# Read-across the targetome
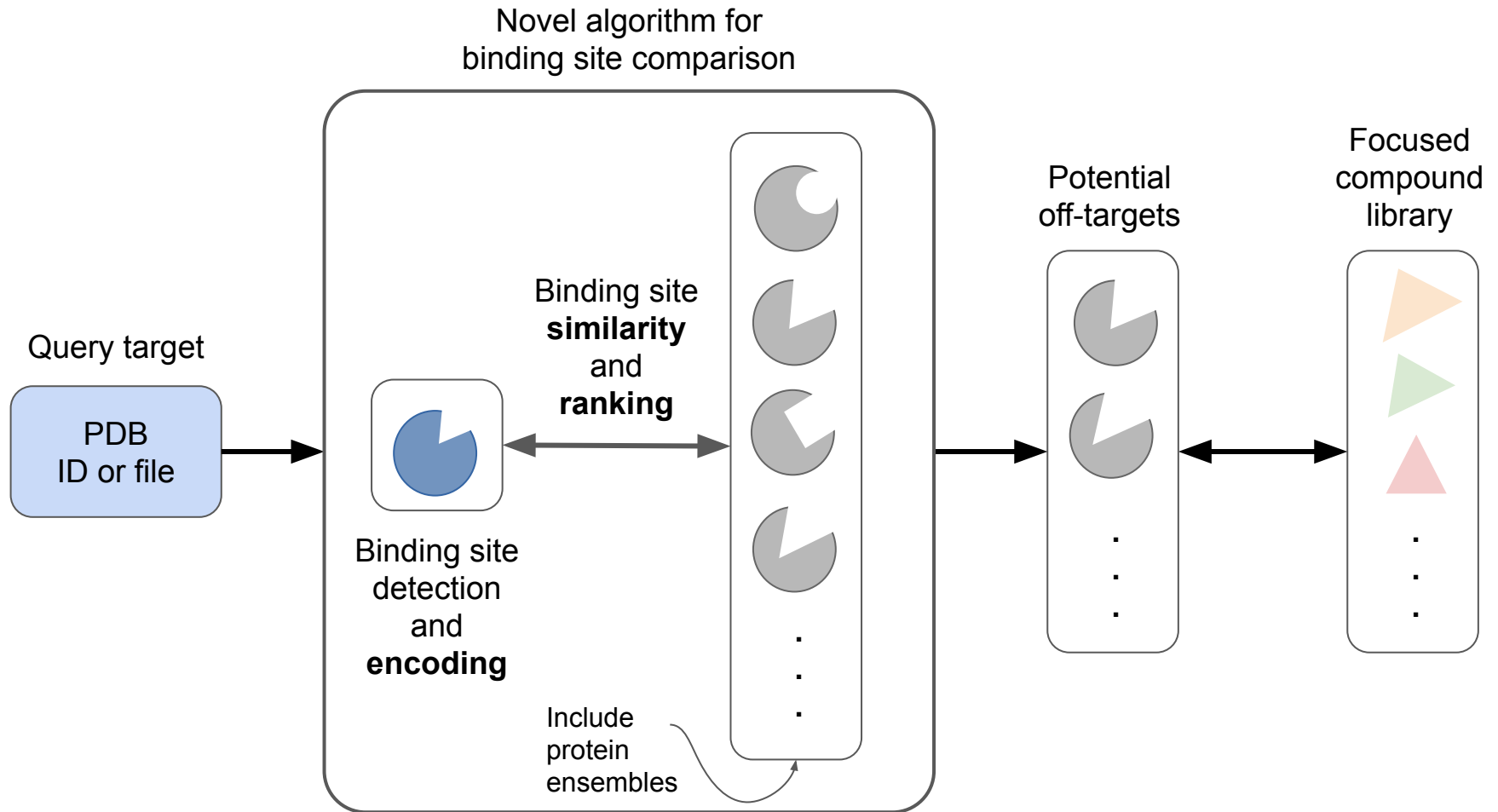
## Binding site comparison
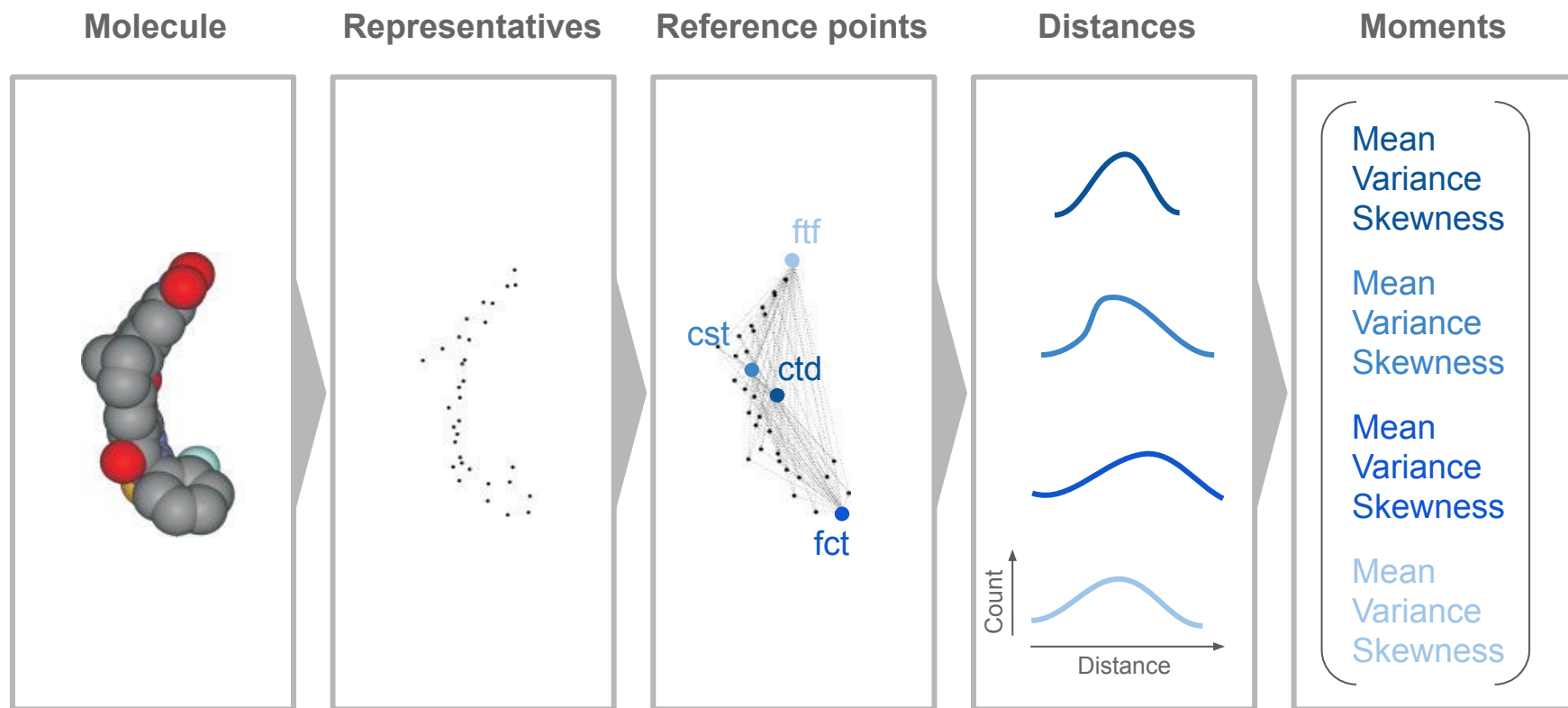
Group meeting

Dominique Sydow
30.04.2019

Novel algorithm for
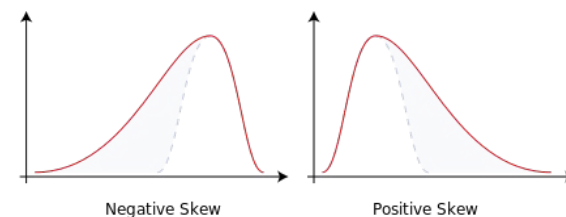binding site comparison

Query target

PDB
ID or file

Binding site
**similarity**
and
**ranking**

Binding site
detection
and
**encoding**

Include
protein
ensembles

Potential
off-targets

Focused
compound
library

# Intro: USR method for ligand encoding



Molecule → Representatives → Reference points → Distances → Moments

Reference points: ftf, cst, ctd, fct

Distances: Count vs Distance

Moments: Mean, Variance, Skewness (×4)

Ultrafast shape recognition (USR) - Ballester et al. 2006

# Intro: USR method for ligand encoding



| Molecule | Representatives | Reference points | Distances | **Moments** |

Reference points: ftf, cst, ctd, fct

Distances: Count / Distance

Moments:
Mean Variance Skewness
Mean Variance Skewness
Mean Variance Skewness
Mean Variance Skewness

Negative Skew          Positive Skew

# Idea: USR method for binding site encoding?



| Molecule | Representatives | Reference points | Distances | Moments |
|----------|-----------------|------------------|-----------|---------|

# Methods: Binding site dimensionality

**Molecule**   **Representatives**   Reference points   Distances   Moments

Cα atoms
Pseudocenters

HO ── C ── OH
      │
      NH₂

Serine

Count
Distance

Mean
Variance
Skewness

Mean
Variance
Skewness

Mean
Variance
Skewness

Mean
Variance
Skewness

How many dimensions
per data point?

**Molecule**

**Representatives**

Cα atoms
Pseudocenters

HO ——— C ——— OH
         |
        NH₂

Serine

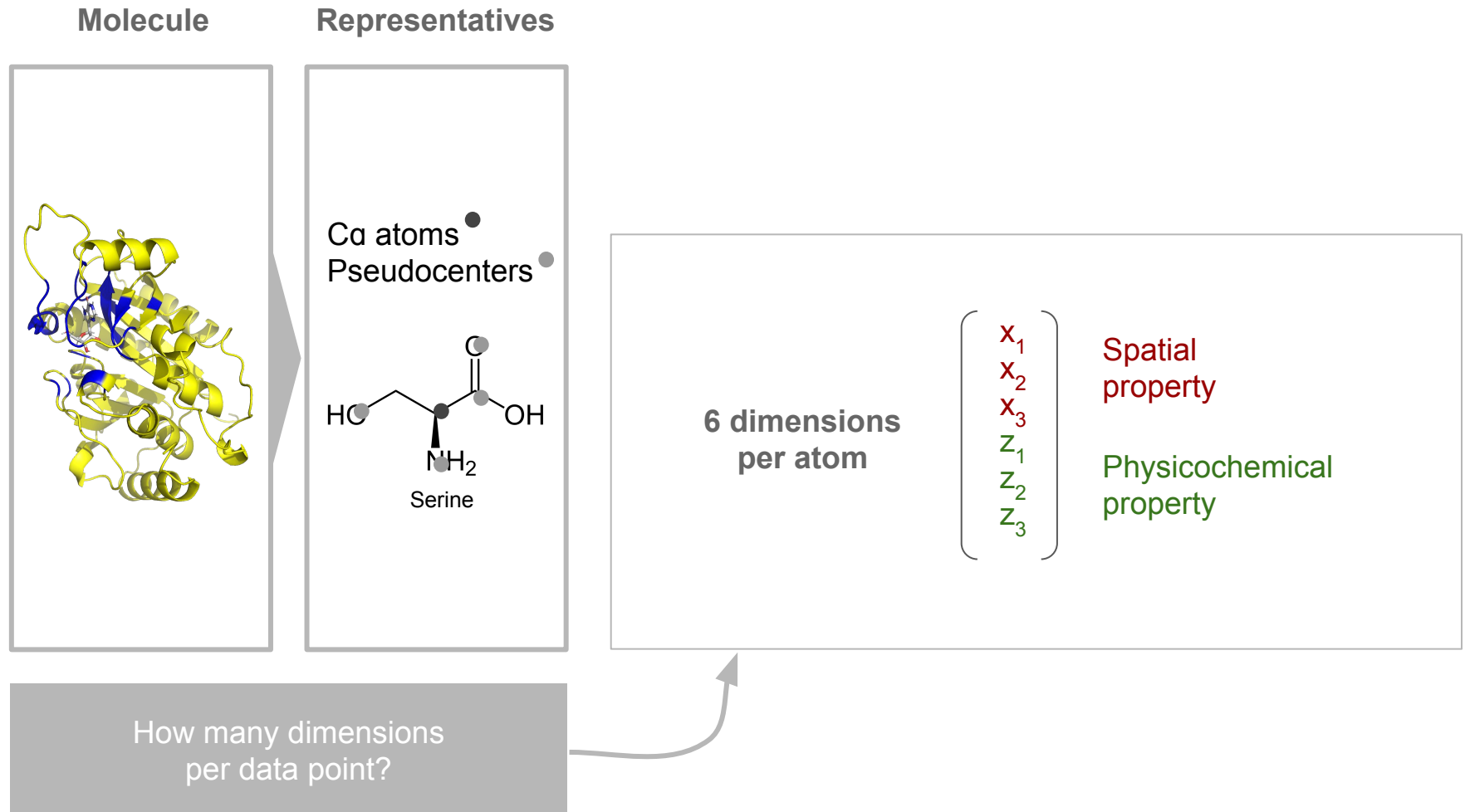**Amino acid descriptors**

molecular weight (g/mol)
TLC % migration on silica gel, ethanol/water (70/30)[a]
TLC, silica gel, 1-butanol/acetic acid/water (40/10/10)
TLC, silica gel, phenol/water (75/25)
TLC, silica gel, butanone/pyridine/acetic acid/water (70/15/2/15)
TLC, cellulose, ethanol/water (70/30)
TLC, cellulose, pyridine/isoamyl alcohol/water (35/30/30)
TLC, kiselguhr, butanone/water/phenol/acetone/ethanol (1/1)
side chain van der Waals volume (cm$^3$/mol)
NMR α-proton shift at pD = 2 (ppm)
NMR α-proton shift at pD = 7 (ppm)
NMR α-proton shift at pD = 12.5 (ppm)
$^{10}$log(octanol/water) partition coefficient
energy of highest occupied molecular orbital (eV)
energy of lowest unoccupied molecular orbital (eV)
heat of formation (kcal)
α-polarizability (Å$^3$)
absolute electronegativity (eV)
absolute hardness (eV)
total accessible molecular surface area (log Å$^2$)
polar accessible molecular surface area (log Å$^2$)
nonpolar accessible molecular surface area (log Å$^2$)
number of hydrogen bond donors
number of hydrogen bond acceptors
indicator of positive charge in side chain
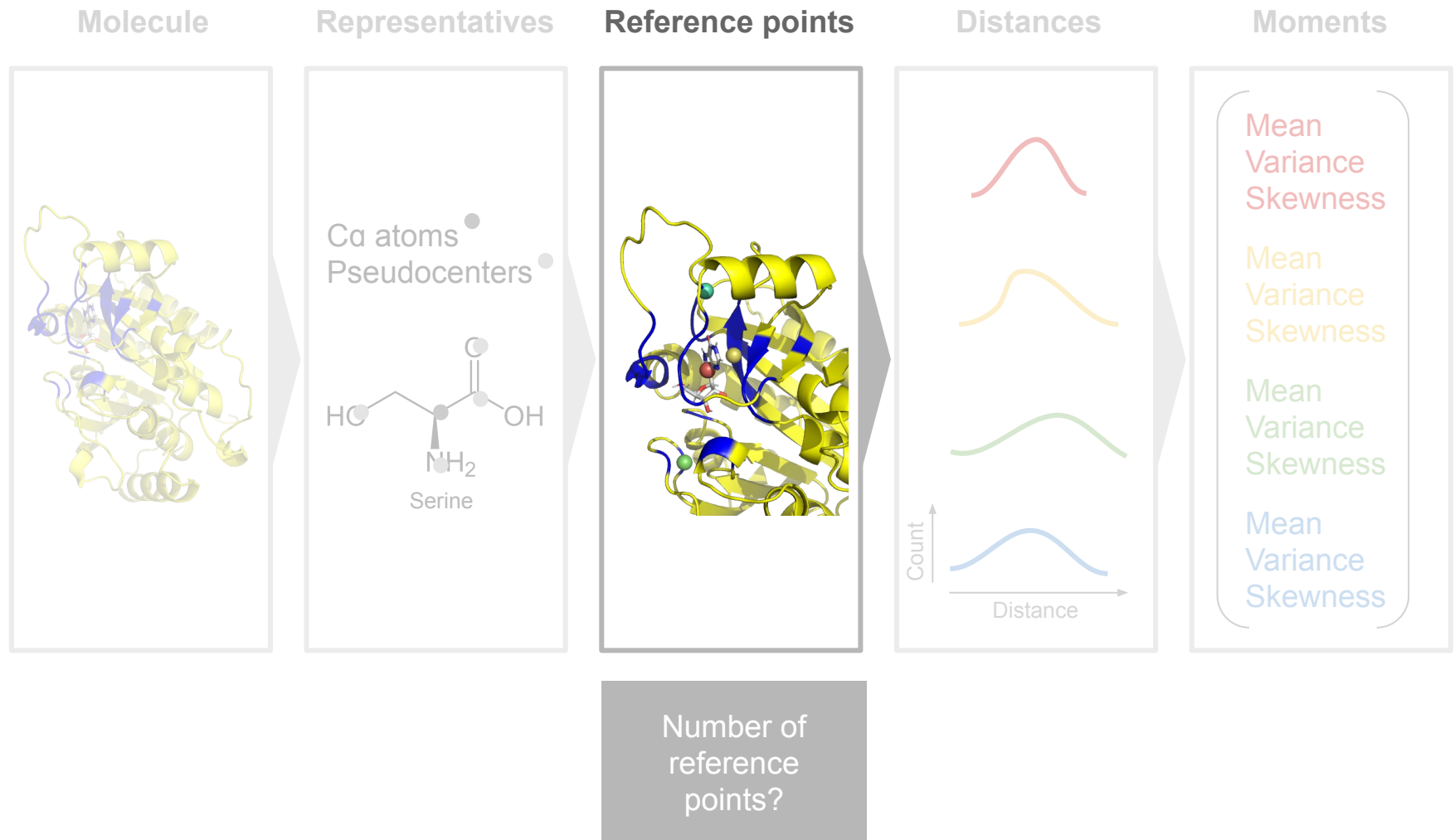indicator of negative charge in side chain

**Z-scales**

|  | $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ |
|---|---|---|---|---|---|
| alanine | 0.24 | −2.32 | 0.60 | −0.14 | 1.30 |
| arginine | 3.52 | 2.50 | −3.50 | 1.99 | −0.17 |
| asparagine | 3.05 | 1.62 | 1.04 | −1.15 | 1.61 |
| aspartic acid | 3.98 | 0.93 | 1.93 | −2.46 | 0.75 |
| cysteine | 0.84 | −1.67 | 3.71 | 0.18 | −2.65 |
| glutamine | 1.75 | 0.50 | −1.44 | −1.34 | 0.66 |
| glutamic acid | 3.11 | 0.26 | −0.11 | −3.04 | −0.25 |
| glycine | 2.05 | −4.06 | 0.36 | −0.82 | −0.38 |
| histidine | 2.47 | 1.95 | 0.26 | 3.90 | 0.09 |
| isoleucine | −3.89 | −1.73 | −1.71 | −0.84 | 0.26 |
| leucine | −4.28 | −1.30 | −1.49 | −0.72 | 0.84 |
| lysine | 2.29 | 0.89 | −2.49 | 1.49 | 0.31 |
| methionine | −2.85 | −0.22 | 0.47 | 1.94 | −0.98 |
| phenylalanine | −4.22 | 1.94 | 1.06 | 0.54 | −0.62 |
| proline | −1.66 | 0.27 | 1.84 | 0.70 | 2.00 |
| serine | 2.39 | −1.07 | 1.15 | −1.39 | 0.67 |
| threonine | 0.75 | −2.18 | −1.12 | −1.46 | −0.40 |
| tryptophan | −4.36 | 3.94 | 0.59 | 3.44 | −1.59 |
| tyrosine | −2.54 | 2.44 | 0.43 | 0.04 | −1.47 |
| valine | −2.59 | −2.64 | −1.54 | −0.85 | −0.02 |

$z_1$ Lipophilicity
$z_2$ Steric bulk/
    polarisability
$z_3$ Polarity

How many dimensions
per data point?

# Methods: Binding site dimensionality

**Molecule**

**Representatives**

Cα atoms
Pseudocenters

HC — C(=O) — OH
       |
      NH$_2$

Serine

**6 dimensions per atom**

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ z_1 \\ z_2 \\ z_3 \end{pmatrix}$$

Spatial property

Physicochemical property

How many dimensions per data point?

# Methods: Number of reference points?

**Molecule**　　**Representatives**　　**Reference points**　　**Distances**　　**Moments**



Cα atoms
Pseudocenters

HO　　OH

NH₂

Serine

Count

Distance

Mean
Variance
Skewness

Mean
Variance
Skewness

Mean
Variance
Skewness

Mean
Variance
Skewness

**Number of reference points?**

**Reference points**



Number of reference points?

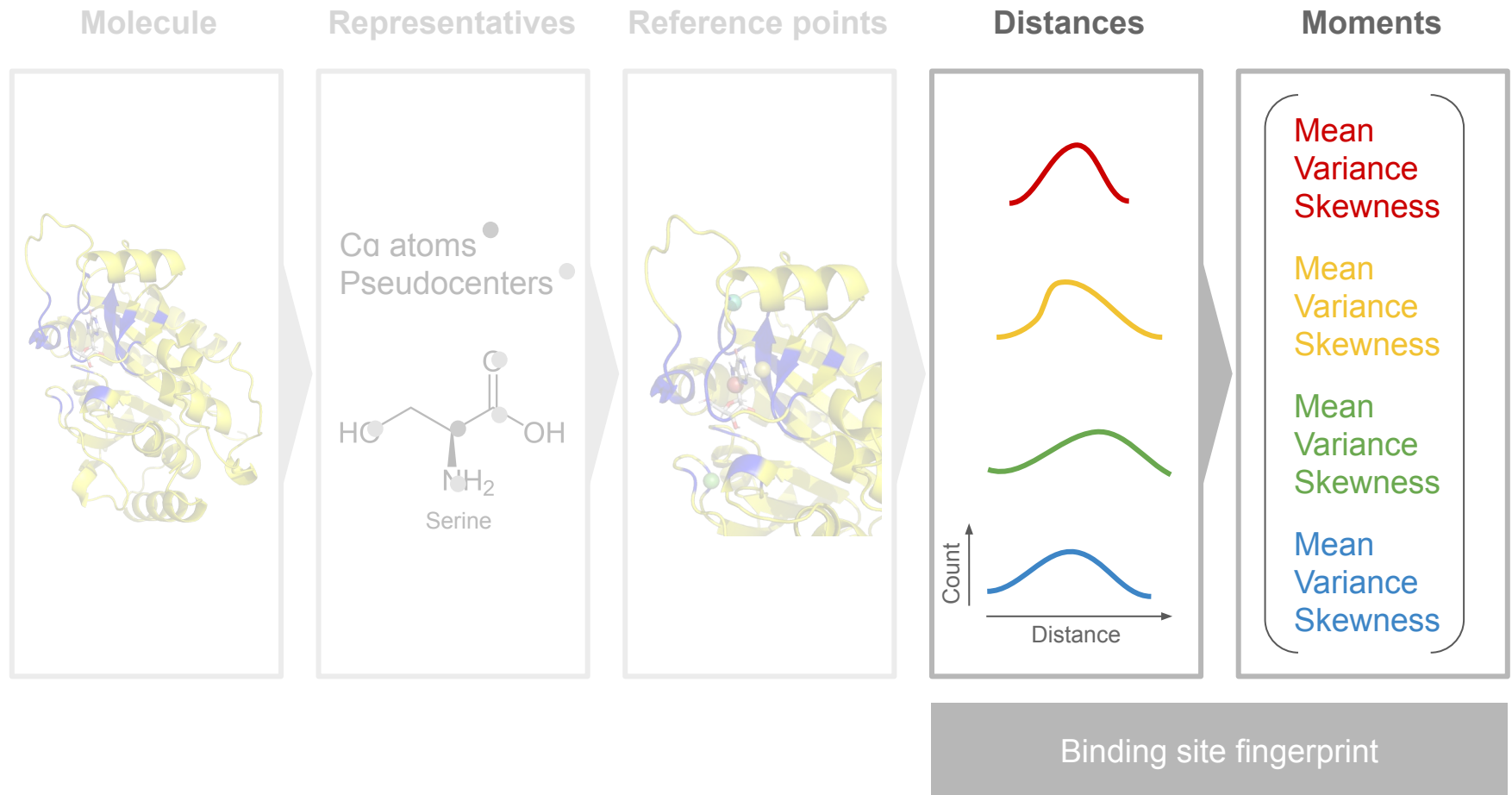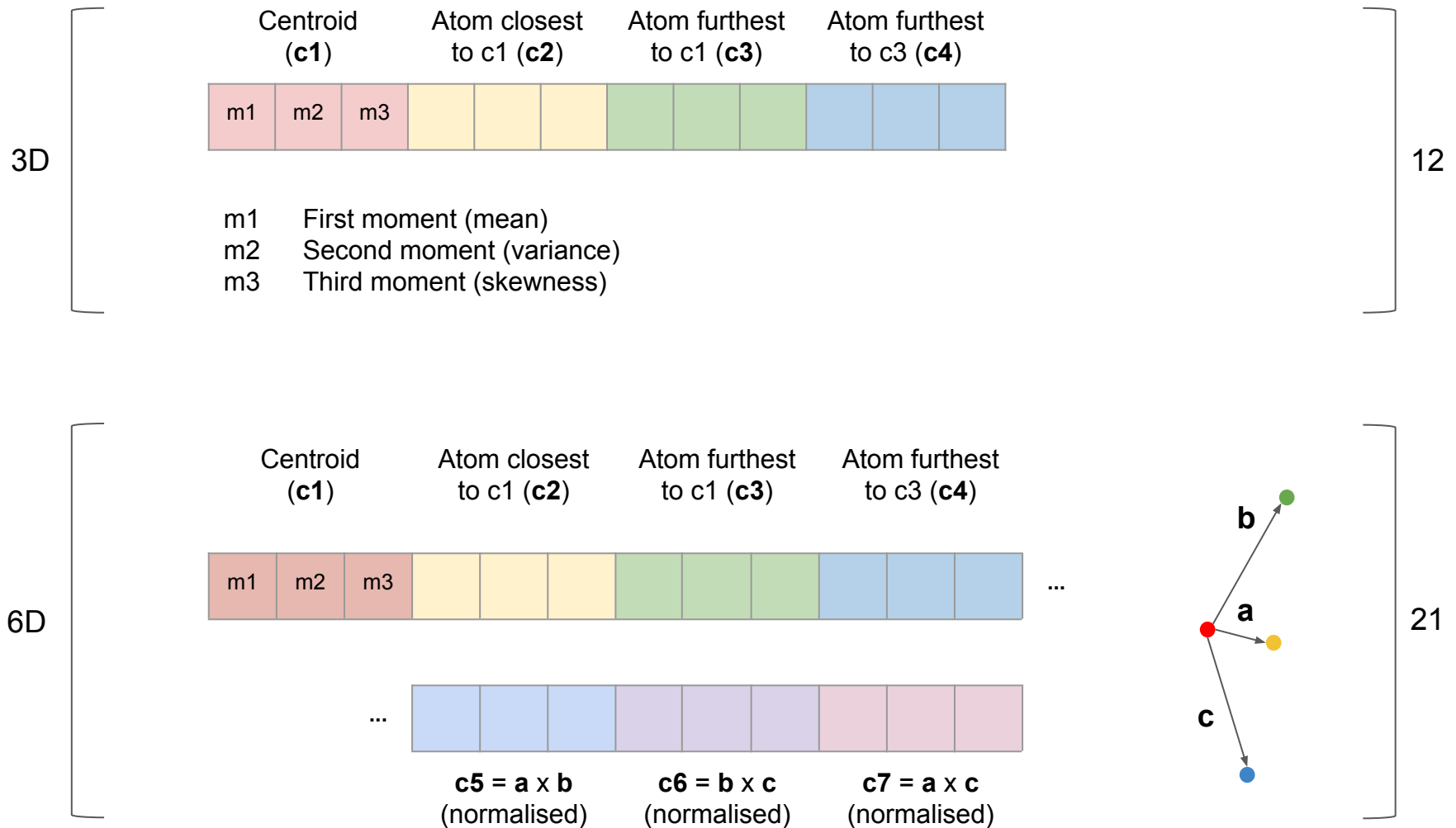| | *n* dimensions | *n* reference points | *n*+1 reference points |
|---|---|---|---|

To find the exact position of a **point** in $\mathbf{R}^n$,
distances to *n*+1 **fixed (reference) points** are needed.
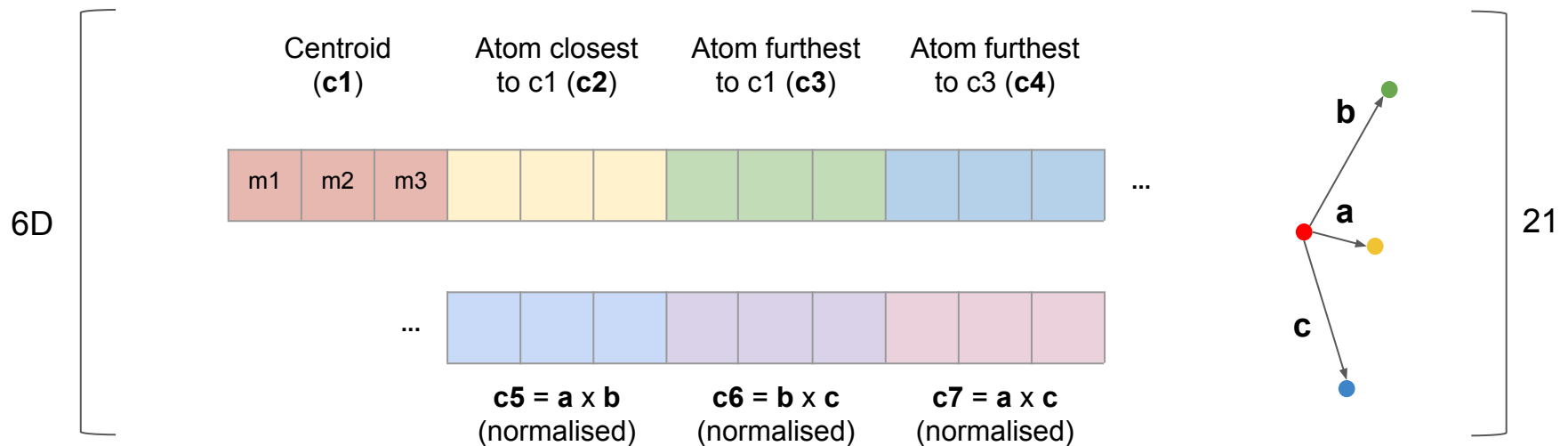
# Methods: Binding site fingerprints

**Molecule**

**Representatives**

**Reference points**

**Distances**

**Moments**

Cα atoms
Pseudocenters

HO

OH

NH$_2$

Serine

Count

Distance

Mean
Variance
Skewness

Mean
Variance
Skewness

Mean
Variance
Skewness

Mean
Variance
Skewness

Binding site fingerprint

3D

| Centroid (**c1**) | | | Atom closest to c1 (**c2**) | | | Atom furthest to c1 (**c3**) | | | Atom furthest to c3 (**c4**) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| m1 | m2 | m3 | | | | | | | | | |

m1    First moment (mean)
m2    Second moment (variance)
m3    Third moment (skewness)

12

6D

| Centroid (**c1**) | | | Atom closest to c1 (**c2**) | | | Atom furthest to c1 (**c3**) | | | Atom furthest to c3 (**c4**) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| m1 | m2 | m3 | | | | | | | | | | ... |

... | | | | | | | | |

**c5** = a x b (normalised)    **c6** = b x c (normalised)    **c7** = a x c (normalised)

21

Outlook: Calculate cross product in 6 dimensions?

6D

| Centroid (**c1**) | | | Atom closest to c1 (**c2**) | | | Atom furthest to c1 (**c3**) | | | Atom furthest to c3 (**c4**) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| m1 | m2 | m3 | | | | | | | | | | ... |

21

... | **c5** = **a** x **b** (normalised) | | | **c6** = **b** x **c** (normalised) | | | **c7** = **a** x **c** (normalised) |

**b**

**a**

**c**

# Methods: Binding site similarity measure

| Centroid (**c1**) | | | Atom closest to c1 (**c2**) | | | Atom furthest to c1 (**c3**) | | | Atom furthest to c3 (**c4**) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| m1 | m2 | m3 | | | | | | | | | |

$$\vec{M}^q = (4.44,\ 2.98,\ 1.04,\ 4.55,\ 4.70,\ 0.23,\ 8.30,\ 16.69,\ -22.97,\ 7.37,\ 15.64,\ 0.51)$$

$$S_{qi} = \frac{1}{1 + \frac{1}{12}\sum_{i=1}^{12}\left|M_i^q - M_i^i\right|} \in (0,1] \qquad S_{qi} = 0.812$$

$$\vec{M}^i = (4.39,\ 3.11,\ 1.36,\ 4.50,\ 4.44,\ 0.09,\ 8.34,\ 16.78,\ -23.20,\ 7.15,\ 16.52,\ 0.13)$$

**Similarity measure**

Inverse of the translated and scaled
**Manhattan distance**

# Methods: Binding site similarity measure

| Centroid (**c1**) | | | Atom closest to c1 (**c2**) | | | Atom furthest to c1 (**c3**) | | | Atom furthest to c3 (**c4**) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| m1 | m2 | m3 | | | | | | | | | |

$$\vec{M}^q = (4.44,\ 2.98,\ 1.04,\ 4.55,\ 4.70,\ 0.23,\ 8.30,\ 16.69,\ -22.97,\ 7.37,\ 15.64,\ 0.51)$$

$$S_{qi} = \frac{1}{1 + \frac{1}{12} \sum_{i=1}^{12} \left| M_i^q - M_i^i \right|} \in (0,1] \qquad S_{qi} = 0.812$$

$$\vec{M}^i = (4.39,\ 3.11,\ 1.36,\ 4.50,\ 4.44,\ 0.09,\ 8.34,\ 16.78,\ -23.20,\ 7.15,\ 16.52,\ 0.13)$$

**Similarity measure**
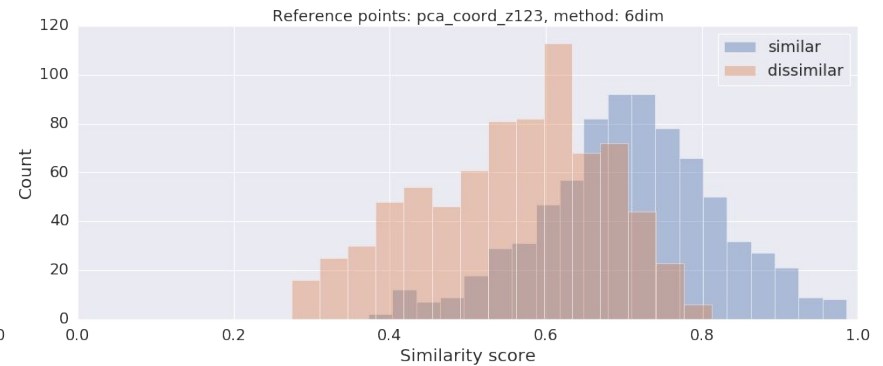
Inverse of the translated and scaled
**Manhattan distance**

# Evaluation: Similar vs. dissimilar pairs (FuzCav)



769 pairs of non-redundant **similar binding sites**

vs.

769 pairs of non-redundant **dissimilar binding sites**

Data set generation

1. Cluster scPDB by UniProt name (911 clusters and 1204 singletons)
2. All-against-all comparison of all active sites within each cluster (SiteAlign)
3. Define *cutoff* for similarity measure discriminating between similar/dissimilar binding sites
4. Choose pairs
   - Similar pairs: select randomly two entries per cluster (considering *cutoff*)
   - Dissimilar pairs: select two entries from clusters differering at the first level of their EC numbers

**3D**

**Ca atoms**

# Evaluation: Similar vs. dissimilar pairs (FuzCav)
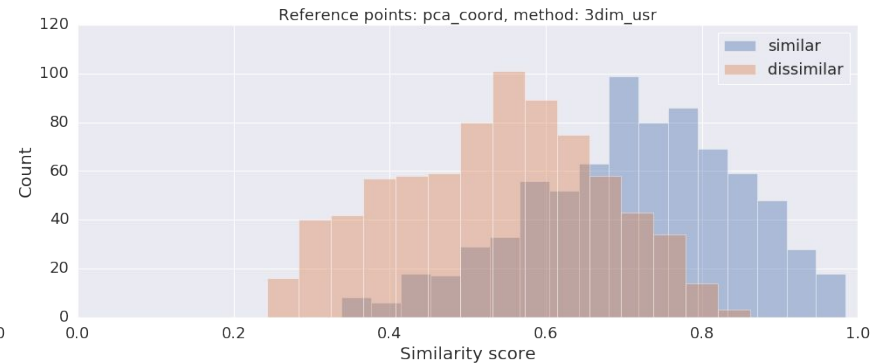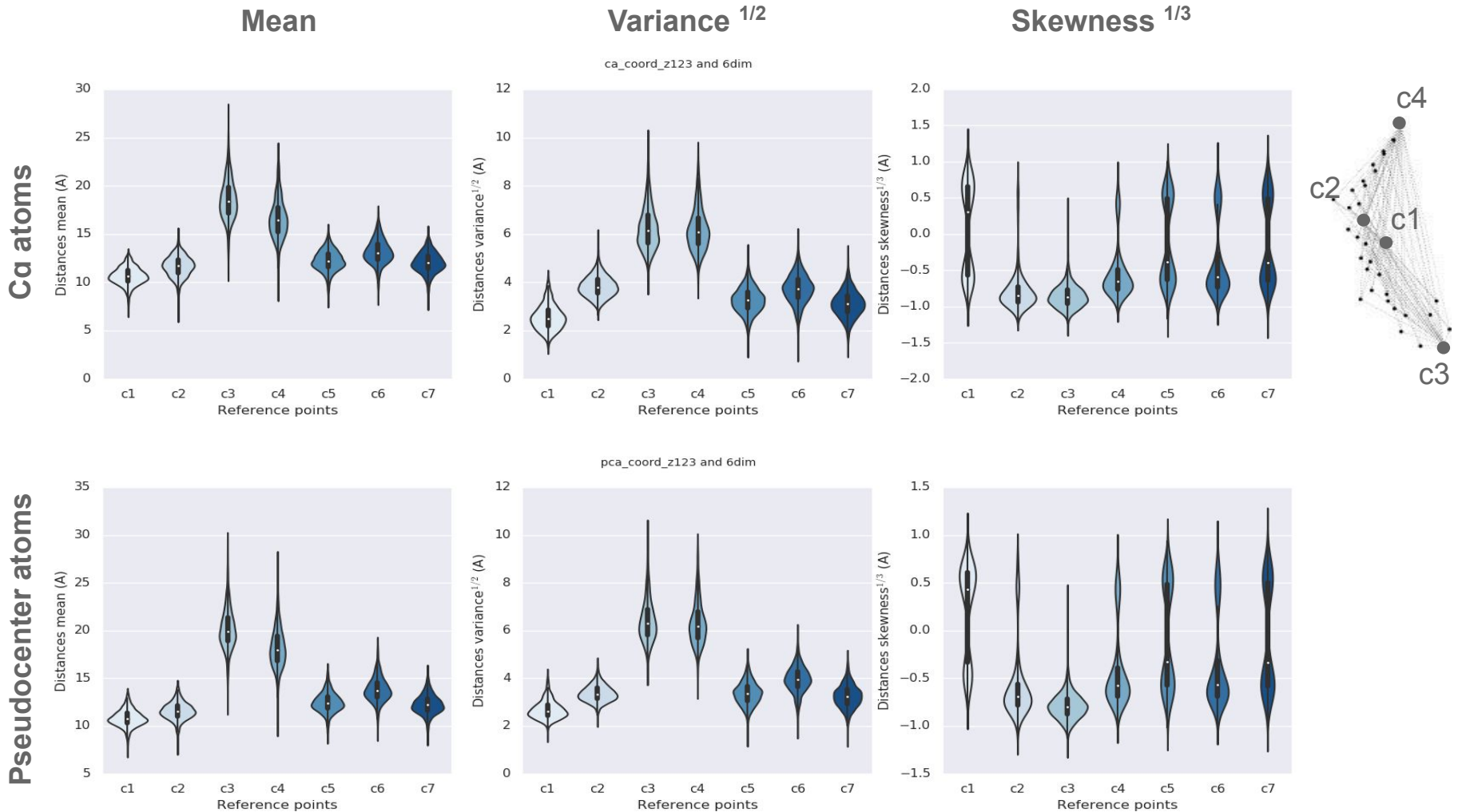
**Ca atoms**

**Pseudocenter atoms**

ROC curves

Legend:
- 0.29: ca_3dim_usr
- 0.58: ca_3dim_csr
- 0.3: ca_z1_4_dim_electroshape
- 0.32: ca_z123_6dim
- 0.29: pca_3dim_usr
- 0.58: pca_3dim_csr
- 0.3: pca_z1_4_dim_electroshape
- 0.32: pca_z123_6dim
- 0.29: pc_3dim_usr
- 0.58: pc_3dim_csr
- 0.31: pc_z1_4_dim_electroshape
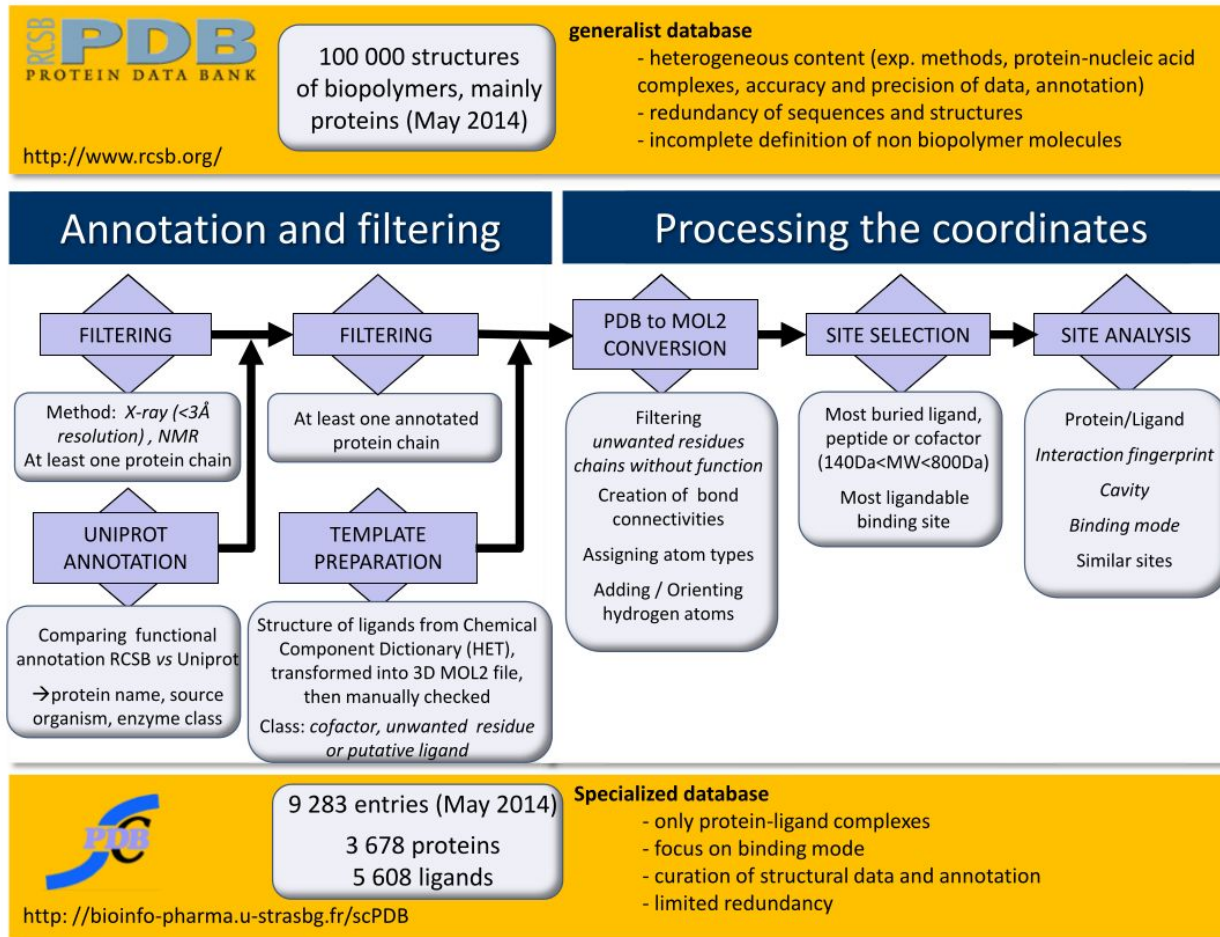- 0.33: pc_z123_6dim

# Outlook

- **Set reference points based on 6 dimensions**
  - Cross product in 6 dimensions
- **Instead of Z-scales (per residue) use**
  - AutoDock partial charges (per atom) and/or
  - Flexibility information via normal mode analysis (per atom)
- **Introduce subpockets/regions**
  - Calculate fingerprints for overlapping regions
  - All-against-all fingerprint comparison between binding site regions
  - Find maximal neighboring matches
- **Use more information from distance histograms than moments for fingerprint**
- **Apply method to benchmarking datasets other than FuzCav dataset**
  - TOUGH-M1 dataset (Govindaraj et al. 2018)
  - ProSPECCTs datasets (Ehrt et al. 2018)
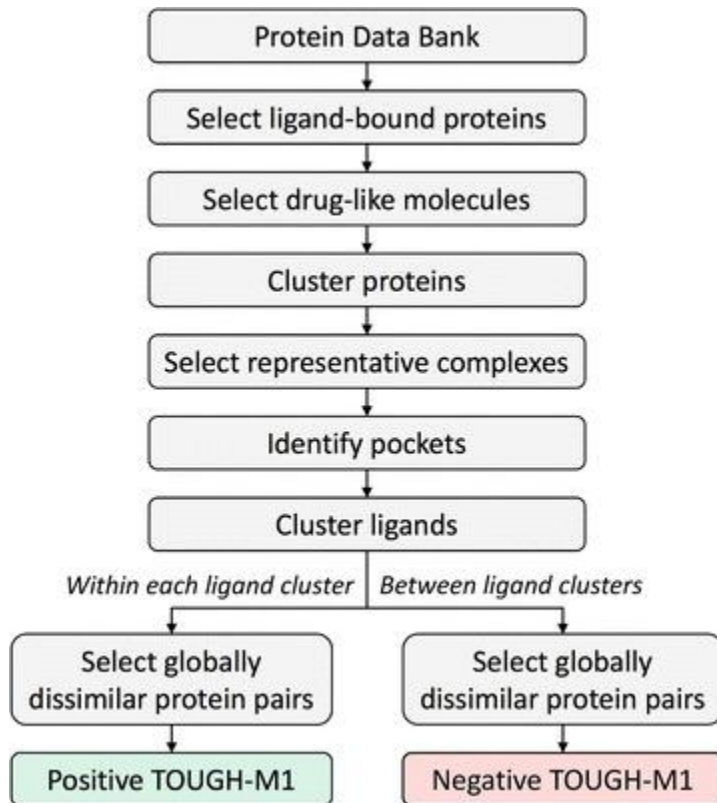


| c1 | c2 | c3 | c4 | c5 | c6 | c7 |

# Code review sessions?

Questions I have...

- How to initialize global variables at start of full program and how to pass them to all associated scripts?
- How to store complex data structures?
  - Dict of dict of Pandas DataFrames
  - Database?
- How to note functions that are only called within other functions (but will not be called by themselves)?
- What is the advantage of Docker over conda - and when is what good to use?
- ...

The end.

# scPDB

CHARITÉ UNIVERSITÄTSMEDIZIN BERLIN

# Benchmarking dataset: TOUGH-M1

# Benchmarking datasets: ProSPECCTs

| goal | number of comparisons (similar or active / dissimilar or inactive pairs) | resolution (mean ± stddev, minimum, maximum) [Å] | R_work (mean ± stddev, minimum, maximum) | average overall G-factor (mean ± stddev, minimum, maximum) |
|---|---|---|---|---|
| structures with identical sequences (data set 1) | | | | |
| sensitivity with respect to the binding site definition, score range for active and inactive pairs | 13,430 / 92,846 (12 groups of structures with identical sequences) | 1.79 ± 0.37, 0.8, 2.71 | 0.174 ± 0.027, 0.091, 0.264 | 0.023 ± 0.23, -1.27, 0.6 |
| structures with identical sequences and similar ligands (data set 1.2) | | | | |
| impact of ligand diversity on binding site comparison | 241 / 1,784 | 1.73 ± 0.37, 0.92, 2.5 | 0.171 ± 0.025, 0.104, 0.232 | 0.019 ± 0.22, -0.57, 0.6 |
| NMR structures (data set 2) | | | | |
| sensitivity with respect to the binding site flexibility | 7,729 / 100,512 (17 structural ensembles of diverse proteins) | n.d. | n.d. | -0.279 ± 0.705, -2.8, 0.21 |
| decoy set 1 (data set 3) | | | | |
| differentiation between binding sites with different physicochemical properties | 13,430 / 67,150 (complete data set) 13,430 / 13,430 (data set with five residue variants) | n.d. | n.d. | n.d. |
| decoy set 2 (data set 4) | | | | |
| differentiation between binding sites with different physicochemical and shape properties | 13,430 / 67,150 (complete data set) 13,430 / 13,430 (data set with five residue variants) | n.d. | n.d. | n.d. |
| Kahraman data set[63] without phosphate binding sites (data set 5) | | | | |
| classification of proteins binding to identical ligands and cofactors | 920 / 5,480 | 2.02 ± 0.37, 0.88, 2.9 | 0.202 ± 0.033, 0.089, 0.265 | 0.166 ± 0.228, -0.56, 0.47 |
| Kahraman data set[63] (data set 5.2) | | | | |
| original data set | 1,320 / 8,680 | 2.02 ± 0.4, 0.88, 2.9 | 0.201 ± 0.031, 0.089, 0.265 | 0.162 ± 0.218, -0.56, 0.47 |
| Barelier data set[64] (data set 6) including cofactors (data set 6.2) | | | | |
| identification of distant relationships between protein binding sites with identical ligands which "observe" a similar environment | 19 / 43 | 2.16 ± 0.44, 0.93, 3.1 | 0.196 ± 0.027, 0.104, 0.25 | 0.117 ± 0.23, -1.46, 0.53 |
| data set of successful applications (data set 7) | | | | |
| recovery of known binding site similarities within a set of diverse proteins | 115 / 56,284 (49 query structures) | 1.98 ± 0.43, 0.8, 3.25 | 0.191 ± 0.029, 0.101, 0.284 | 0.13 ± 0.208, -2.8, 1.35 |

CHARITÉ UNIVERSITÄTSMEDIZIN BERLIN

Dominique Sydow

27

# Binding site encoding: representatives

## Point number

- Cα atoms
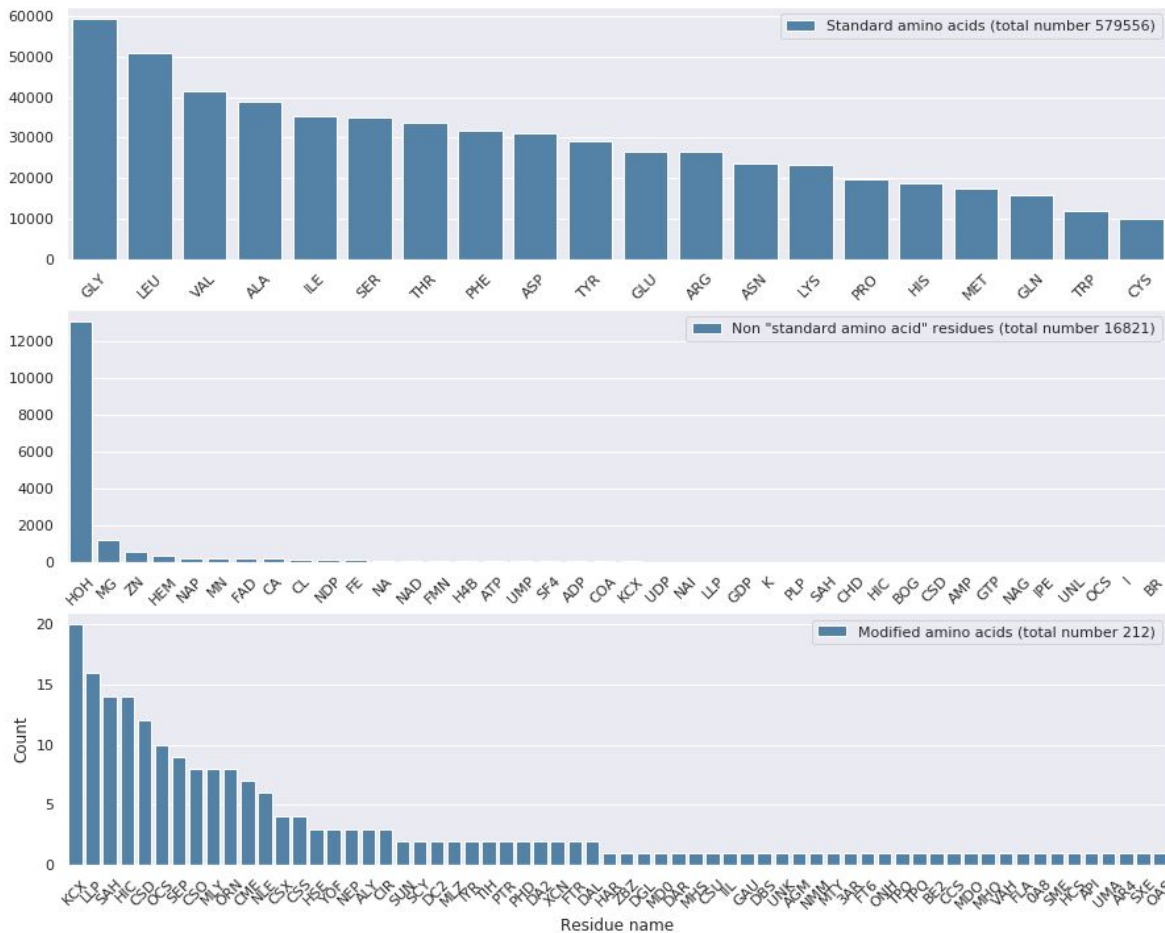- Pseudocenters

## Point dimensions $n$

- Spatial information
  - $x_1$, $x_2$, and $x_3$
- Physicochemical information
  - Z-scales $z_1$, $z_2$, and $z_3$ (lipophilicity, steric bulk/polarisability, and polarity)
  - Physicochemical atom subsets (based on pseudocenters: aliphatic, donor, acceptor, aromatic, or donor/acceptor)

Pseudocenters



| Side-chain | Amino acid | Pseudocenter (type) | Origin atoms |
|---|---|---|---|
| H₃C | Ala | Aliphatic | CB |
| | Arg | Aliphatic | CB, CG, CD |
| | | Donor | NE |
| | | Donor | NH1 |
| | | Donor | NH2 |
| | Asn | Acceptor | OD1 |
| | | Donor | ND2 |
| | Asp | Acceptor | OD1 |
| | | Acceptor | OD2 |
| (R,H)S | Cys | Aliphatic | CB, SG |
| | Gln | Acceptor | OE1 |
| | | Donor | NE2 |
| | Glu | Acceptor | OE1 |
| | | Acceptor | OE2 |
| | His | PI | CG, ND1, CD2, CE1, NE2 |
| | | DON_ACC | NE1 |
| | | DON_ACC | NE2 |

Schmitt et al. 2002 (Cavbase)

# Data set: Residue composition (scPDB)

Overall - how often are standard amino acids and other residues?



| Mod. aa in scPDB & z-scales | # in scPDB |
|---|---|
| NLE | 6 |
| ISE | 1 |
| ORN | 8 |

CHARITÉ UNIVERSITÄTSMEDIZIN BERLIN

# Data set: Residue composition (scPDB)

Per binding site - are standard amino acids somewhat equally distributed?



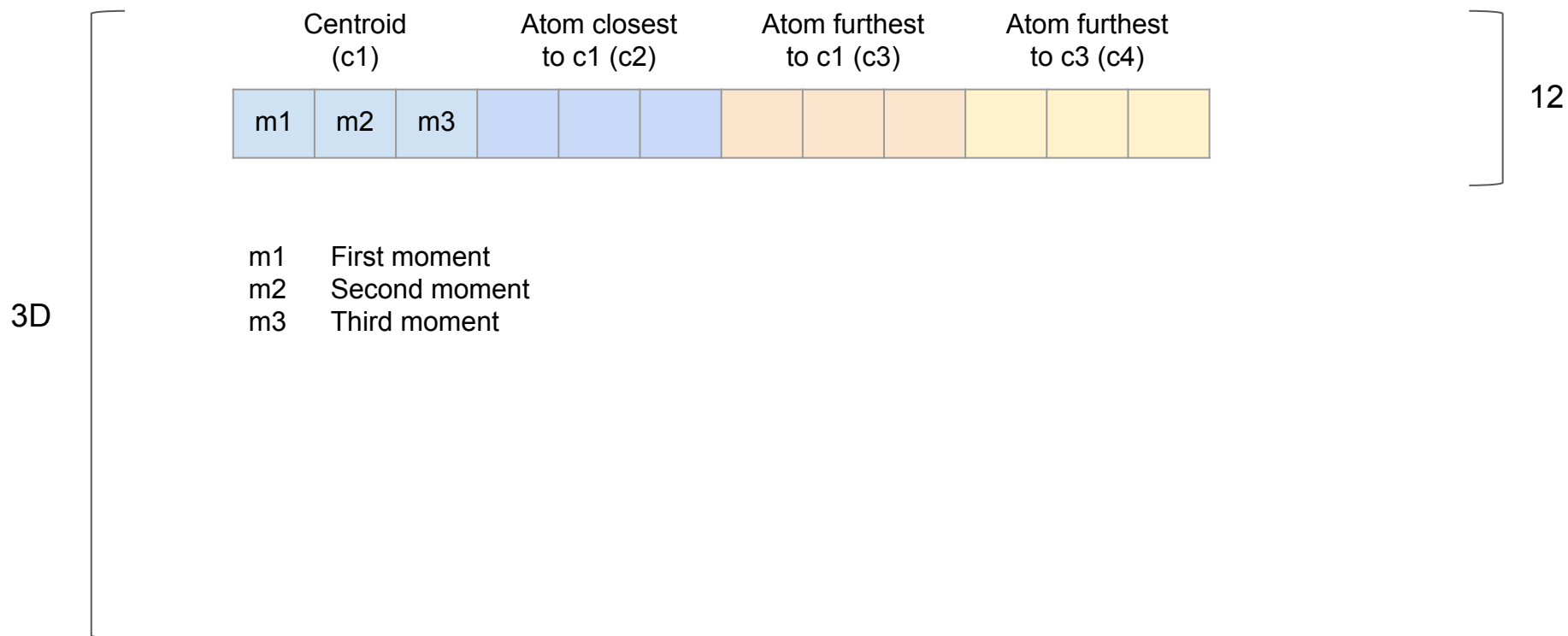Distribution of amino acid ratios in binding sites (scPDB)

# Moments

| Moment number | Name | Measure of | Formula |
|---|---|---|---|
| 1 | Mean | Central tendency | $\bar{X} = \dfrac{\sum_{i=1}^{N} X_i}{N}$ |
| 2 | Variance (Volatility) | Dispersion | $\sigma^2 = \dfrac{\sum_{i=1}^{N}(X_i - \bar{X})^2}{N}$ |
| 3 | Skewness | Symmetry (Positive or Negative) | $Skew = \dfrac{1}{N}\sum_{i=1}^{N}\left[\dfrac{(X_i - \bar{X})}{\sigma}\right]^3$ |
| 4 | Kurtosis | Shape (Tall or flat) | $Kurt = \dfrac{1}{N}\sum_{i=1}^{N}\left[\dfrac{(X_i - \bar{X})}{\sigma}\right]^4$ |

Where X is a random variable having N observations (i = 1,2,…,N).

# Binding site "fingerprints" (implemented)

Centroid (c1)  Atom closest to c1 (c2)  Atom furthest to c1 (c3)  Atom furthest to c3 (c4)

| m1 | m2 | m3 | | | | | | | | | |

12

3D

m1   First moment
m2   Second moment
m3   Third moment

# Binding site "fingerprints" (implemented)

# Binding site "fingerprints" (implemented)

# Binding site "fingerprints" (points with 4 dimensions)

**Points**

ElectroShape
$(x_1, x_2, x_3, q)$

ElectroShape for $z_1$
$(x_1, x_2, x_3, z_1)$

$\mathbf{c1}$      geometric centre, $\mathbb{R}^4$
$\mathbf{c2}$      atom furthest from $\mathbf{c1}$, $\mathbb{R}^4$
$\mathbf{c3}$      atom furthest from $\mathbf{c2}$, $\mathbb{R}^4$

$$\mathbf{a} = \mathbf{c2} - \mathbf{c1} \qquad \mathbf{a_s} \text{ only spatial part}, \mathbb{R}^3$$
$$\mathbf{b} = \mathbf{c3} - \mathbf{c1} \qquad \mathbf{b_s} \text{ only spatial part}, \mathbb{R}^3$$

$$\mathbf{c_s} = \left( \frac{\|\mathbf{a}\|}{2} \right) \frac{\mathbf{a_s} \times \mathbf{b_s}}{\|\mathbf{a_s} \times \mathbf{b_s}\|}$$

$$\mathbf{c4} = \mathbf{c1_s} + \mathbf{c_s} + (0,0,0,\mu q_+)$$
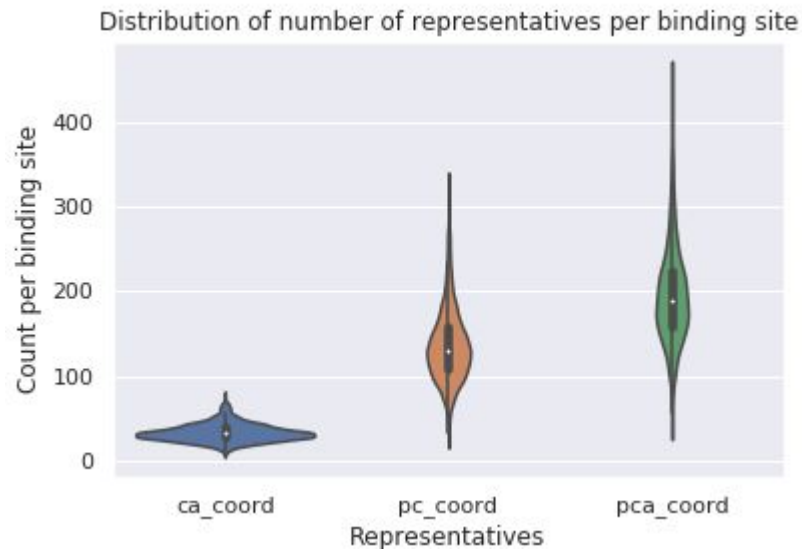$$\mathbf{c5} = \mathbf{c1_s} + \mathbf{c_s} + (0,0,0,\mu q_-)$$

$\mu$      scaling factor
$q_{+/-}$      highest/lowest value of 4th dimension in molecule

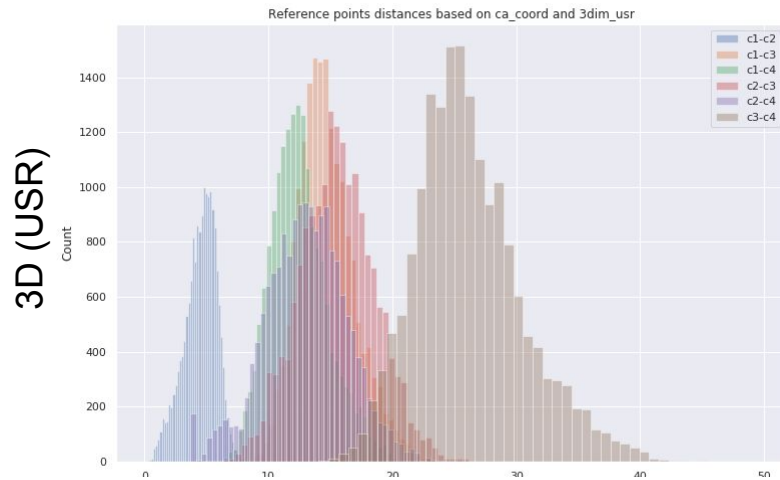| | Centroid (c1) | | | Atom furthest to c1 (c2) | | | Atom furthest to c2 (c3) | | | ~Cross product of c1-c3 (c4) | | | ~Cross product of c1-c3 (c5) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| m1 | m2 | m3 | | | | | | | | | | | | | |

4D     ...     15

# Cross product

$$\begin{pmatrix} a_x \\ a_y \\ a_z \end{pmatrix} \times \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} = \begin{pmatrix} a_y b_z - b_y a_z \\ a_z b_x - b_z a_x \\ a_x b_y - b_x a_y \end{pmatrix}$$

# Reference points distances Full scPDB

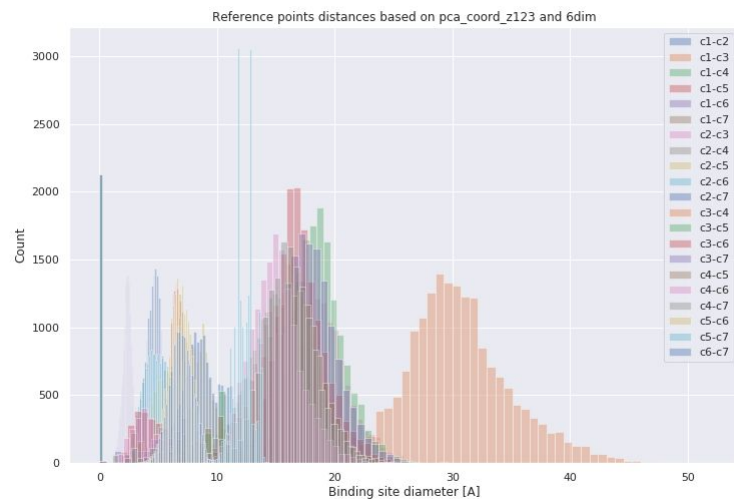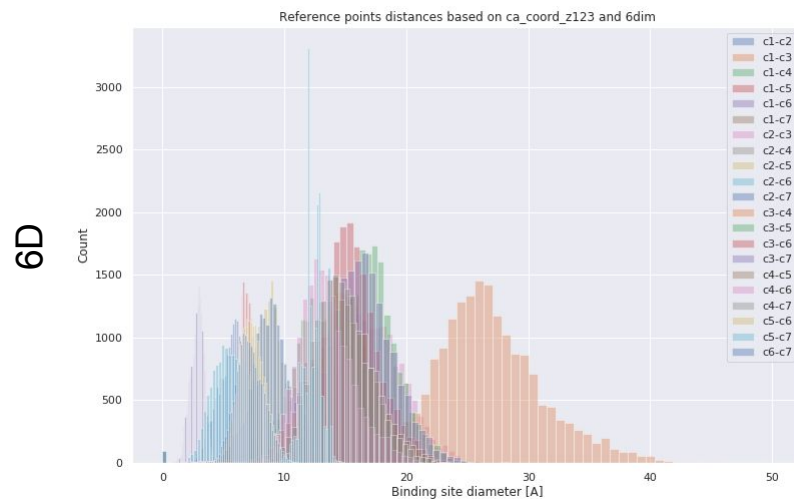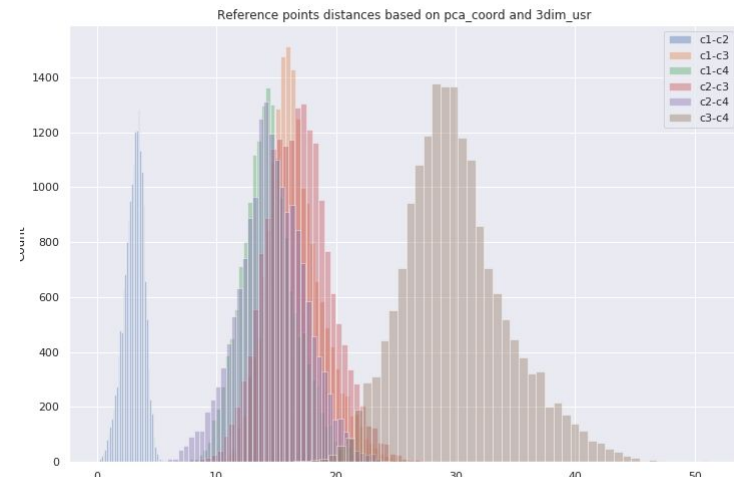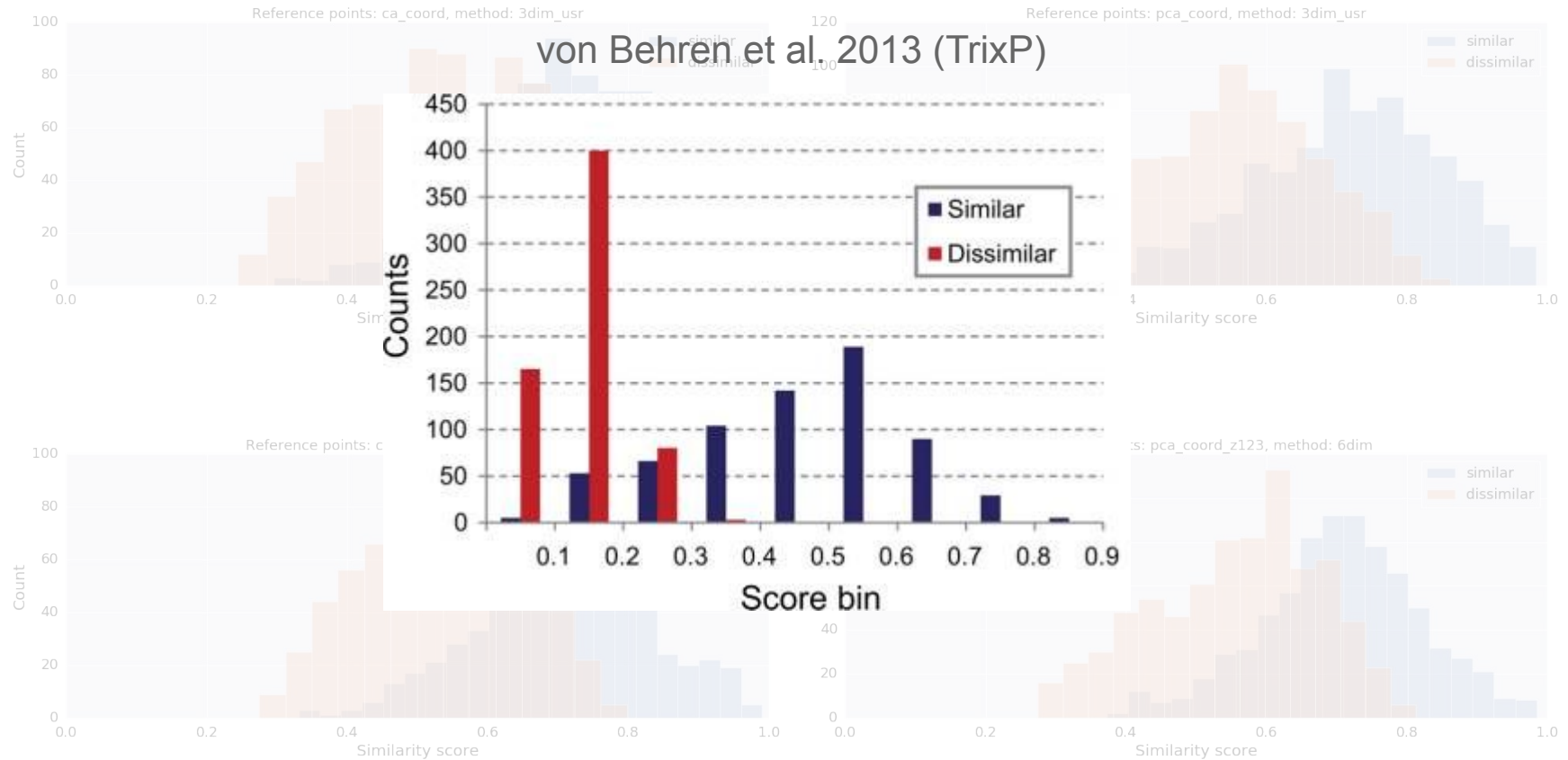# Encoding: PCA for binding sites (FuzCav: 2061 structures)



PCA for all binding sites

# Evaluation: similar vs. dissimilar pairs Weill et al. 2010 (FuzCav)

# Outlook: Introduce subpockets/regions

**Ideas**

- Overlapping/sliding window
  - Adapting von Behren et al. 2013 (TrixP)

- Overlapping subgraphs
  - Adapting Konc et al. 2010 (ProBis)

- Triangulation/Voronoi
  - Adapting Lindow et al. 2011

- Density-based clustering of binding site atoms
  - Adapting Oliver Lempke (group of Bettina Keller)