

## Nematode Gene Sequences, December 2000 Update

JAMES P. McCARTER,<sup>1</sup> DAVID MCK. BIRD,<sup>2</sup> SANDRA W. CLIFTON,<sup>1</sup> AND  
ROBERT H. WATERSTON<sup>1</sup>

High-throughput sequencing is revolutionizing molecular nematology by providing the sequences of thousands of genes never before characterized. The most rapid and cost-effective route to gene discovery for nematode genomes is the generation of expressed sequence tags (ESTs), single-pass reads from random cDNA library clones that provide 300–600 nucleotides of sequence from a gene. Projects are currently under way at Washington University's Genome Sequencing Center that will generate 225,000 ESTs from 14 nematode species by 2003. Additionally, the Sanger Centre and Edinburgh University will produce 80,000 ESTs from seven species. New sequences are immediately submitted to the dbEST (database of expressed sequence tags) division of GenBank. By the completion of these efforts, we anticipate the identification of more than 80,000 new nematode genes.

Here we present a brief progress report on publicly available ESTs from nematodes. *Caenorhabditis elegans* has long been a focus of sequencing efforts, and 111,639 ESTs are available from *Caenorhabditis* species. Cur-

rently, 72,764 ESTs are available from 19 nematode species beyond *Caenorhabditis*, including seven human parasites, five animal parasites, five plant parasites, and two free-living bacteriovores (Table 1). The majority of these ESTs were generated in 1999–2000. EST sequences are publicly available from dbEST at GenBank and a number of parasite-specialized Web sites (Table 2).

A major factor contributing to the popularity of ESTs is the versatility of the data. ESTs can be used for new gene discovery, gene discovery by similarity searching, genomic physical mapping, determination of gene coding regions and splice isoforms, determination of a species' codon bias, design of gene expression biochips and microarrays, determination of evolutionary relationships between species, identification of protein and domain families, and identification of single nucleotide polymorphisms (SNPs). Strategies for using ESTs, as well as discussions of the strengths and weaknesses of EST data, are available from reviews (Marra et al., 1998; Blaxter et al., 1999; McCarter et al., 2000a).

Nematologists also benefit greatly from the availability of the complete genome sequence of *C. elegans* and the annotation of its genes (The *C. elegans* Sequencing Consortium, 1998). The essentially complete sequence of *C. elegans* published in 1998 was composed of 97 megabases with 19,099 predicted protein encoding genes. Gap filling to date has brought the total genome to just over 100 finished megabases (Genome Sequencing Center, unpubl. data) with 19,705 predicted proteins including 405 splice variants (Sanger Centre Wormpep Release 33, 2000). A number of small gaps remain. Ad-

Received for publication .

Parasitic nematode EST sequencing at Washington University is supported by NIH-NIAID research grant AI 46593 to Robert Waterston and NSF Plant Genome award 0077503 to David Bird (PI) and co-PIs Sandra Clifton, Joseph Kieber, Charles Opperman, and Jeffrey Thorne. James McCarter is a Merck Postdoctoral Fellow of the Helen Hay Whitney Foundation.

<sup>1</sup> Genome Sequencing Center, Department of Genetics, Box 8501, Washington University School of Medicine, St. Louis, MO 63108, USA.

<sup>2</sup> Plant Nematode Genetics Group, Department of Plant Pathology, North Carolina State University, Raleigh, NC 27695, USA.

The authors thank members of the Genome Sequencing Center EST lab, especially Deana Pape, John Martin, Todd Wylie, and Brandi Chiapelli.

E-mail: mccarter@genetics.wustl.edu

This paper was edited by B. C. Hyman.

TABLE 1. 21 Nematode species have more than 50 ESTs registered in the GenBank dbEST database, December 2000.

Nematode species	Number of ESTs	Major source	Publications (ESTs analyzed)
<i>Caenorhabditis elegans</i>	109,215	NIG, Japan <sup>1</sup>	Waterston et al., 1992 (1,517) McCombie et al., 1992 (720)
<i>Brugia malayi</i>	22,392	WHO FilGenNet <sup>2,3,6</sup>	Blaxter et al., 1996 (596)
<i>Onchocerca volvulus</i>	13,802	Smith College, MA <sup>3</sup>	Unnasch and Williams, 2000
<i>Strongyloides stercoralis</i>	10,922	GSC, St. Louis, MO <sup>4</sup>	Moore et al., 1996 (55)
<i>Meloidogyne incognita</i>	6,626	GSC, St. Louis, MO <sup>4</sup>	McCarter et al., 2000b (5,713)
<i>Ancylostoma caninum</i>	5,546	GSC, St. Louis, MO <sup>4</sup>	
<i>Pristionchus pacificus</i>	4,989	GSC, St. Louis, MO <sup>4</sup>	
<i>Caenorhabditis briggsae</i>	2,424	GSC, St. Louis, MO <sup>5</sup>	
<i>Haemonchus contortus</i>	2,399	Univ. Edinburgh <sup>6</sup> & Sanger Centre, UK, <sup>7</sup> Wash. State, WA <sup>8</sup>	
<i>Heterodera glycines</i>	1,506	GSC, St. Louis, MO <sup>4</sup>	
<i>Meloidogyne javanica</i>	1,208	GSC, St. Louis, MO <sup>4</sup>	
<i>Globodera rostochiensis</i>	894	Scottish Crop & Wageningen <sup>9</sup>	Popeijus et al., 2000 (894)
<i>Ascaris suum</i>	588	Univ. Edinburgh, UK <sup>6</sup>	
<i>Toxocara canis</i>	519	Univ. Edinburgh, UK <sup>10</sup>	Tetteh et al., 1999 (519)
<i>Zeldia punctata</i>	378	GSC, St. Louis, MO <sup>4</sup>	
<i>Trichuris muris</i>	301	Univ. Edinburgh, UK <sup>6</sup>	
<i>Necator americanus</i>	211	Univ. Edinburgh, UK <sup>6</sup>	Daub et al., 2000 (211)
<i>Litomosoides sigmodontis</i>	198	Univ. Edinburgh, UK <sup>6</sup>	
<i>Wuchereria bancrofti</i>	131	Smith College, MA <sup>3</sup>	
<i>Globodera pallida</i>	94	Scottish Crop & Wageningen <sup>9</sup>	Popeijus et al., 2000 (94)
<i>Onchocerca ochengi</i>	60	Smith College, MA <sup>3</sup>	
Total ESTs	184,403		
Total Non- <i>Caenorhabditis</i>	72,764		

<sup>1</sup> Korhara et al., National Institute of Genetics, Mishima, Japan.

<sup>2</sup> World Health Organization Filarial Genome Network.

<sup>3</sup> Williams et al., Dept. of Biology, Smith College, Northampton, MA USA.

<sup>4</sup> McCarter et al., Genome Sequencing Center, Washington Univ. School of Medicine, St. Louis, MO USA.

<sup>5</sup> Marra et al., Genome Sequencing Center, Washington Univ. School of Medicine, St. Louis, MO USA.

<sup>6</sup> Blaxter et al., Institute of Cell, Animal, and Population Biology, Univ. of Edinburgh, Edinburgh, UK.

<sup>7</sup> Sanger Centre, Hinxton, UK.

<sup>8</sup> Jasmer et al., Dept. of Veterinary Microbiology and Pathology, Washington State Univ., Pullman, WA.

<sup>9</sup> Jones et al., Nematology Dept., Scottish Crop Research Inst., Dundee, UK, and Laboratory of Nematology, Wageningen Univ., Wageningen, The Netherlands.

<sup>10</sup> Maizels et al., Institute of Cell, Animal, and Population Biology, Univ. of Edinburgh, Edinburgh, UK.

TABLE 2. Selected Web resources for nematode EST access.

EST resources	URLs
GenBank dbEST	<a href="http://www.ncbi.nlm.nih.gov/dbEST">www.ncbi.nlm.nih.gov/dbEST</a>
Blaxter Lab ESTs & Clusters	<a href="http://www.nematodes.org">www.nematodes.org</a>
Genome Sequencing Center ESTs & NemaGene	<a href="http://www.nematode.net">www.nematode.net</a>
EMBL Parasite Genome Server	<a href="http://www.ebi.ac.uk/parasites/parasite-genome.html">www.ebi.ac.uk/parasites/parasite-genome.html</a>
The Filarial Genome Network	<a href="http://nema.cap.ed.ac.uk/fgn/ests.html">nema.cap.ed.ac.uk/fgn/ests.html</a> or <a href="http://circuit.neb.com/fgn/ests.html">circuit.neb.com/fgn/ests.html</a>
More Extensive Links	<a href="http://elegans.swmed.edu/Nematodes/">elegans.swmed.edu/Nematodes/</a>

The Blaxter Lab EST page provides access to EST data sets from *Haemonchus contortus*, *Ascaris suum*, *Necator americanus*, *Litomosoides sigmodontis*, and *Trichuris muris* as well as data from *Trichinella spiralis* and *Loa loa* not yet available in dbEST. A BLAST server with data from 18 species is also available. Nematode.net, a collaboration between the Washington University Genome Sequencing Center and North Carolina State University is a project to provide comprehensive access to nematode ESTs, including search capacity, ftp downloads, and sequence trace files. Nematode.net will also provide access to the NemaGene Index of nematode genes built by clustering overlapping ESTs. Currently available information includes ESTs and NemaGene version 1.0 clusters from *Meloidogyne incognita*. The European Molecular Biology Lab (EMBL) Parasite Genome Server provides BLAST search access to the *Brugia malayi* ESTs as well as several protozoan parasites. *Brugia malayi* data, as well as information about the *B. malayi* EST project, is also provided by the Filarial Genome Network page.

ditionally, 10.1 megabases of the *C. briggsae* genome is in finished form (Genome Sequencing Center, unpubl. data). The majority of genes identified to date in parasitic

TABLE 3. Selected Web resources for *C. elegans* genome access.

EST resources	URLs
Wormbase	<a href="http://www.wormbase.org">www.wormbase.org</a>
Wormpep, Sanger Centre	<a href="http://www.sanger.ac.uk/Projects/C_elegans/wormpep">www.sanger.ac.uk/Projects/C_elegans/wormpep</a>
WormPD	<a href="http://www.proteome.com/databases/index.html">www.proteome.com/databases/index.html</a>
<i>C. elegans</i> Project & BLAST Server at Sanger Centre	<a href="http://www.sanger.ac.uk/Projects/C_elegans">www.sanger.ac.uk/Projects/C_elegans</a>
or Genome Sequencing Center	or <a href="http://genome.wustl.edu/gsc/C_elegans/elegans.shtml">genome.wustl.edu/gsc/C_elegans/elegans.shtml</a>
<i>C. elegans</i> WWW Server	<a href="http://elegans.swmed.edu">elegans.swmed.edu</a>

Wormbase, an outgrowth of Acedb, is a database of all genetic and genomic information for *C. elegans*, including genetic and physical maps. Wormpep is an up-to-date list of all predicted proteins from *C. elegans* maintained by the Sanger Centre. WormPD from Proteome Inc. is a *C. elegans* protein database with extensive literature annotation. Both the Sanger Centre and the Washington University Genome Sequencing Center (GSC) maintain information on the *C. elegans* genome project including a *C. elegans*-specific BLAST server and access to ftp downloads. The GSC maintains similar information for *C. briggsae*. The *C. elegans* WWW Server provides comprehensive links for *C. elegans*, including a literature search engine.

nematodes have homologues in *C. elegans*. For example, BLAST analysis revealed that 66% of *Meloidogyne incognita* genes have a *C. elegans* homologue ( $E < 10^{-5}$ ) (McCarter et al., 2000a). Key *C. elegans* genome resources are shown in Table 3.

An awareness of available nematode sequence data can save time and effort in the lab as well as greatly affect plans for experimental design. Periodic updates on the status of nematode gene sequence will be provided over the next several years as nematode EST data sets rapidly expand.

#### LITERATURE CITED

- Blaxter, M., M. Aslett, D. Guiliano, J. Daub, and the Filarial Genome Project. 1999. Parasitic helminth genomics. *Parasitology* 118:S39–S51.
- Blaxter, M. L., N. Raghavan, I. Ghosh, D. Guiliano, W. Lu, S. A. Williams, B. Slatko, and A. L. Scott. 1996. Genes expressed in *Brugia malayi* infective third-stage larvae. *Molecular and Biochemical Parasitology* 77:77–93.
- Daub, J., A. Loukas, D. I. Pritchard, M. Blaxter. 2000. A survey of genes expressed in adults of the human hookworm, *Necator americanus*. *Parasitology* 120:171–184.
- Marra, M. A., L. Hillier, and R. H. Waterston. 1998. Expressed sequence tags—ESTablishing bridges between genomes. *Trends in Genetics* 14:4–7.
- McCarter, J., P. Abad, J. T. Jones, and D. Bird. 2000a. Rapid gene discovery in plant-parasitic nematodes via expressed sequence tags. *Nematology*, in press.
- McCarter, J., D. McK. Bird, U. Rao, A. Kloek, S. Clifton, D. Pape, J. Martin, T. Wylie, S. Eddy, R. Waterston, and the Washington University GSC EST Team. 2000b. Progress toward high throughput gene discovery in parasitic nematodes: Initial findings from *Meloidogyne incognita* EST sequencing. *Journal of Nematology*, in press.
- McCombie, W. R., M. D. Adams, J. M. Kelley, M. G. Fitzgerald, T. R. Utterback, M. Khan, M. Dubnick, A. R. Kerlavage, J. C. Venter, and C. Fields. 1992. *Caenorhabditis elegans* expressed sequence tags identify gene families and potential disease gene homologues. *Nature Genetics* 1:124–131.
- Moore, T. A., S. Ramachandran, A. A. Gam, F. A. Neva, W. Lu, L. Saunders, S. A. Williams, and T. B. Nutman. 1996. Identification of novel sequences and codon usage in *Strongyloides stercoralis*. *Molecular and Biochemical Parasitology* 79:243–248.
- Popeijus, H., V. C. Blok, L. Cardle, J. Bakker, M. S. Phillips, J. Helder, G. Smant, J. T. Jones. 2000. Analysis of genes expressed in second-stage juveniles of the potato cyst nematodes *Globodera rostochiensis* and *G. pallida* using the expressed sequence tag approach. *Nematology*, in press.
- Tetteh, K. K., Loukas, C. Tripp, and R. M. Maizels. 1999. Identification of abundantly expressed novel and conserved genes from the infective larval stage of *Toxocara canis* by an expressed sequence tag strategy. *Infection and Immunity* 67:4771–4779.
- The *C. elegans* Sequencing Consortium. 1998. Genome sequence of the nematode *C. elegans*: A platform for investigating biology. *Science* 282:2012–2018.
- Unnasch, T. R., and S. A. Williams. 2000. The genomes of *Onchocerca volvulus*. *International Journal of Parasitology* 30:543–452.
- Waterston, R., C. Martin, M. Craxton, C. Huynh, A. Coulson, L. Hillier, R. Durbin, P. Green, R. Showkeen, N. Metzstein, T. Hawkins, R. Wilson, M. Berks, Z. Du, K. Thomas, J. Thierry-Mieg, and J. Sulston. 1992. A survey of expressed genes in *Caenorhabditis elegans*. *Nature Genetics* 1:114–123.