

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.

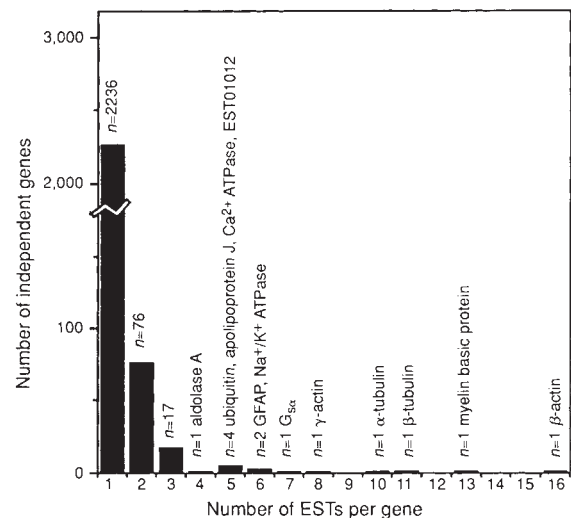


FIG. 1. Redundancy of sequencing of ESTs. The number of putatively identified EST clones plus the groups of ESTs that form contigs are plotted against the number of independent genes represented. The number of genes is given above each bar, with the names of the genes for redundancies of 4–16. We define 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 1 Gene composition of human brain cDNA libraries

EST	Putative Identification	Species	DB	Accession	Length	%ID	EST	Putative Identification	Species	DB	Accession	Length	%ID
EST02245	14.3.3 protein (PKI regulator)	H	(G)	H1M433	272	98.9	EST01651	Lactate dehydrogenase A	H	(G)	H1MDHA	378	99.5
EST01757	2',3'-cyclic nucleotide phosphodiesterase H	(G)	H1MCPDEA	272	98.5	EST01670	Lactate dehydrogenase B	H	(G)	H1MDHBR	401	99.5	
EST01784	60K filarial antigen	A	(P)	A28209	88	50.6	EST01764	Lamin B receptor	C	(P)	A36247	76	71.4
EST01455	80K-H protein	H	(G)	H1MG19P1A	263	97.3	EST01724	Ion protease	E	(P)	J00901	103	41.3
EST01982	ADP-ribosylation factor 1	H	(P)	B33283	84	41.2	EST02413	Long-chain-fatty-acid-CoA ligase	R	(P)	A36275	36	62.2
EST02077	ADP/ATP translocase	H	(G)	H1M1LCA	200	96.5	EST02418	MARCKS homology	M	(E)	MFE52	237	92.2
EST01620	AMP deaminase, brain	R	(P)	A37056	57	100.0	EST01952	MHC class I HLA-Cw2 heavy chain	H	(G)	H1M6KWB	231	97.8
EST01504	Actin, beta, cytoskeletal	H	(G)	H1MCCYBA	247	98.8	EST01640	MHC class III HSP70-1 gene (HLA)	H	(G)	H1M6HSP	181	100.0
EST01543	Actin, gamma, cytoskeletal	H	(G)	H1M6CTGM	359	98.9	EST02505	Matrin 3	R	(G)	R1MTRIN3	137	93.5
EST01891	Actinin, alpha	H	(G)	H1M6ACTAR	315	81.6	EST01918	Metalloproteinase inhibitor	H	(G)	H1M6MET	236	99.2
EST01801	Adaptin, beta	H	(G)	H1M6ADPTA	380	99.5	EST01865	Microtubule-associated protein 1B	R	(G)	R1M6ADU	299	86.4
EST01710	Adenosine deaminase	H	(G)	H1M6ADG	316	100.0	EST01473	Microtubule-associated protein 4	H	(G)	H1M6M4	249	97.6
EST01625	Ahrin	(G)	H1M6ADR	103	84.6	EST01678	Milk fat globule membrane protein	M	(P)	A36479	48	61.2	
EST02113	Ala	I	(H)	R1C0023G_1	38	59.0	EST01704	Monamine oxidase A (MAOA)	H	(G)	H1M6MOA	255	98.8
EST00675	Alcohol dehydrogenase	H	(G)	H1M6ALD	346	98.0	EST01769	Myelin basic protein	H	(G)	H1M6MBP	314	99.4
EST01660	Aldolase A	H	(G)	H1M6ALD	317	99.0	EST01580	Myeloid differentiation primary response	M	(H)	M1SKYD118	176	88.3
EST01761	Aldolase C	H	(G)	H1M6ALD	317	99.0	EST02585	Myoblast cell surface antigen 24.1D5	H	(G)	H1M641D5	177	98.3
EST02688	Alpha-Enolase	H	(G)	H1M6ENO	345	98.8	EST01614	Mycosin heavy chain, non-muscle	H	(G)	H1M6MYM	291	99.3
EST01635	Alpha-2-Macroglobulin	H	(G)	H1M6AZM	319	99.4	EST01744	N-formylpeptide receptor	H	(G)	H1M6NFR	237	99.2
EST01664	Amyloid A	H	(P)	A29030	52	54.7	EST01744	NAD(P)+ transhydrogenase (H-specific)	B	(P)	DEBOM	86	93.1
EST01700	Anion exchanger homolog AE3	M	(P)	A33638	95	97.9	EST01805	NF-kappa-B transcription factor	H	(G)	H1M6NFKB	98	98.0
EST01585	Apolipoprotein J	H	(G)	H1M6APOJ	317	99.0	EST01805	Na+/K+ ATPase, alpha subunit	H	(G)	H1M6ATP2	277	99.9
EST01825	Aspartate aminotransferase, cytosolic	H	(G)	H1M6ASAM	309	98.7	EST02610	Neural cell adhesion molecule L1	M	(P)	S05479	82	43.4
EST02671	Aspartate aminotransferase, mitochondria	H	(G)	H1M6ASAM	231	99.1	EST01471	Neuraxin	R	(P)	S06017	120	84.3
EST01634	Axonal glycoprotein TAG-1	R	(P)	A34695	69	87.1	EST01519	Neurofibromatosis type 1	H	(G)	H1M6NF1	369	98.6
EST02530	B cell-specific Mo-MLV integration site	(G)	M1M6S11A	111	87.5	EST01749	Neurofilament heavy chain	H	(G)	H1M6NFH	356	99.7	
EST02306	Bib protein	D	(P)	S09699	57	53.4	EST02453	Neurofilament subunit M (NF-M)	H	(G)	H1M6NFM	164	99.4
EST01443	CDP-30-p1	E	(P)	JF0168	33	41.2	EST02609	Neutrophil oxidase factor	H	(P)	A34855	43	47.7
EST01800	Ca2+-transporting ATPase	H	(G)	H1M6CA	380	99.2	EST02632	Neutrophil oxidase factor	H	(P)	A24400	63	39.1
EST02146	Calbindin D28	R	(G)	R1M6CALB28	81	87.8	EST01643	Neutrophil protein 1	K	(P)	A24400	63	39.1
EST01823	Calcineurin A2	H	(G)	H1M6CNAB	287	99.7	EST01961	Neutrophil protein 1	K	(P)	A24400	63	39.1
EST02055	Calcium channel	L	(P)	S05054	33	67.6	EST02429	Notch/Notch homologue	H	(H)	H1M6NTN1	85	57.0
EST01849	Calmodulin	H	(G)	H1M6CALM	327	98.5	EST01573	Nuclear factor I-like protein (NF1)	H	(H)	H1M6NF1A	111	92.0
EST01466	Calmodulin-dependent PKI, BII	R	(P)	A26464	93	98.9	EST01657	Nuclear factor I-like protein (NF1)	H	(P)	A33386	71	52.8
EST02378	cAMP-dependent PKI inhibitor	M	(G)	M1M6PKI	234	91.5	EST01657	Osteonectin/SPARC	H	(G)	H1M6SPARC	348	99.7
EST01644	cAMP-dependent PKI regulatory subunit RIIB	H	(G)	H1M6RIIB	198	97.5	EST01822	Osteopontin	H	(G)	H1M6OP	170	96.5
EST01628	cAMP-dependent PKI regulatory subunit RIIB	H	(G)	H1M6RIIB	198	97.5	EST01828	Oxid homeotic protein	D	(P)	A35912	35	52.8
EST01041	cAMP-regulated phosphoprotein	B	(P)	B33508	21	86.4	EST01486	Pancreatic tumor-related protein	H	(G)	H1M6PANTR	422	99.3
EST02447	cAMP-specific phosphodiesterase	H	(G)	H1M6DEAA	363	69.0	EST01798	Peptidylprolyl isomerase (cyclophilin)	H	(G)	H1M6CYC	382	99.5
EST01536	Cannabinoid receptor	H	(S)	C1M6NRHMAN	97	93.9	EST01751	P-4,5-BPL1	R	(P)	A28807	40	90.2
EST01979	Carboxypeptidase E	H	(G)	X11405	330	99.4	EST01656	Phosphoglycerate kinase	H	(G)	H1M6GK1	374	98.7
EST01606	Casein kinase II alpha subunit	H	(G)	H1M6CKII	186	97.3	EST00992	Polymyxin B resistance	Y	(P)	A32714	20	76.2
EST01733	Catalase	H	(G)	H1M6CAT	270	96.3	EST01806	Prohibitin	H	(P)	R1M6PHB1	120	97.5
EST01810	Cathepsin D	H	(G)	H1M6CATD	246	98.4	EST01775	Prothymosin cleavage enzyme	M	(H)	M1M6CSLA	91	93.5
EST01799	Cell surface glycoprotein MUC18	H	(G)	H1M6MUC18	413	97.8	EST01461	Prothymosin cleavage enzyme	M	(H)	M1M6CSLA	91	93.5
EST01487	C/D PG40	H	(G)	H1M6PG40	261	100.0	EST02087	Protein kinase C, zeta	H	(G)	H1M6PKCZ	382	98.7
EST01913	Clathrin coat assembly protein homologue	Y	(H)	Y1M6YAP54	62	63.5	EST01650	Protein phosphatase 2A beta subunit	H	(G)	H1M6PP2AB	288	76.8
EST01796	Coagulation factor VII	H	(G)	H1M6CFVII	227	97.8	EST01584	PR65 (alpha)	H	(G)	H1M6PR65	242	98.8
EST01676	Cofilin	P	(G)	P1C0051L	132	89.5	EST01786	Protein-tyrosine kinase (clone JTK1)	H	(G)	C18269	38	100.0
EST01774	Cysteine proteinase inhibitor	H	(G)	H1M6CYSICR	163	97.6	EST01572	Protochlorophyllide reductase	W	(P)	S04783	34	57.1
EST01824	Cysteine-rich intestinal protein	R	(P)	GYR1	56	66.7	EST01538	Pyruvate dehydrogenase alpha subunit	H	(G)	H1M6PDHA	322	99.7
EST01502	Cytochrome P-450IIE1	H	(G)	H1M6CYPIIE	175	97.2	EST01587	Pyruvate kinase isozyme M2	H	(G)	H1M6PKM2	395	98.5
EST01951	Cytochrome c	H	(G)	H1M6CYCA	172	96.0	EST02683	Rab1	H	(H)	H1M6RASP1	78	97.7
EST01808	DNA binding protein YAVR367	H	(G)	H1M6YAVR367	429	98.8	EST01572	Rab5	H	(P)	F34323	91	82.6
EST01721	DNA repair helicase (ERCC3)	H	(G)	H1M6ERCC3A	334	98.5	EST01072	Rac protein kinase	H	(G)	H1M6RACPC	91	100.0
EST01621	DNF1552 (lung) mRNA	H	(G)	H1M6DNF1552	137	96.4	EST01389	Radial spoke protein 3	F	(P)	S05962	58	52.5
EST01257	Diacylglycerol kinase, lymphocyte	P	(P)	S09156	44	42.2	EST01787	Retinaldehyde-binding protein	H	(G)	H1M6RALBP	129	100.0
EST02477	Dianion acetyltransferase	H	(S)	A1M6ASHMAN	74	45.3	EST01579	Retrovirus-related gag polyprotein	H	(P)	H0F4E2	99	77.1
EST01508	Dihydroallopurinol acetyltransferase	H	(S)	H1M6DHAA	254	98.4	EST02550	Retrovirus-related pol polyprotein	H	(P)	GNLGL	50	54.9
EST02062	Dilute (myosin heavy chain)	M	(H)	M1M6DILUTE_1	27	100.0	EST01578	Riboflavin H	H	(G)	H1M6RIBH	289	99.2
EST01719	Disc-large tumor suppressor	D	(H)	D1M6DLG1	53	63.0	EST01928	Ribosomal phosphoprotein P0	H	(G)	H1M6RPP0	273	98.2
EST02627	Elongation factor 1 alpha	H	(G)	H1M6EF1A	361	99.2	EST01583	Ribosomal phosphoprotein P1	H	(G)	H1M6RPP1	126	91.3
EST01165	Elongation factor 1 beta	N	(P)	A24806	36	64.9	EST01583	Ribosomal protein L18a	R	(P)	R1M6R18A	155	94.9
EST02596	Enolase, gamma-2, neuron specific	H	(G)	H1M6ENOG	345	98.5	EST01627	Ribosomal protein L19a	X	(P)	A24579	75	63.1
EST01743	Epoxide hydrolase	H	(G)	H1M6EH	363	99.7	EST01667	Ribosomal protein L3	O	(P)	J00771	74	80.0
EST01946	Ezrin	H	(G)	H1M6EZRN	153	99.3	EST01826	Ribosomal protein S10	S	(P)	R3MY10	36	51.4
EST01971	Familial adenomatous polyposis coli	H	(G)	H1M6APPC	316	98.1	EST01459	Ribosomal protein Y10	Y	(P)	S11581	40	68.3
EST01325	Fatty acid synthase	R	(G)	H1M6FATP	96	79.8	EST02608	Ribosomal protein S3	H	(G)	H1M6RMS3	279	95.7
EST01476	Fibrillarin	H	(G)	H1M6FIBA	201	99.5	EST01442	Seven in absentia	D	(P)	A36195	46	80.8
EST01790	Fibroblast growth factor receptor	H	(G)	H1M6FGFR	376	99.2	EST01960	Spectrin, beta	H	(G)	H1M6SPTB	268	67.7
EST01967	Fibronectin	H	(G)	H1M6FNCT	120	98.3	EST01699	Sperm membrane protein	R	(P)	A35981	52	58.5
EST02186	Fodrin, alpha	H	(G)	H1M6FOD	282	98.2	EST01760	Spermidine/spermine N1-acetyltransferase	H	(G)	H1M6SPMAT	102	97.1
EST02428	Fumarate hydratase, mitochondrial	H	(G)	H1M6FH	96	96.9	EST01545	Sphingomyelin phosphodiesterase	H	(G)	H1M6SMAS	159	97.5
EST01665	G(1) alpha	H	(G)	H1M6G1A	310	98.1	EST01984	Stathmin (p18)	H	(G)	H1M6STH	180	100.0
EST02389	G(1) alpha	H	(G)	H1M6G1A	275	99.3	EST01987	Succinate dehydrogenase flavoprotein	B	(H)	B1M6SDHFP1	44	100.0
EST02362	Ga binding protein, beta subunit	M	(H)	M1M6GAPC	86	90.8	EST00742	Synaptobrevin (p65)	H	(S)	S1M6SVB	27	53.6
EST01745	GTP-binding protein beta chain homologue	H	(G)	H1M6GALB23	319	99.1	EST01586	T-cell surface glycoprotein E2	H	(G)	X1G996	277	99.6
EST00825	GABA-aminotransferase alpha transporter	R	(P)	A35918	26	59.3	EST01809	TAB-1, 26-kDa cell surface protein	H	(G)	H1M6TAB1	271	96.5
EST01738	Gelatin factor AB-280	H	(P)	A37098	74	80.0	EST01955	TCTE1	H	(G)	H1M6TCTE1	162	98.8
EST01776	Geisolin precursor, plasma	H	(G)	H1M6GSR	171	99.4	EST01575	TFBF transcription factor	H	(G)	H1M6TFBF	369	98.1
EST01649	Glial fibrillary acidic protein	H	(G)	H1M6GFAP	399	99.2	EST02402	Talin	M	(H)	M1M6TALINR	79	81.2
EST01965	Globin, gamma	H	(G)	H1M6GLB3	147	100.0	EST01601	Thiosulfate sulfurtransferase (rhodanase)	B	(P)	ROB	65	81.8
EST02192	Glutamate decarboxylase	H	(G)	H1M6GAD	213	100.0	EST01435	Threonyl-tRNA synthetase	H	(G)	H1M6TRSYNT	228	99.6
EST01702	Glutamate dehydrogenase	H	(G)	H1M6GDH	223	96.4	EST02420	Thrombospondin precursor	H	(G)	H1M6THSP	98	98.3
EST02446	Glutamate-aspartate carrier protein	E	(P)	J0092	57	37.9	EST01783	Thy-1 glycoprotein	H	(G)	H1M6THY1A	353	98.6
EST02034	Glutaminase	T	(S)	G1M6SRAT	34	74.3	EST02455						

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.