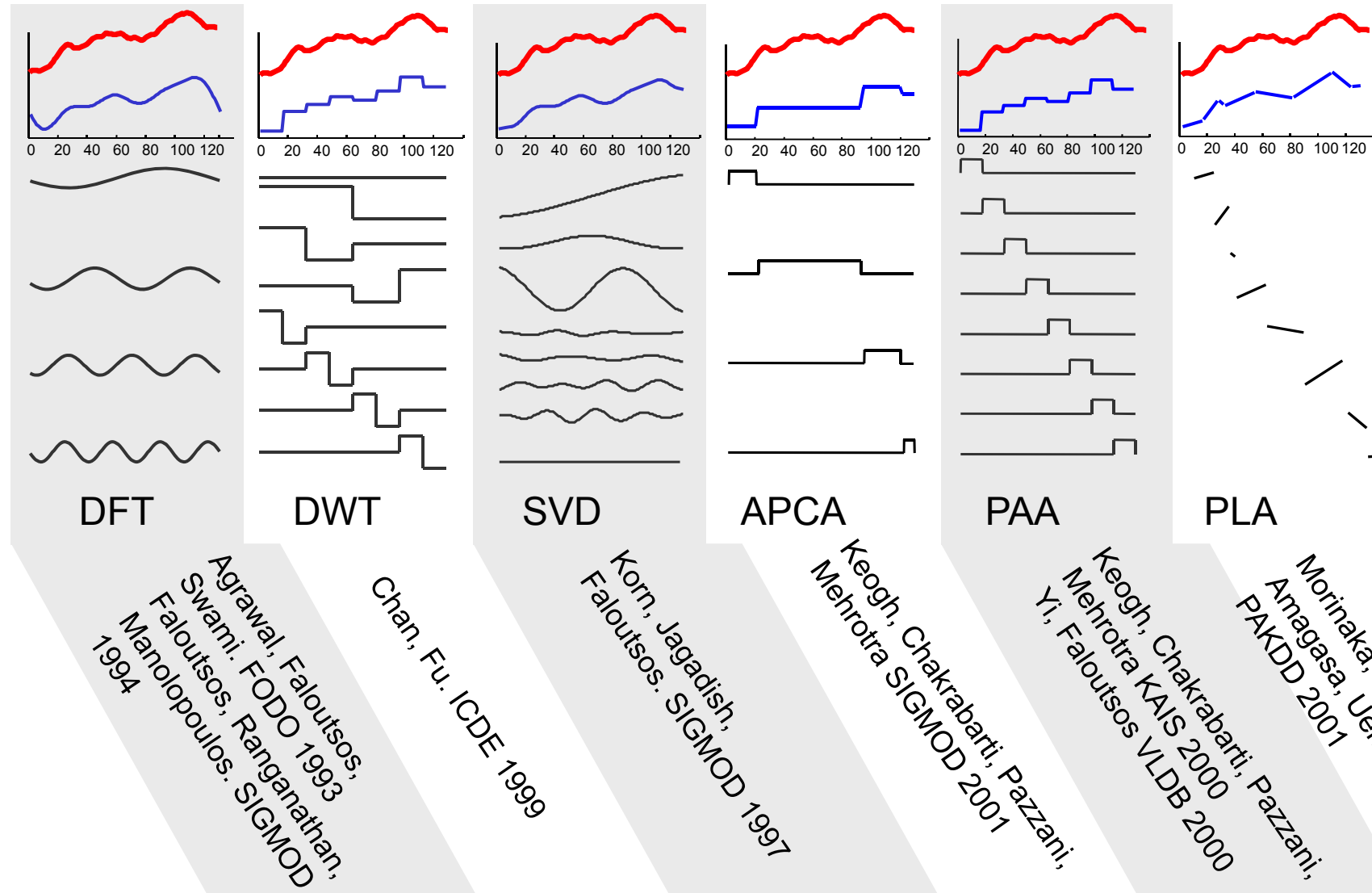
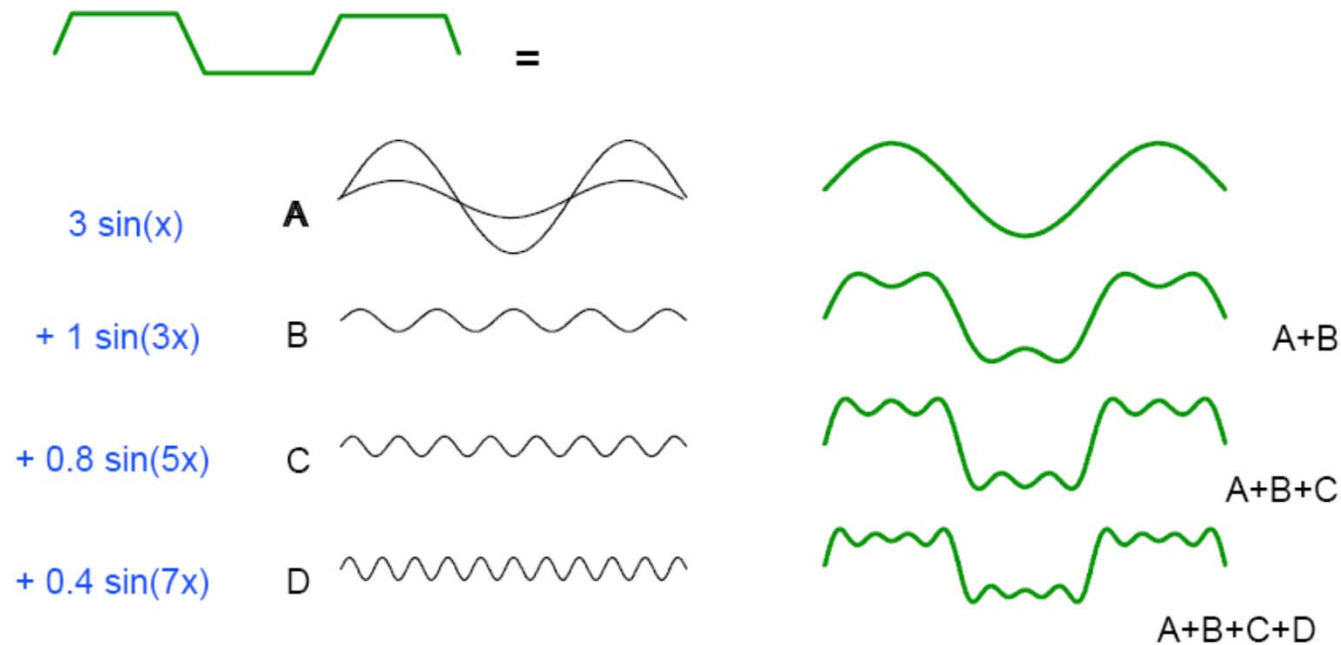


- Instead of directly working with the entire time series, we can also extract features from them
- Many feature extraction techniques exist that basically follow two different purposes:
 - Many of them aim at representing time series in a compact way (e.g. as a “shorter” approximation of the original time series) with minimum loss of modelling error
 - => this is mostly done for performance considerations
 - => approach is closely related to dimensionality reduction/feature selection
 - Examples covered here: DFT, DTW, SVD, APCA, PAA, PLA
 - Other model specific properties of the time series relevant to a given application
 - Example covered here: threshold-based modelling

Compact Representations: Overview



- Discrete Fourier Transformation (DFT)
 - Idea:
 - Describe a periodic function as a weighted sum of periodical base functions with varying frequency
 - Base functions here: sin und cos
 - Example:



– Basic foundation: [Fouriers Theorem](#)

Any periodic function can be represented by a sum of sin- and cos-functions of different frequency

- DFT does not „change“ the function but simply finds a different equivalent representation (and DFT can be reversed)
- Formally:
 - Let $x = [x_t]$, $t = 0, \dots, n - 1$ be a time series of length n
 - DFT transforms x into $X = [X_f]$ of n complex numbers with frequencies $f = 0, \dots, n - 1$ such that

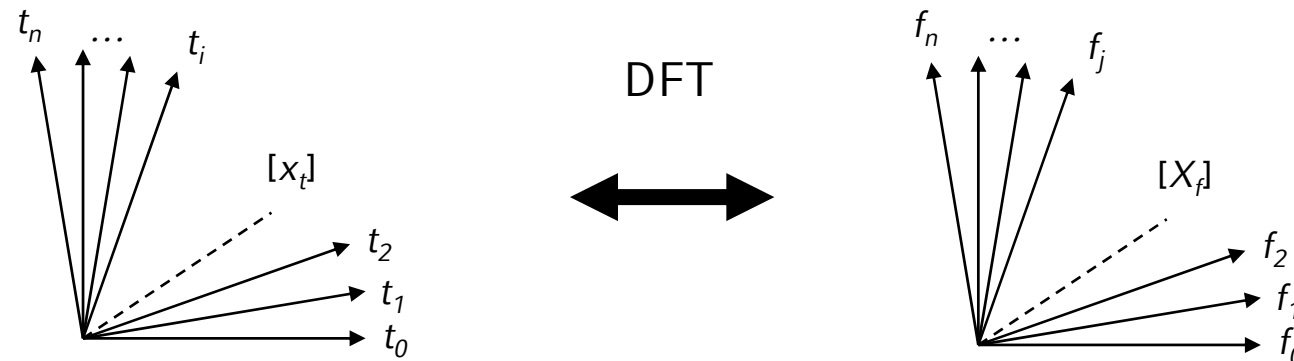
$$X_f = \frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} x_t \cdot e^{\frac{-j2\pi ft}{n}} =$$

$$\underbrace{\frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} x_t \cos\left(\frac{2\pi ft}{n}\right)}_{\text{Realteil}} - j \cdot \underbrace{\frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} x_t \sin\left(\frac{2\pi ft}{n}\right)}_{\text{Imaginärteil}}$$

where $j^2 = -1$.

- » Realteil is the portion of cosine in frequency f
- » Imaginärteil is the portion of sinus in frequency f

- DFT can be interpreted as a transformation of the basis vectors (like e.g. PCA):



- The new axis represent the frequencies

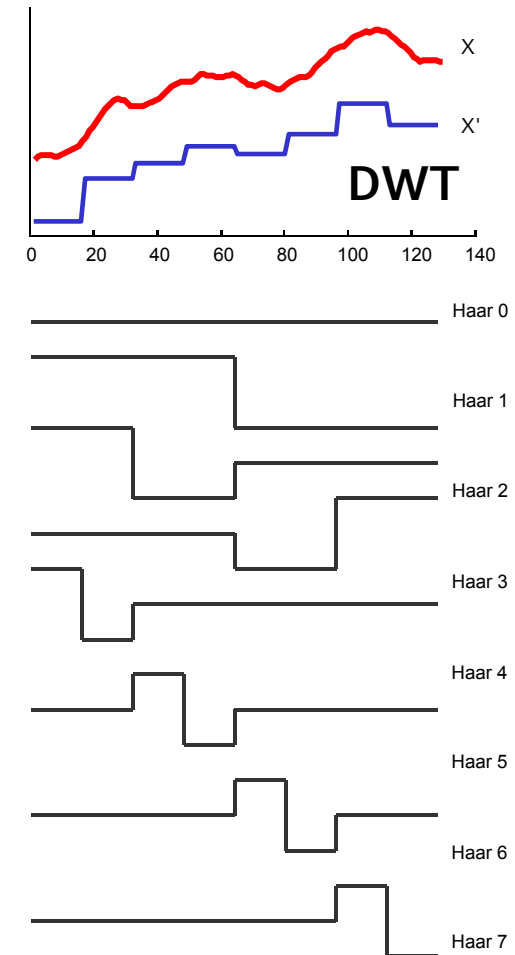
- But how does that help?
 - So far, we transformed an n -dimensional time series into an n -dimensional vector ...
- Well first of all, it holds that the euclidean distance is preserved after DFT, i.e. $||x - y||^2 = ||X - Y||^2$

- This follows from Parseval's theorem (and the linearity of DFT) which states that the energy of a sequence (= sum of squared amplitudes $E(x) = \|x\|^2 = \sum_{t=0}^{n-1} |x_t|^2$) is preserved, i.e.:

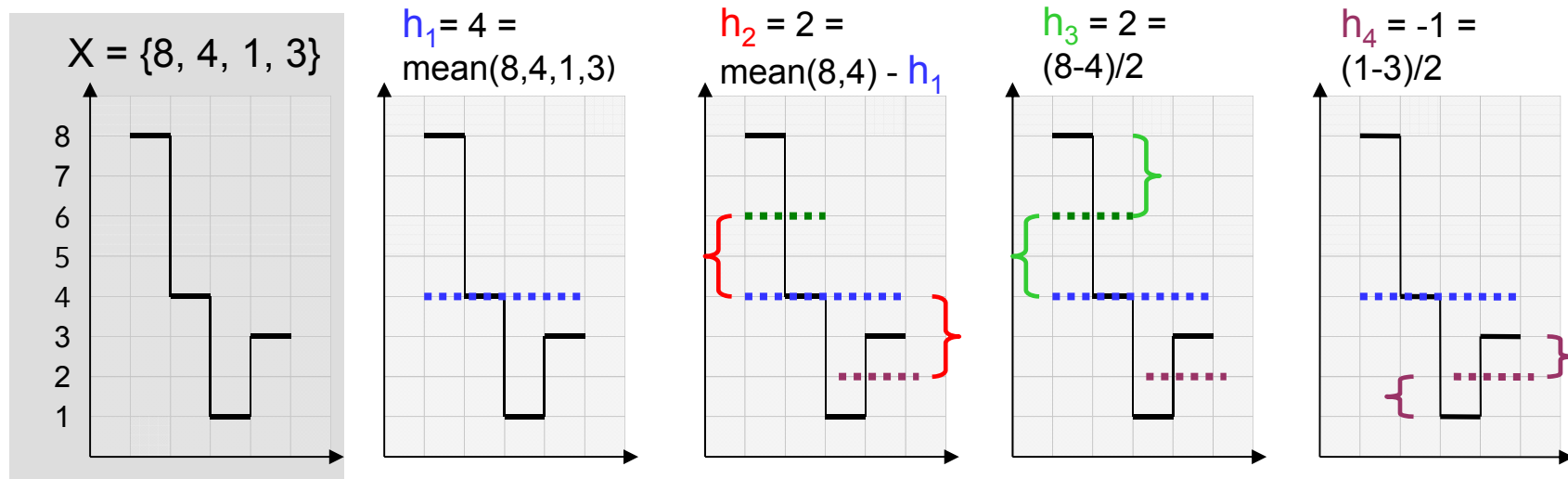
$$\sum_{t=0}^{n-1} |x_t|^2 = \sum_{f=0}^{n-1} |X_f|^2$$

- Now comes the important trick: in practice, the low frequencies (first components) have the highest impact, i.e. contain the most information
- Focusing on the first c coefficients is a good choice if we want to reduce the „dimensionality“ of a sequence
- Since $\|x - y\|^2 = \|X - Y\|^2$ holds, using only c components instead of n yields a lower bounding approximation of the Euclidean Distance
- This approximation will be better when using DFT components instead of original time stamps

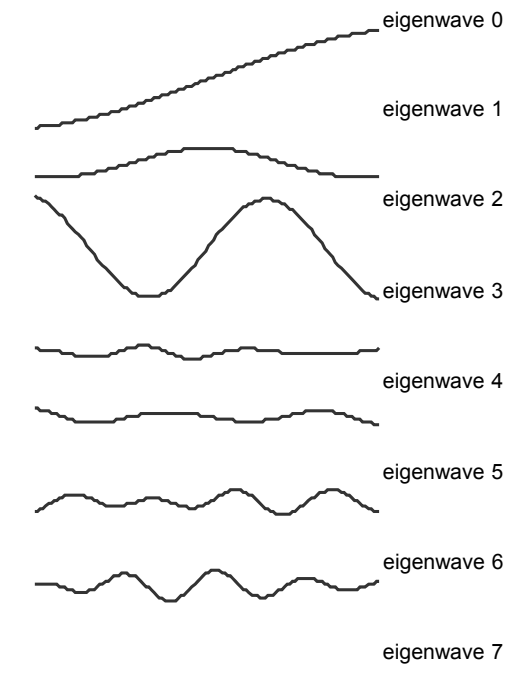
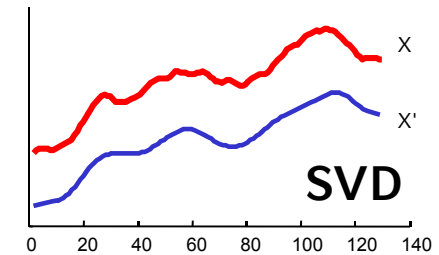
- Discrete Wavelet Transformation (DWT)
 - Idea:
 - Represent a time series as a linear combination of base functions (Wavelet-functions)
 - Typically, Haar-Wavelets are used
 - Properties:
 - The more stationary the time series is, the better is the approximation with fewer components
 - Distance on DWT components also lower bounds Euclidean and DTW distance on original time series
 - Time series are restricted to be of length 2^i (for any i)



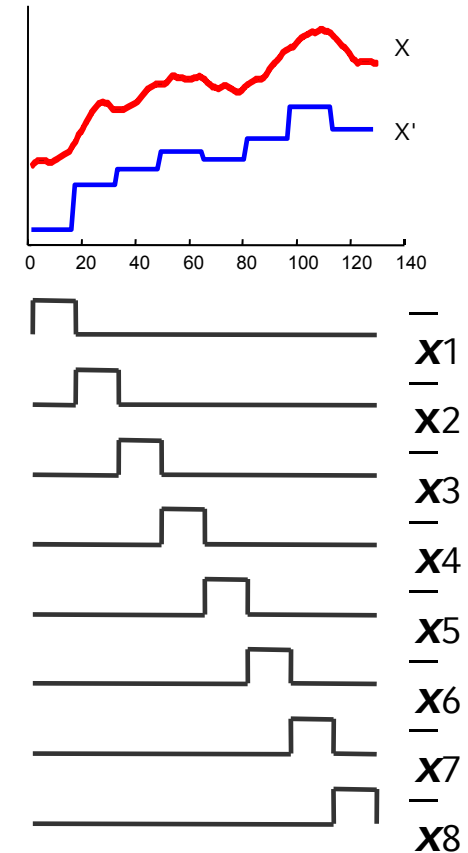
- Example:
 - Stepwise transformation of time series $x = \{8, 4, 1, 3\}$ into Haar Wavelet representation $H = [4, 2, 2, -1]$



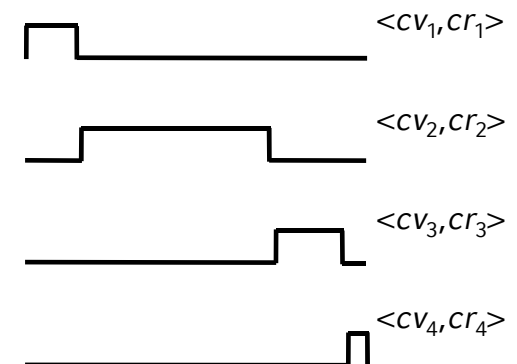
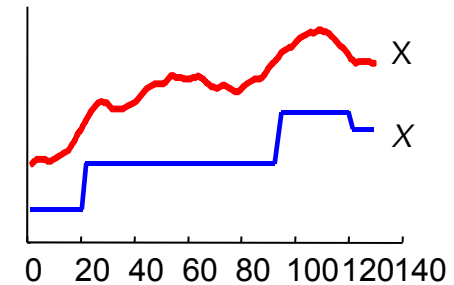
- Singular Value Decomposition (SVD)
 - Idea:
 - Instead of sinus/cosine use Eigen Waves
 - Properties:
 - Minimizes the quadratic approximation error (like PCA and SVD on high dimensional data)
 - The semantics of the components of SVD depends on the actual data while DFT (sin/cos) and DWT (const) are not data dependent
 - In text mining and Information Retrieval, SVD as a feature extraction technique is also known as „Latent Semantic Indexing“



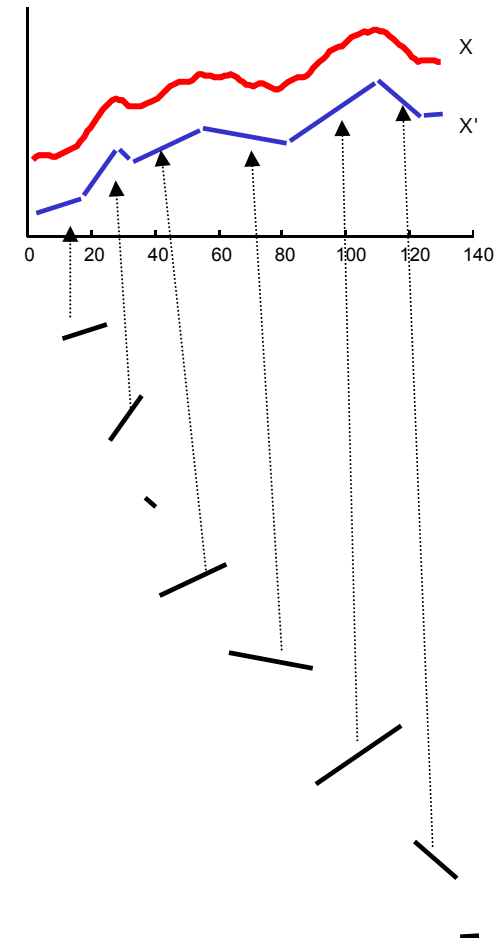
- Piecewise Aggregate Approximation (PAA)
 - Idea:
 - Transform time series into a sequence of box-functions
 - Each box has the same length and approximates the interval by the mean.
 - Properties
 - Lower bounding property
 - Time series may have arbitrary length



- Extension: Adaptive Piecewise Constant Approximation (APCA)
 - Motivation
 - Time series may have time intervals with a small amount details (small amplitude) and intervals with a large amount of details (large amplitude)
 - PAA cannot account for varying amounts of detail
 - Idea
 - Use boxes of variable length
 - Each segment now requires 2 parameters



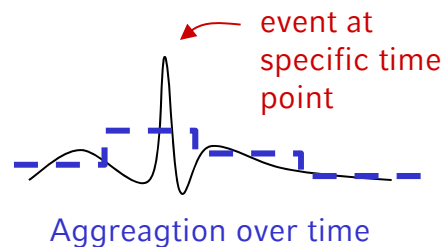
- Piecewise Linear Approximation (PLA)
 - Idea
 - Transform time series into a sequence of line segments $s = (\text{length}, \text{height}_{\text{start}}, \text{height}_{\text{end}})$
 - Two consecutive segments need not to be connected
 - Properties
 - Good approximation depends on #segments
 - Each component (segment) is a rich approximation but requires more parameters
 - Lower bounds Euclidean and DTW



- An example of a specific feature transformation to model a special notion of similarity of time series is „threshold-based similarity“

[Assfalg, Kriegel, Kröger, Kunath, Pryakhin, Renz. Proc. 10th Int. Conf. on Extending Database Technology (EDBT), 2006]

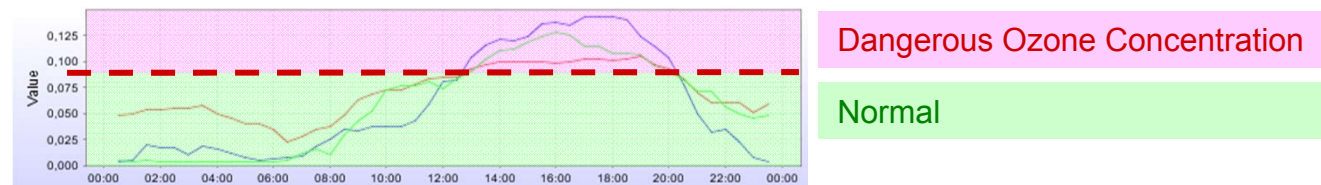
- Basic Idea:
 - In some applications, only significant „events“ that are defined by certain amplitudes (or amplitude values) are interesting
 - So far, the feature extraction extracts features modeling certain properties of time intervals but not of amplitude intervals



- Sample Applications

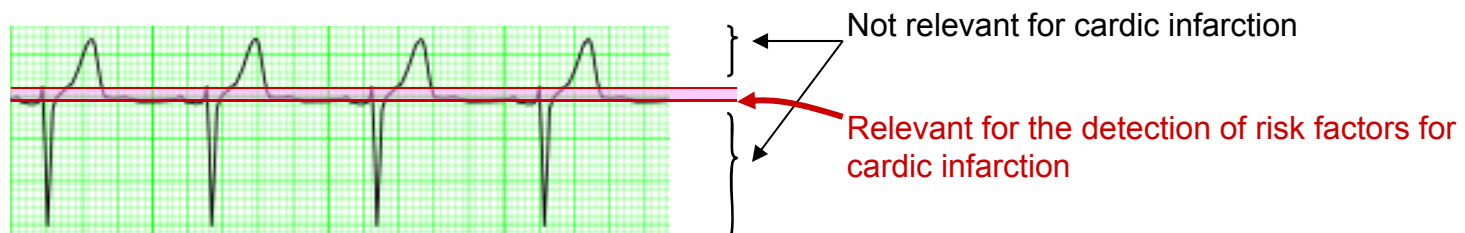
- Environmental Science: analyzing critical ozone concentrations?

- Find cluster of regions (time series) that exceed the allowed threshold in similar time intervals



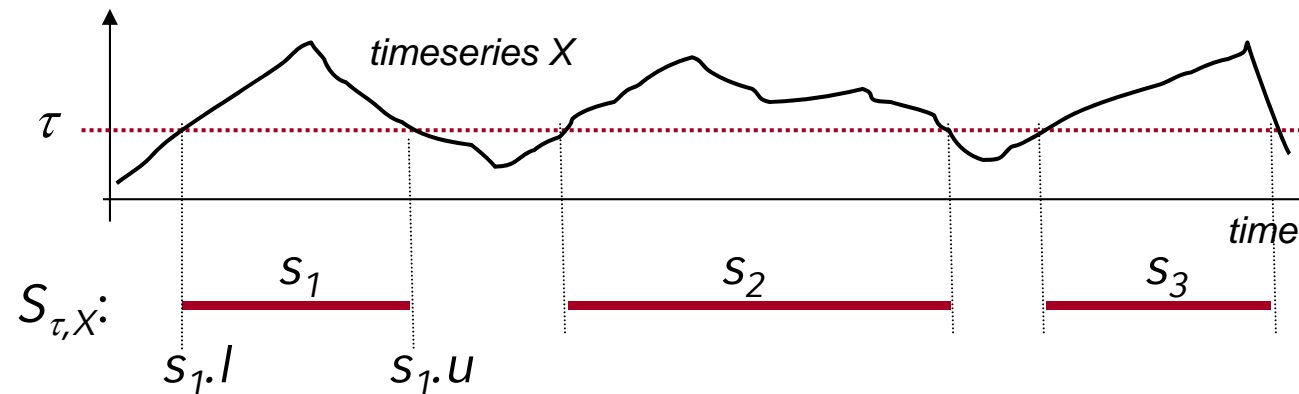
- Medical diagnosis: potential for cardiac infarction?

- Find clusters of heart rates by focusing on the relevant amplitude intervals



– Similarity Model

- Time series $X = \langle (x_i, t_i) : i = 1..N \rangle$ is transformed into a sequence of intervals $S_{\tau, X} = \{s_j : j = 1..M\}$, such that: $\forall t \in T : (\exists s_j \in S_{\tau, X} : s_j.l < t < s_j.u) \Leftrightarrow x(t) > \tau$.



- Similarity of time series = similarity of sequences of intervals

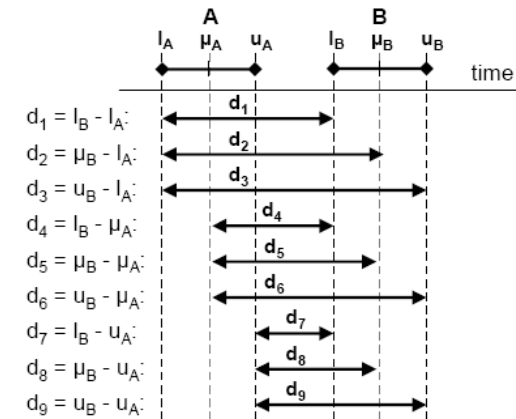


Threshold-based similarity

- Similarity between sequences of intervals?
- First: distance on intervals?
 - Euclidean distance on l- and u-values:

$$d_{\text{int}}(s_1, s_2) = \sqrt{(l_1 - l_2)^2 + (u_1 - u_2)^2}$$

where $s_1 = (l_1, u_1)$ and $s_2 = (l_2, u_2)$



- Use sum of minimum distance between two sequences of intervals S_X and S_Y

$$d_{TS}(S_X, S_Y) = \frac{1}{2} \cdot \left(\underbrace{\frac{1}{|S_X|} \cdot \sum_{s \in S_X} \min_{t \in S_Y} d_{\text{int}}(s, t)}_{S_X \dashrightarrow S_Y} + \underbrace{\frac{1}{|S_Y|} \cdot \sum_{t \in S_Y} \min_{s \in S_X} d_{\text{int}}(t, s)}_{S_X \dashleftarrow S_Y} \right)$$

- Round-up:
 - Feature extraction method serve the purpose of
 - Finding a compact representation of the original time series (mostly for performance reasons)
 - Compact representations can be used for approximate similarity computations
 - Some have bounding properties (e.g. lower bounding the exact distance/similarity) that can be used for indexing/pruning
- or
- Modeling a specialized notion of similarity of a time series for a given application

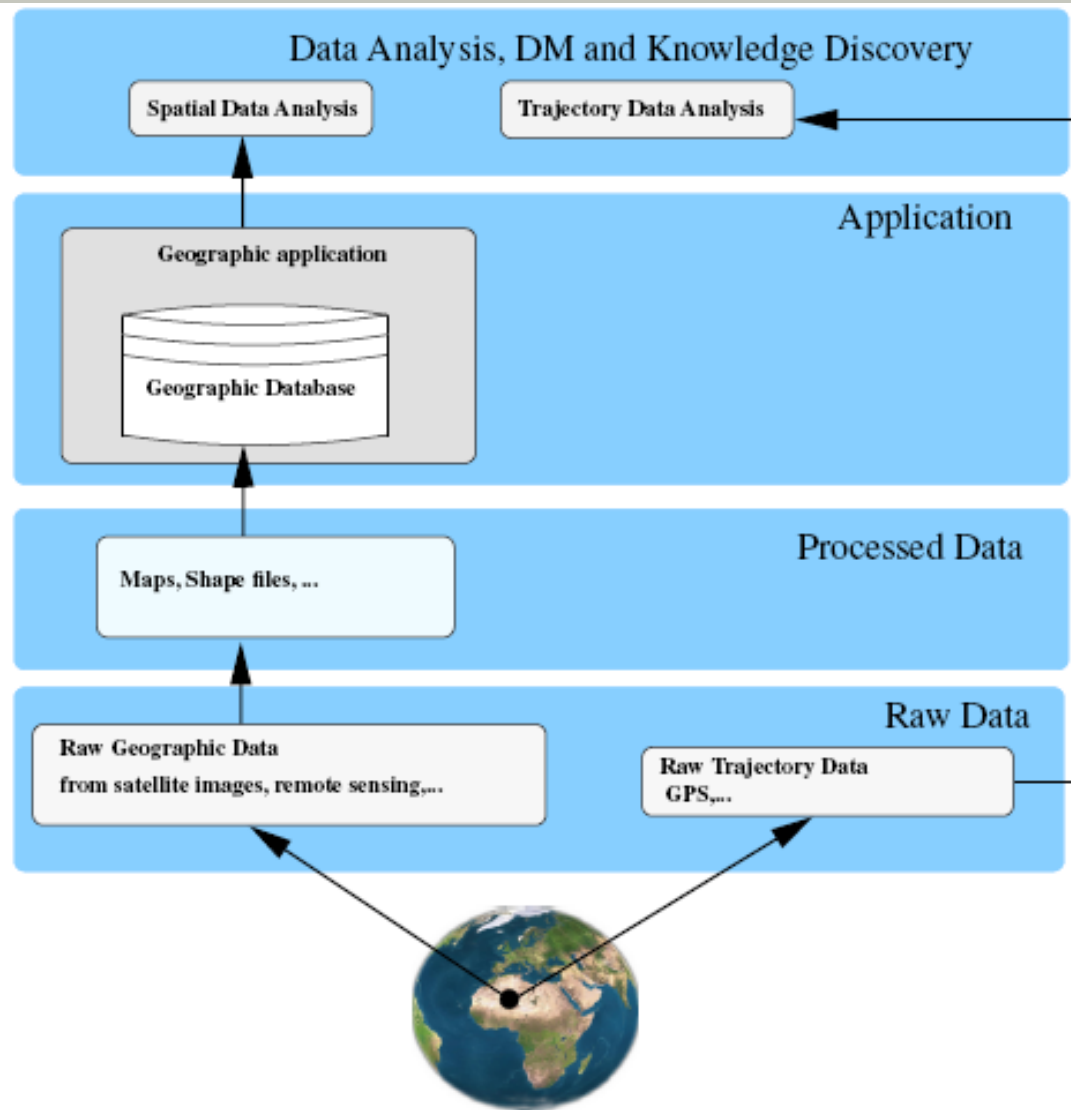
1. Introduction
2. Sequence Data
3. Time Series Data
4. Spatial Temporal Data^[1]

1. BOGORNY, V., and S. SHEKHAR. "Tutorial on Spatial and Spatio-Temporal Data Mining." *Part ii-Trajectory*.

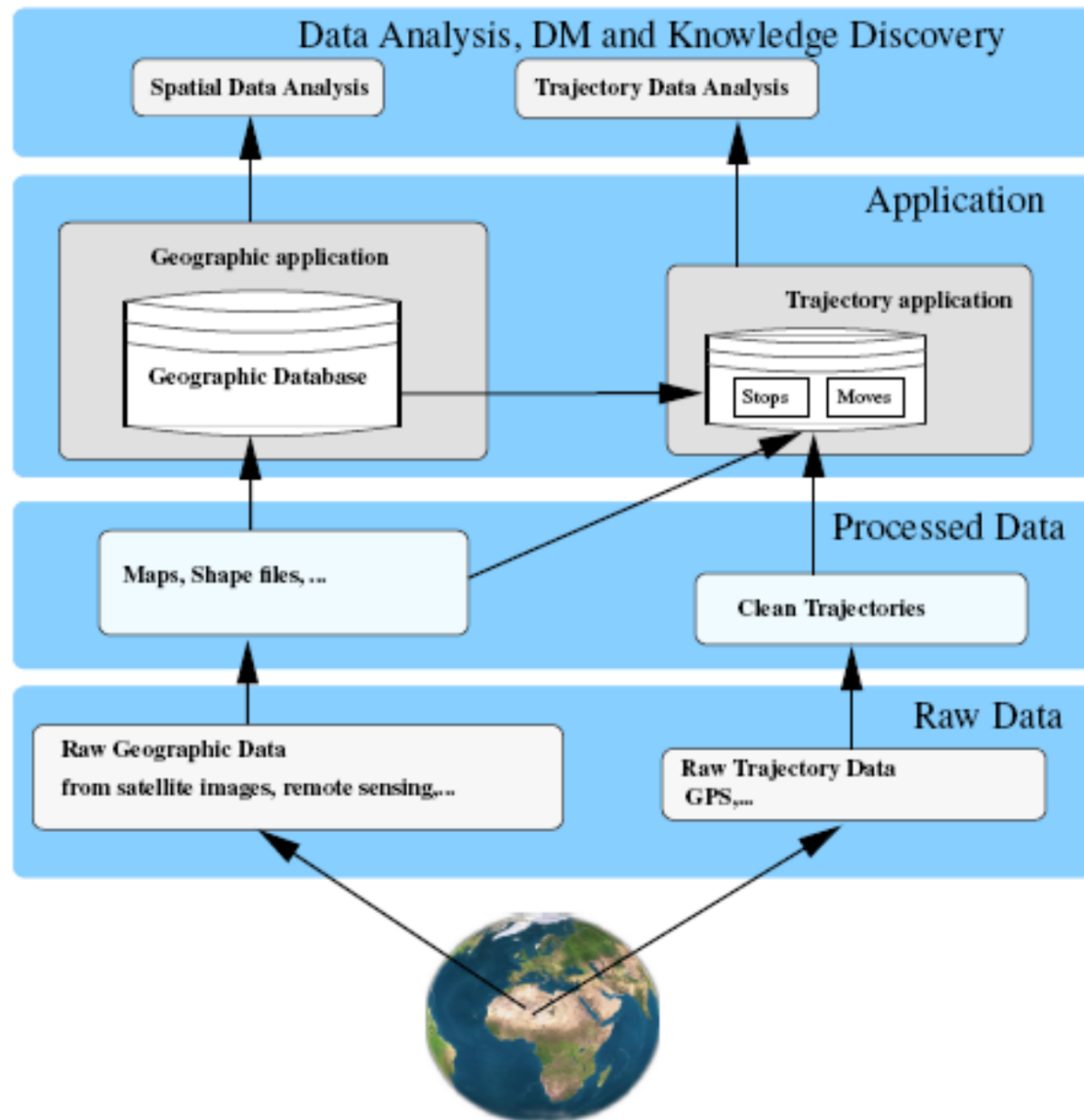
- Spatial-temporal data is a special case of time series where (one of) the information recorded at each time point is the location of an object
- A time series over spatial locations is also called “trajectory”
- Often, there is additional information on time slots (e.g. semantic information on the location such as “museum” or “airport” ...)
- We review the some of the recent trends in mining spatial-temporal (aka: spatio-temporal) data

- In general, there are two major approaches to trajectory mining:
 - **Geometry-based methods** consider only geometrical properties of trajectories; they focus on “location-based” similarity
 - **Semantic-based methods** compute patterns based on the semantics of the data and are somewhat independent of the specific spatial locations

Geometry-based approach

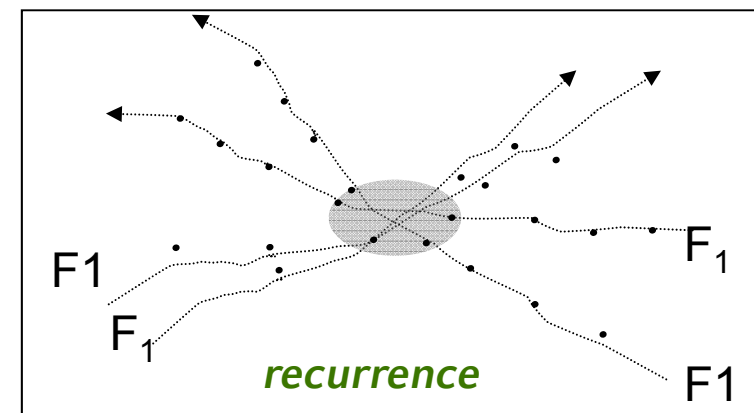
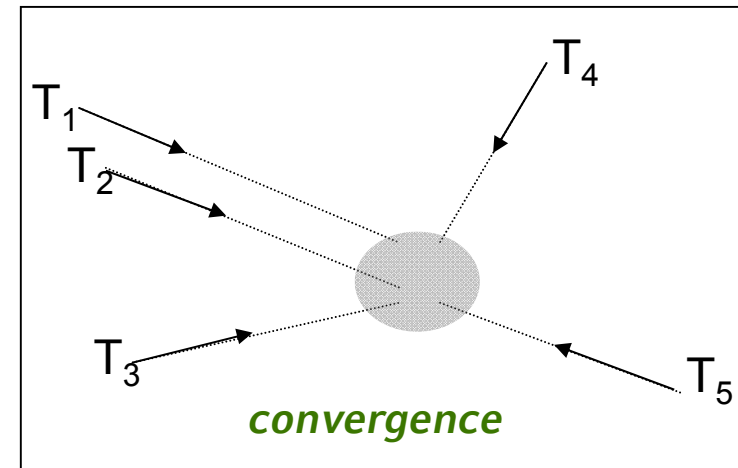


Semantic-based approach

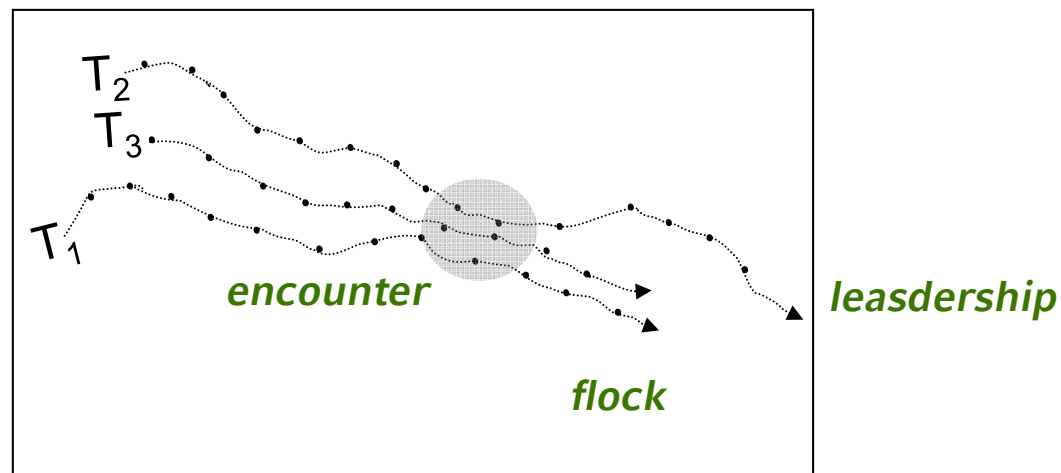


- Laube et al. 2004 proposed five patterns based on location, direction, and/or movement:

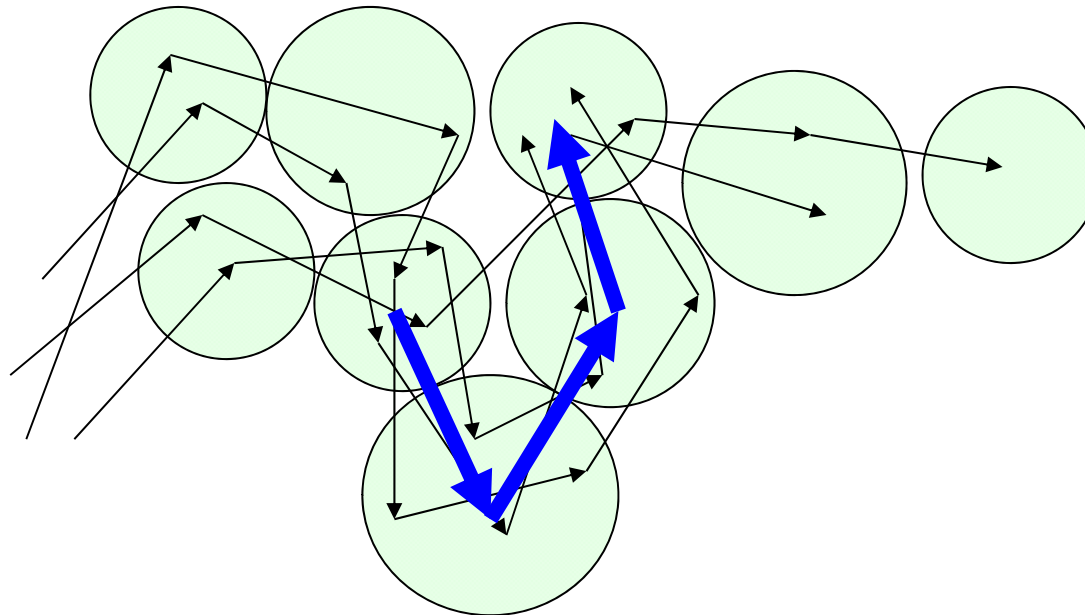
- Convergence:** At least m entities pass through the same circular region of radius r (regardless of time)
- Recurrence:** at least m entities visit a circular region at least k times



3. **Flock pattern:** At least m entities are within a region of radius r and move in the same direction during a time interval $\geq s$ (e.g. traffic jam)
4. **Leadership:** At least m entities are within a circular region of radius r , they move in the same direction, and at least one of the entities is heading in that direction for at least t time steps. (e.g. bird migration)
5. **Encounter:** At least m entities will be concurrently inside the same circular region of radius r , assuming they move with the same speed and direction. (e.g. traffic jam at some moment if cars keep moving in the same direction)



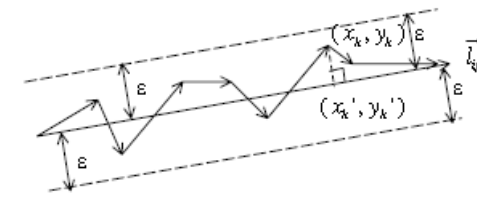
- Frequent patterns: frequent followed paths/frequent sequential patterns



- Computing frequent sequential patterns (e.g. Cao 2005):

1. Transforms each trajectory in a line with several segments

- A distance tolerance measure is defined
- All trajectory points inside this distance are summarized in one segment



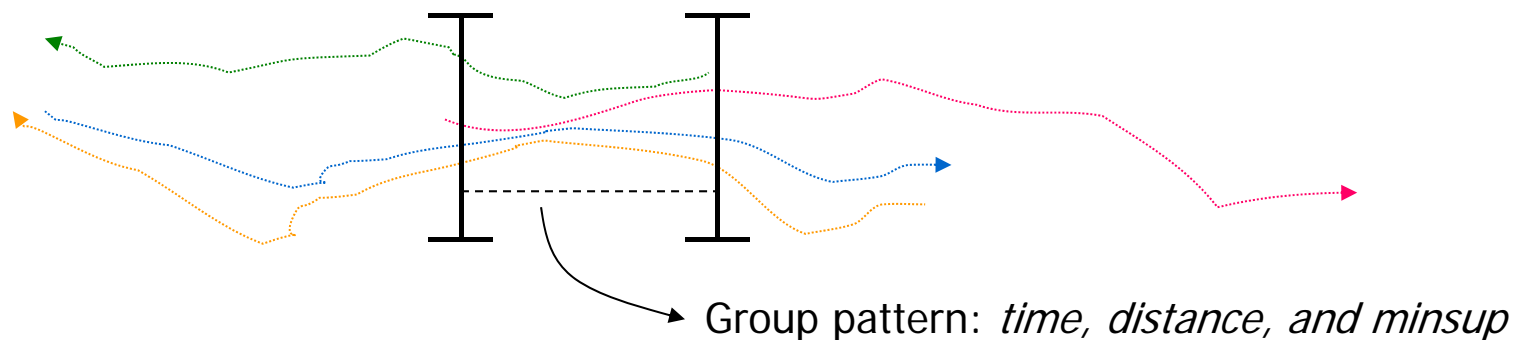
(a) Segment complies with $l_{ij}^{\vec{r}}$

2. Similar segments are grouped

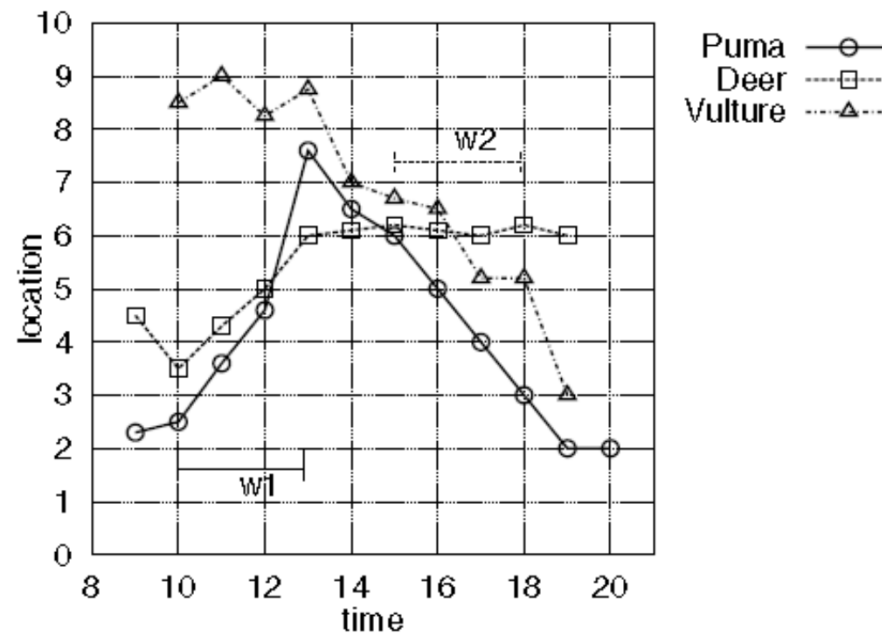
- Similarity is based on the angle and the spatial length of the segment
 - Segments with same angle and length have their distance checked based on a given distance threshold
- From the resulting groups, a medium segment is created
 - From this segment a region is created

3. Frequent sequences of regions are computed considering a minSup threshold

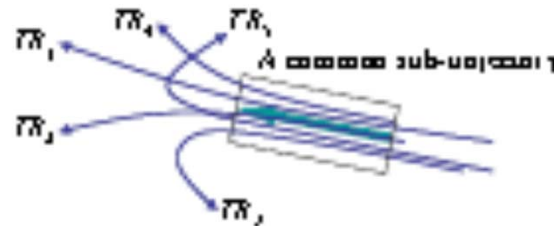
- Frequent mobile group patterns (Hwang 2005):
 - A group pattern is a set of trajectories close to each other (with distance less than a given *minDist*) for a minimal amount of time (*minTime*)
 - Direction is not considered
 - Use Apriori algorithm to compute frequent groups



- Co-location Patterns (Cao 2006):
 - Co-location episodes in spatio-temporal data
 - Trajectories are spatially close in a time window and move together



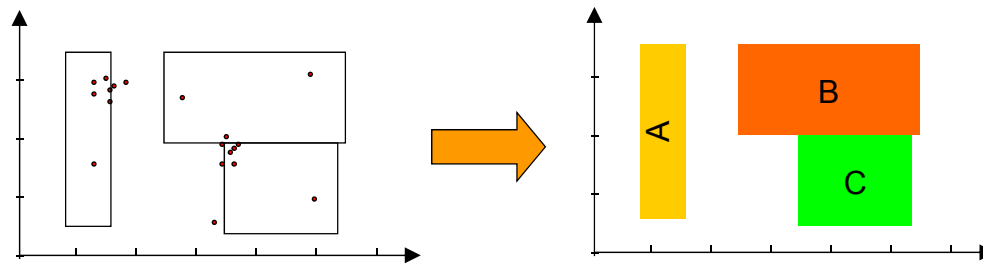
- Trajectory Clustering (Han 2007):
 - Algorithm TraClus: Group sub-trajectories using a density based clustering algorithm
 - 2 step approach
 1. Partition each trajectory in line segments with a user defined length L
 2. Cluster similar line segments based on spatial proximity of the time points
 - Similarity of line segments: Euclidean distance between segments (sub-trajectories); in theory: could be anything else
 - => however, time is not considered in this approach



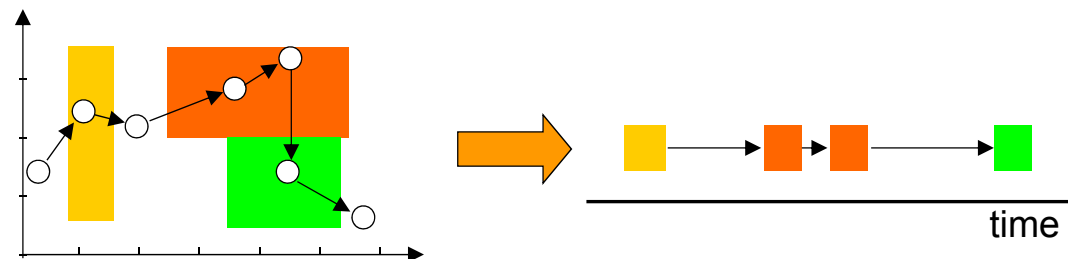
- Sequential Trajectory Pattern Mining (T-Patterns; Giannotti 2007):
 - Considers both space and time
 - Describes frequent behavior in terms of visited regions (ROIs)
 - Three-step approach
 1. Compute regions of interest (ROIs), i.e., regions with many trajectories (regardless of time)
 2. Transform trajectory into sequence of ROIs: select trajectories intersecting at least two regions in a sequence and annotate the time traveled between regions
 3. Compute T-Patterns, i.e., sequences of regions visited during the same time intervals

– Visualization of the idea of T-Patterns:

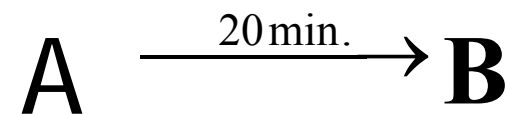
- Regions of interest (ROIs)



- Transform trajectory into a sequence of ROIs

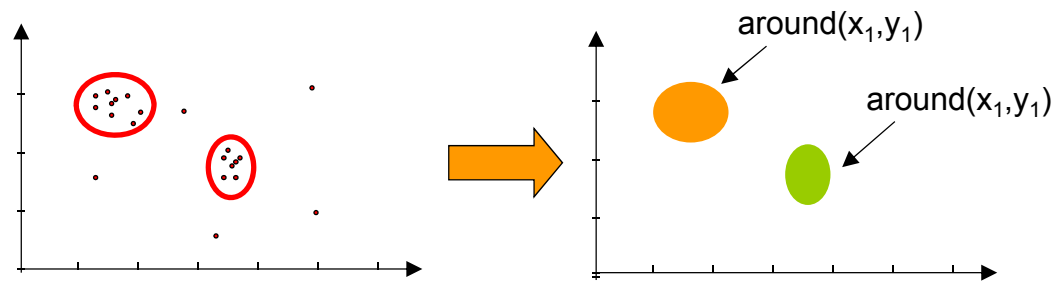


- Sample pattern:

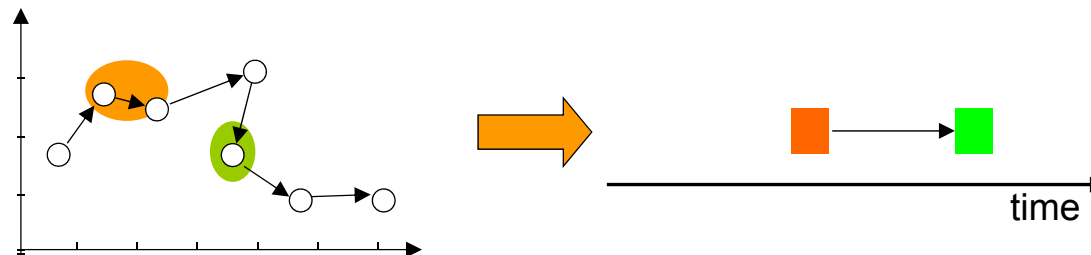


– Visualization of the approach

- Step 1: detection of ROIs



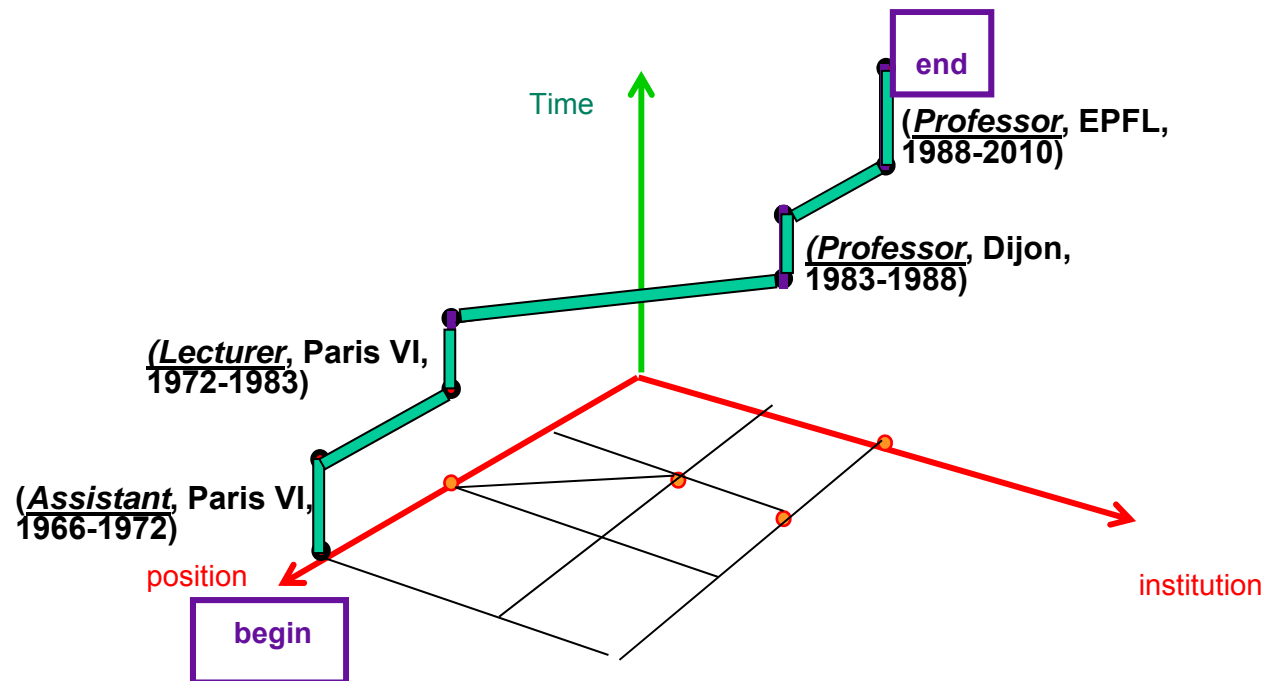
- Step 2: transformation



- Compute pattern:

$$\textit{around}(x_1, y_1) \xrightarrow{20 \text{ min.}} \textit{around}(x_2, y_2)$$

- A Conceptual View on Trajectories (Spaccapietra 2008)
 - Trajectory is a spatio-temporal object that has generic features (independent of the application) and *semantic* features (depend on the application)
 - Trajectory = travel in abstract space, e.g. 2D career space:



- Semantic trajectories = geo data + trajectory data



Trajectory Samples (x,y,t)



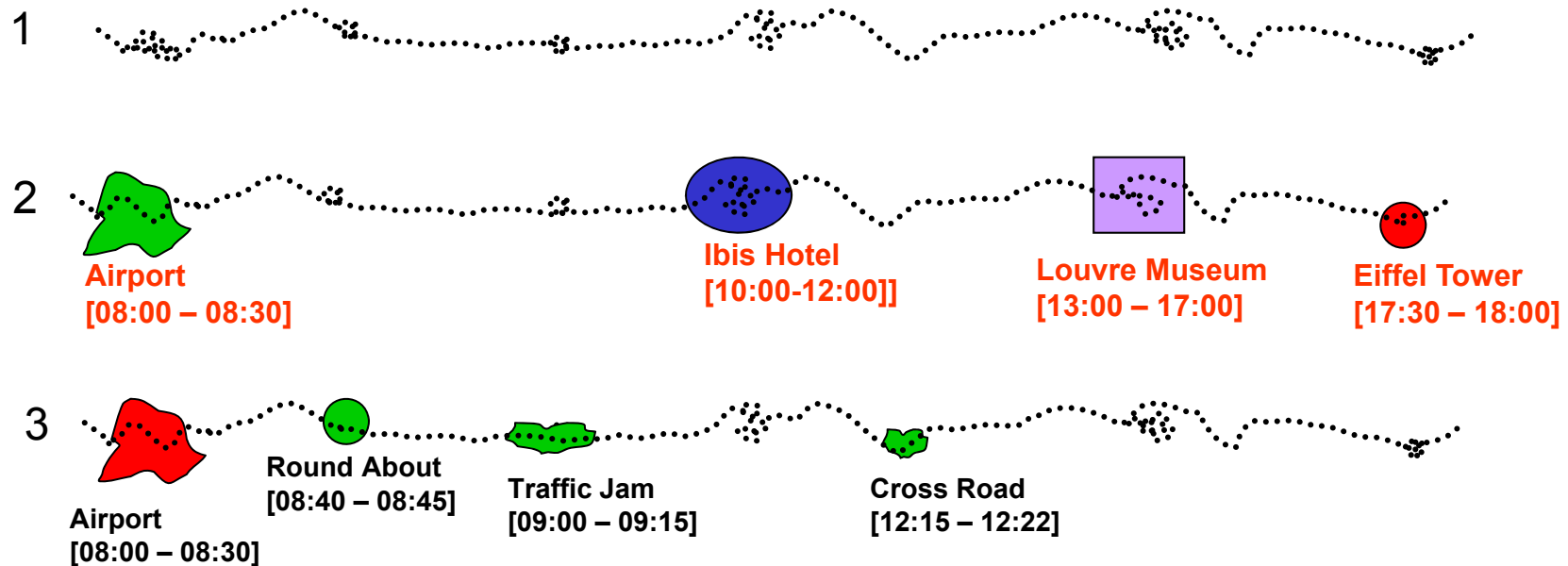
Geographic Data



Geographic Data +
Trajectory Data =
Semantic Trajectories

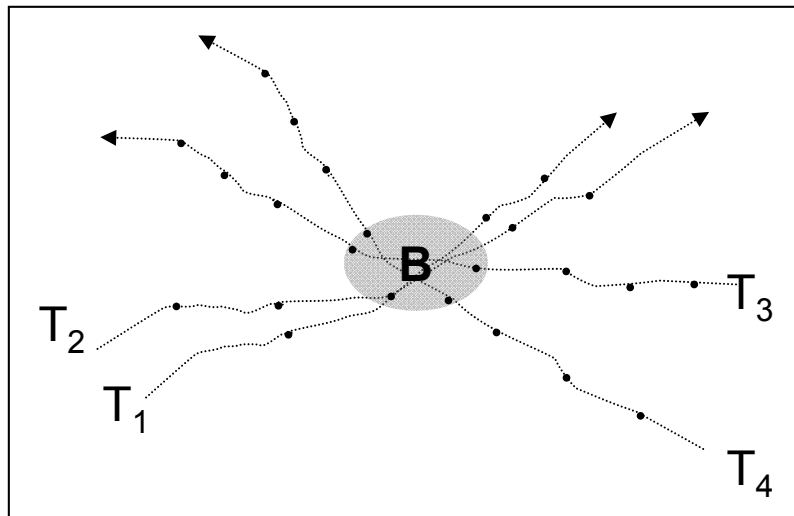
- Difference between stops and moves
 - STOPS
 - Important parts of trajectories
 - Where the moving object has stayed for a minimal amount of time
 - Stops are application dependent
 - Tourism application: Hotels, touristic places, airport, ...
 - Traffic Management Application: Traffic lights, roundabouts, big events...
 - MOVES
 - Are the parts that are not stops

- Stops and moves are independent of the application

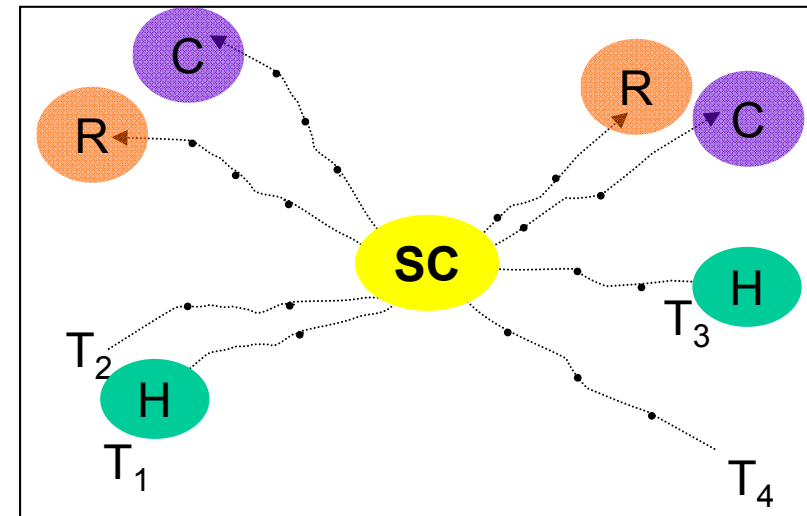


- Geometric Patterns enriched by semantics (Bogorny 2008):
 - Very little semantics in most trajectory mining approaches (geometry-based approaches)
- Thus:
- Patterns are purely geometrical
 - Hard to interpret
- Thus:
- Enrich geometric patterns with semantic information
(stimulated many approaches on how to add semantics to trajectories)

- Semantic Enrichment (Example):



Geometric Pattern



● H Hotel
 ● R Restaurant
 ● C Cinema



Semantic trajectory Pattern

- (a) Hotel to Restaurant, passing by SC
- (b) go to Cinema, passing by SC

- Stop and Move computation: SMoT (Alvares 2007a)
 - A *candidate stop* C is a tuple (R_C, Δ_C) , where
 - R_C is the geometry of the candidate stop (spatial feature type)
 - Δ_C is the *minimal time duration*

E.g. [Hotel - 3 hours]

- An *application* A is a finite set

$A = \{C_1 = (R_{C_1}, \Delta_{C_1}), \dots, C_N = (R_{C_N}, \Delta_{C_N})\}$ of *candidate stops* with non-overlapping geometries R_{C_1}, \dots, R_{C_N}

E.g. [Hotel - 3 hours, Museum – 1 hour]

- Stop and Move computation: SMoT (Alvares 2007a)

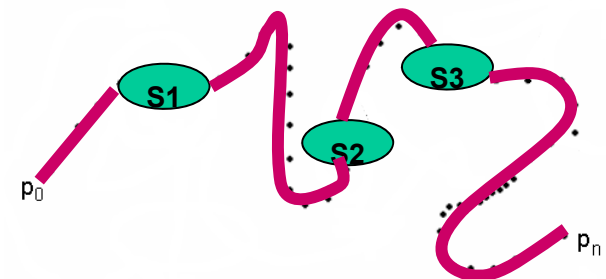
A *stop* of a trajectory T with respect to an *application* A is a tuple (R_{C_k}, t_j, t_{j+n}) , such that a maximal subtrajectory of

$$T \{(x_i, y_i, t_i) \mid (x_i, y_i) \text{ intersects } R_{C_k}\} = \{(x_j, y_j, t_j), (x_{j+1}, y_{j+1}, t_{j+1}), \dots, (x_{j+n}, y_{j+n}, t_{j+n})\}$$

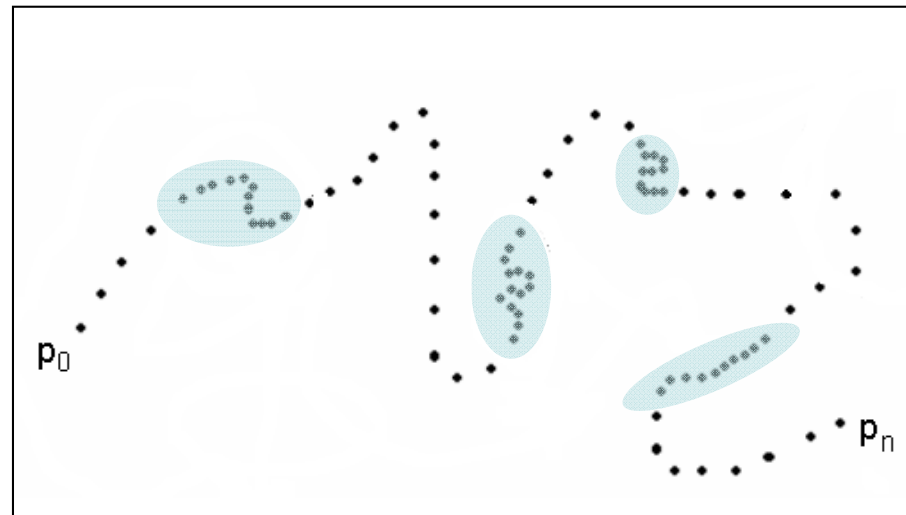
where R_{C_k} is the geometry of C_k and $|t_{j+n} - t_j| \geq \Delta_{C_k}$

A *move* of T with respect to A is:

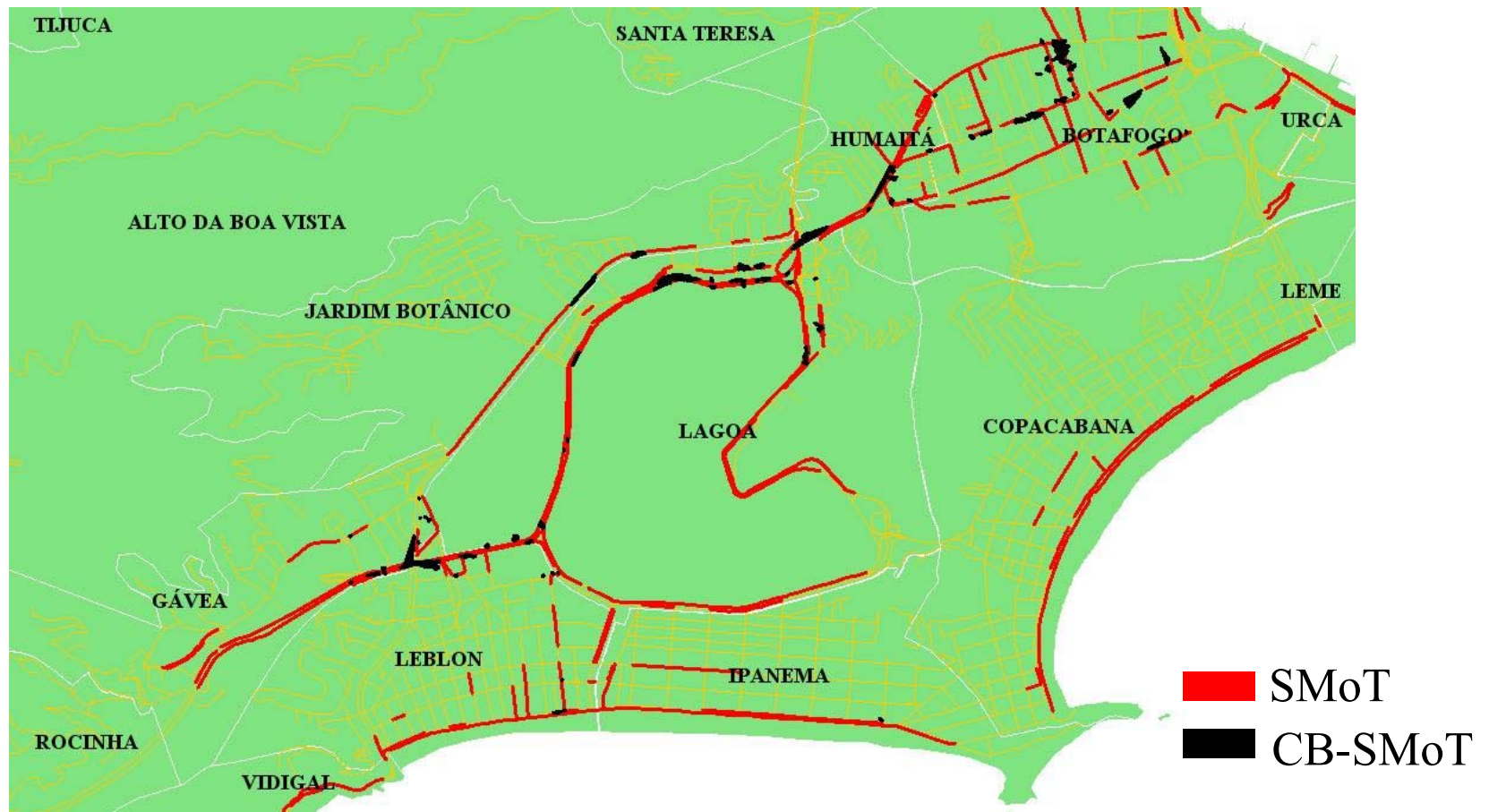
- ❖ a maximal contiguous subtrajectory of T :
 - ❖ between the starting point of T and the first stop of T ; OR
 - ❖ between two consecutive stops of T ; OR
 - ❖ between the last stop of T and the ending point of T ;
- ❖ or the trajectory T itself, if T has no stops.



- Improvement: CB-SMoT (Palma 2008)
 - Cluster based: cluster trajectories based on speed
 - Low speed => important place
 - Algorithm similar to SMoT but clusters trajectory points first and adds semantics to clusters



- Comparison: SMOt vs. CB-SMOt (Application: transportation)



- Geometric based methods

Huiping Cao, Nikos Mamoulis, David W. Cheung: Discovery of Periodic Patterns in Spatiotemporal Sequences. *IEEE Trans. Knowl. Data Eng.* 19(4): 453-467 (2007)

Panos Kalnis, Nikos Mamoulis, Spiridon Bakiras: On Discovering Moving Clusters in Spatio-temporal Data. *SSTD*, 364-381 (2005)

Florian Verhein, Sanjay Chawla: Mining spatio-temporal patterns in object mobility databases. *Data Min. Knowl. Discov.* 16(1): 5-38 (2008)

Florian Verhein, Sanjay Chawla: Mining Spatio-temporal Association Rules, Sources, Sinks, Stationary Regions and Thoroughfares in Object Mobility Databases. *DASFAA*, 187-201 (2006)

Changqing Zhou, Dan Frankowski, Pamela J. Ludford, Shashi Shekhar, and Loren G. Terveen. Discovering personally meaningful places: An interactive clustering approach. *ACM Trans. Inf. Syst.*, 25(3), 2007.

Cao, H., Mamoulis, N., and Cheung, D. W. (2005). Mining frequent spatio-temporal sequential patterns. In *ICDM '05: Proceedings of the Fifth IEEE International Conference on Data Mining*, pages 82–89, Washington, DC, USA. IEEE Computer Society.

- Geometric based methods (cont.)

Laube, P. and Imfeld, S. (2002). Analyzing relative motion within groups of trackable moving point objects. In Egenhofer, M. J. and Mark, D. M., editors, GIScience, volume 2478 of Lecture Notes in Computer Science, pages 132–144. Springer.

Laube, P., Imfeld, S., and Weibel, R. (2005a). Discovering relative motion patterns in groups of moving point objects. *International Journal of Geographical Information Science*, 19(6):639–668.

Laube, P., van Kreveld, M., and Imfeld, S. (2005b). Finding REMO: Detecting Relative Motion Patterns in Geospatial Lifelines. Springer.

Lee, J.-G., Han, J., and Whang, K.-Y. (2007). Trajectory clustering: a partition-and-group framework. In Chan, C. Y., Ooi, B. C., and Zhou, A., editors, SIGMOD Conference, pages 593–604. ACM.

Li, Y., Han, J., and Yang, J. (2004). Clustering moving objects. In KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, pages 617–622, New York, NY, USA. ACM Press.

Nanni, M. and Pedreschi, D. (2006). Time-focused clustering of trajectories of moving objects. *Journal of Intelligent Information Systems*, 27(3):267–289.

- Geometric based methods (cont.)

Verhein, F. and Chawla, S. (2006). Mining spatio-temporal association rules, sources, sinks, stationary regions and thoroughfares in object mobility databases. In Lee, M.-L., Tan, K.-L., and Wuwongse, V., editors, DASFAA, volume 3882 of Lecture Notes in Computer Science, pages 187–201. Springer.

Gudmundsson, J. and van Kreveld, M. J. (2006). Computing longest duration flocks in trajectory data. In [de By and Nittel 2006], pages 35–42.

Gudmundsson, J., van Kreveld, M. J., and Speckmann, B. (2007). Efficient detection of patterns in 2d trajectories of moving points. *GeoInformatica*, 11(2):195–215.

Hwang, S.-Y., Liu, Y.-H., Chiu, J.-K., and Lim, E.-P. (2005). Mining mobile group patterns: A trajectory-based approach. In Ho, T. B., Cheung, D. W.-L., and Liu, H., editors, PAKDD, volume 3518 of Lecture Notes in Computer Science, pages 713–718. Springer.

Cao, H., Mamoulis, N., and Cheung, D. W. (2006). Discovery of collocation episodes in spatiotemporal data. In ICDM, pages 823–827. IEEE Computer Society.

- Semantic based method

Bogorny, V. ; Wachowicz, M. A Framework for Context-Aware Trajectory Data Mining. In: Longbing Cao, Philip S. Yu, Chengqi Ahang, Huaifeng Zhang. (Org.). Domain Driven Data Mining: Domain Problems and Applications. 1 ed. : Springer, 2008a.

Bogorny, V., Kuijpers, B., and Alvares, L. O. (2008b). St-dmql: a semantic trajectory data mining query language. International Journal of Geographical Information Science. Taylor and Francis, 2008.

Palma, A. T; Bogorny, V.; Kuijpers, B.; Alvares, L.O. *A Clustering-based Approach for Discovering Interesting Places in Trajectories*. In: 23rd Annual Symposium on Applied Computing, (ACM-SAC'08), Fortaleza, Ceara, 16-20 March (2008) Brazil. pp. 863-868.

Spaccapietra, S., Parent, C., Damiani, M. L., de Macedo, J. A., Porto, F., and Vangenot, C. (2008). A conceptual view on trajectories. Data and Knowledge Engineering, 65(1):126–146.

Alvares, L. O., Bogorny, V., de Macedo, J. F., and Moelans, B. (2007a). Dynamic modeling of trajectory patterns using data mining and reverse engineering. In Twenty-Sixth International Conference on Conceptual Modeling - ER2007 - Tutorials, Posters, Panels and Industrial Contributions, volume 83, pages 149–154. CRPIT.

Alvares, L. O., Bogorny, V., Kuijpers, B., de Macedo, J. A. F., Moelans, B., and Vaisman, A. (2007b). A model for enriching trajectories with semantic geographical information. In ACM-GIS, pages 162–169, New York, NY, USA. ACM Press.