

# Overview

- **Introduction and problem statement**
- Case-study datasets
- Comparative assessment of clustering techniques
- Two-stage clustering
- Conclusion and future directions

# High-resolution building energy data

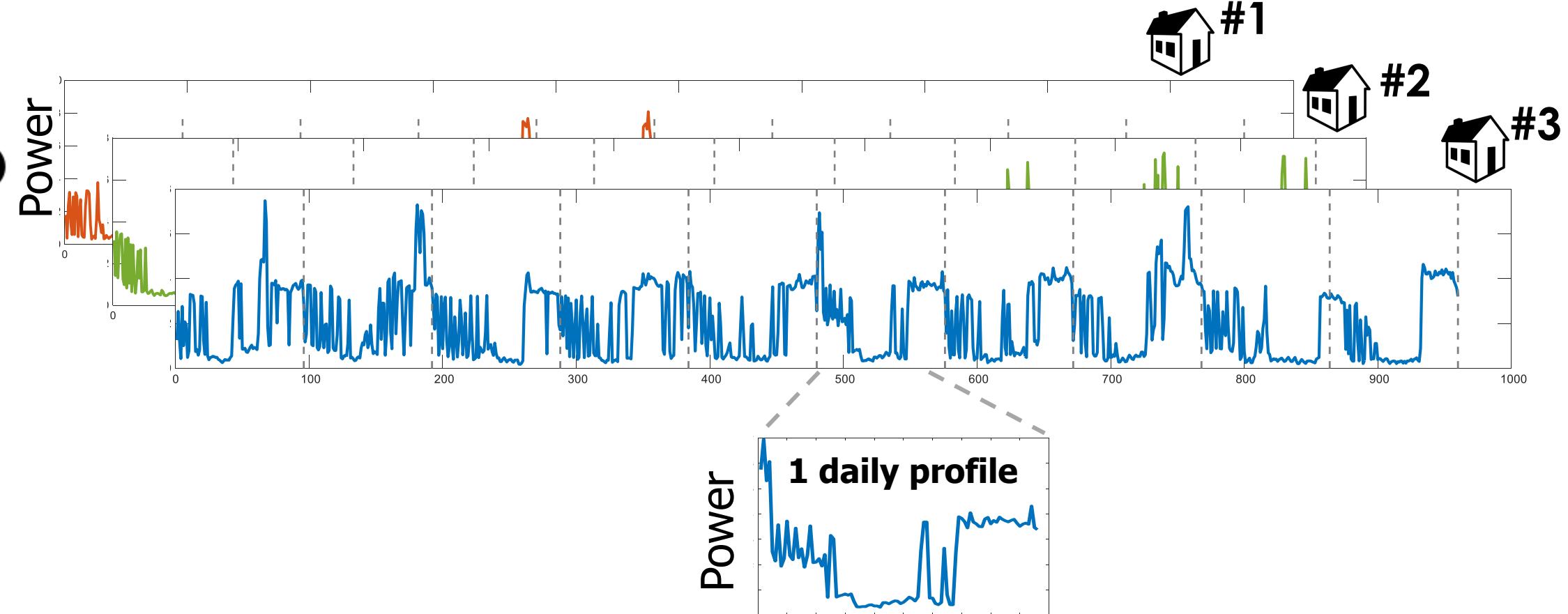


- Employing massive smart meter data promotes the efficiency of the power grid.
- Provides insight on consumption behaviors and lifestyles of the consumers at neighborhood scale.

# Complex and varied energy patterns

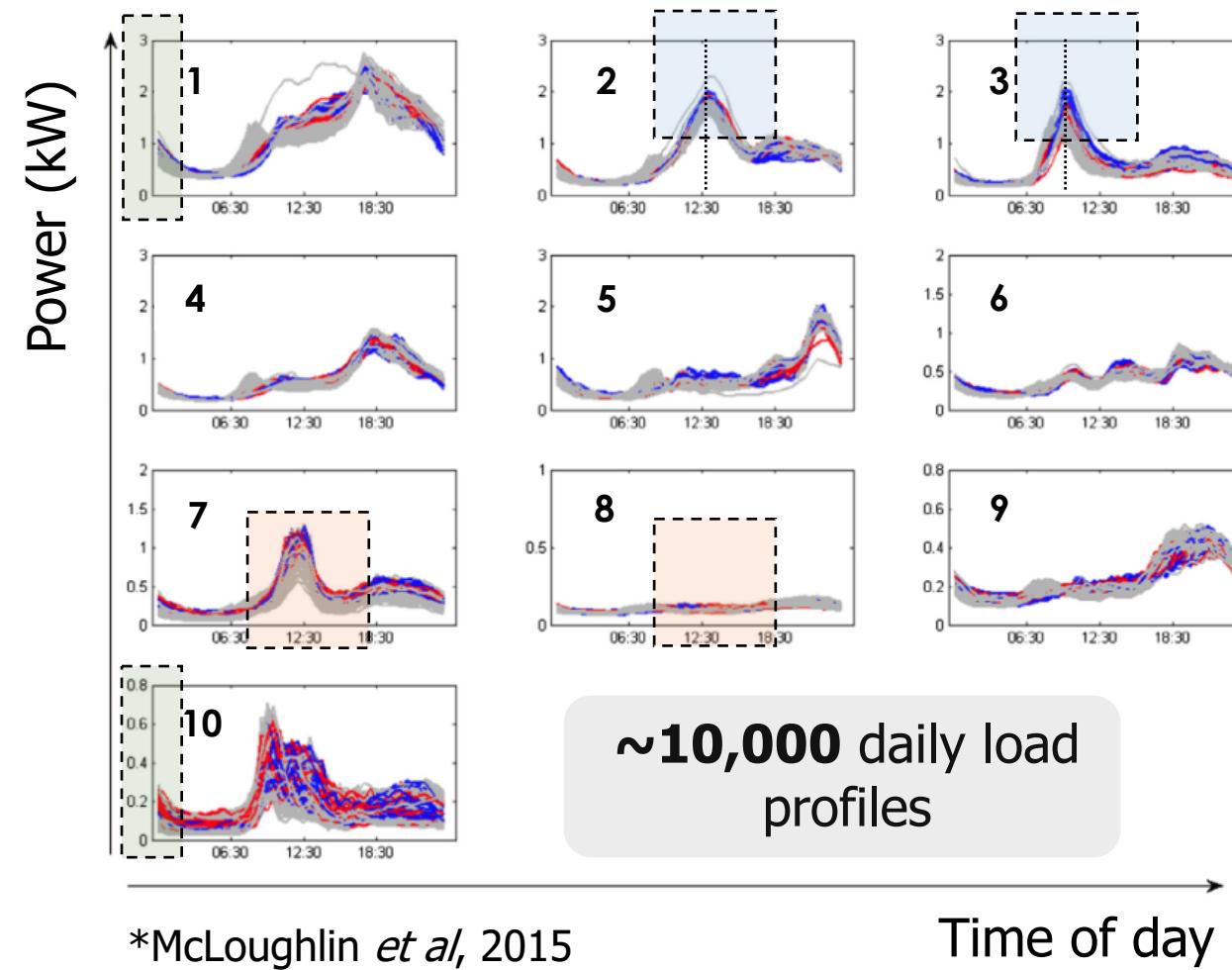


Time-series



- Decision-makers need the distribution of **typical consumption patterns**. However, daily load profiles show **highly volatile behaviors**.

# Daily profile segmentation (*Clustering*)

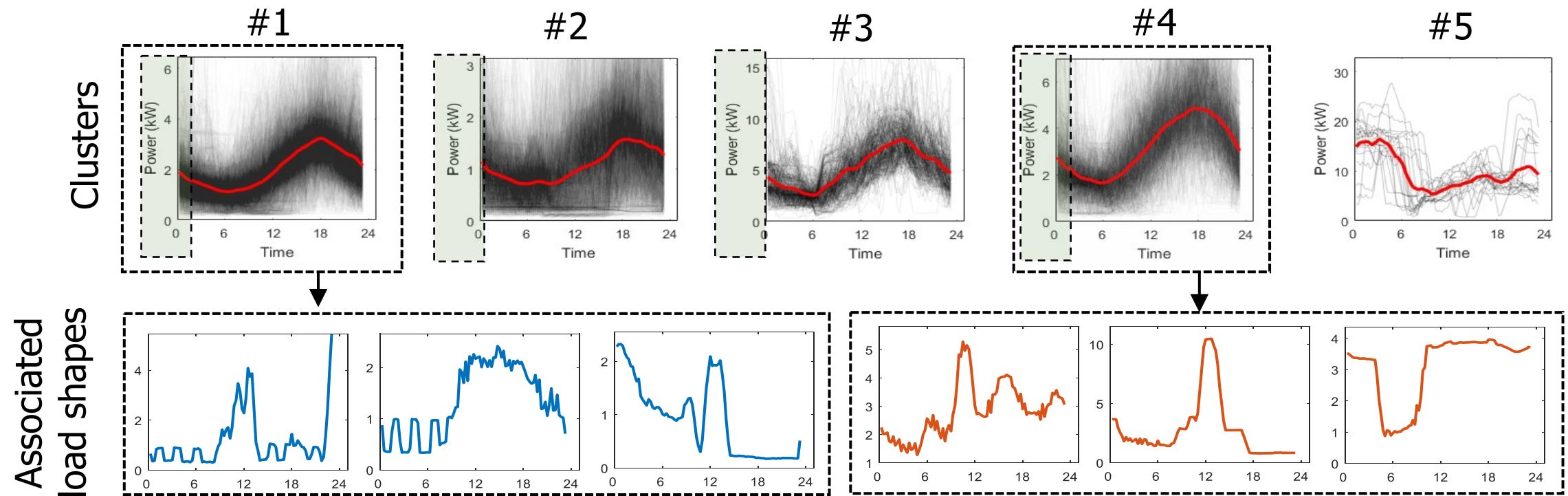


Important features in clustering:

1. Peak detection
2. Peak timing
3. Energy volume

# Problem statement

- Clustering based on validation indices can result in simplified patterns that do not represent the shape or energy volume of their associated load profiles.

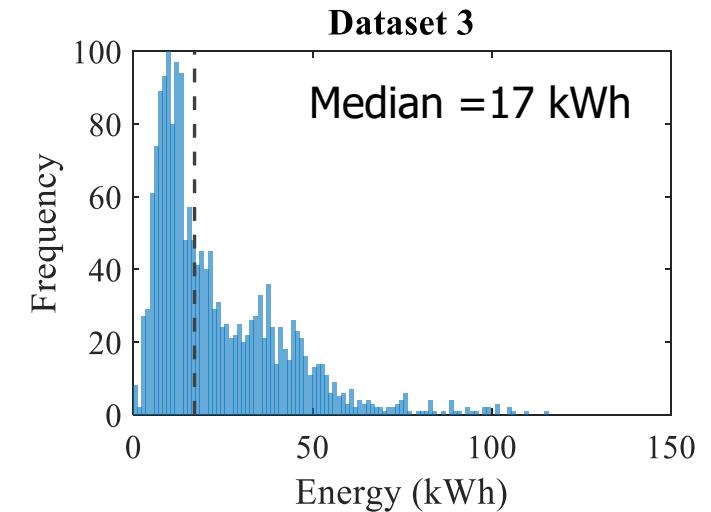
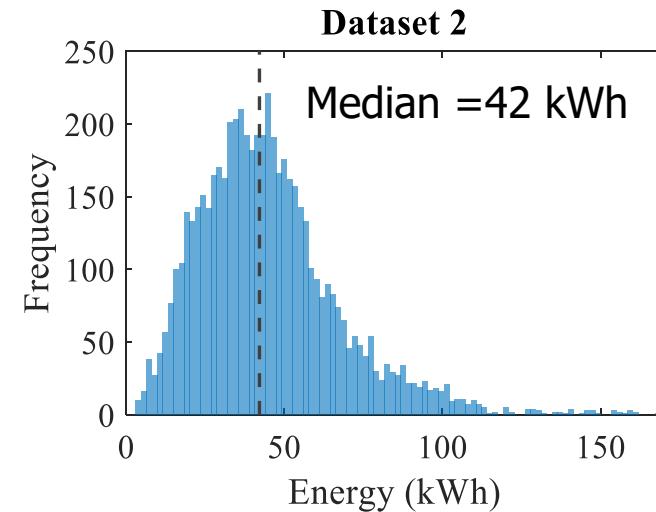
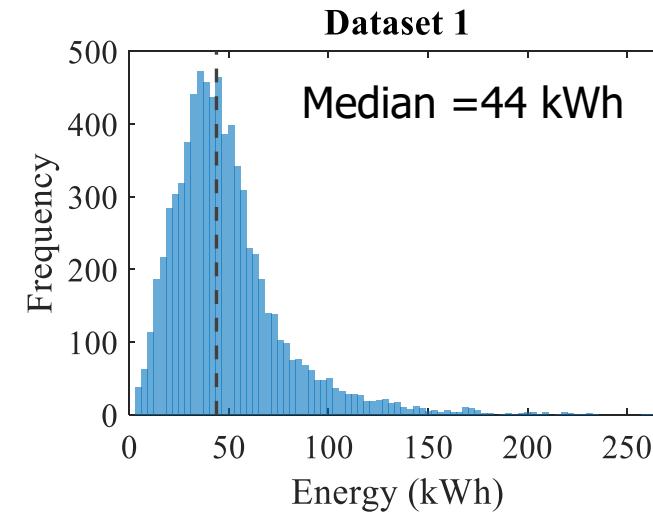


# Overview

- Introduction and problem statement
- **Case-study datasets**
- Comparative assessment of clustering techniques
- Two-stage clustering
- Conclusion and future directions

# Case-study datasets

| Dataset   | Location    | Number of buildings | Duration | # of daily profiles | # of samples per profile |
|-----------|-------------|---------------------|----------|---------------------|--------------------------|
| Dataset 1 | Austin, TX  | 129                 | 60 days  | 7535                | 96                       |
| Dataset 2 | Austin, TX  | 100                 | 60 days  | 5676                | 96                       |
| Dataset 3 | Boulder, CO | 31                  | 60 days  | 1790                | 96                       |

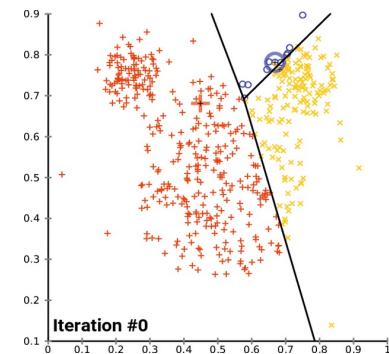


# Overview

- Introduction and problem statement
- Case-study datasets
- **Comparative assessment of clustering techniques**
- Two-stage clustering
- Conclusion and future directions

# Clustering techniques

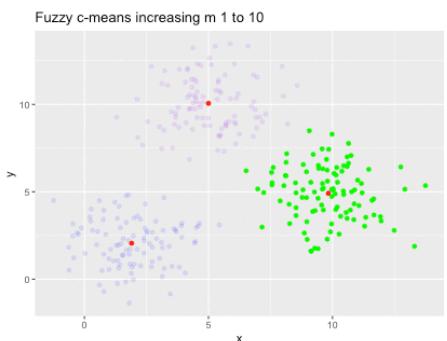
## K-means



- Distance measure-based

$$SSE = \sum_{k=1}^K \sum_{x \in C_k} \|x - c_k\|^2$$

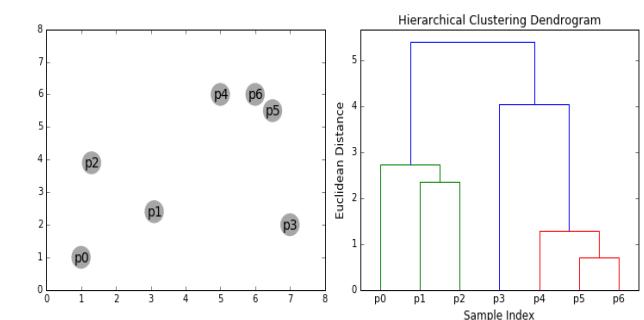
## Fuzzy c-means



- Degree of fuzziness

$$\sum_{i=1}^N \sum_{k=1}^K \mu_{ik}^m \|x_i - c_k\|^2$$

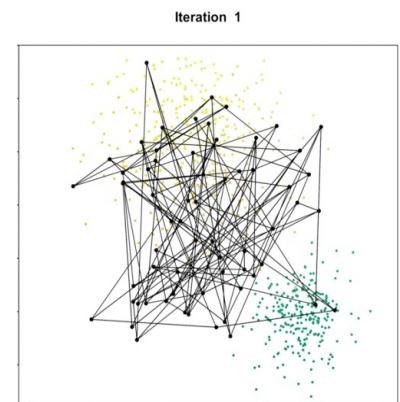
## Hierarchical clustering



- Tree structure

$$d_{ij} = \|X_i - X_j\|^2$$

## SOM



- Neural network-based

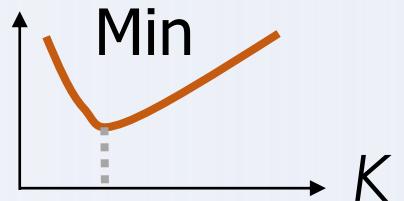
$$n_i(t+1) = n_i(t) + \alpha(t) \theta_{bi}(x(t) - n_i(t))$$

Ref for images: algobeans.com, wikipedia.com

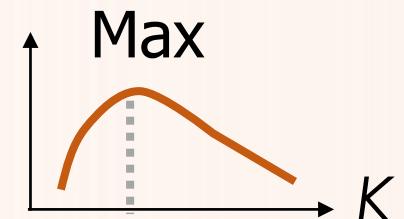
# Cluster validation index (CVI)

**Davies-Bouldin:**  $DBI = \frac{1}{K} \sum_i \max_{j, j \neq i} [\frac{1}{\|C_i\|} \sum_{x \in C_i} d(x, \mu_i) + \frac{1}{\|C_j\|} \sum_{x \in C_j} d(x, \mu_j)]$

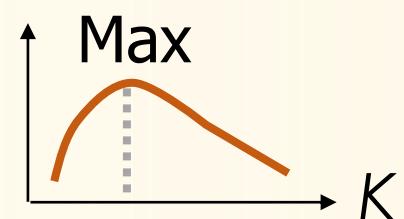
Selection criteria



**Silhouette:**  $SIL = \frac{1}{K} \sum_i \left\{ \frac{1}{\|C_i\|} \sum_{x \in C_i} \frac{b(i) - a(i)}{\max(b(i), a(i))} \right\}$

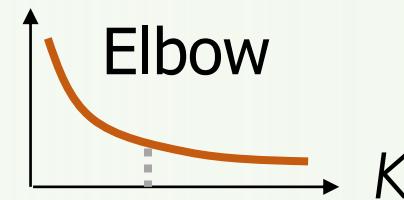


**Calinski-Harabasz:**  $CHI = \frac{\sum_i \|C_i\| * d^2(\mu_i, \mu_j) / (K - 1)}{\sum_i \sum_{x \in C_i} d^2(x, \mu_j) / (N - K)}$



**Within cluster sum of square error:**

$$WCSSE = \sum_i \sum_{x \in C_i} \|x - \mu_i\|^2$$

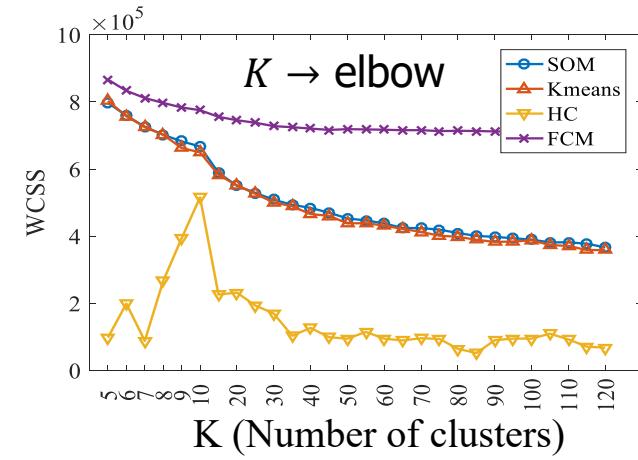
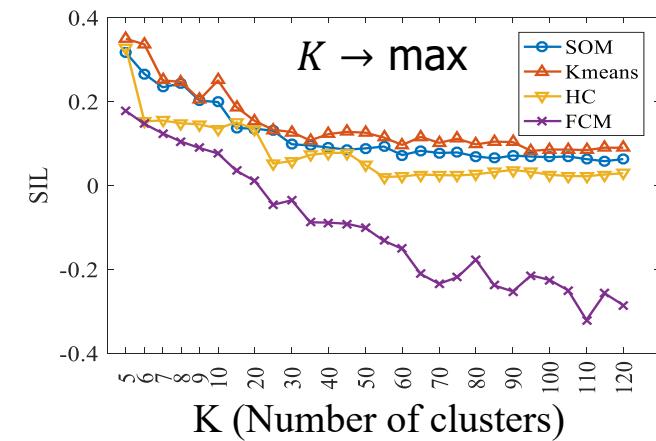
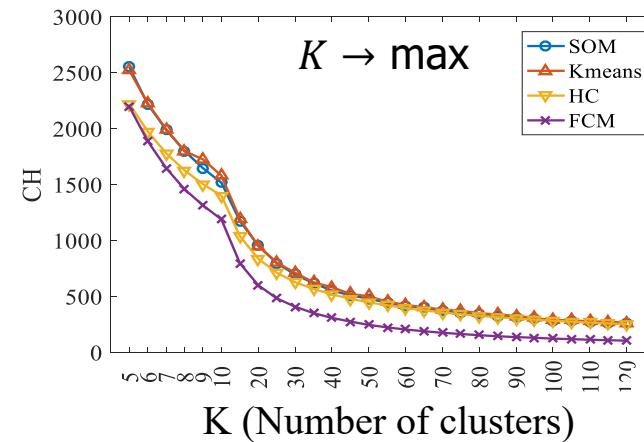
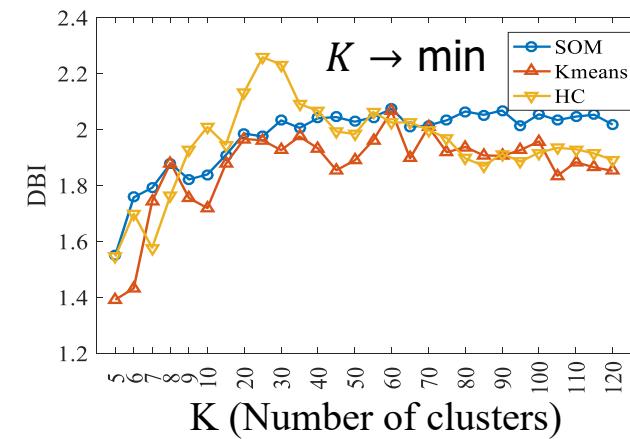


# Comparative assessment

- Different combinations of Clustering techniques and Validation Index for evaluation

| Clustering technique | Validation index |             |                   |      |   |
|----------------------|------------------|-------------|-------------------|------|---|
|                      | Davies-Bouldin   | Sillhouette | Calinski-Harabasz | WCSS |   |
| K-means              | ?                | ?           | ?                 | ?    | ? |
| Hierarchical         | ?                | ?           | ?                 | ?    | ? |
| Fuzzy c-means        | ?                | ?           | ?                 | ?    | ? |
| SOM                  | ?                | ?           | ?                 | ?    | ? |

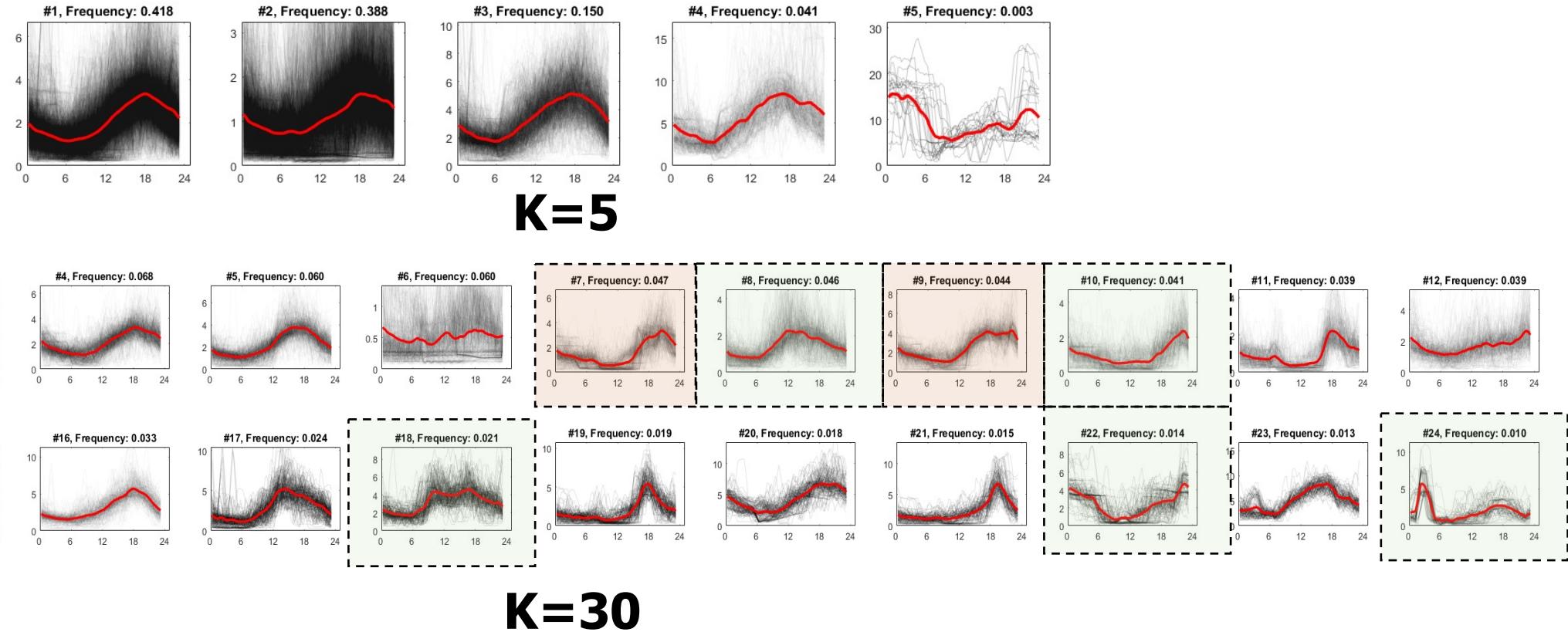
# CVI results



- Three out of four metrics estimate a low number of clusters (**K=5**).
- Only WCSS selects a higher number (~**K=30**).

Compared two scenarios  
**K=5** versus **K=30**

# Empirical investigation



Going with higher number reveals distinct patterns.

However, it results in some similar (redundant) clusters.

# Comparative assessment

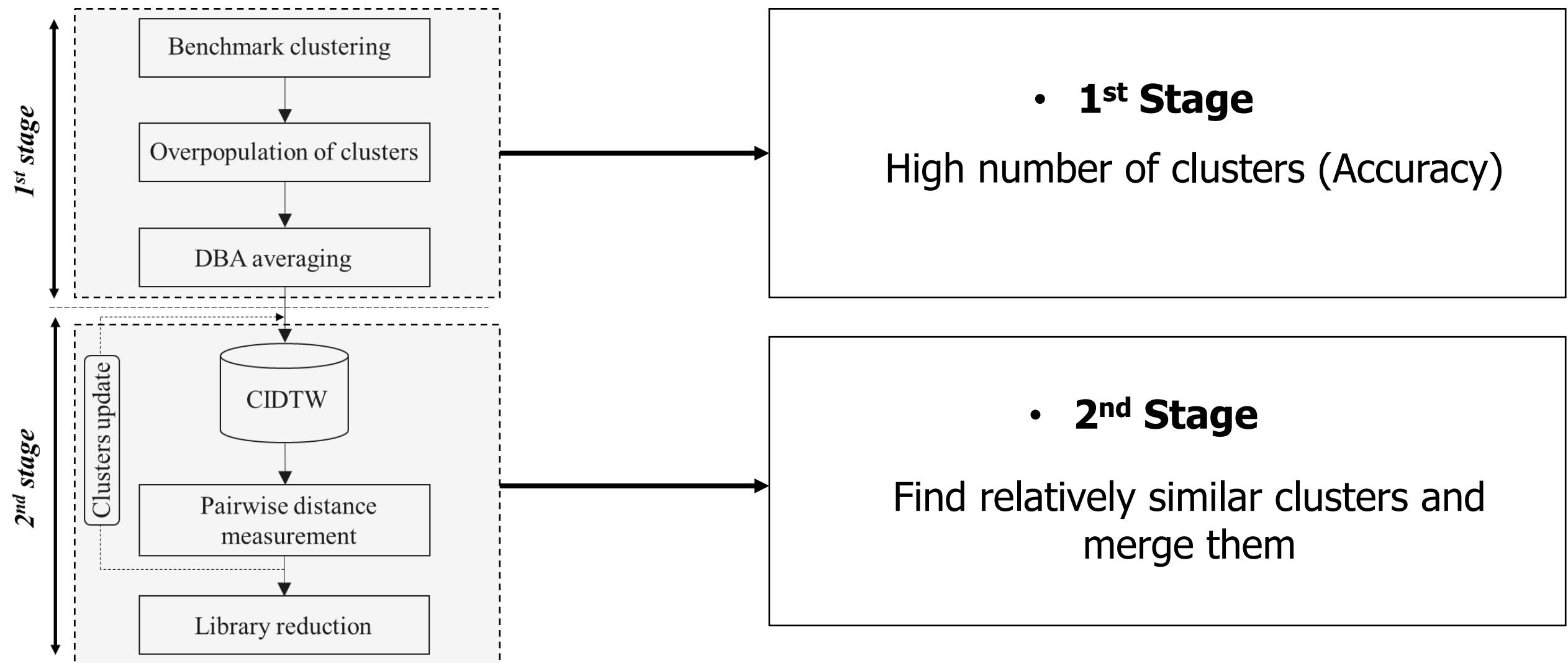
|               | Validation index |             |                   |      |
|---------------|------------------|-------------|-------------------|------|
|               | Davies-Bouldin   | Sillhouette | Calinski-Harabasz | WCSS |
| K-means       | ✓                | ✓           | ✓                 | ✓    |
| Hierarchical  | ✓                | ✓           | ✓                 | ✓    |
| Fuzzy c-means | ✓                | ✓           | ✓                 | ✓    |
| SOM           | ✓                | ✓           | ✓                 | ✓    |

**Limitation:** Low number of clusters can't retrieve realistic energy profiles. High number of clusters could solve this problem but results in similar (redundant) clusters.

# Overview

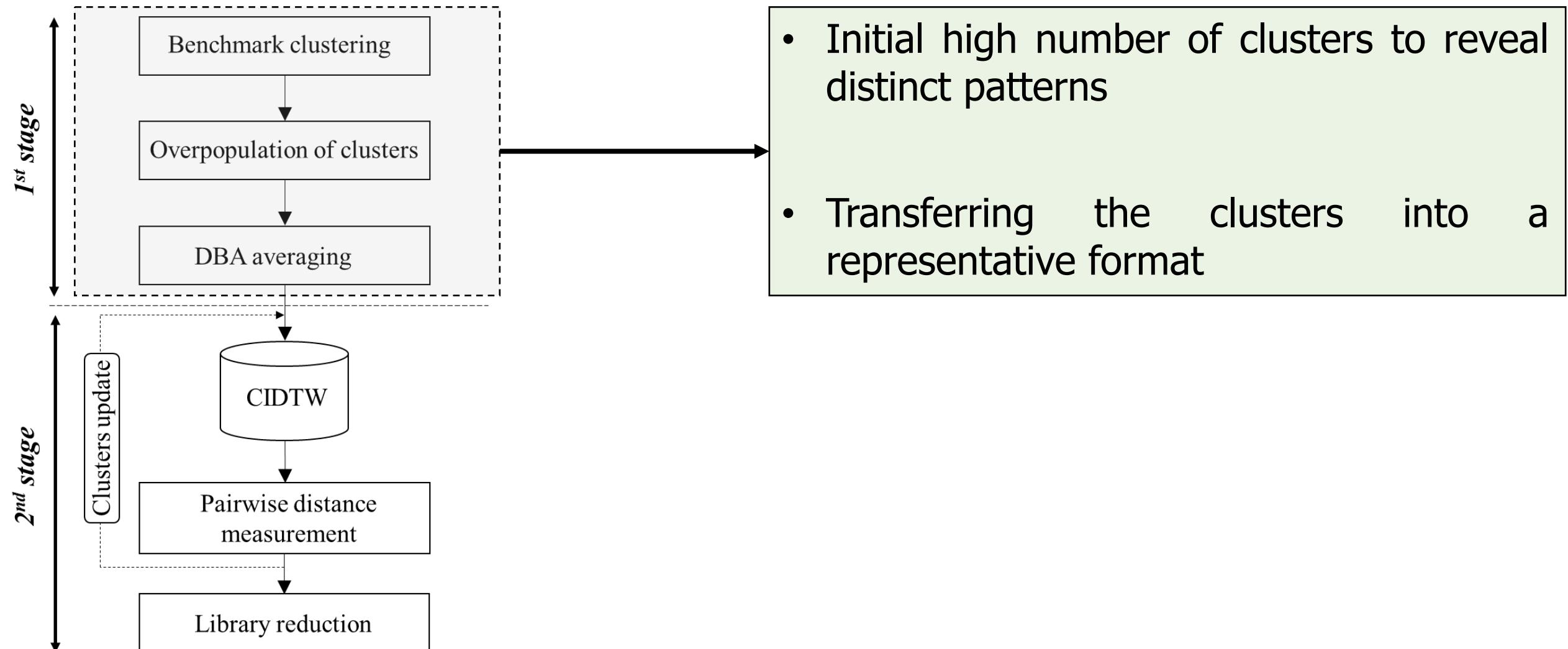
- Introduction and problem statement
- Case-study datasets
- Comparative assessment of clustering techniques
- **Two-stage clustering**
- Conclusion and future directions

# Proposed two-stage clustering



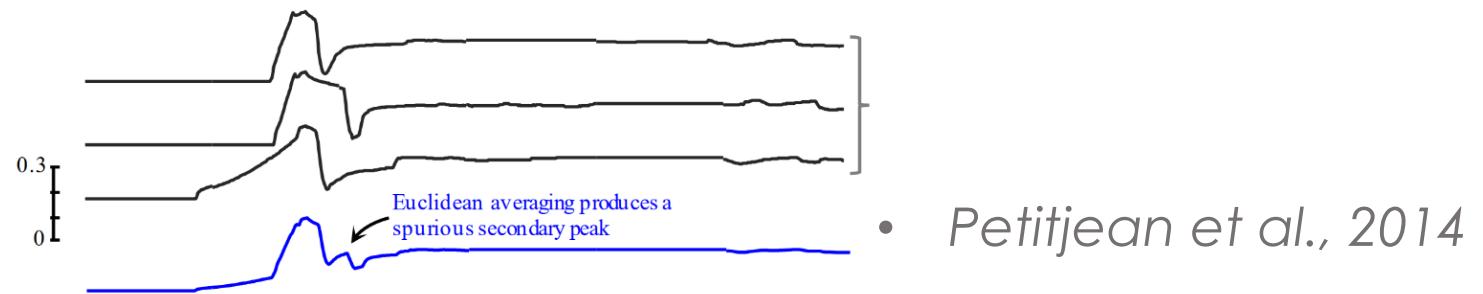
# Proposed two-stage clustering

***1<sup>st</sup> Stage***



# DBA averaging

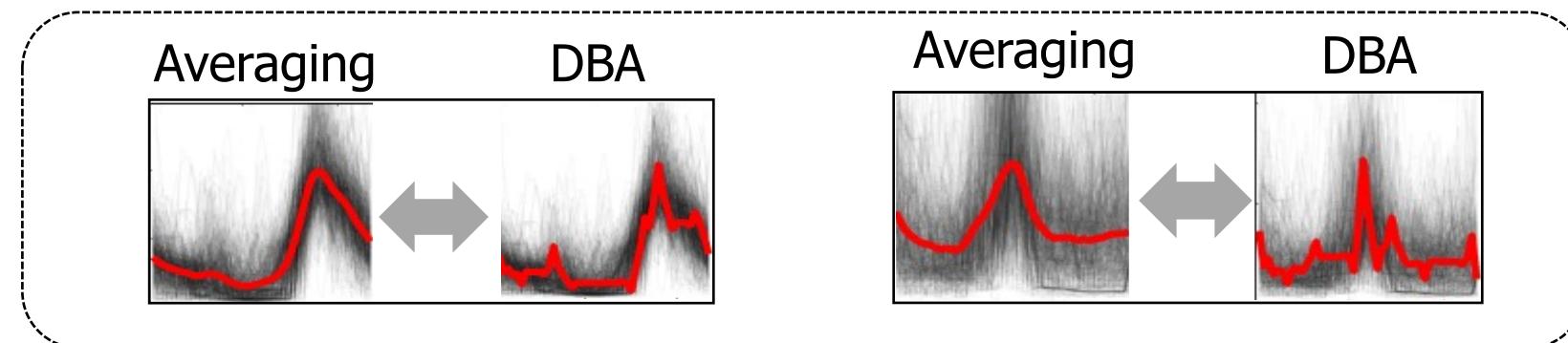
- Conventional averaging problem: Centroid is different from the profiles.



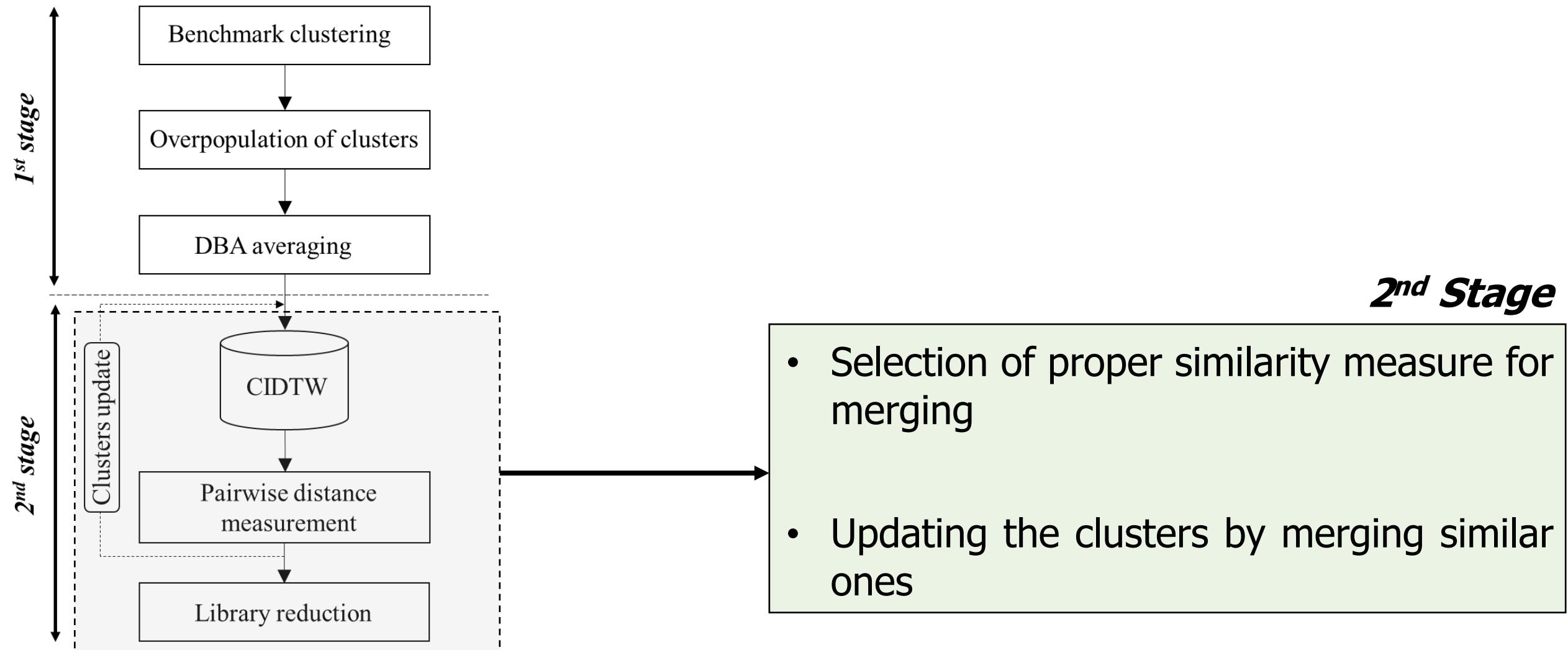
**DBA:** Using an expectation maximization approach to refine the medoid of each group through finding the best set of alignments through iterations.

$$C'_i(t) = \text{barycenter}(\text{assoc}(C_i(t)))$$

$$\text{barycenter}(X_1, X_2, \dots, X_m) = \frac{X_1 + X_2 + \dots + X_m}{m}$$

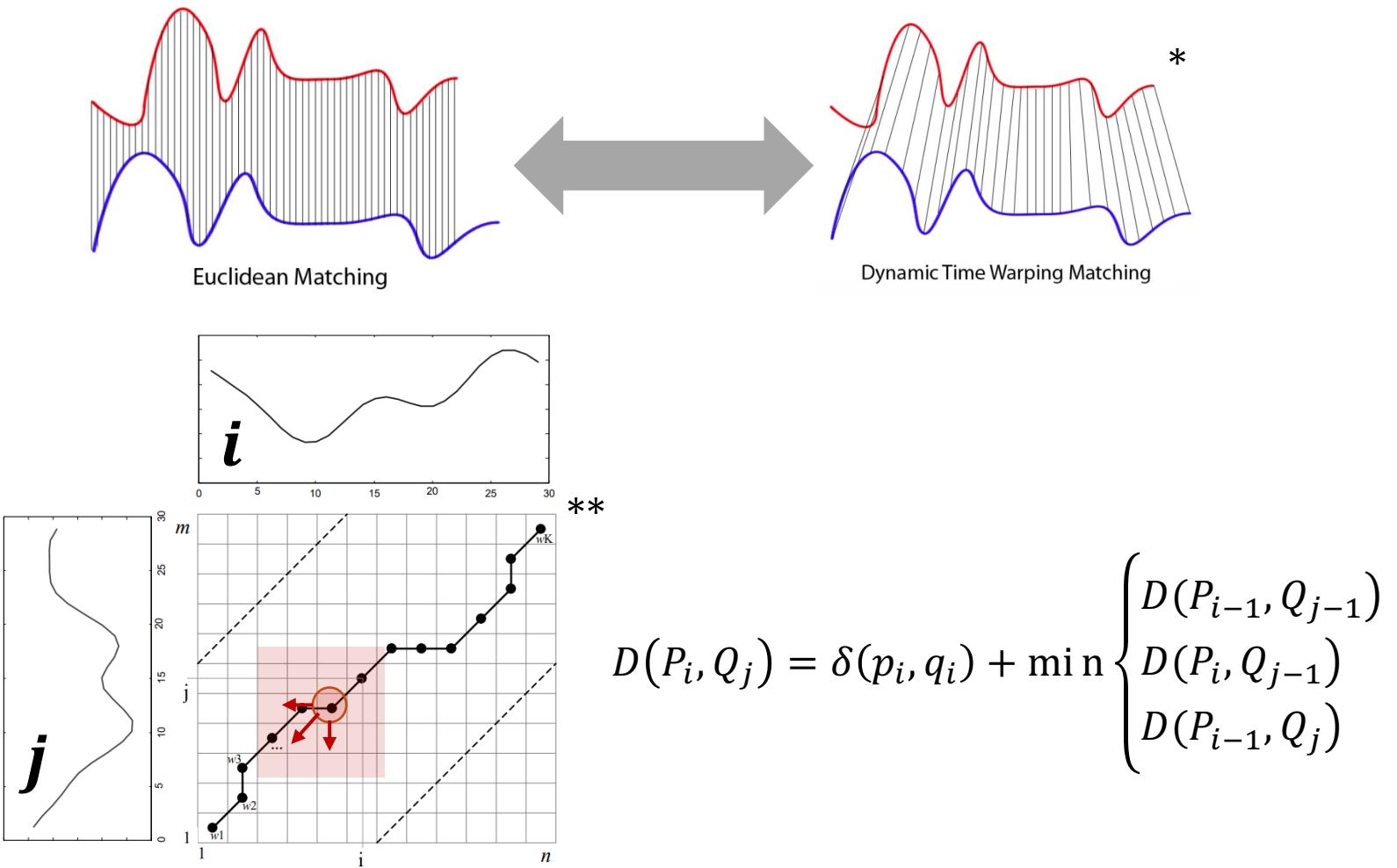


# Proposed two-stage clustering



# CI-DTW (*complexity invariant dynamic time warping*)

Using a dynamic programming-based distance measure (CI-DTW) to find the similarity between pairs of clusters.



\* <https://towardsdatascience.com/>

\*\* Keogh and Pazzani, 2001



# Cluster merging

$K'$  initial clusters at the first stage is transferred into  $K$  final clusters. ( $K < K'$ ).

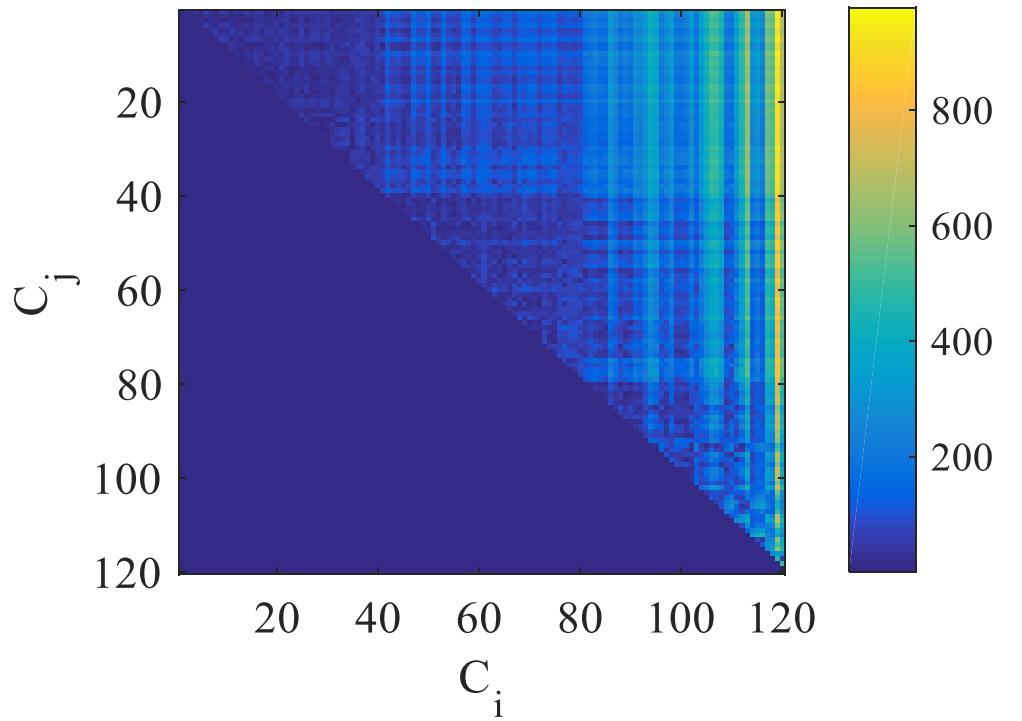
---

**Algorithm 1.** Merging cluster centroids
 

---

**Input:** Overpopulated clustering results, cluster centroids with DBA, initial cluster number ( $K'$ )

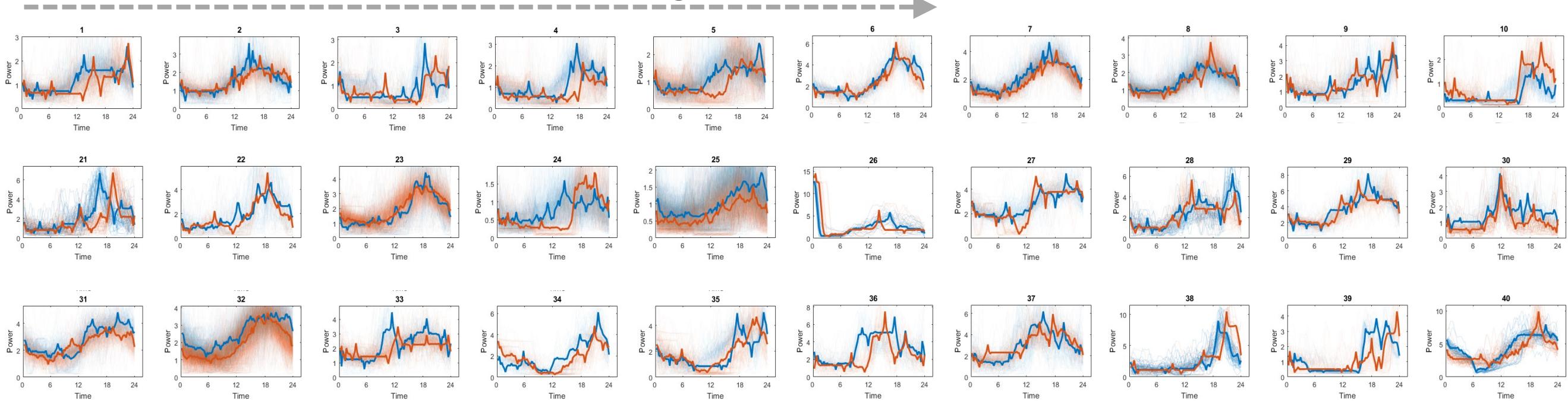
- 1: Set the target cluster number  $K$ .
  - 2: While  $K' > K$ :
  - 3:     Find the closest cluster centroids  $C_i$  and  $C_j$  based on CI-DTW metric. →
  - 4:     While  $\|C_i\| + \|C_j\| > \tau * n$ :
  - 5:         Find the next set of closest  $C_i$  and  $C_j$ .
  - 6:         Set  $C_i = (n_i C_i + n_j C_j) / (n_i + n_j)$ .
  - 7:         Delete  $C_j$ .
  - 8:         Update cluster index from cluster  $j + 1$  to the last one.
  - 9:      $K' = K' - 1$
- 



Pairwise distance between all pairs of clusters ( $i$  and  $j$ ) at one iteration.

# Cluster merging results

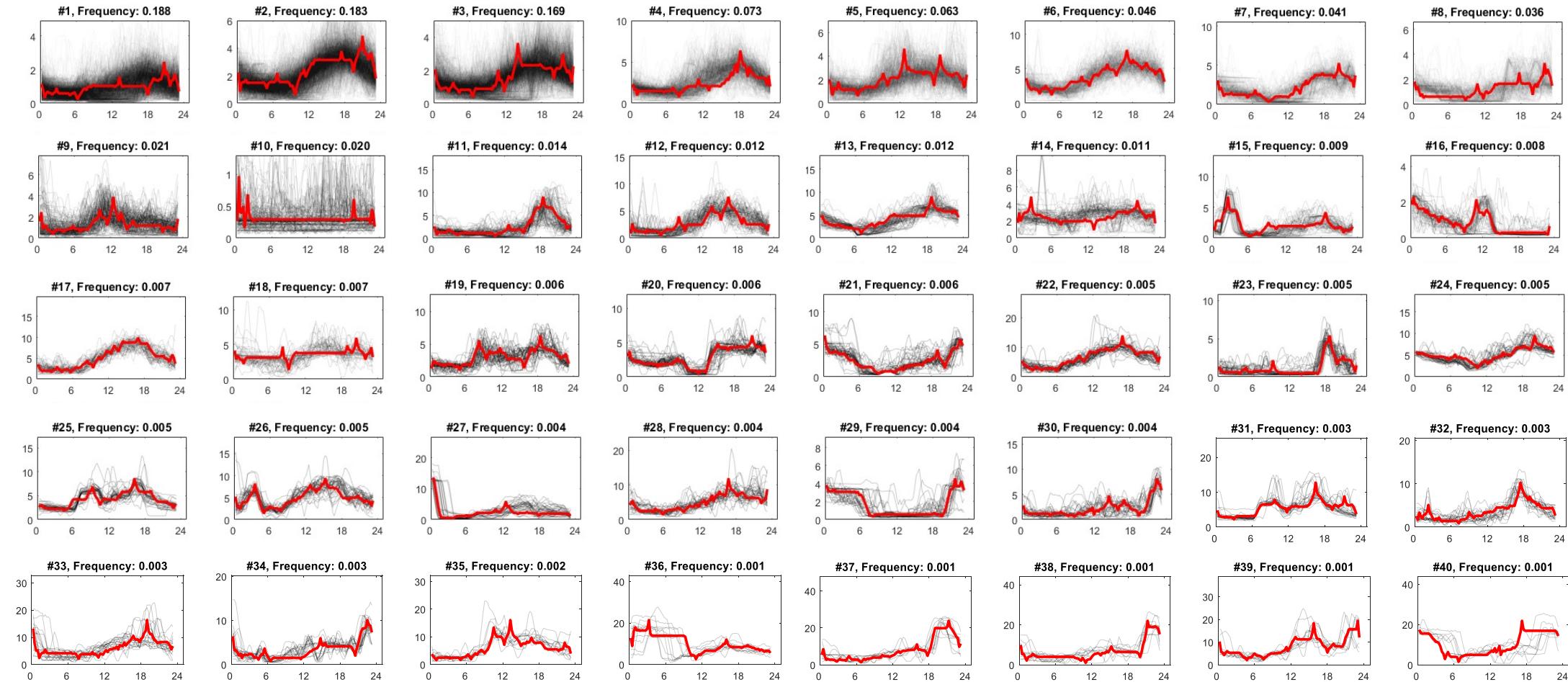
*Higher iteration*



Dataset1:  $K'=90$  and  $K=40$ . First 40 iteration is shown in the plot.

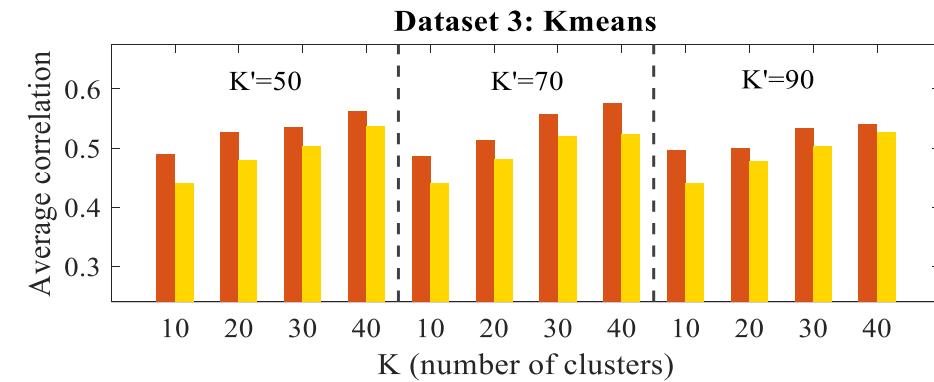
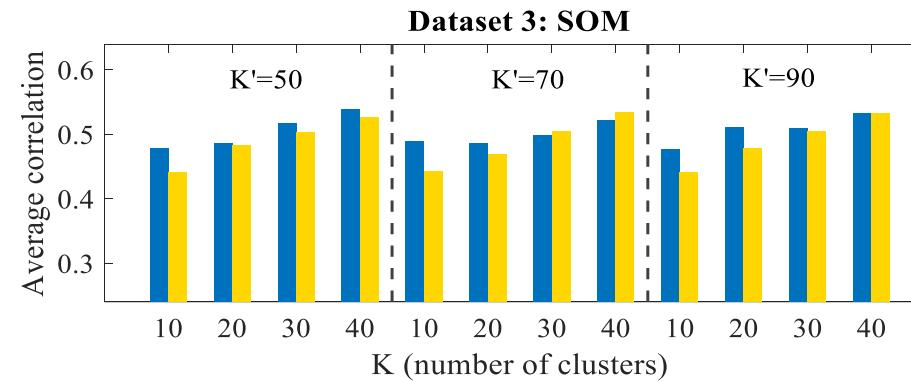
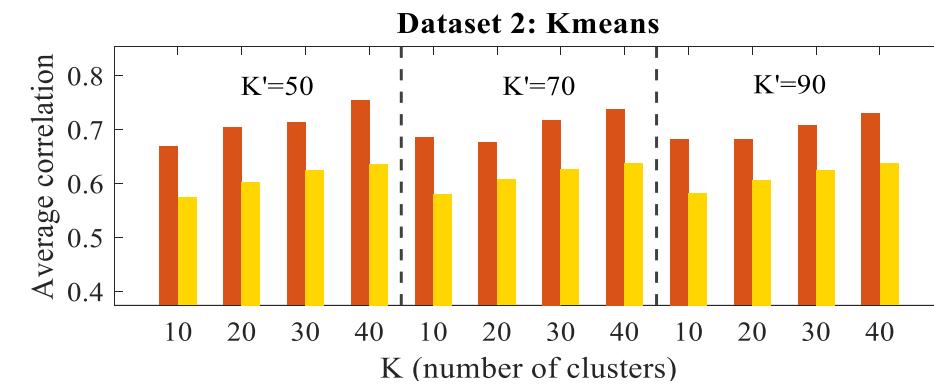
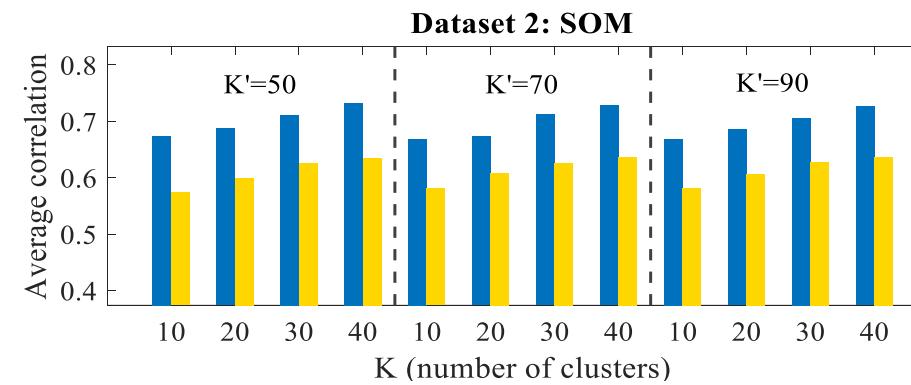
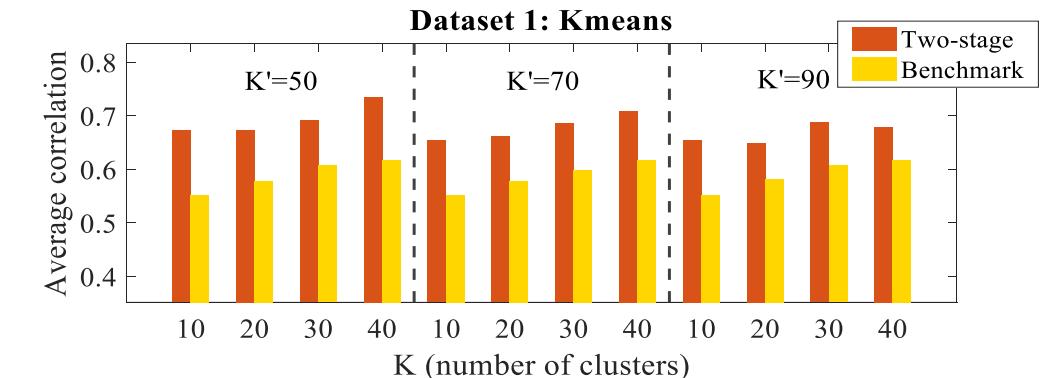
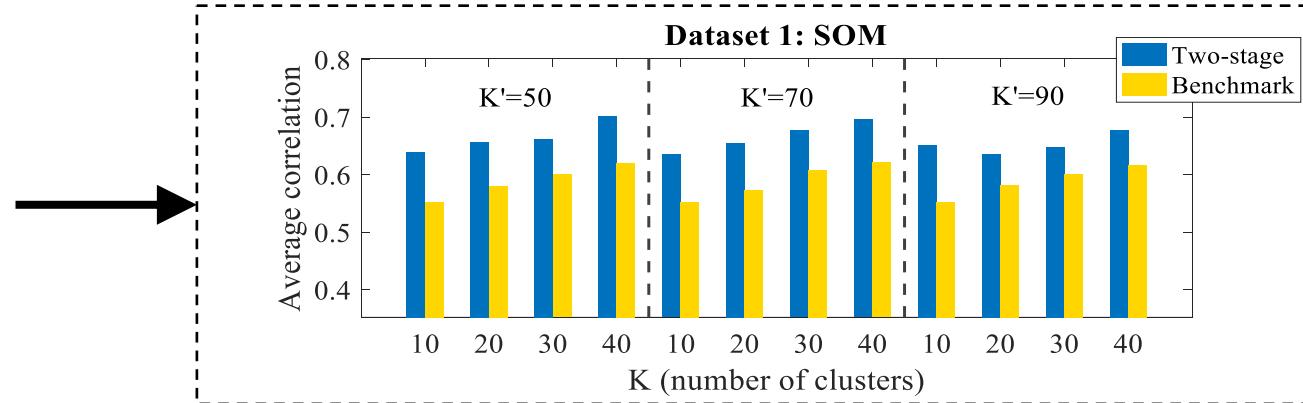
Merged clusters are subjectively close both in temporal patterns and peak magnitudes.

# Two-stage clustering results



- Well-separated clusters
- Important features (peak detection, peak timing, and energy volume) are accentuated

# Quantified investigation (*Correlation*): higher is better



# Quantified investigation (*WCSS*): lower is better

