

Evaluation of stereo-array isotope labeling (SAIL) patterns for automated structural analysis of proteins with CYANA

Teppei Ikeya,^{1,2} Tsutomu Terauchi,² Peter Güntert³ and Masatsune Kainosho^{2*}

¹ Japan Biological Informatics Consortium (JBIC), Japan

² CREST/JST and Graduate School of Science, Tokyo Metropolitan University, 1-1 Minami-ohsawa, Hachioji, 192-0397, Japan

³ Tatsuo Miyazawa Memorial Program, RIKEN Genomic Sciences Center, 1-7-22 Suehiro-cho, Tsurumi, Yokohama, 230-0045, Japan

Received 4 December 2005; Revised 17 February 2006; Accepted 28 February 2006

Recently we have developed the stereo-array isotope labeling (SAIL) technique to overcome the conventional molecular size limitation in NMR protein structure determination by employing complete stereo- and regiospecific patterns of stable isotopes. SAIL sharpens signals and simplifies spectra without the loss of requisite structural information, thus making large classes of proteins newly accessible to detailed solution structure determination. The automated structure calculation program CYANA can efficiently analyze SAIL-NOESY spectra and calculate structures without manual analysis. Nevertheless, the original SAIL method might not be capable of determining the structures of proteins larger than 50 kDa or membrane proteins, for which the spectra are characterized by many broadened and overlapped peaks. Here we have carried out simulations of new SAIL patterns optimized for minimal relaxation and overlap, to evaluate the combined use of SAIL and CYANA for solving the structures of larger proteins and membrane proteins. The modified approach reduces the number of peaks to nearly half of that observed with uniform labeling, while still yielding well-defined structures and is expected to enable NMR structure determinations of these challenging systems. Copyright © 2006 John Wiley & Sons, Ltd.

KEYWORDS: NMR; ¹H; ¹⁵N; ¹³C; SAIL; protein structure; CYANA

INTRODUCTION

In this era of structural genomics the number of protein structures deposited in the Protein Data Bank (PDB)¹ is increasing rapidly, and has reached a total of about 36 000 as of March 2006. In the past two decades, NMR spectroscopy has become the second accepted method, after X-ray crystallography, for determining the three-dimensional structures of proteins. NMR provides information about the structures and dynamic properties of proteins in solution, and offers an approach to solve the structures of proteins that fail to crystallize. However, NMR still has limitations, as compared to X-ray crystallography. Statistical data from the PDB indicate that NMR was used to solve only about 15% of the protein structures in the PDB, and the percentage of NMR-based structures among all newly deposited protein structures has not increased in recent years. One likely reason for this situation is that the NMR method for protein structure determination is not efficient. Because of the complexity of the spectra, the analysis of protein NMR data requires much time and expertise. In addition, the PDB statistics revealed that there are very few NMR structures of proteins

over 25 kDa, while for proteins smaller than 10 kDa, there are rather more NMR than X-ray structures. This reflects the molecular size limitation of the NMR method. Thus far, it has been very difficult to solve the structures of molecules larger than 30 kDa by NMR, because of spectral crowding and line broadening of the signals due to fast transverse relaxation. Therefore, for the structure determinations of larger proteins we must address two main challenges in terms of accuracy and efficiency; improved quality and automatic analysis of the NMR spectra.

Recently we developed the stereo-array isotope labeling (SAIL) method for surmounting these problems by applying a complete stereo- and regiospecific pattern of stable isotopes.² In the SAIL method, the 20 proteinaceous amino acids (Fig. 1) are prepared by chemical and enzymatic syntheses, based on the following design concepts:² (i) stereo-selective replacement of one ¹H in methylene groups by ²H, (ii) replacement of two ¹H in methyl groups by ²H, (iii) stereo-selective labeling of the isopropyl groups of Leu and Val, such that one methyl is ⁻¹²C(²H)₃ and the other is ⁻¹³C¹H(²H)₂, and (iv) labeling of six-membered aromatic rings by alternating ¹²C – ²H and ¹³C – ¹H moieties. These principles are designed to provide an optimal labeling pattern for structure determination, without the loss of relevant structural information. In this way, the nuclear Overhauser effects (NOEs) that are additionally present in

*Correspondence to: Masatsune Kainosho, Graduate School of Science, Tokyo Metropolitan University, 1-1 Minami-ohsawa, Hachioji, 192-0397, Japan.
E-mail: kainosho@nmr.chem.metro-u.ac.jp

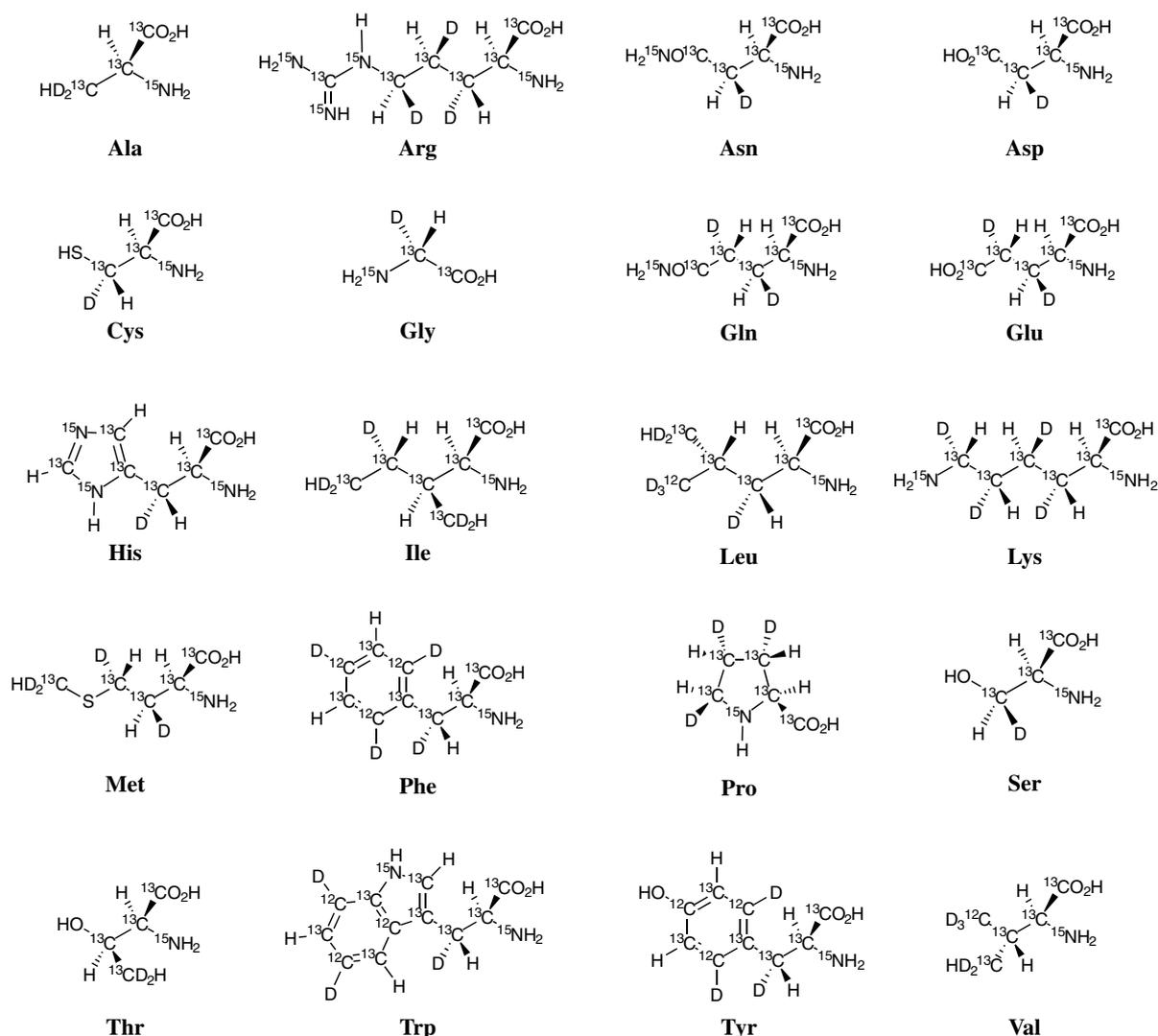


Figure 1. Chemical structures of the SAIL amino acids. H denotes ^1H ; D denotes ^2H .

uniformly labeled proteins contributes to spectral overlap but are, in the absence of stereo-specific assignments, virtually redundant with those observed with SAIL. In addition, SAIL makes it possible to derive more accurate distance restraints from NOE measurements, owing to a decrease in the spin diffusion effects. The viability and power of the SAIL approach have recently been shown² by the high-quality structure determinations of the SAIL-calmodulin (17 kDa) and SAIL-maltodextrin-binding protein (MBP) proteins (41 kDa), using automated NOESY assignment and structure calculation by the CYANA program.

Ideally, the SAIL technique could be applied to the NMR structure determinations of proteins larger than 50 kDa and membrane proteins. For this, it will presumably be necessary to optimize the isotope labeling patterns further, to cope with the extensive crowding and line broadening that are characteristic of the spectra of such proteins. In an overlap and relaxation optimized version of the SAIL approach, the number of ^1H nuclei is reduced further to enable the observation of well-shaped and separated signals, even in the cases of proteins beyond 50 kDa. However, this approach carries the potential risk that crucial structural information

will be lost for proteins with ^1H densities that are too low. The sparser distribution of protons may also hamper the performance of network anchoring in the automated NOE assignment algorithm of CYANA, which utilizes the partial redundancy of the NOE distance restraint network to ensure the reliability of the NOE assignments.³ Therefore, to facilitate the design of overlap and relaxation optimized SAIL amino acids, in terms of the requisite structural information and the suitability for automatic assignments, we performed test structure calculations for three model proteins, calmodulin, LpxC, and OmpA using either uniform labeling UL, original SAIL, or its modified version, and evaluated the relationship between the reduction of the proton density and the accuracy of the structures.

EXPERIMENTAL

Proteins

Calmodulin is a protein for which experimental SAIL NMR data were collected and a high-quality solution structure was solved previously.² The LpxC deacetylase from *Aquifex aeolicus* (31 kDa) is the only protein larger than 25 kDa for

which the BioMagResBank (BMRB)⁴ has more than 90% of the commonly observable ¹H chemical shifts.⁵ A solution structure of LpxC was solved by NMR using uniform ¹³C and ¹⁵N labeling.^{5,6} As a test system for membrane protein structure determination, we chose OmpA, a 19 kDa β -barrel membrane protein that was previously studied by NMR.⁷ Using uniform deuteration complemented by selective protonation of methyl groups, a low-resolution structure of the structurally related integral membrane protein OmpX has recently been determined by solution NMR.⁸ In addition, high-resolution X-ray structures are available for all three test proteins.^{9–11}

Overlap and relaxation optimized SAIL patterns

SAIL patterns can be characterized by the percentage of ¹H atoms relative to the uniformly labeled protein in an H₂O solution. The original SAIL method retains 64% of all protons, and 44% of all side-chain protons.² The modified SAIL pattern studied in this paper comprises 53% of all protons, and only 28% of all side-chain protons (Table 1). The modified SAIL ¹H labeling pattern is essentially a subset of standard SAIL. For instance, a uniformly ¹³C labeled Leu side chain contains four ¹³C and nine ¹H nuclei. The original SAIL pattern reduces the NMR-active nuclei to three ¹³C and three ¹H (Fig. 1). In the case of modified SAIL, the H ^{γ} methine proton, which tends to be overlapped and to provide only largely redundant NOEs, is additionally replaced by ²H (Fig. 2). C ^{γ} can be replaced by ¹²C for a further reduction of the relaxation, or be kept as ¹³C to enable the assignment of the ¹³C¹H(²H)₂ methyl group by through-bond experiments (Fig. 2). The only instance in which modified SAIL is not a subset of SAIL occurs with the aromatic ring of Phe, which has ¹H exclusively in the two ϵ positions in SAIL, but only at the H ^{ζ} position in modified SAIL (Fig. 2).

Chemical shifts and NOESY peak lists

¹³C- and ¹⁵N-edited NOESY peak lists for modified SAIL calmodulin were simulated by modification of the experimental NOESY peak lists from the original SAIL-calmodulin NMR structure determination.² The NOESY peak lists for SAIL and modified SAIL are identical except for the peaks

that were removed for modified SAIL, due to the incorporation of additional ²H nuclei. For the exceptional case of Phe H ^{ζ} , the NOESY peaks were simulated based on the X-ray structure (see below) and added to the peak lists.

For LpxC, the resonance assignments deposited in the BMRB database were used. No side-chain assignments are available for OmpA. Chemical shift values for OmpA were therefore simulated randomly assuming normal distributions with the average and standard deviation taken from the chemical shift statistics overall proteins in the BMRB database.⁴ NOE lists of LpxC and OmpA were then simulated on the basis of the X-ray structures, which were first regularized to adhere to the Engh and Huber standard geometry¹² used by CYANA for the covalent bond lengths, bond angles, and planar groups. Regularization was carried out with CYANA in the presence of a large number of distance constraints extracted from the original X-ray structures, and led only to minimal changes in the overall RMSD. The CYANA-regularized X-ray structures were used for the simulation of the NOESY peak lists and as reference structures for the evaluation of the present test calculations. NOESY cross peak volumes, V , were calculated from the corresponding distance, r , in the reference structure by assuming a $1/r^6$ relationship and random fluctuations to model the fact that actual experimental NOE data do not strictly follow the theoretical distance-to-volume relationship for an isolated spin pair in an internally rigid molecule. Using three random numbers, f_1, f_2 , and f_3 , between 0 and 1, the cross peak volumes were set to $V = f_1 a / r^6$ if $f_3 > 0.3$, and to $V = f_1(1 + f_2) a / r^6$ otherwise, using an arbitrary value $a = 3 \times 10^9$. Additionally, strong NOEs that correspond to a short distance can in practice be detected with much higher probability than weak ones. Simulating all NOEs up to the cutoff distance of 5.5 Å would result in an unrealistically large number of weak NOEs. Therefore, we considered only a subset of all NOEs for distances below 5.5 Å. NOEs were chosen randomly according to a Gaussian probability distribution, with a mean value of 2.4 Å and a standard deviation of 2.36 Å. The parameters of the Gaussian function were adjusted for optimal agreement of the simulated distance distribution with the experimental NOE data for SAIL calmodulin.

Table 1. Calmodulin, LpxC, and OmpA proteins with conventional uniform labeling (UL), original SAIL and modified SAIL

	Calmodulin			LpxC			OmpA		
	UL	SAIL	Modified ^b	UL	SAIL	Modified ^b	UL	SAIL	Modified ^b
¹ H atoms in H ₂ O	1094	695	584	2302	1423	1194	1252	802	709
¹ H atoms in ² H ₂ O	851	452	340	1808	929	689	960	802	709
Side-chain ¹ H atoms in ² H ₂ O	692	304	192	1499	647	407	762	338	241
NOESY peaks	–	4814	3457	11 079	7615	6111	8893	5435	4197
Assigned NOESY peaks	–	4568	3159	10 941	7532	6036	8574	5232	3941
Backbone RMSD (Å) ^a	–	0.37	0.61	0.40	0.50	0.80	0.37	0.58	0.64
Heavy atom RMSD (Å) ^a	–	0.69	1.04	0.95	1.03	1.37	0.58	0.81	0.94
Backbone RMSD to X ray (Å) ^a	–	0.95	0.87	0.89	0.95	0.97	1.01	0.94	1.24

^a RMSDs were calculated for the backbone atoms N, C ^{α} , C ^{γ} , or for all heavy atoms of residues 82–146 (calmodulin), 1–255 (LpxC), or 1–172 (OmpA). The value given is either the average over the 20 CYANA conformers that represent the solution structure of the RMSDs to the mean coordinates, or the single RMSD between the mean coordinates and the reference X-ray structure.

^b The overlap and relaxation optimized version of SAIL, modified SAIL.

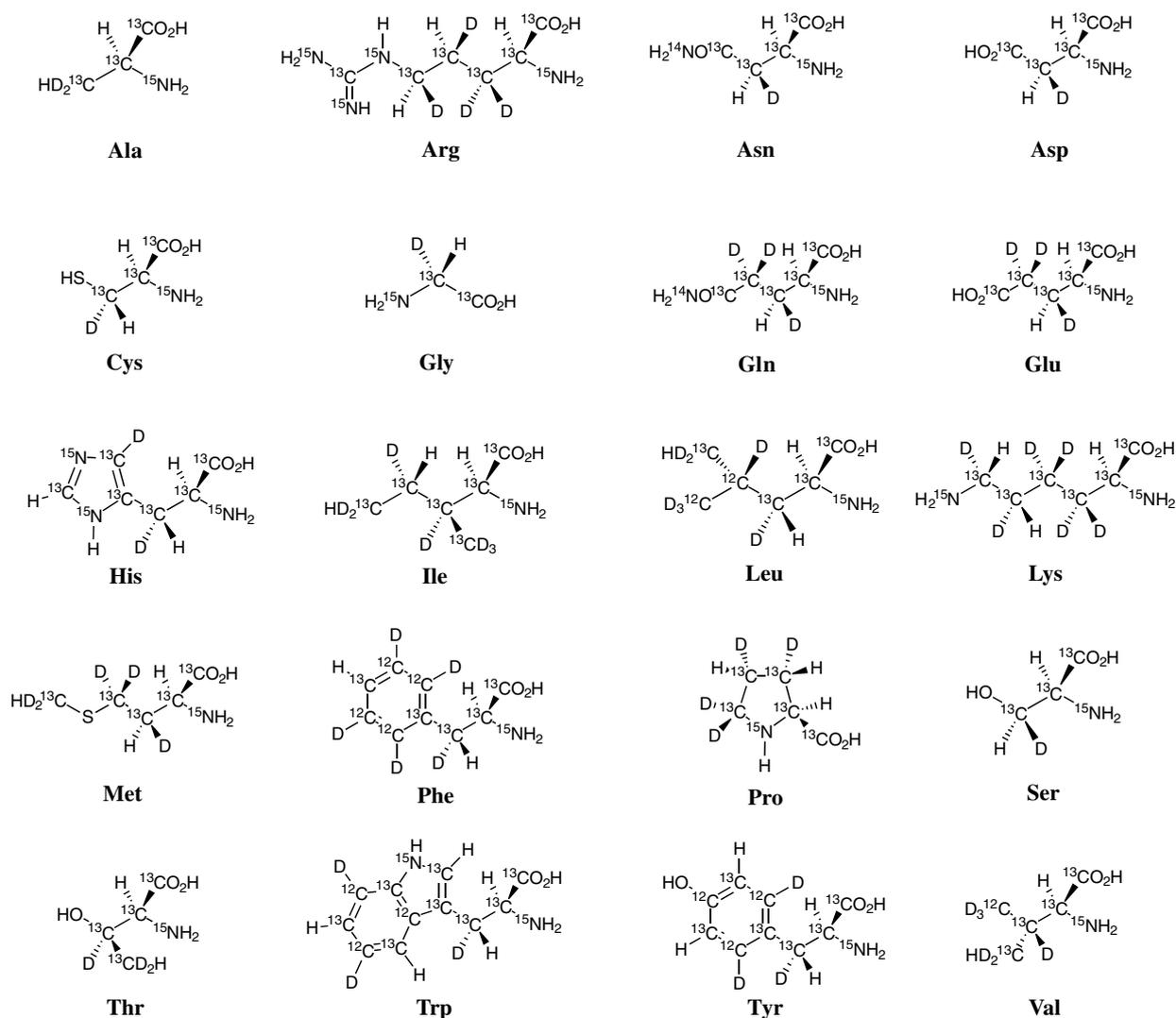


Figure 2. Modified, overlap and relaxation optimized SAIL patterns. H denotes ^1H ; D denotes ^2H .

Automated NOE assignment and structure calculation using CYANA

The calmodulin, LpxC, and OmpA structures were computed with the program CYANA,¹³ using automated NOE assignment³ and torsion angle dynamics for the structure calculation,¹⁴ which was started from 100 conformers with random torsion angle values. The strictly probability-based NOE assignment algorithm of CYANA 2.2¹⁵ was used in place of the earlier CANDID³ algorithm. The standard CYANA simulated annealing schedule was applied with 15 000 torsion angle dynamics steps. Backbone torsion angle restraints obtained from chemical shifts with the program TALOS¹⁶ were added to the input for CYANA in the case of calmodulin, but not for LpxC or OmpA. Hydrogen bond restraints were not used. No stereo-specific assignments were assumed in the calculations with UL. The 20 conformers with the lowest final CYANA target function values were analyzed and compared with the reference structure. All calculations were carried out five times, with different random number generator seed values and otherwise identical input parameters and data. The data reported are averages over five calculations.

RMSD values were calculated for the C-terminal domain of residues 82–146 for calmodulin, for residues 1–255 of LpxC, and for all residues of OmpA.

RESULTS AND DISCUSSION

Evaluation of overlap and relaxation optimized SAIL with calmodulin

A series of CYANA structure calculations were performed for calmodulin using in turn the modified SAIL pattern for each individual amino acid type together with standard SAIL for the other 19 amino acid types. This allowed us to evaluate the relationship between the precision and accuracy of the structures and the reduction of the proton density in each amino acid type (Table 2). Since the SAIL and modified SAIL patterns differ from each other for 12 out of the 20 amino acid types, 12 calculations were performed. As compared with the reference calculation using the original SAIL pattern for all amino acids, the number of NOESY peaks was reduced by 6–216 (Table 2), or by 6–21 per modified SAIL amino acid residue. The most pronounced decrease in the number of NOEs was observed for Lys, which occurs eight times in the amino acid sequence of calmodulin. Nevertheless, the

Table 2. Structural statistics using single modified SAIL amino acid types in SAIL calmodulin

	SAIL	Arg	His	Ile	Leu	Lys	Met	Pro	Phe	Thr	Val	Gln	Glu	no H α
NOESY peaks	4814	4738	4808	4681	4692	4647	4710	4777	4665	4698	4700	4719	4598	3030
Assigned NOESY peaks	4568	4487	4567	4428	4440	4400	4459	4521	4411	4446	4451	4493	4336	2724
Backbone RMSD (Å) ^a	0.37	0.44	0.40	0.43	0.42	0.39	0.44	0.41	0.38	0.40	0.42	0.45	0.41	0.92
Heavy atom RMSD (Å) ^a	0.69	0.79	0.72	0.74	0.76	0.74	0.76	0.72	0.46	0.73	0.75	0.78	0.80	1.38
RMSD to reference (Å) ^a	0.95	1.08	0.95	1.00	1.06	0.97	0.94	1.08	1.03	1.00	0.88	0.89	0.93	1.31

^a RMSDs are calculated for the backbone atoms N, C α , C', or for all heavy atoms of the C-terminal domain of residues 82–146. The value given is either the average over the 20 CYANA conformers that represent the solution structure of the RMSDs to the mean coordinates, or the single RMSD between the mean coordinates and the reference X-ray structure.

structure obtained with modified SAIL-Lys did not differ significantly from the reference structure. The RMSDs in the presence of modified SAIL amino acids were only slightly higher than those obtained with the original SAIL. In practice, modified SAIL is actually expected to yield better results, because its higher resolution and decreased overlap can compensate for the small reduction in the theoretical number of restraints.

An additional calculation was performed in which all of the H α assignments and their corresponding NOESY peaks were eliminated. Such an approach would be expected to increase the spectral quality further, because H α is a major source of deleterious transversal relaxation in larger proteins. It could be imagined that the NOEs of H α are partially redundant with those from nearby H N and H β . However, the results of the test calculations (Table 2) clearly showed that the resulting structures are of much lower quality than the original SAIL ones, in terms of both precision (0.92 vs 0.37 Å) and accuracy (1.31 vs 0.95 Å). Thus, we conclude that maintaining the H α protons is crucial for obtaining a high-quality structure.

Next, structure calculations with complete modified SAIL calmodulin were carried out. In comparison with the SAIL reference structure, these structures showed only a slight increase of the backbone RMSD from 0.37 to 0.61 Å and of the side-chain RMSD, from 0.69 to 1.04 Å. Therefore, the structures still provide detailed information on the side-chain conformations. The deviations from the X-ray structure were 0.95 Å for SAIL and 0.87 Å for modified SAIL (Fig. 3 and Table 1). This result shows that it is possible to determine high-quality structures with modified SAIL, despite the decreased number and density of peaks. The automated NOE assignment algorithm in CYANA was able to assign 95% of the 4814 NOESY peaks with SAIL, and 91% of the 3457 NOESY peaks with modified SAIL. This indicates that automatic NOE assignment also works well with a lower proton density.

Structure calculations of LpxC and OmpA

For the structure calculations of the 31 kDa protein LpxC 92% of the ¹H chemical shifts were available. Structure calculations based on UL, SAIL, and modified SAIL consistently yielded assignments for more than 98.7% of the (simulated) NOESY peaks (Table 1). The resulting structures (Fig. 4) deviate from the X-ray structure by 0.89, 0.95, and 0.97 Å in terms of the backbone RMSD, respectively, and

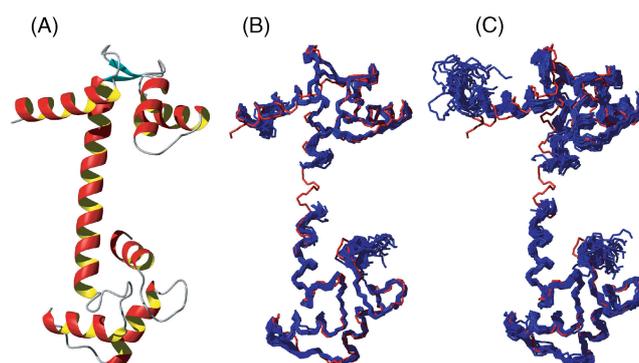


Figure 3. Calmodulin structures obtained with the original SAIL approach and with modified, overlap and relaxation optimized SAIL. (A) SAIL, ribbon diagram, (B) SAIL, structure bundle, and (C) modified SAIL. For comparison, the regularized X-ray structure is shown in red in (B) and (C).

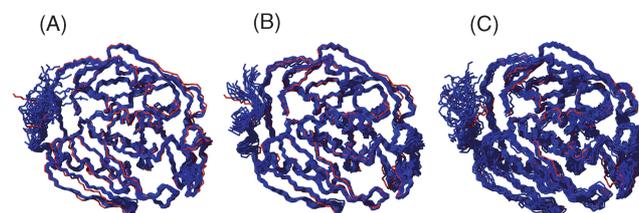


Figure 4. LpxC structures obtained from the experimental chemical shifts and simulated NOESY peaks assuming (A) UL, (B) SAIL, and (C) modified SAIL. For comparison, the regularized X-ray structure is shown in red.

thus have equivalent accuracy. The corresponding RMSDs within the structure bundles were 0.40, 0.50, and 0.80 Å for the backbone, and 0.95, 1.03, and 1.37 Å for all heavy atoms, respectively, for UL, SAIL, and modified SAIL (Table 1). This suggests equivalent precision for the UL and original SAIL methods, and somewhat lower apparent precision with modified SAIL. However, the number of NOEs in the simulated peak lists is almost twice as large with UL than with modified SAIL. In practice, it would be exceedingly difficult to evaluate more than 11 000 peaks in the crowded NOESY spectrum of uniformly labeled LpxC. (Only 4502 NOE restraints were used in the experimental NMR structure determination of LpxC.⁶) On the other hand, it is feasible to identify the approximately 6000 peaks generated in the modified SAIL pattern. Furthermore, due to decreased spectral overlap and relaxation, a much shorter, yet more complete, list of the

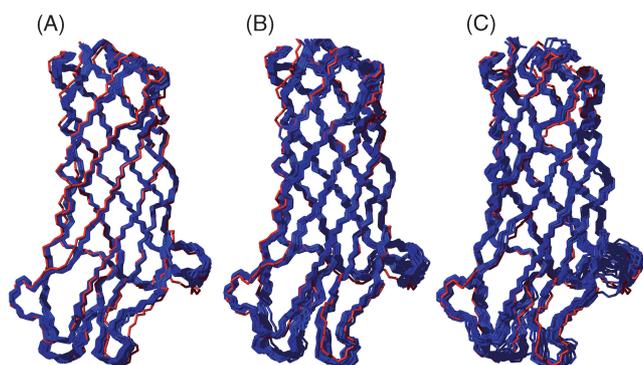


Figure 5. OmpA structures obtained from simulated chemical shifts and simulated NOESY peaks assuming (A) uniform labeling, (B) SAIL, and (C) modified SAIL. For comparison, the regularized X-ray structure is shown in red.

chemical shift assignments can be expected from protein samples with modified SAIL labeling.

Membrane proteins tend to exhibit severe line broadening in solution NMR, even if their molecular size is below 20 kDa. The detergents required to solubilize membrane proteins make the rotational motion of the protein–detergent system slow and the transverse relaxation short. To investigate the potential of the SAIL method for the NMR structure determination of membrane proteins, we carried out automated NOESY assignment and structure calculations for OmpA. From previous studies of OmpA only partial resonance assignments are available. Our test calculations based on UL, SAIL, and modified SAIL yielded assignments for 94–96% of the simulated NOESY peaks (Table 1) and structures (Fig. 5), which deviated from the X-ray structure by 1.01, 0.94, and 1.24 Å in terms of the backbone RMSD, respectively. The corresponding RMSDs within the structure bundles are 0.37, 0.58, and 0.64 Å for the backbone, and 0.58, 0.81, and 0.94 Å for all heavy atoms, respectively, for UL, SAIL, and modified SAIL (Table 1). The modified SAIL approach showed only a moderate reduction in the structural precision and accuracy relative to the highly unrealistic, idealized case of exhaustively analyzed NOESY spectra for the uniformly labeled protein.

CONCLUSIONS

Our simulation results suggest that the combined use of optimized SAIL patterns and automated structure calculation algorithms with CYANA has high potential for the

automated structure determinations of higher molecular weight proteins as well as membrane proteins. The precision of the structure using the new, overlap and relaxation optimized modified SAIL pattern was slightly lower than that obtained by uniform labeling and the original SAIL method because modified SAIL reduces the number of peaks expected under ideal conditions. However, with the actual experimental data better results can be expected for modified SAIL, because it facilitates the analysis of many otherwise broadened or overlapped signals. Our test calculations further showed that the automatic assignment algorithm in CYANA also works well with the lower density of protons in modified SAIL. We expect that the overlap and relaxation optimized SAIL patterns will contribute decisively to the determination of high-quality solution structures of proteins in the 30–100 kDa size range and membrane proteins.

Acknowledgements

Financial support by CREST/JST and by the Japan Biological Informatics Consortium (JBIC) is gratefully acknowledged.

REFERENCES

- Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. *Nucleic Acids Res.* 2000; **28**: 235.
- Kainosho M, Torizawa T, Iwashita Y, Terauchi T, Ono AM, Güntert P. *Nature* 2006; **440**: 52.
- Herrmann T, Güntert P, Wüthrich K. *J. Mol. Biol.* 2002; **319**: 209.
- Seavey BR, Farr EA, Westler WM, Markley JL. *J. Biomol. NMR* 1991; **1**: 217.
- Coggins BE, Li X, McClerren AL, Hindsgaul O, Raetz CRH, Zhou P. *Nat. Struct. Biol.* 2003; **10**: 645.
- Coggins BE, McClerren AL, Jiang L, Li X, Rudolph J, Hindsgaul O, Raetz CRH, Zhou P. *Biochemistry* 2005; **44**: 1114.
- Arora A, Abildgaard F, Bushweller JH, Tamm LK. *Nat. Struct. Biol.* 2001; **8**: 334.
- Fernández C, Hilty C, Wider G, Güntert P, Wüthrich K. *J. Mol. Biol.* 2004; **336**: 1211.
- Babu YS, Bugg CE, Cook WJ. *J. Mol. Biol.* 1988; **204**: 191.
- Whittington DA, Rusche KM, Shin H, Fierke CA, Christianson DW. *Proc. Natl. Acad. Sci. U.S.A.* 2003; **100**: 8146.
- Pautsch A, Schulz GE. *Nat. Struct. Biol.* 1998; **5**: 1013.
- Engh RA, Huber R. *Acta Crystallogr., Sect. A* 1991; **47**: 392.
- Güntert P. *Prog. NMR Spectrosc.* 2003; **43**: 105.
- Güntert P, Mumenthaler C, Wüthrich K. *J. Mol. Biol.* 1997; **273**: 283.
- Güntert P. Automated NMR structure calculation. *NMR Techniques in Structural Biology—From Liquid to Solid State*. Springer: New York, 2006.
- Cornilescu G, Delaglio F, Bax A. *J. Biomol. NMR* 1999; **13**: 289.