

PROTOGENE: turning amino acid alignments into bona fide CDS nucleotide alignments

Sébastien Moretti, Frédéric Reinier, Olivier Poirot, Fabrice Armougom, Stéphane Audic, Vladimir Keduas and Cédric Notredame*

Information Génomique et Structurale, CNRS UPR2589, Institute for Structural Biology and Microbiology (IBSM), Parc Scientifique de Luminy, 163 Avenue de Luminy, FR 13288, Marseille cedex 09, France

Received February 14, 2006; Revised and Accepted March 20, 2006

ABSTRACT

We describe Protogene, a server that can turn a protein multiple sequence alignment into the equivalent alignment of the original gene coding DNA. Protogene relies on a pipeline where every initial protein sequence is BLASTed against RefSeq or NR. The annotation associated with potential matches is used to identify the gene sequence. This gene sequence is then aligned with the query protein using Exonerate in order to extract a coding nucleotide sequence matching the original protein. Protogene can handle protein fragments and will return every CDS coding for a given protein, even if they occur in different genomes. Protogene is available from <http://www.tcoffee.org/>.

INTRODUCTION

Although they constitute the material with which primary biological databases are made of, nucleotide sequences are rarely used when it comes to analyzing proteins. The reason is that evolutionary models designed for comparing nucleic acids are often too simplistic and almost never take into account the constraints associated with the coding nature of gene sequences. In practice, biologists dealing with proteins are often encouraged to use protein databases and associated tools to build their models. However, the transposition of these results onto the bona fide nucleotide sequences is time consuming. This limitation can be an issue, especially when reconstructing phylogenetic trees of closely related species or when looking for conserved nucleotide patterns within multiple coding sequences.

In theory, the task of turning a protein multiple sequence alignment (MSA) into the associated CDS (CoDing Sequence) MSA is trivial. Yet in practice things can prove more

complicated for a variety of simple reasons: unknown gene names, MSAs of domain or partial protein sequence, incomplete database annotation, and the ever faster evolution of genomic resources. Of course, each of these problems can usually be solved manually, on a case by case basis, but altogether they tend to hamper the establishment of automatic procedures for seamlessly connecting the protein and the nucleotide worlds.

When comparing protein and nucleotide sequences, efficient methods exist to either align CDSs using their coding potential (1–3) or thread CDSs onto a pre-established protein sequences (protal2dna: <http://bioweb.pasteur.fr/seqanal/interfaces/protal2dna.html> and pal2nal (<http://www.bork.embl-heidelberg.de/pal2nal/>)) but all these tools require the user to preprocess the data, gather the appropriate CDSs and make sure these are compatible with any subdomain extracted from the original protein sequence. To the best of our knowledge, no tool is available online to automatically identify the nucleotide sequence (genomic or transcript) associated with a protein partial sequence (domain or fragment) and process it to replace that protein with its bona fide CDS while retaining the original alignment.

We developed a fully automated program named Protogene (PROtein TO GENE) that when given a protein sequence alignment (pairwise or multiple) returns the corresponding CDS alignment. Protogene searches RefSeq (4) and NR with BLAST (5) in order to identify the transcript or genomic sequence(s) most likely to be associated with the original protein sequences. The sequences thus identified are processed by Exonerate (6) in order to extract a portion of CDS matching perfectly the original protein sequence. This CDS is re-introduced within the MSA to replace the original protein. Although every attempt is made to identify the genuine protein CDS, a conservative post-filtering step is still needed to eliminate any sequence that may not have been correctly processed. Protogene is not a gene finding tool and depends entirely on the pre-established proteomes found in RefSeq and NR. We have put the emphasis on robustness and reliability rather than

*To whom correspondence should be addressed. Tel: +33 491 825 427; Fax: +33 491 825 420; Email: cedric.notredame@europe.com

exhaustiveness. In our hands, Protogene returns a bona fide CDS for 95% of the sequences and fails on ~5%. Whenever the missed sequences matter, it is up to the user to explore the labyrinth of nucleotide databases and discover the glitch that breaks the information chain between the protein and its CDS. Protogene is available on <http://www.tcoffee.org/>.

METHODS

Server pipeline

Figure 1 shows Protogene's flow chart. Each sequence within the provided MSA is treated individually. The first step is a BLASTp (5) against the RefSeq (4) protein database. Matches against RefSeq are only accepted when associated with an alignment that displays 100% identity and 100% coverage with the query sequence (i.e. 100% of the query sequence residues aligned with identical residues). Multiple hits with

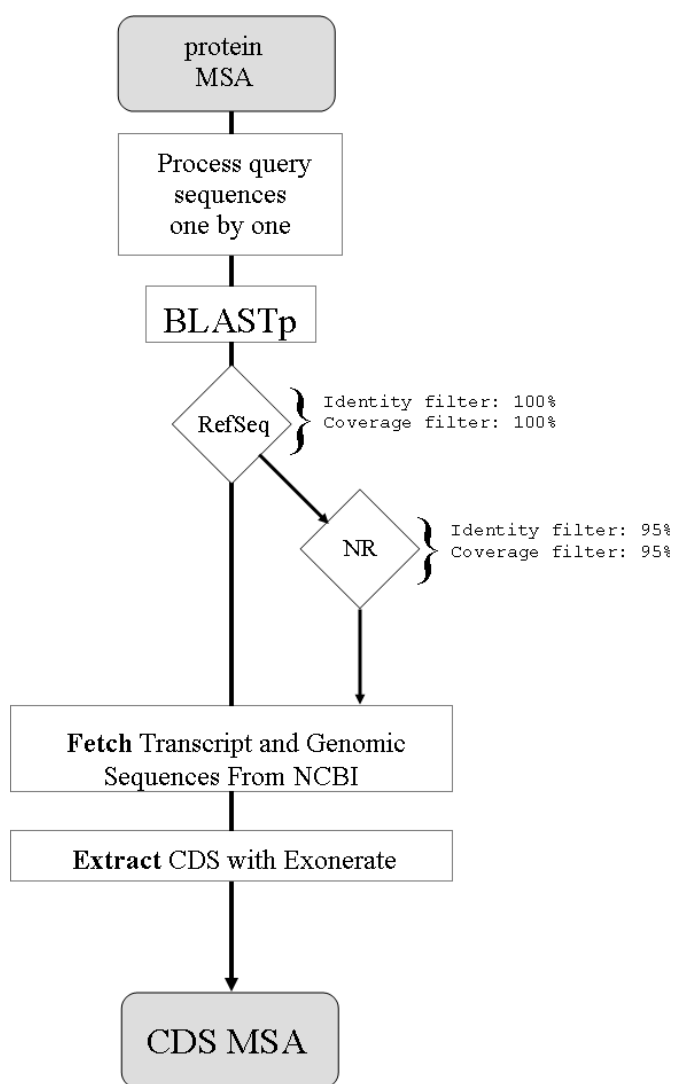


Figure 1. Protogene flow chart sequences are first BLASTed against RefSeq. If no match is found, they are then BLASTed against NR. Nucleotide sequences are fetched from NCBI and processed with Exonerate to yield CDSs that perfectly match the original protein.

the same score are all kept. If no suitable hit is found in RefSeq, the query sequence is BLASTed against the NR non-redundant protein sequence database, with a lowered acceptance threshold (95% identity, 95% coverage). NCBI EFetch and EUtils utilities (http://eutils.ncbi.nlm.nih.gov/entrez/query/static/eutils_help.html) are then used to fetch the nucleotide sequences associated with the query sequence. Exonerate is finally used to extract a CDS matching perfectly the original protein sequence (be it a full-length protein or just a fragment). Exonerate is able to splice introns and can handle genomic and transcript sequences alike. Mismatches between the query and the CDS are indicated with NNN codons. It is up to the user to decide whether these are sequences errors, polymorphisms or identification errors. The user may also request automatic back-translations using the IUC ambiguity code for unmatched amino acids. CDSs are returned along with some basic annotation including the nucleotide sequence accession number, the source organism and the RefSeq/NR accession number. CDSs where >5% of the nucleotides have had to be replaced with Ns are considered unreliable and discarded with an explicit mention in the output.

It is important to point out that when several proteins from different organisms have a perfect identity with the original query, each corresponding nucleotide sequence is integrated within the final alignment, leaving it to the user to remove the nucleotide sequences he is not interested in (see the Histone example in the next section). To ease this selection, the final alignment is reported in FASTA format. Nucleotide sequences gathered from the NCBI are kept in cache for 2 weeks, thus insuring faster second runs when reanalyzing a dataset with minor modifications.

Distribution

The Protogene server is available from <http://www.tcoffee.org/>. It relies on a collection of Perl scripts and two external programs: BLAST and Exonerate. BLAST searches against NR and RefSeq are made using the gigablast service (<http://www.igs.cnrs-mrs.fr/>).

USING PROTOGENE

We provide three simple examples of how Protogene can be used to rapidly and efficiently ask simple questions regarding protein sequence conservation at the nucleotide level. The first one is an analysis of the CLP Serine protease family. Serine is the only amino acid coded by two sets of codons (UCN and AGY) that cannot be interconverted by a single point mutation. Class switches, appear however to be very frequent and the question whether they arise from double mutations has been the subject of intense scrutiny and debate (7). Given a collection of protein sequences, Protogene makes it straightforward to analyze the codon conservation of the serine. Figure 2 shows the output obtained after cutting and pasting the Pfam (8) seed MSA of the CLP Serine Protease family (PF00574). This alignment is not made of complete proteins but restricted to the serine protease domain. For each domain, Protogene managed to identify at least one corresponding gene and also reported identical protein sequences coming from closely related genomes. Figure 2 represents the portion of this alignment containing the conserved Serine that is part of the catalytic triad. In a second test, we evaluated Protogene's

```

Q8W605_9CAUD_17-189_NC_003291      ---CCAGCTAAAA TCAGTATTTATATTGAT ---GCTCTTGC TGGAGTGGTGCATCTGT
Q9B0E1_9CAUD_16-186_NC_007053      ---CCTGCAAAAA TTAATATCTATGTCGAT ---GCCTTAGC GGCATCAATTGCTAGTGT
Q9B0E1_9CAUD_16-186_NC_003923      ---CCTGCAAAAA TTAATATCTATGTCGAT ---GCCTTAGC GGCATCAATTGCTAGTGT
Q9B0E1_9CAUD_16-186_NC_002951      ---CCTGCAAAAA TTAATATCTATGTCGAT ---GCTTTAGC GGCATCAATTGCTAGTGT
Q8ZQG9_SALTY_47-225_NC_003197      ---CCGGCGAGAAA AAGTGGTGTATGTGGAC ---GGTGTGGC CTGTCGATGGCGTCCGGT
Q9EYD3_ECO57_6-171_NC_002695      ---GGGGCGACCA TTACCGTGTATGTGGAT ---GGCGTTGC CGCCTCGATGGCATCTGT
Q9EYD3_ECO57_6-171_NC_002655      ---GGGGCGACCA TTACCGTGTATGTGGAT ---GGCGTTGC CGCCTCGATGGCATCTGT
Q8P2G1_STRP3_13-175_NC_003485      ---CCAGGCAATG TTGAAGTTGTTATCACA ---GGACTAGC TGCAAGCATTGCTAGCAT
Q8P2G1_STRP3_13-175_NC_004584      ---CCAGGCAATG TTGAAGTTGTTATCACA ---GGACTAGC TGCAAGCATTGCTAGCAT
Q8P2G1_STRP3_13-175_NC_004070      ---CCAGGCAATG TTGAAGTTGTTATCACA ---GGACTAGC TGCAAGCATTGCTAGCAT
Q8P2G1_STRP3_13-175_NC_004606      ---CCAGGCAATG TTGAAGTTGTTATCACA ---GGACTAGC TGCAAGCATTGCTAGCAT
Q9CFR9_LACLA_7-176_NC_004746       ---AATGGTAAAC CTGTAAC TGTAAATATTCAAGGGTTGGCAGCATCTGCAGCATCGGT
Q9CFR9_LACLA_7-176_NC_002667       ---AATGGTAAAC CTGTAAC TGTAAATATTCAAGGATTGGCAGCATCTGCAGCTTCAGT
Q9CFR9_LACLA_7-176_NC_002662       ---AATGGTAAAC CTGTAAC TGTAAATATTCAAGGATTGGCAGCATCTGCAGCTTCAGT
Q9XJT4_BPD3_38-211_NC_002484       ---AAGGGCAAGG TCACGGTGAACATCATC ---GGCCTGGC TGCCCTGCCCGCTCTTT
Q8LTH1_9CAUD_11-173_NC_005356      ---AAAGGCAAGG TGAATGTTGATCACA ---GCAATAGC AGCAAGTGGCGCATCGCT
CLPP_PSINU_13-197_NC_003386         ---GTACCAGATG TACATACTATTTGCATG ---GGATTAGC TGCTCAATGGGATCTTT
CLPP_CHAGL_13-197_NC_004115         ---GAACCAGAAA TTAGAACAATATGTATG ---GGAGTTGC TGCTCAATGGGTTCTTT
CLPP_SPIOL_13-196_NC_002202         ---CGACCAGATG TACATACTATTTGCATG ---GGATTAGC TGCTCAATGGGATCTTT
CLPP_WHEAT_13-197_NC_002762        ---ACACCAGATA TATATACAATATGCTCT ---GGAATAGC CGCCTCCATGGCATCCTT
CLPP_PINCO_13-197_L28807            ---GTACCAGATG TAAATACAATATGCTATG ---GGGGTAGC TGCTCAATGGGATCTTT
CLPP_CHLVU_13-196_NC_001865         ---AAATCCGAAG TCACAACGATTTGTGTT ---GGAACAGC AGCTTCGATGGCTCTTT
CLPP_NEPOL_13-196_NC_000927         ---ATGCAAGATG TGACGACGATTTGCGTA ---GGTATTGC CAGCATCTATGGCTCCCT
CLPP_MESVI_13-196_NC_002186         ---AACGTTGATG TTAACAATATTTGTATG ---GGTTAGC  TGCTCTATGGCTCTTT
CLPP_OENHO_19-204_NC_002693         ---GCACCTCCTG TGTATACACTAGGCCTG ---GGGGTACT CGCTCAATGGCATCCTT
CLPP1_MYCTU_11-193_NZ_AAIX01000031 ---CCCTGTGACA TCGCCACCTACGCGATG ---GGCATGGC CGCCTCCATGGCGAGTT
CLPP1_MYCTU_11-193_NZ_AAKR01000078 ---CCCTGTGACA TCGCCACCTACGCGATG ---GGCATGGC CGCCTCCATGGCGAGTT
CLPP2_CORGL_13-195_NC_006958        ---CCATGCGACA TCGCAACCTACGGCATG ---GGCCTGGC AGCATCCATGGGCCAGTT
CLPP1_STRCO_32-208_NC_003888        ---AAGAACGACG TGGTGACGATCGCGATG ---GGTCTCGC GGCCTCCATGGGACAGTT
CLPP_BUCAP_16-197_NC_004061         ---AAACCCGATG TTAACAATATTTGTATA ---GGACAAGC GTGTCAATGGCCATT
Q8VQM6_BACTU_12-193_AF454758       ---AAACCGATG  TGCAAAACGCTGTGCATG ---GGCTTTGC GGCATCAATTGGTGCAAT
CLPP2_BACHD_13-194_NC_002570        ---AAGCCGTGCA TTCATACGATTTGCACA ---GGTATGGC TGCTCCATTGGCCGCAAT
CLPP_LISIN_12-193_NC_003212         ---AAAGCTGACG TGCAAAC TATCGGTATG ---GGGATGGC TGCTCCATGGGCTCATT
CLPP1_BACHD_12-193_NC_002570        ---AAACCAACG TCTCAACCATTTGCATC ---GGGATGGC CGCTCAATGGGAGCCTT
CLPP_STAAW_12-193_NC_002953         ---AAACCTGATG TCAAAACAATTTGTATC ---GGTATGGC TGCTCAATGGGATCATT
CLPP_STAAW_12-193_NC_002951         ---AAACCTGATG TCAAAACAATTTGTATC ---GGTATGGC TGCTCAATGGGATCATT
CLPP_STAAW_12-193_NC_007622         ---AAACCTGATG TCAAAACAATTTGTATC ---GGTATGGC TGCTCAATGGGATCATT
CLPP_CLOAB_12-193_NC_003030         ---AAACCCGATG TCAAAACAATTTGTATA ---GGAATGGC TGCTCAATGGGGTCTATT
CLPP_CLOPE_13-194_NC_003366         ---AAGCCTGACG TATCTACAATCTGTATA ---GGTATGGC TGCTCTATGGGAGCATT
CLPP_THETN_12-193_NC_003869         ---AAGCCGGACG TTGTGACACTTTGTGTG ---GGCATGGC AGCATCTATGGCTGCTTT
CLPP_CAMJE_11-192_NC_003912         ---AAACCTGATG TTTGTACGATTTGCATA ---GGACAAGC TGCTCTATGGGAGCATT
CLPP_CAMJE_11-192_NC_002163         ---AAACCTGATG TTTGTACGATTTGTATA ---GGACAAGC TGCTCTATGGGAGCATT
CLPP_HELPJ_12-193_NC_000921         ---CGCCCTGATG TTTCCACGATTTGCATC ---GGTCAAGC GGCTCTATGGGGCGT
CLPP_THEMA_21-202_NC_000853         ---AAGTGTGATG TCTCAACCATATGTGTA ---GGACAGGC GGCTCCATGGCGGCTGT
CLPP_AQUAE_19-200_NC_000918         ---AAACCCGACG TGGTTACTATATGCTATG ---GGACAGGC GGCTCCATGGGAGCAAT
CLPP1_MYXXA_11-192_AF013216         ---AAGTGTCCGG TGTCCACCATCTGTGTG ---GGGCAGGC GGCTCCATGGGCGCGCT
CLPP_YEREN_25-206_U55059           ---AAGCCGGATG TCAGCAGCATTGTATG ---GGCCAGGC ATGTCAATGGGTGCATT
CLPP1_PSEAE_28-209_NC_002516        ---AAGCCCAACG TCTCGACCACCTGTATC ---GGCCAGGC GTGCAGCATGGGTGCCCT
CLPP1_PSEAE_28-209_NZ_AABQ07000002 ---AAGCCCAACG TCTCGACCACCTGTATC ---GGTCAAGC GTGCAGCATGGGTGCCCT
CLPP_XYLFA_19-200_NC_002488         ---AAACCTGCTG TACAGTACTATCTGTGT ---GGTCAAGC TGCTCTATGGGGCGT
CLPP_HAEIN_12-193_NC_000907         ---AAGCCAGATA TTCGCACTCTTTGTATT ---GGTCAAGC TTGTCAATGGGCGCATT
CLPP_HAEIN_12-193_NC_007146         ---AAGCCAGATA TTCGCACTCTTTGTATT ---GGTCAAGC TTGTCAATGGGCGCATT
CLPP_RALSO_34-215_NC_003295         ---AAGCCCAGCG TGTCCACGTTGTGCATG ---GGCATGGC CGCCAGCATGGGCGCGTT
CLPP_RALSO_34-215_NZ_AAKL01000002 ---AAGCCCAGCG TGTCCACGCTGTGCATG ---GGCATGGC CGCTAGCATGGGCGCGTT
CLPP_NEIMA_16-199_NC_002946         ---AAGCCCCGAT GATCGACTTTGTGCTTG ---GGGCAGGC GGCAAGTATGGGCGGTT
CLPP_NEIMA_16-199_NC_003116        ---AAGCCCCGAT GATCGACTTTGTGCTTG ---GGGCAGGC GGCAAGTATGGGCGGTT
CLPP1_SYNP7_10-191_NC_006576        ---CGACCCGATG TCTCGACCGTTTGTGTC ---GGGCTGGC TGCCAGCATGGGCGCCTT
CLPP1_SYNP7_10-191_NC_007604        ---CGACCCGATG TCTCGACCGTTTGTGTC ---GGGCTGGC TGCCAGCATGGGCGCCTT
CLPP1_ANASP_11-192_NC_003272       ---CGCCCTGATG TTTGTACCATCTGTACA ---GGATTGGC GGCAAGTATGGGTGCTTT
CLPP2_SYNY3_32-213_NC_000911       ---CGTCCCAGAT GGGTCACCATCTGTTT ---GGTCTGGC TGCCAGCATGGGGCTTT

```

Figure 2. Protogene output on the CLP Serine Protease family. The Seed MSA of the PFAM profile entry (PFAM PF00574) was processed by Protogene. The portion of the alignment containing the Serine active site classes are indicated in yellow (UCN) and green (AGY).

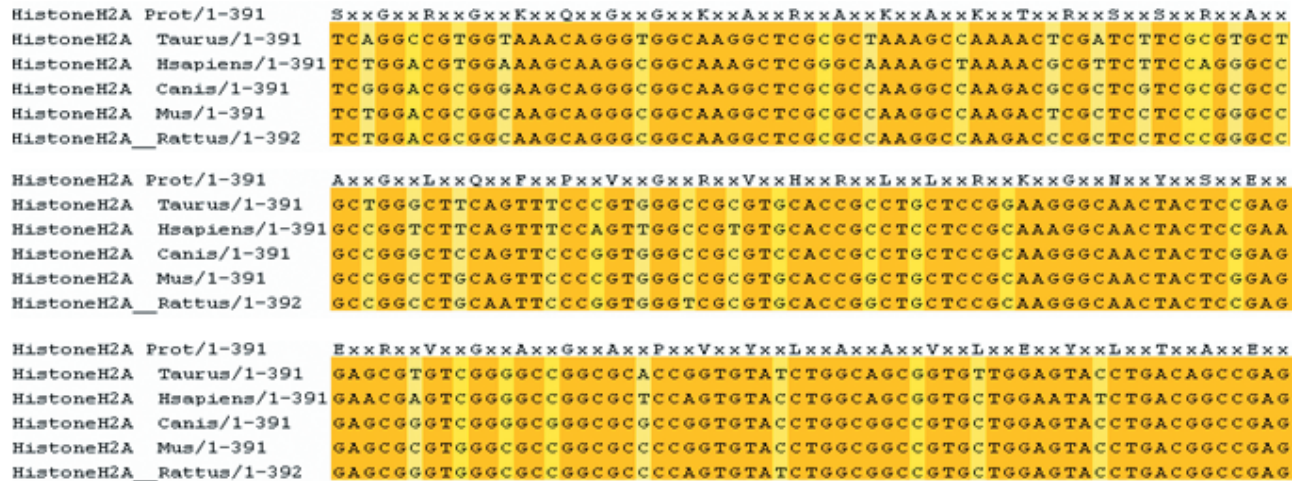


Figure 3. Protogene output on the Human H2A Histone protein. The original protein sequence is indicated on the top. Light coloured columns are those not entirely conserved.

ability to process a complete eukaryotic domain sequence dataset. For that purpose we selected the 209 human trypsin like serine proteases listed in the SMART database (9) (SM00020). These are protein domains processed using a SMART Hidden Markov Model. Protogene returned the CDSs associated with 253 protein sequences found in RefSeq and NR (157 in RefSeq and 96 in NR). Out of the 209 original human sequences, 4 matched equally well a chimpanzee protein (thus prompting the return of the associated Chimpanzee CDSs) and 66 matched two distinct human entries (thus prompting the return of 66 extra Human CDSs). Of the original sequences 26 could not be associated with an acceptable CDS: 20 did not pass the BLAST step (i.e. no suitable match was found in RefSeq or NR) and the 6 remaining could not be properly processed by exonerate against the nucleotide sequences indicated by the database annotation.

Our third example (Figure 3) addresses the question of nucleotide sequence conservation in the Histone H2A family. Histones are notoriously conserved proteins and in the present case, launching a Protogene analysis on the Human H2A sequence (SwissProt P28001) returned 5 perfect matches in RefSeq, resulting in 5 CDSs being reported: Human, Cow, Rat, Mouse and Dog. The alignment is shown of Figure 3. Such a nucleotide alignment of perfectly conserved protein sequences is ideal for phylogenetic studies or motif discoveries. It is worth pointing out that although it is identical, the Chimpanzee Histone was not reported by Protogene because it is not included in RefSeq. This finding reveals the heavy bias of Protogen toward model systems included in RefSeq. The systematic use of NR rather than RefSeq could help solve this problem, but this would come at the cost of a more complex output.

CONCLUSION

In this paper we describe Protogene, a web server that makes it possible to turn a protein MSA into the corresponding CDS MSA, using bona fide genomic or transcriptome data. Protogene is meant to be a simple yet powerful data exploration

tool. Its purpose is to rapidly ask simple questions, with an emphasis on accuracy and robustness rather than sensitivity.

ACKNOWLEDGEMENTS

We thank Guy Slater for his advice on the use of Exonerate, NCBI team for their help with EFetch and EUtils tools. We thank Prof. Jean-Michel Claverie (head of IGS) for stimulating discussions and support. The development was supported by CNRS (Centre National de la Recherche Scientifique), Sanofi-Aventis Pharma SA., Marseille-Nice Génopole and the French National Genomic Network (RNG). Funding to pay the Open Access publication charges for this article was provided by CNRS.

Conflict of interest statement. None declared.

REFERENCES

- Bininda-Emonds, O.R. (2005) transAlign: using amino acids to facilitate the multiple alignment of protein-coding DNA sequences. *BMC Bioinformatics*, **6**, 156.
- Stocsits, R.R., Hofacker, I.L., Fried, C. and Stadler, P.F. (2005) Multiple sequence alignments of partially coding nucleic acid sequences. *BMC Bioinformatics*, **6**, 160.
- Wernersson, R. and Pedersen, A.G. (2003) RevTrans: multiple alignment of coding DNA from aligned amino acid sequences. *Nucleic Acids Res.*, **31**, 3537–3539.
- Wheeler, D.L., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Chetverin, V., Church, D.M., DiCuccio, M., Edgar, R., Federhen, S. *et al.* (2006) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.*, **34**, D173–D180.
- 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.
- Slater, G.S. and Birney, E. (2005) Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics*, **6**, 31.
- Averof, M., Rokas, A., Wolfe, K.H. and Sharp, P.M. (2000) Evidence for a high frequency of simultaneous double-nucleotide substitutions. *Science*, **287**, 1283–1286.
- Finn, R.D., Mistry, J., Schuster-Bockler, B., Griffiths-Jones, S., Hollich, V., Lassmann, T., Moxon, S., Marshall, M., Khanna, A., Durbin, R. *et al.* (2006) Pfam: clans, web tools and services. *Nucleic Acids Res.*, **34**, D247–D251.
- Letunic, I., Copley, R.R., Pils, B., Pinkert, S., Schultz, J. and Bork, P. (2006) SMART 5: domains in the context of genomes and networks. *Nucleic Acids Res.*, **34**, D257–D260.