**Tim Hempel**

**Noé Group**

**Department of Mathematics and Computer Science**

Freie Universität Berlin

# Molecular Dynamics Data Input and Featurization in PyEMMA
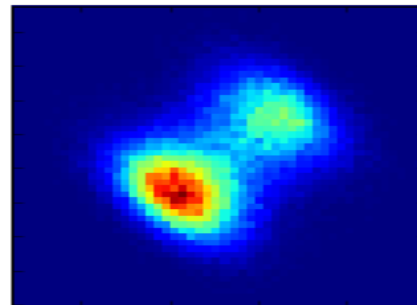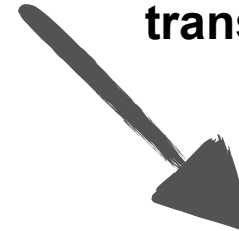
# The classical MSM analysis pipeline



"MD data"

**Featurization**
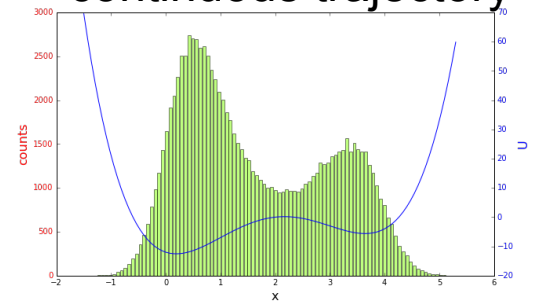"picking observables",
e.g. backbone
torsions

high dimensional
continuous trajectory

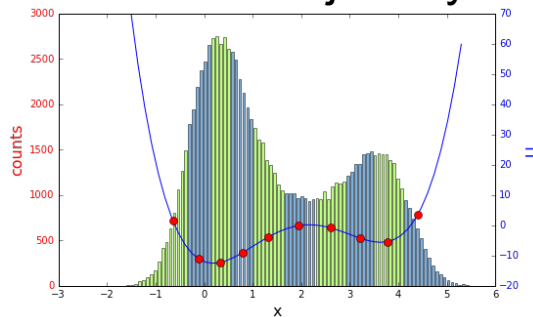**Coordinate
transform**
e.g. PCA,
TICA

shortcut

Discrete trajectory

low dimensional
continuous trajectory

**Markov
Model**

**clustering**
e.g. k-means

# The classical MSM analysis pipeline

"MD data"

**Featurization**
"picking observables",
e.g. backbone torsions

high dimensional continuous trajectory

a) "what is the best description of my system?"
b) "what do I want to model?"
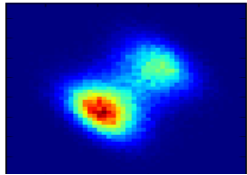
**PyEMMA natively supported features:**

- coordinates: all, heavy, Ca, selection
- angles:
  - backbone torsions
  - sidechain torsions
  - dihedrals
- distances or contacts between
  - all atom
  - Ca
  - heavy atom
- minimum distances
  - between residues or groups
- custom features



16.80

19.55.87
16.33

5.80

17.17.8
13.31

12.28

17.12
13.16

6.50

# The classical MSM analysis pipeline

high dimensional
continuous trajectory

low dimensional
continuous trajectory

**Coordinate
transform**



*"What is the minimum
dimensionality that still
represents all of the
important processes?"*
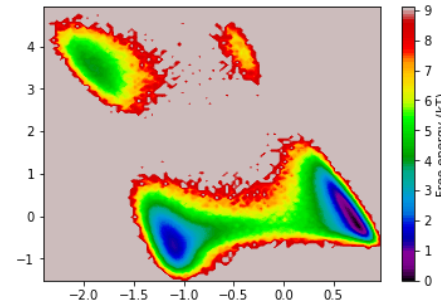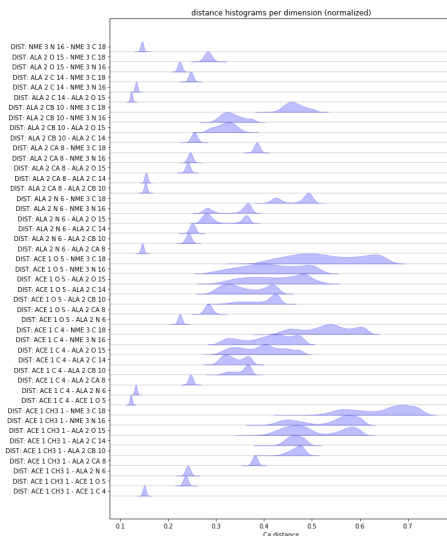
**PyEMMA natively supported coordinate transforms:**
- TICA (time-lagged independent component analysis)
  - strongly recommended
- PCA (principal component analysis)

# The classical MSM analysis pipeline



**clustering**

low dimensional continuous trajectory

Discrete trajectory

**PyEMMA natively supported clustering algorithms:**
- k-means
- regular space
- uniform time

*"What discretization resolves my processes best?"*



continuous trajectory

discrete trajectory