

Figure 2. Histograms of IDIs between sets of identical (upper left), isosteric (upper right), near isosteric (lower left) and non-isosteric (lower right) base pair instances from the 3D structures in the reduced-redundancy dataset having better than 3.0 Å resolution. Upper left: IDIs calculated between identical base pairs (i.e. GC cWW with GC cWW, UA tWH with UA tWH, etc.). Upper right: IDIs between 200 GC cWW and 200 UA cWW pairs. Lower left: IDIs between 200 GC cWW and 200 GU cWW pairs. Lower right: IDIs between 200 GU cWW and 200 UG cWW pairs.

base pairs that are germane to their ability to substitute for one another in 3D structures. We evaluated the IDI by checking how it handles four distinct cases. First, the IDI should be lowest between instances of the same base pair. The upper-left panel of Figure 2 shows a histogram of the IDI calculated between each distinct pair of base pair instances of the same kind from 3D structures (i.e. GC cWW with GC cWW, UA tWH with UA tWH, etc.). The base pairs were drawn from the crystal structures in the reduced-redundancy dataset having a resolution better than 3.0 Å. At most 200 instances of each base combination from each geometric family were used to prevent the cWW base pairs from dominating the histogram. When more than 200 instances were available, 200 were selected randomly. The distribution peaks at IDI = ~0.6 and is narrowly distributed as it should be when comparing identical base pairs, regardless of geometric family.

Second, a quantitative measure of isostericity should be comparably low between base pairs classified as isosteric by qualitative criteria. In the upper-right panel of Figure 2 we show the IDI calculated between combinations of 200 GC and 200 UA cWW base pairs from the same 3D structures as described above. As in the first panel, the vast majority (over 96%) of comparisons result in an

IDI below 2.0 and the distribution peaks below 1.0 and is similar in shape. Third, the IDI should be larger for base pairs that are near isosteric and known to occasionally substitute for one another. The lower-left panel of Figure 2 shows the IDI between 200 GC cWW and 200 GU cWW base pairs. The peak of the histogram occurs to the right of 2.0 and the distribution is largely non-overlapping with the distributions for isosteric or identical base pairs in the upper panels of Figure 2. Finally, the IDI should be largest for base pairs which are geometrically dissimilar (non-isosteric). The lower-right panel of Figure 2 shows the IDI between 200 GU cWW pairs and 200 UG cWW pairs. This distribution peaks at IDI = ~4.5 and is largely non-overlapping with the others in Figure 2. When base pairs from different geometric families are compared, even larger IDI values are obtained, ranging up to 20. Histograms were also made using base pairs extracted from structures with 2.0 Å or better resolution. We observed that the corresponding IDI distributions from the 2.0 Å and 3.0 Å data peak within ~0.2 Å of each other. As expected, the 2.0 Å IDI distributions were narrower, with full width at half height ~0.5 Å vs. ~0.8 Å for the 3.0 Å data. Based on these and similar analyses for non-WC base pairs, we chose IDI threshold

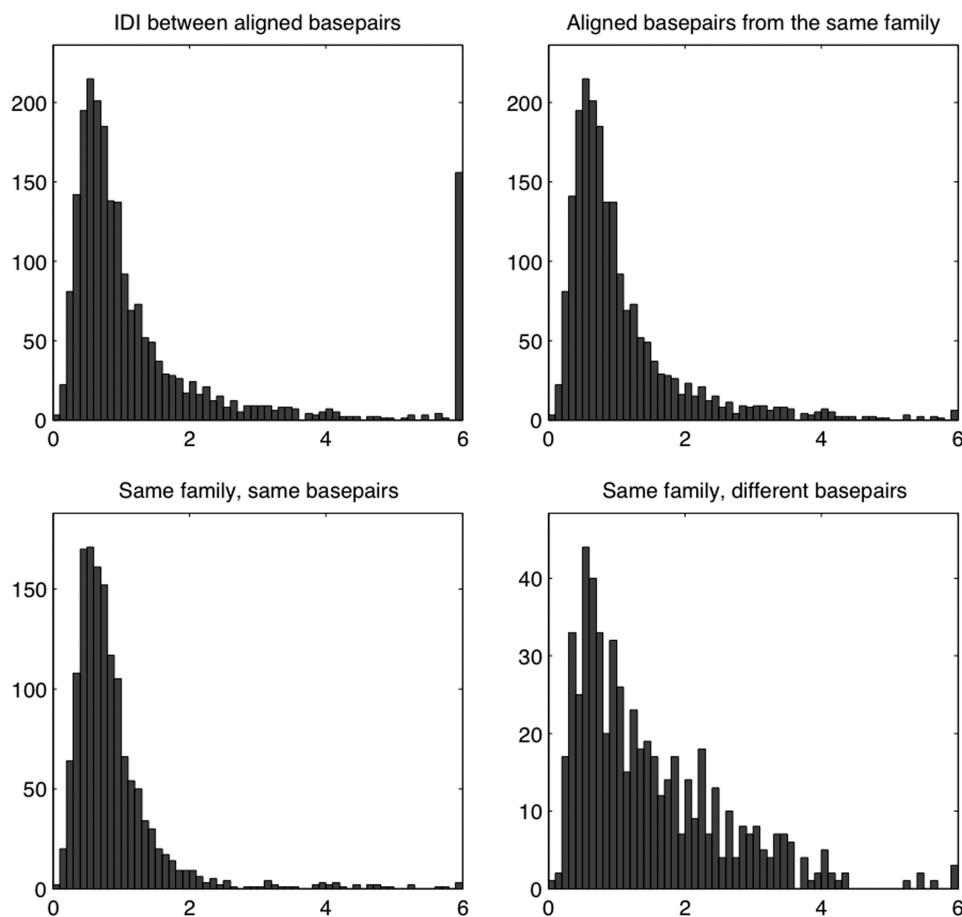


Figure 4. Histograms of IDIs between actual base pairs in the 3D–3D alignment of *E. coli* and *T. thermophilus* 5S, 16S and 23S rRNAs. The IDIs used in these histograms were calculated before the revision of the 3D structures to correct syn-anti errors. The upper-left panel shows the IDI between all aligned base pairs, whether in the same geometric family or not. The base pairs with IDI > 6.0 are discussed in section ‘Base pair discrepancies between aligned positions in the rRNA 3D structural alignments’. The upper-right panel shows the IDI between aligned base pairs that belong to the same geometric family, and the lower panels subdivide these into two cases, those in which with identical base combinations (lower left) and those with different base combinations (lower right). All IDI values above 6 are placed in the rightmost bin in each histogram.

Supplementary Table S8 provides data for all the other base pair families. The Materials and methods section should be consulted on how we removed redundant sequences from the multiple sequence alignments used to obtain base pair frequencies; how base pairs in the 3D structure were selected for analysis and how we estimated confidence intervals for the base pair occurrence frequencies we obtained from the sequences and the 3D structures.

The reliability of the base pair frequencies determined from the sequence alignments depends critically on the quality of the sequence alignments. Misleading results can be obtained if a sequence does not in fact contain the base pair type inferred from the 3D structures or if the sequence is not aligned correctly to the structure, so the wrong base or a gap (‘-’) is placed in one of the columns. We can estimate the extent to which the data are affected by such errors by examining the frequencies of base combinations in the sequences which cannot form an allowed base pair in the geometric base pair family that occurs at the corresponding site of the 3D structure. For example, the base combination GG is not allowed at

cWW base pairing positions, while CG, GA, GC, GU, UC and UU cannot occur at tHS sites. For the base pairs of the conserved core, the frequencies of non-allowed base combinations was <0.6% for all base combinations in all base pair families, with the exception of cWH AA (2.4%). In addition, very few gaps (<0.7%) occurred in the alignments at the positions included in the frequency analysis.

Base pair frequencies from rRNA. The sequences provide far more data than the 3D structures and so potentially provide more reliable estimates and narrower confidence intervals of base pair frequencies. These data, however, need some care in their interpretation. First, there are different numbers of sequences in the different multiple sequence alignment, 101 sequences for 5S rRNA, 717 sequences for 16S rRNA and 136 sequences for 23S rRNA. Second, the columns in the alignment corresponding to each base pair family in the conserved core may have gaps or letters other than A, C, G, U and these are counted as missing data for the purposes of this table.

We explain how we calculated the base pair frequencies from the sequence alignments using the tWH base pair

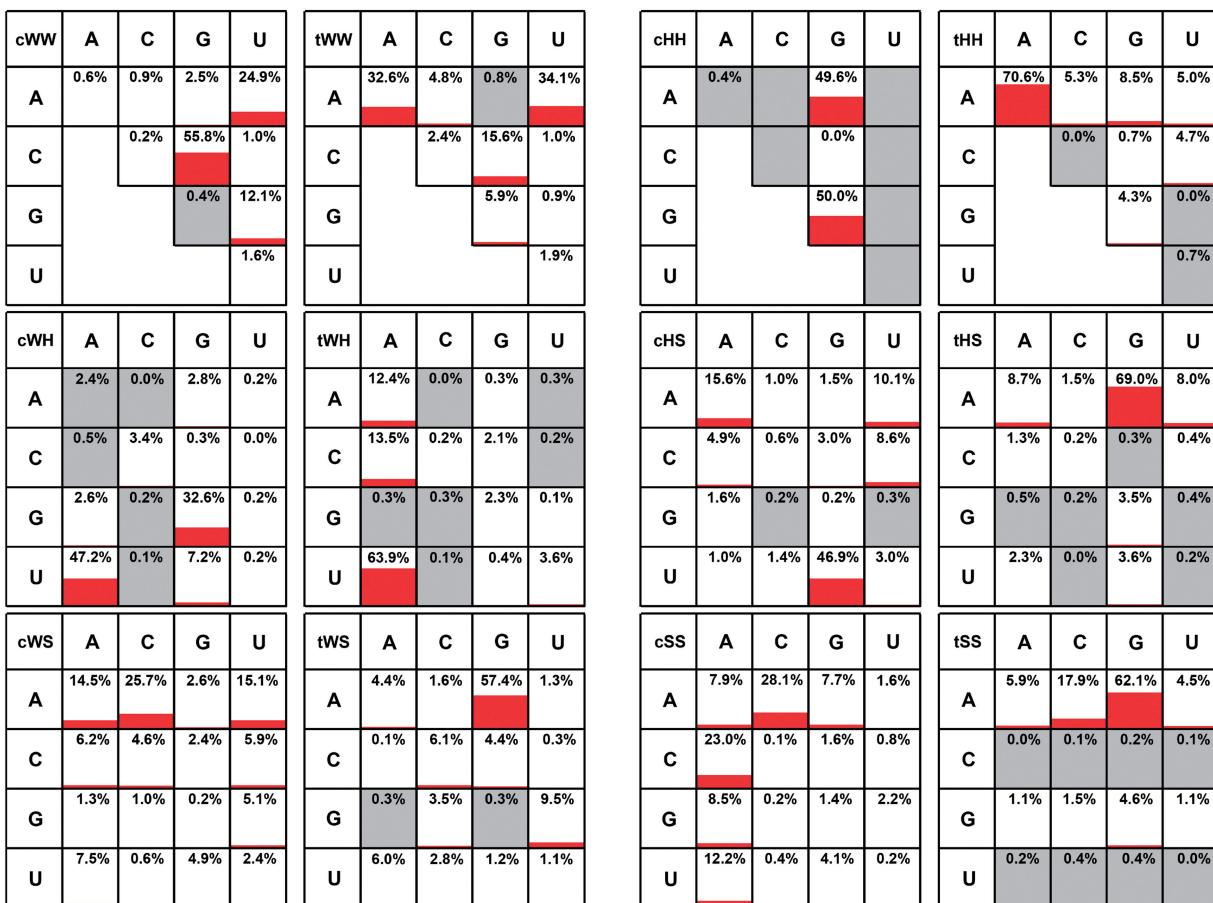


Figure 5. A graphical summary of the base pair occurrence frequencies within each base pair family, obtained from rRNA sequence data (data from Supplementary Table S8). For cWW, tHH, tWH, tHS, tWS and tSS, one base combination accounts for >50% of instances. The gray boxes in each matrix indicate base combinations that do not form that type of base pair. For example, there is no GG cWW base pair.

family as an example: For each of the 95 instances of tWH base pairs in the conserved core of the 5S, 16S and 23S rRNA 3D alignments, we calculated the frequency (as a percentage) of each base combination in the corresponding two columns of the multiple sequence alignment. For each location in the 3D alignment, these frequencies add to 100%. Then we averaged the 95 sets of frequencies, thus giving equal weight to each location of a tWH base pair in the 3D structures, rather than weighting the data by the number of sequences available at each location. Even though we have dramatically reduced the redundancy within the aligned sequences, statistical dependence exists between the base combinations in the alignment corresponding to each particular instance of a tWH base pair in the conserved core. The simultaneous 95% confidence intervals derived from sequences (Table 8, row 1) are somewhat narrower than the confidence intervals calculated from the *E. coli* rRNA structures (Table 8, row 2) or the *T. thermophilus* rRNA structures (Table 8, row 3), but not as narrow as those obtained from the reduced-redundancy set of structures (Table 8, row 4). This indicates that using data from the multiple sequence alignment raises the *effective* number of observations above the total number of base pairs of a particular family in the rRNA 3D structures (i.e. 105 tWH base

pairs in *T. thermophilus* 5S, 16S and 23S rRNA), but not as high as the *total* number of instances of that family in the reduced-redundancy 3D database (i.e. 519 tWH instances), or anywhere near the *observed* number of base combinations from the multiple sequence alignments (i.e. 8139 for tWH).

We provide a graphical summary of the base pair occurrence frequencies within each family, obtained from rRNA sequences, in Figure 5. The cWW, tWW, cHH and tHH families have symmetric base pairs; for example, each instance of GC cWW is also an instance of CG cWW. For this reason, we only display the data on upper right half of the matrices for these families. It is interesting to note that across the ribosomal structures, none of the base combinations in these four families show a 5' to 3' asymmetry due to order in the nucleotide sequence. For example, ~50% of the GC cWW base pairs have G occurring earlier in the nucleotide sequence than C, and 50% have C first in the nucleotide sequence.

DISCUSSION

We have defined the IDI to quantify base pair isostericity and to evaluate the usefulness of the isostericity concept for understanding non-WC base pairs and RNA 3D

19. Yang,H., Jossinet,F., Leontis,N., Chen,L., Westbrook,J., Berman,H. and Westbrook,E. (2003) Tools for the automatic identification and classification of RNA base pairs. *Nucleic Acids Res.*, **31**, 3450–3460.
20. Leontis,N.B., Stombaugh,J. and Westbrook,E. (2002) Motif prediction in ribosomal RNAs Lessons and prospects for automated motif prediction in homologous RNA molecules. *Biochimie*, **84**, 961–973.
21. Leontis,N.B. and Westbrook,E. (2003) Analysis of RNA motifs. *Curr. Opin. Struct. Biol.*, **13**, 300–308.
22. Nissen,P., Ippolito,J.A., Ban,N., Moore,P.B. and Steitz,T.A. (2001) RNA tertiary interactions in the large ribosomal subunit: the A-minor motif. *Proc. Natl Acad. Sci. USA*, **98**, 4899–4903.
23. Lescoute,A., Leontis,N.B., Massire,C. and Westbrook,E. (2005) Recurrent structural RNA motifs, Isostericity Matrices and sequence alignments. *Nucleic Acids Res.*, **33**, 2395–2409.
24. Sheridan,P.P., Freeman,K.H. and Brenchley,J.E. (2003) Estimated minimal divergence times of the major bacterial and archaeal phyla. *Geomicrobiol. J.*, **20**, 1–14.
25. Leontis,N.B., Altman,R.B., Berman,H.M., Brenner,S.E., Brown,J.W., Engelke,D.R., Harvey,S.C., Holbrook,S.R., Jossinet,F., Lewis,S.E. et al. (2006) The RNA Ontology Consortium: an open invitation to the RNA community. *RNA*, **12**, 533–541.