# BioMagResBank database with sets of experimental NMR constraints corresponding to the structures of over 1400 biomolecules deposited in the Protein Data Bank

Jurgen F. Doreleijers, Steve Mading, Dimitri Maziuk, Kassandra Sojourner, Lei Yin§, Jun Zhu*, John L. Markley, and Eldon L. Ulrich

BioMagResBank, UW-Madison, 433 Babcock Dr., Madison, WI, 53706, U.S.A.

§Keithley Instruments Inc. 28775 Aurora Road, Cleveland, OH, 44139, U.S.A.

*Department of Animal Health and Biomedical Sciences, UW-Madison, 1656 Linden Dr., Madison, WI, 53706, U.S.A.

## INTRODUCTION

This project[1] incorporates all the NMR data currently available from the Protein Data Bank (PDB)[2,3] into the BioMagResBank[4] (BMRB - http://bmrb.wisc.edu). The BMRB is a repository for information on biological macromolecules derived from NMR spectroscopy. The PDB has collected and archived experimental NMR data in so-called MR files (Magnetic Resonance). In these files, the different types of information are concatenated together with a PDB-like header and an often-large nomenclature mapping table at the end. There are three main problems when using this information on a large scale. First, sections (blocks) in an MR file correspond to different files that the authors provided, but the boundaries between individual files are not indicated. Second, the data are not in a standard format, which makes it hard if not impossible to use the data if the format cannot be read by the software available to a particular user. Lastly, the data have been archived "as is" and have not been validated[5].

All three problems have been addressed In this project and the BMRB is making a weekly updated archive available to the scientific community. We are currently working on matching the experimental data described here to the coordinate data available from the PDB. The experimental data themselves can then potentially be used to further refine the structures and help to reduce the time needed to solve homologous structures.



Fig. 1: Representation of the query interface showing the selections available to the user. For brevity, the listing of available data types was omitted here, but can be found in full in the Table below.

## RESULTS

At the time of writing[1] (November 21, 2002) the PDB contained 1410 entries with associated MR files. They were classified into 8383 sections. Thus on average each MR file consists of approximately 6 sections or data blocks. Almost all MR files start with a header such as that seen in PDB formatted coordinate files and end with a nomenclature mapping table (ignored in the following discussion), so on average about 4 data blocks constitute an entry.

The results are available from the BMRB web site under "Features/Data Access" and from the PDB web site where individual PDB entries have been linked to the results from this project. A new query interface at the BMRB (Fig. 1) allows a user to select files or blocks of data from individual PDB entries or to browse all of the PDB entries that have associated MR files. An overview of all classified files with raw data (Table 1) with further classification by program and data type is shown.

In total, 31 specific types of data and 16 programs were observed in addition to many formats for which a particular program has not been identified. In most file formats, the number of text lines in a block corresponds roughly to the number of constraints. The histogram of the number of lines in blocks (shown in Fig. 2) peaks for the bin containing blocks with 129 through 256 lines, but many blocks have substantially more lines. The most common data type is the distance restraint with 2511 blocks in total (combining the observed simple and ambiguous formats for NOEs, hydrogen bonds, etc.).

From the total number of MR file blocks annotated as constraints, it proved possible to parse 84% (3337/3975). The constraint lists that were parsed correspond to three data types (2511 distance, 788 dihedral angle, and 38 residual dipolar couplings lists) from the three most popular software packages used in NMR structure determination: X-PLOR/CNS (2520 lists), DISCOVER (412 lists), and DYANA/DIANA (405 lists). In parsing, invalid statements and comments are captured and put into a separate tables for future analysis. These constraints were then mapped to a developmental version of the BMRB data model.
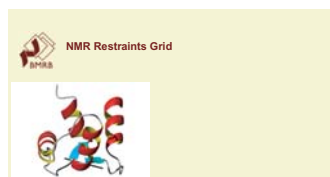
A specific set of files (Fig. 3) or blocks (Fig. 4) can be investigated further by following specific links showing the number of text lines in the block. At the level of an individual data block: the actual data text, a molecular image (prepared using MOLMOL[6]), and links are presented (Fig. 5). Tables from any level in the hierarchy can be downloaded as a comma-separated-value file, which enables easy upload to most spreadsheets. The data themselves can be downloaded as a zipped archive of a set of data blocks.



Fig. 3: A view of a set of files sharing a particular property that was previously selected. In this case, the files all have one or more blocks with dipolar couplings (not shown).



Fig. 4: The set of NOE distance restraint blocks with some additional information such as the number of lines in each block. Links to the MR file, PDB file, and block specific for each block are shown as underscored items.



Fig. 5: Individual block containing parsed RDCs (truncated). The example is from the structure of calmodulin[7].

## METHODS

All software used for this project, with the exception of the molecular image generation, is written in Java. The data files are stored in an Oracle 8i database using the entity-relationship model shown in Fig. 6. The images are actually served directly from a regular file system, because that is faster by far than by any Java servlet technique. The text files are streamed in a zipped format directly from the database.

The PDB archive including the restraints and coordinates is weekly mirrored using the mirror software (http://wuarchive.wustl.edu/packages/mirror) and the classifications and parsed data blocks will be kept up to date at BMRB.

The Java compiler compiler (javacc) program, which is freely available from http://www.webgain.com/products/java_cc, but owned by Sun Inc., allows us to express complicated grammar (for example that used by XPLOR/CNS files) in a more human readable fashion. Javacc then generates the actual parser program in Java.



Fig. 6: Diagram showing the entity-relationship data model used by the Wattos computer software package. For simplicity, the scheme shown omits three temporary tables (having the same definitions as the three tables shown) used for batched updates.

## REFERENCES

1. Doreleijers, J.F. et al., J.Biomol.NMR. accepted (2003).
2. Berman, H.M. et al., Nucleic Acids Res. 28, 235 (2000).
3. Bernstein, F.C., et al., J.Mol.Biol. 112, 535 (1977).
4. Ulrich, E.L. et al., XVIIth Int. Conf. Magn. Res. Biol. Systems. Tokyo, Japan (1998).
5. Doreleijers, J.F.., et al., J.Mol.Biol. 281, 149 (1998).
6. Koradi, R. et al., J.Mol.Graph. 14, 51 (1996).
7. Chou, J.J. et al., Nat.Struct.Biol. 8, 990 (2001).

Table. 1: Overview of the number of PDB MR files as sorted by data type (columns 2 through 4) and separated by the associated computer program (top row in columns 6 through 23).
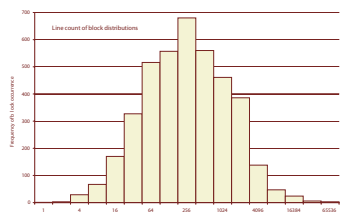


Fig. 2: Histogram of the distribution of the number of lines in blocks other than comments and mapping tables. The x-axi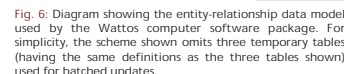s is shown on a logarithmic scale.