# Artemis: sequence visualization and annotation

Kim Rutherford[1], Julian Parkhill[1], James Crook[2], Terry Horsnell[2], Peter Rice[1], Marie-Adèle Rajandream[1] and Bart Barrell[1]

[1]The Sanger Centre, Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SA, UK and [2]The MRC Laboratory of Molecular Biology, Hills Road, Cambridge, CB2 2QH, UK

## Abstract

*Summary: Artemis is a DNA sequence visualization and annotation tool that allows the results of any analysis or sets of analyses to be viewed in the context of the sequence and its six-frame translation. Artemis is especially useful in analysing the compact genomes of bacteria, archaea and lower eukaryotes, and will cope with sequences of any size from small genes to whole genomes. It is implemented in Java, and can be run on any suitable platform. Sequences and annotation can be read and written directly in EMBL, GenBank and GFF format.*

*Availability: Artemis is available under the GNU General Public License from http://www.sanger.ac.uk/Software/Artemis*

*Contact: kmr@sanger.ac.uk*

Visualization of features within a sequence is an essential component of modern molecular biology. Some commercial programs exist to meet this need, but most freely available programs are either designed for large eukaryotic genomes (such as AceDB; http://www.acedb.org/) or have limited capabilities. We have developed Artemis, a Java-based sequence visualization tool, for our own in-house analysis. It is used extensively in our laboratories for the annotation of bacterial and lower eukaryotic genomes, most recently for the annotation of the complete genomes of *Campylobacter jejuni* and *Neisseria meningitidis*, and the chromosomes of *Plasmodium falciparum*. Artemis grew out of the concepts developed in DIANA (for DIsplay and ANAlyse; J.C., T.H. & P.R., *unpublished*), which was previously used at the MRC-LMB and the Sanger Centre for sequence analysis.

Artemis can be used simply as a sequence viewer, on any platform running Java (Unix-, Macintosh- or PC-based), and will allow the visualization of sequence and annotation taken directly from EMBL (Baker *et al.*, 2000) and GenBank (Benson *et al.*, 2000) format files. Properties of the sequence, such as G+C content, G/C skew (Lobry, 1996), frame-specific G+C content (Bibb *et al.*, 1984), codon usage etc. can be directly plotted against the sequence and its features. Each plot allows dynamic modification of the window size used for the calculation, and the sequence and plots can be zoomed together into the single base level or out to the complete genome. Two sequence windows exist which can be used to view the same sequence at different zoom levels simultaneously. Properties of individual protein features, such as hydrophobicity and hydrophilicity can also be viewed in the same way. Artemis can also be run as a WWW applet, allowing the serving of sequence and annotation over the web in a interactive format.
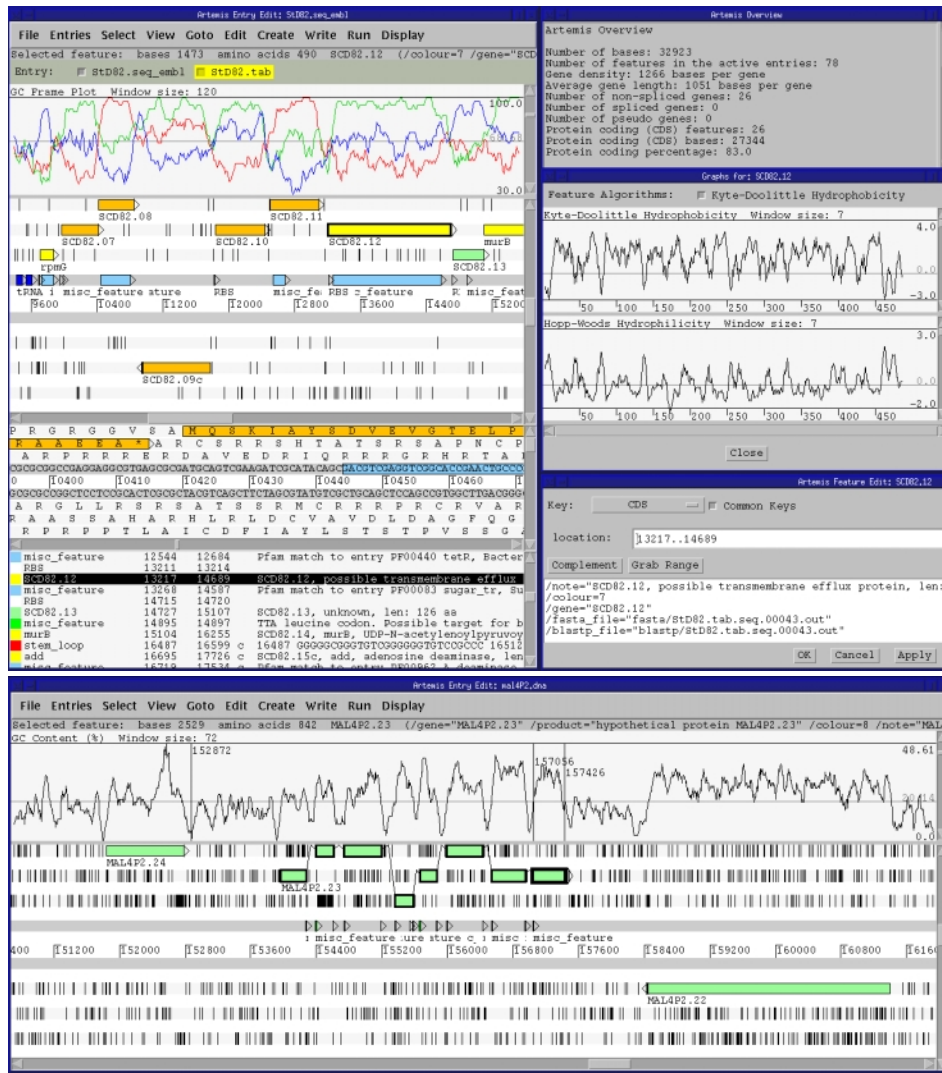
Artemis is specifically designed to be used as an annotation tool. In this mode, the results of any external analysis program can be parsed into the correct format (EMBL/GenBank feature table format; http://www.ebi.ac.uk/embl/Documentation/User_manual/format.html) and be overlaid on the sequence. These features can then be modified and moved between files to build up a single annotation table. External programs, such as FASTA (Pearson and Lipman, 1988) and BLAST (Altschul *et al.*, 1990) can be run on any number of individual features, and the results viewed from within the program. The qualifiers attached to each feature can be edited directly, and users can add any number of custom qualifiers. Because Artemis is designed to be as portable as possible, and to display the results of as wide a variety of analyses as possible, complex algorithms have not been built into the program itself. Sample parsing scripts for commonly used programs are available as part of the distribution.

## Acknowledgements

## References

Altschul,S.F., Gish,W., Miller,W., Myers,E.W. and Lipman,D.J. (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.

**Fig. 1.** Top left; the main Artemis windoow showing a section of the *Streptomyces coelicolor* cosmid EMBL:SCD82. The upper window is a frame-specific G+C plot, the second is the main sequence display with three forward reading frames, forward and reverse DNA lines separated by a scale bar, and three reverse reading frames. Vertical bars represent stop codons, open boxes show features. The third window is the secondary sequence display zoomed into the nucleotide level. The bottom window is a feature list. Top right; the sequence overview, showing the results of global calculations on the sequence. Middle right; protein property plots for one coding sequence. Bottom right; the feature edit window: Artemis makes all legitimate keys and qualifiers available from drop-down menus. Bottom; G+C plot and main sequence display for a highly spliced gene from *Plasmodium falciparum* chromosome 4 (EMBL:PFMAL4P2).

Baker,W., van den Broek,A., Camon,E., Hingamp,P., Sterk,P., Stoesser,G. and Tuli,M.A. (2000) The EMBL nucleotide sequence database. *Nucleic Acids Res.*, **28**, 19–23.

Benson,D.A., Karsch-Mizrachi,I., Lipman,D.J., Ostell,J., Rapp,B.A. and Wheeler,D.L. (2000) GenBank. *Nucleic Acids Res.*, **28**, 15–18.

Bibb,M.J., Findlay,P.R. and Johnson,M.W. (1984) The relationship between base composition and codon usage in bacterial genes and its use for the simple and reliable identification of protein-coding sequences. *Gene*, **30**, 157–166.

Lobry,J.R. (1996) Asymmetric substitution patterns in the two DNA strands of bacteria. *Mol. Biol. Evol.*, **13**, 660–665.

Parkhill,J., Wren,B.W., Mungall,K., Ketley,J.M., Churcher,C., Basham,D., Chillingworth,T., Davies,R.M., Feltwell,T., Holroyd,S., Jagels,K., Karlyshev,A.V., Moule,S., Pallen,M.J., Penn,C.W., Quail,M.A., Rajandream,M.A., Rutherford,K.M., van Vliet,A.H., Whitehead,S. and Barrell,B.G. (2000) The genome sequence of the food-borne pathogen *Campylobacter jejuni* reveals hypervariable sequences. *Nature*, **403**, 665–668.

Pearson,W.R. and Lipman,D.J. (1988) Improved tools for biological sequence comparison. *Proc. Natl. Acad. Sci. USA*, **85**, 2444–2448.