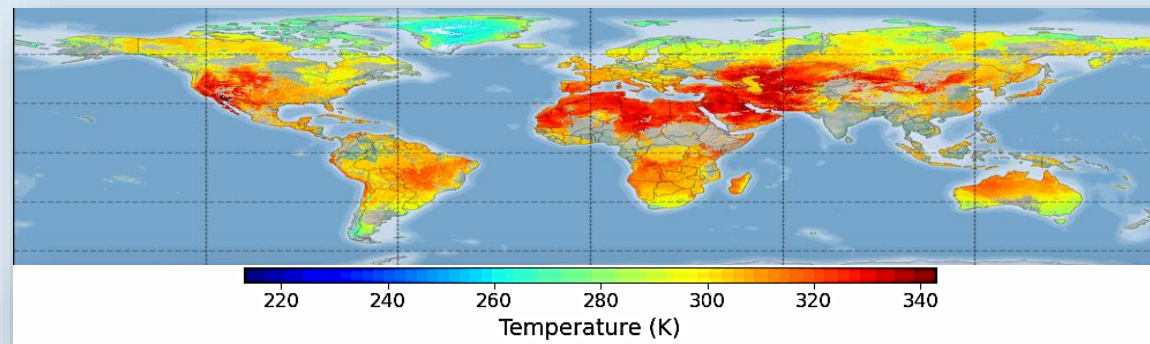


NDVI
-0.05 0.05 0.15 0.25 0.35 0.45 0.55 0.65

4

**Examples:**

- **CI** – (Biological) **C**rust **I**ndex. Soil crust contain cyanobacteria, mosses, algae, etc. that are essential components of arid and semi-arid eco-systems. There is a unique feature of the pigment found in soil crusts, resulting in a relatively higher reflectance in the **blue** spectral region compared to soil without cyanobacteria. A **higher** value of CI indicates a higher content of cyanobacteria in the soil.
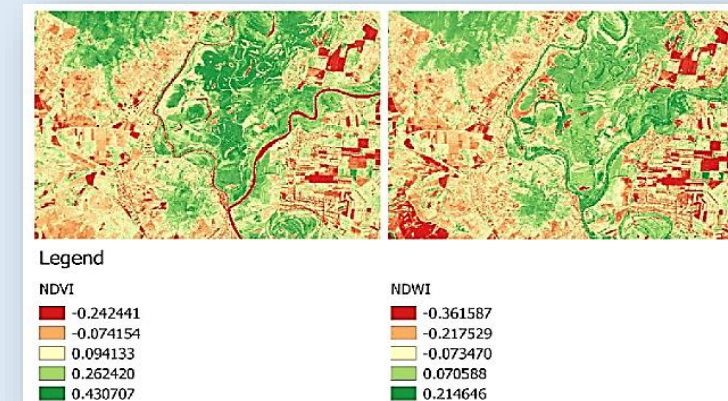
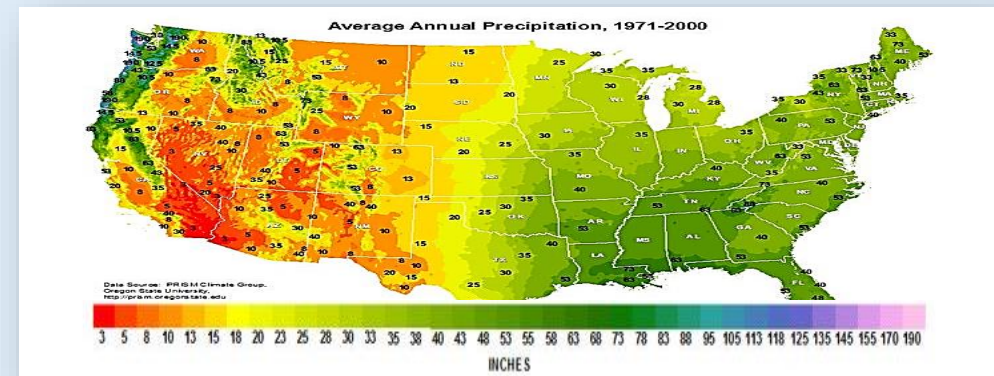$$CI = 1 - \frac{RED-BLUE}{RED+BLUE}$$

- **LST** – **L**and **S**urface **T**emperature. It is estimated by the **Infra-Red** spectral channels of the satellites' sensors. It depends on the **albedo** (fraction of light that is reflected by a body or surface), the vegetation cover and soil moisture. They all respond rapidly to changes in solar radiation due to cloud cover, aerosol concentration, and daily variations of illumination, which affect the LST, too.

- **NDVI** – **N**ormalized **D**ifference **V**egetation **I**ndex. Useful for vegetation monitoring. It describes the difference between visible and near-infrared reflectance of **vegetation cover**, thus, it can be used to estimate the **density of green** on an area of land.

$$NDVI = \frac{NIR-RED}{NIR+RED}$$

- **NDWI** – **N**ormalized **D**ifference **W**ater **I**ndex. Reflects moisture content in plants and soil.

$$NDWI = \frac{NIR-SWIR}{NIR+SWIR}$$



Legend

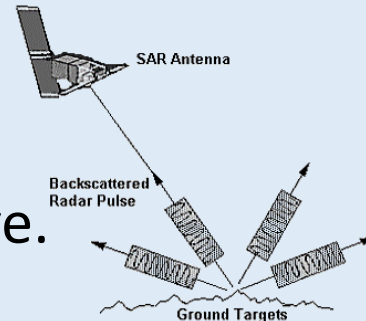| NDVI | | NDWI | |
|---|---|---|---|
| | -0.242441 | | -0.361587 |
| | -0.074154 | | -0.217529 |
| | 0.094133 | | -0.073470 |
| | 0.262420 | | 0.070588 |
| | 0.430707 | | 0.214646 |

- **Precipitation** – Any liquid or frozen water that forms in the atmosphere and falls back to the Earth. It comes in many forms, like: rain, hail, snow. Along with evaporation and condensation, precipitation is one of the three major parts of the global water cycle. Can be measured either by ground-based instruments, or satellites' sensors which estimate the electro-magnetic radiation as reflected from the top of the clouds, rain droplets.



Average Annual Precipitation, 1971-2000

- **RADAR** Backscatter Coefficient – The portion of the outgoing radar signal that the target redirects directly back towards the radar antenna. In general, it is computed as the **ratio** between the **received** energy by the sensor, to the **transmitted** energy by the source. Usually measured in [dB].
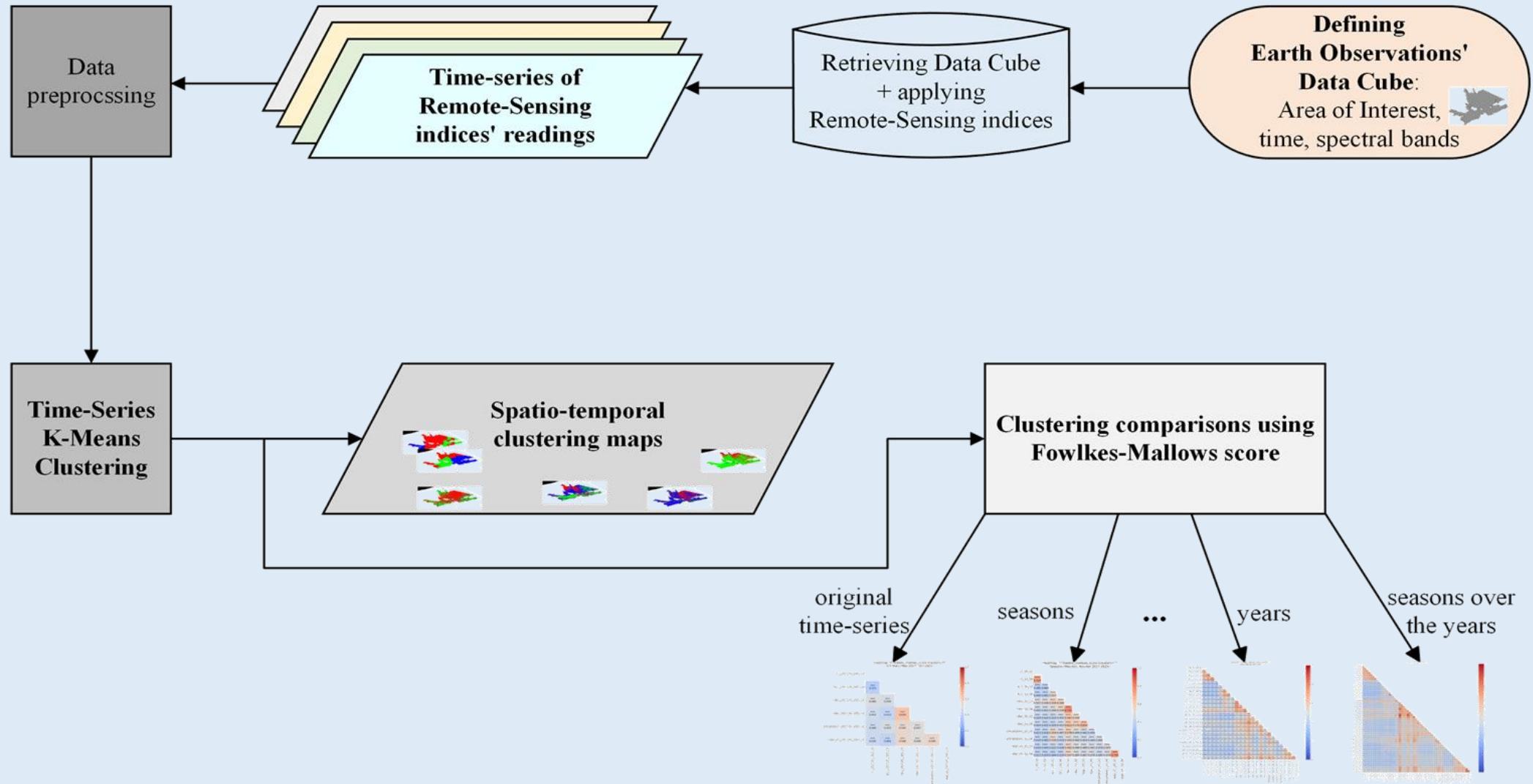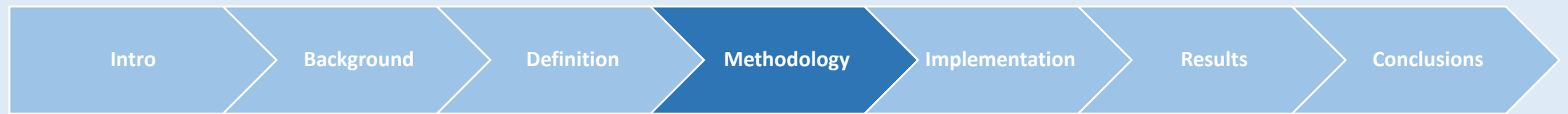
  → Example of polarization mode: VH – **V**ertical Transmit-**H**orizontal Receive.

  → SAR **– S**ynthetic **A**perture **R**adar. A form of radar system that is used to create 2D/3D reconstructions of objects, such as landscapes, independent of weather, and solar radiation. The applications include **topography**, geology mapping.
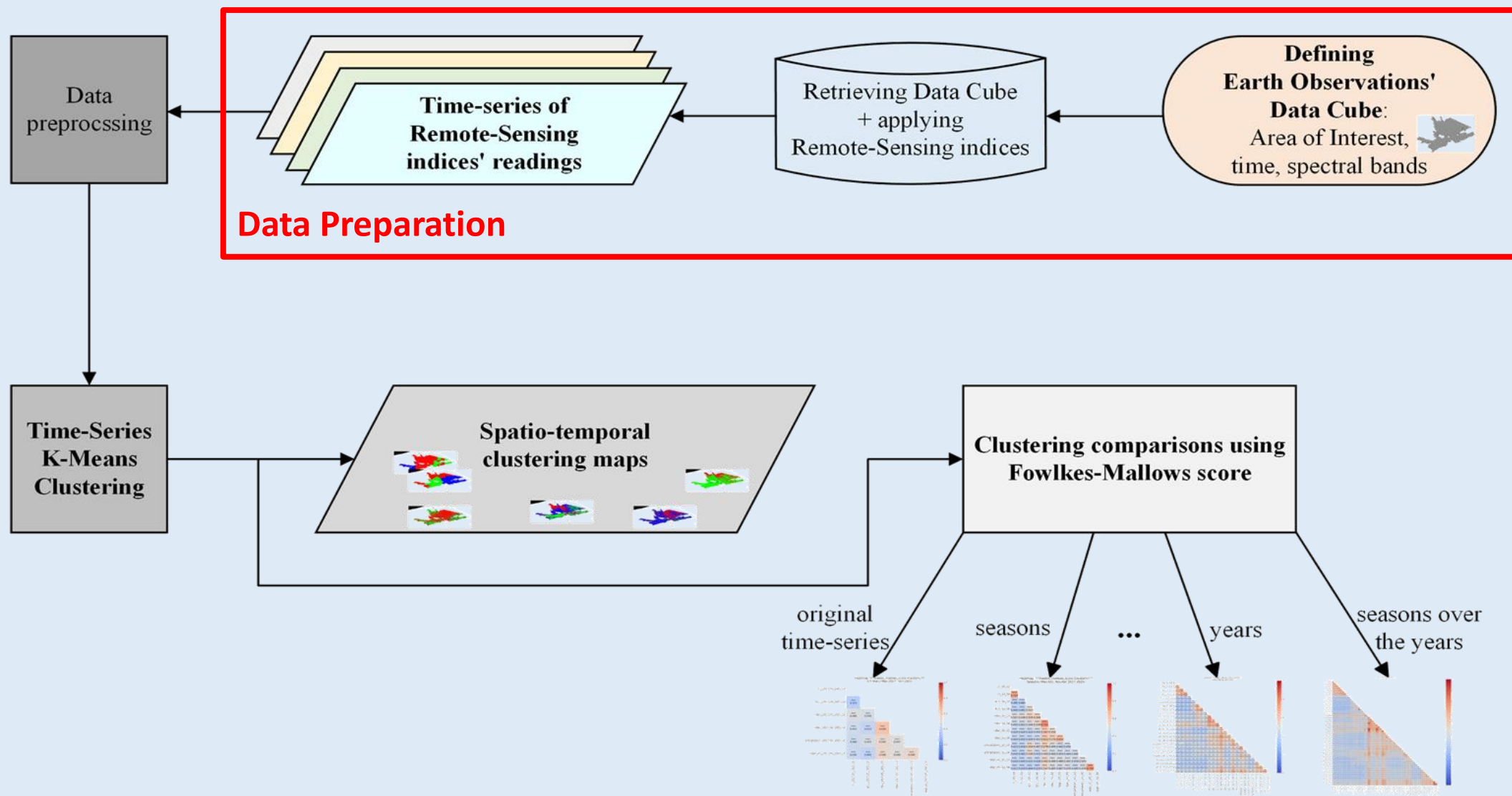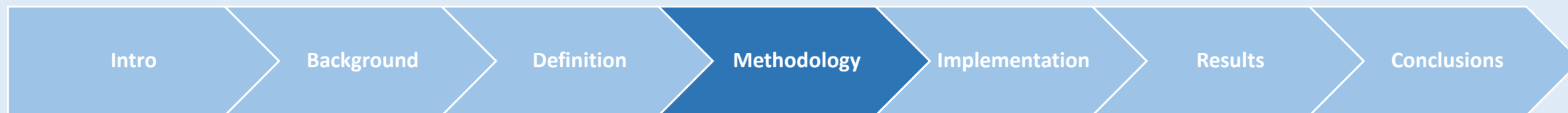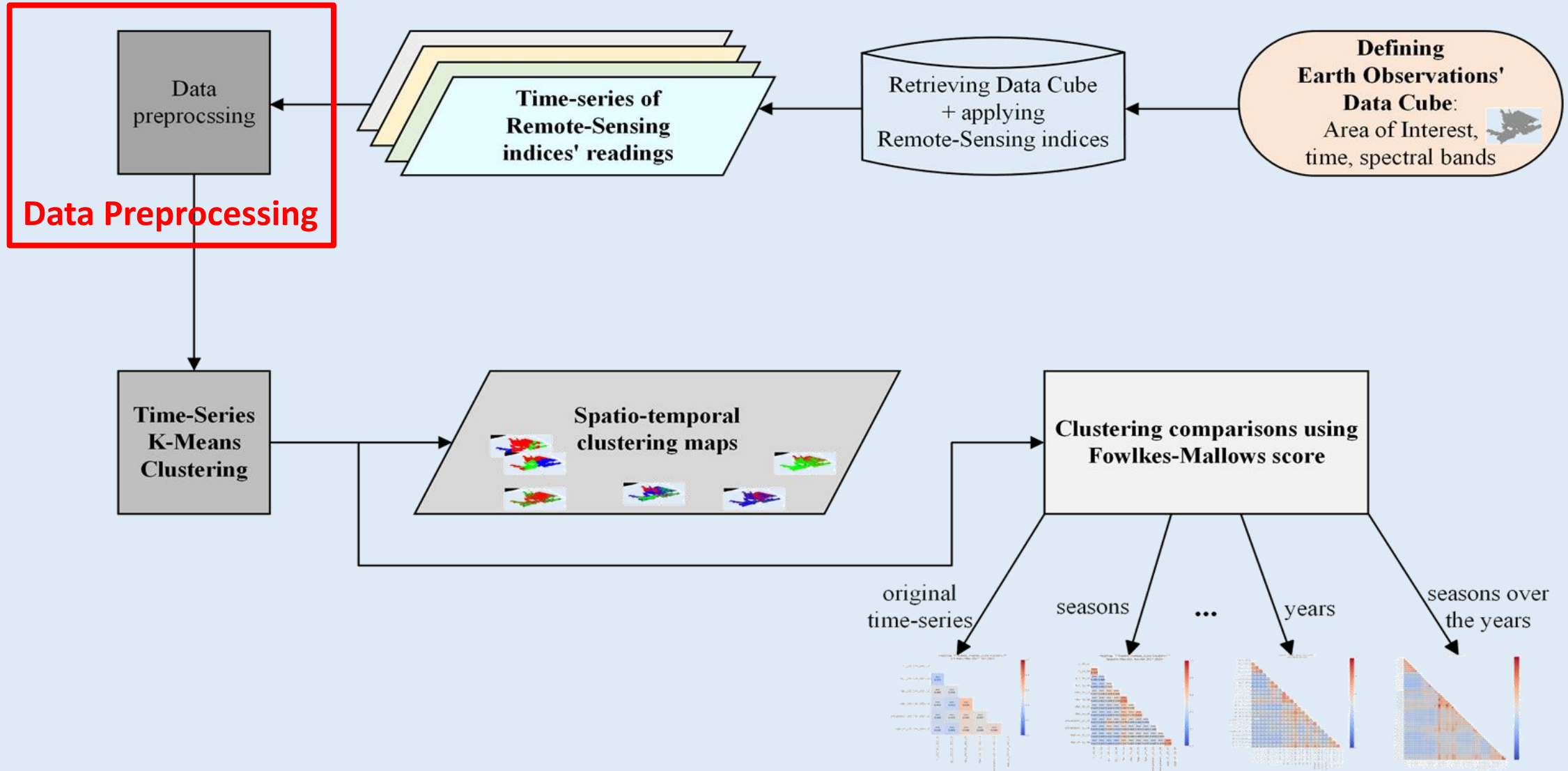


**9**

- **Research Goal**: evaluating the similarities between:

  - *Different* time intervals' **clusterings**, of the *same* RS index.

  - *Same* time intervals' **clusterings**, of *different* RS indices.

  - *Different* time intervals' **clusterings**, of *different* RS indices.

- **Motivation**: areas representing natural and anthropogenic land transformation, require monitoring, understanding the changes, and responding to them.
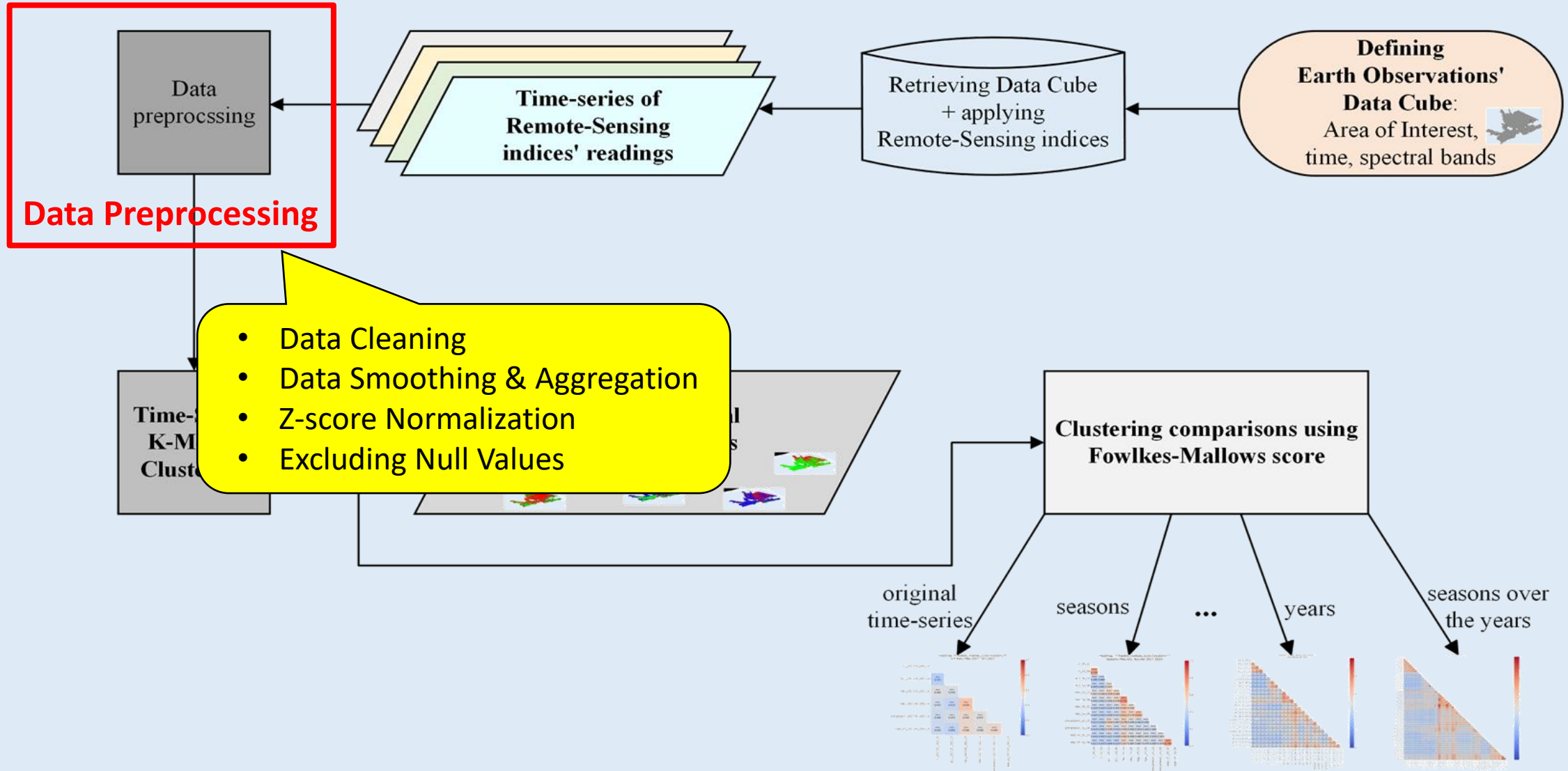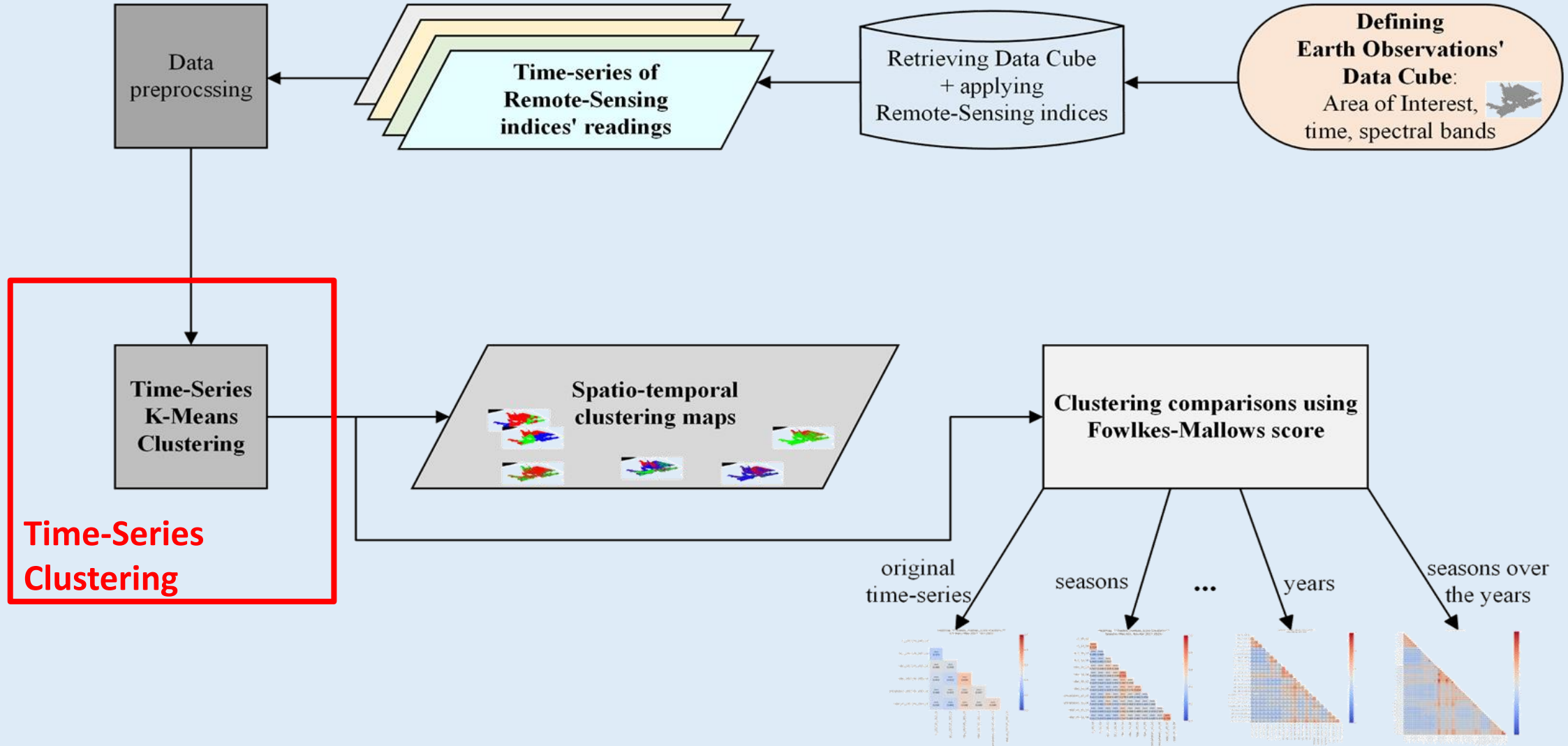
- **Current limitations:** Existing studies mainly deal with EO data related to a **limited** number of RS indices (~1-4), with low frequency time-series analysis (e.g. yearly, monthly sampling).

- **Contribution**: This research suggests a seasonality detection model which is based on data related to **multiple** RS indices (6), which together cover a **wide** range of the electromagnetic spectrum, with **high** frequency time-series analysis (daily/ weekly sampling), in order to capture the subtle optical changes/ anomalies in time.

**12**

**Data Preparation**

Data preprocsing

Time-series of Remote-Sensing indices' readings

Retrieving Data Cube + applying Remote-Sensing indices

Defining Earth Observations' Data Cube: Area of Interest, time, spectral bands

Time-Series K-Means Clustering

Spatio-temporal clustering maps

Clustering comparisons using Fowlkes-Mallows score

original time-series

seasons

...

years

seasons over the years

**13**

**Data Preprocessing**

**Data Preprocessing**

- Data Cleaning
- Data Smoothing & Aggregation
- Z-score Normalization
- Excluding Null Values

**Time-Series Clustering**

Data preprocssing

Time-series of Remote-Sensing indices' readings

Retrieving Data Cube + applying Remote-Sensing indices

**Defining Earth Observations' Data Cube**: Area of Interest, time, spectral bands

Time-Series K-Means Clustering

Spatio-temporal clustering maps

Clustering comparisons using Fowlkes-Mallows score

original time-series

seasons
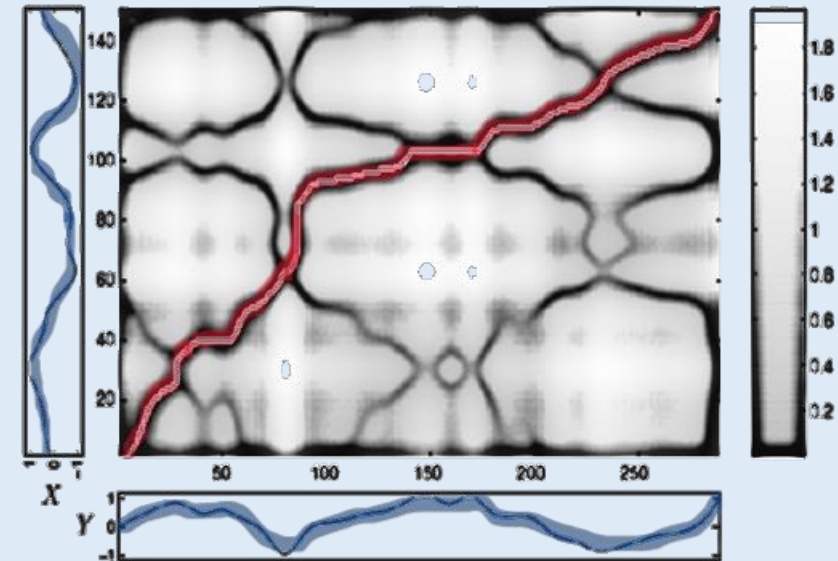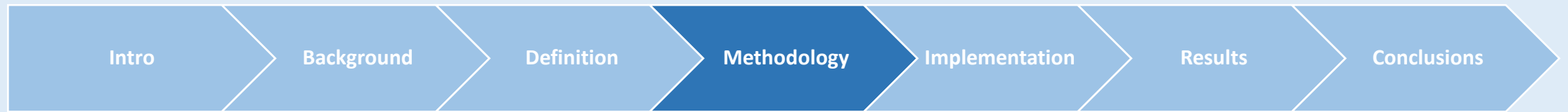
...

years

seasons over the years

**Time Series Clustering**:

- Each RS index original time-series of readings (i.e. during the entire defined time period) is split into the following time **intervals**: seasons (dry/ wet), years (cycles), season each year.

- A **K-Means Clustering** is performed on all n points' time-series, for each time-interval and RS index, using **Dynamic Time Warping (DTW).**



17

• **Dynamic Time Warping (DTW)**: measuring similarity between two temporal sequences (**X**, **Y**), which may vary in speed and length. A cost/ distance measure is applied on each pair of samples, resulting in a **cost matrix**. The goal is to find an alignment between the sequences having a minimal overall cost ("**warping path**").

- Here, the similar-shaped time intervals (as determined by the DTW) are grouped into **clusters**. The number of clusters, **K**, is selected beforehand, out of a range of numbers of clusters, if it reached the highest average **Silhouette** score, for all time intervals, of all RS indices.
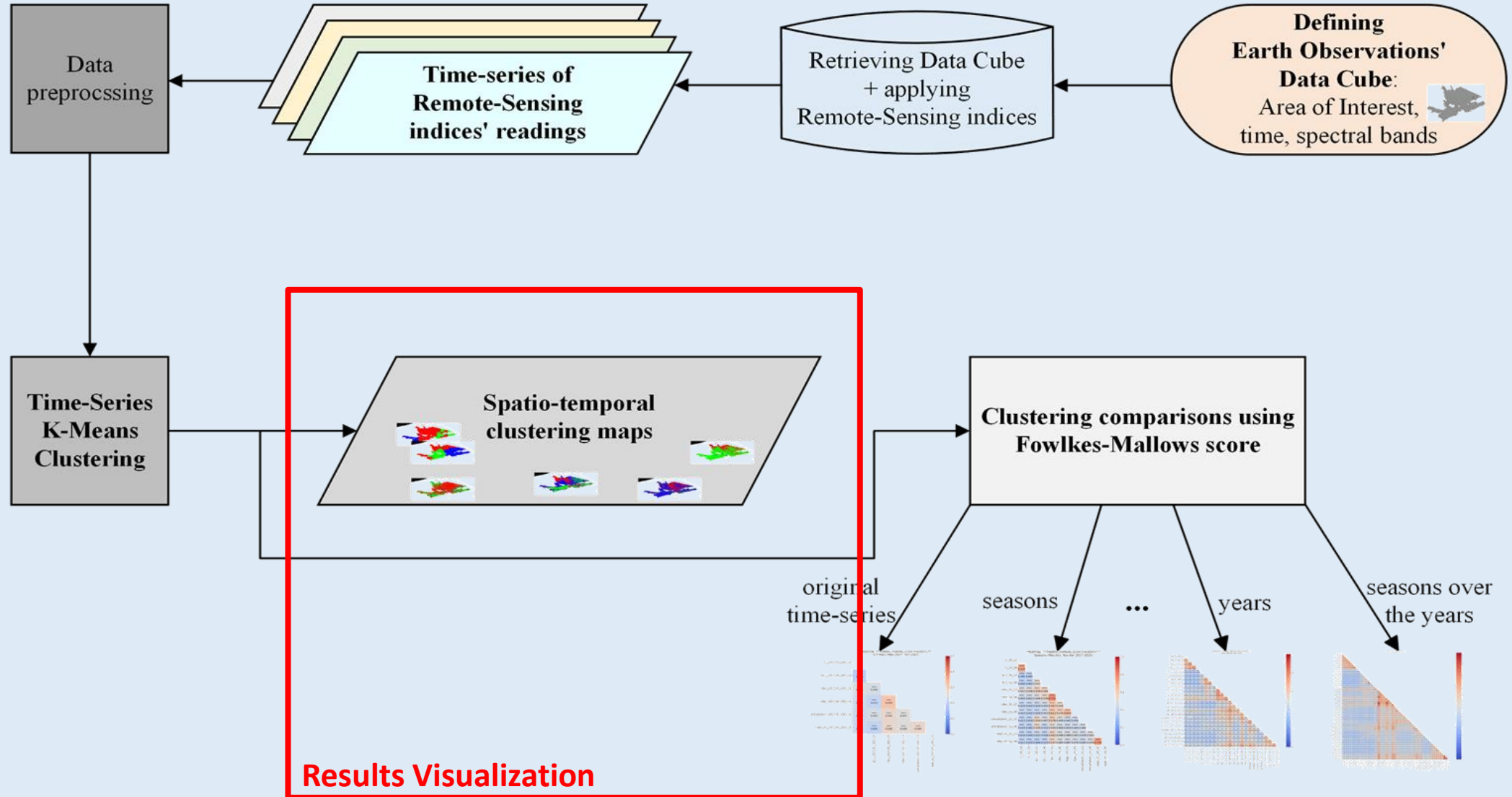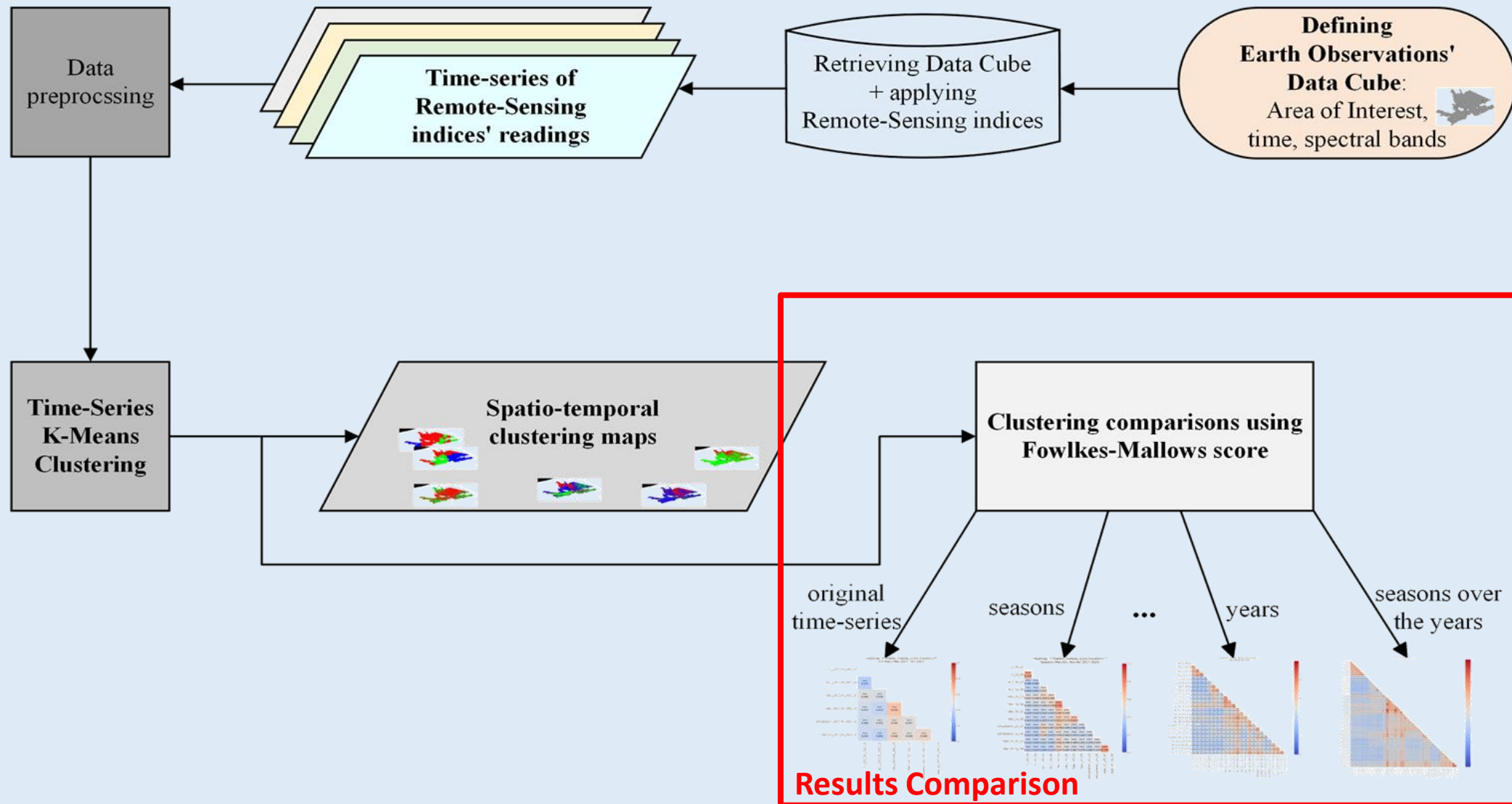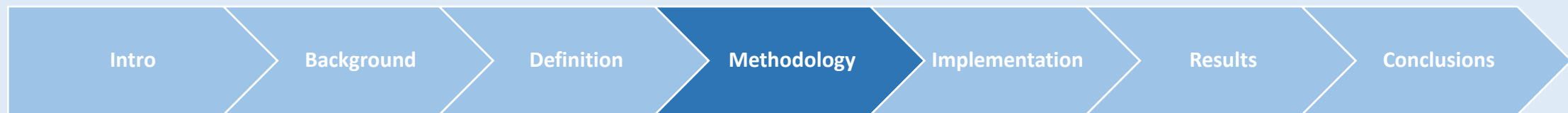
- **Silhouette Score**: a metric used to evaluate the quality of the clustering configuration. Its value ranges from -1 to 1. "1" indicates that the clusters are well separated. "0" means that the clusters are overlapping. "-1" indicates that the clusters are wrongly assigned.

$$s = \frac{b-a}{\max(a, b)}$$

intra-cluster distance

inter-cluster distance

**Results Visualization**

21

Data preprocsing

Time-series of Remote-Sensing indices' readings

Retrieving Data Cube + applying Remote-Sensing indices

**Defining Earth Observations' Data Cube**: Area of Interest, time, spectral bands

Time-Series K-Means Clustering

Spatio-temporal clustering maps

Clustering comparisons using Fowlkes-Mallows score

original time-series

seasons

...

years

seasons over the years

**Results Comparison**

- **Fowlkes-Mallows (FM) Score**: measures similarity between two **clusterings**. The FM can be used to compare either two cluster label sets: S1, S2 (**symmetrical**) or a cluster label set with a true label set. It is defined as the geometric mean between of the precision and recall.

- The score ranges from 0 (totally random) to 1 (identical).

$$FM = \sqrt{\frac{TP}{TP+FP}} * \sqrt{\frac{TP}{TP+FN}}$$

# pairs of points belong to the **same** cluster in both sets

# pairs of points that **don't** belong to S2

# pairs of points that **don't** belong to S1

- **Earth Observations' Data Cube:**

  - **AOI (Area of Interest):** The Israeli-Egyptian Sandfield (n=20,000 points; **uniformly random** within the area).

  - **Time:** May 2017-Oct 2021 (4.5 years).

  - **RS Indices:** CI, LST, NDVI, NDWI, precipitation, RADAR backscatter.

- **Data Retrieval Platform:** GEE (Sentinel 1+2, MODIS, CHIRPS datasets).

- **Data Preprocessing:** filtering high-quality, cloud-free, not null values, MA, aggregation, averaging, normalization.

- **Earth Observations' Data Cube:**

  - **AOI (Area of Interest):** The Israeli-Egyptian Sandfield (n=20,000 points; **uniformly random** within the area).

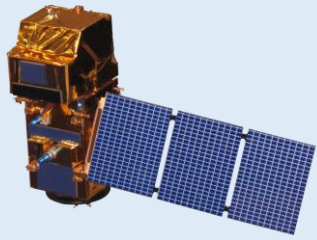  - **Time**:  May 2017-Oct 2021 (4.5 years).

  - **RS Indices**: CI, LST, NDVI, NDWI, precipitation, RADAR backscatter.

- **Data Retrieval Platform**: GEE (Sentinel 1+2, MODIS, CHIRPS datasets).

- **Data Preprocessing**: filtering high-quality, cloud-free, not null values, MA, aggregation, averaging, normalization.

- **Case Study Area:** *The Israeli-Egyptian Sandfield* (total area: 1243.86 [km$^2$], coordinates (center): 31°0'18.9288"N, 34°21'7.6206"E) – is an example of an area which demonstrates particular seasonal optical dynamics in various spectral regions, ranging from **Visible** to **Microwave**, in various **temporal** scales (seasonal, annually, etc.), and the effects of overgazing vs. nature conservation from the two sides of the Israeli-Egyptian border.

- **Earth Observations' Data Cube:**

  - **AOI (Area of Interest):** The Israeli-Egyptian Sandfield (n=20,000 points; **uniformly random** within the area).

  - **Time**:  May 2017-Oct 2021 (4.5 years).

  - **RS Indices**: CI, LST, NDVI, NDWI, precipitation, RADAR backscatter.

- **Data Retrieval Platform**: GEE (Sentinel 1+2, MODIS, CHIRPS datasets).

- **Data Preprocessing**: filtering high-quality, cloud-free, not null values, MA, aggregation, averaging, normalization.

• **Datasets' details:**

| RS Index [unit] | Satellite/ Source | Spectral Region(s), Band(s) | Spatial Res. [m] | Revisit Time [Days] | # Samples |
|---|---|---|---|---|---|
| CI | Sentinel-2 (A, B) | VIS [nm]:<br>Band 2 (Blue): 496.6 (A)/ 492.1 (B).<br>Band 4 (Red): 664.5 (A)/ 665 (B). | 10 | 5 | 275 |
| NDVI | | VIS [nm]:<br>Band 4 (Red): 664.5 (A)/ 665 (B).<br><br>NIR [nm]:<br>Band 8: 835.1 (A)/ 833 (B). | 10 | 5 | 275 |
| NDWI | | NIR [nm]:<br>Band 8: 835.1 (A)/ 833 (B).<br><br>SWIR [nm]:<br>Band 12: 2202.4 (A)/ 2185.7 (B). | 10 (NIR),<br>20 (SWIR) | 5 | 275 |
| (Nightly) LST [Average/week][C] | MODIS | Calculation based on<br>7 TIR MODIS bands [$\mu$m]:<br>1. Band 20: 3.660–3.840.<br>2. Band 22: 3.929–3.989.<br>3. Band 23: 4.020–4.080.<br>4. Band 29: 8.400–8.700.<br>5. Band 31: 10.780–11.280.<br>6. Band 32: 11.770–12.270.<br>7. Band 33: 13.185–13.485. | 1000 | 1-5 | 118 |
| precipitation [mm/week] | CHIRPS | Estimation based on:<br>multi-satellites TIR data +<br>rain gauges measurements. | 5566 | 1 | 236 |
| RADAR VH back-scatter [dB] | Sentinel-1 (B) | Microwave [cm]:<br>5.547 (= 5.405 [GHz]). | 10 | 5 | 261 |

- **Earth Observations' Data Cube:**

  - **AOI (Area of Interest):** The Israeli-Egyptian Sandfield

    (n=20,000 points; **uniformly random** within the area).

  - **Time**:  May 2017-Oct 2021 (4.5 years).

  - **RS Indices**: CI, LST, NDVI, NDWI, precipitation, RADAR backscatter.

- **Data Retrieval Platform**: GEE (Sentinel 1+2, MODIS, CHIRPS datasets).

- **Data Preprocessing**: filtering high-quality, cloud-free, not null values, MA, aggregation, averaging, normalization.

**Data preprocessing steps:**

1. **Filtering only high-quality readings**: relevant to: **LST** (using MODIS' quality band: "QC_Night"; 954 samples out of 975; 97.846%; before weekly aggregation) and **RADAR** (268 samples).

2. **Filtering only cloud-free readings**: relevant to: **CI**, **NDVI**, **NDWI**. Information provided by the "QA60" cloud mask band of Sentinel-2 sensor. There were 244.342 cloud-free samples in average out of 301 (81.177%).

3. **Excluding null values**: in all datasets.

**Data preprocessing steps (cont.):**

4.  **Synchronization of dates**: in all RS indices, for all 20,000 time-series.

5.  **Applying Moving Average (MA) to remove noise**: e.g. of two-weeks time in **CI**, **NDVI** and **NDWI**, to overcome missing data due to cloudness, and null values.

6.  **Weekly aggregation**: relevant to **precipitation** dataset (1645 → 236 samples).

7.  **Weekly averaging**: relevant to **LST** dataset (954 → 232 samples).

8.  **Z-score normalization**: of all datasets.

**Time Intervals (#):**

- **Original time-series**: 4.5 years – May 2017-Oct 2021 (1).

- **Seasons**: dry (May-Oct) and wet (Nov-Apr) seasons (2).

- **Years** (Cycles): May-Apr (4).

- **Season each year**: dry/ wet season of each year cycle (9).

→ Total # of time intervals for each RS index = **16** (=1+2+4+9).

**Time Series Clustering**:

- **K**-Means Clustering of each interval, for all 20,000 points in AOI, using **DTW**.

- **K=3** (=highest average Silhouette Score out of a range of **3-5** clusters):



**3 3**

## Results Visualization – Spatio-Temporal Clustering Maps:

- **96** clustering maps were generated out of all clustering results of the 20,000 points (= **6** RS indices * **16** time intervals).

- **Example**: clustering map of CI during 2017's dry season (May-Oct):



34

**Results Comparison using Fowlkes-Mallows (FM) Score between**:

- *Different* time intervals' **clusterings**, of the *same* RS index.

- *Same* time intervals' **clusterings**, of *different* RS indices.

- *Different* time intervals' **clusterings**, of *different* RS indices.

- ➔ FM heat-maps.

**Software and hardware information:**

- **Javascript** (datasets generation), **Python** (clustering: *tslearn*, analysis + maps: *sklearn, seaborn, folium*).

- **Total datasets size** (tabular; after preprocessing): ~600 [MB].

- Total time-series clustering **running time** (of all 96 intervals): ~6-7 [days].

- **Hardware specifications**:

| Machine Type/ Name | CPU Model Name | CPU [GHz] | # Cores | RAM [GB] |
|---|---|---|---|---|
| Google Compute Engine (Colab) | AMD EPYC | 2.2 | 2 | 12.69 |
| e2-highmem-4 (GCP) | Intel Xeon | 2.2 | 4 | 32 |

**36**

- **Evaluation**: The were no ground-truth labelling. The whole clustering process was **unsupervised**. The estimation of the reliability of the results was done by Prof. Karnieli & Dr. Micha Silver.

- **Statistical Significance**: of each FM result was calculated using "**Bootstrapping**" procedure with 1000 permutations of each clustering result.

## Examples:

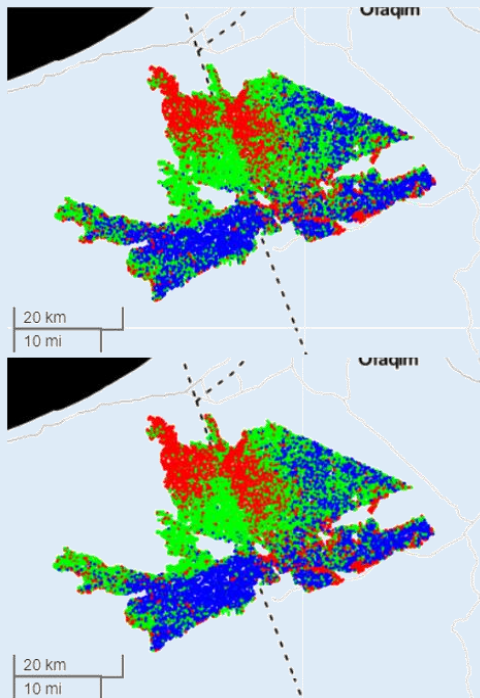- 4.5 years (May 2017-Oct 2021):



NDVI

NDWI

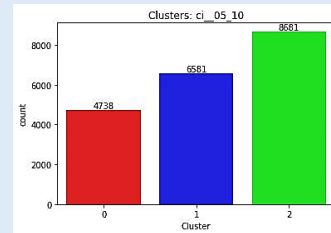Heatmap: ** fowlkes_mallows_score (Clusters) **
4.5 Years (May 2017 – Oct 2021)

- The highest similarity was between the **NDVI** and **NDWI** clusterings (together they reached the highest average FM scores). It fits the assumption that both clusterings will be similar as they reflect **vegetation health**.

- In terms of a single RS index, the: **RADAR**, **CI**, and **NDVI**'s clustering behaviors maintained high level of similarity between the dry and the wet seasons (FM ≥ 0.719). This indicates that the surface topography, biocrust, and vegetation cover remained stable regardless the season.

**Future work**:

- Writing a paper (Remote Sensing).

- **Implementation**: apply the model on other AOIs; more/ other RS indices; longer period of data (= detect long-term patterns/ changes).

- **Methodology**: using/ comparing with other clustering algorithms (e.g. Hierarchical, deep learning); different approach – spatial clustering (e.g. DBSCAN); ensemble of different clustering algorithms/ results.

# Thank you!