

#Todos Check tenses? Past or present Stick with atomistic or all-atom

Combining MD simulations with sPRE calculations for protein structure prediction

Abstract

The following pages contain the analysis basis (available MD data), the methods and, results of applying XPLORES to three differently linked (K6, K29, K33) ubiquitin dimers. These sections are meant to appear in the supplementary information of an upcoming paper. The figures can be used and included in the main part as one sees fit.

Available data

Previous work

CG data

MD data at coarse-grained (CG) resolution was taken from [1]. This dataset comprised the ubiquitination sites K6, K29 and K33 where the data for every protein is comprised of 10 replicas of 10 ms each. The MARTINI v2.2 force field was applied to the non-bonded interaction, whereas bonded interactions have been modeled using the IDEN method. These forcefields had to be altered to accommodate the isopeptide bond between the proximal subunit's lysine and the distal subunit's C-terminal glycine (Figure 1). The CG data of all 3 ubiquitylation sites encompass approximately 300 ms.

Back-mapped atomistic data

These simulations were used as input-data for a dimensionality reduction algorithm called sketch-map. Berg et al. used suitable high-dimensional collective variables to project the protein's phase space into 2D. In such representations similarity of any two conformations is represented in the distance in sketch-space. Naturally, basins of higher density represent statistically significant conformational sub-states of low free energy. Subsequently, new all-atom simulations were started by selecting points in sketch-space from both the center of these basins (4 simulations per ubiquitination site with 10 ns each) and from randomly chosen points in the vicinity of the basin (40 simulations per ubiquitination site with 3 ns each). Martini's **backward** script, paired with energy minimization to settle initial high-energy structures, yielded atomistic starting structures for these simulations.

Simulations started from extended structures

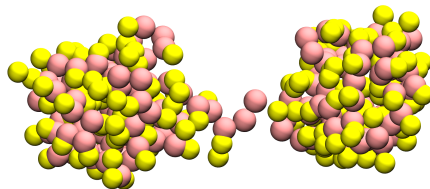


Figure 1: Coarse grained representation of K6-diUbi. Pink beads are unified backbone beads (C_α , N, and C). Yellow beads are unified sidechain beads (one bead for all non-GLY residues). The proximal subunit (right) offers its LYS6 residue to be ubiquitinated by the distal (left) unit’s GLY76 residue.

Additionally, to the back-mapped atomistic simulations, further atomistic data were taken from Berg et al. The starting structures of these simulations were created by slightly altering the position between the proximal and distal subunit of the di-ubiquitin proteins. As a reference for the monomer, the 1UBQ crystal structure from the protein database was used [2]. For every ubiquitination site, 12 simulations with 50 ns each are available. The atomistic simulations employed the GROMOS54a7 forcefield which was adapted to accommodate for the isopeptide bond, which was parametrized from regular peptide bonds.

MD simulations conducted for this work

The dataset of atomistic simulations was complemented with additional simulations started from starting structures created by rotating the sidechain χ_3 angle of the ubiquitinated lysine residue in 40 degree steps. The rotamer simulations were intended to accelerate the exploration of the phase space available to the diubiquitin proteins. This dataset encompasses 9 simulations per ubiquitination site with approximately 50 ns each.

Combination of CG and AA data

This might be interesting info for the main part

The simultaneous usage of CG and AA data for this project has been proven to be especially advantageous. Normally, all-atom simulations sample a protein’s phase space at a fraction of the speed available to CG simulations (approx. $3500 \frac{ns}{day}$ for CG and approx. $8 \frac{ns}{day}$ for AA in the case for diUbi). This disadvantage

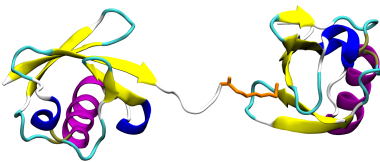


Figure 2: All-atom representation of K6-diUbi. The structure in this figure was used as a basis for Figure 1 and the location of proximal and distal units are identical. Secondary structure has been visualized as cartoon representation and colored accordingly. The residues of the isopeptide bonds (GLY76 (left) and LYS6 (right)) are colored orange.

led to sophisticated methods to accelerate sampling. Some of these so-called advanced sampling methods introduce a bias in the form of added potentials, which leads to unavoidable skew in the MD ensemble. Even starting multiple short AA simulations introduces a bias towards the starting structures. One could run long simulations to reduce this starting structure bias, but sooner or later will find themselves in a Catch-22 situation, where any efforts to accelerate the sampling will lead to additional bias. Pure CG simulations will not be able to break free from this problem, as they simply do not offer the resolution of AA simulations. We broke this vicious cycle by combining CG and AA data. Long CG simulations pave the way and define the general shape of the free energy surface and short AA simulations increase the resolution in the interesting parts. This can be seen in Figure 4.

All in all a single ubiquitination site is characterized by 100 ms of CG data and 1.21 ms of AA data.

Methods

Generally, the analysis steps were performed as follows: 1. Extract high-dimensional collective variables. 2. Project into 2D space using Encodermap. 3. Use HDBSCAN to extract statistically-significant sub-states of ensemble. 4. Calculate sPRE values for the all-atom conformations using XPLOR NIH (task has been parallelized). 5. Normalize the sPRE data from XPLOR NIH. 6. Identify some clusters that represent the ensemble available to di-ubiquitin using sPRE and HDBSCAN .

Extact high-dimensional collective variables

Understanding the dynamics of a protein can be difficult. Of course, its possible to look at a nicely rendered movie of how the protein might move. However, making assessments general conformations i.e. states is impossible from hour-long MD movies. For that we chose to use Encodermap to easily visualize overall protein dynamics as a 2D map. But before creating this map a set of collective variables needed to be defined. Collective variables represent a type of data that is somehow aligned with the raw xyz positional data of a molecular dynamics trajectories. With CVs complex processes can be broken down into meaningful variables. One such example might be the simplification of receptor-ligand docking processes by using a receptor-ligand distance as a collective variable. The torsion angles of the backbone are another widely used example of collective variables. From the ϕ and ψ dihedral angles a ramachandran plot can be created to visualize and simplify the complex structure of proteins. In this paper we faced the challenge to unify the CG and AA data, meaning CVs needed to be identical for both types of MD data. Torsional angles could not be calculated from CG data as the backbone N, C, and CA atoms are unified within one CG bead. Distance matrices between all CA atoms would yield a $\binom{15}{2} = 11476$ dimensional dataset. We chose to use the approach also used by Berg et al. employing residue-wise minimal distances, which makes the input data 152-dimensional.

Encodermap dimensionality reduction

The high-dimensional dataset consists of the 152 residue-wise minimal distances for each MD frame (CG or atomistic). The approximately XXXX frames per ubiquitination site, made it unfeasable to employ non-machine-learning dimensionality algorithms. They could not cope with the number of data (for example: multidimensional-scaling would require a $n_{frames} \times n_{frames}$ distance matrix, which would require ubnfeasable amounts of system memory. Sketch-map would require to pre-select points and reduce the amount of data actually used in dimensionality reduction). We used Encodermap’s auto-encoder neural network to retrieve a dimensionally reduced representation of all aforementioned simulations [3]. The high-dimensional input data was fed into a dense, fully-connected sequential neural network comprised 250, 250, 125, 2, 125, 250, and 250 neurons. Three main cost functions were minimized in the training phase to ensure a good 2D representation.

- Center cost: Mean absolute distance between all points and the coordinate origin. This cost function keeps the 2D points at the coordinate origin. It becomes higher, the further away points are from the coordinate origin.

$$CenterLoss = \frac{\sum_{i=1}^n \hat{y}_i^2}{n}$$

- Auto cost: Compare the mean absolute difference between input and output. This function becomes higher the greater the difference between the encoder input and decoder output are.

$$AutoLoss = \frac{\sum_{i=1}^n \|\sqrt{y_i^2 - \hat{y}_i^2}\|}{n}$$

- Encodermap cost: Compare the input and the latent space by transforming it with a sigmoid function. This function becomes higher the greater the sigmoid-weighted difference between the pairwise distances of encoder input and encoder output (i.e. latent space) are:

$$EncoderMapLoss = \frac{1}{m} \sum_{i \neq j} [SIG_h(R_{ij}) - SIG_l(r_{ij})]^2$$

,

where SIG is a sigmoid function originally devised by Ceriotti et al. for their sketch-map dimensionality reduction algorithm [4].

$$SIG_{\sigma,a,b}(r) = 1 - \left(1 + \left(2^{\frac{a}{b}} - 1\right) \left(\frac{r}{\sigma}\right)^a\right)^{-\frac{b}{a}}$$

R_{ij} and r_{ij} are the pairwise distance between all points in the high-dimensional input space and the 2-dimensional latent space. They are defined using the pairwise distance matrices *boldsymbol* $R_{i,j}$:

$$\mathbf{R}_{i,j} = \begin{pmatrix} r_{1,1} & r_{1,2} & \dots & r_{1,n} \\ r_{2,1} & r_{2,2} & & \\ \vdots & & \ddots & \\ r_{n,1} & & & r_{n,n} \end{pmatrix}$$

without the diagonal elements (which are always 0) and vectorizing them:

$$R_{ij} = \text{vec}(\mathbf{R}_{i,j,i \neq j}) = \begin{pmatrix} r_{1,2} \\ r_{1,3} \\ \vdots \\ r_{1,n} \\ r_{2,1} \\ \vdots \\ r_{n,n-1} \end{pmatrix}$$

The sigmoid function which transforms the high-dimensional input space and the low-dimensional latent space was using the following parameters:

σ_{highd}	A	B	σ_{lowd}	a	b
5.9	12	4	5.9	2	4

The resulting 2D projection could then be used as input to HDBSCAN clustering.

HDBSCAN clustering

The hierarchical density-based clustering algorithm HDBSCAN has been tried and tested for molecular dynamics systems and especially the low-dimensional projections of such [5]. Thus, this established method was chosen. After the 2D projections were obtained for every ubiquitination site, the xy-values were fed into HDBSCAN using the ‘leaf’ algorithm, with a minimal cluster size of 2500, resulting in 13, 18, and 24 selected clusters for K6, K29, and K33, respectively. The clusters are enumerated based on the number of included points/conformations, with cluster 0 being the largest cluster. Another notable point is that HDBSCAN offers a great advantage by distinguishing between clusters and noise. Noise - speaking from an ensemble viewpoint - is comprised of unique structures that carry no significance and thus can be disregarded.

Look at a specific cluster

As an example cluster 0 of K6 diUbi will be further investigated. This cluster consists of 131969 conformations sampled by XXX independent CG and XXX independent AA simulations and thus makes up roughly 12% of the complete K6 diUbi ensemble. Only 3% of the conformations in this cluster are from all-atom simulations (Figure 5). Visualizing 10 of the approximately 4000 all-atom conformations shows, that they are indeed structurally similar and this task succeeded (Figure 6).

XPLOR calculations

Check again.sentences are odd.

The NMR prediction and refinement program XPLOR NIH in version 3.3 was used to calculate the sPRE and 15N-relaxation time NMR observables. XPLOR’s new python functionality was used to call XPLOR functions from within python. XPLOR’s intended workflow follows the steps of structure loading, setting up potentials, setting up the dynamics using the IVM (internal variable module) and module and refining the loaded structure by minimizing the potentials iteratively. As we already had an extensive set of structures we dropped the refinement steps and let XPLOR calculate the initial potential values for the starting structure. To accelerate this undertaking, we parallelized this computation by using python’s joblib process-based multitasking. Every worker was provided with a number

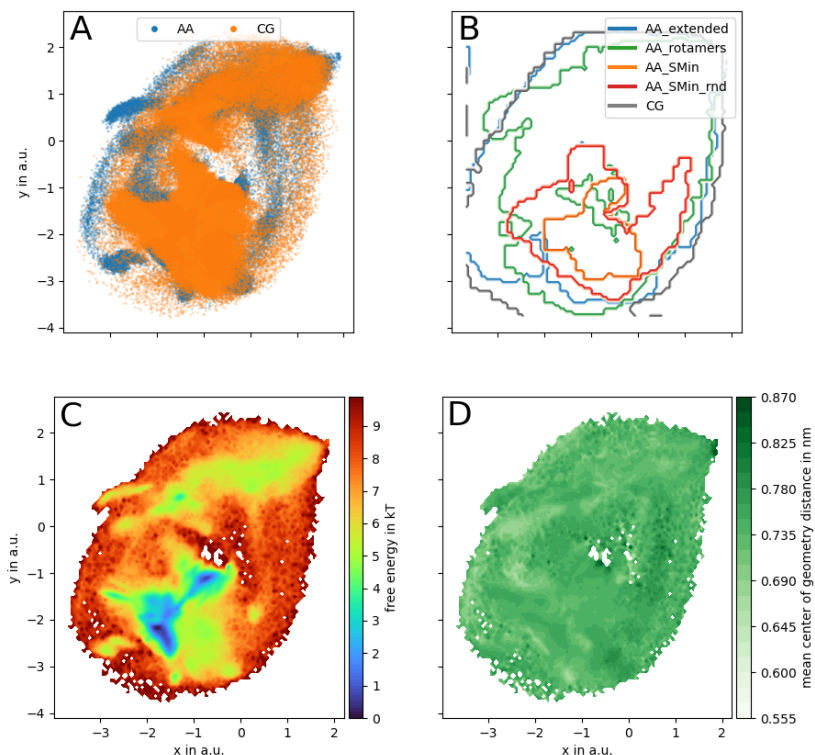


Figure 3: 2D projections of all CG and AA conformations of K6 diUbi. When the low-dimensional coordinates are plotted as a scatter plot, every point can be associated with one conformation sampled by a MD simulation (A). Orange points belong to CG simulations, blue points to AA. The closeness of two points. In (B) the areas occupied by specific MD simulation trajectories is visualized. Here, it can be seen, that the long AA extended simulations sampled more unique conformations, when compared to the limited regions of the sketch-map minima and sketch-map minima random sims. It should be noted, that the AA_rotamers were a more successful strategy in creating unique conformations, when compared to SMin and SMin_rnd. However, the AA_rotamer simulations are inherently starting-structure biased and might not yield a representative i.e. biased ensemble. That's why the long (and with than more unbiased) CG simulations are used to obtain a more correctly weighted ensemble. The Encodermap projections correlate to many other structural properties of the system. So can the distance between the center of geometries of the Ub subunits be used to color the projection and distinct regions can be observed (D).

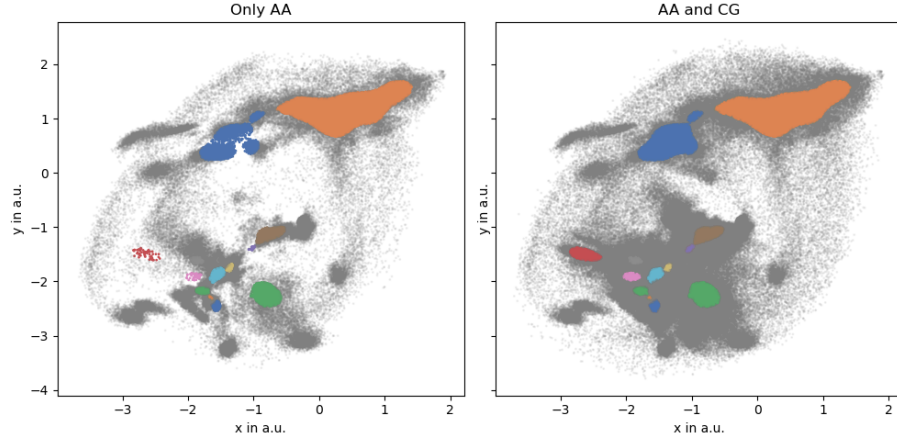


Figure 4: The Encodermap projection of K6 colored according to cluster membership. The advantage of using combined CG and AA data becomes apparent, when the red cluster is investigated. In a combined dataset this cluster contains more than 2500 structurally similar conformations in a dense basin (i.e. statistically significant sub-state of the K6 diUbi system), but multiple AA simulations were not able to sample these structures sufficiently, which would have removed these structures from the list of candidates for further analysis.

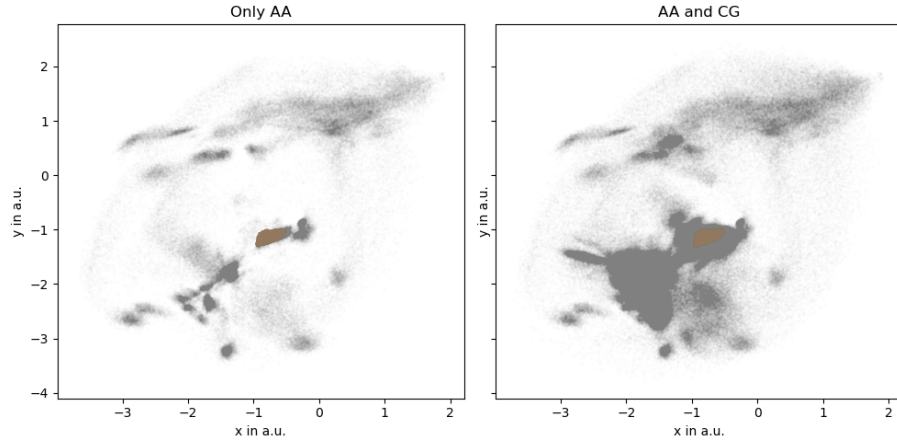


Figure 5: Encodermap projections of K6 diUbi, with cluster 0 highlighted.

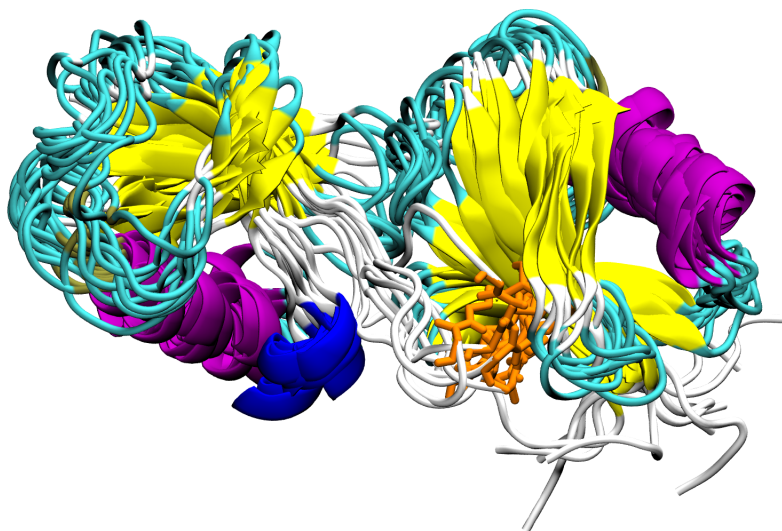


Figure 6: 10 conformations from cluster 0 of K6 diUbi. The structures are rendered as cartoon representation and colored accordingly. The isopeptide LYS and GLY residues are visualized as orange licorice (proximal subunit on the right, distal subunit on the left). It is notable, that the structural difference in the residues taking part in the isopeptide bond is larger, when compared to the protein backbone. The α -helices seem to overlap more precisely.

of simulation frames it should run XPLOR on and then carried out the task of saving structure files and calling XPLOR.

15-N relaxation times

As XPLOR NIH did not recognize the isopeptide bond in the provided pdb files (XPLOR’s pdbTool does not possess, contrary to the online documentation, a `readConect` method), we created a workaround that involved creating psf files for the three proteins with XPLOR’s pdb2psf program and manually adding the required bond. Some atoms in the pdb files then needed to be manually adjusted as they did not follow conventional naming (e.g. the hydrogens at the N-terminus are most commonly labelled as H1, H2, and H3. XPLOR’s pdb2psf program labelled these atoms as HT1, HT2, HT3). With these manual adjustments XPLOR’s relaxratiopot potential produced different output, compared to the old protocol, which did not use the psf files.

IMAGE

Code availability

The code developed for this publication will be made publicly available at: https://github.com/kevinsawade/xplor_functions

Normalization of XPLOR values The sPRE values from these calculations needed to be normalized before they could be quantitatively compared to experimental findings. We adjusted the normalization method described by Gong et al. for our system as follows [6]:

The simulated XPLOR values for the proximal and distal subunit build a $76 \times N$ matrix, respectively. Here, N is the number of aggregated frames over all trajectories for a specific di-ubiquitin protein. As ubiquitin has 76 residues, that number is fixed.

$$sPRE_{sim}^{76 \times N} = \begin{bmatrix} sPRE_{MET1,1} & sPRE_{MET1,2} & \dots & sPRE_{MET1,n} \\ sPRE_{GLN2,1} & sPRE_{GLN2,2} & & \\ sPRE_{ILE3,1} & sPRE_{ILE3,2} & & \\ \vdots & & \ddots & \\ sPRE_{GLY76,1} & sPRE_{GLY76,2} & \dots & sPRE_{GLY76,n} \end{bmatrix}$$

From New-Mexico experiments we identified some residues to exhibit a fast proteon exchange rate with the solvent and thus do not yield representative sPRE values.

This section needs to be checked and maybe extended with more info about new mexico

These residues have been removed from the considerations for the normalizations. Furthermore, the proximal and distal unit were considered separately, resulting

in 6 factors, two for K6, K29, and K33-linked di-ubiquitin, respectively. The per-residue variance over the simulation frames was calculated yielding a vector of variance:

$$\text{var} \left(\mathbf{sPRE}_{sim}^{76 \times N} \right) = \begin{bmatrix} \text{var}(sPRE_{MET1,1}) & sPRE_{MET1,2} & \dots & sPRE_{MET1,n} \\ \text{var}(sPRE_{GLN2,1}) & sPRE_{GLN2,2} & & \\ \text{var}(sPRE_{ILE3,1}) & sPRE_{ILE3,2} & & \\ \vdots & & \ddots & \\ \text{var}(sPRE_{GLY76,1}) & sPRE_{GLY76,2} & \dots & sPRE_{GLY76,n} \end{bmatrix}$$

Instead of selecting one singular deeply buried - and thus yielding most reliable sPRE values - residue, we chose to select the 10 residues with the smallest variance. For each of these 10 residues a normalization factor $f_{i,ubq,pos}$, where i is the residue, ubq the ubiquitination site and pos either proximal or distal, was calculated:

$$f_{i,ubq,pos} = \frac{v_{i,ubq,pos,exp}}{v_{i,ubq,pos,sim}}$$

where $v_{i,ubq,pos,exp}$ is the experimental sPRE value for residue i and $v_{i,ubq,pos,sim}$ is the mean simulated sPRE for residue i over N simulation frames. The normalization factors $F_{ubq,pos}$, which were used to normalize all simulated sPRE values were obtained as the mean of the 10 normalization factors $f_{i,ubq,pos}$:

$$F_{ubq,pos} = \frac{\sum_i^{10} f_{i,ubq,pos}}{10}$$

Metric of comparison: mean abs difference As a metric of difference between one set of sPRE-values and another we chose the mean absolute difference, which can in general be obtained as:

$$MAD = \frac{\sum_i^n \|x_i - \hat{x}_i\|}{n}$$

We exclusively used the MAD to compare experimental and simulated sPRE values, so the equation can be re-written as:

$$MAD = \frac{\sum_i^{152} \|x_{exp,i} - x_{sim,i}\|}{152}$$

With this measure, we could compare any simulation frame from the atomistic simulations to the experimental values (Figure 7).

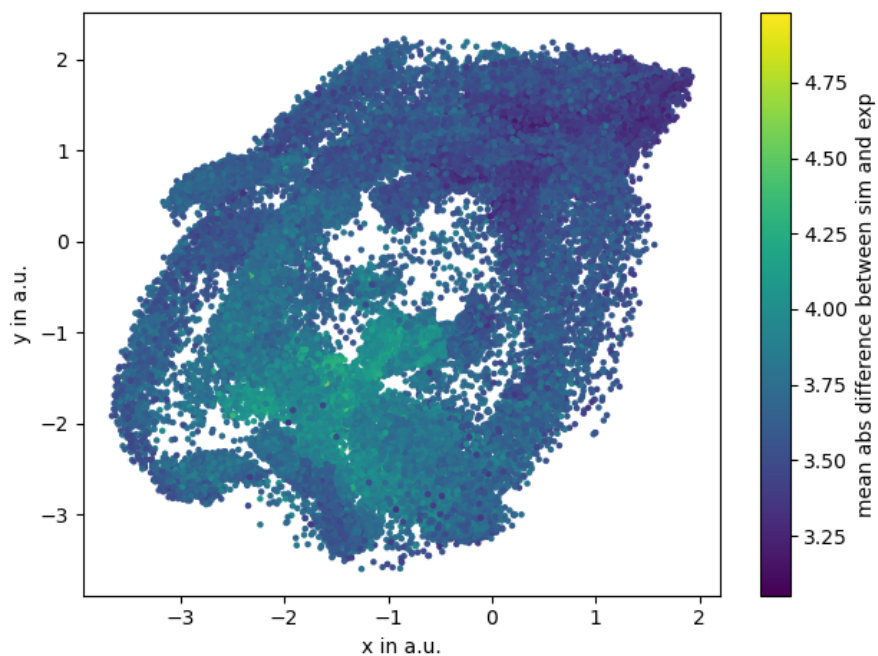


Figure 7: The low-dimensional Encodemap projection (x and y values) colored by the mean abs difference between a conformation from MD simulations and the experimental sPRE values. Certain regions with lower and higher mean abs difference can clearly be distinguished.

Identification of clusters

As a final step we conducted a linear combination to identify only a few clusters per ubiquitination site that should be able to represent the overall conformation of the di-ubiquitin proteins in a dynamic system. Starting from the experimental sPRE-values ($sPRE_{exp}$), we set up this system of equations:

$$\begin{bmatrix} sPRE_{exp,MET1,prox} \\ sPRE_{exp,GLN2,prox} \\ \vdots \\ sPRE_{exp,GLY76,dist} \end{bmatrix} = x_1 \cdot \begin{bmatrix} sPRE_{clu_1,MET1,prox} \\ sPRE_{clu_1,GLN2,prox} \\ \vdots \\ sPRE_{clu_1,GLY76,dist} \end{bmatrix} + \dots + x_n \cdot \begin{bmatrix} sPRE_{clu_n,MET1,prox} \\ sPRE_{clu_n,GLN2,prox} \\ \vdots \\ sPRE_{clu_n,GLY76,dist} \end{bmatrix}$$

Here, the sPRE-values are represented as column-vectors with proximal and distal values unified. The residues identified by the New Mexico experiment to exchange fast have been excluded from these considerations, which makes the column vectors shorter than 152 entries. For every residue there is one experimental value ($sPRE_{exp,MET1,prox}$) and multiple cluster-values. For every cluster the mean sPRE value of all conformations assigned to that cluster was taken, so that:

$$sPRE_{clu_1,MET1,prox} = \frac{\sum_n^i sPRE_{clu_1,MET1,prox,i}}{n_{\text{points in cluster}}}$$

Summarising the column vectors this equation can be written as.

$$sPRE_{exp} = x_1 \cdot sPRE_{clu_1} + x_2 \cdot sPRE_{clu_2} + \dots x_n \cdot sPRE_{clu_n}$$

The objective function, that was minimized can be written as:

$$f(\mathbf{x}) = \|\sum (\mathbf{x} \times sPRE_{sim}) - sPRE_{exp}\|$$

Here, \mathbf{x} is a vector of coefficients and $sPRE_{sim}$ is a matrix of the simulated sPRE values of the n clusters. The minimization was carried out with boundary conditions, so that any $x_i \in \mathbb{R}_{0 \leq x \leq 1}$ and $\sum_n^i x_i = 1$.

$$\min_{\mathbf{x} \in \mathbb{R}_{\geq 0}^n} f(\mathbf{x})$$

The results of this linear combination can be taken from Figure 8.

- [1] A. Berg, O. Kukharensko, M. Scheffner, and C. Peter, "Towards a molecular basis of ubiquitin signaling: A dual-scale simulation study of ubiquitin dimers," *PLoS computational biology*, vol. 14, no. 11, p. e1006589, 2018.

ubq site	cluster id	N	ensemble %	aa %	σ^2		coefficient	mean abs diff to exp w/ coeff	mean abs diff to exp w/o coeff
		frames			internal RMSD				
k33	0	119908	10	2	0.02	0.41		6.27	7.00
	1	95897	8	1	0.02	0.33		6.43	6.92
	2	31411	3	1	0.02	0.00		7.89	6.99
	3	30973	3	47	0.14	0.00		7.89	7.53
	4	29716	3	2	0.08	0.00		7.89	7.07
	5	29250	2	2	0.03	0.00		7.89	7.05
	6	22921	2	1	0.01	0.00		7.89	6.97
	7	21385	2	18	0.22	0.00		7.89	7.21
	8	21253	2	4	0.05	0.00		7.89	7.35
	9	20462	2	9	0.06	0.05		7.62	7.09
	10	9469	1	1	0.00	0.00		7.89	7.19
	11	9425	1	34	0.17	0.00		7.89	7.45
	12	9010	1	100	0.16	0.00		7.89	7.57
	13	7593	1	81	0.20	0.00		7.89	7.43
	14	6179	1	11	0.17	0.00		7.89	7.27
	15	6174	1	3	0.08	0.00		7.89	7.08
	16	5877	0	3	0.29	0.00		7.89	7.59
	17	5140	0	6	0.13	0.00		7.89	7.19
	18	4838	0	100	0.38	0.00		7.89	7.34
	19	4735	0	1	0.01	0.22		6.72	6.97
	20	4326	0	1	0.25	0.00		7.89	7.23
	21	3231	0	4	0.04	0.00		7.89	7.15
	22	3050	0	100	0.05	0.00		7.89	7.33
	23	2584	0	100	0.00	0.00		7.89	7.96
sum and combination	final linear Σ 504807	Σ 43				Σ 1.0	final combination: 6.7		

Figure 8: Solving the linear combination on the example of K29

		N	ensemble %	aa	σ^2	coefficient	mean abs diff to exp w/ coeff	mean abs diff to exp w/o coeff
ubq site	cluster id	frames		%	internal RMSD			
k29	0	186366	17	3	0.19	0.68	4.49	4.69
	1	99803	9	1	0.06	0.00	5.41	4.99
	2	48921	4	16	0.07	0.00	5.41	4.96
	3	44759	4	2	0.05	0.00	5.41	4.85
	4	28701	3	5	0.01	0.00	5.41	4.91
	5	22909	2	13	0.26	0.00	5.41	4.73
	6	16151	1	7	0.06	0.00	5.41	4.92
	7	14024	1	2	0.01	0.00	5.41	4.86
	8	12681	1	5	0.19	0.00	5.41	4.95
	9	11039	1	15	0.08	0.00	5.41	4.99
	10	7428	1	3	0.25	0.00	5.41	4.86
	11	7142	1	7	0.04	0.26	4.83	4.82
	12	6436	1	0	0.00	0.05	5.27	4.93
	13	5148	0	7	0.06	0.00	5.41	4.81
	14	4662	0	2	0.14	0.00	5.41	4.84
	15	4606	0	36	0.10	0.00	5.41	4.93
	16	4112	0	1	0.00	0.01	5.38	4.80
	17	3379	0	1	0.01	0.00	5.41	4.99
sum final linear Σ and combination 528267		Σ 48		Σ 1.0		final combination: 4.64		

Figure 9: Same results fro K29.

		N	ensemble %	aa	σ^2	coefficient	mean abs diff to exp w/ coeff	mean abs diff to exp w/o coeff
ubq site	cluster id	frames		%	internal RMSD			
k6	0	131969	12	3	0.01	0.21	4.92	4.55
	1	96726	9	1	0.10	0.47	4.56	4.71
	2	43256	4	4	0.21	0.06	5.29	4.64
	3	39818	4	1	0.10	0.00	5.43	4.77
	4	27491	2	74	0.12	0.00	5.43	4.78
	5	17686	2	0	0.14	0.00	5.43	4.70
	6	15884	1	32	0.03	0.00	5.43	4.87
	7	13093	1	17	0.11	0.00	5.43	4.67
	8	12241	1	3	0.02	0.00	5.43	4.73
	9	7014	1	1	0.03	0.00	5.43	4.82
	10	6479	1	18	0.21	0.00	5.43	4.75
	11	5848	1	0	0.05	0.26	4.82	4.63
	12	2955	0	3	0.03	0.00	5.43	4.80
sum final linear Σ and combination 420460		Σ 38		Σ 1.0		final combination: 4.58		

Figure 10: Saem results for K6.

- [2] S. Vijay-Kumar, C. E. Bugg, and W. J. Cook, "Structure of ubiquitin refined at 1.8 Å resolution," *Journal of molecular biology*, vol. 194, no. 3, pp. 531–544, 1987.
- [3] T. Lemke and C. Peter, "EncoderMap: Dimensionality reduction and generation of molecule conformations," *Journal of chemical theory and computation*, vol. 15, no. 2, pp. 1209–1215, 2019.
- [4] M. Ceriotti, G. A. Tribello, and M. Parrinello, "Simplifying the representation of complex free-energy landscapes using sketch-map," *Proceedings of the National Academy of Sciences*, vol. 108, no. 32, pp. 13023–13028, 2011.
- [5] L. McInnes, J. Healy, and S. Astels, "Hdbscan: Hierarchical density based clustering," *Journal of Open Source Software*, vol. 2, no. 11, p. 205, 2017.
- [6] Z. Gong, C. D. Schwieters, and C. Tang, "Theory and practice of using solvent paramagnetic relaxation enhancement to characterize protein conformational dynamics," *Methods*, vol. 148, pp. 48–56, 2018.