

Implementation of a Bayesian Secondary Structure Estimation Method for the SESCO Circular Dichroism Analysis Package

Gabor Nagy¹ & Helmut Grubmüller¹

October 2020

¹: Department of Theoretical and Computational Biophysics, Max Planck Institute for Biophysical Chemistry, Am Fassberg 11, 37077 Göttingen, Germany

Keywords: CD spectroscopy, SS estimation, Bayesian statistics, model validation, accuracy improvement, SESCO

1 **Abstract**

2 Circular dichroism spectroscopy is a structural biology technique frequently
3 applied to determine the secondary structure composition of soluble proteins.
4 Our recently introduced computational analysis package SESCA aids the in-
5 terpretation of protein circular dichroism spectra and enables the validation of
6 proposed corresponding structural models. To further these aims, we present
7 the implementation and characterization of a new Bayesian secondary structure
8 estimation method in SESCA, termed SESCA_bayes. SESCA_bayes samples
9 possible secondary structures using a Monte Carlo scheme, driven by the like-
10 lihood of estimated scaling errors and non-secondary-structure contributions of
11 the measured spectrum. SESCA_bayes provides an estimated secondary struc-
12 ture composition and separate uncertainties on the fraction of residues in each
13 secondary structure class. It also assists efficient model validation by providing
14 a posterior secondary structure probability distribution based on the measured
15 spectrum. Our presented study indicates that SESCA_bayes estimates the sec-
16 ondary structure composition with a significantly smaller uncertainty than its
17 predecessor, SESCA_deconv, which is based on spectrum deconvolution. Fur-
18 ther, the mean accuracy of the two methods in our analysis is comparable, but
19 SESCA_bayes provides more accurate estimates for circular dichroism spectra
20 that contain considerable non-SS contributions.

21 **1 Introduction**

22 Circular dichroism (CD) spectroscopy in the far ultraviolet (UV) range (175-
23 260 nm) is an established method to study the structure of proteins in solution
24 [1, 2], because of the conformation-dependent characteristic CD signal of pep-
25 tide bonds that comprise the backbone of all proteins and oligo-peptides. In
26 particular, the CD spectrum is known to change with the secondary structure

27 (SS) of proteins, and markedly different spectra are observed for proteins rich in
28 α -helices, β -sheets, and disordered regions [3, 4]. Because of these characteristic
29 signals, it is common to interpret CD spectra by decomposing them into a set
30 of basis spectra that each represent the average CD signal of pure (secondary)
31 structure elements.

32 The CD analysis package SESCOA (Structure-based Empirical Spectrum Cal-
33 culation Approach) [5] allows for using several empirical basis spectrum sets in
34 two methods. The first method predicts a theoretical CD spectrum from a
35 proposed SS composition, which is typically obtained from a model structure
36 or structural ensemble. The second method fits a measured CD spectrum to
37 estimate the protein SS composition. Both methods can be used to validate
38 protein structural models. The accuracy and precision of validation methods
39 is mainly limited by scaling errors due to the uncertainty of the measured pro-
40 tein concentration and non-SS contributions that are not represented in the
41 basis spectra. We have quantified the uncertainty caused by these deviations
42 between measured CD spectra and their predicted SS signals previously [6].

43 The same study also revealed a potential caveat in the current SS estimation
44 method used in SESCOA. In this deconvolution method, a linear combination of
45 selected basis spectra is used to approximate a measured CD spectrum of the
46 protein of interest. The coefficients of the approximation with the smallest
47 deviation are used to estimate the fraction of SS elements in the protein under
48 the measurement conditions. Unfortunately, the interference caused by non-SS
49 contributions may increase the deviation from the measured spectrum for some
50 SS compositions and decrease it for others, which may lead to significant errors
51 in deconvolution-based SS estimates.

52 To alleviate this problem, we developed and implemented a new SS esti-
53 mation method for SESCOA. The Python module, SESCOA_bayes determines the

54 likelihood of putative SS compositions using a Bayesian inference framework for
55 a given measured CD spectrum and a basis spectrum set. This method uses
56 the expected joint probability distribution of deviations caused by scaling errors
57 and non-SS contributions, and thus fully accounts for the uncertainty caused by
58 these two experimental factors. Here, we describe the theoretical background,
59 general workflow, as well as input and output parameters of this implementa-
60 tion. Further, we will assess the accuracy and precision of this method through
61 a series of sample applications.

62 **2 Theory: Bayesian SS probabilities**

63 Our goal using this method is to determine the conditional probability $P(SS|CD)$
64 of SS compositions given a previously measured CD spectrum. According to
65 Bayes' rule [7], this probability can be inferred according to

$$P(SS_j|CD) \propto P(CD|SS_j) \cdot P(SS_j), \quad (1)$$

66 where $P(CD|SS_j)$ is the probability of observing the measured spectrum for a
67 protein with a given SS composition j (*i.e.*, the likelihood function) and $P(SS_j)$
68 is the prior probability of the given SS composition of the protein. As shown
69 in Fig. 1 (top), the likelihood $P(CD|SS_j)$ is determined in five steps. First,
70 the SS signal is predicted from the SS composition of interest (C_{ji}) using an
71 appropriate basis spectrum set (B_{il}), as discussed in our previous study [5].
72 Second, if the basis set provides side chain corrections based on the protein
73 sequence, they are added to the predicted spectrum. Third, the measured CD
74 spectrum is rescaled to minimize the root-mean-square deviation (RMSD) from
75 the predicted spectrum. The obtained scaling factor α_j quantifies and eliminates
76 deviations from scaling errors of the measured spectrum, whereas the RMSD

77 from the rescaled spectrum ($RMSD_j$) quantifies the average deviation due to
78 unaccounted non-SS contributions. Once $RMSD_j$ and α_j (collectively CD de-
79 viations) are computed, the likelihood of such deviations is determined from
80 the joint probability distribution ($P_{RMSD,\alpha}$, see below), which also estimates
81 the likelihood of observing the measured CD spectrum for the given SS compo-
82 sition $P(CD|SS_j)$. Finally, to compute the posterior probability $P(SS_j|CD)$
83 of SS composition j , the CD spectrum likelihood is multiplied by the prior SS
84 probability.

85 **3 Methods**

86 **3.1 Joint probability distributions**

87 We computed discrete joint-probability distribution functions for `SESCA_bayes`
88 that can be used to determine CD spectrum likelihoods. These probability
89 distributions were computed from CD deviations extracted from SS estimations
90 of previously measured CD spectra. Reference CD spectra were taken from the
91 SP175 reference set [8], which contains 71 synchrotron radiation CD (SR-CD)
92 spectra of globular proteins with varying SS compositions. The CD spectrum
93 of Jacalin (SP175/41) was discarded from the data set due to issues reported
94 during the measurement and its unusually large estimated CD deviations.

95 The joint probability distribution functions of CD deviations were con-
96 structed as the sum of 70 two-dimensional Gaussian functions, each representing
97 the estimated scaling factors and non-SS contributions of a reference spectrum
98 from the SP175 set. The mean and the variance of these Gaussian functions
99 was determined by averaging over multiple $RMSD_j$ and α_j values obtained for
100 each CD spectrum from SS estimations using four different basis spectrum sets.
101 This approach yielded likelihood functions that were defined for a wide range

102 of possible CD deviations, and took the uncertainty due to discretization errors
103 of the basis spectrum determination into account.

104 In SESCO there are two types of basis sets, those that are solely based on
105 SS compositions, and those that also include side chain corrections. Because the
106 average size of CD deviations differ for these two basis set types, we determined
107 two probability distributions shown in Fig. 2. The joint probability distribution
108 function for basis set without side-chain corrections (left) was calculated from
109 CD deviations estimated using the basis sets DS-dT, DSSP-1, HBSS-3, and DS5-
110 4. For basis sets including side-chain corrections, the joint probability of CD
111 deviations (right) were computed using the basis sets DS-dTSC3, DSSP-1SC3,
112 HBSS-3SC1, and DS5-4SC1. For clarity, the Figure shows both a linear (top
113 row) as well as logarithmic (bottom row) representation of the CD deviation
114 likelihood. For both likelihoods, the one-dimensional probability distribution
115 of $RMSD_j$ was also calculated, which can be used to estimate the secondary
116 structure from CD spectra without regards to the applied scaling factors, albeit
117 these estimates naturally have a lower precision.

118 **3.2 Synthetic spectra**

119 To test the accuracy of the Bayesian SS estimation method, six synthetic CD
120 spectra were created using a linear combination of the three basis spectra from
121 the DS-dT basis set (as discussed in our previous study [5]). To this aim,
122 the coefficients shown in Table 1 for the basis spectra α -helix, β -strand, and
123 Other for each spectrum were used. For five of six synthetic spectra ($k= 1$
124 to 5), random coefficients were generated from uniformly distributed random
125 numbers between zero and one, subsequently normalized to sum up to one. For
126 the sixth synthetic spectrum ($k= 6$), the coefficients 0.3, 0.4, and 0.3 as well as
127 the non-SS contributions (see below) were adopted from our previous study [6]

128 for comparison.

129 To model the effects of experimental deviations from the ideal SS signal, the
130 CD spectra were modified by adding non-SS signals and scaling errors. The size
131 of these CD deviations for each synthetic spectrum was quantified by the scaling
132 factors α_k and the root-mean-squared intensities of non-SS signals $RMSI_k$ listed
133 in Table 1. Synthetic spectrum 1 ($k=1$) was a positive control without any
134 CD deviations ($\alpha_k=1.0$, $RMSDI_k=0.0$ kmRE), spectra 2 and 6 included
135 small (0.4 kmRE) and large (3.5 kmRE) non-SS deviations, respectively, but
136 no scaling errors. CD deviations for spectra 3, 4, and 5 were drawn from the
137 marginal distributions of experimentally observed scaling factors and non-SS
138 contributions using rejection sampling.

139 The shapes of the non-SS signals were chosen as sums of Gaussian functions

$$S_{jl}^{\text{nonSS}} = \sum_{g=1}^G \frac{I_g}{\sqrt{2\pi\sigma_g^2}} \times e^{-\frac{(\lambda_l - \mu_g)^2}{2\sigma_g^2}}, \quad (2)$$

140 where the non-SS signal S_{jl} of protein j at wavelengths λ_l from 178 to 269
141 nm was computed from the following randomly chosen parameters. The number
142 of Gaussians G was chosen from the range 1 to 5, the relative peak intensity for
143 Gaussian g I_g was chosen between -20.0 and 20.0, with a peak position μ_g chosen
144 from 178 to 241 nm, and peak half-widths σ_g chosen between 2 and 37 nm. Once
145 the parameters were determined, the non-SS signal at every wavelength (using
146 1 nm spacing) was calculated, and the non-SS signal intensity was rescaled to
147 match the previously defined RMSI values in Table 1.

148 The final synthetic spectra were computed by determining the SS signals
149 first, by adding the appropriately scaled non-SS signal contributions in a second
150 set, and finally by rescaling the resulting CD spectrum according to the indicated
151 scaling factor.

152 4 Algorithm overview

153 Our newly implemented Python module `SESCA_bayes.py` performs a Monte-
154 Carlo (MC) sampling in SS space to determine the most probable SS compo-
155 sition of a protein based on its measured CD spectrum. Figure 3 shows the
156 flowchart of the algorithm that is divided into three phases: preparation, sam-
157 pling, and evaluation.

158 4.1 Preparation and input parameters

159 In the preparation phase, input, output, and run parameters are read based
160 on the user-provided command line arguments. If `SESCA_bayes.py` is used as
161 a Python module, an array of arguments can be processed by the function
162 `Read_Args` and passed to the `Main` function to run the algorithm. Arguments in
163 `SESCA` are identified by preceding command flags (marked by the "@" character
164 in the first position. There are four input files – shown as blue parallelograms
165 in Fig. 2 – that `SESCA_bayes` accepts, each read in white-space separated data
166 blocks stored as simple *ascii* text files.

167 The CD spectrum file (specified using the `@spect` flag) should contain two
168 columns, wavelength in nanometers (nm) and CD signal intensity in 1000 mean
169 residue ellipticity (kMRE) units. This file must be specified for `SESCA_bayes`,
170 and if no command flags are provided, the first argument is automatically rec-
171 ognized as a CD spectrum file.

172 The side-chain correction file (specified by `@corr`) is an optional file to add
173 baseline or sequence-dependent side-chain correction to the predicted CD spec-
174 trum, which are independent of the SS composition. If the basis spectrum
175 set has basis spectra to calculate side-chain contributions, these signals can be
176 computed before running `SESCA_bayes`, and added as a correction.

177 The Bayesian parameter file (`@par`) contains several data blocks, most im-

178 portantly, the binned probability distribution function of CD deviations $P_{RMSD,\alpha}$
179 (likelihood function), prior SS probability distributions for the SS composition
180 $P(SS_j)$ and scaling factors $P(\alpha_j)$, as well as the MC step parameters. If no
181 parameter file is provided by the user, `SESCA_bayes.py` uses one of two default
182 parameter files (`Bayes_2D_SC.dat` and `Bayes_2D_noSC.dat`) found in the "libs"
183 sub-directory of `SESCA`, depending on whether a side chain correction file was
184 provided or not. These files contain one of the two likelihood functions shown
185 in Fig. 2, and uniform prior SS probability distributions. A more detailed
186 description of the parameter blocks is provided in the `examples` sub-directory
187 (`examples 5`).

188 The basis set file (`@lib`) contains several data blocks for CD spectrum cal-
189 culations, including a block where the CD intensity of 3-6 basis spectra at each
190 wavelength (175-269 nm) is provided. Several derived basis sets are available
191 in `libs` sub-directory, and a detailed description of the data blocks is given in
192 `example 1`.

193 In addition to the input files, `SESCA_bayes` recognizes several additional
194 command flags to modify program behavior. The number of initial SS composi-
195 tions for MC sampling phase is specified by `@size`. The number of MC steps per
196 initial SS composition is set by `@iter`. The `@scale` flag allows the user to control
197 whether the measured CD spectrum is rescaled before determining the deviation
198 from the predicted CD spectra or not. In the absence of these command flags,
199 the values 100, 500, and 1 (yes) are used for the SS estimation.

200 Finally, three command flags control the output of `SESCA_bayes.py`; provid-
201 ing a "0" argument to any of these flags disables writing the associated output.
202 The command flag `@write` specifies the file name for the primary output, and
203 if no command flags are given, `SESCA_bayes` automatically recognizes the sec-
204 ond argument as primary output file. This file contains a summary of the

205 input parameters, binned posterior probability distributions for the SS compo-
206 sitions and scaling factors, as well as the most probable SS fractions and their
207 uncertainties. The command flags @proj and @data allows the user to print
208 secondary output files. The @proj flag specifies a file name for heatmap-style
209 two-dimensional projection of the posterior SS distribution. The projection is
210 made along two SS fractions selected using the @pdim flag Finally, the flag
211 @data specifies a file name for printing all the sampled SS compositions the
212 primary output is computed from, along with their estimated CD deviations,
213 prior and posterior probabilities. By default, only the primary output file is
214 printed into 'SS_est.out', and no secondary output is written.

215 **4.2 Monte Carlo sampling**

216 To determine the most probable SS composition of the protein based on its CD
217 spectrum, sampling of the SS space is required. To this aim, SESCO_bayes uses
218 a MC sampling scheme starting from N (set by @size) initial SS compositions,
219 drawn from the prior SS distribution using rejection sampling. As the center
220 part of Fig. 3 shows, at every step t of the MC sampling phase, a change
221 on each of the SS compositions ($C_{ji,t}$) is attempted. The change is realized
222 by transferring a given SS fraction between two randomly chosen SS classes,
223 yielding a new SS composition ($C'_{ji,t}$). The amount of the transferred SS fraction
224 from the donor class to the acceptor class is determined based on the distribution
225 specified in the Bayesian parameter file. If no distribution is provided, the
226 fraction is drawn from a Gaussian distribution with a mean of 0.05 and variance
227 of 0.1. To remain in the space of possible SS compositions, the transferred SS
228 fraction cannot exceed the current fraction assigned to the donor class, and
229 classes that currently have a fraction of zero assigned to them cannot be selected
230 as donors.

231 After the changes are attempted, the posterior probabilities P'_{jt} of the new
232 SS compositions are calculated (see Section 2) and compared to the posterior
233 probabilities (P_{jt}) of the SS compositions before the change. The attempted
234 change is accepted or rejected by applying the Metropolis criterion to the ratio
235 of posterior probabilities, *i.e.* the change is accepted if the ratio P'_{jt}/P_{jt} is
236 larger than a randomly generated number between zero and one. If the change
237 is accepted, $C'_{ji,t}$ is added to the sampled SS distribution and used as the initial
238 SS composition $C_{ji,t+1}$ in the next MC step, otherwise $C_{ji,t}$ is added to the
239 sampled SS distribution (again) and is used in the next MC step. This procedure
240 is repeated until the specified number of MC attempts is reached, and yields
241 $N \times t_{max}$ sampled SS compositions. The sampled SS compositions resemble the
242 prior SS distribution during the initial MC steps but converge towards an SS
243 distribution weighted by the posterior SS probabilities.

244 4.3 Sample evaluation

245 The sampled SS distribution is analysed in the evaluation phase, as shown in
246 the bottom part of Fig 3. To avoid the over-representation of very low posterior
247 probability SS compositions, a fraction of the initially sampled SS compositions
248 may be discarded from final SS distribution. This fraction can be set by the
249 user through the @discard flag, otherwise, the initial 5% of SS compositions is
250 discarded. The remaining probability-weighted ensemble of possible SS compo-
251 sitions is used to compute the estimated SS composition C_{ji}^{est} for the protein, the
252 estimated scaling factor α_j^{est} , as well as to approximate the discrete posterior
253 probability distribution for both quantities.

254 The estimated SS composition is determined by computing the mean and
255 standard deviation (SD) of each SS fraction over the sampled SS compositions.
256 Similarly, the most probable scaling factor is computed as the mean and SD of

257 scaling factors estimated for the sampled SS compositions. The discrete prob-
258 ability distribution for both scaling factors and SS compositions are computed
259 by binning all sampled SS compositions and scaling factors using the parame-
260 ters extracted from the prior distributions provided in the Bayesian parameter
261 file. The number of sampled SS compositions and scaling factors in each bin is
262 normalized by the final sample size to obtain the discrete probability distribu-
263 tions. The computed estimates, their uncertainties and the discrete probability
264 distributions are all written in the primary output file (defined by the @write
265 flag) and returned as output by the SESCO.bayes module. If requested (@proj
266 flag), the sampled SS compositions can be printed in a separate file. Finally,
267 the two-dimensional projection of posterior SS distribution along two chosen SS
268 fractions can also written into a separate output file (@proj flag), formatted as
269 a human readable heat map, that can be easily processed into images using *e.g.*
270 Python’s Matplotlib module [9] or external visualization programs.

271 5 Testing the Algorithm

272 5.1 Accuracy and precision

273 The accuracy and precision of the Bayesian SS estimation was tested using the
274 10 CD spectra listed in Table 1. Six of these spectra ($k= 1-6$) are synthetic spec-
275 tra that were generated from a given SS composition, but modified by adding
276 artificial non-SS signals and scaling errors (see Section 3.2) to emulate CD devi-
277 ations in real measured spectra. The remaining four CD spectra ($k= 7-10$) are
278 measured spectra from the SP175 set [8], for which the estimated SS composi-
279 tions are compared to those extracted from the (protein data bank) structure
280 of the reference protein. Table 1 also lists the (estimated) CD deviations of all
281 ten CD spectra, quantified by the scaling factors α_k and the root-mean-square

282 intensity ($RMSI_k$) of non-SS signals in each spectrum.

283 To test the accuracy of `SESCA.bayes`, we estimated the SS composition
284 of the above ten CD spectra using the same DS-dT basis set with three SS
285 classes (α -helix, β -strand, and Other) that was used to generate the synthetic
286 spectra. The obtained Bayesian estimates for the test set are summarized in
287 Table 2. This table includes the mean and SD (in parentheses) of SS fractions
288 of the sampled posterior distributions, as well as the total SS deviation from
289 the reference SS compositions, computed according to

$$\Delta SS_k = \frac{1}{2} \sum_{i=1}^F |C_{ki}^{est} - C_{ki}^{ref}|, \quad (3)$$

290 where C_{ki}^{est} are the estimated SS fractions and C_{ki}^{ref} are the reference SS fractions
291 listed in Table 1.

292 The obtained SS fractions show a fairly consistent 0.03 to 0.06 uncertainty.
293 As expected, 27 of 30 SS fractions are within two SD of their reference value,
294 with no significant difference in accuracy between synthetic and measured CD
295 spectra. In addition, the calculated total SS deviations (ΔSS) from the reference
296 structures range between 0.03 and 0.12, and eight of ten values are also smaller
297 than the estimated uncertainty of the estimation (two SD) that was calculated
298 from the individual SD of SS fractions (σ_{ki}) according to

$$\sigma_k = \frac{1}{2} \sqrt{\sum_{i=1}^F \sigma_{ki}^2}. \quad (4)$$

299 5.2 Comparison to deconvolution

300 Next, we compare the accuracy and precision of the Bayesian estimates to that
301 of SS estimates obtained through spectrum deconvolution. To this aim, we esti-
302 mated SS compositions with the deconvolution module of `SESCA` (`SESCA.deconv`)

303 for the same ten CD spectra (Table 1), using the same DS-dT basis spectrum set.
304 The deconvolution was carried out using the most accurate protocol (method
305 D2) tested previously [6]. This method constrains the basis spectrum coefficients
306 to positive values, but normalizes them to unity only after the search for the best
307 approximation. The SS compositions obtained using `SESCA_deconv` are listed
308 in Table 3, along with the total SS deviations from reference SS compositions
309 (found in Table 1). The total SS deviation of deconvolution estimates (ΔSS_k)
310 ranges from 0.0 to 0.29. The mean SS deviation for the whole set (0.08) is
311 similar to that of the Bayesian estimates (0.07), but shows a significantly larger
312 scatter (0.9 vs. 0.03). All three CD spectra with larger than average SS devi-
313 ations ($k= 3,4,8$) have large non-SS contributions (2.0-2.9 kMRE), which is in
314 line with our previous findings that non-SS contributions may be detrimental
315 to the accuracy of deconvolution methods.

316 Although the `SESCA_deconv` module does not provide information on the
317 uncertainty of individual SS fractions, many `SESCA` basis sets (including DS-
318 dT) include a calibration curve to estimate the expected total SS deviation if
319 the true SS composition is unknown. This curve was computed from $4.9 \times$
320 10^5 synthetic spectrum-structure combinations, which were binned according to
321 their estimated non-SS contributions ($RMSD_j$), to provide an expected mean
322 and SD of SS deviations for a given (rescaled) RMSD. Comparing the true SS
323 deviations of the deconvolution results with their estimated values shows that
324 these estimates correctly describe the precision of the deconvolution method:
325 six of ten ΔSS_k values are within 1 SD of the estimated total deviation, and
326 all ten fall within 2 SD. However, the average uncertainty of the deconvolution
327 (0.09) is again considerably larger than that of the Bayesian SS estimates (0.04),
328 and it increases with increasing non-SS contributions.

329 In summary, Bayesian SS estimation and spectrum deconvolution provides

330 SS estimates that – in most cases – have a similar accuracy. However, Bayesian
331 SS estimates are considerably more precise when significant non-SS contribu-
332 tions are present in the measured spectrum. Further, the Bayesian approach
333 provides uncertainties for each individual SS fraction as well as the optimal scal-
334 ing factor for the measured CD spectrum, which is an additional advantage of
335 the new method.

336 5.3 Example spectrum analysis

337 To further investigate the differences between the two methods, we analysed the
338 SS estimates for the CD spectrum with the largest difference between the de-
339 convolution and Bayesian SS estimates. Figure 4A shows the obtained posterior
340 SS probability distribution for synthetic spectrum 3, which contains larger than
341 average non-SS contributions (2.02 kMRE). The heatmap shown in Fig. 4A
342 illustrates that the most likely SS compositions are indeed clustered around the
343 SS composition the synthetic spectrum was created from (shown as a red cross),
344 with the highest posterior probability regions (shown in dark green) located in
345 the immediate ($\Delta SS_k < 0.05$) vicinity of correct SS composition. However, the
346 SS composition determined by deconvolution (purple cross) has a much higher
347 α -helix content and it is not in a high-probability region in the Bayesian SS
348 estimation.

349 To examine why the two algorithms evaluate the proposed SS compositions
350 differently, in Fig. 4B we computed the predicted CD signals of the two es-
351 timated SS compositions, rescaled them, and compared them to the synthetic
352 spectrum, as is done during the deconvolution process. The figure shows that
353 with the proper scaling factor both SS compositions approximate the synthetic
354 spectrum well, but the deconvolution estimate (purple dashed line, $RMSD_j$:
355 1.31 kMRE) fits slightly better than the Bayesian estimate (blue dashed lines,

356 $RMSD_j$: 1.71 kMRE).

357 In contrast, the Bayesian SS estimation rescales the synthetic CD spectrum
358 to match the predicted spectra, and evaluates the likelihood of the SS com-
359 positions based on the joint probability of their non-SS contributions $RMSD_j$
360 and scaling factor α_j , as shown in Fig 4C. Although the two estimates have a
361 comparable RMSD in this method as well, the deconvolution estimate requires
362 a scaling factor (α_j : 1.99) to achieve a good agreement that is shown to be very
363 unlikely according to the joint-probability map in Fig. 3. Comparing the two
364 estimated SS signals (dashed lines) to the SS signal of the true SS composition
365 (in red) illustrates how considering scaling factors improves the precision of the
366 SESCO.bayes. In this case, eliminating SS compositions with unlikely scaling
367 factors from the sampled distribution allowed a fairly accurate (RMSD: 0.99
368 kMRE) approximation of the true SS signal.

369 **Declarations**

370 **Funding:** This research project was funded and supported by the Alexander
371 von Humboldt Foundation and the Max Planck Society.

372 **Conflicts of interest:** The authors declare no conflict of interest.

373 **Code availability:** the new SESCO implementation based on this study is
374 available at: <https://www.mpibpc.mpg.de/sesca>

375 **License:** SESCO is free available under GNU general public license 3 (GPLv3)

376 **Authors' contributions:** G.N. designed and performed the computational
377 analysis, and implemented code improvements. H.G. is the corresponding au-
378 thor, and assisted with the conceptualisation. Both authors contributed to
379 writing the manuscript.

380 **Acknowledgements:** The authors would like to thank K. Blom and P. Kellers
381 for the aid in editing the manuscript.

References

- [1] Gerald D. Fasman, editor. *Circular Dichroism and the Conformational Analysis of Biomolecules*. Springer US, Boston, MA, 1996.
- [2] Norma J. Greenfield. Methods to estimate the conformation of proteins and polypeptides from circular dichroism data. *Analytical biochemistry*, 235(1):1–10, 1996.
- [3] Sharon M. Kelly, Thomas J. Jess, and Nicholas C. Price. How to study proteins by circular dichroism. *Biochimica et Biophysica Acta (BBA) - Proteins and Proteomics*, 1751(2):119–139, August 2005.
- [4] S. Brahms and J. Brahms. Determination of Protein Secondary Structure in Solution by Vacuum Ultraviolet Circular Dichroism. *Journal of Molecular Biology*, 138:147–178, 1980.
- [5] Gabor Nagy, Maxim Igaev, Nykola C. Jones, Søren V. Hoffmann, and Helmut Grubmüller. SESCO : Predicting Circular Dichroism Spectra from Protein Molecular Structures. *Journal of Chemical Theory and Computation*, August 2019.
- [6] Gabor Nagy and Helmut Grubmüller. How accurate is circular dichroism-based model validation? *European Biophysics Journal*, 49(6):497–510, September 2020.
- [7] Andrew Gelman and John B. Carlin. *Bayesian Data Analysis*. Texts in Statistical Science. Chapman and Hall/CRC, 3rd edition, 2014.
- [8] J. G. Lees, A. J. Miles, F. Wien, and B. A. Wallace. A reference database for circular dichroism spectroscopy covering fold and secondary structure space. *Bioinformatics*, 22(16):1955–1962, August 2006.
- [9] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007.

Tables

Table 1: SS compositions and CD deviations of model proteins. Columns show the index and name of the respective model protein, the fraction of its amino acids assigned to the SS classes α -helix, β -strand, and Other, as well as scaling factors α_k and root-mean squared intensities $RMSI_k$ of non-SS signals to quantify scaling errors and non-SS contributions in the protein CD spectrum, respectively. Synth denotes synthetic spectrum in the proteins name, whereas Lysm, Dqd-1, and Sub-C abbreviate Lysozyme, Dehydroquinase dehydratase I, and Subtilisin Carlsberg, respectively. Note that SS fractions, scaling factors, and non-SS contributions for all synthetic proteins (k= 1-6) were parameters used to generate their CD spectrum, whereas for real reference proteins (k= 7-10), all values were computed based on their measured spectra and protein data bank structures (193L, 2DHQ, 1KU8, and 1SCD, respectively).

k	protein	α -helix	β -strand	Other	α_k	$RMSI_k$
1	Synth-1	0.11	0.40	0.49	1.0	0.0
2	Synth-2	0.41	0.20	0.39	1.0	0.4
3	Synth-3	0.43	0.10	0.47	1.2	2.0
4	Synth-4	0.27	0.26	0.47	1.6	2.7
5	Synth-5	0.00	0.33	0.67	1.4	0.7
6	Synth-6	0.30	0.40	0.30	1.0	3.6
7	Lysm	0.35	0.03	0.62	1.1	1.0
8	Dqd-1	0.43	0.18	0.39	1.1	2.9
9	Jacalin	0.01	0.28	0.71	0.3	3.2
10	Sub-C	0.30	0.12	0.58	0.4	1.2

Table 2: Bayesian secondary structure estimates. The table lists the index and name of the model protein, the estimated fraction of its amino acids assigned to SS classes α -helix, β -strand, and Other, as well as the total SS deviation ΔSS_k from the reference SS compositions shown in Table 1. The uncertainty (standard deviation) of each SS fraction and deviation is given in parentheses. Estimates that are more than 2 SD away from their reference value are highlighted in red.

k	protein	α -helix	β -strand	Other	ΔSS_k
1	Synth-1	0.14 (0.05)	0.44 (0.06)	0.43 (0.04)	0.06 (0.04)
2	Synth-2	0.49 (0.06)	0.19 (0.06)	0.32 (0.06)	0.08 (0.05)
3	Synth-3	0.45 (0.06)	0.12 (0.05)	0.43 (0.05)	0.03 (0.04)
4	Synth-4	0.22 (0.04)	0.21 (0.05)	0.57 (0.04)	0.09 (0.04)
5	Synth-5	0.03 (0.04)	0.26 (0.03)	0.71 (0.05)	0.07 (0.04)
6	Synth-6	0.36 (0.05)	0.32 (0.04)	0.31 (0.06)	0.08 (0.04)
7	Lysm	0.38 (0.05)	0.04 (0.05)	0.57 (0.05)	0.05 (0.04)
8	Dqd-2	0.48 (0.06)	0.06 (0.05)	0.47 (0.05)	0.12 (0.05)
9	Jacalin	0.01 (0.04)	0.31 (0.06)	0.68 (0.06)	0.03 (0.05)
10	Sub-C	0.26 (0.05)	0.13 (0.04)	0.61 (0.04)	0.04 (0.04)

Table 3: Secondary structure estimates based on spectrum deconvolution. The table lists the index and name of the model protein, the estimated fraction of its amino acids assigned to SS classes α -helix, β -strand, and Other, as well as the total SS deviation ΔSS_k from the reference SS compositions shown in Table 1. The values in parentheses after ΔSS_k show the mean and SD of the estimated total SS deviation computed from the rescaled RMSD between the measured (generated) CD spectrum and predicted spectrum of the SS estimate.

k	protein	α -helix	β -strand	Other	ΔSS_k
1	Synth-1	0.11	0.40	0.49	0.00 (0.00 \pm 0.02)
2	Synth-2	0.41	0.20	0.39	0.00 (0.05 \pm 0.03)
3	Synth-3	0.72	0.03	0.24	0.29 (0.16 \pm 0.09)
4	Synth-4	0.19	0.22	0.59	0.12 (0.08 \pm 0.05)
5	Synth-5	0.01	0.33	0.66	0.01 (0.06 \pm 0.04)
6	Synth-6	0.31	0.31	0.37	0.08 (0.14 \pm 0.09)
7	Lysm	0.34	0.06	0.60	0.03 (0.07 \pm 0.05)
8	Dqd-2	0.51	0.04	0.45	0.14 (0.09 \pm 0.06)
9	Jacalin	0.00	0.35	0.65	0.07 (0.20 \pm 0.09)
10	Sub-C	0.25	0.13	0.62	0.05 (0.07 \pm 0.05)

Figures

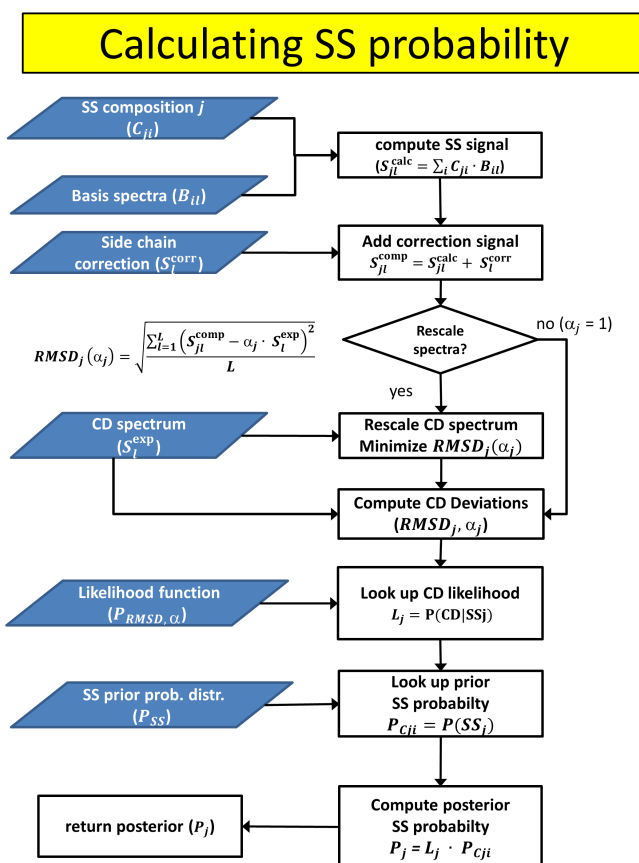


Figure 1: Secondary structure probability calculation scheme. The figure depicts the algorithm to compute the posterior probability of a given secondary structure j , based on its prior probability, and the deviations between its predicted CD signal and a given measured CD spectrum. Input data are depicted as blue parallelograms, operations as white rectangles, and decisions as white diamonds.

A - Without Side Chain Corrections **B - With Side Chain Corrections**

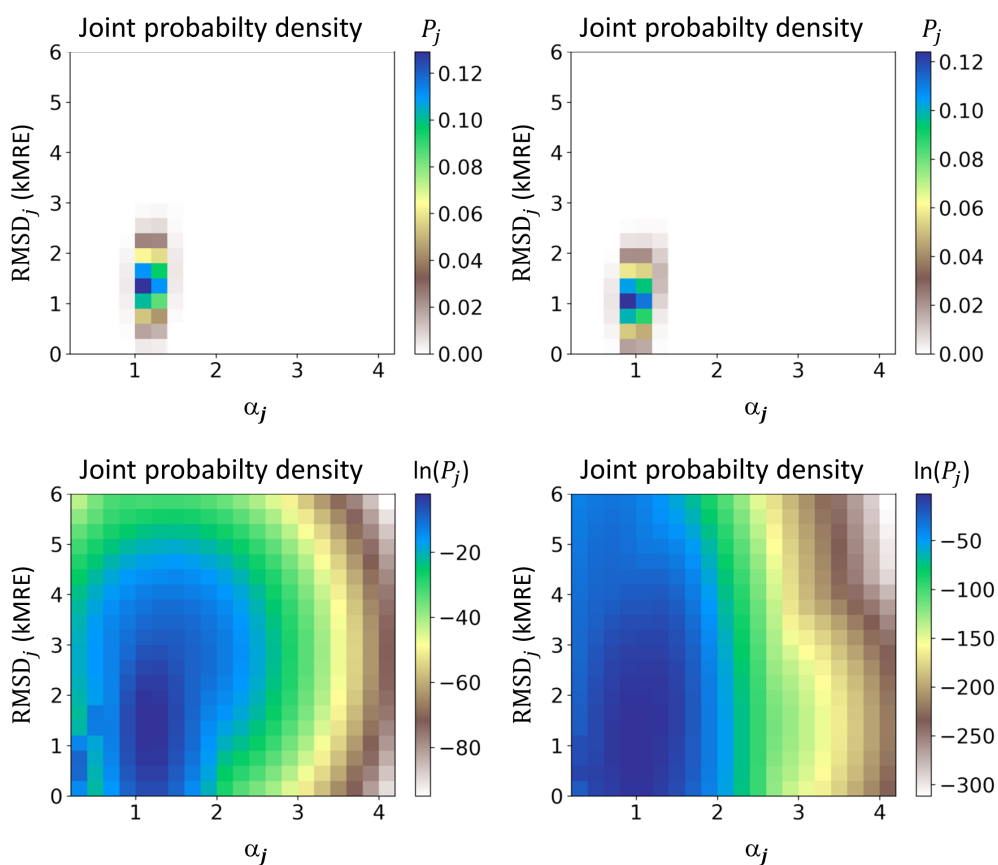


Figure 2: The panels depict the heat map representation of two likelihood functions provided for Bayesian SS estimation with SESCA. The estimated joint-probability distributions are shown for basis spectra that A) predict CD signals solely from SS information (left) and B) also include CD corrections from sequence-based side-chain information (right). Panels on the top and bottom show the same probability distributions using a linear and logarithmic color scale, respectively.

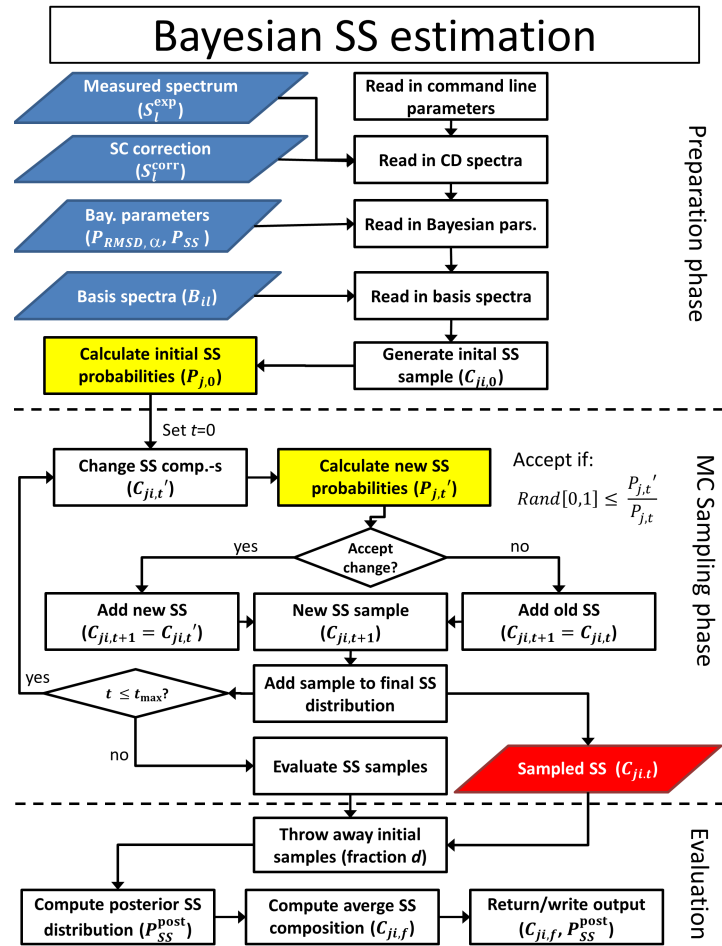


Figure 3: Schematic workflow of the Bayesian secondary structure estimation module in SESCO. The scheme depicts input data files as blue parallelograms, data on the sampled SS compositions are shown as a red parallelogram. Operations are depicted as white rectangles, and decisions are shown as white diamonds. Posterior probability calculation operations (see Fig. 1) are highlighted as yellow rectangles on the scheme.

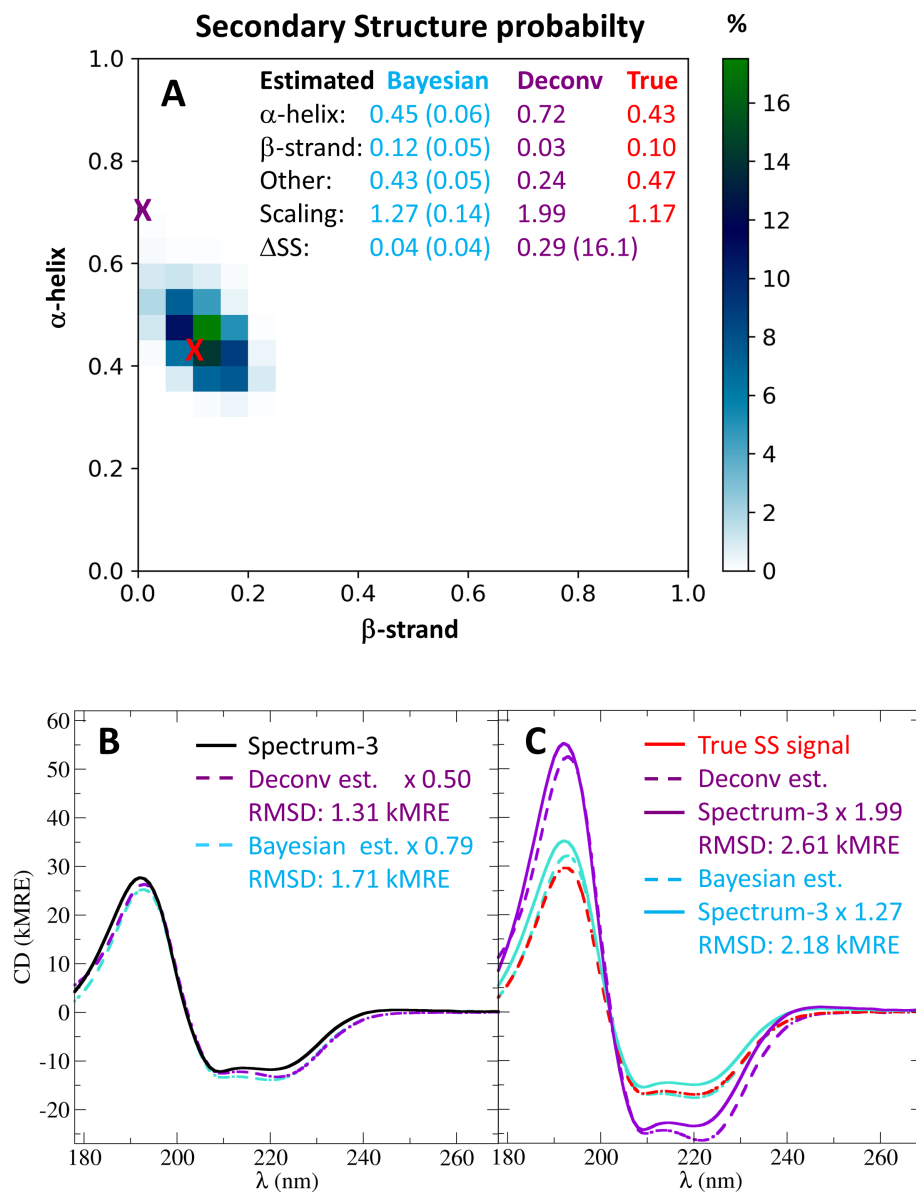


Figure 4: SS estimation for a synthetic CD spectrum. The figure compares the true SS composition (shown in red) with SS compositions obtained from Bayesian SS estimation (in blue) and spectrum deconvolution (in purple). A) shows the posterior probability distribution of sampled SS compositions in a heat map representation and indicates the true SS composition and the deconvolution SS estimate as crosses. The SS compositions, estimated scaling factors, and SS deviations are also listed in a tabulated format on the top. The difference on how the two estimates are evaluated by B) the deconvolution and C) Bayesian SS estimation are also shown. During deconvolution, the predicted CD signal of SS estimates is rescaled to match the measured CD spectrum, and the measure of quality is solely the RMSD. In the Bayesian approach, the measured spectrum is rescaled to match the predicted SS signals, and both the RMSD-s and the scaling factors are used to determine the most likely SS composition.