

Automated projection spectroscopy (APSY)

Sebastian Hiller*, Francesco Fiorito*, Kurt Wüthrich†, and Gerhard Wider

Institut für Molekularbiologie und Biophysik, Eidgenössische Technische Hochschule Zürich, CH-8093 Zürich, Switzerland

Contributed by Kurt Wüthrich, June 9, 2005

This work presents the automated projection spectroscopy (APSY) method for the recording of discrete sets of j projections from N -dimensional ($N \geq 3$) NMR experiments at operator-selected projection angles and automatic identification of the correlation cross peaks. The result from APSY is the fully automated generation of the complete or nearly complete peak list for the N -dimensional NMR spectrum from a geometric analysis of the j experimentally recorded, low-dimensional projections. In the present implementation of APSY, two-dimensional projections of the N -dimensional spectrum are recorded by using techniques developed for projection–reconstruction spectroscopy [Kupče, E. & Freeman, R. (2004) *J. Am. Chem. Soc.* 126, 6429–6440]. All projections are peak-picked with the available automated routine ATNOS. The previously undescribed algorithm GAPRO (geometric analysis of projections) uses vector algebra to identify subgroups of peaks in different projections that arise from the same resonance in the N -dimensional spectrum, and from these subgroups it calculates the peak positions in the N -dimensional frequency space. Unambiguous identification thus can be achieved for all cross peaks that are not overlapped with other peaks in at least one of the N dimensions. Because of the correlation between the positions of corresponding peaks in multiple projections, uncorrelated noise is efficiently suppressed, so that APSY should be quite widely applicable for correlation spectra of biological macromolecules, which have intrinsically low peak density in the N -dimensional spectral space.

GAPRO | multidimensional | NMR | peak picking

In NMR studies of biological macromolecules in solution (1–4), multidimensional NMR data are commonly acquired by sampling the time domain in all dimensions equidistantly at a resolution adjusted to the populated spectral regions (5). With recent advances in sensitivity, due to high field strengths and/or cryogenic detection devices, the time required to explore the time domain in the conventional way typically exceeds by far the time needed for sensitivity considerations, so that the desired resolution in the indirect dimensions determines the duration of the experiment. In this situation of the “sampling limit,” which is common in 3D and higher-dimensional experiments with small and medium-size proteins (6), the desired chemical shift information has been collected by using “unconventional” experimental schemes, such as nonuniform sampling of the time domain (7, 8) or combination of two or more indirect dimensions (9, 10).

The concept of combining indirect dimensions has led to reduced-dimensionality experiments (9) and G-matrix Fourier transform NMR (11, 12). In G-matrix Fourier transform NMR, several evolution periods of a multidimensional NMR experiment are combined, the data are processed by using a G-matrix, and the resulting set of spectra is analyzed jointly to identify the peaks that arise from the same spin system and to calculate their resonance frequencies (11). In another approach, projection–reconstruction NMR (13–16), the projection–cross-section theorem (17, 18) is combined with reconstruction methods from imaging techniques (19, 20). In particular, a scheme for quadrature detection along tilted planes in the time domain allows the direct recording of orthogonal projections of any multidimensional experiment at arbitrary projection angles (15). In projec-

tion–reconstruction NMR, the full multidimensional spectrum then is reconstructed from the projections of the multidimensional spectral data (13–16).

The analysis of complex NMR spectra typically involves intensive human interaction, and automation of NMR spectroscopy with macromolecules is still in development. Thereby the distinction of real peaks from random noise and spectral artifacts as well as peak overlap represent major challenges (21–23). On grounds of principle, automated analysis benefits from higher dimensionality of the spectra (24, 25), because the peaks are then more widely separated, and hence peak overlap is substantially reduced (Fig. 1).

In the present work, we combine technologies to record projections of high-dimensional NMR experiments described by Kupče and Freeman (15) and automated peak-picking using a scheme of Herrmann *et al.* (22) with a previously undescribed algorithm, GAPRO (geometric analysis of projections). Based on geometrical considerations, GAPRO identifies peaks in the projections that arise from the same resonance in the N -dimensional frequency space and subsequently calculates the resonance frequencies in the N -dimensional spectrum without ever considering the high-dimensional data set itself. This automated analysis of projected spectra, APSY (automated projection spectroscopy), yields a peak list of the original multidimensional experiment without any human interaction. In the following sections, the foundations of APSY are introduced, and characteristic properties of APSY are discussed. Two examples of APSY are a 4D HNCOCOA experiment (26, 27) with the 63-residue protein 434-repressor(1–63) (28) and a 5D HACACONH experiment (11) with the 116-residue protein TM1290 (29).

Theoretical Background

Recording of Projection Spectra. The projection–cross-section theorem by Bracewell (17), which was introduced into NMR by Nagayama *et al.* (18), states that an m -dimensional ($m < N$) cross section, $c_m(t)$, through N -dimensional time domain data is related by an m -dimensional Fourier transformation, \mathcal{F}_t , and its inverse, \mathcal{F}_ω , to an m -dimensional orthogonal projection of the N -dimensional NMR spectrum, $P_m(\omega)$, in the frequency domain. Thereby, $P_m(\omega)$ and $c_m(t)$ are oriented by the same angles with regard to their corresponding coordinate systems (Fig. 2). On this basis, Kupče and Freeman proposed to record projections $P_m(\omega)$ by sampling the corresponding time domain data, $c_m(t)$, along a straight line (dashed line in Fig. 2). Quadrature detection is obtained from corresponding positive and negative projection angles for the subsequent hypercomplex Fourier transformation (15).

Projections of Cross Peaks. We describe here 2D projections, $P_2(\omega)$, of an N -dimensional spectrum ($N > 2$). $P_2(\omega)$ represents spectral data in a 2D plane, which is spanned by an indirect dimension with unit vector \hat{p}_1 , and the direct dimension, with \hat{p}_2 .

Abbreviations: APSY, automated projection spectroscopy; GAPRO, geometric analysis of projections.

*S.H. and F.F. contributed equally to the work.

†To whom correspondence should be addressed. E-mail: wuthrich@mol.biol.ethz.ch.

© 2005 by The National Academy of Sciences of the USA

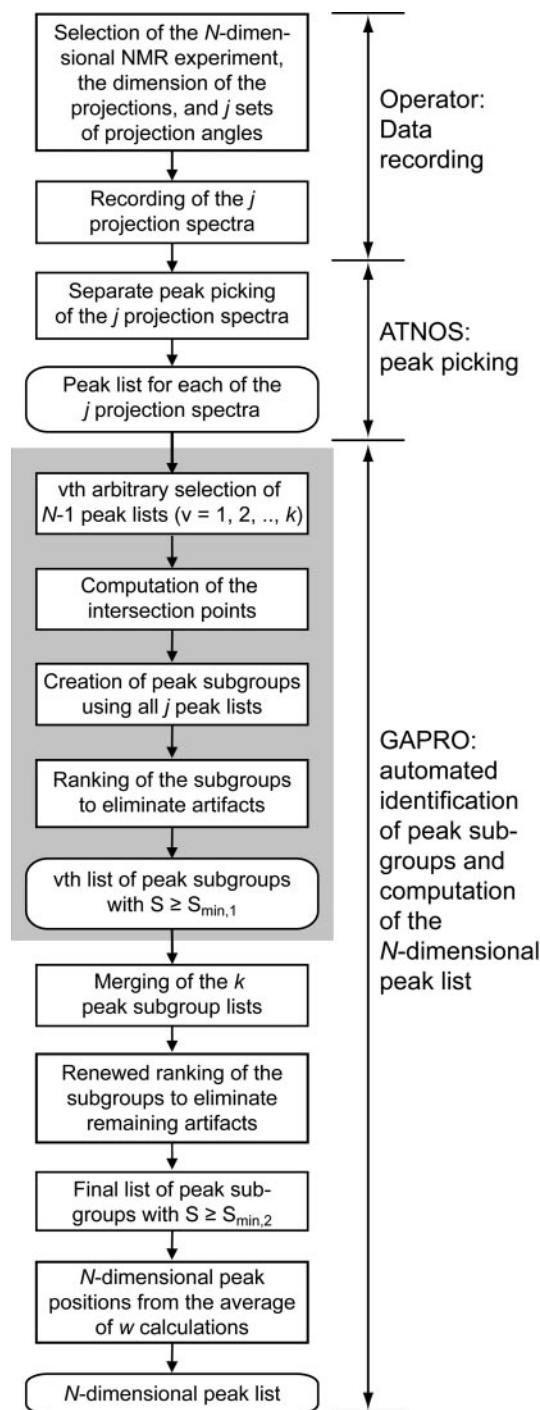


Fig. 3. Flowchart of APSY. Square boxes indicate processes, and boxes with rounded corners denote intermediate or final results. The steps underlayed in gray are repeated k times and thus generate k lists of peak subgroups.

remaining subspaces (Fig. 4c). This procedure is repeated until the value of S for all remaining candidate points falls below a user-defined threshold, $S_{\min,1}$, at which point a list of peak subgroups is generated. The subgroup identification is repeated with k different, randomly chosen starting combinations of $N - 1$ projections (the user-defined parameter k is a small fraction of the total number of possible combinations of $N - 1$ projections), and k peak subgroup lists are thus obtained (gray box in Fig. 3). These lists are merged into a single list, which is again subjected

to the same type of ranking procedure, so that all subgroups with $S < S_{\min,2}$ are eliminated. From the resulting “final” list of subgroups, the peak positions in the N -dimensional space are calculated (Fig. 4d). Below, the computational techniques used for individual steps in Fig. 3 are described.

Intersection of Subspaces. To simplify the mathematical treatment, we describe a $(N - l)$ -dimensional subspace L ($1 < l < N$) by a point Q^L in this subspace and a set of orthonormal vectors, $\{\vec{p}_1^L, \dots, \vec{p}_l^L\}$, orthogonal to L . To intersect, for example, four 3D subspaces in 5D frequency space, two of the 3D subspaces can intersect to a 2D subspace, which can then intersect with one of the remaining 3D subspaces to a 1D subspace, which can intersect with the fourth 3D subspace to a point.

L and M are two subspaces of dimensionality $(N - l)$ and $(N - m)$, with $(1 < l < N)$ and $(1 < m < N)$. L is described by $\{\vec{p}_1^L, \dots, \vec{p}_l^L\}$ and the point Q^L ; M is described by $\{\vec{p}_1^M, \dots, \vec{p}_m^M\}$ and Q^M . Both L and M are orthogonal to the direct dimension, and therefore both $\{\vec{p}_1^L, \dots, \vec{p}_l^L\}$ and $\{\vec{p}_1^M, \dots, \vec{p}_m^M\}$ include the unit vector of the direct dimension. If Eqs. 3 and 4 are satisfied, the subspaces L and M intersect in a subspace K of dimensionality $(N - k)$, with $k = l + m - 1$ as follows:

$$|Q^L(N) - Q^M(N)| \leq \Delta v_{\min} \quad [3]$$

$$\dim\{\vec{p}_1^L, \dots, \vec{p}_l^L, \vec{p}_1^M, \dots, \vec{p}_m^M\} = l + m - 1. \quad [4]$$

$Q^L(N)$ and $Q^M(N)$ are the N th coordinates of Q^L and Q^M , respectively; Δv_{\min} is a user-defined intersection tolerance in the direct dimension; and \dim stands for “dimension of.” Eq. 4 implies that $\{\vec{p}_1^L, \dots, \vec{p}_l^L\}$ and $\{\vec{p}_1^M, \dots, \vec{p}_m^M\}$ share only the direct dimension. The subspace K is then described by the orthonormal basis $\{\vec{p}_1^K, \dots, \vec{p}_k^K\}$, and by a point Q^K with its coordinates 1 to $(N - 1)$ given by the $l + m$ scalar products of Eq. 5

$$\begin{aligned} \vec{p}_z^L \cdot \left(\overrightarrow{Q^K Q^L} \right) &= 0 \quad z = 1, \dots, l \\ \vec{p}_z^M \cdot \left(\overrightarrow{Q^K Q^M} \right) &= 0 \quad z = 1, \dots, m. \end{aligned} \quad [5]$$

The N th coordinate of Q^K is the arithmetic average of the N th coordinates of Q^L and Q^M .

Distance Between a Point and a Subspace. The distance r between a point Q and a $(N - l)$ -dimensional subspace L , as described by a point Q^L and an orthonormal set of vectors orthogonal to L , $\{\vec{p}_1^L, \dots, \vec{p}_l^L\}$, is given by

$$r = \sqrt{\sum_{z=1}^l \left(\vec{p}_z^L \cdot \left(\overrightarrow{Q Q^L} \right) \right)^2}. \quad [6]$$

Peak Positions in the N -Dimensional Space. The fact that the peak positions are generally overdetermined by the experimental data are used to refine the peak coordinates. From each subgroup, $N - 1$ elements are arbitrarily chosen, and their associated subspaces are intersected to yield the position of the N -dimensional peak (Fig. 4). This procedure is repeated w times, where w is a user-defined parameter. Because of the limited precision of the individual chemical shift measurements, this procedure results in w slightly different peak positions, which then are averaged in each dimension to obtain the final positioning of the N -dimensional peak Q^i .

Materials and Methods

Sample Preparation. [U- ^{13}C , ^{15}N]-labeled 434-repressor(1–63) was produced following procedures described in refs. 28 and 30. For

$\pm 60^\circ$), $(90^\circ, \pm 30^\circ)$, $(90^\circ, \pm 60^\circ)$, $(\pm 30^\circ, \pm 30^\circ)$, $(\pm 60^\circ, \pm 30^\circ)$, and $(\pm 45^\circ, \pm 60^\circ)$. The projection spectra were peak picked with ATNOS (22) to generate the input for GAPRO. The 4D peak list that resulted after ≈ 10 min of GAPRO computation time contained 59 peaks, which is to be compared with a total of 60 peaks expected from the chemical structure of the molecule. Although on average 18 ± 9 noise artifacts were picked in each projection, the final 4D peak list generated by the GAPRO algorithm contained 59 cross peaks and not a single artifact. Only the peak that would correlate the residues of the N-terminal dipeptide was missing. It had a signal intensity below the noise level in all projections. The precision of the chemical shifts in the final APSY peak list has been estimated to be 1 Hz in the direct dimension and 8 Hz in each of the three indirect dimensions.

The 5D APSY-HACACONH experiment was recorded with the 12.4-kDa protein TM1290. The pulse sequence used is described in *Supporting Materials and Methods*, and further experimental details are given in *Materials and Methods*. In this experiment, 28 2D projections were recorded in 11 h. The projection angles and the spectral widths are listed in Table 2, and the 28 2D projection spectra are shown in Fig. 9. The final 5D peak list produced from the 5D APSY-HACACONH experiment with the protein TM1290 contained all of the peaks that were expected from the chemical structure of the molecule and the previously published NMR assignments (31), and there were no artifacts contained in the final peak list.

Discussion

In this work, we presented the foundations of APSY and introduced the algorithm GAPRO for automated spectral analysis. We then implemented APSY for high-dimensional heteronuclear correlation NMR experiments with proteins. In two applications without any human intervention after the initial set-up of the experiments, we obtained complete peak lists with high-precision chemical shifts for 4D and 5D triple resonance spectra. For the future, we anticipate that APSY will be the first step, after protein preparation, in a fully automated process of protein structure determination by NMR. In addition to providing automated peak picking and computation of the corresponding chemical shift lists, as described in this work, APSY is expected to support automated sequential resonance assignment. For these envisaged goals, APSY has the promise of being a valid alternative to related NMR techniques that have recently been introduced for similar purposes. Thus, when compared with projection-reconstruction NMR (13–16), APSY has the advantage of relying exclusively on the analysis of experimental low-dimensional projection spectra, with no need to ever reconstruct the parent high-dimensional spectrum. When compared with G-matrix Fourier transform NMR (11), APSY differs in that there are no restrictions on the selection of the number of projections or the combinations of projection angles. The strongest asset of APSY, however, is that the algorithm GAPRO enables fully automated analysis of the experimental projection spectra. As a primary result, complete peak picking and computation of high-precision chemical shift lists are obtained without any bias that could result from human intervention.

APSY and Protein Size. APSY has so far been applied to the 6.9-kDa protein 434-repressor(1–63) and the 12.4-kDa protein TM1290. To obtain an estimate of possible limitations for APSY applications with larger proteins due to spectral overlap, we analyzed the peak separations in 4D and 5D triple resonance spectra of a sample of 54 proteins with sizes from $n = 50$ to 300 residues, which were simulated from the BioMagResBank (32) chemical shift deposits. Considered were the average of the distances from each peak to its nearest-by neighbor, d_{av} , and the distance between the two most closely spaced peaks in the entire spectrum, d_{min} . Fig. 5 represents the data for the two experi-

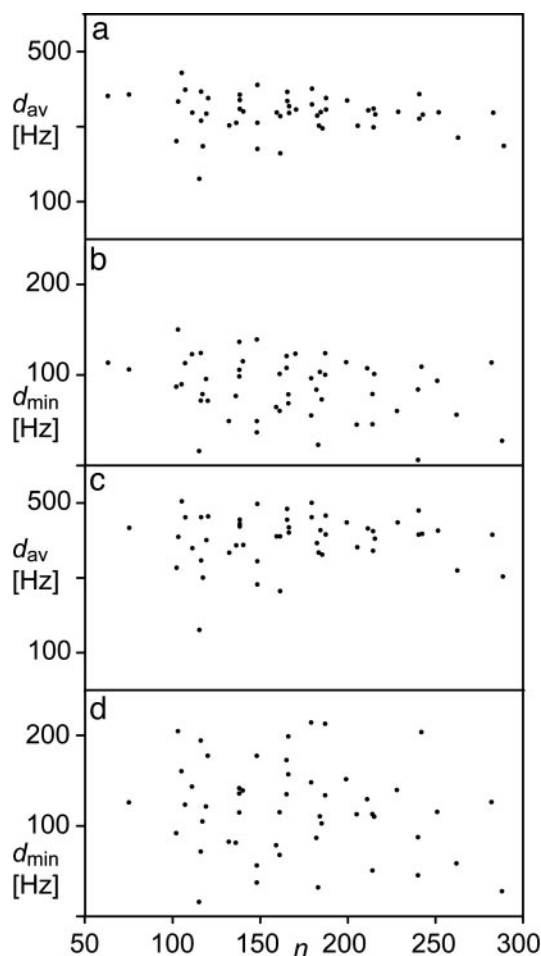


Fig. 5. Plots of peak separation in hertz vs. protein size (n is the number of residues). d_{av} is the average distance to the closest peak, and d_{min} is the distance between the closest pair of peaks. (a and b) 4D HNCOCA. (c and d) 5D HACACONH. The calculations were based on the BioMagResBank chemical shift deposits of 54 proteins (see Table 3, which is published as supporting information on the PNAS web site) and assumed a ^1H frequency of 750 MHz. Gly residues are included only in a and b.

ments 4D HNCOCA and 5D HACACONH. There is no obvious correlation between n and either d_{av} or d_{min} , indicating that close approach of peaks is distributed statistically and depends on particular properties of the protein, irrespective of its size. In Fig. 5, the statistical probability to encounter pairs of peaks that could not be resolved by APSY is $<1\%$ for protein sizes up to at least 300 residues, which is representative for 4D and 5D triple resonance data sets that contain one peak per residue. Foreseeably, sensitivity of signal detection therefore will be a more stringent limitation for APSY applications than spectral crowding. From our experience to-date, projection spectra with a signal-to-noise ratio of $\approx 3:1$ are required for efficient use of APSY with automated peak picking.

For the few expected closely spaced pairs of peaks, APSY is in a good position to resolve potential difficulties, because it is not required that a given N -dimensional resonance is found in all projection spectra. Peaks with overlap in one or several projections usually will be resolved in many other projections (Fig. 4). Similar to G-matrix Fourier transform NMR (11), APSY is also well prepared to deal with inaccurate peak positions from automated peak picking, which may arise from peak overlap. Because the final N -dimensional APSY peak list is computed as the average of a large number of measurements (Fig. 3),

