

# Accuracy and robustness of three-way decomposition applied to NMR data

Tu Luan<sup>a</sup>, Vladislav Yu Orekhov<sup>a</sup>, Aleksandras Gutmanas<sup>b,1</sup>, Martin Billeter<sup>b,\*</sup>

<sup>a</sup> Swedish NMR Centre, Göteborg University, Box 465, 405 30 Göteborg, Sweden

<sup>b</sup> Biophysics Group, Department of Chemistry, Göteborg University, Box 462, 405 30 Göteborg, Sweden

Received 27 September 2004; revised 10 February 2005

Available online 10 March 2005

## Abstract

Three-way decomposition is a very versatile analysis tool with applications in a variety of protein NMR fields. It has been used to extract structural data from 3D NOESYs, to determine relaxation rates in large proteins, to identify ligand binding in screening for lead compounds, and to complement non-uniformly recorded (sparse) spectra. All applications so far concerned experimental data sets; it thus remains to address questions of accuracy and robustness of the method using simulated data where the correct answer is known. Systematic tests are presented for relaxation and NOESY data sets. Mixtures of real and synthetic data are used to allow control of various parameters and comparisons with correct reference data, while working with input that is as realistic as possible. The influence of the following parameters is evaluated: signal-to-noise, overlap of signals and the use of a regularization procedure within the algorithm. The main criteria used for the evaluation are accuracy and precision. It is shown that deterioration of accuracy is indicated by internal checks such as decrease of precision. Both with relaxation data and when interpreting NOESY spectra, three-way decomposition exhibits a robust behavior in situations with severe signal overlap and/or poor signal-to-noise, e.g., by avoiding false positives in the NOE shapes of NOESY decompositions. As a complement to this study, three-way decomposition is compared to other methods that achieve the same type of results.

© 2005 Elsevier Inc. All rights reserved.

**Keywords:** MUNIN; <sup>15</sup>N-NOESY-HSQC; Three-way decomposition; Relaxation; Tikhonov regularization

## 1. Introduction

Three-way decomposition (TWD) is an analysis tool applicable to various types of NMR data sets [1]. Its underlying model is tightly coupled to multidimensional NMR spectroscopy and the basic decomposition formula can in fact be derived from a general pulse sequence describing a NMR experiment [2]. A consequence of this model is that the signals found in a NMR spectrum are grouped into “natural” subsets called components. Consider the example of a <sup>15</sup>N-NOESY-HSQC [1]. All NOE

signals in this spectrum are automatically, i.e., without specific user instruction, grouped so that each component exactly describes the NOEs to one amide proton HN. During this grouping, overlap among signals is very efficiently resolved. The decomposition avoids to a large extent the inclusion of noise and certain artifacts into the components, collecting them in a residual term. Similarly, in relaxation or ligand binding studies, which are typically based on a set of 2D <sup>15</sup>N-HSQC-type spectra, the decomposition is again performed for individual amide groups. This allows proper separation of overlapping signals, providing correct signal intensities. Another property of the TWD model is the absence of any requirement on signal shapes. Thus, TWD works equally well on signals with the appearance of Lorentzian or similar curves, of decaying sine functions describing FIDs, of exponential decays

\* Corresponding author. Fax: +46 31 773 3910.

E-mail address: [martin.billeter@chem.gu.se](mailto:martin.billeter@chem.gu.se) (M. Billeter).

<sup>1</sup> Present address: Clinical Genomics Centre, MBRC, 200 Elizabeth Street, Toronto, Ont., Canada M5G 2C4.

coupled to relaxation processes, or of any other shape. Furthermore, the line forms of the signals along the various dimensions in multidimensional spectra can be an arbitrary combination of different shape types. Finally, a recently introduced modification of the algorithm allows the application of TWD to sparsely recorded time-domain data sets in order to computationally fill the gaps and obtain full data sets [3].

TWD has been applied to a wide range of NMR data sets demonstrating its usefulness in various situations. From a  $^{15}\text{N}$ -NOESY–HSQC for the 128 residue long protein azurin a complete set of NOEs was extracted and it was shown that this set coincides closely with short proton-proton distances observed in the crystal structure of the same protein [4]. The TWD analysis thus provided a highly complete and reliable input for structure calculations. Thorough TWD analyses were performed for relaxation data including a demonstration of the method using  $T_{1\rho}$  data for azurin [5], and the determination of  $T_1$  relaxation times for all 341 assigned backbone amide groups of the 370 residue long protein maltose binding protein (MBP) [6]. Overlapping signals from up to three different amide groups were resolved. More recently, TWD has been applied in a routine manner in several relaxation studies [7–12]. A TWD application similar to the analysis of relaxation data is the use of decomposition for the screening of potential ligands for a target protein [13], since again a series of 2D  $^{15}\text{N}$ -HSQC-type spectra form the input. One difference is that the number of spectra is much larger, reaching several hundred 2D spectra; another is that instead of extracting a relaxation curve, TWD detects spectra where certain peaks have changed position. This screening is achieved in a single-step, avoiding peak picking with the necessity to characterize uncertain peaks. A rather different usage of TWD concerns the reconstruction of time-domain data sets from experiments performed in a sparse mode [14]. This allows time savings on expensive instruments by factors of three to five. It may be noted that for many of the above applications, there are few feasible alternatives to the analysis with TWD. This includes in particular proper separation of intensities for highly overlapped peaks in relaxation data or reliable reconstruction from sparse data in the case of NOESY spectra, which contain many peaks with widely varying intensities.

All the above TWD applications concern experimental data sets for the demonstration of the general applicability of TWD or simply for its routine use. Indications for the correctness of the TWD results could sometimes be derived by comparing these to independent data, e.g., a crystal structure in the case of extraction of distance information from NOESY spectra [4]. In other cases, e.g., with relaxation measurements, no truly independent data is available and the only comparison possible was by applying another method to the same NMR spectra [6]. For a thorough description of

any novel approach, its behavior should be characterized under controlled conditions using synthetic input with a priori known answers. Simulation studies have for example shown the advantages and limitations of methods such as linear prediction [15] or maximum entropy reconstruction [16].

The main goal of the present study is thus to systematically analyze the potential of TWD in terms of absolute accuracy, robustness in difficult situations and the extent to which internal error calculations, i.e., precision, correlate with accuracy. Parameters varied in the analysis are the extents of signal overlap and noise, and the use of a regularization procedure. Both extraction of data from a 3D spectrum, a  $^{15}\text{N}$ -NOESY–HSQC, and a set of 2D spectra,  $^{15}\text{N}$ -HSQC, are investigated. The resulting data provide guidelines for the reliability of TWD in difficult situations. While this goal can only be achieved by comparison of the TWD output with a priori known results, comparisons with alternative methods can provide additional information. Thus, for the case of relaxation data, comparisons are also performed with other tools that yield the same type of results, namely relaxation times; for NOESY decompositions no other tools have, to our knowledge, been described that yield shapes corresponding to the TWD output.

## 2. Methods

The application of TWD to a 3D NMR data set is based on the model assumption that this data set can be approximated by a sum of components [1]. Each component is a 3D entity of the same size as the original data set, but containing only a subset of the signals. Components are in turn defined as direct products of three one-dimensional vectors called shapes; because all shapes are normalized, their direct product is further multiplied by an amplitude. Mathematically, this idea can be expressed as follows:

$$S_{ijk} \approx \sum_{m=1}^M a^m \cdot F1_i^m \cdot F2_j^m \cdot F3_k^m, \quad (1)$$

where  $S_{ijk}$  is a matrix element describing the experimental 3D NMR data set, and the indices  $i$ ,  $j$ , and  $k$  cover the entire spectrum  $\mathcal{S}$ . Each of the  $M$  terms in the sum represents one component defined by the normalized 1D shape vectors  $F1^m$ ,  $F2^m$ , and  $F3^m$  and the overall amplitude  $a^m$ . The TWD algorithm will find shape vectors  $F1^m$ ,  $F2^m$ , and  $F3^m$  and amplitudes  $a^m$  that optimally approximate a given experimental spectrum  $\mathcal{S}$  (thus the  $\approx$  sign). More specifically, it will minimize the penalty function defined by the following expression:

$$\min \sum_{ijk} \left| S_{ijk} - \sum_{n=1}^N a^n \cdot F1_i^n \cdot F2_j^n \cdot F3_k^n \right|^2 + \lambda \sum_{n=1}^N (a^n)^2, \quad (2)$$

where  $S_{ijk}$  and the shapes  $F1^n$ ,  $F2^n$ , and  $F3^n$  and amplitudes  $a^n$  have the same meaning as in Eq. (1). The number  $N$  is input to the algorithm and represents an estimate of the true number of components  $M$ . In most applications a good estimate is easy to find, e.g., for a  $^{15}\text{N}$ -NOESY-HSQC each amide group will define one component [1,4]. Over-estimating  $M$  by about 10% will not affect the performance of the algorithm. The new sum at the end of expression (2) is used for regularization, i.e., it will ensure that all component amplitudes are of comparable size. The parameter  $\lambda$  refers to the Tikhonov regularization factor [17].

The flowchart in Fig. 1 summarizes the procedure used for testing TWD implemented as the program MUNIN [4]. Various input spectra are created (left side of flowchart) and decomposed with MUNIN according to expression (2). Input spectra are constructed according to Eq. (1) using shapes that can be either synthesized or extracted from experimental spectra, and user-defined amplitudes. In addition, synthetic or real noise, the latter typically extracted from empty regions of experimental spectra, is added according to predefined signal-to-noise ratios (S/N). The output shapes from MUNIN can be compared to the input shapes for the determination of accuracy, and subjected to other analyses.

### 2.1. Simulations using relaxation data

A series of calculations were designed to analyze the influence of noise and overlap on the accuracy and precision of relaxation times extracted with TWD. Relaxation times are usually obtained from a set of 2D spectra, where the intensity of a peak with given frequencies along the two axes follows an exponential decay. Typically,  $^{15}\text{N}$ -HSQC spectra are used for this purpose [18]. Although TWD works with any line forms, ideal shapes were used to construct an input spectrum in order to control the extent of overlap between two peaks that describe two atom groups with different relaxation behavior and to determine the accuracy of the resulting relaxation times. Shapes along the two frequency dimensions (e.g., HN and  $^{15}\text{N}$ ) were defined as absorption Lorentzian line-shapes according to:

$$f(\Omega) = W \frac{\alpha}{1 + \alpha^2(\Omega - \Omega_0)^2}, \quad (3)$$

where  $W$  is a normalization factor, the inverse of  $\alpha$  describes the line width of a Lorentzian, and  $\Omega_0$  indicates the center of the peak. Shapes in the relaxation dimension followed a simple exponential function with a normalization factor  $A$  and a relaxation time  $T^0$ :

$$f(t) = Ae^{-t/T^0}. \quad (4)$$

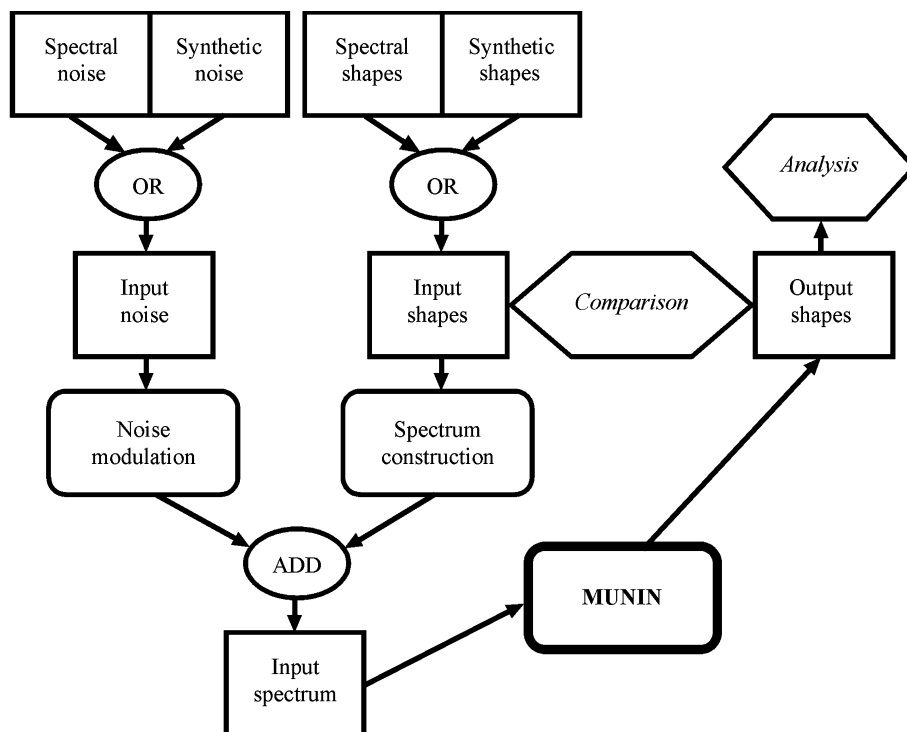


Fig. 1. Flowchart depicting the simulation procedure. The left part describes the preparation of the input spectra consisting of signal data either determined synthetically or constructed from experimental data, and noise data either obtained from experimental spectra or created with the help of a random number generator. Noise modulation includes scaling and reshuffling (see Section 2). The combination of signal and noise data forms the simulated input spectrum (bottom of flowchart). The right side of the figure presents the decomposition by MUNIN and the evaluation of the results. Shapes of the resulting components can be compared with the input shapes or they can be further analyzed, e.g., by deriving relaxation times.

Planes with frequency axes of the size  $35 \times 35$  points were used. Two peaks were placed onto the diagonal by applying Eq. (3) along each axis. The initial intensity ratio between the two peaks was set to 3:4. For both peaks,  $\alpha$  was set to  $1/3$  yielding a half width at half height of 3 points. Variations of  $\Omega_0$  allowed to move the two peaks along the diagonal towards each other and thus to control the overlap between the two peaks in the frequency dimensions. Using Eq. (1) and with a third shape defined by Eq. (4), an input data set with 12 planes was created for each calculation. Values of  $t$  in Eq. (4) were varied between 0 and 50 ms (with values of 0, 1, 2, 3, 5, 8, 12, 17, 24, 31, 40, and 50 ms), and the relaxation times  $T^0$  for the strong and weak peaks were set to 28 and 40 ms, respectively. Noise was then added at various, user-defined S/N ratios. For this purpose, a pool with 400 noise planes was created as follows to ensure spectrum-like noise. FIDs containing random Gaussian noise were created and processed in the same way as spectral FIDs, i.e., zero-filled, multiplied by a window function and Fourier transformed. Randomly selected noise planes from the pool were added to each of the 12 planes with the two peaks and this yielded the input spectra for MUNIN (Fig. 1). For each choice of overlap and S/N, calculations were repeated 50 times with different noise planes to obtain statistical results for accuracies and precision.

MUNIN was applied to various input spectra with different choices for signal overlap, S/N and the regularization factor  $\lambda$ . Relaxation times from the MUNIN output were extracted by a least square fit of the output shapes along the relaxation dimension utilizing a routine from the DASHA package [19]. The corresponding fitting errors  $\Delta err$  were estimated by covariance matrix calculation. For each choice of overlap and S/N, the average measurement precision over 50 runs with different noise is reported as a logarithm:

$$\log(\text{prec}) = \log(\langle \Delta err_i \rangle / T^0). \quad (5)$$

The brackets describe averaging over the  $i = 1, \dots, 50$  runs. Similarly, accuracy of the resulting relaxation times  $T_i$  from the 50 decompositions is reported with the following logarithm:

$$\log(\text{acc}) = \log \sqrt{\frac{\sum (T_i - T^0)^2}{50 \cdot T^{0^2}}}, \quad (6)$$

where  $T^0$  is the true value, which is used for normalization also in the definition of a precision entity (Eq. (5)) to allow better comparisons between the numbers obtained from Eqs. (5) and (6).

## 2.2. Comparisons with other methods

Complementing the above analysis on the extraction of relaxation data using TWD, comparisons were made

with other methods using the same simulated spectra. A first comparison is with a simple (but widely used) approach, in the following referred to as conventional procedure. In the input spectra, i.e., after addition of noise to the constructed spectra, intensities at the locations of the peak center in the input shapes were evaluated and an exponential function was fitted to the resulting curve using the same routine as above from the DASHA package. Again, results were obtained by averaging over 50 runs with different noise and overlap according to Eqs. (5) and (6). A second comparison is with a tool from the software package nmrPipe, the routine “nlinLS” [20]. As input, this method requires start values for the peak positions, the line widths and the relaxation times as well as a choice between Gaussian and Lorentzian line forms. Again, calculations of 50 runs for each combination of overlap and noise were performed with the following start values derived from the simulation input (i.e., from the correct values): Initial peak positions were shifted towards each other by 25% of their correct distance, line widths were set to those used for simulating the input, the initial relaxation times for both peaks were set to the average of the two true relaxation times, and both Gaussian and Lorentzian line forms were tested.

## 2.3. Simulations using NOESY data

Simulations on NOESY spectra were designed to test the ability of TWD to reproduce the input components for various choices of S/N and signal overlap. For this purpose, the resulting shapes were compared to the input shapes for all three dimensions. All input data were constructed from an experimental spectrum, a  $^{15}\text{N}$ -NOESY-HSQC recorded for the protein azurin. This spectrum was Fourier transformed in the HN and  $\text{H}_{\text{NOE}}$  dimensions, and a region covering the HN frequency interval from 8.64 to 8.94 ppm was extracted. This region was decomposed with MUNIN into 20 components as described earlier [4], and two of the resulting components were selected. Each selected component consisted of 21 points covering 0.3 ppm along the HN dimension, 470 points covering the entire interval of 13.43 ppm along the  $\text{H}_{\text{NOE}}$  dimension and 80 time domain points in the non-transformed  $^{15}\text{N}$  dimension. In order to vary the degree of overlap, one of the components was kept fixed, while the other one was modified by shifting its HN frequency shape. The  $\text{H}_{\text{NOE}}$  shape had to be shifted simultaneously, in order to preserve the diagonal peak as such. In addition, a section of the same size as the selected region was chosen from a part of the spectrum that contained no signals to provide the basic noise for the simulation. This noise was scaled to achieve various signal-to-noise ratios. For this test on NOESY data, TWD calculations were performed with a factor  $\lambda$  of zero in expression (2), i.e., without Tikhonov

nov regularization. Because the input shapes originated from an experimental spectrum, no absolute discrimination of weak peaks and noise is possible. Direct comparisons of the input and output shapes, e.g., with dot products between corresponding shapes as used earlier [1,4], tend to be insensitive to variations of small peaks. Therefore, the analysis was based on a one-dimensional peak picker [4], which was used with default parameters. While the overall result of any peak picker depends on its inherent properties and the run-time parameters chosen, a comparison of input and output is nonetheless informative provided that the same peak picker is used with the same run-time parameters. However, comparisons with alternative methods based on other peak pickers (e.g., peak picking in three dimensions) would mainly reflect inherent features of different peak pickers.

### 3. Results

#### 3.1. Three-way decomposition of relaxation data

In a first test the simulation protocol of Fig. 1 was applied to relaxation data. Input spectra consisting of 12 planes were constructed as described in Section 2, and the following three parameters were varied. Overlap between two peaks was changed from an initial separation in each dimension of the peak centers of 8 spectral points down to 2 points in steps of 2 points; this can be compared to a half width at half height of 3 spectral points in each frequency dimension. The extent of overlap for these four simulations is also illustrated in Fig. 2, which shows the first planes with  $t = 0$  (Eq. (4)). S/N was varied in 15 steps by scaling the input S/N ratios between 10 and 100, calculated for the stronger peak in the first of the 12 planes. Finally, calculations were performed with a regularization factor  $\lambda$  of 0.001, or without regularization ( $\lambda = 0$ ) [14]. For each combination of overlap, S/N and regularization factor, 50 individual calculations were performed by noise reshuffling as described in Section 2 to ensure significance of the statistical analysis.

Fig. 3 reports logarithms of accuracy according to Eq. (6) as a function of S/N for the MUNIN calculations without (thick solid lines) and with regularization (thin solid lines). Results for the different overlap situations described above and illustrated in Fig. 2 are reported for both peaks (left and right panels, respectively) in different pairs of panels. Fig. 4 provides precision values as determined with Eq. (5) for the same MUNIN calculations. The difference between using regularization or not is marginal for cases with small overlap. However, regularization may help in difficult situations with strong overlap and low S/N. With high S/N, regularization rather represents an unnecessary distortion of the model. The precision should follow

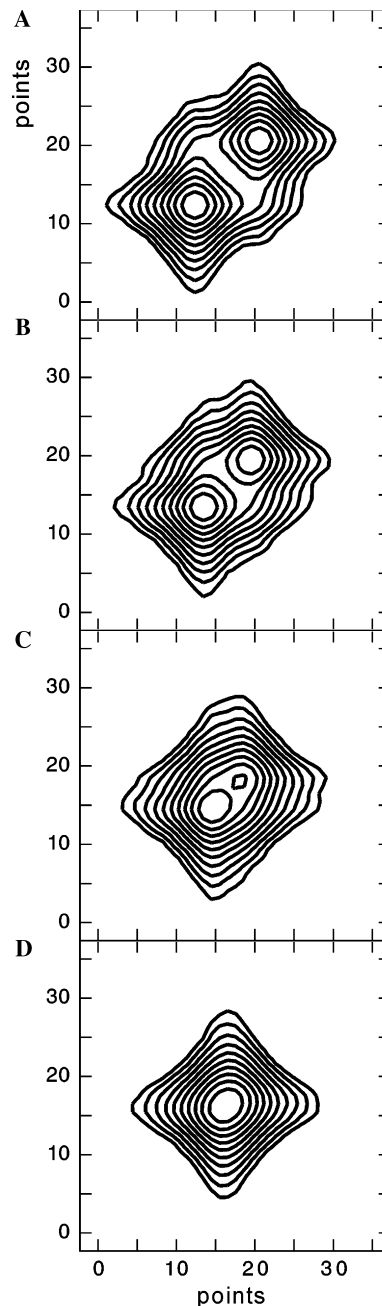


Fig. 2. First planes of the input spectra for simulations of relaxation data sets. (A–D) The extent of overlap of two peaks for four different simulations. The total input for each simulation consists of 12 planes.

the accuracy; otherwise one may for example obtain a too optimistic impression of the reliability of the resulting relaxation times. For TWD calculations without regularization on cases with limited overlap, there is good correspondence between precision and accuracy and the difference is not significant (Figs. 3 and 4, A–C). With high overlap, the over-estimation of the quality of the data becomes more pronounced and one should be more careful with the interpretation. With high S/N, the above mentioned model distortions caused by

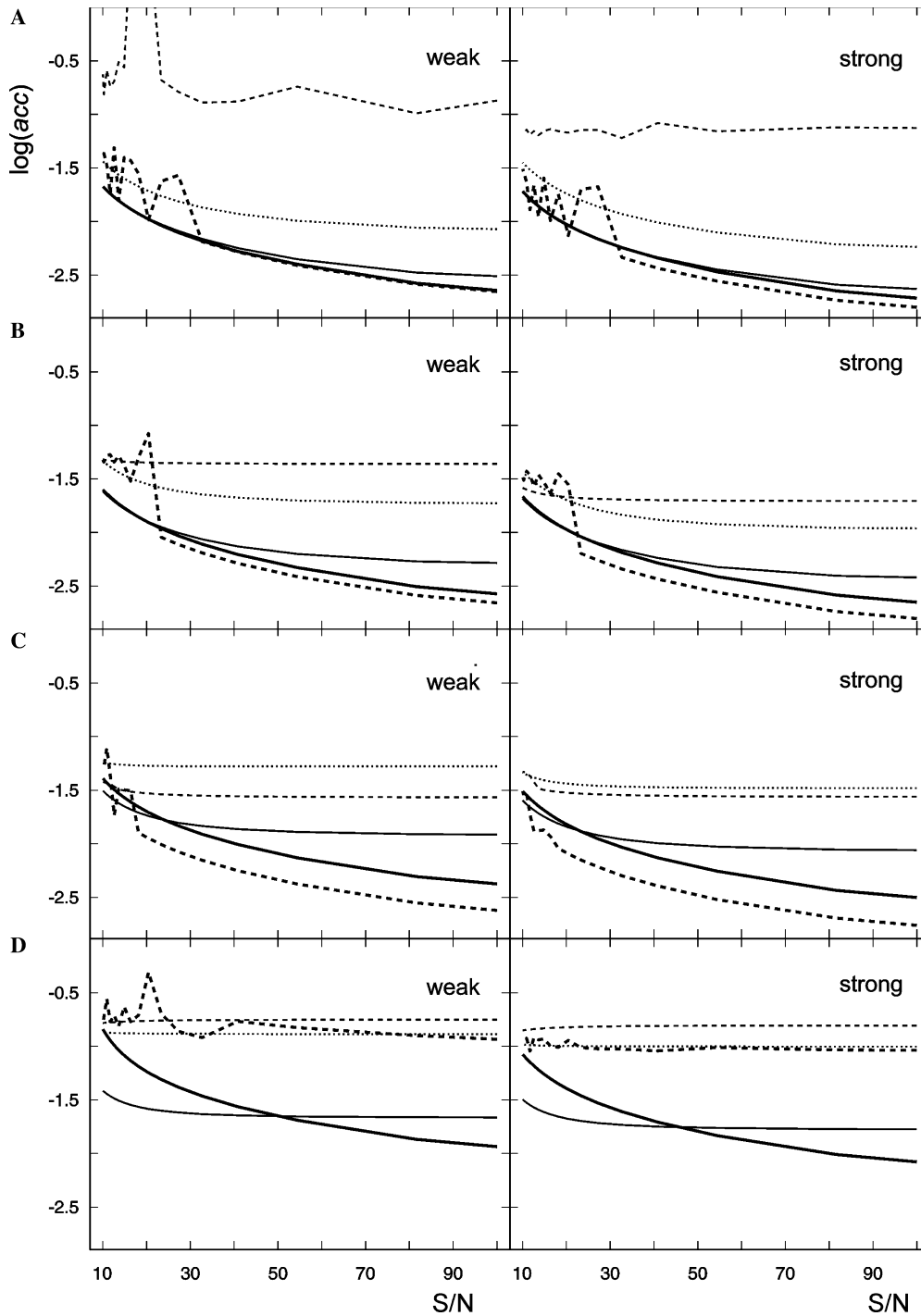


Fig. 3. Accuracy according to Eq. (6) as a function of S/N for various overlap situations. The panels on the left side, labeled “weak” correspond to results for the weak peak in the overlap situations A–D in Fig. 2. Similarly, the panels on the right side, labeled “strong” correspond to results for the strong peak in the overlap situations A–D in Fig. 2. The thin and thick solid lines correspond to MUNIN calculations with and without regularization, respectively. The thin and thick dashed lines correspond to calculations with the “nlinLS” routine from the nmrPipe package using Gaussian and Lorentzian line forms, respectively, in both frequency dimensions [20]. The dotted line stands for the use of a conventional procedure.

regularization effect mostly the accuracy, thus increasing the difference between accuracy and precision.

Accuracies and precisions of the TWD results are also compared to those of other methods. The conventional method, based on peak intensities alone, gives a

stable impression with little dependence of the accuracy on S/N (dotted lines in Figs. 3 and 4). More worrying here is the prominent discrepancies between accuracy and precision, where with maximal overlap the precision is among the best whereas the accuracy is among the

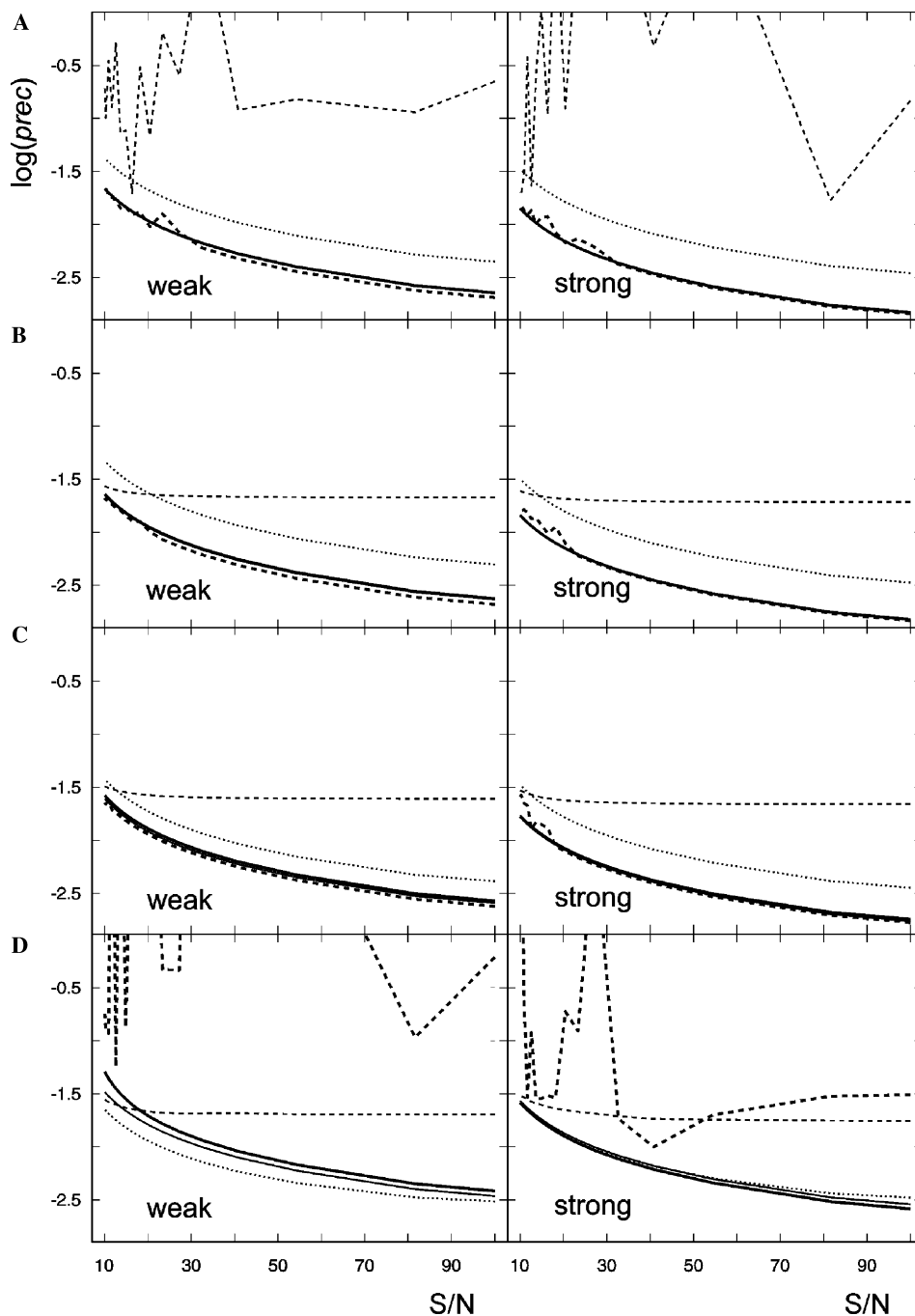


Fig. 4. Precision according to Eq. (5) as a function of  $S/N$  for various overlap situations. Panel labeling and definition of lines are as in Fig. 3.

worst (Figs. 3D and 4D). In the case of strong overlap, the approach implemented in the routine “nlinLS” of the nmrPipe [20] package represents a more realistic tool (dashed lines in Figs. 3 and 4). Indeed, when assuming pure Lorentzian lines (thick dotted lines), this method follows the lines for MUNIN both for accuracy and precision; with moderate to sizeable overlap (B- and C-panels) it exceeds the MUNIN accuracy. However, it breaks down with strong overlap (D-panels). Furthermore, it shows a somewhat erratic behavior with low  $S/N$ ,

caused by poor convergence of some of the 50 runs with varied noise (see Section 2). In contrast to the simulated input used here, experimental peaks are not strictly Lorentzian, but are often approximated by a mixture of Lorentzian and Gaussian curves [21,22]. Thus, the striking differences observed in Figs. 3 and 4 between the use of Lorentzian and Gaussian curves (thick and thin dotted lines) indicate another potential instability.

Fig. 5 further analyzes the calculations with varying overlap using a fixed  $S/N$  ratio of 32 (this value corre-

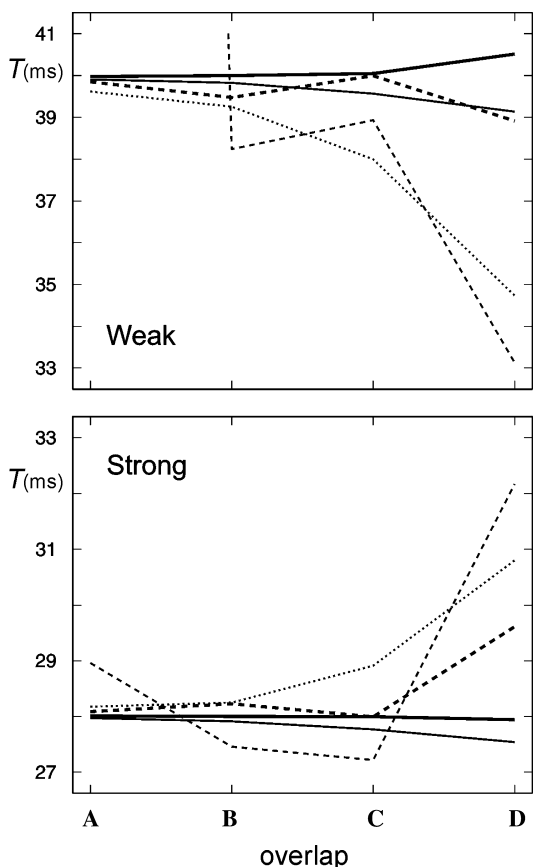


Fig. 5. Plots of average values of relaxation times from 50 runs with different noise for varying signal overlap as defined in Fig. 2. S/N was set to 32. The top and the bottom panel correspond to the weak and strong peaks, respectively. Overlap is indicated along the horizontal axis by referring to the different panels in Fig. 2. Similarly to Figs. 3 and 4, thin and thick solid lines correspond to MUNIN calculations with and without regularization, thin and thick dashed lines to the use of “nlinLS” with Gaussian and Lorentzian line forms, and the dotted line to the use of a conventional procedure.

sponds to the S/N observed in an application to the protein MBP; see Section 4) by plotting the average over 50 runs of the relaxation times,  $\langle T \rangle$ , as a function of overlap for both peaks. Note that this average differs from the one used in Fig. 3 where squares of deviations from the correct value  $T^0$  are averaged. The MUNIN calculations yield values close to the correct ones of 28 and 40 ms over the entire range of signal overlap (solid lines). While the calculations without regularization appear to perform better than those using regularization, the differences lie within the error range. Again, comparisons to other methods may reveal intrinsic features of individual approaches. As expected (and indicated in [6]) the conventional method shows strong averaging of the two relaxation times with increasing overlap (dotted line). The “nlinLS” routine from nmrPipe [20] results in some averaging for the strongest overlap situation when a Lorentzian line form is used, while with a Gaussian line form a rather instable behavior is observed (dashed lines).

### 3.2. Three-way decomposition of NOESY data

The simulation protocol of Fig. 1 was also used to investigate the influence of overlap and S/N on the extraction of data from a 3D NOESY. Two components obtained by applying MUNIN to an experimental  $^{15}\text{N}$ -NOESY-HSQC of azurin describe the NOEs involving two HNs. Fig. 6A shows the shapes along the  $H_{\text{NOE}}$  dimension for the two components. The shapes along the  $^{15}\text{N}$  dimension were kept fixed and in time-domain. For illustration, the extent of overlap in this dimension after Fourier transform is shown in Fig. 6B. The shapes in the two proton dimensions for the second component (thick lines in Figs. 6 and 7) were shifted simultaneously to preserve the diagonal peak, and different overlap situations ranging from 0.0 to 0.056 ppm were constructed. Noise taken from an empty region of the same spectrum was scaled to achieve user-defined levels of S/N and added to the input components as described in Section 2. S/N is defined here as the average intensities of the diagonal peaks of the two components divided by the variance of the noise. Note that the intensity of a cross-peak seldom exceeds 5% of that of the diagonal peak, and that therefore the signal-to-noise ratio for cross-peaks is about 20 times smaller than the S/N reported. The S/N in the original, unmodified  $^{15}\text{N}$ -NOESY-HSQC was calculated as 326, and the S/N was varied between 41 and 326.

Table 1 summarizes results for calculations for five different degrees of overlap and varying S/N, and selected situations indicated in the table are illustrated in Fig. 7. It should be noted here that the primary output of MUNIN are shapes; their interpretation by peak picking yields a result that also depends on the characteristic of the peak picker used. This is in particular true for small peaks that naturally disappear in the noise when S/N is reduced. However, an important question concerns the possible appearance of false positives, i.e., new peaks created by distortions of the shapes in the  $H_{\text{NOE}}$  dimension. Thus, both the input and output shapes were subjected to a one-dimensional peak picker [4] along this dimension. The two numbers listed in Table 1 for each combination of overlap and S/N report the false positives detected for the two components, i.e., peaks picked in the output shapes but not in the input shape (the latter are indicated in Fig. 6A by arrows). The MUNIN approach is thus successful even with full overlap provided that S/N is high enough. With increasing separation of the two components, proper shapes are extracted even with very low S/N. When both strong overlap is present and S/N is low, the method fails by being no longer able to correctly separate the two components. This can be seen from the fact that false positives in one component correspond to correct peaks from the other component. Therefore, one output peaks component tends to describe



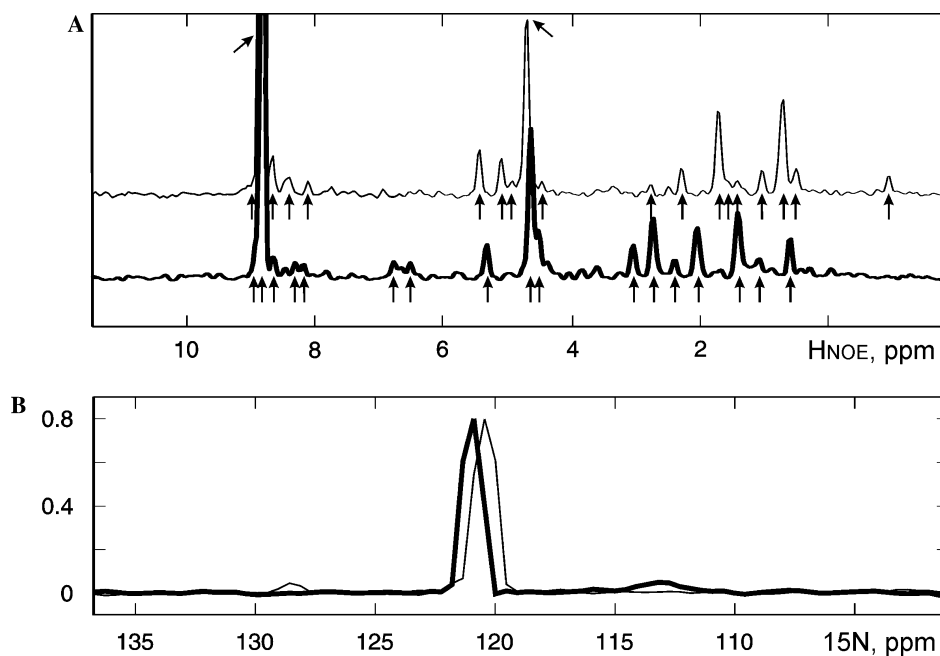


Fig. 6. Two components used for simulations of the decomposition of a  $^{15}\text{N}$ -HSQC–NOESY. (A) Shapes along the  $H_{\text{NOE}}$  dimension used as input for the simulation. Peaks detected by peak picking are indicated by arrows. (B) Input shapes along the  $^{15}\text{N}$  dimension after Fourier transform. Thin lines in both panels correspond to the first component, which was kept fixed, and thick lines to the second, moving component.

both input components, while the other output component contains merely noise.

Fig. 7 illustrates three critical situations with views of the input and output. Panels A1 and A2 characterize the case of complete overlap and  $S/N = 163.1$  (labeled “(A)” in Table 1), panels B1 and B2 the case with overlap of 0.014 ppm and  $S/N = 65.2$  (labeled “(B)” in Table 1), and panels C1, C2, and C3 the case with overlap of 0.056 ppm and  $S/N = 40.8$  (labeled “(C)” in Table 1). Panels A1, B1, and C1 (the latter is plotted at higher contour levels than the former two) of Fig. 7 show the input spectrum for the MUNIN decomposition, while panels A2, B2, and C2 of the same figure show the resulting shapes along the  $H_{\text{NOE}}$  dimension. When compared to the shapes used for construction of the input spectra shown in Fig. 6A, the output shapes in panels A2 and B2 contain all but the very smallest peaks, and no artifact peaks are observed; in the panel C2 only the large peaks can be detected in the output due to the decreased  $S/N$ , but still no false positives are observed. Panel C3 of Fig. 7 is a reconstruction of the spectrum according to Eq. (1) using the output shapes of panel C2 and the corresponding output shapes in the other two dimensions. It shows that, even with the high level of noise present in the input, the larger peaks can be detected. Comparison of panel C3 with C1 (plotted at the same contour levels) illustrates that the decomposition removes a significant part of the noise, some of which with relatively strong intensity.

#### 4. Discussion and conclusions

The above systematic analyses of TWD calculations applied to simulated input describing NOESY and relaxation data provide a basis for the application of TWD to various experimental data sets. Questions of accuracy, robustness, and the relation between accuracy and precision were addressed. In earlier applications of MUNIN to relaxation data, reliability of the results was indicated by internal criteria such as the residual of the optimization (the value of expression (2) at the end of the minimization), variation of the number of components used ( $N$  in expression (2)), and comparison to other methods [5,6]. In the analysis of relaxation times for all assigned backbone amide groups of the protein MBP, it was proposed that TWD performs a proper separation of strongly overlapped peaks, resulting in accurate relaxation times. This explanation is confirmed by the present calculations, in particular by Fig. 5 that is based on the same  $S/N$ , 32, as observed in the MBP spectra. Similarly, the power of TWD to separate NOEs in a  $^{15}\text{N}$ -NOESY–HSQC observed for two amide protons with extensive overlap in both the HN and  $^{15}\text{N}$  dimensions could be demonstrated quantitatively. False peaks that would yield wrong distance restraints in a structural study appear only with both strong overlap and poor  $S/N$ . Furthermore, the first false positives to appear are actually correct peaks but from the overlapped component, meaning that a NOE is indeed present but it would receive a wrong assignment. For both

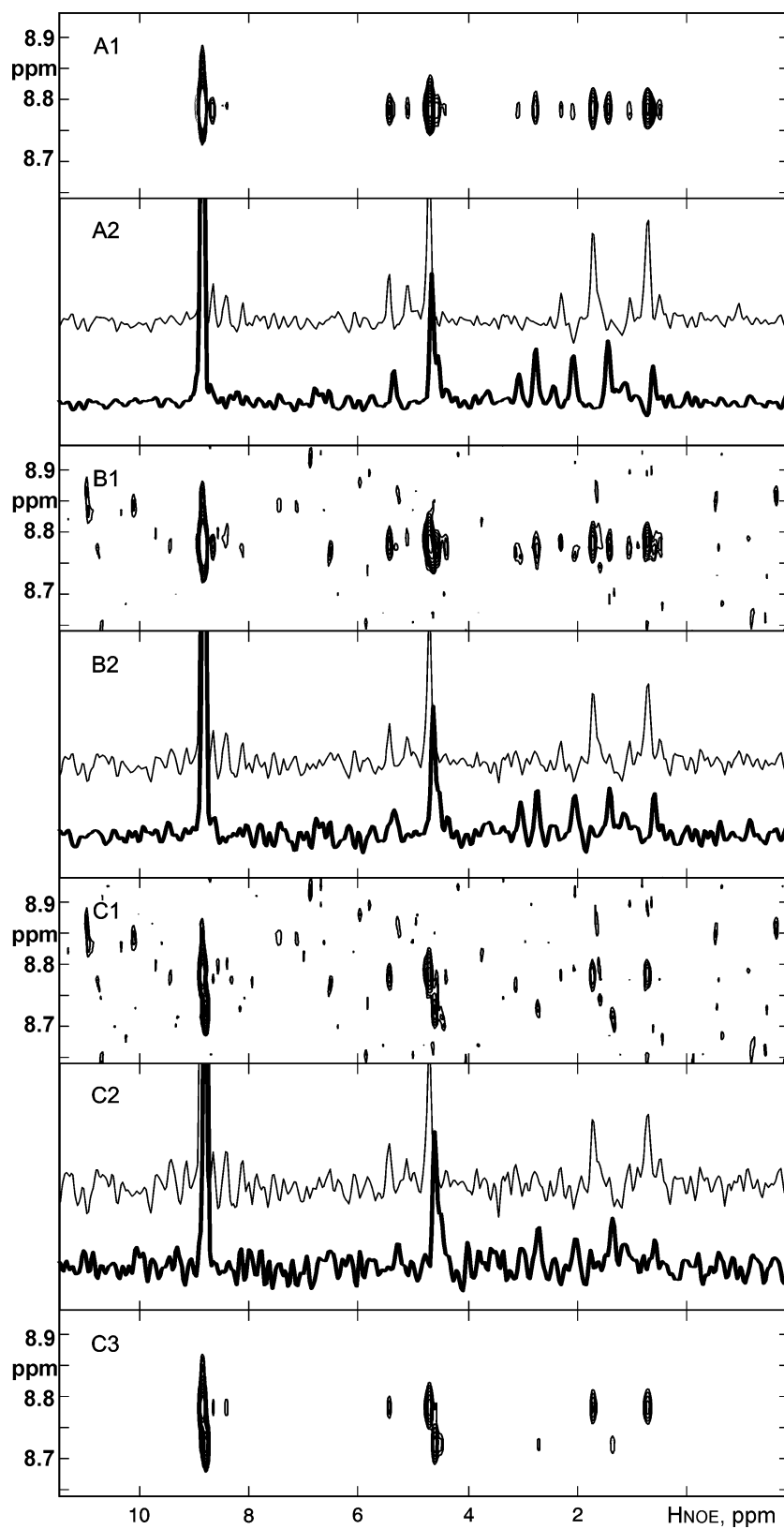


Fig. 7. Input and output of selected MUNIN decompositions for a  $^{15}\text{N}$ -HSQC-NOESY. The letters in the panel identifications (A, B or C) correspond to overlap and S/N situations indicated in Table 1 by the same letters (see footnote d to this table). Panels with identifications containing the number 1 provide the input spectra for the simulation. Panels with the number 2 show the output shapes along the  $H_{\text{NOE}}$  dimension; these can be compared to Fig. 6A. In addition, the panel with identification C3 shows a reconstruction of the spectrum from the output shapes according to Eq. (1). Panels C1 and C3 are plotted at contour levels 1.5 times higher than panels A1 and B1.

Table 1  
Performance of TWD when applied to 3D NOESY data for various choices of overlap and S/N<sup>a</sup>

S/N <sup>b</sup>	Peak separation (ppm) <sup>c</sup>				
	0.0	0.014	0.028	0.042	0.056
326.2	0 0	0 0	0 0	0 0	0 0
163.1	0 0 (A) <sup>d</sup>	0 0	0 0	0 0	0 0
108.7	(4) 4	0 0	0 0	0 0	0 0
81.5	(3) 4	0 0	0 0	0 0	0 0
65.2	(3) 4	0 0 (B) <sup>d</sup>	0 0	0 0	0 0
54.4	(3) 4	(3) 4	(3) 3	(2) 3	0 0
46.6	(3) 4	(3) 4	(2) 3	(1) 3	0 0
40.8	(2) 4	(2) 4	(1) 3	(1) 3	0 0 (C) <sup>d</sup>

<sup>a</sup> The two numbers in each entry indicate for both components how many of the peaks detected in the output shapes along the H<sub>NOE</sub> dimension have no matches in the input shapes. Numbers in parentheses represent peaks whose matches in the input shapes were found in the other component (because of mixing).

<sup>b</sup> S/N was calculated based on the intensity of the diagonal peak after Fourier transformation in the <sup>15</sup>N dimension.

<sup>c</sup> Overlap is indicated by the separation in ppm of the HN frequencies of the two components.

<sup>d</sup> Calculations for the entries with labels (A), (B), and (C) are further illustrated in Fig. 7.

types of applications, also the limitations set by strong overlap and/or poor S/N are characterized.

Another parameter that was investigated is the use of regularization according to Tikhonov [17] (last term in expression (2)). This regularization was shown earlier to speed up the optimization process, in particular in difficult situations [6,14]. However, regularization is mostly helpful in decompositions involving many components. In the present simulations with only two components involved, the motivation of using regularization is limited. We have tested its influence for relaxation data because this part of the simulation study provided more information in terms of accuracy and robustness. Our conclusion is that regularization with a small value of  $\lambda$ , for example half the value used here, has no significant negative consequences. With very high S/N, regularization is less necessary, but in many experimental situations and when many components are present, the earlier described advantages may well motivate some regularization.

Besides the investigation on the absolute reliability of TWD by means of accuracy and robustness, the MUNIN results regarding relaxation data were also compared to results obtained with an alternative method. While the use of a routine from the nmrPipe package [20] provided in many cases a similar or even better accuracy than MUNIN, the latter proved more robust with low S/N and/or strong overlap. Furthermore, the results for the former are strongly influenced by the assumptions made on line forms (Lorentzian or Gaussian), which are known to be problematic. For the present simulations, assuming Lorentzian line forms is ideal

as this corresponds to the simulated data, but in a real case mixtures of Lorentzian or Gaussian expressions have been shown to work best [21,22].

Another issue that has been discussed in the context of TWD is the observation of mixing between components [1,4]. This occurs when the shapes from two components along one dimension are very similar. Mixing implies that the shapes in the other two dimensions contain features of both true components. These true components can be obtained as linear combinations of the output components of TWD, a process referred to as demixing. In the simulations presented here, we chose to control peak overlap by moving the input components in a diagonal manner (see for example Fig. 2). This avoids mixing as long as possible when signals approach each other. While demixing could have improved the final results in a few cases, we chose to not apply this procedure and rather report a situation with significant mixing as a failure (e.g., in Table 1).

In conclusion, TWD is able to extract accurate data from NMR spectra and to escape interference by noise. The method is robust, avoiding false positives and providing in general trustworthy precision data.

## Acknowledgments

This work was supported by grants from the Swedish Foundation for Strategic Research (A3 04:160d) and the Swedish Research Council (621-2003-4048).

## References

- [1] V.Y. Orekhov, I.V. Ibraghimov, M. Billeter, MUNIN: a new approach to multi-dimensional NMR spectra interpretation, *J. Biomol. NMR* 20 (2001) 49–60.
- [2] M. Billeter, V.Y. Orekhov, Three-way decomposition and nuclear magnetic resonance, in: P.M.A. Sliot, D. Abramson, A.V. Bogdanov, J.J. Dongarra, A.Y. Zomaya, Y.E. Gorbachev (Eds.), *Computational Science—ICCS 2003*, Springer, Berlin, 2003, pp. 15–24.
- [3] I.V. Ibraghimov, Application of the three-way decomposition for matrix compression, *Numer. Linear Algebra Appl.* 9 (2002) 551–565.
- [4] A. Gutmanas, P. Jarvoll, V.Y. Orekhov, M. Billeter, Three-way decomposition of a complete 3D <sup>15</sup>N-NOESY-HSQC, *J. Biomol. NMR* 24 (2002) 191–201.
- [5] D.M. Korzhnev, I.V. Ibraghimov, M. Billeter, V.Y. Orekhov, MUNIN: application of three-way decomposition to the analysis of heteronuclear NMR relaxation data, *J. Biomol. NMR* 21 (2001) 263–268.
- [6] A. Gutmanas, L. Tu, V.Y. Orekhov, M. Billeter, Accurate relaxation parameters for large proteins, *J. Magn. Reson.* 167 (2004) 107–113.
- [7] D.M. Korzhnev, K. Kloiber, V. Kanelis, V. Tugarinov, L.E. Kay, Probing slow dynamics in high molecular weight proteins by methyl-TROSY NMR spectroscopy: application to a 723-residue enzyme, *J. Am. Chem. Soc.* 126 (2004) 3964–3973.
- [8] D.M. Korzhnev, K. Kloiber, L.E. Kay, Multiple-quantum relaxation dispersion NMR spectroscopy probing millisecond time-

- scale dynamics in proteins: theory and application, *J. Am. Chem. Soc.* 126 (2004) 7320–7329.
- [9] D.M. Korzhnev, V.Y. Orekhov, F.W. Dahlquist, L.E. Kay, Off-resonance R-1 rho relaxation outside of the fast exchange limit: an experimental study of a cavity mutant of T4 lysozyme, *J. Biomol. NMR* 26 (2003) 39–48.
- [10] D.M. Korzhnev, B.G. Karlsson, V.Y. Orekhov, M. Billeter, NMR detection of multiple transitions to low-populated states in azurin, *Protein Sci.* 12 (2003) 56–65.
- [11] A. Zhuravleva, D.M. Korzhnev, E. Kupce, A.S. Arseniev, M. Billeter, V. Orekhov, Gated electron transfers and electron pathways in azurin: a NMR dynamic study at multiple fields and temperatures, *J. Mol. Biol.* 342 (2004) 1599–1611.
- [12] V. Tugarinov, W.-Y. Choy, V.Yu. Orekhov, L.E. Kay, *Proc. Natl. Acad. Sci. USA* 102 (2005) 622–627.
- [13] C.S. Damberg, V.Y. Orekhov, M. Billeter, Automated analysis of large sets of heteronuclear correlation spectra in NMR-based drug discovery, *J. Med. Chem.* 45 (2002) 5649–5654.
- [14] V.Y. Orekhov, I.V. Ibraghimov, M. Billeter, Optimizing resolution in multidimensional NMR by three-way decomposition, *J. Biomol. NMR* 27 (2003) 165–173.
- [15] P. Koehl, Linear prediction spectral analysis of NMR data, *Prog. Nucl. Magn. Reson. Spectrosc.* 34 (1999) 257–299.
- [16] A.S. Stern, K.B. Li, J.C. Hoch, Modern spectrum analysis in multidimensional NMR spectroscopy: comparison of linear-prediction extrapolation and maximum-entropy reconstruction, *J. Am. Chem. Soc.* 124 (2002) 1982–1993.
- [17] A.N. Tikhonov, A.A. Samarskij, *Equations of Mathematical Physics*, Dover, New York, 1990.
- [18] A.G. Palmer, J. Williams, A. McDermott, Nuclear magnetic resonance studies of biopolymer dynamics, *J. Phys. Chem.* 100 (1996) 13293–13310.
- [19] V.Y. Orekhov, D.E. Nolde, A.P. Golovanov, D.M. Korzhnev, A.S. Arseniev, Processing of heteronuclear NMR relaxation data with the new software DASHA, *Appl. Magn. Reson.* 9 (1995) 581–588.
- [20] F. Delaglio, S. Grzesiek, G.W. Vuister, G. Zhu, J. Pfeifer, A. Bax, Nmrpipe—a multidimensional spectral processing system based on Unix pipes, *J. Biomol. NMR* 6 (1995) 277–293.
- [21] W. Denk, R. Baumann, G. Wagner, Quantitative evaluation of cross-peak intensities by projection of two-dimensional NOE spectra on a linear-space spanned by a set of reference resonance lines, *J. Magn. Reson.* 67 (1986) 386–390.
- [22] R. Koradi, M. Billeter, M. Engeli, P. Güntert, K. Wüthrich, Automated peak picking and peak integration in macromolecular NMR spectra using AUTOPSY, *J. Magn. Reson.* 135 (1998) 288–297.