

Representing and understanding patterns in materials and molecules

Michele Ceriotti
EPFL/IMX/COSMO

MPG-EPFL Summer School



Acknowledgements



Piero Gasparotto
Sandip De



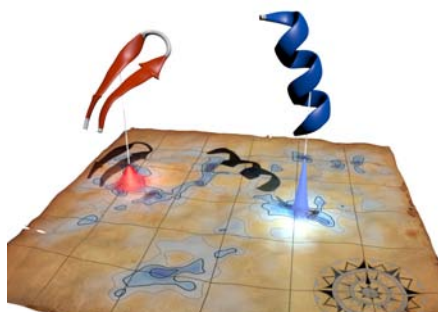
Michele Parrinello
Gareth Tribello



Outline



- Analysis of molecular data from simulations: big data and high dimensionality
- Cluster analysis and recognition of molecular patterns
 - Hydrogen bonds, and secondary structure patterns
- Mapping high-dimensional data in low dimension
 - Linear methods: Principal Components Analysis
 - Non-linear methods: ISOMAP, LLE, Sketch-map
 - From proteins to clusters



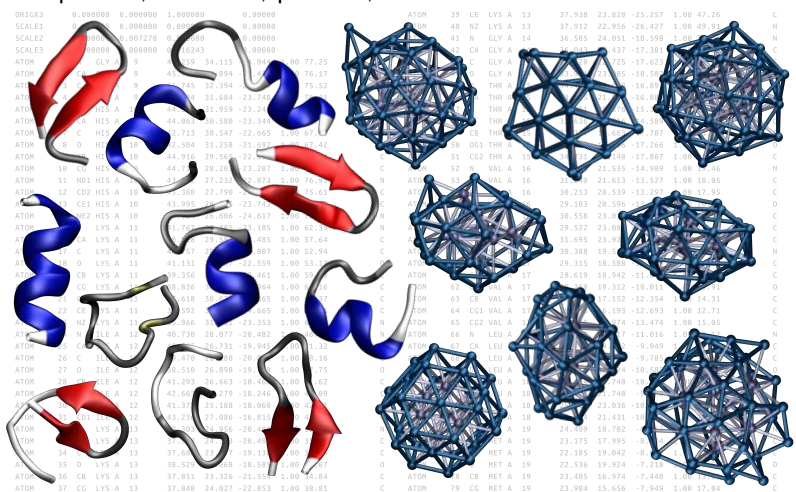
High dimensional data in atomistic simulation



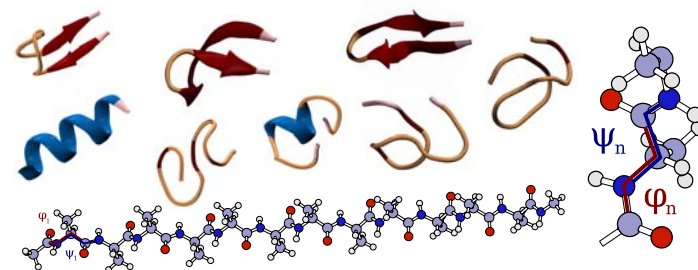
- Atomistic simulations provide **too much** information
- It is hard to decipher the essential features in structurally-complex compounds, materials, proteins, etc.

| | | | | | | | | | | | | | | | | |
|--------|----------|----------|----------|----------|------|--------|--------|---------|------|-------|--------|--------|---------|------|--------|---|
| ORIGX3 | 0.000000 | 0.000000 | 1.000000 | 0.000000 | ATOM | 39 | CE | LVS | A | 13 | 37.938 | 23.028 | -25.257 | 1.00 | 47.26 | C |
| SCALE1 | 0.008865 | 0.000000 | 0.000054 | 0.000000 | ATOM | 40 | NZ | LVS | A | 13 | 37.912 | 22.956 | -26.427 | 1.00 | 49.91 | N |
| SCALE2 | 0.000000 | 0.007278 | 0.000000 | 0.000000 | ATOM | 41 | N | GLY | A | 14 | 36.505 | 24.051 | -18.598 | 1.00 | 30.23 | N |
| SCALE3 | 0.000000 | 0.000000 | 0.010243 | 0.000000 | ATOM | 42 | CA | GLY | A | 14 | 36.043 | 23.437 | -17.393 | 1.00 | 27.87 | C |
| ATOM | 1 | N | GLY | A | 9 | 47.259 | 34.115 | -24.004 | 1.00 | 77.25 | 34.728 | 22.725 | -17.623 | 1.00 | 24.00 | C |
| ATOM | 2 | CA | GLY | A | 9 | 45.958 | 33.894 | -23.454 | 1.00 | 76.17 | 33.931 | 23.085 | -18.582 | 1.00 | 24.28 | N |
| ATOM | 3 | C | GLY | A | 9 | 45.745 | 32.394 | -23.509 | 1.00 | 75.52 | 34.467 | 21.648 | -16.893 | 1.00 | 23.77 | N |
| ATOM | 4 | O | GLY | A | 9 | 46.734 | 31.684 | -23.748 | 1.00 | 76.41 | 33.165 | 21.088 | -16.808 | 1.00 | 23.99 | C |
| ATOM | 5 | N | HIS | A | 10 | 44.513 | 31.959 | -23.246 | 1.00 | 75.46 | 33.851 | 20.876 | -15.366 | 1.00 | 20.67 | C |
| ATOM | 6 | CA | HIS | A | 10 | 44.063 | 30.580 | -23.344 | 1.00 | 70.95 | 33.636 | 20.316 | -14.589 | 1.00 | 20.91 | N |
| ATOM | 7 | C | HIS | A | 10 | 42.713 | 30.547 | -22.665 | 1.00 | 67.35 | 33.307 | 19.663 | -17.787 | 1.00 | 21.11 | C |
| ATOM | 8 | O | HIS | A | 10 | 42.504 | 31.250 | -21.692 | 1.00 | 67.42 | 32.445 | 18.636 | -17.266 | 1.00 | 20.61 | C |
| ATOM | 9 | CB | HIS | A | 10 | 44.916 | 29.565 | -22.586 | 1.00 | 72.64 | 34.721 | 19.148 | -17.867 | 1.00 | 21.03 | C |
| ATOM | 10 | CG | HIS | A | 10 | 44.712 | 28.203 | -23.207 | 1.00 | 74.51 | 31.766 | 21.535 | -14.989 | 1.00 | 19.46 | N |
| ATOM | 11 | ND1 | HIS | A | 10 | 43.872 | 27.232 | -22.873 | 1.00 | 76.95 | 31.240 | 21.613 | -13.527 | 1.00 | 18.85 | C |
| ATOM | 12 | CD2 | HIS | A | 10 | 45.388 | 27.790 | -24.318 | 1.00 | 75.61 | 38.213 | 20.839 | -13.297 | 1.00 | 17.95 | C |
| ATOM | 13 | CE1 | HIS | A | 10 | 43.995 | 26.255 | -23.743 | 1.00 | 77.59 | 29.183 | 20.396 | -13.979 | 1.00 | 20.19 | O |
| ATOM | 14 | NE2 | HIS | A | 10 | 44.904 | 26.606 | -24.617 | 1.00 | 77.84 | 38.558 | 23.011 | -13.268 | 1.00 | 17.93 | C |
| ATOM | 15 | N | LVS | A | 11 | 41.767 | 29.785 | -23.185 | 1.00 | 62.39 | 29.522 | 23.084 | -12.325 | 1.00 | 14.31 | C |
| ATOM | 16 | CA | LVS | A | 11 | 40.518 | 29.556 | -22.465 | 1.00 | 57.64 | 31.605 | 23.027 | -12.855 | 1.00 | 19.85 | C |
| ATOM | 17 | C | LVS | A | 11 | 40.757 | 28.260 | -21.987 | 1.00 | 52.94 | 38.388 | 19.552 | -12.431 | 1.00 | 16.80 | N |
| ATOM | 18 | O | LVS | A | 11 | 41.131 | 27.357 | -22.559 | 1.00 | 53.17 | 29.335 | 18.588 | -12.315 | 1.00 | 16.32 | C |
| ATOM | 19 | CB | LVS | A | 11 | 39.336 | 29.487 | -23.461 | 1.00 | 59.88 | 28.619 | 18.942 | -11.811 | 1.00 | 16.99 | C |
| ATOM | 20 | CC | LVS | A | 11 | 38.836 | 28.831 | -23.964 | 1.00 | 61.16 | 29.259 | 19.312 | -10.911 | 1.00 | 15.91 | C |
| ATOM | 21 | CD | LVS | A | 11 | 37.618 | 30.622 | -24.865 | 1.00 | 62.37 | 29.896 | 17.352 | -12.354 | 1.00 | 14.31 | C |
| ATOM | 22 | CE | LVS | A | 11 | 36.592 | 31.732 | -24.665 | 1.00 | 64.46 | 28.785 | 16.393 | -12.693 | 1.00 | 12.71 | C |
| ATOM | 23 | NZ | LVS | A | 11 | 35.966 | 31.661 | -23.353 | 1.00 | 65.85 | 38.980 | 16.974 | -13.474 | 1.00 | 15.85 | C |
| ATOM | 24 | N | ILE | A | 12 | 40.738 | 28.977 | -28.462 | 1.00 | 46.69 | 27.289 | 18.817 | -11.014 | 1.00 | 16.12 | N |
| ATOM | 25 | CA | ILE | A | 12 | 40.865 | 26.731 | -19.945 | 1.00 | 41.32 | 26.438 | 19.322 | -9.949 | 1.00 | 10.883 | C |
| ATOM | 26 | C | ILE | A | 12 | 39.470 | 26.180 | -20.031 | 1.00 | 38.16 | 25.212 | 18.444 | -9.785 | 1.00 | 17.44 | C |
| ATOM | 27 | O | ILE | A | 12 | 38.510 | 26.890 | -19.759 | 1.00 | 37.75 | 25.012 | 17.524 | -10.585 | 1.00 | 18.53 | O |
| ATOM | 28 | CB | ILE | A | 12 | 41.293 | 26.663 | -18.468 | 1.00 | 41.62 | 28.161 | 20.768 | -10.891 | 1.00 | 16.68 | C |
| ATOM | 29 | CG1 | ILE | A | 12 | 42.663 | 27.279 | -18.246 | 1.00 | 43.69 | 26.438 | 19.322 | -9.949 | 1.00 | 10.883 | C |
| ATOM | 30 | CG2 | ILE | A | 12 | 41.377 | 25.188 | -18.060 | 1.00 | 42.11 | 25.617 | 23.036 | -10.575 | 1.00 | 17.40 | C |
| ATOM | 31 | CD1 | ILE | A | 12 | 43.236 | 27.086 | -16.815 | 1.00 | 42.37 | 23.754 | 21.431 | -10.378 | 1.00 | 16.91 | C |
| ATOM | 32 | N | LVS | A | 13 | 39.303 | 24.916 | -20.484 | 1.00 | 36.89 | 24.409 | 18.982 | -8.762 | 1.00 | 15.77 | N |
| ATOM | 33 | CA | LVS | A | 13 | 37.973 | 24.393 | -20.494 | 1.00 | 35.27 | 23.175 | 17.995 | -8.584 | 1.00 | 14.82 | C |
| ATOM | 34 | C | LVS | A | 13 | 37.680 | 23.777 | -19.338 | 1.00 | 33.46 | 22.185 | 19.042 | -8.012 | 1.00 | 17.55 | C |
| ATOM | 35 | O | LVS | A | 13 | 38.529 | 23.060 | -18.585 | 1.00 | 34.67 | 22.536 | 19.924 | -7.218 | 1.00 | 17.69 | O |
| ATOM | 36 | CB | LVS | A | 13 | 37.851 | 23.311 | -21.554 | 1.00 | 34.84 | 23.405 | 16.974 | -7.440 | 1.00 | 13.63 | C |
| ATOM | 37 | CC | LVS | A | 13 | 37.848 | 24.027 | -22.853 | 1.00 | 38.81 | 23.984 | 15.656 | -7.949 | 1.00 | 17.84 | C |

- Atomistic simulations provide **too much** information
- It is hard to decipher the essential features in structurally-complex compounds, materials, proteins, etc.



- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!



- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!

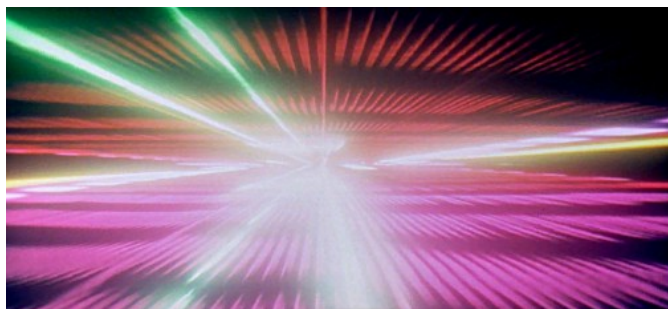


Image from: 2001, A Space Odyssey

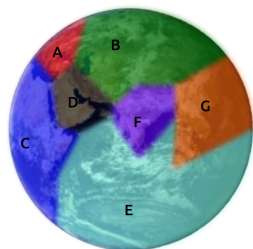
- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!



Pattern Recognition vs Nonlinear Maps



- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!



5

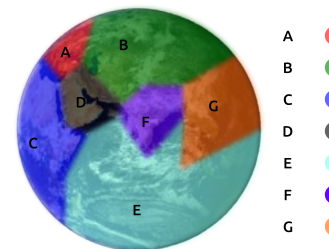
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Pattern Recognition vs Nonlinear Maps



- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!



5

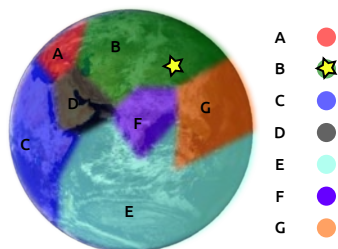
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Pattern Recognition vs Nonlinear Maps



- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!



5

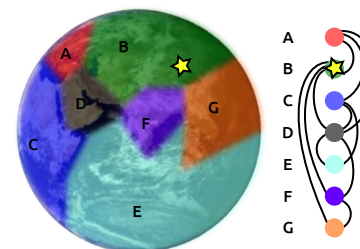
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Pattern Recognition vs Nonlinear Maps



- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!



5

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Pattern Recognition vs Nonlinear Maps



- We can describe a complex molecular structure as a point in a high-dimensional space.
- Clustering/pattern recognition partitions configuration space into regions that can be assigned to (meta) stable structures
- (Non-linear) dimensionality reduction corresponds to making a low-dimensional map: more informative!



5

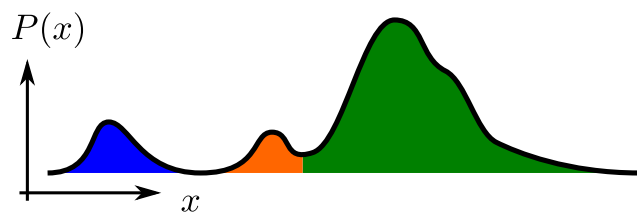
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Mode Analysis of a Distribution



- A natural way of recognizing patterns in a distribution is to identify its modes, and the basin of attraction of each mode.
- One can then fit a simple Gaussian model (with fixed centers), and use posteriors to assign fingerprints to each cluster



Piero Gasparotto & Michele Ceriotti, JCP 174110, 141 (2014)

6

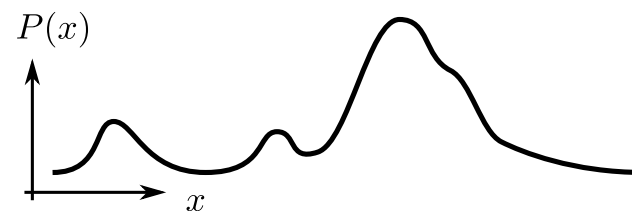
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Mode Analysis of a Distribution



- A natural way of recognizing patterns in a distribution is to identify its modes, and the basin of attraction of each mode.
- One can then fit a simple Gaussian model (with fixed centers), and use posteriors to assign fingerprints to each cluster



Piero Gasparotto & Michele Ceriotti, JCP 174110, 141 (2014)

6

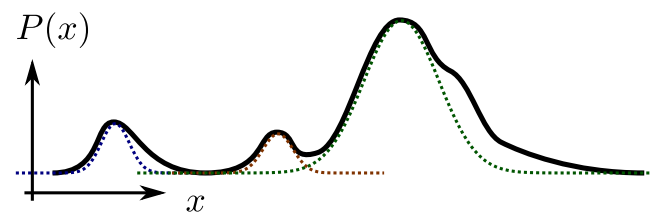
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Mode Analysis of a Distribution



- A natural way of recognizing patterns in a distribution is to identify its modes, and the basin of attraction of each mode.
- One can then fit a simple Gaussian model (with fixed centers), and use posteriors to assign fingerprints to each cluster



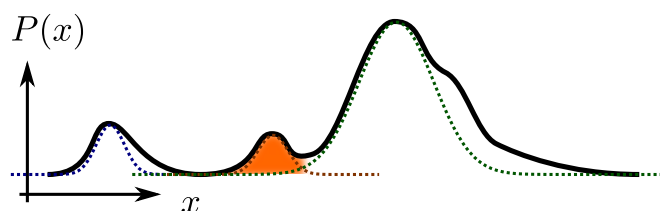
Piero Gasparotto & Michele Ceriotti, JCP 174110, 141 (2014)

6

Michele Ceriotti EPFL/IMX/COSMO

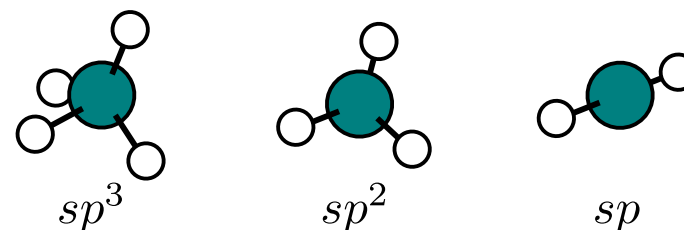
Representing and understanding patterns in materials and mol

- A natural way of recognizing patterns in a distribution is to identify its modes, and the basin of attraction of each mode.
- One can then fit a simple Gaussian model (with fixed centers), and use posteriors to assign fingerprints to each cluster

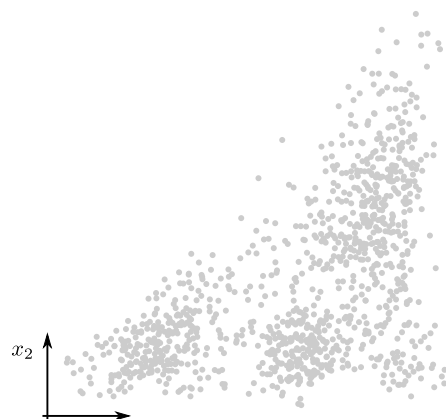


Piero Gasparotto & Michele Ceriotti, JCP 174110, 141 (2014)

- We still need an effective high-dimensional description to start with
- “Chemical intuition” builds on recognizing recurring patterns in atomic configurations
- Automatic scheme to single out structural motifs in atomistic simulations

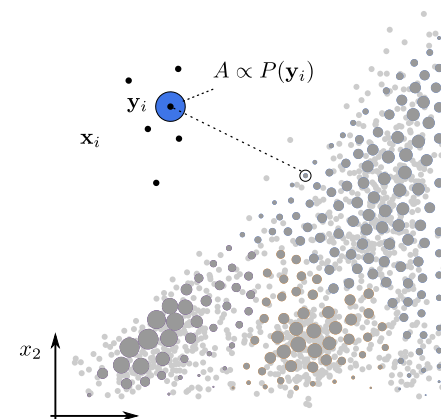


- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



Piero Gasparotto & Michele Ceriotti, JCP 174110, 141 (2014)

- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space

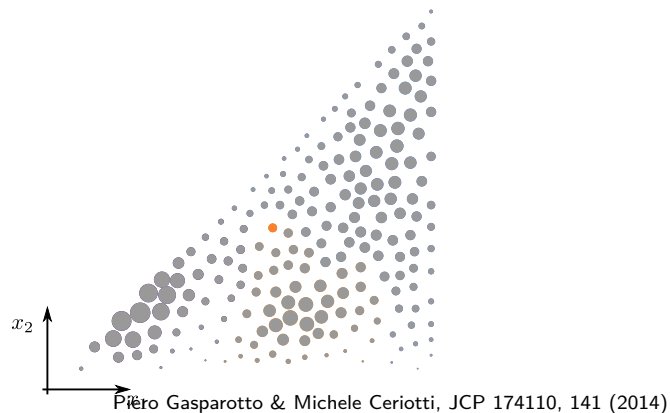


Piero Gasparotto & Michele Ceriotti, JCP 174110, 141 (2014)

Probabilistic Analysis of Molecular Motifs



- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



8

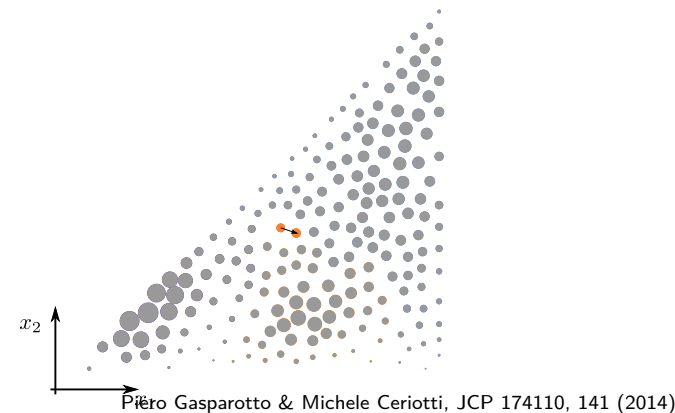
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Probabilistic Analysis of Molecular Motifs



- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



8

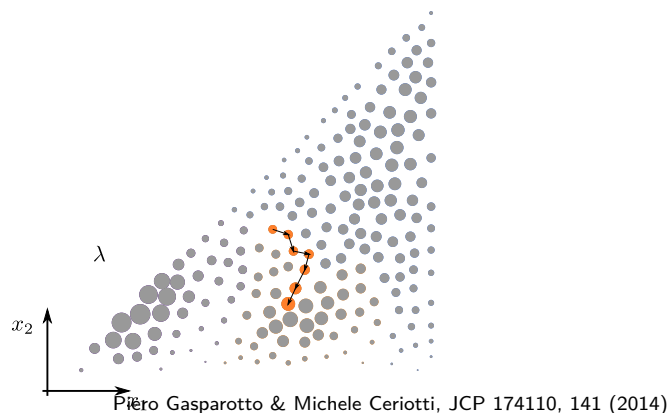
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Probabilistic Analysis of Molecular Motifs



- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



8

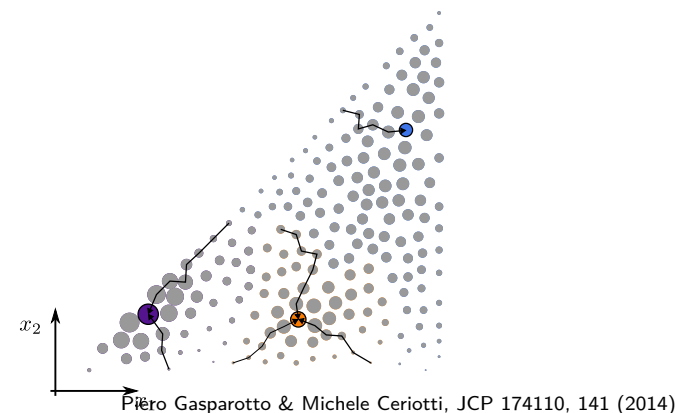
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Probabilistic Analysis of Molecular Motifs



- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



8

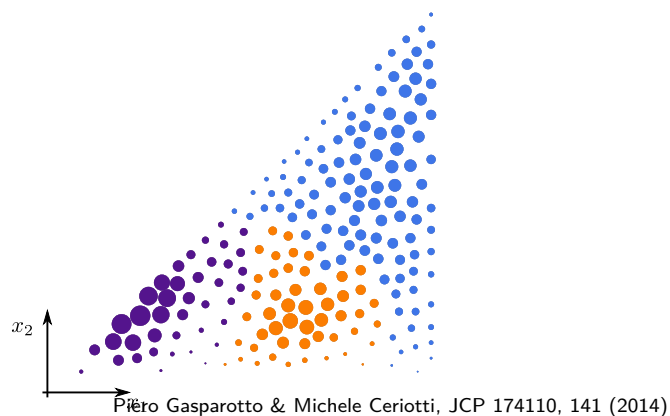
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Probabilistic Analysis of Molecular Motifs



- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



8

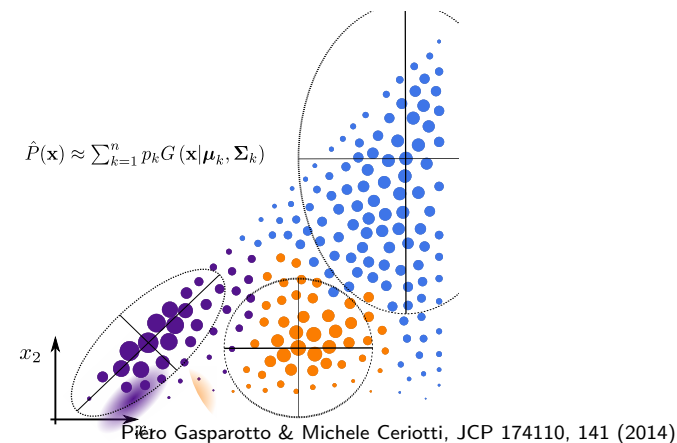
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Probabilistic Analysis of Molecular Motifs



- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



8

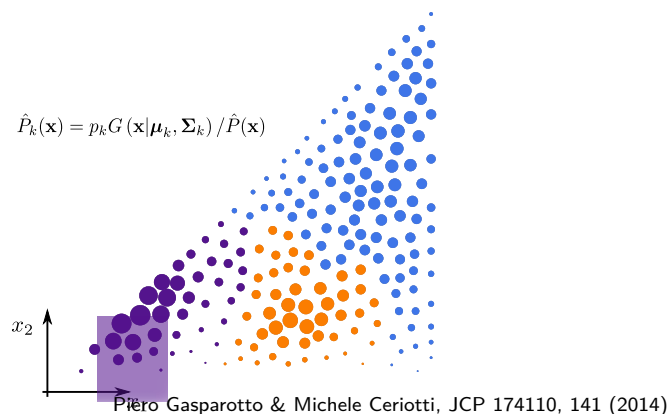
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Probabilistic Analysis of Molecular Motifs



- Evaluate the probability distribution of molecular structures
- Cluster it around the modes of the distribution
- Naturally gives a fuzzy and continuous partitioning of configuration space



8

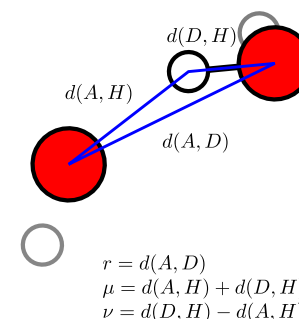
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

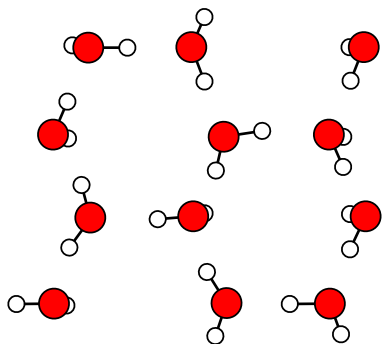
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

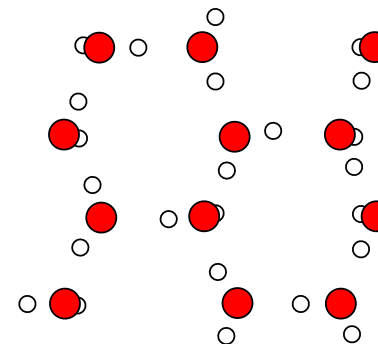
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

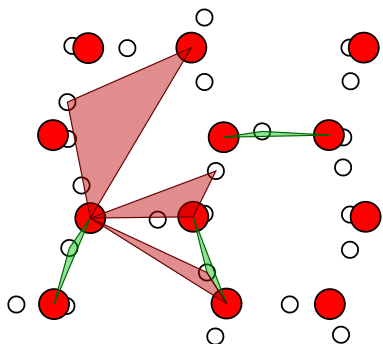
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

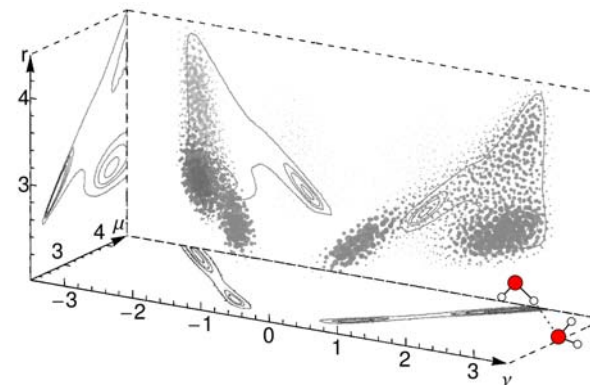
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

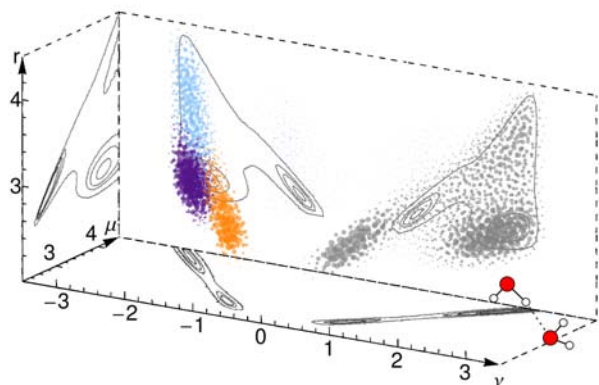
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

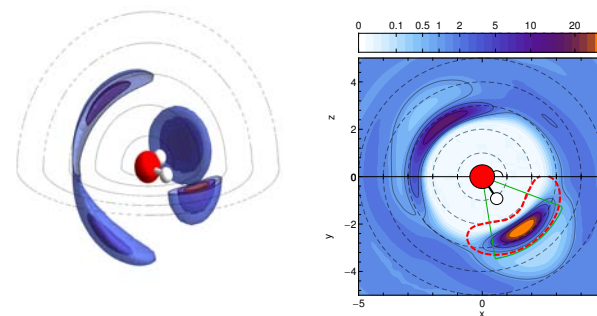
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

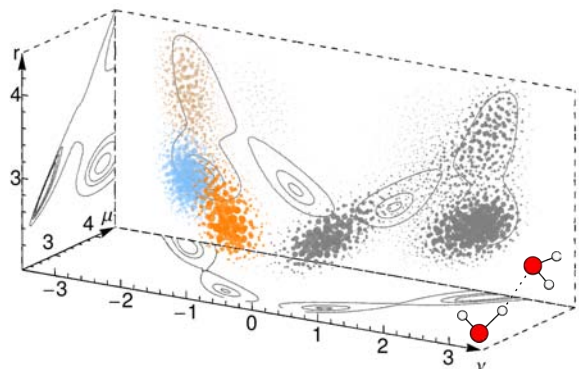
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy



9

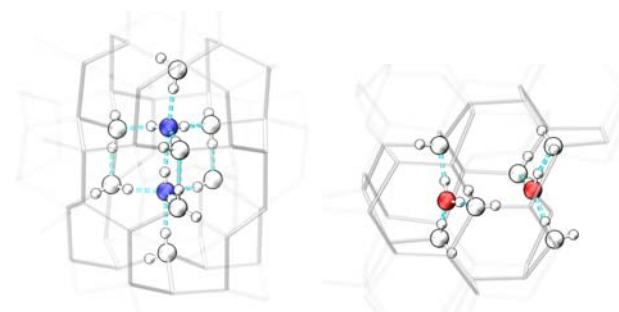
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

An Agnostic Definition of the H-Bond



- Most general description of a H-bond geometry: 3 distances
- PAMM recognizes multiple modes - one corresponds to the H-bond
- PAMM H-bond fingerprints can be used as HB counts, but are adaptive, unbiased and fuzzy

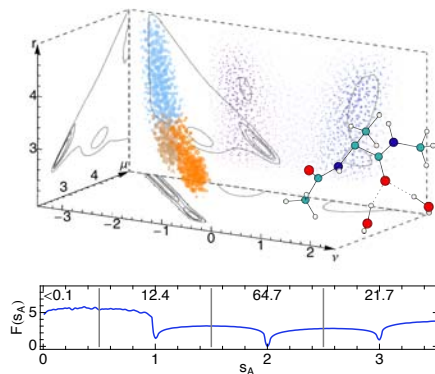


9

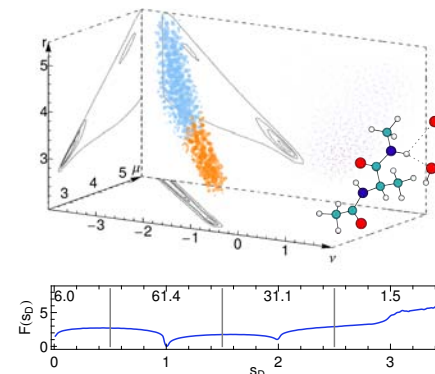
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

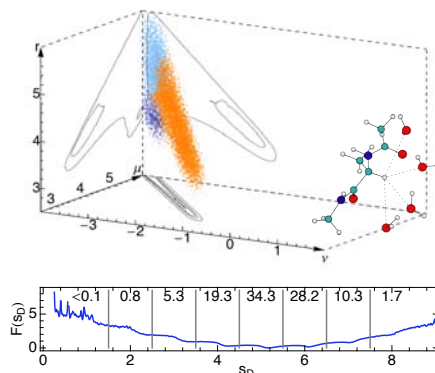
- Different groups should be treated with a different geometric definition of HB
- PAMM provides data-driven, unbiased procedure to determine the structures that can be labeled as bonded



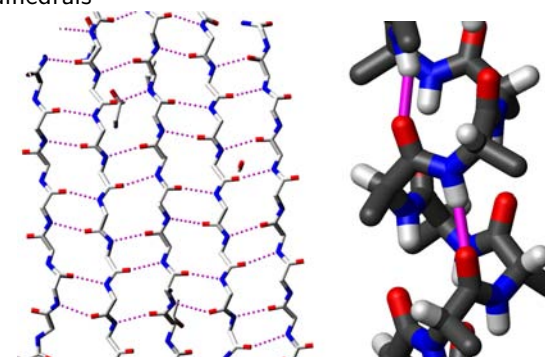
- Different groups should be treated with a different geometric definition of HB
- PAMM provides data-driven, unbiased procedure to determine the structures that can be labeled as bonded



- Different groups should be treated with a different geometric definition of HB
- PAMM provides data-driven, unbiased procedure to determine the structures that can be labeled as bonded



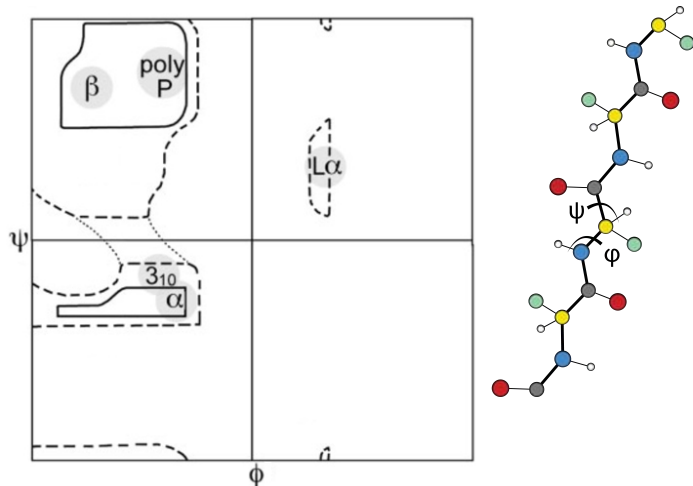
- Secondary structure is induced by H-bonds, but correlates strongly with backbone dihedrals



Machine-learning the Ramachandran plot



- Secondary structure is induced by H-bonds, but correlates strongly with backbone dihedrals



11

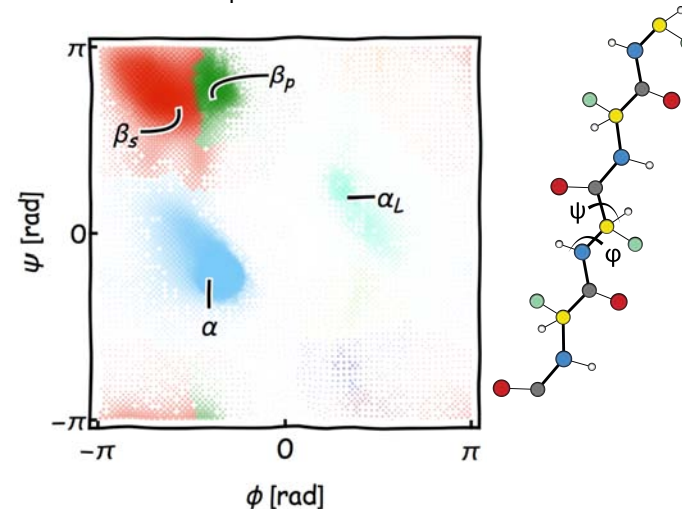
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Machine-learning the Ramachandran plot



- Use data from the PDB, and "learn" with PAMM the stable patterns of proteins in dihedral space



11

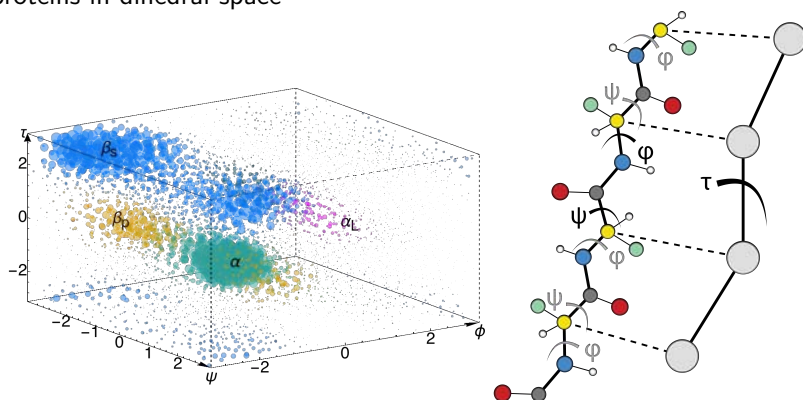
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Machine-learning the Ramachandran plot



- Use data from the PDB, and "learn" with PAMM the stable patterns of proteins in dihedral space



11

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Describing structural complexity



- We are looking for **collective variables** that can describe structural complexity **globally**
 - Discriminate between different structures
 - Follow the system across transitions
- This is not only important for post-processing
 - Good CVs make for better transition-state approximation to the rate
 - Biased MD requires coarse-grained but thorough description of the problem
- Finding these variables is time-consuming and error-prone: can we **automate** the process?

12

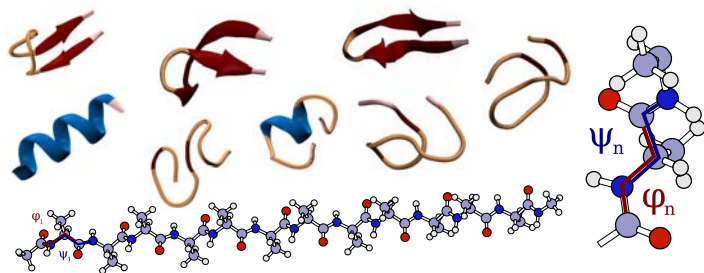
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Dimensionality reduction



- We can describe a complex atomistic structure as a point in a high-dimensional space. Then finding CVs means finding a low-dimensional **map** to describe the accessible configurations!
 - Take a set of configurations \Rightarrow high-dim. **landmark points**
 - Define a measure of dissimilarity between the points
 - Arrange low-dim. points so that the dissimilarities are preserved
 - Locate other configurations with an **out-of-sample embedding**



13

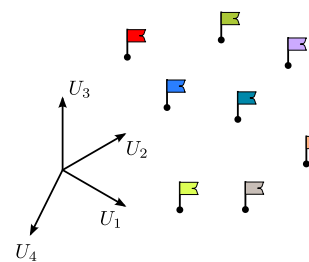
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Dimensionality reduction



- We can describe a complex atomistic structure as a point in a high-dimensional space. Then finding CVs means finding a low-dimensional **map** to describe the accessible configurations!
 - Take a set of configurations \Rightarrow high-dim. **landmark points**
 - Define a measure of dissimilarity between the points
 - Arrange low-dim. points so that the dissimilarities are preserved
 - Locate other configurations with an **out-of-sample embedding**



13

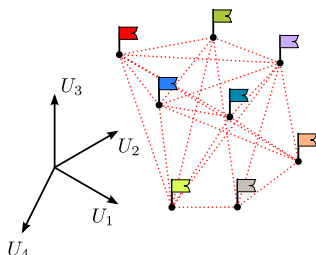
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Dimensionality reduction



- We can describe a complex atomistic structure as a point in a high-dimensional space. Then finding CVs means finding a low-dimensional **map** to describe the accessible configurations!
 - Take a set of configurations \Rightarrow high-dim. **landmark points**
 - Define a measure of dissimilarity between the points
 - Arrange low-dim. points so that the dissimilarities are preserved
 - Locate other configurations with an **out-of-sample embedding**



13

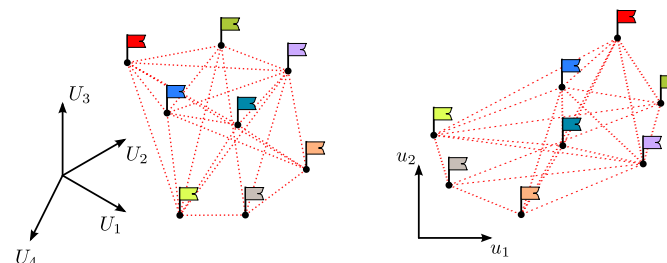
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Dimensionality reduction



- We can describe a complex atomistic structure as a point in a high-dimensional space. Then finding CVs means finding a low-dimensional **map** to describe the accessible configurations!
 - Take a set of configurations \Rightarrow high-dim. **landmark points**
 - Define a measure of dissimilarity between the points
 - Arrange low-dim. points so that the dissimilarities are preserved
 - Locate other configurations with an **out-of-sample embedding**

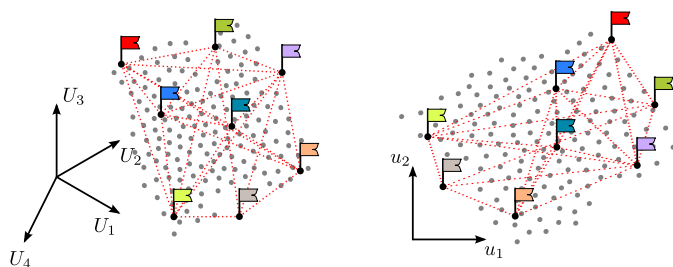


13

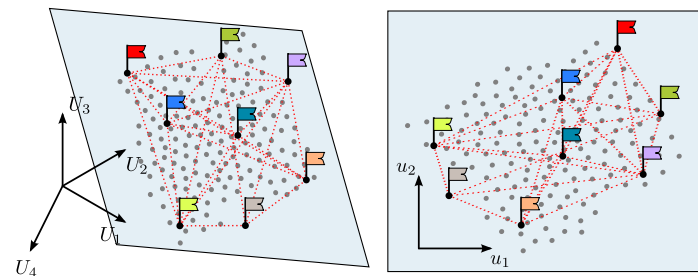
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

- We can describe a complex atomistic structure as a point in a high-dimensional space. Then finding CVs means finding a low-dimensional **map** to describe the accessible configurations!
 - Take a set of configurations \Rightarrow high-dim. **landmark points**
 - Define a measure of dissimilarity between the points
 - Arrange low-dim. points so that the dissimilarities are preserved
 - Locate other configurations with an **out-of-sample embedding**



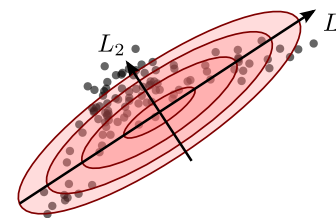
- We can describe a complex atomistic structure as a point in a high-dimensional space. Then finding CVs means finding a low-dimensional **map** to describe the accessible configurations!
 - Take a set of configurations \Rightarrow high-dim. **landmark points**
 - Define a measure of dissimilarity between the points
 - Arrange low-dim. points so that the dissimilarities are preserved
 - Locate other configurations with an **out-of-sample embedding**



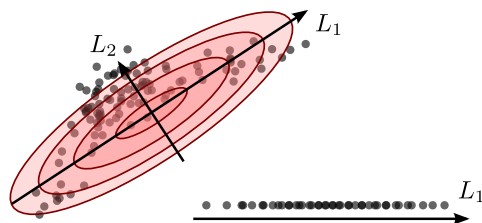
- Principal component analysis: assumes that the “important” coordinates are the linear combinations with the largest variance
 - $\{X_i\}$ are N vectors in D dimensions. Let \mathbf{X} be the $N \times D$ matrix with the X_i as rows.
 - Define the $N \times N$ centering matrix $H_{ij} = \delta_{ij} - \frac{1}{n}$
 - Define the covariance matrix $\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{H} \mathbf{X}$,
 $C_{ij} = \frac{1}{n} \sum_k (X_k - \bar{X})_i (X_k - \bar{X})_j$
 - Pick the d eigenvectors P_i associated with the largest eigenvalues λ_i and use them as the rows of a **linear projector P**.
 - The low-dimensional projections are $x_i = \mathbf{P}^T X_i$



- Principal component analysis: assumes that the “important” coordinates are the linear combinations with the largest variance
 - $\{X_i\}$ are N vectors in D dimensions. Let \mathbf{X} be the $N \times D$ matrix with the X_i as rows.
 - Define the $N \times N$ centering matrix $H_{ij} = \delta_{ij} - \frac{1}{n}$
 - Define the covariance matrix $\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{H} \mathbf{X}$,
 $C_{ij} = \frac{1}{n} \sum_k (X_k - \bar{X})_i (X_k - \bar{X})_j$
 - Pick the d eigenvectors P_i associated with the largest eigenvalues λ_i and use them as the rows of a **linear projector P**.
 - The low-dimensional projections are $x_i = \mathbf{P}^T X_i$



- Principal component analysis: assumes that the “important” coordinates are the linear combinations with the largest variance
 - $\{X_i\}$ are N vectors in D dimensions. Let \mathbf{X} be the $N \times D$ matrix with the X_i as rows.
 - Define the $N \times N$ centering matrix $H_{ij} = \delta_{ij} - \frac{1}{n}$
 - Define the covariance matrix $\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{H} \mathbf{X}$,
 $C_{ij} = \frac{1}{n} \sum_k (X_k - \bar{X})_i (X_k - \bar{X})_j$
 - Pick the d eigenvectors P_i associated with the largest eigenvalues λ_i and use them as the rows of a **linear projector P**.
 - The low-dimensional projections are $x_i = \mathbf{P}^T X_i$



- Measures of four morphological features (petal/sepal width/length) of a total of 150 samples of three species of iris
- 4D dataset, strong correlation between the indicators, but also spread within one species
- Apply PCA and classical MDS – equivalent modulo a rotation.
 - Clearly clustered in 2D.
 - *Versicolor* and *Virginica* are pretty close...



Iris Setosa



Iris Versicolor



Iris Virginica

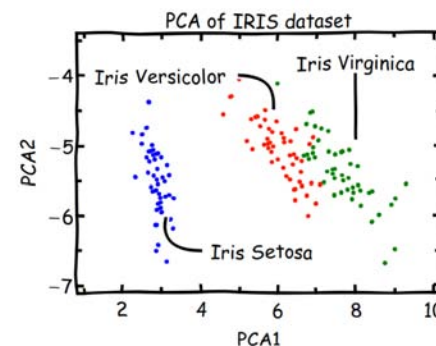
- A literal implementation of the general idea of dimensionality reduction
 - define $\Delta_{ij} = \Delta(X_i, X_j)$ where $\Delta(X, Y)$ is a measure of similarity between points in D dimensions
 - find d -dimensional projections $\{x_i\}$ minimizing

$$\chi^2 = \sum_{ij} (\Delta_{ij} - |x_i - x_j|)^2$$

- *Classical* MDS turns this iterative optimization in an eigenvalue problem
 - Define $S_{ij} = \Delta_{ij}^2$ and $\mathbf{B} = -\frac{1}{2} \mathbf{H} \mathbf{S} \mathbf{H}$. Note that $\mathbf{B} = (\mathbf{H} \mathbf{X})(\mathbf{H} \mathbf{X})^T$
 - Compute the largest d eigenvalues of \mathbf{B} , λ_i and the eigenvectors \mathbf{V}_i
 - Make the $n \times d$ matrix whose columns are $\sqrt{\lambda_i} \mathbf{V}_i$. The rows are the x_i low-dimensional projections
- If $\Delta(X_i, X_j)$ is the Euclidean norm, classical MDS is the best *linear* projection preserving the squared distances. It corresponds to PCA, but it is more easily generalized to different dissimilarities

Cox & Cox, *Multidimensional Scaling* (CRC Press, 2010)

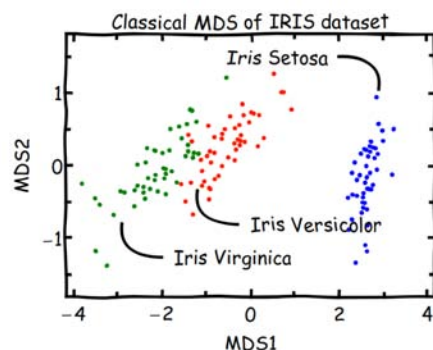
- Measures of four morphological features (petal/sepal width/length) of a total of 150 samples of three species of iris
- 4D dataset, strong correlation between the indicators, but also spread within one species
- Apply PCA and classical MDS – equivalent modulo a rotation.
 - Clearly clustered in 2D.
 - *Versicolor* and *Virginica* are pretty close...



The IRIS dataset



- Measures of four morphological features (petal/sepal width/length) of a total of 150 samples of three species of iris
- 4D dataset, strong correlation between the indicators, but also spread within one species
- Apply PCA and classical MDS – equivalent modulo a rotation.
 - Clearly clustered in 2D.
 - *Versicolor* and *Virginica* are pretty close...



16

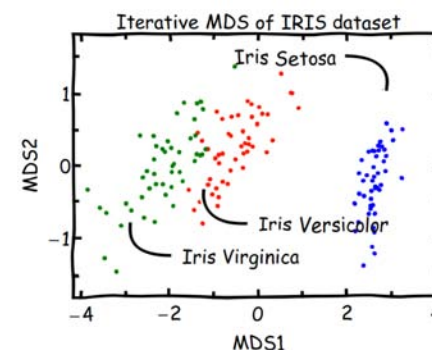
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The IRIS dataset



- Measures of four morphological features (petal/sepal width/length) of a total of 150 samples of three species of iris
- 4D dataset, strong correlation between the indicators, but also spread within one species
- Apply PCA and classical MDS – equivalent modulo a rotation.
 - Clearly clustered in 2D. Iterative MDS does not make a big difference
 - *Versicolor* and *Virginica* are pretty close...



16

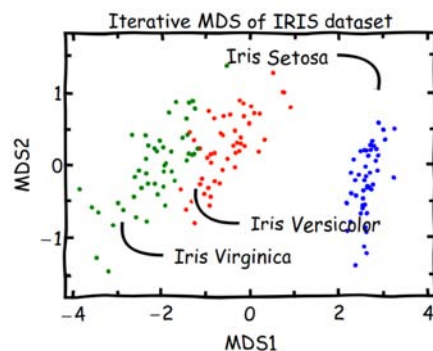
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The IRIS dataset



- Measures of four morphological features (petal/sepal width/length) of a total of 150 samples of three species of iris
- 4D dataset, strong correlation between the indicators, but also spread within one species
- Apply PCA and classical MDS – equivalent modulo a rotation.
 - Clearly clustered in 2D. Iterative MDS does not make a big difference
 - *Versicolor* and *Virginica* are pretty close...



16

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The IRIS dataset



- Measures of four morphological features (petal/sepal width/length) of a total of 150 samples of three species of iris
- 4D dataset, strong correlation between the indicators, but also spread within one species
- Apply PCA and classical MDS – equivalent modulo a rotation.
 - Clearly clustered in 2D. Iterative MDS does not make a big difference
 - *Versicolor* and *Virginica* are pretty close... and they look quite similar!



Iris Setosa



Iris Versicolor



Iris Virginica

16

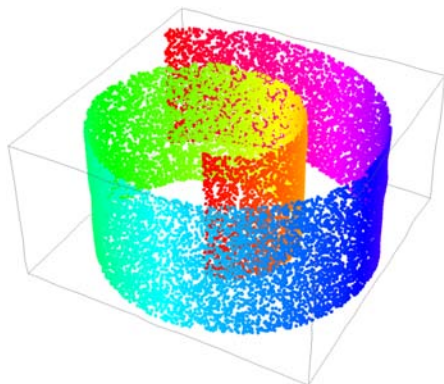
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The Swiss roll dataset



- A 2D manifold embedded in three-dimensional space
- PCA cannot capture the low-dimensional structure of the manifold, because it is just a **linear** projection!
 - Linear methods work when data lie (almost) on a plane
 - One would need a method that can deal with a curved manifold which is only **locally** linear



17

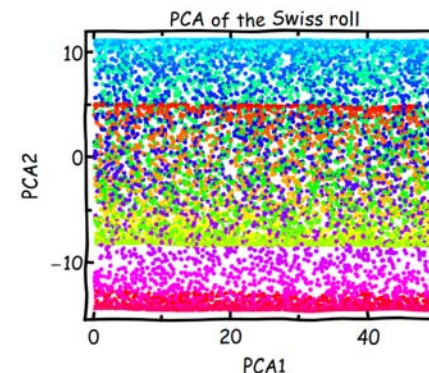
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The Swiss roll dataset



- A 2D manifold embedded in three-dimensional space
- PCA cannot capture the low-dimensional structure of the manifold, because it is just a **linear** projection!
 - Linear methods work when data lie (almost) on a plane
 - One would need a method that can deal with a curved manifold which is only **locally** linear



17

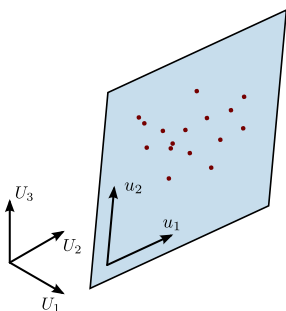
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The Swiss roll dataset



- A 2D manifold embedded in three-dimensional space
- PCA cannot capture the low-dimensional structure of the manifold, because it is just a **linear** projection!
 - Linear methods work when data lie (almost) on a plane
 - One would need a method that can deal with a curved manifold which is only **locally** linear



17

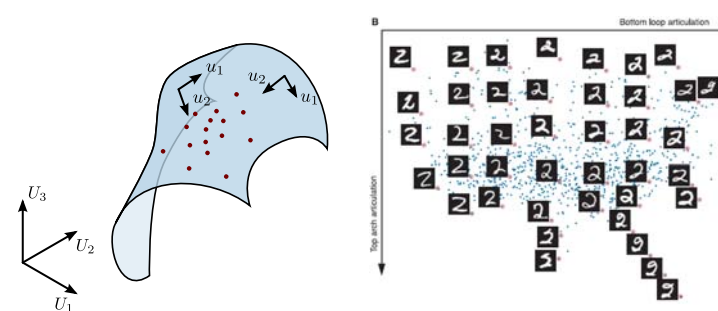
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The Swiss roll dataset



- A 2D manifold embedded in three-dimensional space
- PCA cannot capture the low-dimensional structure of the manifold, because it is just a **linear** projection!
 - Linear methods work when data lie (almost) on a plane
 - One would need a method that can deal with a curved manifold which is only **locally** linear



17

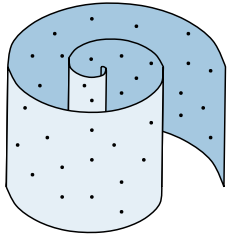
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Non-Linear DR: ISOMAP



- A family of methods introduces non-linearity in the dissimilarity metric
- ISOMAP defines point-point distances based on geodesics
 - Approximate geodesics by hopping between neighbours
 - Problem: very sensitive to uneven sampling and noise. One “wrong” neighbor detection can mess up geodesics completely.



Tenenbaum et al., Science (2000)

18

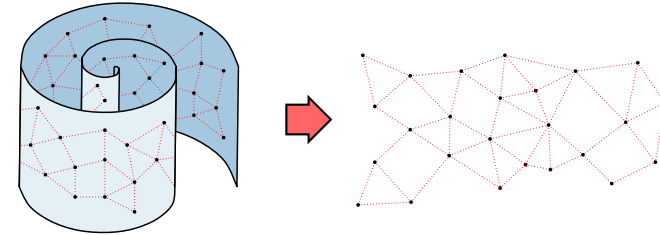
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Non-Linear DR: ISOMAP



- A family of methods introduces non-linearity in the dissimilarity metric
- ISOMAP defines point-point distances based on geodesics
 - Approximate geodesics by hopping between neighbours
 - Problem: very sensitive to uneven sampling and noise. One “wrong” neighbor detection can mess up geodesics completely.



Tenenbaum et al., Science (2000)

18

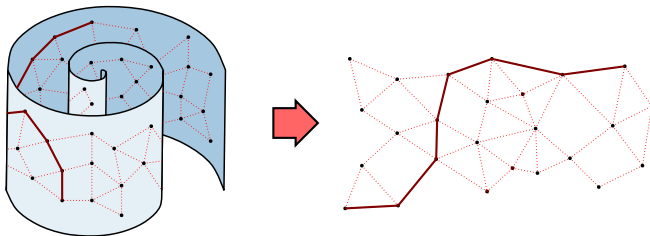
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Non-Linear DR: ISOMAP



- A family of methods introduces non-linearity in the dissimilarity metric
- ISOMAP defines point-point distances based on geodesics
 - Approximate geodesics by hopping between neighbours
 - Problem: very sensitive to uneven sampling and noise. One “wrong” neighbor detection can mess up geodesics completely.



Tenenbaum et al., Science (2000)

18

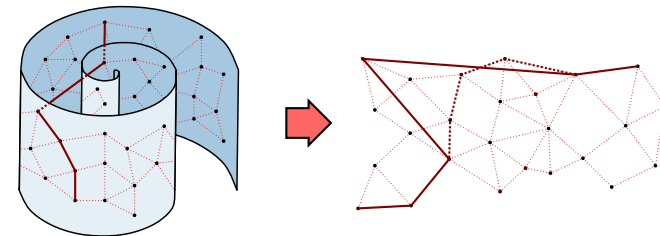
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Non-Linear DR: ISOMAP



- A family of methods introduces non-linearity in the dissimilarity metric
- ISOMAP defines point-point distances based on geodesics
 - Approximate geodesics by hopping between neighbours
 - Problem: very sensitive to uneven sampling and noise. One “wrong” neighbor detection can mess up geodesics completely.



Tenenbaum et al., Science (2000)

18

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Non-Linear DR: ISOMAP



- A family of methods introduces non-linearity in the dissimilarity metric
- ISOMAP defines point-point distances based on geodesics
 - Approximate geodesics by hopping between neighbours
 - Problem: very sensitive to uneven sampling and noise. One “wrong” neighbor detection can mess up geodesics completely.
- ① Define neighbour relations between points (k nearest neighbours or points within ϵ)
- ② Compute the graph distance matrix as an approximant to geodesics
- ③ Run classical MDS.

Tenenbaum et al., Science (2000)

18

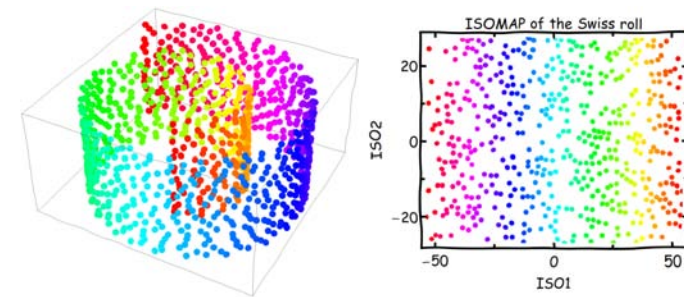
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

ISOMAP and the Swiss roll



- ISOMAP works very well for the Swiss roll, as it identifies beautifully the manifold directions
- It is however very sensitive to noise, and to uneven sampling. When it fails, it fails dramatically!



19

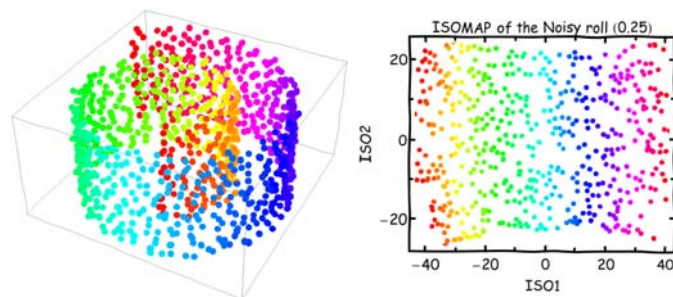
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

ISOMAP and the Swiss roll



- ISOMAP works very well for the Swiss roll, as it identifies beautifully the manifold directions
- It is however very sensitive to noise, and to uneven sampling. When it fails, it fails dramatically!



19

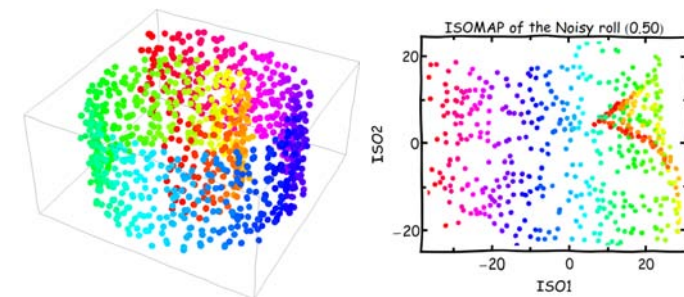
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

ISOMAP and the Swiss roll



- ISOMAP works very well for the Swiss roll, as it identifies beautifully the manifold directions
- It is however very sensitive to noise, and to uneven sampling. When it fails, it fails dramatically!

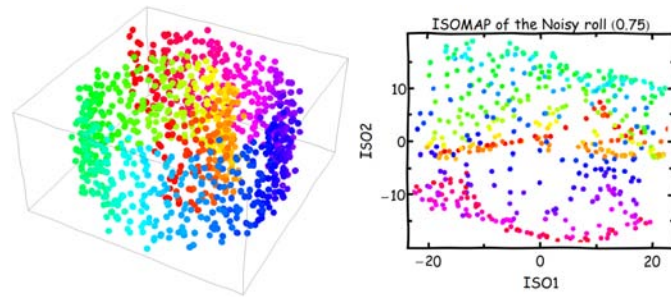


19

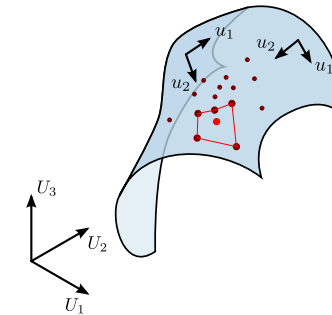
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

- ISOMAP works very well for the Swiss roll, as it identifies beautifully the manifold directions
- It is however very sensitive to noise, and to uneven sampling. When it fails, it fails dramatically!

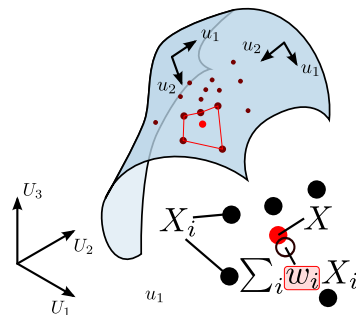


- Use the fact that if the manifold is locally flat each point can be expressed as a combination of its neighbors.
 - 1 Determine a neighborhood of each point X , and the weights w_i that best match X and its embedding
 - 2 Determine the low-dimensional points such that for each point, x an its embedding are as close as possible *keeping the weights fixed*
 - 3 This is again formulated as an eigenvalue problem



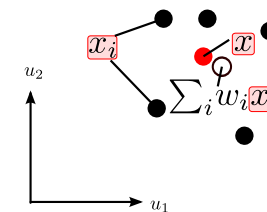
Roweis, Saul, Science (2000)

- Use the fact that if the manifold is locally flat each point can be expressed as a combination of its neighbors.
 - 1 Determine a neighborhood of each point X , and the weights w_i that best match X and its embedding
 - 2 Determine the low-dimensional points such that for each point, x an its embedding are as close as possible *keeping the weights fixed*
 - 3 This is again formulated as an eigenvalue problem



Roweis, Saul, Science (2000)

- Use the fact that if the manifold is locally flat each point can be expressed as a combination of its neighbors.
 - 1 Determine a neighborhood of each point X , and the weights w_i that best match X and its embedding
 - 2 Determine the low-dimensional points such that for each point, x an its embedding are as close as possible *keeping the weights fixed*
 - 3 This is again formulated as an eigenvalue problem

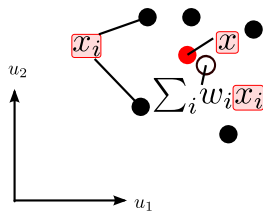


Roweis, Saul, Science (2000)

Locally Linear Embedding



- Use the fact that if the manifold is locally flat each point can be expressed as a combination of its neighbors.
 - 1 Determine a neighborhood of each point X , and the weights w_i that best match X and its embedding
 - 2 Determine the low-dimensional points such that for each point, x and its embedding are as close as possible *keeping the weights fixed*
 - 3 This is again formulated as an eigenvalue problem



Roweis, Saul, Science (2000)

20

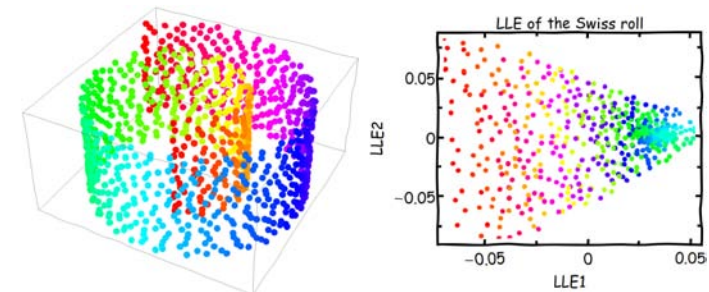
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

LLE and the Swiss roll



- LLE seems not to work ... 2-but some higher variables do capture the manifold structure
- Growing noise destabilizes the embedding, and shifts to even higher-order LLE vectors the reasonable map



21

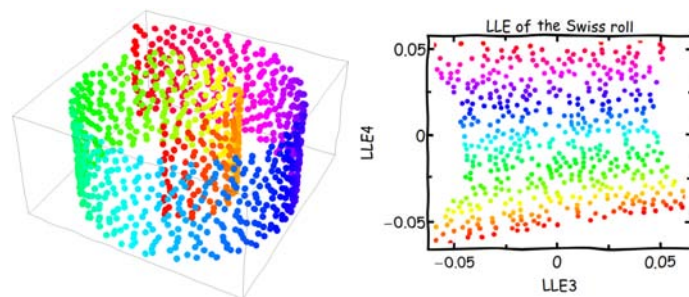
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

LLE and the Swiss roll



- LLE seems not to work ... 2-but some higher variables do capture the manifold structure
- Growing noise destabilizes the embedding, and shifts to even higher-order LLE vectors the reasonable map



21

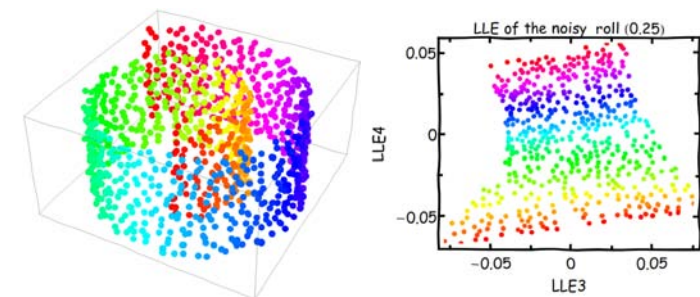
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

LLE and the Swiss roll



- LLE seems not to work ... 2-but some higher variables do capture the manifold structure
- Growing noise destabilizes the embedding, and shifts to even higher-order LLE vectors the reasonable map

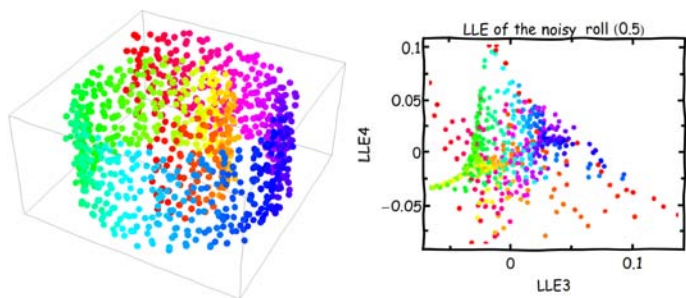


21

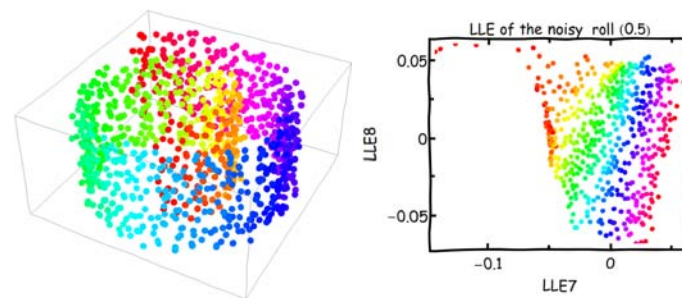
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

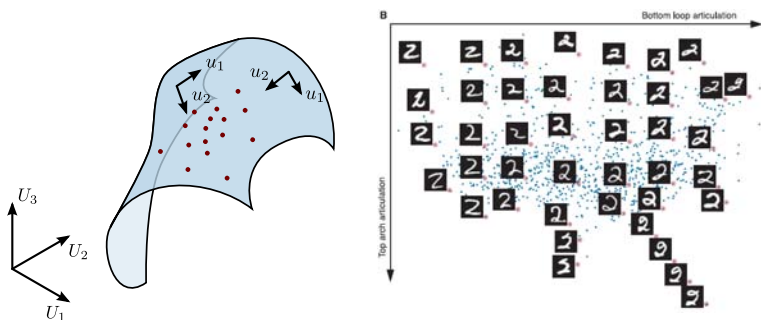
- LLE seems not to work ... 2-but some higher variables do capture the manifold structure
- Growing noise destabilizes the embedding, and shifts to even higher-order LLE vectors the reasonable map



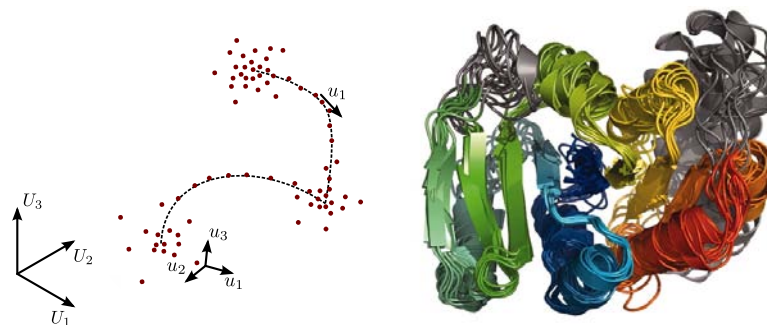
- LLE seems not to work ... 2-but some higher variables do capture the manifold structure
- Growing noise destabilizes the embedding, and shifts to even higher-order LLE vectors the reasonable map



- Non-linear dimensionality reduction algorithms:
 - Describe curved, "locally-flat" manifolds
 - Developed by the CS community (image recognition)
 - Attempts to apply to chemical problems (PCA, ISOMAP, LLE, ...)
- Atomistic simulations are harder:
 - **Thermal fluctuations** are high-dimensional
 - A network of transition pathways with a complex topology



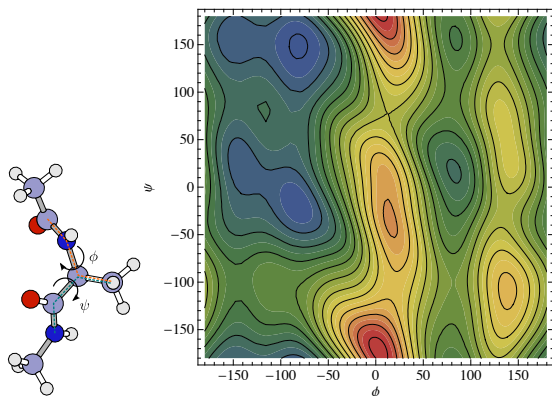
- Non-linear dimensionality reduction algorithms:
 - Describe curved, "locally-flat" manifolds
 - Developed by the CS community (image recognition)
 - Attempts to apply to chemical problems (PCA, ISOMAP, LLE, ...)
- Atomistic simulations are harder:
 - **Thermal fluctuations** are high-dimensional
 - A network of transition pathways with a complex topology



Features of a folding landscape



- The free energy landscape for ala₂ contains low-energy basins and a spider web of transition pathways
- Reconnaissance metadynamics¹ for ala₁₂: similar distribution of points for any pair of dihedrals



¹Tribello, Ceriotti, Parrinello, PNAS 2010

23

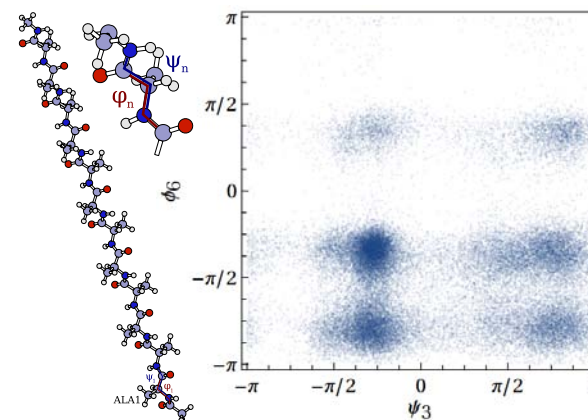
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Features of a folding landscape



- The free energy landscape for ala₂ contains low-energy basins and a spider web of transition pathways
- Reconnaissance metadynamics¹ for ala₁₂: similar distribution of points for any pair of dihedrals



¹Tribello, Ceriotti, Parrinello, PNAS 2010

23

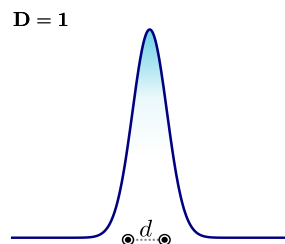
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Matching fluctuations: a space odyssey



- Inherent problem when projecting **full-dimensional features**
- Take for instance the distribution of distances between points taken from a D -dimensional Gaussian
- This is a disaster for distance matching! It is *impossible* to match the distances for a 24-dimensional Gaussian using a 3d Gaussian!



24

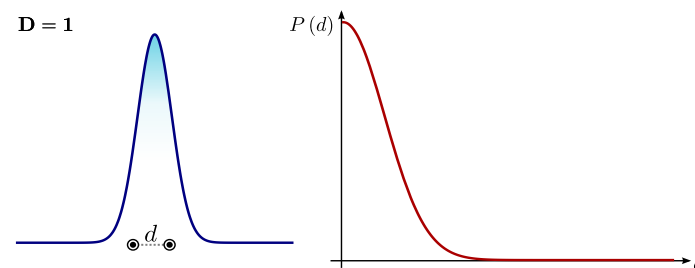
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Matching fluctuations: a space odyssey



- Inherent problem when projecting **full-dimensional features**
- Take for instance the distribution of distances between points taken from a D -dimensional Gaussian
- This is a disaster for distance matching! It is *impossible* to match the distances for a 24-dimensional Gaussian using a 3d Gaussian!



24

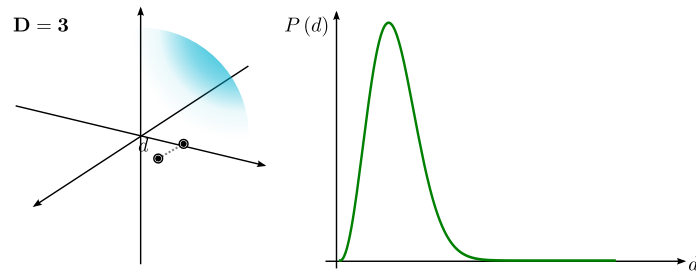
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Matching fluctuations: a space odyssey



- Inherent problem when projecting **full-dimensional features**
- Take for instance the distribution of distances between points taken from a D -dimensional Gaussian
- This is a disaster for distance matching! It is *impossible* to match the distances for a 24-dimensional Gaussian using a 3d Gaussian!



24

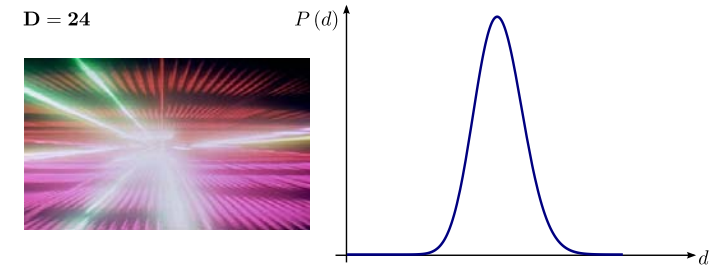
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Matching fluctuations: a space odyssey



- Inherent problem when projecting **full-dimensional features**
- Take for instance the distribution of distances between points taken from a D -dimensional Gaussian
- This is a disaster for distance matching! It is *impossible* to match the distances for a 24-dimensional Gaussian using a 3d Gaussian!



24

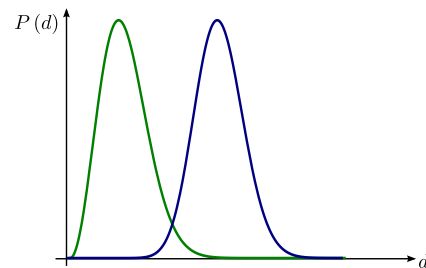
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Matching fluctuations: a space odyssey



- Inherent problem when projecting **full-dimensional features**
- Take for instance the distribution of distances between points taken from a D -dimensional Gaussian
- This is a disaster for distance matching! It is *impossible* to match the distances for a 24-dimensional Gaussian using a 3d Gaussian!



24

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

No need for a perfect map



- Developing a NLDR method which is more robust and suited for trajectory data
 - Basic idea: we don't need a precise, isometric map.
 - We need the computational equivalent of a hand sketched map



25

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

No need for a perfect map



- Developing a NLDR method which is more robust and suited for trajectory data
 - Basic idea: we don't need a precise, isometric map.
 - We need the computational equivalent of a hand sketched map



25

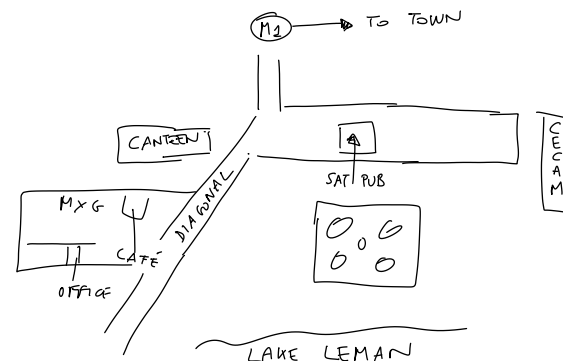
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

No need for a perfect map



- Developing a NLDR method which is more robust and suited for trajectory data
 - Basic idea: we don't need a precise, isometric map.
 - We need the computational equivalent of a hand sketched map



25

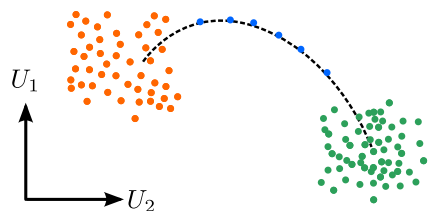
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Proximity matching



- We would like to capture the low-dimensional structure of complex transitions
 - How to deal with full-dimensional thermal fluctuations? Portions of the landscape cannot be projected by matching high and low-dimensional distances.
 - Idea: simpler task, aim for proximity matching: close→close, far→far



Ceriotti, Tribello, Parrinello, PNAS (2011); JCTC (2013)

26

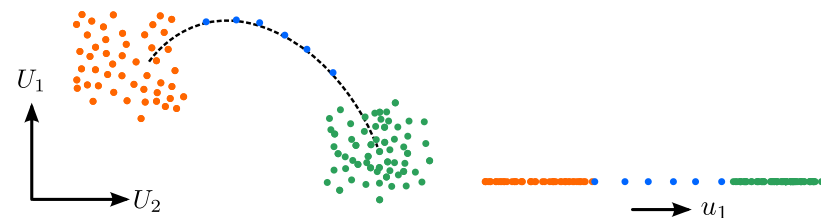
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Proximity matching



- We would like to capture the low-dimensional structure of complex transitions
 - How to deal with full-dimensional thermal fluctuations? Portions of the landscape cannot be projected by matching high and low-dimensional distances.
 - Idea: simpler task, aim for proximity matching: close→close, far→far



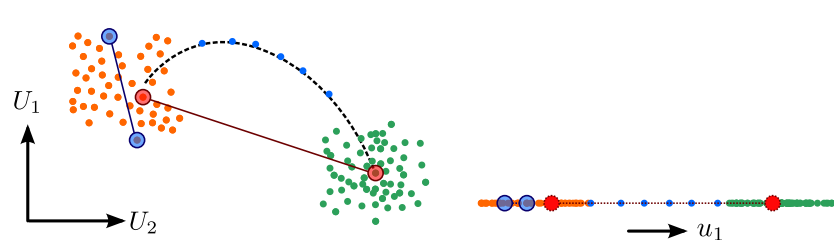
Ceriotti, Tribello, Parrinello, PNAS (2011); JCTC (2013)

26

Michele Ceriotti EPFL/IMX/COSMO

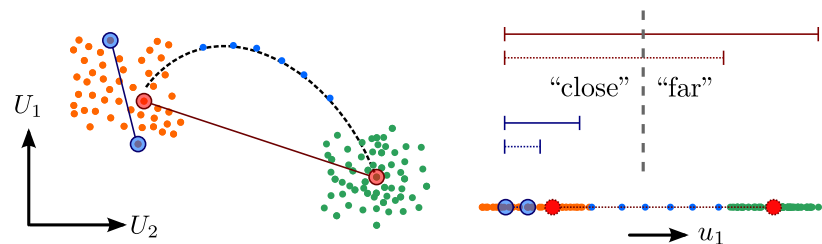
Representing and understanding patterns in materials and mol

- We would like to capture the low-dimensional structure of complex transitions
- How to deal with full-dimensional thermal fluctuations? Portions of the landscape cannot be projected by matching high and low-dimensional distances.
 - Idea: simpler task, aim for **proximity matching**: close↔close, far↔far



Ceriotti, Tribello, Parrinello, PNAS (2011); JCTC (2013)

- We would like to capture the low-dimensional structure of complex transitions
- How to deal with full-dimensional thermal fluctuations? Portions of the landscape cannot be projected by matching high and low-dimensional distances.
 - Idea: simpler task, aim for **proximity matching**: close↔close, far↔far



Ceriotti, Tribello, Parrinello, PNAS (2011); JCTC (2013)

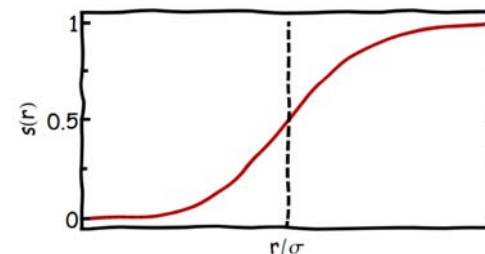
- In “metric” MDS a stress function that measures how well distances are reproduced is minimized
- Modify the objective function to aim for proximity matching
 - Distances are transformed by **sigmoid functions** in both high and low dimension

$$\chi^2 = \sum_{i,j=1}^N [||X_i - X_j| - |x_i - x_j|]^2$$

Ceriotti, Tribello, Parrinello, PNAS (2011); JCTC (2013)

- In “metric” MDS a stress function that measures how well distances are reproduced is minimized
- Modify the objective function to aim for proximity matching
 - Distances are transformed by **sigmoid functions** in both high and low dimension

$$\chi^2 = \sum_{i,j=1}^N [s(|X_i - X_j|) - s(|x_i - x_j|)]^2$$

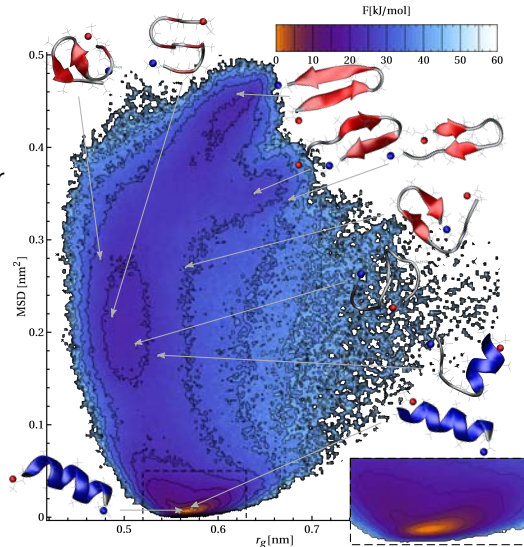


Ceriotti, Tribello, Parrinello, PNAS (2011); JCTC (2013)

The folding landscape of ala₁₂



“Conventional” CVs recognize the folded state, but many meta-stable structures overlap with each other



Ceriotti, Tribello, Parrinello, PNAS (2011); Tribello, Ceriotti, Parrinello PNAS (2012)

28

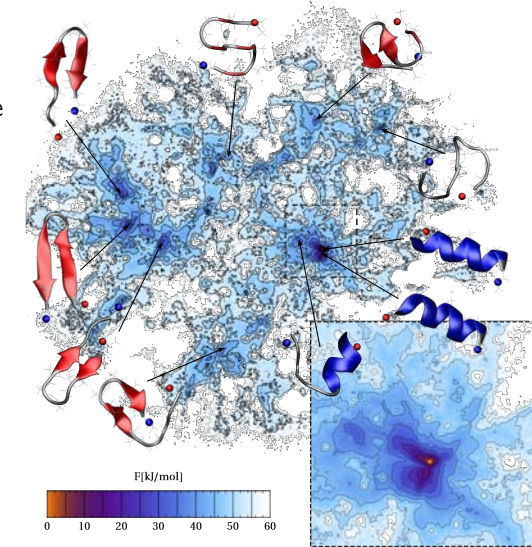
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

The folding landscape of ala₁₂



Sketch-map CVs give a very detailed picture, where each meta-stable configuration is clearly singled out¹



Ceriotti, Tribello, Parrinello, PNAS (2011); Tribello, Ceriotti, Parrinello PNAS (2012)

28

Michele Ceriotti EPFL/IMX/COSMO

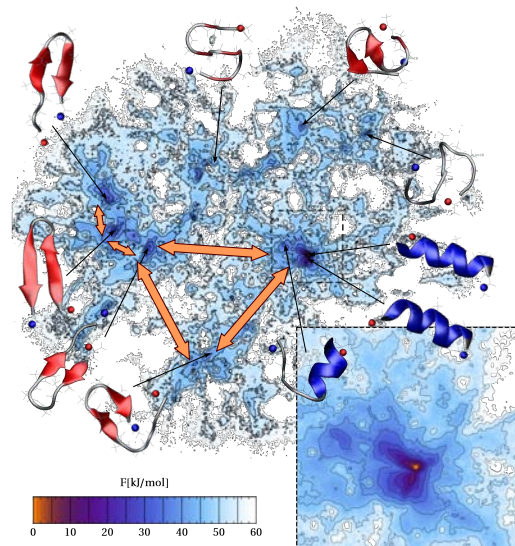
Representing and understanding patterns in materials and mol

The folding landscape of ala₁₂



Sketch-map CVs give a very detailed picture, where each meta-stable configuration is clearly singled out¹

Can be used effectively for accelerated dynamics:
field-overlap
metadynamics



Ceriotti, Tribello, Parrinello, PNAS (2011); Tribello, Ceriotti, Parrinello PNAS (2012)

28

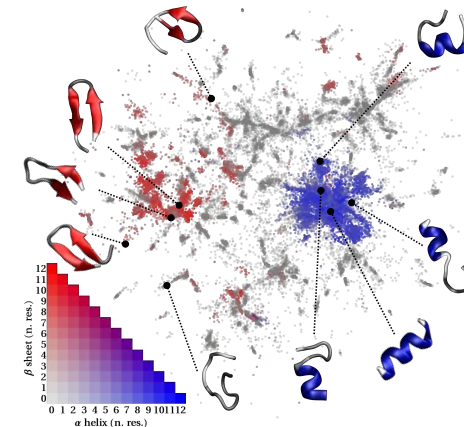
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Sketch-map and secondary structure



- Same qualitative features of the hi-D description (basins + network of transitions)
 - Sketch-map CVs correlate nicely with the **secondary structure**
 - Qualitatively similar picture if using contact-maps distance



29

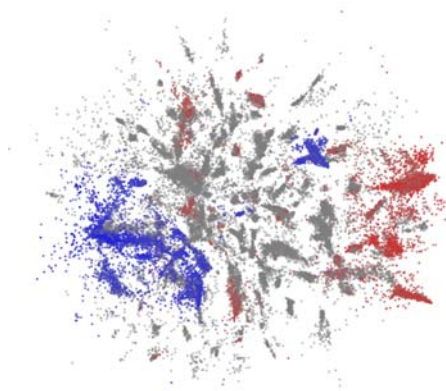
Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Sketch-map and secondary structure



- Same qualitative features of the hi-D description (basins + network of transitions)
 - Sketch-map CVs correlate nicely with the **secondary structure**
 - Qualitatively similar picture if using contact-maps distance



29

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Out-of-sample embedding



- In order to use (N)LDR as CV, one needs a way to project an arbitrary point X to low dimension.
 - PCA has a very natural linear projector solution: $x^T = \mathbf{L}X^T$
 - There are specialized solutions for different NLDR methods
- A general approach: “generalized” path coordinates

$$x(X) = \frac{\sum_i x_i e^{-|X-x_i|/\lambda}}{\sum_i e^{-|X-x_i|/\lambda}}$$

- Problem: this is a convex embedding so points away from everything will map to the center of the landmark projections

Spiwok, Králová, J. Chem. Phys. 2011

30

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

Out-of-sample embedding



- In order to use (N)LDR as CV, one needs a way to project an arbitrary point X to low dimension.
 - PCA has a very natural linear projector solution: $x^T = \mathbf{L}X^T$
 - There are specialized solutions for different NLDR methods
- A general approach: “generalized” path coordinates

$$x(X) = \frac{\sum_i x_i e^{-|X-x_i|/\lambda}}{\sum_i e^{-|X-x_i|/\lambda}}$$

- Problem: this is a convex embedding so points away from everything will map to the center of the landmark projections

Spiwok, Králová, J. Chem. Phys. 2011

30

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

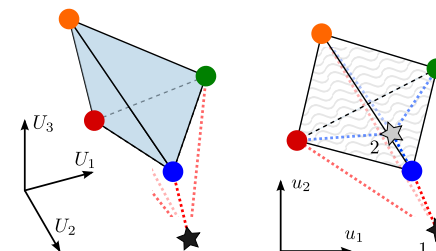
Out-of-sample embedding



- In order to use (N)LDR as CV, one needs a way to project an arbitrary point X to low dimension.
 - PCA has a very natural linear projector solution: $x^T = \mathbf{L}X^T$
 - There are specialized solutions for different NLDR methods
- A general approach: “generalized” path coordinates

$$x(X) = \frac{\sum_i x_i e^{-|X-x_i|/\lambda}}{\sum_i e^{-|X-x_i|/\lambda}}$$

- Problem: this is a convex embedding so points away from everything will map to the center of the landmark projections



Spiwok, Králová, J. Chem. Phys. 2011

30

Michele Ceriotti EPFL/IMX/COSMO

Representing and understanding patterns in materials and mol

- Sketch-map only provides projections x_i for the **landmark points** X_i
- One can work out a very natural **out-of-sample embedding** for a new point X
 - Introduce a “stress function” based on a set of landmarks and their projections

$$\chi^2(x, X) = \sum_{i=1}^N [s(|X - X_i|) - s(|x - x_i|)]^2$$

- The d -dimensional projection of the point X can be defined as the \bar{x} which minimizes $\chi^2(x, X)$

$$\bar{x}(X) = x : \min \chi^2(x, X)$$

- Sketch-map only provides projections x_i for the **landmark points** X_i
- One can work out a very natural **out-of-sample embedding** for a new point X
 - Introduce a “stress function” based on a set of landmarks and their projections

$$\chi^2(x, X) = \sum_{i=1}^N [s(|X - X_i|) - s(|x - x_i|)]^2$$

- The d -dimensional projection of the point X can be defined as the \bar{x} which minimizes $\chi^2(x, X)$

$$\bar{x}(X) = x : \min \chi^2(x, X)$$

- Sketch-map only provides projections x_i for the **landmark points** X_i
- One can work out a very natural **out-of-sample embedding** for a new point X
 - Introduce a “stress function” based on a set of landmarks and their projections

$$\chi^2(x, X) = \sum_{i=1}^N [s(|X - X_i|) - s(|x - x_i|)]^2$$

- The d -dimensional projection of the point X can be defined as the \bar{x} which minimizes $\chi^2(x, X)$

$$\bar{x}(X) = x : \min \chi^2(x, X)$$

- Sketch map can describe configurations that are in “no man’s land”, far from any landmark point!
 - We can build a useful map from rough preliminary sampling.
 - We can compare different systems using the same map.

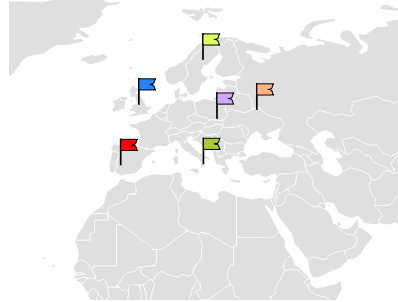
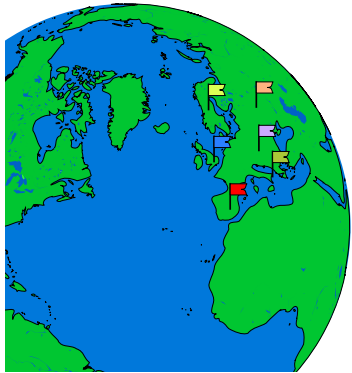


Ceriotti, Tribello, Parrinello, JCTC (2013)

Mapping in no man's land



- Sketch map can describe configurations that are in “no man's land”, far from any landmark point!
 - We can build a useful map from rough preliminary sampling.
 - We can compare different systems using the same map.

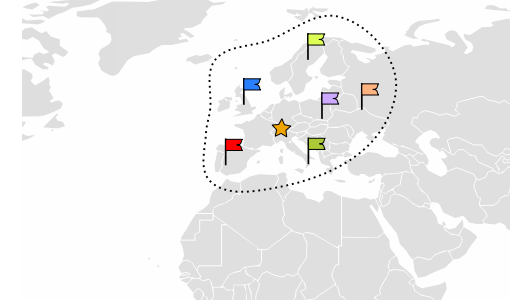
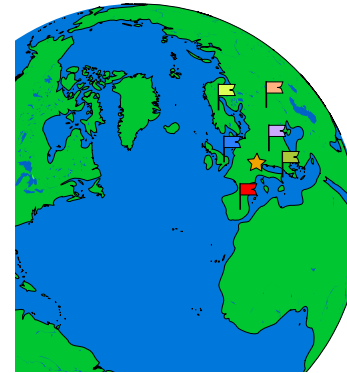


Ceriotti, Tribello, Parrinello, JCTC (2013)

Mapping in no man's land



- Sketch map can describe configurations that are in “no man's land”, far from any landmark point!
 - We can build a useful map from rough preliminary sampling.
 - We can compare different systems using the same map.

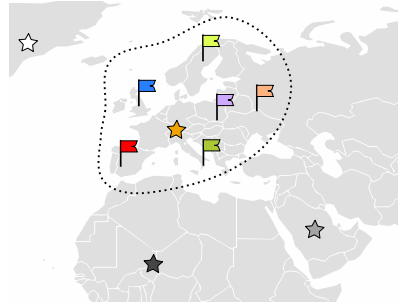
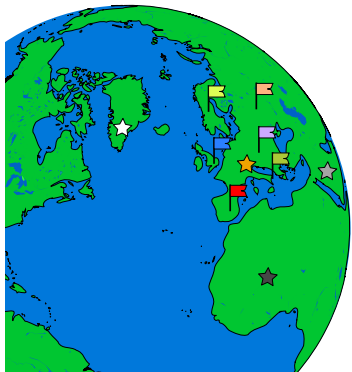


Ceriotti, Tribello, Parrinello, JCTC (2013)

Mapping in no man's land



- Sketch map can describe configurations that are in “no man's land”, far from any landmark point!
 - We can build a useful map from rough preliminary sampling.
 - We can compare different systems using the same map.

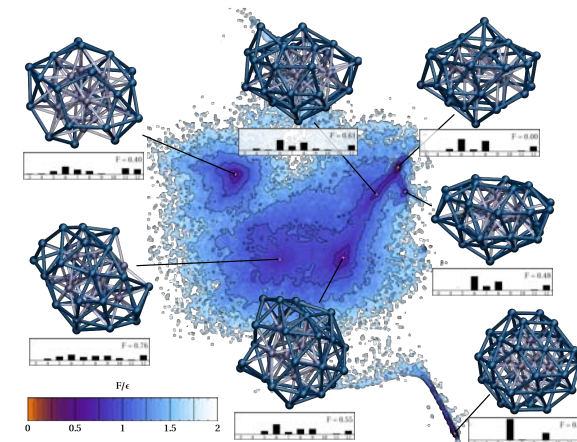


Ceriotti, Tribello, Parrinello, JCTC (2013)

From clusters to defects in the bulk

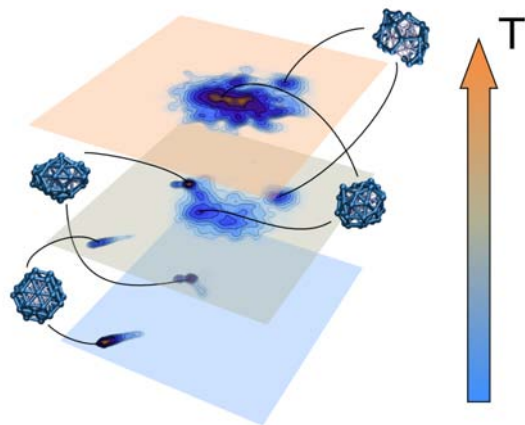


- Start building a map for a Lennard-Jones cluster
 - The same map describes the cluster across phase transitions
 - ... and can even be used to identify defects in a bulk system!



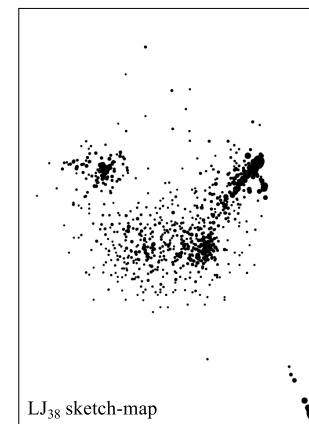
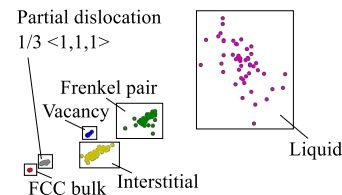
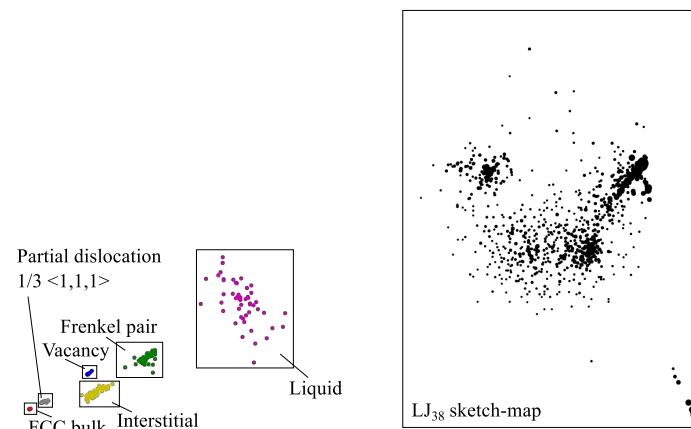
Ceriotti, Tribello, Parrinello, JCTC (2013)

- Start building a map for a Lennard-Jones cluster
- The same map describes the cluster across phase transitions
- ... and can even be used to identify defects in a bulk system!



Cerriotti, Tribello, Parrinello, JCTC (2013)

- Start building a map for a Lennard-Jones cluster
- The same map describes the cluster across phase transitions
- ... and can even be used to identify defects in a bulk system!



Cerriotti, Tribello, Parrinello, JCTC (2013)

- A problem that is common in atomistic simulations (but also data analysis in general) is how to
 - Recognize recurring patterns, appearing more often than expected
 - Perform **dimensionality reduction**, to describe a complex problem with few order parameters
- One can perform these analyses automatically
 - Mode analysis of molecular patterns by **PAMM**
 - PCA/Classical MDS are robust but **linear** techniques
 - ISOMAP, works well for “locally flat”, densely sampled data. Very sensitive to **noise**!
 - **Sketch map** targets specifically the features of atomistic simulation data.
- Another problem one should keep in mind: out-of-sample embedding. Should be *continuous* and *predictive*

Bibliography

Cox & Cox, “Multidimensional Scaling”; Tenenbaum et al., Science (2000);
 Roweis, Saul, Science (2000); Cerriotti, Tribello, Parrinello, PNAS (2011);
<http://epfl-cosmo.github.io/sketchmap>