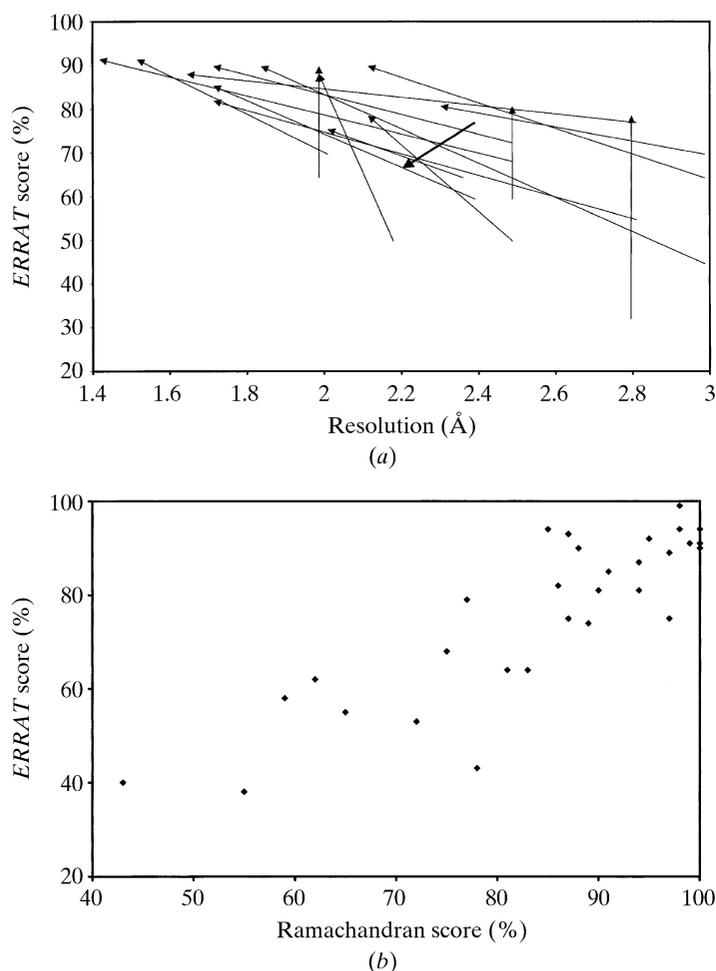21.3. DETECTION OF ERRORS IN PROTEIN MODELS





Fig. 21.3.5.5. Evaluation of old and revised models in a database survey by *ERRAT*. (*a*) Percentage of residues within the 95% confidence limit given by *ERRAT* as a function of the resolution to which the structure was determined. Each arrow represents an obsolete structure (arrow tail) and the revised structure that replaced it (arrow head). The revised structures (typically analysed at finer resolution) show markedly improved *ERRAT* scores. (*b*) Correlation between the percentage of a structure within the 95% confidence limit according to *ERRAT*, and the percentage of residues in the most favoured regions of a Ramachandran diagram according to *PROCHECK*.

by the San Diego Supercomputer group (http://pdbobs.sdsc.edu) served as a benchmark for evaluating the *ERRAT* program. For testing, 17 pairs of protein models were selected. Each pair comprised an obsolete entry and the revised model that replaced it. Using *ERRAT*, the overall quality of each model was expressed as a single number according to the fraction of the structure falling below the 95% confidence limit for rejection. The overall scores are significantly better for the revised structures, most of which were analysed at improved resolution (Fig. 21.3.5.5*a*). This result further demonstrates the utility of *ERRAT* for monitoring the model-

building process. Furthermore, a strong correlation is found between the percentage of residues within the 95% confidence limit given by *ERRAT* and the percentage of residues in the most favoured regions of the Ramachandran plot of *PROCHECK* (Fig. 21.3.5.5*b*). In general, the problematic regions detected by the two programs agree with each other.

### 21.3.6. Summary

In order to ensure the quality of the growing protein structure databases, models must be evaluated carefully during and after the structure determination process. Model evaluation can incorporate two types of measures: agreement between the model and the experimental diffraction data, and agreement between the model and the database of known structures. The latter types use the atomic coordinates of the final model, but do not rely on the diffraction data. In recent years, powerful methods of this type have been developed.

The most informative and reliable model-evaluation criteria are those that measure properties not optimized as part of the automatic refinement procedure. The free $R$ value has become important for monitoring the progress of atomic refinement for the same reason: it is based on reflections not included in refinement. We have focused here on two programs, *VERIFY*3D and *ERRAT*, which both evaluate high-level geometric properties not optimized during atomic refinement. Each offers the convenience of a single score over a sliding window along the protein sequence. Because *VERIFY*3D operates on the level of amino-acid residues, it is sensitive to errors on that scale, particularly those that affect the distribution of polar and nonpolar residues. *ERRAT* operates on the atomic level and has proven to be particularly useful for pinpointing local regions of protein models that require further adjustments. When used in combination, these methods and others can help crystallographers produce more accurate structural models of proteins.

### 21.3.7. Availability of software

The programs *ERRAT* and *VERIFY*3D are available on the World Wide Web for non-commercial applications. The URL for *VERIFY*3D is http://www.doe-mbi.ucla.edu/Services/Verify3D.html and the URL for *ERRAT* is http://www.doe-mbi.ucla.edu/Services/Errat.html. *VERIFY*3D and *ERRAT* expect a coordinate file in PDB format. The programs return plots of the type shown in this chapter.

### Acknowledgements

references