

well as their overall quantity as a subfraction of DOC, is only summarily known^{2,3,8,28}. Likewise, we lack precise knowledge about the real straining (and aggregation or agglutination) capacity of oikopleurid food-concentrating filters¹². But even a very conservative estimate of 10% of DOC as grazable by oikopleurids, means that this source of food is as important for the oikopleurids as is total POC. This agrees well with previous findings that POC accounts for a maximum of 30% of the energy needs of *O. dioica*^{19,20}.

Oikopleurid tunicates often occur in high densities in discrete strata at various depths^{21,22}, and may under such conditions clear 30–60% of the water mass in 24 h^{1,4}. Obviously, they may remove and repack colloidal DOC (>0.2 µm particle size) rapidly under such conditions. On the basis of filter parameters, this ability to graze on colloidal DOC is probably shared by caddisfly larvae²³, pedal worms²⁴, ascidians²⁵, salps²⁶ and amphioxus⁹. □

Received 10 September; accepted 8 November 1991.

1. Alldredge, A. L. *Limnol. Oceanogr.* **26**, 247–257 (1981).
2. Sugimura, Y. & Suzuki, Y. *Mar. Chem.* **24**, 105–131 (1988).
3. Koike, I., Hara, S., Terauchi, K. & Kogure, K. *Nature* **345**, 242–244 (1990).

4. Knoechel, R. & Steel-Flynn, D. *Mar. Ecol. Prog. Ser.* **53**, 257–266 (1989).
5. Shelbourne, J.E. *J. mar. biol. Ass. U.K.* **42**, 243–252 (1962).
6. Gadowski, D. M. & Boelert, G.W. *Mar. Ecol. Prog. Ser.* **20**, 1–12 (1984).
7. Keats, D. W., Steele, D. H. & South, G. R. *Canad. J. Zool.* **65**, 49–53 (1987).
8. Flood, P. R. *Experientia* **34**, 173–175 (1978).
9. Flood, P. R. *Biomed. Res. Suppl.* **2**, 49–53 (1981).
10. Deibel, D., Dickson, M.-L. & Powell, C. V. L. *Mar. Ecol. Prog. Ser.* **27**, 79–86 (1987).
11. Deibel, D. & Powell, C. V. L. *Mar. Ecol. Prog. Ser.* **39**, 81–85 (1987).
12. Flood, P. R. *Mar. Biol.* **111**, 95–111 (1991).
13. Deibel, D. & Powell, C. V. L. *Mar. Ecol. Prog. Ser.* **35**, 243–250 (1987).
14. Flood, P. R., Deibel, D. & Morris, C. C. *Biol. Bull. Mar. Biol. Lab., Woods Hole* **178**, 118–125 (1990).
15. Johnson, B. D. & Wangersky, P. J. *Limnol. Oceanogr.* **30**, 966–971 (1985).
16. Deibel, D. *Mar. Biol.* **99**, 177–186 (1988).
17. Bagnara, J. T. & Hadley, M. E. *Chromatophores and Color Change, the Comparative Physiology of Animal Pigmentation* 46–50 (Prentice-Hall, Englewood Cliffs, NJ, 1973).
18. Cauwet, G. *Oceanologica Acta* **1**, 99–105 (1978).
19. Gorsky, G. thesis, Univ. de P. et. M. Curie, Paris VI (1980).
20. King, K. R. thesis, Univ. Washington (1981).
21. Youngbluth, M. J., Bailey, T. G. & Jacoby, C. A. in *Man in the Sea* (eds Lin, Y. C. & Shida, K. K.) Vol. 2 191–208. (Best, San Pedro, California, 1990).
22. Magnesen, T., Aksnes, D. L. & Skjoldal, H. R. *Sarsia* **74**, 115–126 (1989).
23. Wallace, J. B. & Malas, D. *Arch. Hydrobiol.* **77**, 205–212 (1976).
24. Flood, P. R. & Fiala-Medioni, A. *Mar. Biol.* **72**, 27–33 (1982).
25. Flood, P. R. & Fiala-Medioni, A. *Acta Zool. Stockh.* **62**, 53–65 (1981).
26. Bone, Q., Braconnot, J.-C. & Ryan, K. P. *Acta Zool. Stockh.* **72**, 55–60 (1991).
27. Harbison, G. R. & McAlister, V. L. *Limnol. Oceanogr.* **24**, 875–892 (1979).
28. Williams, P. M. & Druffel, E. R. M. *Oceanography* **1**, 14–17 (1988).

ACKNOWLEDGEMENTS. We thank E. Bru, B. Hansen, M. Riehl and the scuba divers of the Ocean Science Centre for technical assistance. This work was supported by the Norwegian Fisheries Research Council, the Norwegian Research Council for Science and the Humanities (to P.R.F.) and by the Natural Sciences and Engineering Research Council of Canada (to D.D.).

Sequence identification of 2,375 human brain genes

Mark D. Adams, Mark Dubnick, Anthony R. Kerlavage, Ruben Moreno, Jenny M. Kelley, Teresa R. Utterback, James W. Nagle, Chris Fields & J. Craig Venter*

Receptor Biochemistry and Molecular Biology Section, National Institute of Neurological Disorders and Stroke, National Institutes of Health, Bethesda, Maryland 20892, USA

WE recently described a new approach for the rapid characterization of expressed genes by partial DNA sequencing to generate 'expressed sequence tags'¹. From a set of 600 human brain complementary DNA clones, 348 were informative nuclear-encoded messenger RNAs. We have now partially sequenced 2,672 new, independent cDNA clones isolated from four human brain cDNA libraries to generate 2,375 expressed sequence tags to nuclear-encoded genes. These sequences, together with 348 brain expressed sequence tags from our previous study, comprise more than 2,500 new human genes and 870,769 base pairs of DNA sequence. These data represent an approximate doubling of the number of human genes identified by DNA sequencing and may represent as many as 5% of the genes in the human genome.

Most (83%) of the 2,375 partial cDNA sequences reported here (Table 2) are not related to any previously described sequences. Based on database matches to known genes from humans and from such evolutionarily distant organisms as *Escherichia coli*, yeast, *Caenorhabditis elegans*, *Drosophila*, barley, *Arabidopsis*, rice and green algae, we have putatively identified 217 of the expressed sequence tags (ESTs; Table 1). These include a novel gene similar to *Notch/TAN-1* (refs 1, 2), a new neurotransmitter transporter gene, and a new member of the multidrug resistance gene family. Several genes involved in development or cell differentiation in *Drosophila* are represented by similar human ESTs, including *seven in absentia*³, *big-brain*⁴, the *discs-large* tumour suppressor⁵ and the homeotic gene *orthodenticle*⁶. New members of previously known gene families in humans include a Ca²⁺-transporting ATPase, an ADP ribosylation factor and a new neural-cell adhesion molecule gene.

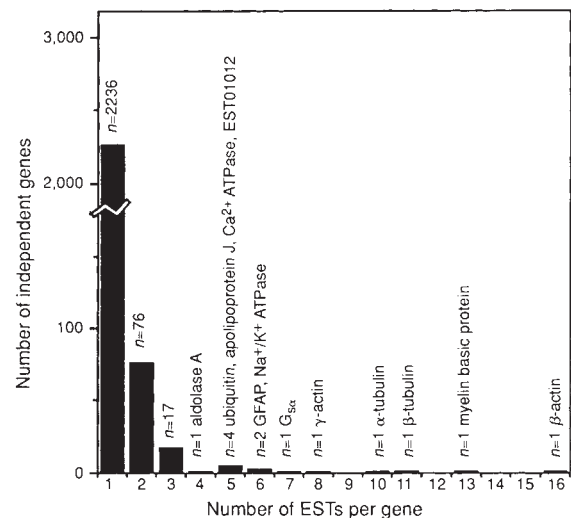


FIG. 1. Redundancy of sequencing of ESTs. The number of putatively identified EST clones plus the groups of ESTs that form contigs are plotted against the number of independent genes represented. The number of genes is given above each bar, with the names of the genes for redundancies of 4–16. We define redundancy as the number of times each gene is represented by an EST; for example β-actin has an EST redundancy of 16. One five-member EST contig was constructed that did not match any known sequences, indicating that at least one common transcript in brain, EST 01012, has not been reported previously. GFAP, glial fibrillary acidic protein.

The 1,971 ESTs without a putative identification were analysed using the coding-region prediction program CRM with the GRAIL server⁷. Some of the unknown ESTs (15%) scored a likely probability of containing protein-coding sequence. Half of the ESTs to known human genes contain protein-coding sequences, so at most half of the unknown ESTs are likely to contain coding sequences. We have found no evidence that genomic DNA or cDNA to unspliced precursor RNA is a major contaminant of either the hippocampus or fetal brain library.

The limited extent of redundancy of EST sequencing is shown in Fig. 1. Of the nuclear-encoded messenger RNAs, the most common ESTs were to the β-actin (0.6% of the EST clones)

* To whom correspondence should be addressed.

TABLE 2 Gene composition of human brain cDNA libraries

	Total	Hippocampus unscreened	Hippocampus prescreened	Fetal brain*	Fetal brain†	Whole brain‡
No database match	1,942	474	394	1,025	32	17
Exact human match	255	76	88	87	0	4
Non-exact human match	51	8	10	33	0	0
Non-human match	99	34	26	35	2	2
Alu repeat	313	70	70	172	1	0
L1 repeat	58	13	5	39	0	1
THE-Itr sequence	17	1	1	14	1	0
Other repeat	9	3	4	2	0	0
Mitochondrial	181	115	33	27	0	6
rRNA	57	29	7	16	5	0
poly(A) insert only	339	171	161	5	0	2
Total	3,321	994	799	1,455	41	32

Four cDNA libraries were used as sources of clones for sequencing. Human hippocampus and fetal brain (*) libraries, plasmid template preparation, sequencing reactions, and automated sequencing were performed as described¹. A pooled probe consisting of inserts from 10 different EST clones with sequences that matched either mitochondrial genes or the 18S or 28S rRNAs was used to prescreen a gridded filter array of the hippocampus library; nonhybridizing clones are referred to as the 'prescreened library'. Another fetal brain library (†) was constructed by and was a gift from B. Soares (Columbia University). A directionally-cloned library (‡) was prepared¹⁴ using human adult brain mRNA from Clontech (Palo Alto). Of 482 clones analysed by restriction-enzyme digestion, 33% contained inserts at least 1,500 base pairs long. Stratagene hippocampus and fetal brain library totals include data from ref. 1. Sequences of nuclear-encoded cDNAs that did not include the interspersed repeats Alu, L1 or THE-Itr¹⁵⁻¹⁷ were searched against GenBank and, in 6-frame translation, against a comprehensive, non-redundant peptide database using the network BLAST¹⁸ server at the National Center for Biotechnology Information. For significant similarities, a putative gene name and protein identification resource (PIR) gene family identification¹⁹ for the EST were assigned. ESTs without significant matches using BLAST were searched in translation against PIR using the program FASTA. Ten additional marginal matches were found. A total of 2,300 new EST sequences comprising 765,505 nucleotides from the current data set have been submitted to GenBank and assigned accession numbers M77851-M79278 and M85308-M86179. All ESTs except those multiply representing actin, tubulin, and myelin basic protein clones were submitted. cDNA clones from which ESTs were derived are available from the American Type Culture Collection (Rockville, Maryland) with accession numbers 77501-78999 and 81000-81756. The Genome Data Base²⁰ expressed D-segment numbers for these clones are DOS1E-DOS2300E. We have developed a database which includes the clone and sequence data, sequence analysis results, physical mapping data, tissue localization and cross-references to the public databases and distribution of the clones, mapping and sequence data using the Sybase relational database management system (Sybase Inc., Emeryville, California). Comprehensive reports on the sequences described in this paper are available in electronic form. A README file describing how to access the EST database reports is available via anonymous file transfer protocol (FTP) to briggs.ninds.nih.gov. Questions on data access and database structure can be addressed via electronic mail to arkerlav@briggs.ninds.nih.gov.

and myelin basic protein genes (0.5% of the clones). Myelin basic protein, a highly expressed structural component of nerve tissue⁸, displays four alternate splicing forms, of which at least two are present among the ESTs reported here. Other common ESTs were G-protein subunit G_{sα}, γ-actin and both α- and β-tubulin.

All of the genes for which four or more ESTs were found have been sequenced in humans, except for one which was matched by five unknown ESTs. Assuming that most brain mRNAs are rare transcripts⁹, the chance of finding a new gene by EST sequencing is fairly high when ribosomal and mitochondrial transcripts are eliminated. Therefore, although normalization may be important as we near closure in sequencing every human gene, it is not necessary at this stage to reduce sequencing redundancy or to increase gene diversity. Furthermore, a certain amount of redundancy is desirable to the extent that it promotes assembly of EST contigs into full-length cDNA sequences.

By matching ESTs to known database sequences, a phenotypic characterization of the tissue begins to emerge. Protein super-families matched by ESTs were grouped into three broad functional categories to assess the biological spectrum represented

by these randomly selected cDNA clones. Structural and metabolic classes comprised about 30% of the ESTs each, 25% were involved in regulatory pathways and the remainder were not classifiable. Eleven of the eighteen enzymes of glycolysis and the citric acid cycle are represented by at least one subunit or isozyme. In addition, several genes not previously known to be expressed in the brain were matched, including spermine/spermidine acetyltransferase¹⁰ and osteopontin¹¹. Isolation of 171 ESTs from mouse testes was recently reported¹², including four with database matches in common with human ESTs.

The genomic mapping of these new human expressed genes is among our highest priorities. Physical mapping of the 2,375 EST clones reported here would provide human chromosome markers spaced an average of 1.2 megabases apart and would roughly double the number of expressed sequences that have been localized to chromosomes¹³. Mapped ESTs are a new resource for identifying candidate genes for the estimated 5,000 single-locus diseases¹³. All the sequences and clones described here are publicly available (Table 2). We shall update EST clone identification and map information through the NIH cDNA database. □

Received 27 November 1991; accepted 3 January 1992.

- Adams, M. D. *et al. Science* **252**, 1651-1656 (1991).
- Ellisen, L. *et al. Cell* **66**, 649-661 (1991).
- Carthew, R. & Rubin, G. *Cell* **63**, 561-577 (1990).
- Rao, Y., Jan, L. & Jan, Y. *Nature* **345**, 163-167 (1990).
- Woods, D. & Bryant, P. *Cell* **66**, 451-464 (1991).
- Finkelstein, R., Smouse, D., Capaci, T., Spradling, A. & Perrimon, N. *Genes Dev.* **4**, 1516-1527 (1990).
- Überbacher, E. & Mural, R. *Proc. natn. Acad. Sci. U.S.A.* **88**, 11261-11265 (1991).
- Kamholz, J., de Ferra, F., Puckett, C. & Lazzarini, R. *Proc. natn. Acad. Sci. U.S.A.* **83**, 4962-4966 (1986).
- Galau, G., Klein, W., Britten, R. & Davidson, E. *Archs Biochem. Biophys.* **197**, 584-599 (1977).
- Casero, R. *et al. J. Biol. Chem.* **266**, 810-814 (1991).
- Young, M. *et al. Genomics* **7**, 491-502 (1990).

- Hoög, C. *Nucleic Acids Res.* **19**, 6123-6127 (1991).
- McKusick, V. *FASEB J.* **5**, 12-20 (1991).
- Rubenstein, J. *et al. Nucleic Acids Res.* **18**, 4833-4842.
- Schmid, C. W. & Jalinek, W. R. *Science* **216**, 1065-1070 (1982).
- Paulson, K. E. *et al. Nature* **316**, 359-361 (1985).
- Fanning, T. G. & Singer, M. F. *Biochim. biophys. Acta* **910**, 203-212 (1987).
- Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. *J. molec. Biol.* **215**, 403-410 (1990).
- Barker, W., George, D., Hunt, L. & Garavelli, J. *Nucleic Acids Res.* **19** (Suppl), 2231-2236 (1991).
- Pearson, P. *Nucleic Acids Res.* **19** (Suppl), 2237-2239 (1991).

ACKNOWLEDGEMENTS. We thank J. Powell and J. Kelley of the Division of Computer Research and Technology at NIH for computer systems support and D. Lipman of the National Center for Biotechnology Information at NIH for access to the network BLAST server and nonredundant peptide database.