

A Reference Data Set for Circular Dichroism Spectroscopy Comprised of Validated Intrinsically Disordered Protein Models

Gabor Nagy¹, Søren Vrønning Hoffman², Nykola C. Jones² and Helmut Grubmüller^{1*}

¹: Department of Theoretical and Computational Biophysics, Max Planck Institute for Multidisciplinary Sciences, Göttingen, Germany

² ISA, Department of Physics & Astronomy, Aarhus University, Aarhus, Denmark

*: corresponding author: hgrubmu@mpinat.mpg.de

Keywords: Intrinsically Disordered Proteins, Circular Dichroism Spectroscopy, Reference Data Set, Secondary Structure Estimation, Protein Ensemble Refinement

1 **Abstract:**

2 Circular Dichroism (CD) spectroscopy is an analytical technique that measures the wavelength-
3 dependent differential absorbance of circularly polarized light, and is applicable to most biologically
4 important macromolecules, such as proteins, nucleic acids, and carbohydrates. It serves to
5 characterize the secondary structure composition of proteins, including intrinsically disordered
6 proteins, by analyzing their recorded spectra. Several computational tools have been developed to
7 interpret protein CD spectra. These methods have been calibrated and tested mostly on globular
8 proteins with well-defined structures, mainly due to the lack of reliable reference structures for
9 disordered proteins. It is therefore still largely unclear how accurately these computational methods
10 can determine the secondary structure composition of disordered proteins.

11 Here, we provide such a required reference data set consisting of model structural ensembles and
12 matching CD spectra for eight intrinsically disordered proteins. Using this set of data, we have assessed
13 the accuracy of several published CD prediction and secondary structure estimation tools, including
14 our own CD analysis package SESCA. Our results show that for most of the tested methods, their
15 accuracy for disordered proteins is generally lower than for globular proteins. In contrast, SESCA,
16 which was developed using globular reference proteins, but was designed to be applicable to
17 disordered proteins as well, performs similarly well for both classes of proteins. The new reference
18 data set for disordered proteins should allow for further improvement of all published methods.

19 **Introduction**

20 Circular Dichroism (CD) spectroscopy measurements serve to estimate the average secondary
21 structure (SS) content of proteins, to monitor protein folding under various experimental conditions,
22 and to determine folding kinetics.¹⁻⁴ Several CD-based SS estimation methods have been developed
23 either as web-based applications like DichroCalc⁵, K2D3⁶, BestSel⁷, and PDB2CD⁸ or as stand-alone
24 bioinformatics tools like SELCON3⁹, CCA¹⁰, and SESCA¹¹. Online tools and repositories such as
25 Dichroweb² and the Protein Circular Dichroism Databank¹² (PCDDDB) also allow easy access to these
26 tools and provide a platform for further development efforts (See Table S1 for available links).

27 CD spectroscopy is also often used to identify intrinsically disordered proteins (IDPs). IDPs form a
28 major class of proteins that fulfil their biological function without adopting a well-defined secondary
29 or tertiary structure under physiological conditions, and thus do not conform to the classical structure-
30 function paradigm¹³. Instead, IDPs often adopt a large number of partially folded transient structures,
31 and this conformational flexibility provides them functional advantages over their well-folded globular
32 counterparts^{14,15}. Rather than forming two distinct classes, the transition between ordered and
33 disordered proteins is continuous, and studies estimate that approximately 30% of human proteins
34 contain flexible or disordered domains.^{14,15} Because of their abundance and functional importance in
35 higher organisms, several tools have been developed to identify IDPs and intrinsically disordered
36 regions (IDRs) in otherwise folded proteins. Most of these methods are based either on protein
37 sequence, or the measured CD spectra of the respective regions.^{13,16,17}

38 The SS composition of proteins strongly affects their CD spectra.^{4,18,19} Structure-based predictions of
39 CD spectra using quantum-mechanical calculations are challenging and computationally
40 demanding^{20,21}; therefore, many CD-based analysis tools use reference data sets (RDS) instead to
41 empirically extract structure-spectrum relationships. For folded proteins, such RDSs are available,
42 consisting of proteins with known structures derived from X-ray crystallography, and respective CD
43 spectra²²⁻²⁴. Information from these data sets is often the basis of current algorithms that predict the
44 CD spectrum of a putative protein structure, or infer the unknown SS composition of a protein based

1 on its measured CD signal. Unfortunately, the conformational flexibility of IDPs and IDRs renders them
2 hard to characterize in terms of their structure both experimentally and computationally. Most IDPs
3 do not form regular crystals, and if they do, *e.g.*, in the presence of a binding partner, their crystal
4 structure usually does not reflect their conformational flexibility in solution. Due to the lack of reliable
5 IDP structural models, disordered proteins are largely absent from currently available RDSs, despite
6 the fact that their CD spectra are often published^{25–29}, and are distinctly different from that of folded
7 proteins.

8 The conformational flexibility of IDPs can be modelled through structural ensembles³⁰. Structural
9 ensembles (or ensemble models) consist of protein conformations and associated weights that
10 collectively describe the average protein structure and its fluctuations over time. Marked
11 improvements in simulation force fields^{31–34} and molecular modelling tools^{35–37} now allows one to
12 construct increasingly realistic ensemble models, which agree with or predict experimental
13 observables. Additionally, recent developments in prediction tools can now process structural
14 ensembles to predict observables such as, fluorescence spectra, nuclear magnetic resonance (NMR),
15 electron paramagnetic resonance (EPR), and small angle X-ray scattering (SAXS). These developments
16 have recently enabled more rigorous validations and further refinements of IDP ensembles.^{38,39}
17 Additionally, online repositories such as the protein ensemble database⁴⁰ (PED), the Biological
18 Magnetic Resonance Databank⁴¹ (BMRB), and the PCDDB¹² compile and link available data to facilitate
19 ensemble model generation. Finally, regarding the prediction of CD spectra, our CD analysis package
20 SESCO can predict CD spectra not only of individual protein structures, but also of structural
21 ensembles¹¹, and estimate the SS composition of proteins based on their measured CD spectrum⁴².
22 Initially, due to the lack of available IDP RDSs, SESCO was parametrized and validated on folded
23 proteins only.

24 These advances, taken together, now enable us to provide a small RDS, namely IDP8, consisting
25 measured CD spectra and structural ensembles for eight disordered proteins. Further, we will use this
26 newly constructed RDS to assess the accuracy of several established modelling tools for either CD
27 prediction or SS estimation of IDPs including the current version of our own SESCO analysis package.
28 Our analysis indicates that the IDP8 RDS offers the opportunity not only to assess the prediction
29 accuracy of CD-based analysis tools regarding disordered proteins, but to further improve their
30 accuracy and precision as well.

31 **Materials and Methods**

32 **Reference data set assembly**

33 The IDP8 RDS consists of eight IDP CD spectra and 14 structural ensembles, which were assembled
34 with the aim of testing the accuracy of CD-based prediction and SS estimation methods. The RDS
35 includes eight disordered protein models : 1) α -synuclein (*asyn*), 2) the measles virus nucleoprotein
36 tail domain (*mevn*), 3) *Saccharomyces cerevisiae* CDK inhibitor N-terminal targeting domain (*sic1*), 4)
37 the human tau protein K18 fragment (*tk18*), 5) the activator of thyroid hormone and retinoid receptor
38 protein activation domain 1 (*actr*), 6) CREB-binding protein nuclear coactivator binding domain (*cbpn*),
39 7) the protein 53 N-terminal transactivation domain (*p53t*), and 8) an RS-repeat peptide (*rsp8*,
40 sequence: GAMGPSYGRSRSRSRSRSRSRSRS). Two of these models are full length IDPs (*asyn* and
41 *rsp8*), and the other six models (*mevn*, *sic1*, *tk18*, *actr*, *cbpn*, *p53t*) are IDRs of larger proteins. All eight
42 disordered models were selected based on the availability of experimental and modelling data.

43 The CD spectra of the proteins included in the IDP8 RDS are shown in Fig. 1. We measured the CD
44 spectra of *actr*, *asyn*, *cbpn*, *p53t*, and *rsp8* using a synchrotron radiation source (SR-CD), which allowed

1 us to determine additional short wavelength information (down to 178 nm). The remaining three CD
2 spectra of *mevn*, *sic1*, and *tk18* depicted in Fig. 1 were measured using conventional CD
3 spectrophotometers. Due to high absorbance and a weaker UV source in conventional CD
4 spectrophotometers, measurements at short wavelengths are unreliable for these spectra and
5 therefore were truncated to the wavelengths provided in Table 1. Further details are provided in the
6 Circular Dichroism measurements section.

7 The 14 structural ensembles of the data set shown in Fig. 2, are organized into three groups A, B, and
8 C based on the experimental data used in their creation. Group A contains four IDP model ensembles
9 for *asyn*, *mevn*, *tk18*, *sic1* that were previously published on the PED⁴⁰ under accession codes provided
10 in Table 2. These ensembles were fitted mainly against results from NMR measurements, partly
11 complemented by data from electron-spin paramagnetic resonance (EPR), small angle X-ray scattering
12 (SAXS), and residual dipolar coupling (RDC) experiments. For these ensembles, CD spectra were not
13 used during the ensemble refinement process. Group B consists of five IDP model ensembles for *mevn*,
14 *actr*, *cbpn*, *p53t*, and *rsp8*. These ensembles were refined from large molecular dynamics (MD)
15 simulation ensembles using the Bayesian Maximum Entropy (BME) approach to fit against measured
16 CD spectra, SAXS curves, and NMR C_α chemical shifts as described below. Finally, Group C contains
17 five model ensembles of the same five IDP domains as group B, but here, the refinement was carried
18 out without the CD information. Separating the ensemble models into three groups allowed us to
19 compare the average accuracy of BME refined ensemble models to established structural ensembles
20 (group A vs. group-C), and to assess the effects of including CD spectra in the refinement process
21 (group B vs. group C).

22 Protein sample preparation

23 The protein samples for four IDP domains were manufactured by the company Karebay and were
24 delivered in a lyophilized form. These samples included *actr*, *cbpn*, *p53t* (13-61), and an RS repeat *rsp8*.
25 Samples for two other variants of *p53t* (1-73 and 1-94), as well as *asyn* were kindly provided by S.
26 Becker, Max Planck Institute for Multidisciplinary Sciences, Department of NMR-based Structural
27 Biology, Göttingen, Germany. All seven listed protein samples were dissolved in a 10 mM sodium-
28 phosphate buffer, pH 7.2, including 50 mM NaF for electrostatic screening. A summary of the
29 sequence details of the IDP8 model proteins is provided in Table 2.

30 Circular Dichroism measurements

31 Circular Dichroism (CD) spectra for seven of the protein samples described above were recorded on
32 the AU-CD beamline of the ASTRID2 synchrotron radiation source, at the Department of Physics &
33 Astronomy, Aarhus University, Denmark. The spectra were measured at 25 °C using a 0.1 mm quartz
34 cuvette under a nitrogen atmosphere. The CD intensities were recorded every 1 nm, with an average
35 of 2 seconds per measurement. The final CD spectrum was calculated as the smoothed average of five
36 independently measured and baseline corrected spectra recorded between 178-280 nm. Spectra
37 were smoothed using a 7 pt Savitzky-Golay filter. The protein samples of *actr*, three variants of *p53t*,
38 *cbpn*, *rsp8*, and *asyn* were measured in the buffer solution as described above. Protein concentrations
39 for CD measurements were between 0.3 – 1.5 g/L, calculated from sample UV absorption at 280 nm
40 as well as 214 nm. The molar extinction coefficient at 214 nm was estimated based on the protein
41 sequence using the method proposed by Kupiers *et al.*⁴³

42 The CD spectrum for *mevn* was kindly provided by Longhi *et al.*⁴⁴. This spectrum was measured in a
43 Jasco-810 dichrograph using a 1 mm quartz cuvette, 7 μM protein sample in a 10 mM sodium-
44 phosphate buffer, pH 7.0, at 20 °C, under nitrogen atmosphere. The CD spectrum of *sic1* was kindly
45 measured and provided by Chong *et al.* (personal communication). It was measured in a Jasco-1500

1 CD spectrophotometer using a 0.1 mm quartz cuvette under nitrogen atmosphere, at 25 °C. The
2 measured sample had a 10 µM protein concentration, dissolved in a 50 mM potassium-phosphate
3 buffer, pH 7, and included 150 mM NaCl and 1 mM EDTA. The CD spectrum of *tk18* was extracted from
4 the work of Braghorn *et al.*⁴⁵.

5 Small Angle X-ray scattering

6 Small-angle X-ray scattering curves were measured for *actr*, three *p53t* variants, *cbpn*, and *asyn* at the
7 European Synchrotron Radiation Facility (ESRF) Grenoble, France, at the BioSAXS beamline BM29 in
8 2018. All measurements were carried out under sample flow to reduce the effects of radiation damage
9 during the measurement. SAXS curves were collected over 10 data frames of 0.3 seconds each. The
10 measured scattering curves were normalized for the protein concentration, corrected for the buffer
11 signal, and averaged to obtain the final scattering curves. Data processing and automated analysis was
12 done using the Edna software package⁴⁶. Samples were measured under similar conditions as
13 described above, with protein concentrations ranging from 2-8 g/L.

14 The SAXS curve for *mevn* was kindly provided by Longhi *et al.*, measured at the European Synchrotron
15 Radiation Facility (ESRF) using a 10 mM Tris/Cl buffer (pH 8) containing 10% glycerol and 600 µM *mevn*
16 at 8 °C. The SAXS curve for *rsp8* was kindly provided by Rauscher *et al.*, which was measured at 25 °C
17 in a 50 mM sodium-phosphate buffer (pH 7), at a concentration of 750 µM *rsp8* and 100 mM NaCl.
18 The SAXS curves for *tk18* and *sic1* were extracted from the studies of Mylonas *et al.*⁴⁷ and Mittag *et*
19 *al.*⁴⁸, respectively.

20 Nuclear Magnetic resonance chemical shifts

21 Backbone chemical shifts for *asyn*, *sic1*, *tk18*, *actr*, *cbpn*, *p53t* were downloaded from the Biomagnetic
22 Resonance database (BMRB)⁴¹: entry numbers 19257, 16657, 19253, 15397, 16363, and 17660,
23 respectively. Backbone chemical shifts for *mevn* were measured by Gely *et al.*⁴⁹ and kindly provided
24 by S. Longhi. The chemical shifts were determined for a 500 µM *mevn* sample in 10 mM sodium-
25 phosphate buffer with 50 mM NaCl, 1 mM EDTA, and 5% D₂O, at 25 °C, pH 6.5. The backbone chemical
26 shifts of *rsp8* were measured and kindly provided by Rauscher *et al.*³⁴. These chemical shifts were
27 measured for a 750 µM peptide sample in a 50 sodium-phosphate buffer, 100 mM NaCl, at 25 °C and
28 pH 7.

29 Molecular Dynamics simulations

30 All-atom molecular dynamics (MD) simulations for *mevn*, *actr*, *cbpn*, *p53t*, and *rsp8* were carried out
31 using the GROMACS 2019 software package⁵⁰. All simulations were performed at a constant
32 temperature of 25 °C, constant pressure of 1 atm in a dodecahedral simulation box filled with explicit
33 water molecules and periodic boundary conditions. To accommodate extended IDP conformations,
34 simulation box radii were chosen to be larger than the expected radius of gyration by at least 2.5 nm,
35 resulting in system sizes of 60 000-300 000 atoms. Sodium and chloride ions were added to all
36 simulation boxes to obtain neutral systems with NaCl concentrations of 50-150 mM. Details on
37 simulation conditions, used force fields, and length of simulation trajectories for individual IDPs are
38 provided in Table 5.

39 The temperature was kept constant by using the velocity rescaling algorithm⁵¹ and a coupling constant
40 of 0.1 ps. Pressure was maintained by the Parrinello-Rahman barostat⁵² using a 0.1 ps coupling
41 constant and the isothermal compressibility of water $4.5 \times 10^{-5} \text{ bar}^{-1}$. Simulations were propagated
42 using a leapfrog integrator⁵³ with 4 fs time steps. To enable such large time steps, fast vibrational
43 degrees of freedom were removed by using the LINCS algorithm⁵⁴ and applying a sixth order iterative
44 restraint on the bond angles. Apolar hydrogen positions were described using virtual atom sites⁵⁰ to

1 eliminate hydrogen bond vibrations. Electrostatic and van der Waals interactions were explicitly
2 calculated within a cutoff distance of 1.0 nm. Electrostatic interactions beyond the cutoff distance
3 were calculated by particle-mesh Ewald summation⁵⁵ with a grid spacing of 0.12 nm. Long-range van
4 der Waals dispersion corrections⁵⁶ were applied to the total energy of the system in all simulations.

5 MD trajectories were generated using six different force fields, for which the accuracy for IDPs has
6 been assessed previously^{34,57,58}. These force fields include the Amber03 force field⁵⁹ with a modified
7 TIP4P water model⁶⁰, an Amber99SB parameter set with a modified TIP4P water model⁶⁰, the
8 Amber99SB-disp force field with a modified TIP4P water model re-parametrized with dispersion
9 corrections⁶¹, the Amber14SB⁶² force field with an Optimal Point Charge (OPC) water model⁶³, the
10 CHARMM22* force field⁶⁴ with a modified TIP3P water model⁶⁵, and the CHARMM36M⁶⁶ force field
11 with an OPC water model.

12 To provide initial conformations for the BME refinement, conformations were taken at 1-100 ns
13 intervals from 5-60 MD simulation trajectories per system amounting to total simulation times of 30-
14 800 μ s. Starting conformations for these simulations were either extended disordered structures or
15 conformations observed in the crystalized complex structures published in the protein data bank⁶⁷
16 (PDB) entries 1KB6, 2L14, and 1ZQO, respectively.

17 Bayesian Maximum Entropy refinement

18 Structural ensembles of group B and C for *mevn*, *actr*, *cbpn*, *p53t*, and *rsp8* were obtained using BME
19 refinement.⁶⁸ Table 6 summarizes the refinement parameters used for each IDP model. For each IDP
20 model an initial ensemble was formed from 5,000 to 50,000 conformations obtained from all-atom
21 MD simulations described above. Uniform prior weights were assigned to each conformation of the
22 initial ensembles. For each conformation, CD spectra, backbone carbon chemical shifts, and SAXS
23 curves were computed using the SESCA (V0.96)¹¹, Sparta+ (V2.6)⁶⁹, and CRY SOL (ATSAS V2.7.2.5)⁷⁰
24 analysis software packages, respectively.

25 All conformations of the initial ensembles were reweighted using the BME approach such that the re-
26 weighted ensemble fits O_{ki} , the i^{th} measured observable of type k , as best as possible, while at the
27 same time minimizing the loss of relative entropy $S_{\text{rel}} = -w_j \cdot \log(w_j/w_j^0)$ from redistributing the
28 conformation weights. Here, k and i denote the type and index of the observable, while j is the index
29 of conformations. The redistributed (posterior) weights w_j were obtained by minimizing

30

$$31 \quad L(O_{ki}, w_j, w_j^0, \theta) = \frac{m}{2} \cdot \sum_k \chi_k^2(O_{ki}, w_j) - \theta \cdot S_{\text{rel}}(w_j, w_j^0), \quad (1)$$

32 where w_j^0 are the initial (uniform) weights of each conformation, θ is the scaling parameter for the
33 entropy loss, and $M = \sum_k M_k$ is the total number of fitted observables, and M_k is the number of fitted
34 observables of type k . The deviation from the observables was quantified by the χ^2 deviations
35 between each observable computed from the reweighted ensemble $O_{kij}^{\text{calc}} = \sum_j w_j \cdot O_{kij}^{\text{calc}}$ and the
36 measured observable O_{ki}

37

$$38 \quad \chi_k^2 = \frac{1}{M_k} \sum_i \left(\frac{O_{kij} - \alpha_k \cdot O_{kij}^{\text{calc}}}{\sigma_{ki}} \right)^2, \quad (2)$$

39

1 where α_k is the uniform scaling factor to match the measured and calculated observables of type k ,
2 and σ_{ki} is the uncertainty of the observable O_{ki} . For group B ensembles, three types of measured
3 observables were used for BME refinement: 1) the intensities of the measured CD spectra, 2) the SAXS
4 intensities, and 3) the resolved C_α backbone chemical shifts of each residue. For group C ensembles,
5 only SAXS intensities and C_α chemical shifts were used. We note that the scaling factor α_k was used
6 to compensate for the machine-dependent beam intensity for SAXS measurements ($\alpha \in R_+$). For the
7 CD spectrum intensities and NMR chemical shifts, no such compensation was required, hence α was
8 set to 1.

9 The uncertainty σ_{ki} for SAXS measurements was defined as the SD of obtained SAXS intensities at
10 scattering vector q_i . The uncertainty of the backbone carbon chemical shift i was set to 0.95, 1.03, and
11 1.13 ppm for C_α , C_β , and carboxylate C shifts, respectively, to reflect the uncertainty of SPARTA+
12 chemical shift predictions⁶⁹. These are conservative estimates that are considerably larger than the
13 0.1-0.4 ppm errors indicated in available BMRB entries of *actr*, *p53t*, and *sic1*. The uncertainty of CD
14 intensities was computed as $\sigma_{ki} = \delta_k \cdot O_{ki} + \sigma_k^0$, where $\delta_k = 0.2$ represents the typical uncertainty
15 of concentration determination⁷¹, and $\sigma_k^0 = 0.75$ kMRE is the machine error of CD measurements,
16 determined from the average SD of obtained CD intensities upon repeated measurements.

17 The refinement parameter θ controls the balance between close agreement with measured
18 observables and reducing the effective ensemble size. To find the optimal θ parameter for each model
19 ensemble, several refinements with $\theta = \{0.1, 1, 2, 5, 10, 20, 50, 100, 200\}$ were carried out while
20 monitoring the computed χ_k^2 values. The refinement with the largest θ and significant improvements
21 to χ_k^2 values was selected and was used to draw sub-ensembles that constitute the final ensemble
22 models.

23 To obtain the final ensemble models, smaller sub-ensembles of 5, 10, 20, 50, 100, and 200
24 conformations were drawn at random by rejection sampling based on the redistributed weights of the
25 conformations after refinement. Conformations with high redistributed weights in the initial ensemble
26 may be included multiple times in the final ensemble to represent their importance. To assess the
27 effect of the sub-ensemble size on the model accuracy, five sub-ensembles were drawn and the
28 deviation from experimental observables was computed and averaged for each size. The sub-
29 ensembles of each size were concatenated to form a combined ensemble model for each IDP.
30 Deviations from the measured observables were calculated for the concatenated ensemble as well.
31 Finally, the ensemble with the smallest size was selected for each IDP that met the two following
32 criteria: 1) increasing the ensemble size further does not improve the average χ_k^2 deviations
33 considerably, and 2) the average χ_k^2 deviations of sub-ensembles are similar to χ_k^2 deviations of the
34 concatenated ensemble within uncertainty. The selected ensemble sizes and θ parameters and for all
35 derived IDP models are summarized in Table 6. This procedure yielded small ensemble models of 100-
36 250 conformations with integer weights for each refined ensemble.

37 Accuracy of the predicted CD spectra

38 To assess the accuracy of the CD spectrum predicted for protein j , the predicted CD spectrum was
39 compared to the measured spectrum by computing the root mean squared deviation (RMSD) of CD
40 intensities,

41

$$42 \quad \text{RMSD}_j^{\text{CD}} = \sqrt{\frac{1}{N} \cdot \sum_{\lambda}^N \left(\alpha_j \cdot I_{j\lambda}^{\text{exp}} - I_{j\lambda}^{\text{calc}} \right)^2}, \quad (3)$$

1

2 expressed in 1000 mean residue ellipticity units (kMRE, 1000 deg cm²/dmol) for each wavelength λ
3 for which both the measured and the predicted spectrum were available. Here, N is the number of
4 available wavelengths, I_{λ}^{exp} and $I_{\lambda}^{\text{calc}}$ are the measured and predicted CD intensities, respectively. The
5 scaling factor α_j minimizes the RMSD described above and accounts for experimental spectrum
6 normalization errors.

7 To assess the accuracy of CD spectrum prediction methods for disordered proteins, the CD spectra of
8 all model IDPs in the IDP8 RDS were predicted from their structural ensembles, the deviation from
9 their measured reference CD spectra were computed, and the resulting $\text{RMSD}_j^{\text{CD}}$ values were
10 averaged to determine the mean accuracy of the respective method. A similar approach was followed
11 to assess the mean accuracy of each studied CD prediction method for globular proteins, using the
12 reference structures and CD spectra of the SP175 RDS²².

13 Accuracy of estimated SS fractions

14 To determine the accuracy of SS estimation methods, the RMSD of SS fractions was computed for each
15 protein j in globular and disordered protein RDSs as follows. For globular proteins of the SP175 RDS²²
16 the SS fractions estimated from their CD spectra were compared to the respective reference
17 structures derived from X-ray diffraction measurements. For the disordered proteins of the IDP8 RDS,
18 estimated SS fractions were compared to those computed from the reference ensemble models. The
19 RMSD between the estimated and reference SS was computed as

20

$$21 \quad \text{RMSD}_j^{\text{SS}} = \sqrt{\frac{1}{M} \cdot \sum_k^M (F_{jk}^{\text{est}} - F_{jk}^{\text{calc}})^2}, \quad (4)$$

22

23 where M is the number of SS classes within the classification method, and F_{jk}^{est} and F_{jk}^{calc} are the
24 estimated and computed fractions of SS class k , respectively.

25 To be able to apply the RMSD determination according to eq. 4, the SS fractions computed from the
26 reference structures/ensembles by an SS classification method have to be grouped and identified with
27 the classes of the SS estimation method. For SESCA, the documented calculation of SS fractions from
28 the protein structure was used¹¹. For basis sets DS-dTSC3 and DS5-4SC1, the SS composition was
29 computed using the DISICL algorithm⁷², and the SS elements were grouped into three and six SS
30 classes, respectively. For the DSSP-1SC3 basis set, the SS composition was determined using the DSSP⁷³
31 algorithm and the obtained SS fractions were grouped into four SS classes. For the HBSS-3SC1 basis
32 set, the HBSS¹¹ algorithm was used, and the obtained SS composition was grouped into five SS classes.
33 To assess the accuracy of the K2D3 algorithm, SS classification was performed the same way as for the
34 DS-dTSC3 SESCA basis set, and the Alpha(-helix) and Beta(-sheet) sheet fractions were compared to
35 the corresponding estimated SS contents. The third SS fraction (Coil) for the K2D3 composition was
36 computed as $F_{j,\text{Coil}} = 1 - (F_{j,\text{Beta}} + F_{j,\text{Alpha}})$.

37 Finally, to assess the accuracy of the BESTSEL SS estimates, the SS fractions of the reference models
38 was determined by the HBSS algorithm, which uses similar helix and advanced β -sheet classifications.
39 The obtained fractions were grouped into six SS classes as follows: The Helix-1 and Helix-2 classes of
40 BESTSEL were grouped into a common Helix class, which was identified with the 4-Helix class in HBSS.
41 The three anti-parallel β -sheet classes (Anti1-3) were kept separate and were identified with the

1 corresponding HBSS classes (left-handed, non-twisted, and right-handed β -strands). All parallel β -
2 strand classes in HBSS were merged and identified with the parallel β -sheet class of BESTSEL. The SS
3 fractions of all other classes in BESTSEL and HBSS were merged and identified with an “Other” SS class,
4 resulting in six SS classes for both algorithms.

5 *Results and discussion*

6 **Model quality assessment**

7 First, we assessed how well the models of the IDP8 RDS agree with SAXS and NMR chemical shift
8 measurements (agreement with CD spectra will be discussed below). Table 3 shows how well the
9 observables predicted from the model ensembles of the RDS agree with the measured SAXS data as
10 well as with $C\alpha$, $C\beta$ and carbonyl-carbon (CO) chemical shifts. These two groups of observables were
11 chosen due to their complementarity; whereas SAXS curves report overall IDP compactness, carbon
12 chemical shifts are sensitive to the local secondary structure.

13 The χ values for SAXS curves shown in the second column of Table 3 are square roots of the χ^2 metric
14 defined by Sevrugun *et al*⁷⁰. This metric is insensitive to any scaling differences between the measured
15 and predicted SAXS intensities, and reports the deviation in units of the experimental uncertainty
16 determined by σ_i , the standard deviation of the scattering intensities. For seven of the 14 ensemble
17 models, the χ values are below one, meaning that predicted SAXS intensities are on average well
18 within the experimental uncertainty. The remaining seven models achieved χ values between one and
19 two, resulting in an overall average χ of 1.14 for the whole RDS. This result suggests that the size
20 distributions of the model ensembles agree with the available experimental data, except for the *actr*
21 and *cbpn* ensembles for which the predicted SAXS curves deviate from the experiment with χ values
22 between 1.8 and 2.1.

23 Columns three to five in Table 3 report the root-mean-squared deviation (RMSD) of carbon chemical
24 shifts for each model ensemble. The average RMSDs of the data set are 0.46 ppm, 0.49 ppm, and 0.46
25 ppm for $C\alpha$, $C\beta$, and CO chemical shifts, respectively. These RMSD values are slightly larger than the
26 0.1 - 0.4 ppm estimated experimental uncertainty reported in BMRB entries, but are considerably
27 smaller than the average 1.14 ppm, 0.94 ppm, and 1.09 ppm backbone chemical shift deviations
28 reported by Shen and Bax obtained in the context of Sparta+ prediction assessments from high-quality
29 crystallographic structures of globular proteins.⁶⁹

30 To assess the effects of using CD spectrum information in ensemble refinement, we compared the
31 average deviations between predicted and measured SAXS and NMR data for ensembles of group A,
32 B, and C separately. In addition, we also computed the SAXS and NMR deviations of the initial MD
33 ensembles (henceforth group 0) group B and C ensembles were refined from. SAXS intensities and $C\alpha$
34 chemical shifts were used as fit variables during both group B and group C ensemble refinements.

35 The average deviation from measured SAXS curves is within the average uncertainty for group A, as
36 shown by a mean χ value of 0.87 ± 0.17 . Refinement reduced the deviation from measured SAXS
37 curves from an initial χ of 1.59 ± 0.39 to 1.19 ± 0.27 for group B and 1.17 ± 0.26 for group C, showing
38 no significant difference between the two groups. The average deviation of $C\alpha$ chemical shifts for
39 group A is 0.62 ± 0.2 ppm, which is very similar to the deviation of 0.63 ± 0.1 ppm for initial group 0
40 ensembles. Ensemble refinement improved the average deviation from measured $C\alpha$ chemical shifts
41 to 0.43 ± 0.1 ppm for both group B and group C ensembles.

42 The deviations from measured $C\beta$ and CO chemical shifts were not used in ensemble refinements, and
43 thus are used for cross-validation. The average deviation of $C\beta$ chemical shifts in group A is 0.63 ppm.

1 In comparison, the C β chemical shifts are accurately reproduced by the initial MD ensembles with an
2 average C β shift deviation of 0.37 ± 0.04 ppm. Apparently, the refinement process did not cause
3 significant changes in the C β chemical shift deviations for group B or C within the uncertainty. In
4 contrast, average deviations from measured CO chemical shifts improved from an initial value of 0.64
5 ± 0.03 ppm in group 0 to 0.55 ± 0.03 ppm for group B ensembles, and to 0.50 ± 0.06 ppm for group C
6 ensembles. The ensembles in group A are similarly accurate in predicting CO chemical shifts with an
7 average deviation of 0.53 ± 0.06 ppm.

8 In summary, our structural ensembles reproduced both the measured SAXS curves and NMR chemical
9 shifts for all model IDPs with deviations from the measurements close to the experimental
10 uncertainty. The average agreement with SAXS curves and NMR chemical shifts indicates that there
11 are only minor differences between the quality of published PED models in group A and the newly
12 refined ensemble models of groups B and C. The ensembles of groups B and C also showed no
13 significant accuracy difference regarding the predicted SAXS curves and NMR chemical shifts,
14 suggesting that they are of similar quality. Based on the presented quality assessment, we consider
15 the model ensembles sufficiently accurate that they can now be used to assess the accuracy of both
16 structure-based CD prediction methods as well as CD-based SS estimation methods regarding IDPs.

17 Testing CD prediction methods

18 Utilizing the new IDP8 RDS, we proceed to determine the accuracy of the three structure-based CD-
19 spectrum prediction methods SESCA¹¹, PDB2CD⁸, and DichroCalc⁵, and compare their mean accuracy
20 separately for IDPs and globular proteins. The predicted CD spectra of all methods for the IDP8 RDS
21 are compared with the measured CD spectra in Figures S1-S6. The accuracy of these algorithms on
22 globular proteins was previously assessed¹¹ using the SP175 RDS, which contains 71 water soluble
23 globular proteins. The same SP175 data set was used as a training set for the two empirical methods
24 SESCA and PDB2CD, with no IDPs involved. The individual RMSDs computed between the measured
25 CD spectra of IDP8 RDS and the CD spectra predicted from the 14 ensemble models of the RDS are
26 shown in Table 4.

27 Figure 3 shows the average RMSD values between measured and predicted CD spectra (RMSD^{CD}, see
28 eq 3.) for both disordered (IDP8, blue) and globular (SP175, orange) proteins. For SESCA predictions,
29 four different basis sets were used: DS-dTSC3, DSSP-1SC3, HBSS-3SC1, DS5-SC1. These basis sets
30 represent 'pure' CD spectra for given SS elements (α -helix, β -sheet etc., see Nagy *et al.*¹¹ for precise
31 definitions), and therefore differ depending on which and how many SS elements have been used, as
32 well as on which SS classification method (e.g., DISICL⁷², DSSP⁷³, or HBSS¹¹) has been applied. All four
33 chosen basis sets contain correction terms for side chain signals for improved accuracy.

34 The average prediction accuracy of SESCA is 2.0 ± 0.1 kMRE for disordered proteins. As shown in Fig.
35 3, the average accuracy is similar for all four chosen basis sets, ranging between 1.9 and 2.2 kMRE
36 with a mean standard deviation (SD) of 1.0 kMRE for RMSD^{CD} values within the IDP8 RDS using the
37 same basis set. The average scatter of RMSD^{CD} values is 0.67 kMRE, when the measured CD spectra
38 are compared to CD predictions from the same ensemble model using different basis sets. In
39 comparison, the average prediction accuracy of SESCA for globular proteins is 2.1 ± 0.05 kMRE units
40 (as determined from the SP175 RDS). The scatter of RMSD^{CD} values for globular proteins is 1.0 kMRE
41 within the RDS using the same basis set and 0.7 kMRE between predictions from the same crystal
42 structure using different basis sets. The obtained RMSD values do not show a significant difference in
43 prediction accuracy between the chosen basis sets. Most importantly, the RMSD^{CD} values support our
44 previous expectations that, by construction, SESCA should yield a similar accuracy for disordered
45 proteins as for globular proteins.

1 Next, we tested the accuracy of the PDB2CD⁸ algorithm and its recent update PDBMD2CD⁷⁴ that allows
2 CD predictions from small structural ensembles. PDB2CD is based on determining the SS composition
3 from the model structure (or ensemble) by the DSSP algorithm and produces predicted spectra by
4 taking a weighted sum of spectra from structurally similar reference proteins. At the time of writing,
5 PDB2CD can utilize two globular RDS: SP175 and SMP180 to predict the CD spectra of protein models.
6 SMP180 includes all SP175 proteins and 11 additional membrane proteins, but neither RDS includes
7 any disordered proteins, which suggests limited accuracy for this class of proteins. PDBMD2CD is
8 based solely on the SMP180. Therefore, we used this RDS for computing CD predictions of both
9 globular and disordered protein spectra in our evaluation (the average accuracy for globular proteins
10 was still determined from the RMSD^{CD} values of SP175 proteins). As can be seen in Fig. 3, the accuracy
11 of PDBMD2CD for globular proteins is slightly better than that of SESCA, with an RMSD^{CD} of 1.6 ± 0.1
12 kMRE (SD 1.0 kMRE). For disordered proteins, however, the prediction accuracy of PDB2CD is
13 markedly reduced, with an average RMSD^{CD} 5.2 ± 0.5 kMRE (SD 1.7 kMRE).

14 In contrast to the other two empirical algorithms, DichroCalc predictions are calculated directly from
15 the three-dimensional protein structure through parameters derived from time-dependent quantum
16 mechanics (QM) calculations.⁵ The obtained average prediction RMSD^{CD} values for DichroCalc are
17 4.8 ± 0.3 kMRE (SD 2.4 kMRE) for globular proteins, and are even larger (8.7 ± 1.0 kMRE, SD 3.4 kMRE)
18 for disordered proteins. The obtained deviations from measured CD spectra indicate that the
19 approximations that allow CD calculations for entire proteins are rather harsh and limit the accuracy
20 of DichroCalc in reproducing the fine spectral features. These limitations are particularly severe for
21 disordered proteins, because the negative peak that defines the shape of their spectra is not
22 reproduced well by the underlying matrix method⁷⁵.

23 Further, to assess the effect of using CD information during ensemble refinement we also compared
24 the average accuracy of CD predictions of group B ensembles with those of group C ensembles shown
25 in Table 4. Here, we will focus on the prediction accuracies of SESCA, because the large mean and
26 scatter of RMSD^{CD} values for PDBMD2CD and DichroCalc renders it difficult to infer statistically
27 relevant statements about model quality using these methods. During the refinement of group B
28 ensembles, the SESCA basis set DS-dTSC3 was used to compute the CD signal of individual
29 conformations for *mevn-B* and *p53t-B*, whereas the DS5-4SC1 basis set was used for *actr-B*, *cbpn-B*,
30 and *rsp8-B*. The individual RMSD^{CD} values (underlined in Table 4) for CD predictions using these
31 ensembles and the corresponding basis sets average to 1.1 ± 0.2 kMRE, which can be considered the
32 best accuracy achievable by directly fitting the ensemble to match the measured CD spectrum. It is
33 also a considerable improvement over the 2.6 ± 0.3 kMRE average CD deviation of the initial MD
34 ensembles (see Table 3). The average deviation of group B ensemble CD predictions using all four
35 chosen SESCA basis sets (lines 5-9 in Table 4) amounted to 1.8 ± 0.2 kMRE. In comparison, the average
36 CD deviation for group C (lines 10-14) is 2.4 ± 0.4 kMRE, which suggests that including CD data in the
37 ensemble refinement process reduces both the mean and the scatter of RMSD^{CD} values to a small but
38 statistically significant extent.

39 In summary, based on the CD predictions for our IDP8 RDS, SESCA consistently predicts the CD spectra
40 of IDPs with an accuracy similar to that of globular proteins. Additionally, SESCA predictions are robust
41 with respect to the choice of basis set both for folded proteins and IDPs. In contrast, PDBMD2CD and
42 DichroCalc predictions are markedly less accurate regarding IDPs than for the folded proteins. Based
43 on our model quality assessments, including CD information during the ensemble refinement process
44 significantly improves CD predictions from the ensemble models, while maintaining the accuracy of
45 predicted SAXS curves and carbon chemical shifts.

1 Testing IDP SS estimation methods

2 Next, we focused on SS estimation, the second main branch of CD-based methods, which infers the
3 average SS composition of proteins from their measured CD spectra. Here, we assessed the SS
4 estimation accuracy of the Bayesian SS estimator *SESCA_bayes*⁴², using the same four basis sets as
5 above, as well as two other widely used methods, namely *BeStSel*⁷ and *K2D3*⁶. The estimated SS
6 fractions of all methods for the IDP8 RDS are shown in Tables S2-S7. To assess the accuracy of
7 estimated SS compositions, we compared them to reference SS compositions (see Methods Section
8 for details). For globular proteins, SS compositions of the NMR/crystallographic structures of the
9 SP175 RDS were used as reference. For disordered proteins, we selected the SS composition of those
10 ensemble models from IDP8 as reference that had the lowest average RMSD^{CD} for *SESCA* predictions,
11 namely *asyn-A*, *mevn-B*, *sic1-A*, *tk18-A*, *actr-B*, *cbpn-B*, *p53t-B*, and *rsp8-B*. The accuracy of the
12 estimated SS content was quantified by the RMSD to the reference SS fractions (RMSD^{SS} , see eq. 4).
13 The summary of all RMSD^{SS} values shown in Table S8 indicates, that the choice of reference ensemble
14 (except for *mevn*) doesn't have a large impact on the average accuracy of SS estimation methods and
15 wouldn't change our conclusions outlined below. For *mevn* all three tested methods estimated SS
16 fractions in better agreement with the *mevn-B* ensemble than *mevn-A* or *mevn-C*.

17 Figure 4 compares the average SS estimation accuracies of these methods for the IDP8 RDS (in blue)
18 of disordered proteins with those obtained for globular proteins of the SP175 RDS (orange). Overall,
19 the tested methods performed more similarly to one another than the CD prediction methods, albeit
20 larger differences are seen between the four *SESCA* basis set variants. All methods achieved average
21 RMSD^{SS} values between 0.07 and 0.12 for globular proteins and slightly larger average RMSD^{SS} values
22 (between 0.07 and 0.14) for disordered proteins. No clear correlation is observed between the SS
23 estimation accuracy and the number of SS classes used for the estimation method, although the
24 precision of *SESCA_bayes* estimates increased monotonically with the number of SS classes in the basis
25 set.

26 For the four *SESCA_bayes* variants using different basis sets, the smallest average RMSD^{SS} is obtained
27 for the DS5-4SC1 basis set (6 SS classes), with 0.07 RMSD^{SS} for both globular and disordered proteins
28 (SD of 0.04 and 0.06, respectively). The largest average RMSD^{SS} for *SESCA_bayes* are seen for the basis
29 set DSSP-1SC3 (four classes), amounting to an average RMSD^{SS} of 0.12 (SD 0.06) and 0.14 (SD 0.04) for
30 globular and disordered RDSs, respectively.

31 The program *K2D3* estimates a three-class SS composition using a neural network that was trained on
32 *DichroCalc* predictions of globular CD spectra based on their structures. *K2D3* estimates globular
33 protein SS fractions with an average RMSD^{SS} of 0.09 (SD 0.05), similar to the RMSD^{SS} *SESCA_bayes*
34 achieved using the DS-dTSC3 basis set with a similar 3-class SS composition. The RMSD^{SS} of *K2D3* for
35 IDPs is 0.12 (SD 0.05), somewhat larger than that for the globular RDS. We note that the obtained SS
36 estimation errors of *K2D3* are typically small for IDPs, despite the fact that the program provides very
37 poor back-calculated CD spectra and warns the user about the potential unreliability of those SS
38 estimates.

39 The *BeStSel* web application provides a detailed SS estimation based on eight SS classes, four of which
40 are associated with different types of β -sheets⁷. An average RMSD^{SS} of 0.08 (SD 0.03) is obtained for
41 globular proteins, and 0.14 (SD 0.05) for IDPs, which is the largest difference amongst the tested SS
42 estimators. We attribute this difference mainly to an observed systematic overestimation of the right-
43 handed antiparallel β -sheet (Anti3) fractions in our model IDPs (Table S7). Indeed, for the globular
44 RDS, the SS fractions are fairly similar for *BeStSel* estimates and the fractions of the reference (crystal)
45 structures. In contrast, almost none of the IDP ensemble models contains residues classified as the

1 Anti3 class, for which BeStSel estimates fractions between 0.2 and 0.3. The only protein in the IDP8
2 RDS for which the Anti3 fraction was not over-estimated was *cbpn*. However, *cbpn* is a molten-globule
3 type IDP with a stable α -helical structure, and thus its CD spectrum is more similar to those of helical
4 proteins.

5 It is worth noting that BeStSel also provides a simple ordered/disordered classification of proteins
6 based on their CD spectrum¹⁷, which our IDP8 RDS also enabled us to assess. Indeed, seven of the
7 eight proteins are correctly classified as disordered, with *cbpn* being classified as ordered. The latter
8 result is not a true misclassification because *cbpn* is a helical molten globule and its disorder is
9 apparent mostly on the tertiary structure level.

10 In contrast to the other available SS estimators, the Bayesian SS estimation method of SESCO
11 additionally provides uncertainties for the estimated SS fractions. To test if these Bayesian
12 uncertainties are realistic, we expressed the observed deviations to the reference SS fractions in units
13 of χ^2 analogously to eq. 2, but without a scaling factor. Similar to the RMSD^{SS} values above, the
14 computed χ^2 deviations also vary with the choice of the basis set. For the four basis sets, SESCO_bayes
15 achieves average χ^2 deviations for the IDP8 set of 0.87 (HBSS-3SC1), 1.03 (DS5-4SC1), 2.15 (DS-dTSC3),
16 and 2.51 (DSSP-1SC3). Obviously, these deviations are largely within one or two Bayes standard
17 deviations, such that the estimated uncertainty can be considered rather accurate. In contrast, the
18 average χ^2 values for the globular SP175 RDS are 1.32 (DS5-4SC3), 2.62 (HBSS-3SC1), 3.04 (DSSP-
19 1SC3), and 5.59 (DS-dTSC3), significantly larger than for the IDP set. As the RMSD^{SS} values for the SP175
20 are not considerably larger than those of our IDP8 RDS, the significantly larger χ^2 deviations indicate
21 that uncertainties of the SS fractions are underestimated for the DS-DTSC3 basis set, and to a lesser
22 extent for the DSSP-1SC3 basis set as well.

23 Overall, the observed RMSD^{SS} values indicate that SESCO basis sets estimate the SS composition of
24 IDPs with a similar accuracy as globular ones, whereas the average deviation of K2D3 and BeStSel SS
25 estimates are somewhat smaller for globular proteins and larger for IDPs. Our results also suggest that
26 SESCO basis sets DS5-4SC1 and Ds-dTSC3 are slightly more accurate for SS estimations than HBSS-3SC1
27 and DSSP-1SC3, but the uncertainties of DS-dTSC3 may be underestimated.

28 **Conclusions**

29 **Current method accuracy**

30 We introduced a new reference data set (RDS) for disordered proteins comprising CD spectra of eight
31 proteins and 14 ensemble models. This RDS, referred to as IDP8, served here to assess existing CD-
32 based biophysical analysis methods and can also support their further development. We first
33 determined the accuracy of the CD prediction methods SESCO, DichroCalc, and PDB2CD and compared
34 it to their accuracy for folded globular proteins using the curated RDS SP175. Overall, the accuracy of
35 these methods was lower (between 2.0 and 9.0 kmRE) for IDPs than for globular proteins (between
36 1.6 to 4.8 kmRE). SESCO predicted the CD spectra of globular and disordered proteins with a similar
37 high accuracy; PDB2CD performed well on globular proteins but was less accurate for IDPs, whereas
38 larger errors were seen for DichroCalc for both folded as well as disordered proteins.

39 Second, we used the IDP8 data set to assess the accuracy of the CD-based secondary structure
40 estimators SESCO_bayes, K2D3, and BESTSEL. Here, the (absolute) error of SS fraction estimates was
41 found between 0.07 and 0.14 for disordered proteins and between 0.07 and 0.12 for globular proteins.
42 Again, the accuracy of SESCO SS estimates was similar for folded and disordered proteins. However,
43 and in contrast to the above-mentioned CD spectrum predictions, it varied depending on the used

1 basis set. Both K2D3 and BESTSEL provided more accurate SS estimates for globular than for
2 disordered proteins.

3 Importantly, the IDP8 data set also enabled us to test if SESCO_bayes provides realistic uncertainty
4 estimates. For the disordered proteins, the uncertainty estimates largely agreed with the actual
5 deviations from the SS of the reference ensembles, whereas for the folded proteins, the uncertainty
6 estimates, particularly for the smaller basis sets, tended to be smaller than the actual errors. None of
7 the other SS estimators provides uncertainty estimates.

8 Over the past years, several methods for the structural characterization of folded proteins by CD
9 spectroscopy – such as CD spectrum predictors or SS estimators – have been established and are now
10 widely used. Their development and optimization has been enabled and driven by high quality RDSs
11 such as SP175²². Similar developments for IDPs, though pressing, have been hampered by the lack of
12 a suitable reference data set. We addressed this obstacle by compiling IDP8, an intrinsically disordered
13 protein RDS. Our subsequent assessments showed that the structural ensembles of IDP8 agree well
14 with SAXS and NMR chemical shift measurements, thus establishing that their quality is sufficient for
15 CD assessment. Using this new RDS, our assessments showed that SESCO CD predictions and SS
16 estimations achieved similarly high accuracy for disordered proteins as we previously determined for
17 globular proteins, which suggests that SESCO should be equally applicable to both protein classes.
18 Further, the assessment of several other CD prediction and SS estimation methods revealed generally
19 lower accuracy for IDPs than for globular proteins. Further, our data indicated that most of the tested
20 methods (including SESCO) would likely benefit from re-parametrization using the IDP8 RDS. We
21 therefore believe that our IDP8 RDS will also drive further methodological improvements in this rapidly
22 growing field.

23 Data set availability

24 All ensemble models and CD spectra will be made publicly available through the protein ensemble
25 database (PED) and the protein circular dichroism database (PCDDDB), respectively. Until then, the CD
26 spectra and ensemble models of the IDP8 RDS are available on request. Supplementary information
27 about computational tool availability, predicted CD spectra and estimated SS fractions are available
28 online free of charge.

29 Acknowledgements

30 Gabor Nagy would like to thank the Alexander von Humboldt foundation for financial support. The
31 authors would like to thank Sonia Longhi for providing support and experimental data on *mevn*. We
32 are thankful to Martha Brennich for the aid in SAXS measurements. We would like to thank Stefan
33 Becker, Christian Griesinger, and Karin Müller for their help in sample preparation. We would like to
34 thank Sarah Rauscher and Reinhardt Klement for providing simulation trajectories, experimental data,
35 and useful discussions regarding *rsp8* and *asyn*, respectively, and Vytautas Gapsys and Tamás Lázár
36 for helpful discussions regarding ensemble refinement and model deposition.

1

2 **References**

- 3 (1) Manavalan, P.; Johnson, W. C. Protein Secondary Structure from Circular Dichroism Spectra. *J. Biosci.* **1985**, *8* (1–2), 141–149.
- 4 (2) Whitmore, L.; Wallace, B. A. Protein Secondary Structure Analyses from Circular Dichroism
5 Spectroscopy: Methods and Reference Databases. *Biopolymers* **2008**, *89* (5), 392–400.
6 <https://doi.org/10.1002/bip.20853>.
- 7 (3) Wallace, B. A. Protein Characterisation by Synchrotron Radiation Circular Dichroism Spectroscopy.
8 *Q. Rev. Biophys.* **2009**, *42* (4), 317–370. <https://doi.org/10.1017/S003358351000003X>.
- 9 (4) *Circular Dichroism and the Conformational Analysis of Biomolecules*; Fasman, G. D., Ed.; Springer
10 US: Boston, MA, 1996.
- 11 (5) Bulheller, B. M.; Hirst, J. D. DichroCalc--Circular and Linear Dichroism Online. *Bioinformatics* **2009**,
12 *25* (4), 539–540. <https://doi.org/10.1093/bioinformatics/btp016>.
- 13 (6) Louis-Jeune, C.; Andrade-Navarro, M. A.; Perez-Iratxeta, C. Prediction of Protein Secondary
14 Structure from Circular Dichroism Using Theoretically Derived Spectra. *Proteins Struct. Funct.*
15 *Bioinforma.* **2012**, *80* (2), 374–381. <https://doi.org/10.1002/prot.23188>.
- 16 (7) Micsonai, A.; Wien, F.; Kernya, L.; Lee, Y.-H.; Goto, Y.; Réfrégiers, M.; Kardos, J. Accurate Secondary
17 Structure Prediction and Fold Recognition for Circular Dichroism Spectroscopy. *Proc. Natl. Acad.*
18 *Sci.* **2015**, *112* (24), E3095–E3103. <https://doi.org/10.1073/pnas.1500851112>.
- 19 (8) Mavridis, L.; Janes, R. W. PDB2CD: A Web-Based Application for the Generation of Circular
20 Dichroism Spectra from Protein Atomic Coordinates. *Bioinformatics* **2017**, *33* (1), 56–63.
21 <https://doi.org/10.1093/bioinformatics/btw554>.
- 22 (9) Sreerama, N.; Yu, S.; Woody, R. W. Estimation of Protein Secondary Structure from Circular
23 Dichroism Spectra: Inclusion of Denatured Proteins with Native Proteins in the Analysis. *Anal.*
24 *Biochem.* **2000**, *287*, 243–251. <https://doi.org/10.1006/abio.2000.4879>.
- 25 (10) Perczel, Andras; Hollósi, M.; Tudnady, G.; Fasman, G. D. Convex Constraint Analysis: A Natural
26 Deconvolution of Circular Dichroism Curves of Proteins. *Protein Eng.* **1991**, *4* (6), 669–679.
- 27 (11) Nagy, G.; Igaev, M.; Jones, N. C.; Hoffmann, S. V.; Grubmüller, H. SESCO : Predicting Circular
28 Dichroism Spectra from Protein Molecular Structures. *J. Chem. Theory Comput.* **2019**.
29 <https://doi.org/10.1021/acs.jctc.9b00203>.
- 30 (12) Whitmore, L.; Woollett, B.; Miles, A. J.; Klose, D. P.; Janes, R. W.; Wallace, B. A. PCDDDB: The Protein
31 Circular Dichroism Data Bank, a Repository for Circular Dichroism Spectral and Metadata. *Nucleic*
32 *Acids Res.* **2011**, *39* (Database), D480–D486. <https://doi.org/10.1093/nar/gkq1026>.
- 33 (13) Quaglia, F.; Mészáros, B.; Salladini, E.; Hatos, A.; Pancsa, R.; Chemes, L. B.; Pajkos, M.; Lazar, T.;
34 Peña-Díaz, S.; Santos, J.; Ács, V.; Farahi, N.; Fichó, E.; Aspromonte, M. C.; Bassot, C.; Chasapi, A.;
35 Davey, N. E.; Davidović, R.; Dobson, L.; Elofsson, A.; Erdős, G.; Gaudet, P.; Giglio, M.; Glavina, J.;
36 Iserte, J.; Iglesias, V.; Kálmán, Z.; Lamborghini, M.; Leonardi, E.; Longhi, S.; Macedo-Ribeiro, S.;
37 Maiani, E.; Marchetti, J.; Marino-Buslje, C.; Mészáros, A.; Monzon, A. M.; Minervini, G.; Nadendla,
38 S.; Nilsson, J. F.; Novotný, M.; Ouzounis, C. A.; Palopoli, N.; Papaleo, E.; Pereira, P. J. B.; Pozzati,
39 G.; Promponas, V. J.; Pujols, J.; Rocha, A. C. S.; Salas, M.; Sawicki, L. R.; Schad, E.; Shenoy, A.;
40 Szaniszló, T.; Tsigos, K. D.; Veljkovic, N.; Parisi, G.; Ventura, S.; Dosztányi, Z.; Tompa, P.; Tosatto,
41 S. C. E.; Piovesan, D. DisProt in 2022: Improved Quality and Accessibility of Protein Intrinsic
42 Disorder Annotation. *Nucleic Acids Res.* **2022**, *50* (D1), D480–D487.
43 <https://doi.org/10.1093/nar/gkab1082>.
- 44 (14) van der Lee, R.; Buljan, M.; Lang, B.; Weatheritt, R. J.; Daughdrill, G. W.; Dunker, A. K.; Fuxreiter,
45 M.; Gough, J.; Gsponer, J.; Jones, D. T.; Kim, P. M.; Kriwacki, R. W.; Oldfield, C. J.; Pappu, R. V.;
46 Tompa, P.; Uversky, V. N.; Wright, P. E.; Babu, M. M. Classification of Intrinsically Disordered
47 Regions and Proteins. *Chem. Rev.* **2014**, *114* (13), 6589–6631.
48 <https://doi.org/10.1021/cr400525m>.
- 49

- 1 (15) Wright, P. E.; Dyson, H. J. Intrinsically Disordered Proteins in Cellular Signalling and Regulation.
2 *Nat. Rev. Mol. Cell Biol.* **2014**, *16* (1), 18–29. <https://doi.org/10.1038/nrm3920>.
- 3 (16) Dass, R.; Mulder, F. A. A.; Nielsen, J. T. ODINPred: Comprehensive Prediction of Protein Order and
4 Disorder. *Sci. Rep.* **2020**, *10* (1), 14780. <https://doi.org/10.1038/s41598-020-71716-1>.
- 5 (17) Micsonai, A.; Moussong, É.; Wien, F.; Boros, E.; Vadász, H.; Murvai, N.; Lee, Y.-H.; Molnár, T.;
6 Réfrégiers, M.; Goto, Y.; Tantos, Á.; Kardos, J. BeStSel: Webserver for Secondary Structure and
7 Fold Prediction for Protein CD Spectroscopy. *Nucleic Acids Res.* **2022**, *50* (W1), W90–W98.
8 <https://doi.org/10.1093/nar/gkac345>.
- 9 (18) Boehm, G.; Muhr, R.; Jaenicke, R. Quantitative Analysis of Protein Far UV Circular Dichroism
10 Spectra by Neural Networks. *Protein Eng.* **1992**, *5* (3), 191–195.
- 11 (19) Johnson Jr., W. C. Protein Secondary Structure and Circular Dichroism: A Practical Guide.
12 *PROTEINS Struct. Funct. Genet.* **1990**, *7*, 205–214.
- 13 (20) Oakley, M. T.; Bulheller, B. M.; Hirst, J. D. First-Principles Calculations of Protein Circular Dichroism
14 in the Far-Ultraviolet and Beyond. *Chirality* **2006**, *18*, 340–347.
15 <https://doi.org/10.1002/chir.20264>.
- 16 (21) Štěpánek, P.; Bouř, P. Multi-Scale Modeling of Electronic Spectra of Three Aromatic Amino Acids:
17 Importance of Conformational Averaging and Explicit Solute–Solvent Interactions. *Phys Chem*
18 *Chem Phys* **2014**, *16* (38), 20639–20649. <https://doi.org/10.1039/C4CP02668C>.
- 19 (22) Lees, J. G.; Miles, A. J.; Wien, F.; Wallace, B. A. A Reference Database for Circular Dichroism
20 Spectroscopy Covering Fold and Secondary Structure Space. *Bioinformatics* **2006**, *22* (16), 1955–
21 1962. <https://doi.org/10.1093/bioinformatics/btl327>.
- 22 (23) Abdul-Gader, A.; Miles, A. J.; Wallace, B. A. A Reference Dataset for the Analyses of Membrane
23 Protein Secondary Structures and Transmembrane Residues using Circular Dichroism
24 Spectroscopy. *Bioinformatics* **2011**, *27* (12), 1630–1636.
25 <https://doi.org/10.1093/bioinformatics/btr234>.
- 26 (24) Evans, P.; Bateman, O. A.; Slingsby, C.; Wallace, B. A. A Reference Dataset for Circular Dichroism
27 Spectroscopy Tailored for the B γ -Crystallin Lens Proteins. *Exp. Eye Res.* **2007**, *84* (5), 1001–1008.
28 <https://doi.org/10.1016/j.exer.2007.01.016>.
- 29 (25) Barghorn, S.; Davies, P.; Mandelkow, E. Tau Paired Helical Filaments from Alzheimer’s Disease
30 Brain and Assembled in Vitro Are Based on β -Structure in the Core Domain. *Biochemistry* **2004**,
31 *43* (6), 1694–1703. <https://doi.org/10.1021/bi0357006>.
- 32 (26) Dogan, J.; Schmidt, T.; Mu, X.; Engström, Å.; Jemth, P. Fast Association and Slow Transitions in the
33 Interaction between Two Intrinsically Disordered Protein Domains. *J. Biol. Chem.* **2012**, *287* (41),
34 34316–34324. <https://doi.org/10.1074/jbc.M112.399436>.
- 35 (27) Lin, C. H.; Hare, B. J.; Wagner, G.; Harrison, S. C.; Maniatis, T.; Fraenkel, E. A Small Domain of
36 CBP/P300 Binds Diverse Proteins: Solution Structure and Functional Studies. *Mol. Cell* **2001**, *8* (3),
37 581–590.
- 38 (28) Wong, T. S.; Rajagopalan, S.; Freund, S. M.; Rutherford, T. J.; Andreeva, A.; Townsley, F. M.;
39 Petrovich, M.; Fersht, A. R. Biophysical Characterizations of Human Mitochondrial Transcription
40 Factor A and Its Binding to Tumor Suppressor P53. *Nucleic Acids Res.* **2009**, *37* (20), 6765–6783.
41 <https://doi.org/10.1093/nar/gkp750>.
- 42 (29) Shu, Y.; Habchi, J.; Costanzo, S.; Padilla, A.; Brunel, J.; Gerlier, D.; Oglesbee, M.; Longhi, S. Plasticity
43 in Structural and Functional Interactions between the Phosphoprotein and Nucleoprotein of
44 Measles Virus. *J. Biol. Chem.* **2012**, *287* (15), 11951–11967.
45 <https://doi.org/10.1074/jbc.M111.333088>.
- 46 (30) Fisher, C. K.; Stultz, C. M. Constructing Ensembles for Intrinsically Disordered Proteins. *Curr. Opin.*
47 *Struct. Biol.* **2011**, *21* (3), 426–431. <https://doi.org/10.1016/j.sbi.2011.04.001>.
- 48 (31) Liu, H.; Song, D.; Zhang, Y.; Yang, S.; Luo, R.; Chen, H.-F. Extensive Tests and Evaluation of the
49 CHARMM36IDPSFF Force Field for Intrinsically Disordered Proteins and Folded Proteins. *Phys.*
50 *Chem. Chem. Phys.* **2019**, *21* (39), 21918–21931. <https://doi.org/10.1039/C9CP03434J>.

- 1 (32) Robustelli, P.; Piana, S.; Shaw, D. E. Developing a Molecular Dynamics Force Field for Both Folded
2 and Disordered Protein States. *Proc. Natl. Acad. Sci.* **2018**, *115* (21).
3 <https://doi.org/10.1073/pnas.1800690115>.
- 4 (33) Hornak, V.; Abel, R.; Okur, A.; Strockbine, B.; Roitberg, A.; Simmerling, C. Comparison of Multiple
5 Amber Force Fields and Development of Improved Protein Backbone Parameters. *Proteins Struct.*
6 *Funct. Bioinforma.* **2006**, *65* (3), 712–725. <https://doi.org/10.1002/prot.21123>.
- 7 (34) Rauscher, S.; Gapsys, V.; Gajda, M. J.; Zweckstetter, M.; de Groot, B. L.; Grubmüller, H. Structural
8 Ensembles of Intrinsically Disordered Proteins Depend Strongly on Force Field: A Comparison to
9 Experiment. *J. Chem. Theory Comput.* **2015**, *11* (11), 5513–5524.
10 <https://doi.org/10.1021/acs.jctc.5b00736>.
- 11 (35) Ozenne, V.; Bauer, F.; Salmon, L.; Huang, J. -r.; Jensen, M. R.; Segard, S.; Bernado, P.; Charavay, C.;
12 Blackledge, M. Flexible-Meccano: A Tool for the Generation of Explicit Ensemble Descriptions of
13 Intrinsically Disordered Proteins and Their Associated Experimental Observables. *Bioinformatics*
14 **2012**, *28* (11), 1463–1470. <https://doi.org/10.1093/bioinformatics/bts172>.
- 15 (36) Teixeira, J. M. C.; Liu, Z. H.; Namini, A.; Li, J.; Vernon, R. M.; Krzeminski, M.; Shamandy, A. A.;
16 Zhang, O.; Haghghatlari, M.; Yu, L.; Head-Gordon, T.; Forman-Kay, J. D. IDPConformerGenerator:
17 A Flexible Software Suite for Sampling the Conformational Space of Disordered Protein States. *J.*
18 *Phys. Chem. A* **2022**, *126* (35), 5985–6003. <https://doi.org/10.1021/acs.jpca.2c03726>.
- 19 (37) Köfinger, J.; Stelzl, L. S.; Reuter, K.; Allande, C.; Reichel, K.; Hummer, G. Efficient Ensemble
20 Refinement by Reweighting. *J. Chem. Theory Comput.* **2019**, *15* (5), 3390–3401.
21 <https://doi.org/10.1021/acs.jctc.8b01231>.
- 22 (38) Salmon, L.; Nodet, G.; Ozenne, V.; Yin, G.; Jensen, M. R.; Zweckstetter, M.; Blackledge, M. NMR
23 Characterization of Long-Range Order in Intrinsically Disordered Proteins. *J. Am. Chem. Soc.* **2010**,
24 *132* (24), 8407–8418. <https://doi.org/10.1021/ja101645g>.
- 25 (39) Ezerski, J. C.; Zhang, P.; Jennings, N. C.; Waxham, M. N.; Cheung, M. S. Molecular Dynamics
26 Ensemble Refinement of Intrinsically Disordered Peptides According to Deconvoluted Spectra
27 from Circular Dichroism. *Biophys. J.* **2020**, *118* (7), 1665–1678.
28 <https://doi.org/10.1016/j.bpj.2020.02.015>.
- 29 (40) Varadi, M.; Kosol, S.; Lebrun, P.; Valentini, E.; Blackledge, M.; Dunker, A. K.; Felli, I. C.; Forman-
30 Kay, J. D.; Kriwacki, R. W.; Pierattelli, R.; Sussman, J.; Svergun, D. I.; Uversky, V. N.; Vendruscolo,
31 M.; Wishart, D.; Wright, P. E.; Tompa, P. pE-DB: A Database of Structural Ensembles of Intrinsically
32 Disordered and of Unfolded Proteins. *Nucleic Acids Res.* **2014**, *42* (D1), D326–D335.
33 <https://doi.org/10.1093/nar/gkt960>.
- 34 (41) Hoch, J. C.; Baskaran, K.; Burr, H.; Chin, J.; Eghbalnia, H. R.; Fujiwara, T.; Gryk, M. R.; Iwata, T.;
35 Kojima, C.; Kurisu, G.; Maziuk, D.; Miyanoiri, Y.; Wedell, J. R.; Wilburn, C.; Yao, H.; Yokochi, M.
36 Biological Magnetic Resonance Data Bank. *Nucleic Acids Res.* **2023**, *51* (D1), D368–D376.
37 <https://doi.org/10.1093/nar/gkac1050>.
- 38 (42) Nagy, G.; Grubmüller, H. Implementation of a Bayesian Secondary Structure Estimation Method
39 for the SESCO Circular Dichroism Analysis Package. *Comput. Phys. Commun.* **2021**, *266*, 108022.
40 <https://doi.org/10.1016/j.cpc.2021.108022>.
- 41 (43) Kuipers, B. J. H.; Gruppen, H. Prediction of Molar Extinction Coefficients of Proteins and Peptides
42 Using UV Absorption of the Constituent Amino Acids at 214 Nm To Enable Quantitative Reverse
43 Phase High-Performance Liquid Chromatography-Mass Spectrometry Analysis. *J. Agric. Food*
44 *Chem.* **2007**, *55* (14), 5445–5451. <https://doi.org/10.1021/jf070337l>.
- 45 (44) Troilo, F.; Bonetti, D.; Bignon, C.; Longhi, S.; Gianni, S. Understanding Intramolecular Crosstalk in
46 an Intrinsically Disordered Protein. *ACS Chem. Biol.* **2019**, *14* (3), 337–341.
47 <https://doi.org/10.1021/acscchembio.8b01055>.
- 48 (45) Barghorn, S.; Davies, P.; Mandelkow, E. Tau Paired Helical Filaments from Alzheimer’s Disease
49 Brain and Assembled in Vitro Are Based on β -Structure in the Core Domain. *Biochemistry* **2004**,
50 *43* (6), 1694–1703. <https://doi.org/10.1021/bi0357006>.

- 1 (46) Incardona, M.-F.; Bourenkov, G. P.; Levik, K.; Pieritz, R. A.; Popov, A. N.; Svensson, O. *EDNA : A*
2 *Framework for Plugin-Based Applications Applied to X-Ray Experiment Online Data Analysis. J.*
3 *Synchrotron Radiat.* **2009**, *16* (6), 872–879. <https://doi.org/10.1107/S0909049509036681>.
- 4 (47) Mylonas, E.; Hascher, A.; Bernadó, P.; Blackledge, M.; Mandelkow, E.; Svergun, D. I. Domain
5 Conformation of Tau Protein Studied by Solution Small-Angle X-Ray Scattering. *Biochemistry* **2008**,
6 *47* (39), 10345–10353. <https://doi.org/10.1021/bi800900d>.
- 7 (48) Mittag, T.; Marsh, J.; Grishaev, A.; Orlicky, S.; Lin, H.; Sicheri, F.; Tyers, M.; Forman-Kay, J. D.
8 Structure/Function Implications in a Dynamic Complex of the Intrinsically Disordered Sic1 with the
9 Cdc4 Subunit of an SCF Ubiquitin Ligase. *Structure* **2010**, *18* (4), 494–506.
10 <https://doi.org/10.1016/j.str.2010.01.020>.
- 11 (49) Gely, S.; Lowry, D. F.; Bernard, C.; Jensen, M. R.; Blackledge, M.; Costanzo, S.; Bourhis, J.-M.;
12 Darbon, H.; Daughdrill, G.; Longhi, S. Solution Structure of the C-Terminal X Domain of the Measles
13 Virus Phosphoprotein and Interaction with the Intrinsically Disordered C-Terminal Domain of the
14 Nucleoprotein. *J. Mol. Recognit.* **2010**, *23* (5), 435–447. <https://doi.org/10.1002/jmr.1010>.
- 15 (50) Abraham, M. J.; Murtola, T.; Schulz, R.; Páll, S.; Smith, J. C.; Hess, B.; Lindahl, E. GROMACS: High
16 Performance Molecular Simulations through Multi-Level Parallelism from Laptops to
17 Supercomputers. *SoftwareX* **2015**, *1–2*, 19–25. <https://doi.org/10.1016/j.softx.2015.06.001>.
- 18 (51) Bussi, G.; Donadio, D.; Parrinello, M. Canonical Sampling through Velocity Rescaling. *J. Chem. Phys.*
19 **2007**, *126* (1), 014101. <https://doi.org/10.1063/1.2408420>.
- 20 (52) Parrinello, M.; Rahman, A. Polymorphic Transitions in Single Crystals: A New Molecular Dynamics
21 Method. *J. Appl. Phys.* **1981**, *52* (12), 7182–7190. <https://doi.org/10.1063/1.328693>.
- 22 (53) Van Gunsteren, W. F.; Berendsen, H. J. C. A Leap-Frog Algorithm for Stochastic Dynamics. *Mol.*
23 *Simul.* **1988**, *1* (3), 173–185. <https://doi.org/10.1080/08927028808080941>.
- 24 (54) Hess, B. P-LINCS: A Parallel Linear Constraint Solver for Molecular Simulation. *J. Chem. Theory*
25 *Comput.* **2008**, *4* (1), 116–122. <https://doi.org/10.1021/ct700200b>.
- 26 (55) Darden, T.; York, D.; Pedersen, L. Particle Mesh Ewald: An $N \cdot \log(N)$ Method for Ewald Sums in
27 Large Systems. *J. Chem. Phys.* **1993**, *98* (12), 10089–10092. <https://doi.org/10.1063/1.464397>.
- 28 (56) Allen, M. P.; Tildesley, D. J. *Computer Simulation of Liquids*; Oxford University Press, USA, 1989.
- 29 (57) Abriata, L. A.; Dal Peraro, M. Assessment of Transferable Forcefields for Protein Simulations
30 Attests Improved Description of Disordered States and Secondary Structure Propensities, and
31 Hints at Multi-Protein Systems as the next Challenge for Optimization. *Comput. Struct. Biotechnol.*
32 *J.* **2021**, *19*, 2626–2636. <https://doi.org/10.1016/j.csbj.2021.04.050>.
- 33 (58) Mu, J.; Pan, Z.; Chen, H.-F. Balanced Solvent Model for Intrinsically Disordered and Ordered
34 Proteins. *J. Chem. Inf. Model.* **2021**, *61* (10), 5141–5151.
35 <https://doi.org/10.1021/acs.jcim.1c00407>.
- 36 (59) Duan, Y.; Wu, C.; Chowdhury, S.; Lee, M. C.; Xiong, G.; Zhang, W.; Yang, R.; Cieplak, P.; Luo, R.; Lee,
37 T.; Caldwell, J.; Wang, J.; Kollman, P. A Point-Charge Force Field for Molecular Mechanics
38 Simulations of Proteins Based on Condensed-Phase Quantum Mechanical Calculations. *J. Comput.*
39 *Chem.* **2003**, *24* (16), 1999–2012. <https://doi.org/10.1002/jcc.10349>.
- 40 (60) Best, R. B.; Zheng, W.; Mittal, J. Balanced Protein–Water Interactions Improve Properties of
41 Disordered Proteins and Non-Specific Protein Association. *J. Chem. Theory Comput.* **2014**, *10* (11),
42 5113–5124. <https://doi.org/10.1021/ct500569b>.
- 43 (61) Robustelli, P.; Piana, S.; Shaw, D. E. Developing a Molecular Dynamics Force Field for Both Folded
44 and Disordered Protein States. *Proc. Natl. Acad. Sci.* **2018**, *115* (21).
45 <https://doi.org/10.1073/pnas.1800690115>.
- 46 (62) Maier, J. A.; Martinez, C.; Kasavajhala, K.; Wickstrom, L.; Hauser, K. E.; Simmerling, C. ff14SB:
47 Improving the Accuracy of Protein Side Chain and Backbone Parameters from ff99SB. *J. Chem.*
48 *Theory Comput.* **2015**, *11* (8), 3696–3713. <https://doi.org/10.1021/acs.jctc.5b00255>.
- 49 (63) Izadi, S.; Anandakrishnan, R.; Onufriev, A. V. Building Water Models: A Different Approach. *J. Phys.*
50 *Chem. Lett.* **2014**, *5* (21), 3863–3871. <https://doi.org/10.1021/jz501780a>.

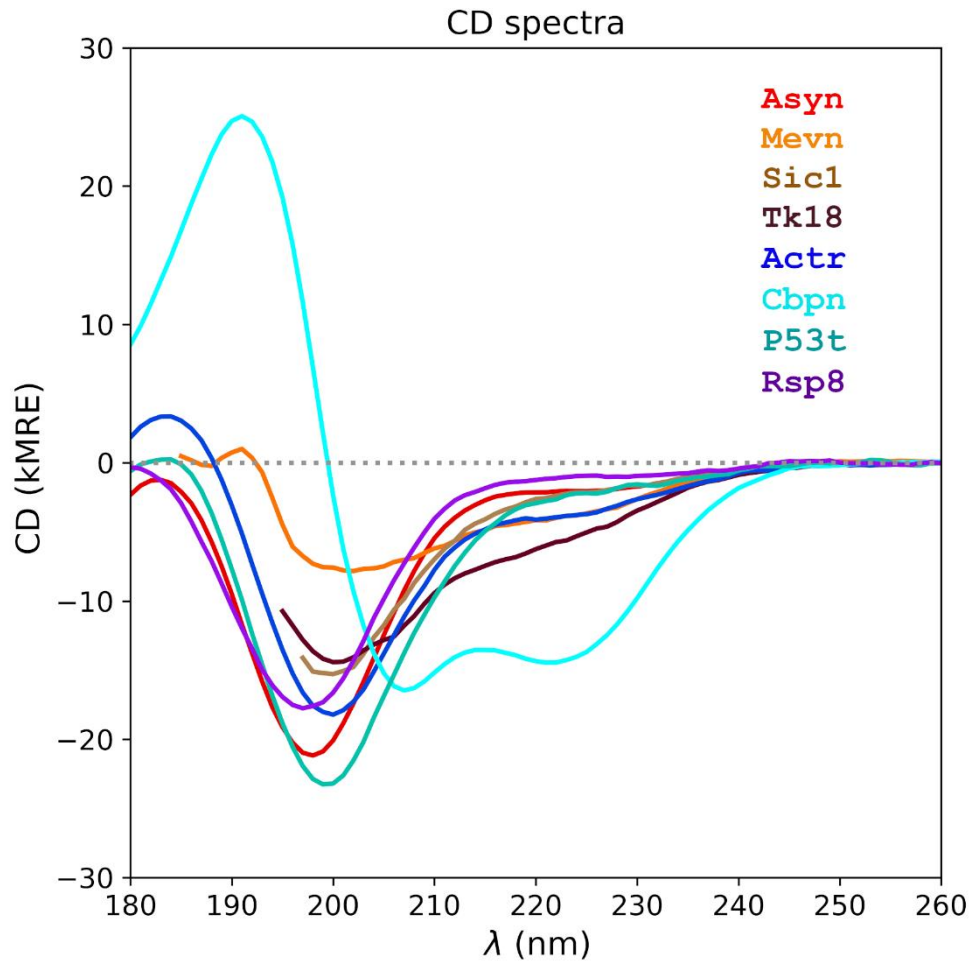
- 1 (64) MacKerell Jr, A. D.; Bashford, D.; Bellott, M.; Dunbrack Jr, R. L.; Evanseck, J. D.; Field, M. J.; Fischer,
2 S.; Gao, J.; Guo, H.; Ha, S. All-Atom Empirical Potential for Molecular Modeling and Dynamics
3 Studies of Proteins. *J. Phys. Chem. B* **1998**, *102* (18), 3586–3616.
- 4 (65) Bjelkmar, P.; Larsson, P.; Cuendet, M. A.; Hess, B.; Lindahl, E. Implementation of the CHARMM
5 Force Field in GROMACS: Analysis of Protein Stability Effects from Correction Maps, Virtual
6 Interaction Sites, and Water Models. *J. Chem. Theory Comput.* **2010**, *6* (2), 459–466.
7 <https://doi.org/10.1021/ct900549r>.
- 8 (66) Huang, J.; Rauscher, S.; Nawrocki, G.; Ran, T.; Feig, M.; de Groot, B. L.; Grubmüller, H.; MacKerell,
9 A. D. CHARMM36m: An Improved Force Field for Folded and Intrinsically Disordered Proteins. *Nat.*
10 *Methods* **2017**, *14* (1), 71–73. <https://doi.org/10.1038/nmeth.4067>.
- 11 (67) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.;
12 Bourne, P. E. The Protein Data Bank. *Nucleic Acids Res.* **2000**, *28* (1), 235–242.
13 <https://doi.org/10.1093/nar/28.1.235>.
- 14 (68) Bottaro, S.; Bengtson, T.; Lindorff-Larsen, K. Integrating Molecular Simulation and Experimental
15 Data: A Bayesian/Maximum Entropy Reweighting Approach. In *Structural Bioinformatics*; Gáspári,
16 Z., Ed.; Methods in Molecular Biology; Springer US: New York, NY, 2020; Vol. 2112, pp 219–240.
17 https://doi.org/10.1007/978-1-0716-0270-6_15.
- 18 (69) Shen, Y.; Bax, A. SPARTA+: A Modest Improvement in Empirical NMR Chemical Shift Prediction by
19 Means of an Artificial Neural Network. *J. Biomol. NMR* **2010**, *48* (1), 13–22.
20 <https://doi.org/10.1007/s10858-010-9433-9>.
- 21 (70) Svergun, D.; Barberato, C.; Koch, M. H. CRY SOL—a Program to Evaluate X-Ray Solution Scattering
22 of Biological Macromolecules from Atomic Coordinates. *J. Appl. Crystallogr.* **1995**, *28* (6), 768–
23 773.
- 24 (71) Nagy, G.; Grubmüller, H. How Accurate Is Circular Dichroism-Based Model Validation? *Eur.*
25 *Biophys. J.* **2020**, *49* (6), 497–510. <https://doi.org/10.1007/s00249-020-01457-6>.
- 26 (72) Nagy, G.; Oostenbrink, C. Dihedral-Based Segment Identification and Classification of Biopolymers
27 I: Proteins. *J. Chem. Inf. Model.* **2014**, *54* (1), 266–277. <https://doi.org/10.1021/ci400541d>.
- 28 (73) Kabsch, W.; Sander, C. Dictionary of Protein Secondary Structure: Pattern Recognition of
29 Hydrogen-Bonded and Geometrical Features. *Biopolymers* **1983**, *22* (12), 2577–2637.
- 30 (74) Drew, E. D.; Janes, R. W. PDBMD2CD: Providing Predicted Protein Circular Dichroism Spectra from
31 Multiple Molecular Dynamics-Generated Protein Structures. *Nucleic Acids Res.* **2020**, *48* (W1),
32 W17–W24. <https://doi.org/10.1093/nar/gkaa296>.
- 33 (75) Bulheller, B. M.; Rodger, A.; Hirst, J. D. Circular and Linear Dichroism of Proteins. *Phys. Chem.*
34 *Chem. Phys.* **2007**, *9* (17), 2020. <https://doi.org/10.1039/b615870f>.
- 35
36

1

2 **Figures**

3

4



5

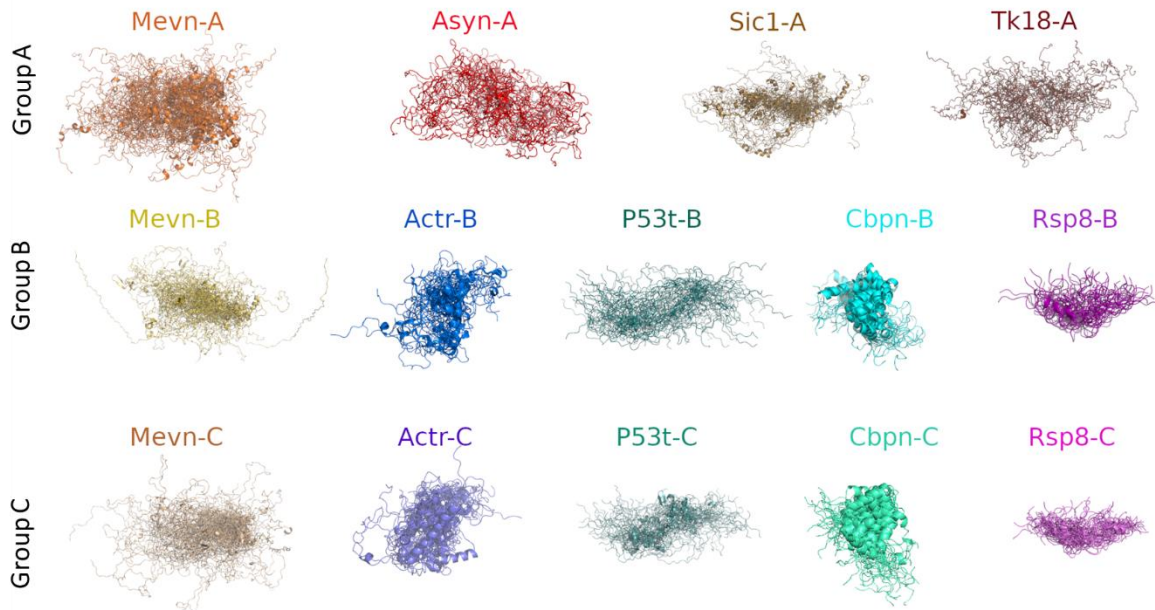
6 **Figure 1: Measured IDP8 CD spectra.** The spectra of eight different IDP domains are shown in different
7 colors. Abbreviations for the name of each domain are shown in the upper right corner (color coded)
8 and are listed in Table 1. The full name of each IDP domain is listed in the Reference data set assembly
9 section of this manuscript. Intensities of the CD spectra are expressed in 1000 mean residue ellipticity
10 units (kMRE or 1000 deg* cm² / dmol). The dotted gray line indicates the CD intensity of 0 kMRE.

11

12

1

2



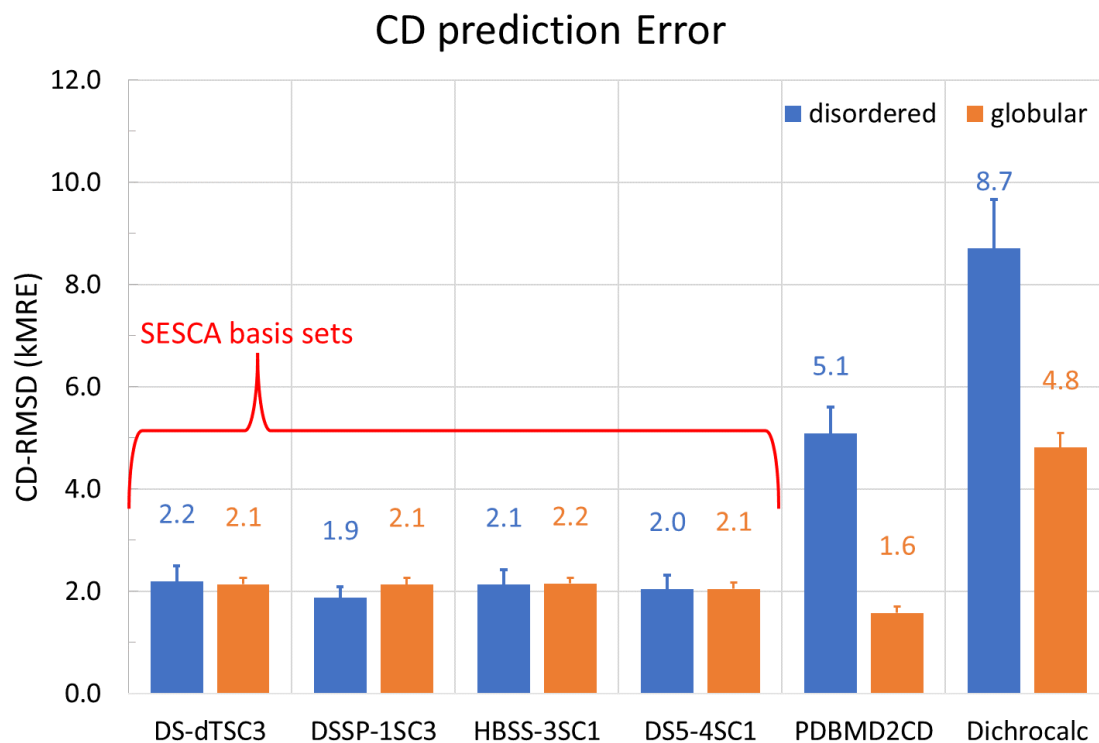
3

4 **Figure 2: IDP8 protein ensemble models.** Each ensemble model is an overlay of 20-50 backbone
5 conformations, shown in cartoon representation, and fitted to the first model of the respective
6 ensemble. The name of each ensemble model is displayed above the model. Group A models were
7 previously published and were obtained from the PED, group B models were derived by the authors
8 using NMR chemical shifts, SAXS, and CD measurements. Group C models were derived similarly as the
9 models of group B but without using CD information.

10

11

1
2

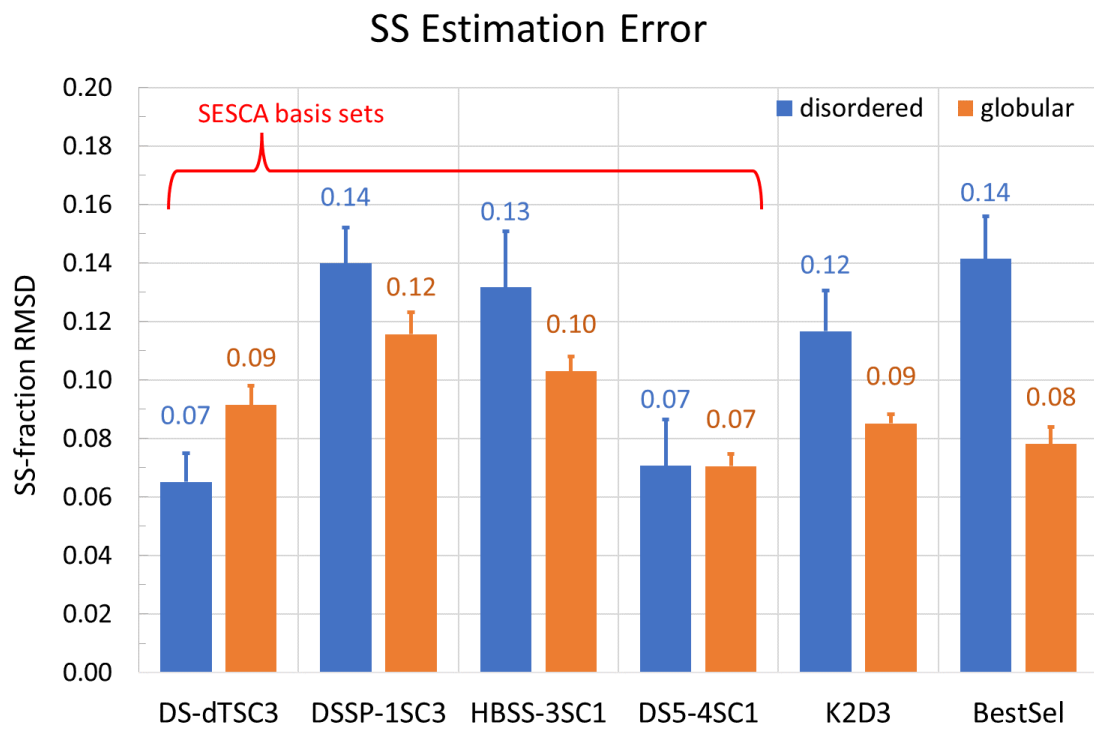


3
4
5
6
7
8
9
10
11

Figure 3 : Accuracy of CD spectrum predictions. Summary of RMSDs of CD spectra predicted from reference model structures relative to measured spectra of the same protein. Shown are RMSD values averaged over all proteins, for the different methods described in the text. Two reference data sets have been used: IDP8 for disordered proteins (blue) and SP175 for folded globular proteins (orange). Tested CD prediction methods are: DichroCalc, PDBMD2CD, and SESCA with four different basis sets (DS-dTSC3, DSSP-1SC3, HBSS-3SC1 and DS5-4SC1).

1

2



3

4 **Figure 4: Accuracy of SS fraction estimates.** Summary of averaged RMSDs of SS fractions estimated from
5 the reference CD spectra by different methods relative to SS fractions computed from the respective
6 reference structure. As in Fig. 3, two RDSs have been used: IDP8 for disordered proteins (blue), and
7 SP175 for folded globular proteins (orange). The tested SS fraction estimators are: K2D3, BESTSEL, and
8 SESCO_Bayes with four different basis sets (DS-dTSC3, DSSP-1SC3, HBSS-3SC1 and DS5-4SC1).

9

10

1 *Tables*

2

3 **Table 1: Measured IDP8 CD spectra.** Properties of the eight CD spectra included within the IDP8 RDS.
4 Columns list the ID and the short code of the protein, their minimum (λ_{\min}) and maximum (λ_{\max})
5 wavelengths (in nm) of the spectra, whether it was recorded on a conventional spectrophotometer
6 (CD) or a synchrotron radiation CD (SR-CD) facility, and the estimated protein concentration (C_{prot} , in
7 μM) of the measured sample.

spectrum-ID	short code	facility	λ_{\min}	λ_{\max}	C_{prot}
1	<i>asyn</i>	SR-CD	178	280	75
2	<i>mevn</i>	CD	185	260	24
3	<i>sic1</i>	CD	200	250	10
4	<i>tk18</i>	CD	195	260	120
5	<i>actr</i>	SR-CD	178	300	75
6	<i>cbpn</i>	SR-CD	178	280	120
7	<i>p53t</i>	SR-CD	178	260	60
8	<i>rsp8</i>	SR-CD	178	300	270

8

9

1

2 **Table 2: IDP8 structural ensembles.** Summary of the 14 model ensembles included within the IDP8 RDS.
3 The columns list the ID and the short code of the models, the PED accession code of the model, residue
4 numbers of the IDP domain in the full protein, the length of peptide sequence (in amino acids), the
5 number of conformations in the model ensemble (ens. size), and experimental data used to construct
6 or refine the model ensemble. Abbreviations of experimental data denote NMR chemical shifts (CS),
7 NMR paramagnetic relaxation enhancement (PRE), NMR residual dipolar coupling (RDC), small angle X-
8 ray scattering (SAXS), and Circular Dichroism (CD).

Group	model ID	short code	PED code	residues	length (aa)	ens. size	exp. Data
A	1	<i>asyn-A</i>	PED:00024-1	1-140	140	567	PRE, SAXS
	2	<i>mevn-A</i>	PED:00020	400-525	132	995	CS, RDC
	3	<i>sic1-A</i>	PED:00160-2	1-90	92	500	CS, RDC, PRE, SAXS
	4	<i>tk18-A</i>	PED:0192	1-130	130	75	CS, RDC, SAXS
B	5	<i>mevn-B</i>	PED:00233	400-525	132	100	CS, SAXS, CD
	6	<i>actr-B</i>	PED:00230	1018-1088	71	100	CS, SAXS, CD
	7	<i>cbpn-B</i>	PED:00228	2059-2117	59	100	CS, SAXS, CD
	8	<i>p53t-B</i>	PED:00229	1-73	73	250	CS, SAXS, CD
	9	<i>rsp8-B</i>	PED:00231	1-24	24	250	CS, SAXS, CD
C	10	<i>mevn-C</i>	PED:00234	400-525	132	100	CS, SAXS
	11	<i>actr-C</i>	PED:00237	1018-1088	71	100	CS, SAXS
	12	<i>cbpn-C</i>	PED:00235	2059-2117	59	100	CS, SAXS
	13	<i>p53t-C</i>	PED:00236	1-73	73	250	CS, SAXS
	14	<i>rsp8-C</i>	PED:00238	1-24	24	250	CS, SAXS

9

10

1

2

3 **Table 3: IDP8 ensemble model assessment.** Summary of IDP8 ensemble model prediction vs. measured
 4 SAXS curves, NMR chemical shifts, and CD spectra. The table lists the ensemble ID, the square root of
 5 the χ^2 deviation of the predicted and measured SAXS curves, the average RMSDs between backbone
 6 NMR chemical shifts (CS) for C α , C β , and carbonyl carbon (CO) atoms, as well as the RMSD of CD
 7 intensities (CD) as predicted using the SESCA basis set DS-dTSC3.

Group	model ID	SAXS χ	CS-C α ppm	CS-C β ppm	CS-CO ppm	CD kMRE
A	<i>asyn-A</i>	1.34	0.36	0.66	0.66	1.9
	<i>mevn-A</i>	0.61	0.28	0.37	0.40	2.0
	<i>sic1-A</i>	0.66	1.30	0.49	NA	3.1
	<i>tk18-A</i>	0.96	0.56	1.01	0.52	1.4
B	<i>mevn-B</i>	0.39	0.34	0.38	0.52	1.2
	<i>actr-B</i>	1.94	0.35	0.35	0.46	0.5
	<i>cbpn-B</i>	1.65	0.69	0.38	0.62	1.6
	<i>p53t-B</i>	0.92	0.36	0.34	0.53	1.3
	<i>rsp8-B</i>	1.03	0.41	NA	0.60	1.0
C	<i>mevn-C</i>	0.37	0.34	0.32	0.37	1.4
	<i>actr-C</i>	1.83	0.28	0.32	0.39	3.6
	<i>cbpn-C</i>	1.64	0.67	0.40	0.62	2.0
	<i>p53t-C</i>	0.91	0.28	0.32	0.47	2.7
	<i>rsp8-C</i>	1.07	0.57	NA	0.64	2.5
0	<i>mevn-0</i>	0.71	0.53	0.35	0.64	2.0
	<i>actr-0</i>	1.72	0.50	0.36	0.53	3.4
	<i>cbpn-0</i>	2.94	0.93	0.49	0.70	2.5
	<i>p53t-0</i>	0.92	0.39	0.29	0.72	2.1
	<i>rsp8-0</i>	1.67	0.83	NA	0.59	2.9

8

9

10

1

2 **Table 4: Accuracy of CD spectrum predictions.** Summary of RMSDs between measured CD spectra, and
3 CD spectra predicted by SESCA from IDP8 reference ensemble models. RMSD values are shown for the
4 four basis sets used for the predictions and described in the text, expressed in 1000 Mean Residue
5 Ellipticity (kMRE) units. The most accurate predictions are indicated in bold, and RMSD values for the
6 basis set used in the ensemble refinement of group B are underlined. The average (avg) and standard
7 deviations (SD) of the RMSD values for each basis set are shown at the bottom of the table.

Group	entry	DS-dTSC3	DSSP-1SC3	HBSS-3SC1	DS5-4SC1	PDBMD2CD	Dichro
A	1	1.92	1.91	1.55	1.92	6.31	8.48
	2	1.95	1.87	1.91	2.88	3.63	5.50
	3	2.67	2.39	0.57	2.69	3.60	3.09
	4	1.36	1.40	1.55	3.18	2.74	11.59
B	5	<u>1.22</u>	1.17	1.60	1.61	4.74	8.76
	6	2.79	2.60	2.66	0.45	6.26	8.33
	7	1.07	1.60	1.06	<u>1.47</u>	2.87	5.42
	8	1.33	1.43	1.66	1.89	6.41	15.61
	9	2.74	2.07	4.29	<u>0.96</u>	6.91	10.49
C	10	1.41	1.20	1.66	2.73	4.61	9.86
	11	4.89	3.95	3.97	3.59	7.02	6.97
	12	1.09	1.85	1.06	1.21	3.23	5.08
	13	2.66	1.35	2.25	0.64	6.98	13.41
	14	3.07	1.83	4.34	2.48	7.08	9.33
	avg	2.2	1.9	2.2	2.0	5.2	8.7
	sd	1.1	0.7	1.2	1.0	1.7	3.4

8

9

1 **Table 5: MD simulations and parameters.** Summary of all MD simulations used for ensemble refinement
 2 of the newly derived IDP8 models. Different columns indicate simulation parameters for different IDP
 3 domains (abbreviations shown in the first row). The subsequent rows indicate simulation parameters
 4 used for all trajectories, temperature (T) in degrees Celsius, pressure (P) in atmospheres, and ion
 5 concentration in mol/dm³ (C_{ion}). The rows further below describe MD trajectories, separated by
 6 horizontal lines, indicating the used force field (FF), the number of calculated trajectories (Ntraj), total
 7 simulation time (tsim), and the number of frames used for ensemble refinement (Nfr). Force fields are
 8 abbreviated as A99SB-disp (Amber99SB with dispersion corrections)⁶¹, A99SB-ws (Amber99SB with
 9 rescaled TIP4P water)⁶⁰, A03-ws (Amber03 with rescaled water interactions)⁵⁹, A14SB-OPC
 10 (Amber14SB⁶² with an OPC water mode⁶³), C22S-TIPS3P (CHARMM22 star with modified TIP3P water)
 11 ⁶⁴, and C36M-OPC (CHARMM36M with OPC water model)⁶⁶.

12

System	mevn	actr	cbpn	p53t	rsp8
T (C°)	25	25	25	25	25
P (atm)	1	1	1	1	1
Cion (M)	150	50	50	150	50
FF	A99SB-disp	A03-ws	A03-ws	A99SB-ws	A03-ws
Ntraj	3	2	2	30	1
tsim (us)	30	9	20	600	10
Nfr	12 000	2400	5000	12000	2000
FF	C36M-OPC	C22S-TIPS3P	A99SB-disp	C36M-OPC	C36M-OPC
Ntraj	6	3	20	20	5
tsim (us)	42	30	100	200	25
Nfr	33000	7500	20000	20000	2500
FF			C36M-OPC		C22S-TIPS3P
Ntraj			20		1
tsim (us)			100		1
Nfr			10000		1000
FF			A14SB-OPC		
Ntraj			20		
tsim (us)			100		
Nfr			10000		

13

14

15

16

1

2 **Table 6: Bayesian Maximum Entropy refinement parameters.** The table summarizes details of Bayesian
3 maximum entropy (BME) refinement of IDP8 ensemble models. Columns denote different disordered
4 models (abbreviations of the model shown in the first row). The subsequent rows show initial ensemble
5 size (N0), Theta scaling parameter (Θ) and final ensemble size (Nf) for group B and group C refinements.
6 Both sets of refinements started from the same respective initial ensemble (described in Methods)
7 using different experimental data.

system	mevn	actr	cbpn	p53t	rsp8
N0	45 000	10000	45000	32000	5500
Θ -B	10	5	5	5	20
Nf-B	5x20	5x20	5x20	1X50	5x50
Θ -C	20	10	10	5	10
Nf-C	5X20	5X20	5x20	5x50	5x50

8