# BPS: a database of RNA base-pair structures

**Yurong Xin and Wilma K. Olson***

Rutgers, the State University of New Jersey, Department of Chemistry & Chemical Biology, BioMaPS Institute for Quantitative Biology, Wright-Rieman Laboratories, 610 Taylor Road, Piscataway, NJ 08854, USA

## ABSTRACT

The BPS (http://bps.rutgers.edu) is a database of RNA base-pair structures, higher-order base interactions and isosteric pairs (base pairs with similar shape). The main functions of the BPS are to find and annotate the structural and chemical features of the Watson–Crick and non-Watson–Crick (noncanonical) base pairs in high-resolution RNA structures, and to provide a user-friendly interface to browse and search for the base pairs. The current database contains 91 265 bp and 3386 higher-order base interactions from 426 RNA crystal structures and 61 819 bp that fall into one of 17 different isosteric classes. The base-pair data can be accessed by searches of base-pair patterns, structure identifiers (IDs) and structural types. The BPS also includes an Atlas with representative images of the various base pairs, higher-order base interactions and isosteric pairs and links to statistical information about these groups of structures.

## INTRODUCTION

The pairing of bases within and between long, stacked helical arrays underlies the unique three-dimensional folding of RNA (1). The RNA bases associate not only through the canonical A·U and G·C Watson–Crick pairs, but also via a wide variety of noncanonical interactions, such as the G·U wobble pair first suggested by Crick (2). Understanding the features of these primary architectural units can help to decipher the principles of RNA folding and to predict the three-dimensional structures of RNA from the primary nucleotide sequence.

The available high-resolution structures of RNA provide a rich resource for characterizing the different types of base pairs, including (i) the number of occurrences of different kinds of pairs, (ii) the sequence context of these pairs, (iii) the structural properties of the pairs, (iv) the structural contexts of the pairs, such as whether the bases are embedded in a double helix or in other parts of a structure and (v) the similar (isosteric) shapes of different pairs that allow them to substitute for one another in an RNA structure, such as the four combinations of Watson–Crick pairs (A·U, U·A, G·C, C·G) that fit in a regular A-RNA helical framework (3,4).

There are currently two databases with some information about RNA base pairs: the NCIR database of noncanonical interactions in RNA developed in the laboratory of George Fox at the University of Houston (5,6) and the Nucleic Database (NDB) (7). The current version of the NCIR lists over 1800 noncanonical base pairs and higher-order base interactions. The NCIR, however, does not provide quantitative information about base pairing, such as hydrogen-bond distances, interbase angles and base-pair statistics. The NDB, by contrast, uses the 3DNA software package (8,9) to extract the six rigid-body parameters that describe the spatial arrangements of paired bases in individual structures. The number of reported base pairs is, nevertheless, incomplete in some structures owing to the choice of settings used in the implementation of 3DNA.

Our newly developed relational database of RNA base-pair structures (BPS) includes both sequential and spatial information about the canonical and noncanonical arrangements of bases found in 426 high-resolution RNA structures. The database is more comprehensive than the aforementioned sites in terms of both the number of noncanonical base pairs and the variety of base-pair types. The searchable web interface, found at http://bps.rutgers.edu, provides a set of user-friendly tools to search for base pairs of different types, in different types of structures and in specific structures. To the best of our knowledge, the BPS is the first database to analyze the structures of all RNA base pairs quantitatively (in terms of parameters generated by 3DNA), to present representative images of the observed spatial patterns, to include isosteric base pairs, to characterize base pairs stabilized by backbone interactions, to allow for searches across structures and to sort the selected data in various ways.

## DATA COLLECTION

### Identification of base pairs

The BPS has compiled over 90 000 canonical and noncanonical hydrogen-bonded base interactions from 426 RNA crystal structures available in the NDB (7) and

*To whom correspondence should be addressed. Tel: +1 732 445 3993; Fax: +1 732 445-5958; Email: wilma.olson@rutgers.edu

PDB (Protein Data Bank) (10,11). We use the 3DNA software package (8,9) to identify canonical and noncanonical base pairs. Specifically, we use the following geometric criteria to find a base pair: (i) the distance between the origins on the standard reference frame (12) embedded on each base must be 15 Å or less; (ii) the magnitude of the vertical offset of the base planes (measured by the rigid-body parameter Stagger, see Figure 1) must be 2.5 Å or less; (iii) the smaller of the two angles between the normals of the base planes (given by the value of Buckle or its supplement, see Figure 1) must be 65° or less; and (iv) the distance between glycosidic base atoms, i.e. the purine N9 and pyrimidine N1, must be 4.5 Å or more. A base pair is accepted if one or more hydrogen bonds (both true donor-acceptor pairs and pseudo hydrogen bonds as described below) fall within a distance of 3.4 Å in the selected configurations. In this way we do not rule out base pairs stabilized by interactions with the sugar-phosphate backbone (9).

## Chemical classification

The base pairs are classified in terms of the chemical identities of the interacting bases, the hydrogen-bonding pattern an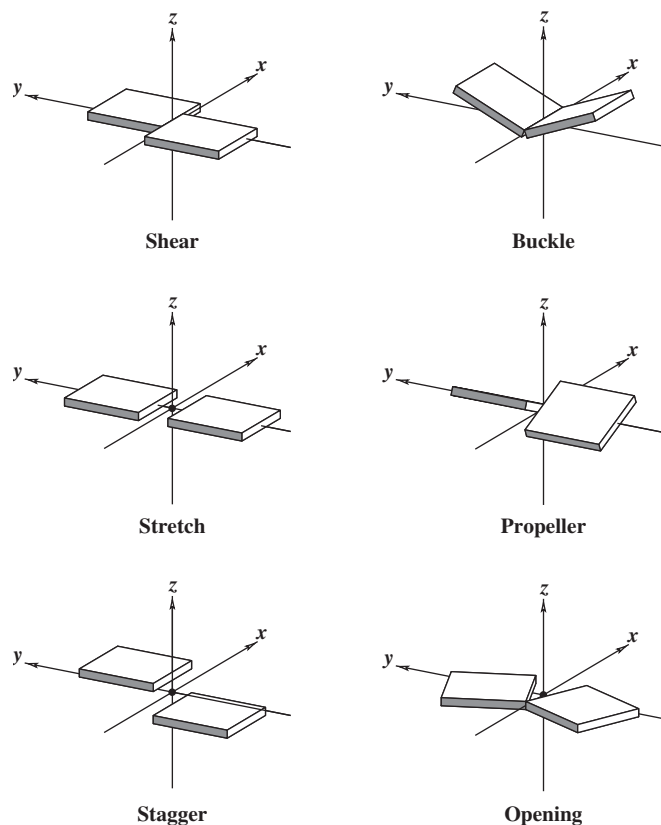d the base-pair orientation. Modified bases are treated by first fitting the reference frames of the most closely related standard bases to the observed chemical species and then computing the relevant spatial parameters. In order to be comprehensive, both true and pseudo hydrogen bonds (e.g. pairs of donor atoms or acceptor atoms that come within the chosen interatomic contact distance) are included in the pattern description. The base-pair orientation is defined by the angle $\gamma$ between the normals ($z$-axes) of the two bases, and described as parallel or antiparallel for values of $\gamma$ in the range $0° \leqslant \gamma \leqslant 90°$ or $90° < \gamma \leqslant 180°$, respectively.

## Structural context

Each base pair is also described by the secondary structural context in which the constituent bases occur, i.e. their locations along the helical stretches of stacked base pairs, identified with the 3DNA software (8,9), or within the intervening single-stranded nucleotide regions. Because no limits are placed on residue number in the search for 'neighboring' base pairs, some of the helices are quasi-continuous in that stacked pairs are not necessarily connected by covalent bonds (Figure 2). Thus, the bases that occur at these 'nicks' in RNA helical structures are distinguished from those that are located at the 5′- and 3′-termini or within the chemically continuous helical regions. The software also identifies single bases inserted within or at the ends of an otherwise continuous helical stretch, isolated base pairs with no immediately stacked neighbors and bases found within or at the ends of single-stranded regions. The various structural categories are illustrated schematically in Figure 2 and described in more detail in Table 1.

## Higher-order interactions

Bases involved in higher-order interactions, found from a 'horizontal' search in the planes of associated residues with 3DNA (9), are placed in groups of the same chemical and hydrogen-bonding type. The current version of the BPS includes 3386 unique higher-order interaction patterns.

## Isosteres

Isosteres are identified with the three virtual parameters used conventionally to characterize Watson–Crick base pairs: the virtual C1′⋯C1′ distance and the angles $\lambda_I$ and $\lambda_{II}$ between the C1′⋯C1′ vector and the glycosidic bonds of nucleotides I and II. These three parameters are sufficient to determine the isostericity of both canonical and noncanonical base pairs and are henceforth termed 'isosteric parameters'. There are currently 17 isosteric classes in the database.



**Figure 1.** Block representation of the six rigid-body parameters used to identify and describe interacting base pairs: three components of displacement called Shear, Stretch and Stagger and three angular parameters termed Buckle, Propeller and Opening (16). The examples illustrate distortions of canonical Watson–Crick base pairs with the shaded minor-groove edge, also called the Sugar edge (3), pointed toward the reader.

## DATABASE CONTENT

The BPS contains two major modules, Atlas and Search, as the top level of the database structure (Figure 3). The Atlas contains images of all base-pair-related classes, namely the base-pairing patterns, higher-order base interactions and isosteric classes, which constitute the middle layer of the database. These three submodules lead,

in turn, to entries for the isosteres, multiplets (i.e. higher-order base interactions) and base-pairing patterns at the lowest layer of the database. The Search module provides access to the base pairs and base-pairing patterns of individual structures, which are linked to each other and to the isosteres and multiplets.

## DATA ACCESS

### Search options and sample output

The database provides four types of search options for identification of base pairs: (i) structure identifier (ID) (i.e. the NDB or PDB ID), (ii) Saenger pattern (i.e. one of the 28 arrangements of base pairs listed in the classic textbook of Wolfram Saenger (13)), (iii) base-pairing pattern (i.e. all currently observed patterns of base association) and (iv) RNA molecule type (i.e. ribozyme, t-RNA, etc.).
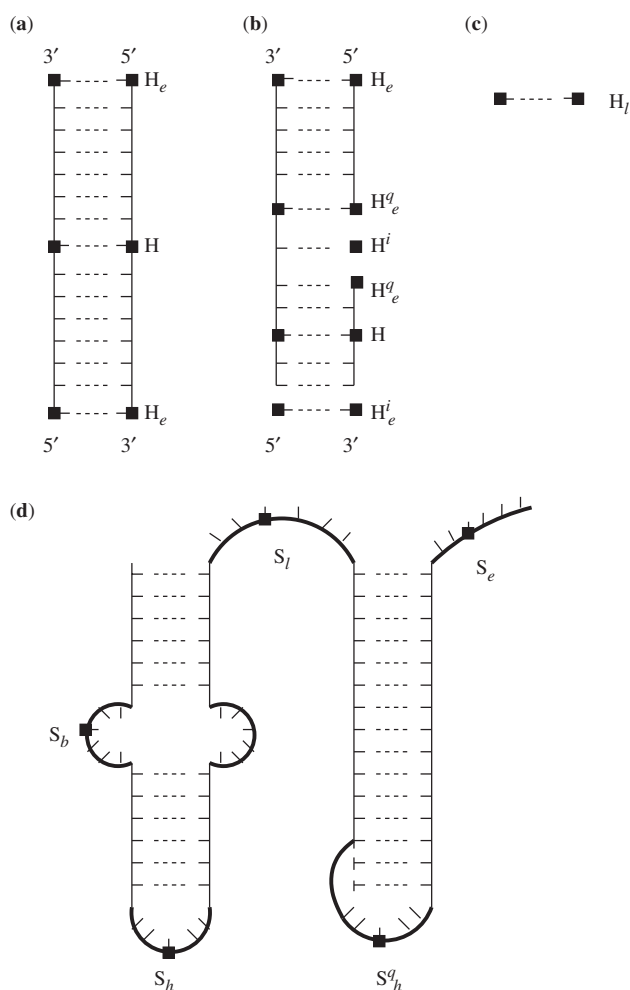


**Figure 2.** Schematic illustration of the terms used to describe the structural context of paired RNA bases. Double helices are represented by pairs of parallel straight lines and single-stranded segments by bold curves. The quasi-continuous helices include gaps along the lines. Ticks connected by dashed lines denote paired bases and ticks without partners represent unpaired bases. The small black boxes illustrate nucleotides in various structural contexts. The symbols near the boxes correspond to the literal information stored in the database.

Figure 4 illustrates representative search results obtained for a given structure ID, here the ID of the P4-P6 group I ribozyme domain (14): URX053 (NDB_ID); 1GID (PDB_ID). The search result gives the number of base pairs and a tabulated list with (i) the type of molecule, (ii) the NDB and 3DNA IDs, (iii) the chemical, residue and chain identities, (iv) the base-pair pattern, (v) the secondary structural form and (vi) the structural context of each base pair in the selected structure (Figure 4A). The result page provides further links to other information: a page that includes a table of the rigid-body parameters describing each base pair (Figure 4B); a summary page for the selected base-pair (URX053_1), here a canonical Watson–Crick G·C base pair, with (i) all stored chemical and structural information, (ii) an atomic-level representation of the interacting nucleotides, (iii) a block image of the sequential context and (iv) a schematic illustration of the structural context, here with G at the 3′-end of a single-stranded helical fragment (denoted by $H_e^q$) and C within a continuous single-stranded helical stretch (denoted by H) (Figure 4C); a page with the rigid-body parameters describing the average arrangement of all base pairs in the same (GC_1) pairing pattern as the selected example

**Table 1.** Structural context of base pairs observed in RNA structures

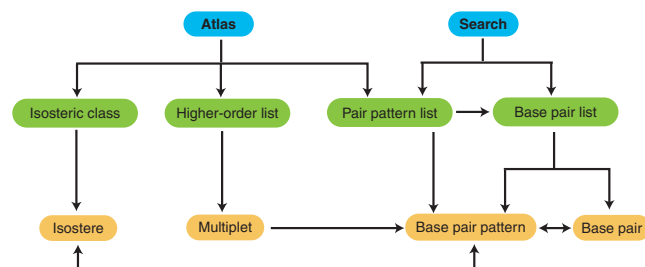| Context ID | Base location |
|---|---|
| H | RNA duplex interior, stacked between both 3′- and 5′-neighbors |
| $H_e$ | RNA duplex terminus, stacked against either 3′- or 5′-neighbor |
| $H_e^i$ | RNA duplex terminus, with sequentially distance base(s) stacked against either 3′- or 5′-neighbor(s) in same quasi-continuous array |
| $H_e^q$ | Terminus (5′- or 3′-end) of a broken single strand within an RNA duplex |
| $H^i$ | Single-base insertion (intercalator) within an RNA duplex |
| $H_l$ | Isolated base pair with no stacking interactions |
| $S_b$ | Single-stranded loop between RNA duplexes |
| $S_h$ | Single-stranded (hairpin) loop at either end of an RNA duplex |
| $S_h^q$ | Single-stranded (hairpin) loop at either end of a quasi-continuous helix |
| $S_l$ | Single-stranded linker between RNA duplexes |
| $S_e$ | Single-stranded chain terminus |



**Figure 3.** Overview of BPS content. The database is organized into two sections, Atlas and Search, which are interconnect by the base-pairing patterns. The Atlas contains galleries with representative images of base-pairing patterns, higher-order base interactions (multiplets) and isosteric base pairs. The Search module provides the interface to access specific base pairs via different input information.
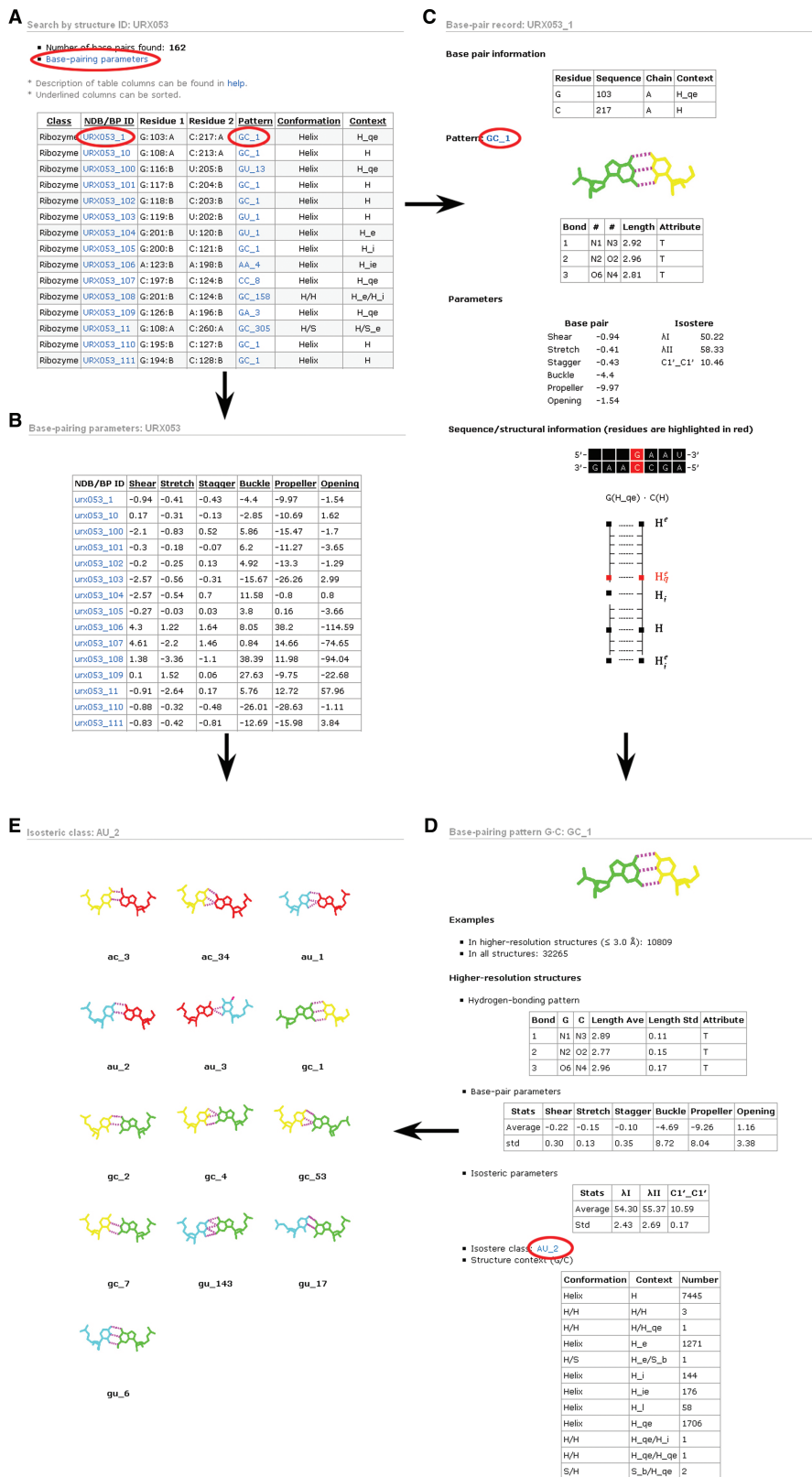
**Figure 4.** Screenshots illustrating representative results of a search for base pairs in the P4-P6 group I ribozyme domain structure (14). (**A**) Overview of chemical and secondary structural information for the first 15 of the 162 bp in the structure. (**B**) Sortable list of the rigid-body parameters characterizing the base pairs, obtained by clicking the parameter link highlighted by the large red ellipse in the upper left corner in (A). (**C**) The base-pair information page obtained by clicking the highlighted base-pair (URX053_1) in (A). (**D**) Summary of the base-pairing pattern (GC_1) of the highlighted residue, which can be accessed from (A) or (C). (**E**) The isosteric class that includes the base-pairing pattern in (D). Not all base-pairing patterns belong to an isosteric class. If a base pair is a member of such a class, the link is found in (D).

(see Figure 4D); and, if applicable, a page with atomic-level images of all base pair types in the same isosteric class as the selected pair (Figure 4E). The results can be sorted by NDB/3DNA ID, base-pairing pattern, conformation type, structural context and each of the base-pair parameters (i.e. the columns with underlined headers in Figure 4A and B). All four base-pair related web pages—namely, the rigid-body parameter page (Figure 4B), the base-pair summary page (Figure 4C), the base-pairing pattern page (Figure 4D) and the isosteric class page (Figure 4E)—can be accessed from the search result (Figure 4A).

### Atlas pages

The Atlas provides representative images and summaries of all classes of canonical and noncanonical base pairs, higher-order base interactions and isosteric base pairs in the database, offering an alternative way to gain access to the data other than performing a search. The base pairs are grouped by the Saenger patterns (13) and by the pattern numbers collected in our analysis of data (i.e. some of the interactions found in available high-resolution structures were not previously anticipated). The base-pairing patterns are of the form shown in Figure 4D. The observed multiplets include a numbers of triplets, quadruplets and quintuplets as well as a few even higher-order interactions. Each multiplet pattern has a single information page showing the base-pair components and the number of occurrences in the current database (Figure 5). The isosteric base pairs fall into the 17 isosteric classes identified with isosteric parameters described above. Each isosteric class has a webpage that illustrates each member with a representative image that links, in turn, to the corresponding page summarizing the base-pairing pattern (Figure 4E).
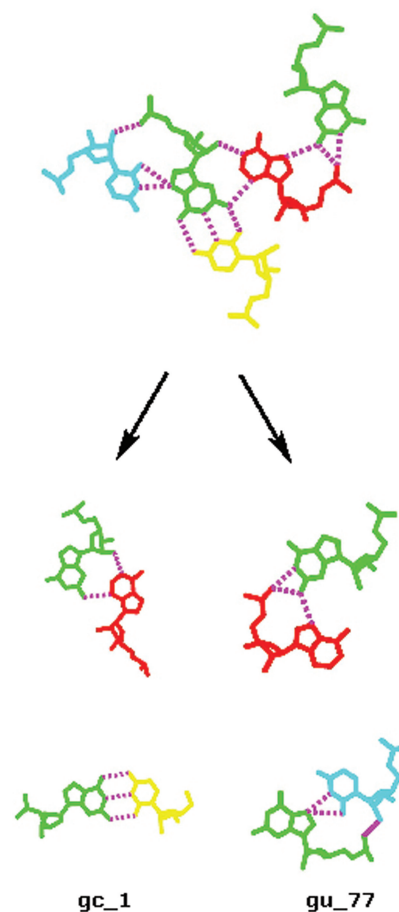
## CONCLUSIONS AND FUTURE WORK

The BPS offers a comprehensive compilation of the canonical and noncanonical base pairs in high-resolution RNA structures. The information includes a variety of features to characterize and use for gaining access to the data. The database combines descriptive and numerical information to annotate the base pairs, and is especially useful for the study of noncanonical interactions. Specifically, the BPS contains (i) quantitative spatial data (i.e. rigid-body parameters) that allow for the exact reconstruction and visualization of each interaction, (ii) the hydrogen-bond identities and distances and characteristic interbase distances and angles, (iii) the identities and quantitative characteristics of all isosteric base pairs, (iv) the corresponding information for all networks of three or more associated bases, (v) the sequential and structural context of each pair, (vi) the mean values and standard deviations of the quantitative parameters associated with each base-pair type and (vii) the capability to sort the data in different ways. The data can be accessed by browsing the Atlas or querying via Search. The base-pair coordinates, raw database files and scripts are available on request.

Our geometric characterization of RNA base pairs complements the well-known Leontis–Westhof (3)

classification of these interactions in terms of the edges (Watson–Crick, Hoogsteen/CH, Sugar) of associated bases and the orientation (*cis* or *trans*) of the glycosidic bonds and the topological groupings of noncanonical pairs introduced by Lee and Gutell (4). We plan in the near future to add the Leontis–Westhof classifiers using an algorithm of our own adapted from the RNAVIEW software (15). The new information will facilitate a more general classification of isosteres and help to group closely related base-pairing patterns. We are currently investigating the feasibility of including the rigid-body 'step' that characterize the spatial disposition of stacked bases and base pairs and thereby quantifying the structural-context information. In addition to the search functions for base pairs, we will include the capability to search multiplets and isosteres. We will update the database at regular



**Figure 5.** Example of a base quintuplet—here a U·G·C·A·G complex found in the crystal structures of the *H. marismortui* large ribosomal subunit (17), its complex with r(CCd)A-p-puromycin (18) and the complex of the same ligand at the peptidyl transferase center of the 50S ribosomal submit (19) (NDB_IDs: RR0011, RR0013 and RR0076, respectively)—and the constituent base pairs illustrated in the BPS Atlas.

intervals as new RNA structures are deposited in the NDB and PDB.

## REFERENCES

1. Holbrook,S.R. (2008) Structural principles from large RNAs. *Ann. Rev. Biophys.*, **37**, 445–464.
2. Crick,F.H.C. (1966) Codon-anticodon pairing: the wobble hypothesis. *J. Mol. Biol.*, **19**, 548–555.
3. Leontis,N.B., Stombaugh,J. and Westhof,E. (2002) The non-Watson-Crick base pairs and their associated isostericity matrices. *Nucleic Acids Res.*, **30**, 3497–3531.
4. Lee,J.C. and Gutell,R.R. (2004) Diversity of base-pair conformations and their occurrence in rRNA structure and RNA structural motifs. *J. Mol. Biol.*, **344**, 1225–1249.
5. Nagaswamy,U., Voss,N., Zhang,Z. and Fox,G.E. (2000) Database of non-canonical base pairs found in known RNA structures. *Nucleic Acids Res.*, **28**, 375–376.
6. Nagaswamy,U., Larios-Sanz,M., Hury,J., Collins,S., Zhang,Z., Zhao,Q. and Fox,G.E. (2002) NCIR: a database of non-canonical interactions in known RNA structures. *Nucleic Acids Res.*, **30**, 395–397.
7. Berman,H.M., Olson,W.K., Beveridge,D.L., Westbrook,J., Gelbin,A., Demeny,T., Hsieh,S.-H., Srinivasan,A.R. and Schneider,B. (1992) The Nucleic Acid Database: a comprehensive relational database of three-dimensional structures of nucleic acids. *Biophys. J.*, **63**, 751–759.
8. Lu,X.J. and Olson,W.K. (2003) 3DNA: a software package for the analysis, rebuilding and visualization of three-dimensional nucleic acid structures. *Nucleic Acids Res.*, **31**, 5108–5121.
9. Lu,X.J. and Olson,W.K. (2008) 3DNA: a versatile, integrated software system for the analysis, rebuilding, and visualization of three-dimensional nucleic-acid structures. *Nat. Protoc.*, **3**, 1213–1227.
10. Bernstein,F.C., Koetzle,T.F., Williams,G.J., Meyer,E.E. Jr, Brice,M.D., Rodgers,J.R., Kennard,O., Shimanouchi,T. and Tasumi,M. (1977) The Protein Data Bank: a computer-based archival file for macromolecular structures. *J. Mol. Biol.*, **112**, 535–542.
11. Berman,H.M., Westbrook,J., Feng,Z., Gilliland,G., Bhat,T.N., Weissig,H., Shindyalov,I.N. and Bourne,P.E. (2000) The Protein Data Bank. *Nucleic Acids Res.*, **28**, 235–242.
12. Olson,W.K., Bansal,M., Burley,S.K., Dickerson,R.E., Gerstein,M., Harvey,S.C., Heinemann,U., Lu,X.-J., Neidle,S., Shakked,Z. *et al.* (2001) A standard reference frame for the description of nucleic acid base-pair geometry. *J. Mol. Biol.*, **313**, 229–237.
13. Saenger,W. (1984) *Principles of Nucleic Acid Structure*. Springer, New York, p. 120.
14. Cate,J.H., Gooding,A.R., Podell,E., Zhou,K., Golden,B.L., Kundrot,C.E., Cech,T.R. and Doudna,J.A. (1996) Crystal structure of a group I ribozyme domain: principles of RNA packing. *Science*, **273**, 1678–1685.
15. Yang,H., Jossinet,F., Leontis,N., Chen,L., Westbrook,J., Berman,H. and Westhof,E. (2003) Tools for the automatic identification and classification of RNA base pairs. *Nucleic Acids Res.*, **31**, 3450–3460.
16. Dickerson,R.E., Bansal,M., Calladine,C.R., Diekmann,S., Hunter,W.N., Kennard,O., von Kitzing,E., Lavery,R., Nelson,H.C.M., Olson,W.K. *et al.* (1989) Definitions and nomenclature of nucleic acid structure parameters. *J. Mol. Biol.*, **208**, 787–791.
17. Ban,N., Nissen,P., Hansen,J., Moore,P.B. and Steitz,T.A. (2000) The complete atomic structure of the large ribosomal subunit at 2.4 Ångstrom resolution. *Science*, **289**, 905–920.
18. Nissen,P., Hansen,J., Ban,N., Moore,P.B. and Steitz,T.A. (2000) The structural basis of ribosome activity in peptide bond synthesis. *Science*, **289**, 920–930.
19. Hansen,J.L., Schmeing,T.M., Moore,P.B. and Steitz,T.A. (2002) Structural insights into peptide bond formation. *Proc. Natl Acad. Sci., USA*, **99**, 11670–11675.