

Background

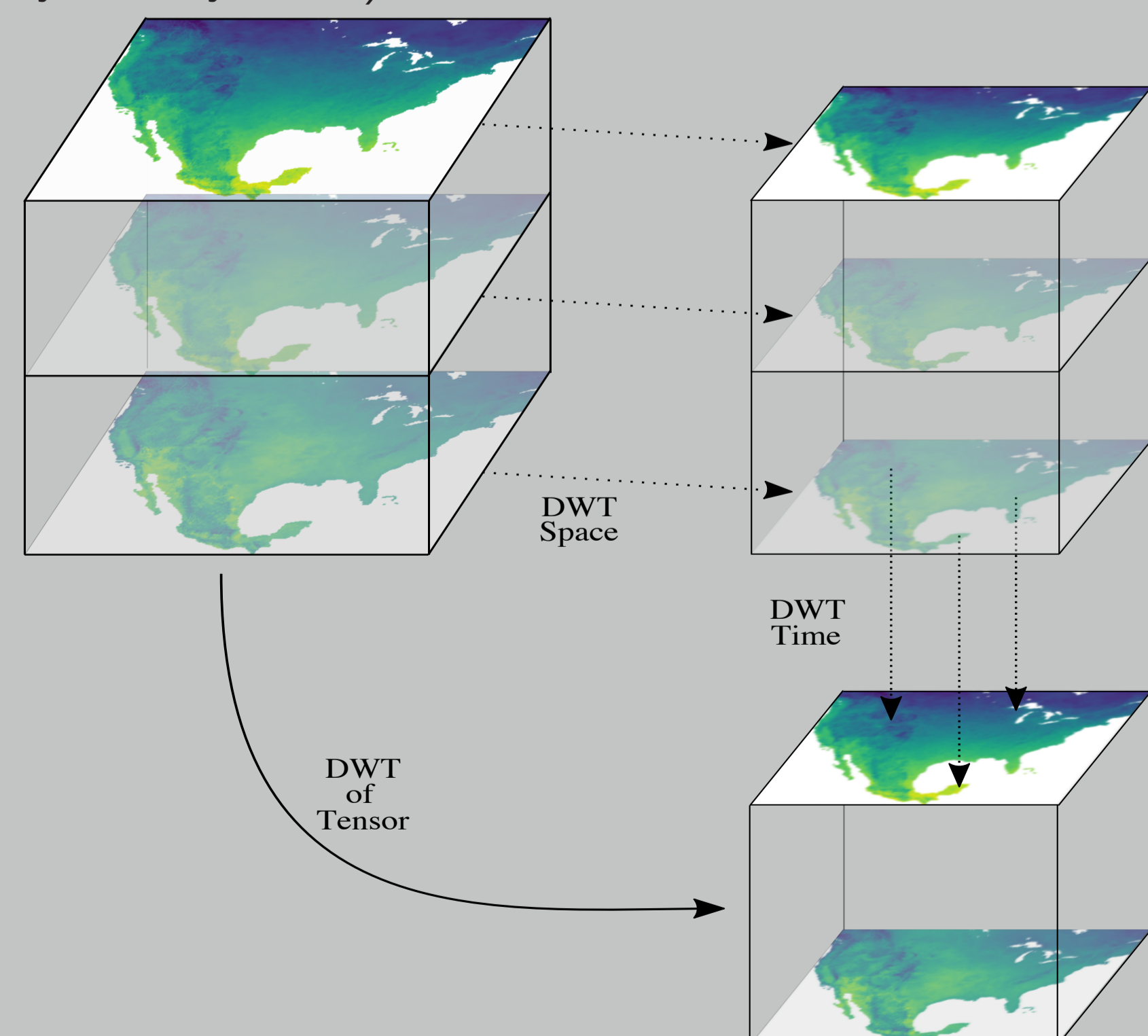
- ▶ Many data structures come in the form of tensors. Frequently, one is interested in the relationship between particular slices of the tensor.
- ▶ Clustering data along a chosen slice provides a method to parse the complexity-complete data into digestible bits of information.
- ▶ It is well known that the results of clustering methods are heavily dependent on the algorithm of choice, as well as the chosen hyperparameters of the algorithm. What isn't understood is the further dependence on "hidden parameters" such as the scale of the data.

Application - Climate Classification

- ▶ Classifying the surface of the Earth into climate biomes is a way to parse the complex micro and macroscopic interdependencies down to provide meaningful diagnostics that more readily relate to physical and biological systems.
- ▶ It is important to analyze both short and long term trends, as well as small and large scale features.
- ▶ Classification of the climate is sensitive to spatial/temporal coarse-graining.

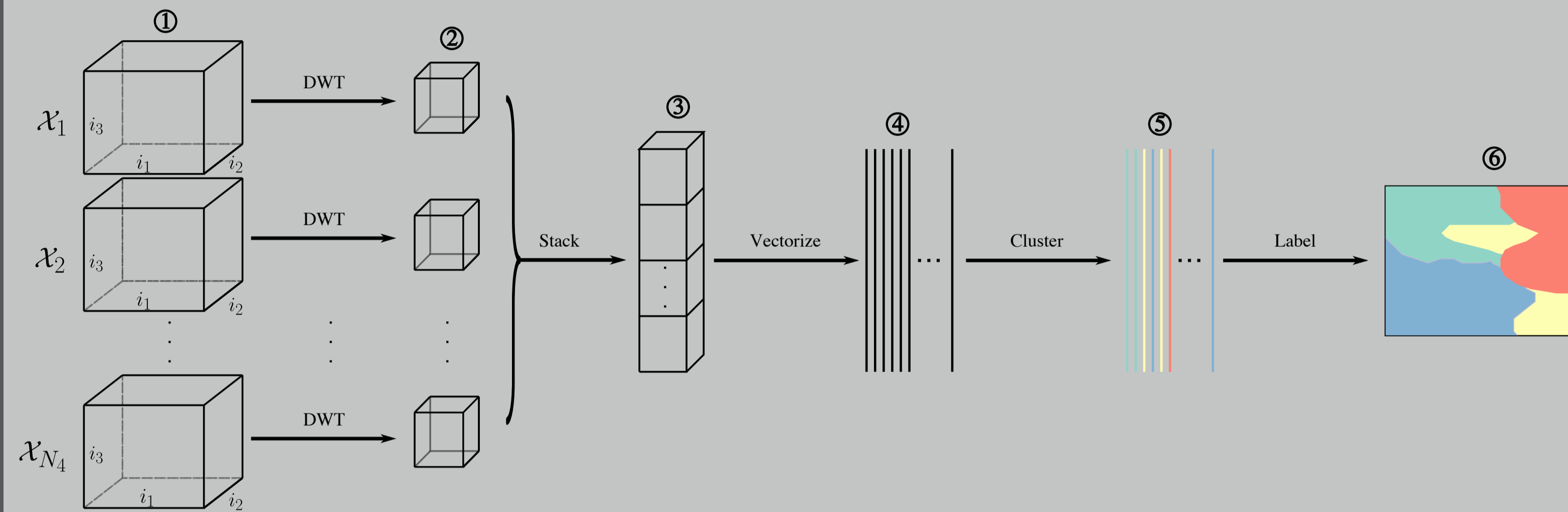
Tool - Discrete Wavelet Transform (DWT)

- ▶ The DWT splits a signal into high and low frequency.
- ▶ Low temporal signal captures climatology (seasons, years, decades), while low spatial signal captures regional features (city, county, state).
- ▶ Each application of DWT reduces size of data by factor of 2 (per dimension).
- ▶ DWT for tensors is computed via repeated 1D DWT on the low frequency signal along each axis.



Algorithm - Coarse Grain Clustering

1. Separate the tensor into 3-way tensors of a single variable (e.g. temperature data).
2. Pick integers l_s, l_t , and wavelets w_s for space and w_t for time and for each tensor of data, take l_s DWT in space, and l_t DWT in time.
3. Stack the new tensors for classification.
4. Vectorize the stack of wavelet coefficients.
5. Cluster the vectors using a chosen clustering algorithm (e.g. K-means).
6. Map the clustering on the DWT back to the initial latitude-longitude points.



Proof of Concept - LOCA Statistical Downscaled climate data

- ▶ Gridded climate data set of North America.
- ▶ Grid cell is monthly data from 1950-2013, six kilometers across.
- ▶ Available variables used: precipitation, maximum temperature, minimum temperature.
- ▶ Haar wavelets for space, db2 in time.
- ▶ Max spatial scale $l_s = 4$ (~ 100 km), max temporal scale $l_t = 6$ (~ 5 years).
- ▶ Clustering algorithm used: K-means.

Results - Effect of Coarse-Graining

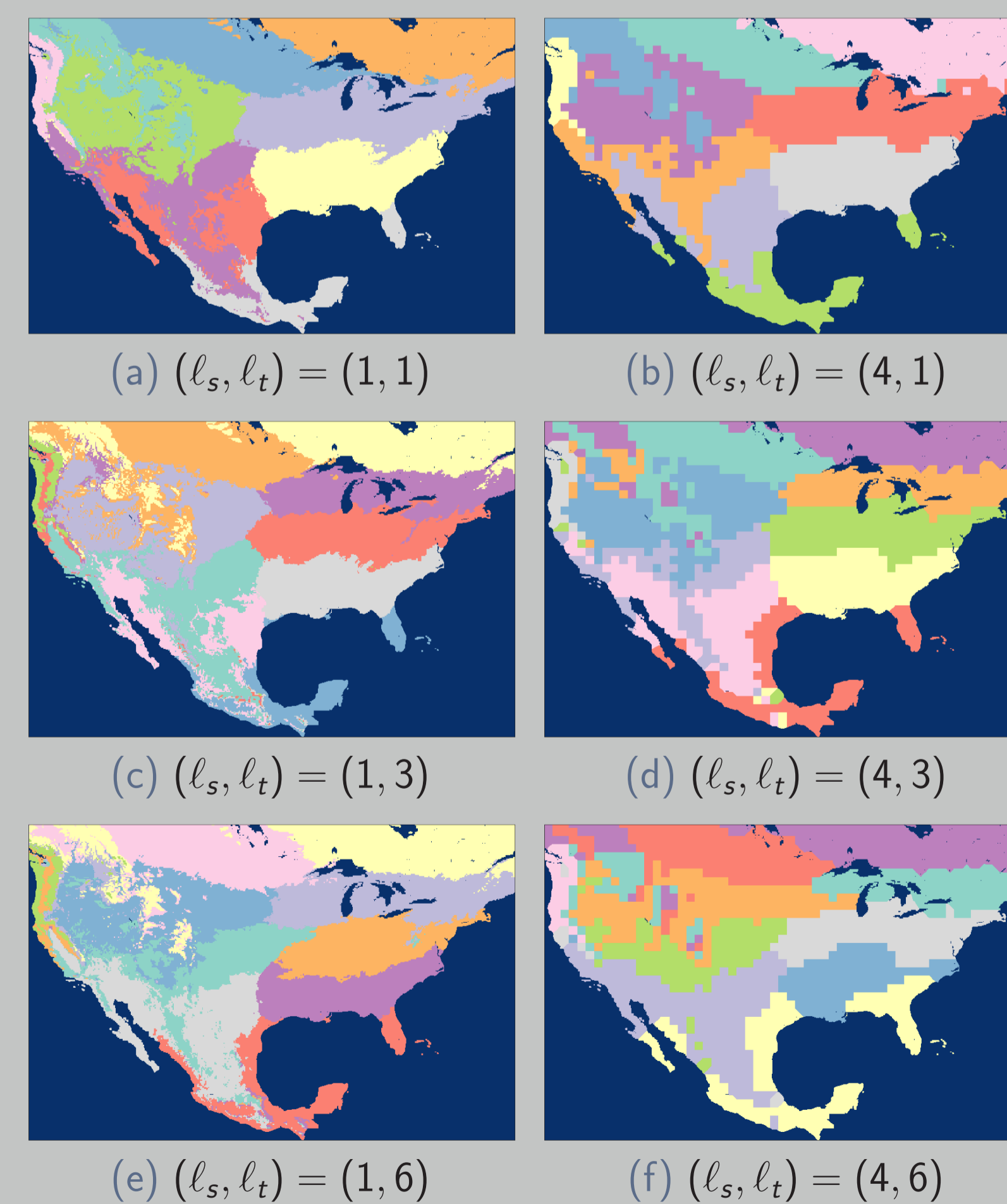


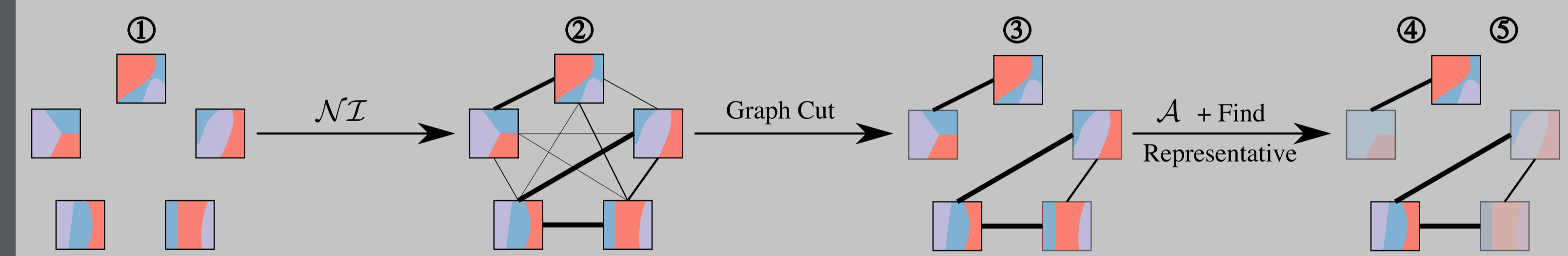
Figure: Clusterings obtained for $k = 10$ at resolution (l_s, l_t) .

Tool - Mutual Information

Given partitions of data $U = \{U_j\}_{j=1}^k, V = \{V_j\}_{j=1}^l$, the **Mutual Information** $\mathcal{NI}(U, V)$ measures how knowledge of one clustering reduces our uncertainty of the other.

Algorithm - Mutual Information Ensemble Reduce

1. Cluster across range of resolutions.
2. Compute pairwise mutual information - creating the mutual information graph.
3. Perform a graph cut to find the most similar connected components.
4. Compute average mutual information between each cluster and all other members of its component.
5. In each component, select cluster with highest expected mutual information.



Results - Key Resolutions

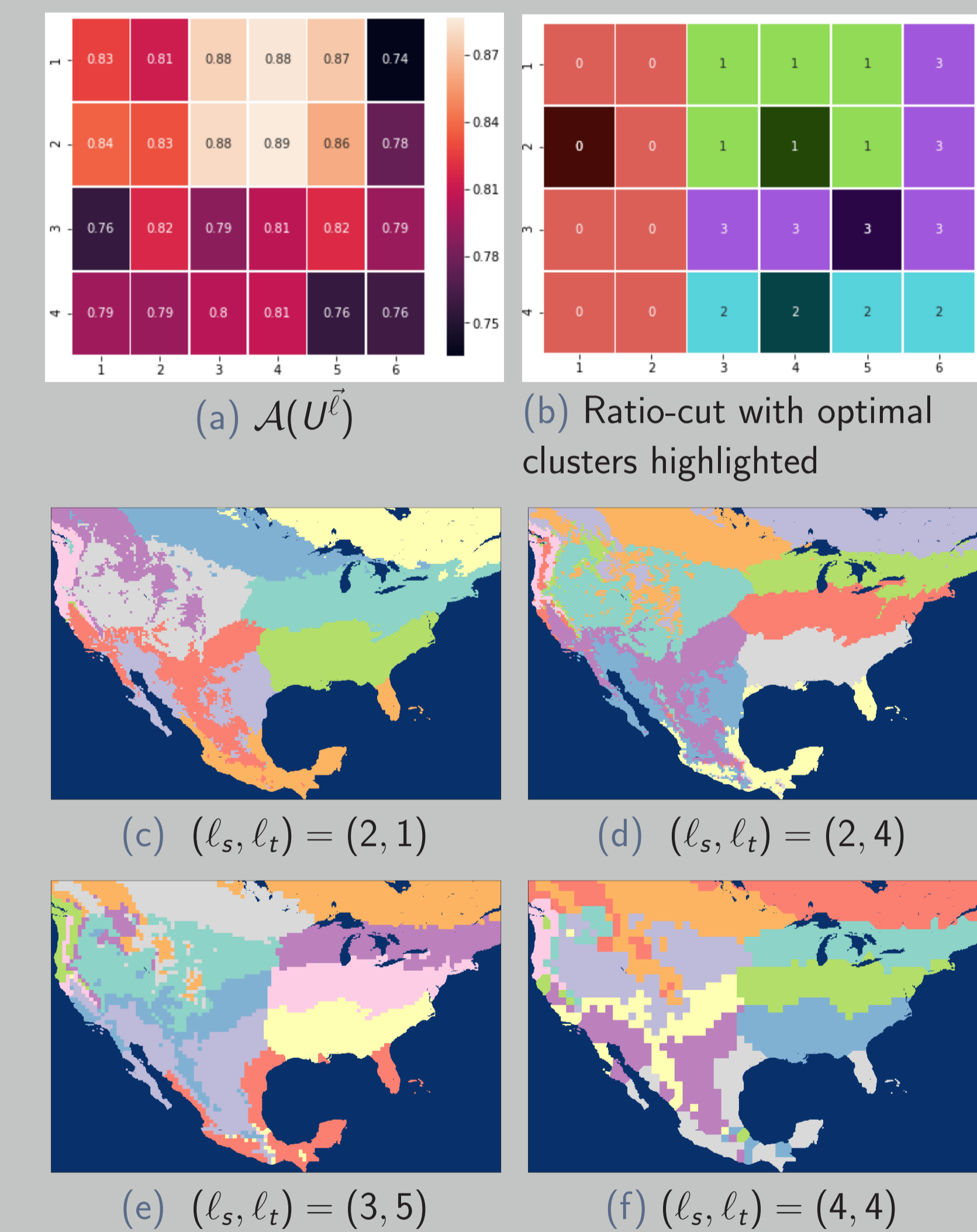


Figure: Output of the MIER algorithm for $k = 10$. a) $\mathcal{A}(U^{\vec{l}})$ for each resolution $\vec{l} \in \mathcal{L}$, b) results of the ratio cut algorithm and key resolutions, c-f) plots of key resolutions.

Discussion and Future Work

- ▶ Clustering the DWT illustrates additional structure within the data.
- ▶ Hence, an **optimal clustering** should be an ensemble of the wavelet clusterings to best characterize the impact of spatio-temporal variability in the climate data.