

*Sequence analysis***Clustal W and Clustal X version 2.0**

M.A. Larkin¹, G. Blackshields¹, N.P. Brown³, R. Chenna³, P.A. McGettigan¹, H. McWilliam⁴, F. Valentin⁴, I.M. Wallace¹, A. Wilm¹, R. Lopez⁴, J.D. Thompson², T.J. Gibson³ and D.G. Higgins^{1,*}

¹The Conway Institute of Biomolecular and Biomedical Research, University College Dublin, Belfield, Dublin 4, Ireland, ²Laboratoire de Biologie et Genomique Structurales, Institut de Génétique et de Biologie Moléculaire et Cellulaire, Illkirch, France, ³European Molecular Biology Laboratory, Meyerhofstrasse 1, 69117 Heidelberg, Germany and ⁴EMBL Outstation-European Bioinformatics Institute, Wellcome Trust Genome Campus Hinxton, Cambridge, CB10 1SD, UK

Received on June 27, 2007; revised on August 3, 2007; accepted on August 3, 2007

Advance Access publication September 10, 2007

Associate Editor: Alex Bateman

ABSTRACT

Summary: The Clustal W and Clustal X multiple sequence alignment programs have been completely rewritten in C++. This will facilitate the further development of the alignment algorithms in the future and has allowed proper porting of the programs to the latest versions of Linux, Macintosh and Windows operating systems.

Availability: The programs can be run on-line from the EBI web server: <http://www.ebi.ac.uk/tools/clustalw2>. The source code and executables for Windows, Linux and Macintosh computers are available from the EBI ftp site <ftp://ftp.ebi.ac.uk/pub/software/clustalw2/>

Contact: clustalw@ucd.ie

1 INTRODUCTION

Multiple sequence alignments are now one of the most widely used bioinformatics analyses. They are needed routinely as parts of more complicated analyses or analysis pipelines and there are several very widely used packages, e.g. Clustal W (Thompson *et al.*, 1994), Clustal X (Thompson *et al.*, 1997), T-Coffee (Notredame *et al.*, 2000), MAFFT (Katoh *et al.*, 2002) and MUSCLE (Edgar, 2004). Clustal is also the oldest of the currently most widely used programs having been first distributed by post on floppy disks in the late 1980s. It was initially written in Microsoft Fortran for MS-DOS and originally ran on IBM compatible personal computers as four separate executable programs, Clustal1–Clustal4 (Higgins and Sharp, 1988, 1989). These were later rewritten in C and merged into a single program, Clustal V (Higgins *et al.*, 1992), that was distributed for VAX/VMS, Unix, Apple Macintosh and IBM compatible PCs. These programs were distributed from the EMBL File server (Stoehr and Omond, 1989), an e-mail and FTP server, based at the EMBL in Heidelberg, Germany.

The current Clustal programs all derive from Clustal W (Thompson *et al.*, 1994), which incorporated a novel position-specific scoring scheme and a weighting scheme for down weighting over-represented sequence groups. The ‘W’

stands for ‘weights’. These programs have been amended and added to many times since 1994 in order to increase functionality and to increase sensitivity. The user-friendliness has also been greatly enhanced by the addition, in 1997, of a full graphical user interface (Thompson *et al.*, 1997). This has made the code complicated to maintain and develop, as the graphical interface must be constantly modified and recompiled for new operating systems and desktop environments (Windows, Macintosh, VMS, Unix and Linux).

By the late 1990s, Clustal W and Clustal X were the most widely used multiple alignment programs. They were able to align medium-sized data sets very quickly and were easy to use. The alignments were of sufficient quality not to require manual editing or adjustment very often. This situation changed greatly with the appearance of the first custom made benchmark test set for multiple alignment programs, BALiBASE (Thompson *et al.*, 1999). This was followed by the appearance of T-Coffee which was able to make very accurate alignments of very divergent proteins but only for small sets of sequences, given its high computational cost. With the increase in processing speed of desktop computers, and subsequent optimisation of the T-Coffee code, the latter is now practical for routine use on moderately sized alignment problems. More recently, MAFFT and MUSCLE appeared; which were, initially, at least as accurate as Clustal, in terms of alignment accuracy, but which were also extremely fast; and able to align many thousands of sequences. Over the past 4 or 5 years, these programs have also gradually become more and more accurate with difficult alignments. Nonetheless, Clustal W and Clustal X continue to be very widely used, increasingly on websites. The EBI Clustal site, gets literally millions of multiple alignment jobs per year.

It is in this context that we developed Clustal W 2.0 and Clustal X 2.0. These programs were rewritten in C++ with a simple object model in order to make it easier to maintain the code and more importantly, to make it easier to modify or even replace some of the alignment algorithms. We have produced two new programs which are very similar in look and feel to the older version 1.83 programs but which can now be managed more easily. We have also made some minor adjustments to the

*To whom correspondence should be addressed.

