

# EST sequencing and usage

## Introduction

Expressed sequence tags, or ESTs, are single DNA sequencing reads made from complementary DNA (cDNA) clone libraries constructed from a known tissue source. Sequencing a large number of these clones from such a library allows one to sample the set of expressed genes, or transcripts, in that particular tissue and experimental state, thus providing a snapshot of the tissue's active genes under those defined conditions.

ESTs provide an excellent resource both for novel gene discovery, and for confirmation of in silico gene predictions from genomic DNA sequence. The physical cDNA clones themselves are useful resources in many types of functional biology studies, including gene inhibition experiments, their potential uses being dependent upon the type of cloning vector employed to construct them.

## Biological material

cDNA clones are the biological material from which EST sequences and finished cDNA sequences are produced. A set (or library) of cDNA clones are made from a biological tissue (or pooled tissues) by using routine molecular biological techniques. Cellular mRNA is extracted from the tissue, it is reverse transcribed to produce DNA complementary to the initial mRNA, and its incorporated into a plasmid, with appropriate vectors and linkers. The plasmids are then grown up in a suitable host such as *E. coli*.

Ideally the cDNA clones produced contain a single DNA insert derived from a gene transcript (mRNA) that was being expressed by the cell at the time the tissue was harvested. The tissue is generally selected to sample a particular cell type or organ of interest, and the time the tissue is taken is chosen to correspond to a developmental stage, and could be from a 'normal' individual/organism, one in disease state, or one having been subjected to experimental manipulation by: drug, environmental stress, or other means.

## Use of EST sequences

Large scale EST sequencing has become a comparatively inexpensive method for novel gene discovery. Prior to elucidating the full genome sequence of an organism, it provides a shortcut to the transcribed portions of the genome, allowing one to rapidly screen for novel genes, based on the various types of bioinformatic sequence searches available.

Consequently millions of human ESTs have been sequenced, as well as many from mouse, zebrafish and other species. However, with the ever-falling cost of automated DNA sequencing, many organisms have now had their whole genome sequenced. ESTs from these organisms remain useful, providing a means to locate the genes within the genome, by sequence similarity search, and also to lend supporting evidence to the genes predicted by other bioinformatic methods, in fact, they play a vital role in decoding and annotating each genome sequence as

it becomes available to the scientific community.

Analysis process and pipeline

## **Process outline**

The est\_db analysis process typically consists of:

- Loading cDNA library and sequencing primer information
- Batch-loading of EST sequences
- Public databank sequence submission and tracking
- EST clustering
- Similarity search analysis
- Protein sequence prediction
- Viewing and extracting results

### **Loading cDNA library and sequencing primer information**

At the initiation of a sequencing project, information describing the cDNA libraries being sequenced needs to be loaded into est\_db, this includes library names, descriptions, tissue source, developmental stage, etc. The sequencing primers to be used in the project also need to be specified.

### **Batch-loading of EST sequences**

The EST sequences themselves are loaded into the est\_db following base-calling and appropriate base-quality and vector clipping. Input is accepted as suitably-formatted FASTA files, and data is validated to confirm the library, clone name and sequencing primers and direction are recognised for each of the sequences.

### **Public databank sequence submission and tracking**

est\_db incorporates an automated system for submission of processed EST sequences to public databanks, tracking of the submission process, and subsequent updates to the databanks should they prove necessary. A final quality control (QC) check is generally performed on all EST sequences prior to their potential submission. This takes the form of one or more project-specific BLASTN similarity searches for any vector, adaptor, host, or contaminating organism sequences that might possibly be encountered in the laboratory environment(s) where the cDNA library construction and sequencing took place. Results of the searches are used to mark each cDNA clone stored in the database as clean or contaminated. EST reads from contaminated clones are not subsequently be submitted.

During the process of sequencing of a particular library, it is likely that some clones will be sequenced more than once, potentially duplicating EST reads. This might occur following resequencing of plate, following partial failure of the initial sequencing reads from the plate, or as a consequence of steps carried out to optimise the sequencing reaction conditions, etc. est\_db incorporates the facility to remove this redundancy, that might be introduced with every additional batch of ESTs sequenced and loaded from a particular library, ensuring the best set of ESTs (based on length and other controllable criteria) are selected for submission,

and thus subsequently made available to the public.

Using this battery of checks, and other user chosen parameters, the submission system then decides to submit a new entry, perform a sequence update, or withdraw an existing entry for each EST sequence. Tracking of the sequences is by their clone and library names, accession numbers received from the public databanks, and sequence checksum.

### **EST clustering and assembly**

Clustering of EST sequences, followed by their assembly is absolutely necessary to make sense of EST sequence data, removing the considerable redundancy always present in EST libraries, and allowing long transcripts to be reconstructed *in silico*. Rather than (re)implement a system to carry out these tasks, we chose to design the est\_db API and relational database schema to store, retrieve, and manipulate EST clustering, and contig assembly results in a generic fashion, independent of their source. Currently est\_db directly accepts results from the StackPACK analysis system, which has become a widely accepted package for EST clustering and assembly.

### **Similarity search**

Similarity searches can be performed to help identify individual ESTs, and also consensus EST sequences derived from EST clustering and assembly. Any number of nucleotide and protein sequences databases can be installed and used by est\_db in similarity searches with the WuBLAST family of programmes. Typically these databases would be specific to the aims of a particular project, based on species from which the cDNA library was constructed, the search for known or novel sequences, of a particular class, etc. est\_db uses a concise format to store and reconstruct the gapped pairwise sequence alignments resulting from such searches, allowing very large numbers of alignments to be stored, should this prove necessary, though a number of options are available to prune the hits stored to those most relevant.

The est\_db pipeline itself has features to allow lengthy analysis processes, such as some BLAST searches, to be split over many CPUs or machines, speeding their execution. These automated job creation and management features can reduce the time to completion for a given BLAST search from weeks to hours, and have been tested to more than three hundred machines simultaneously running BLAST analysis against a single est\_db and its database server.

### **Protein sequence prediction**

Protein sequences are predicted from consensus EST sequences as part of the analysis process. We have incorporated ESTScan into the est\_db system for this purpose.

### **Viewing and extracting results**

All results from the varied analyses are easily accessible through the est\_db web interface, with both textual and graphic views being used to render the

information. A single search form allows one to retrieve information pertaining to a particular cDNA clone, or EST read, as well as view EST clustering information pertaining to it, such as multiple sequence alignments of the cluster it belongs to, etc.

BLAST results can be searched by protein or gene name or description, with the BLAST hits, predicted protein sequence, etc being show together in a graphical panel, hyperlinked to display the gapped-alignments from the original searches. A step by step guide is available on [the searching the est\\_db analysis page](#) showing many of the available views.

Statistics derived from the analysis of a whole library can also be viewed. The total numbers of clones and ESTs sequenced from the library is displayed, with a breakdown available by sequencing direction, etc. A summary of EST clustering results is also included, allowing one to browse resultant clusters by size etc. A full description of the information available is provided on the [Viewing cDNA library analysis statistics page](#)

### Viewing cDNA library analysis statistics

#### Introduction

Detailed statistics are available from the analysis of each of the cDNA libraries sequenced and loaded into est\_db on the LibraryStatisticsView pages. These can be reached by clicking the library names on the est\_db web SearchPage.

Three types of information are provided:

- Total numbers of clones and ESTs
- EST clustering summary
- Analysis of abundant transcripts

#### Total numbers of clones and ESTs

Counts of the numbers of Clones and ESTs sequenced, broken down by sequencing direction, etc, are the first statistics presented (see Figure 1).

<b>Library name</b>	TAl1
<b>Description</b>	X. tropicalis ESTs from TNeu, TGas, TEgg & TTp clone libraries
<b>Clones</b>	160423
<b>5' ESTs</b>	153780
<b>3' ESTs</b>	79372
<b>Clone-paired ESTs</b>	72729
<b>Duplicate reads</b>	0
<b>Total ESTs</b>	233152

Figure 1. Clone and EST counts from the EST sequencing of a cDNA library, in this case over 230,000 *X. tropicalis* ESTs pooled from a number of embryonic, egg and

tadpole tissue sources.

### Clones

The number of complementary DNA (cDNA) clones that were successfully sequenced from the library i.e. yielding at least one EST sequence, and potentially more if sequenced from both ends, or duplicate reads were made from clones.

### 5' ESTs and 3' ESTs

DNA sequencing technology allows one to read from one or both ends of a DNA sequence. Generally the method of cloning mRNA transcripts to produce cDNA sequences, then insert them into vectors to produce cDNA clones, is a directional one. Therefore when the clone is sequenced, one can infer from which end of the cellular mRNA the sequence is derived.

Thus 5' and 3' ESTs are counts of the number of ESTs sequenced from the library, corresponding to the expected 5' and 3' ends, respectively, of the original mRNA transcripts that were used as the library source material.

### Clone-paired ESTs

The number of clones sequenced in the library that yielded both 5' and 3' ESTs.

### Duplicate reads

The number of duplicate reads (ESTs from the same clone, with the same direction) within the library.

### EST clustering summary

Statistics from EST clustering analysis of the library (if performed) are displayed (see Figure 2).

<a href="#">SuperClusters</a>	15803
<a href="#">Clusters</a>	17733
<a href="#">Consensei</a>	27260
<a href="#">Clustered ESTs</a>	207388
<a href="#">Singletons</a>	25764
<a href="#">Rev.Compded ESTs</a>	61234
<a href="#">Library contents</a>	

[VIEW ABUNDANT TRANSCRIPTS AND PROTEIN MATCHES](#)

Figure 2. EST clustering summary for the library. Shown are the numbers of SuperClusters, Clusters, ESTs clustered, singletons, etc. Note in this example counts and programmes which produced them (d2\_cluster, PHRAP, CRAW) are

components of the StackPACK package. See text for full description.

### **Superclusters**

Superclusters are sequence clusters of similar nucleotide sequence produced by the 'loose' clustering programme d2\_cluster. Each supercluster likely represents transcripts derived from the same gene, or several closely related genes. Named: library\_S\_x Where x is a number starting from 1. Example: TNeu\_S\_1

### **Clusters**

Each supercluster is further processed by the PHRAP programme, potentially splitting the supercluster to produce one or more clusters. Named: library\_C\_x\_y Where x is the supercluster number, and y is a number starting from 1. Example: TNeu\_C\_1\_1

### **Consensus sequences (consenseni)**

Clusters are analysed for sequence diversity with the CRAW programme. This results in the generation of one or more alignments, each with a consensus sequence. Multiple consenseni may suggest alternative-splicing of transcripts. Named: Library\_x\_y\_z, where x, y are the supercluster and cluster numbers from which the consensus is derived and z is number starting from 1. Example: TNeu\_1\_1\_1

### **Clustered ESTs**

The total number of ESTs clustered by d2\_cluster, every one thus belonging to a supercluster. EST names are of the form clone\_name.sp6 (sp6 being the sequencing primer)

### **Singletons**

**The total number of ESTs that were not clustered.**

### **Rev.Comped ESTs**

The number of ESTs that analysis suggested they should be included in clusters as the reverse complement of their original sequence. Usually these are 3' sequenced ESTs within the library. Occasionally they result from incorrect direction of insertion of the cDNA into the clone during library construction. Indicated by (R.C.) in alignments.

A breakdown of the results of the clustering analysis, counting how many superclusters and clusters were generated that fall in a certain size range is also given. Hotlinks allow the binned clusters to be reviewed (See Figure 3).

<u>Cluster Sizes</u>	SuperCluster	Count	Cluster	Count
	2	<a href="#">5142</a>	2	<a href="#">5740</a>
	3	<a href="#">2203</a>	3	<a href="#">2534</a>
	4	<a href="#">1522</a>	4	<a href="#">1717</a>
	5	<a href="#">1019</a>	5	<a href="#">1145</a>
	6	<a href="#">751</a>	6	<a href="#">845</a>
	7-8	<a href="#">1084</a>	7-8	<a href="#">1214</a>
	9-10	<a href="#">708</a>	9-10	<a href="#">789</a>
	11-20	<a href="#">1599</a>	11-20	<a href="#">1809</a>
	21-30	<a href="#">671</a>	21-30	<a href="#">734</a>
	31-40	<a href="#">334</a>	31-40	<a href="#">365</a>
	41-50	<a href="#">218</a>	41-50	<a href="#">233</a>
	51-60	<a href="#">104</a>	51-60	<a href="#">116</a>
	61-70	<a href="#">75</a>	61-70	<a href="#">89</a>
	71-80	<a href="#">70</a>	71-80	<a href="#">74</a>
	81-90	<a href="#">49</a>	81-90	<a href="#">49</a>
	91-100	<a href="#">26</a>	91-100	<a href="#">45</a>
	101-200	<a href="#">157</a>	101-200	<a href="#">164</a>
	201-300	<a href="#">35</a>	201-300	<a href="#">35</a>
	301-400	<a href="#">12</a>	301-400	<a href="#">12</a>
	401-500	<a href="#">1</a>	401-500	<a href="#">10</a>
	501-1000	<a href="#">13</a>	501-1000	<a href="#">11</a>
	1001+	<a href="#">10</a>	1001+	<a href="#">3</a>
<b><u>Available Searches</u></b>	<a href="#">Ensembl mouse (WuBLASTX)</a> <a href="#">Unigene X. laevis (WuBLASTN)</a> <a href="#">Swall (WuBLASTX)</a> <a href="#">TIGR XGI (WuBLASTN)</a> <a href="#">Xenopus sp. mRNA (WuBLASTN)</a> <a href="#">Ensembl human (WuBLASTX)</a>			

Figure 3. Supercluster and cluster sizes summarised for the library.

One can notice that there are a large number of both clusters and superclusters made up of only 2 ESTs (more than 5000 in each case), but considerably lower numbers of each as the absolute sizes increase, for example there are only 10 superclusters, which are composed of 1000 ESTs or more. Also shown are the available (BLAST) search results for the library.

### **Analysis of abundant transcripts**

One can investigate the abundant transcripts in a particular library, together with their analysis results by following the 'VIEW ABUNDANT TRANSCRIPTS AND PROTEIN MATCHES' link provided on each of the LibraryStatisticsView pages (see Figure 2). These show the consensus sequences, sorted by transcript abundance, as assessed by the clustering analysis.

View results page: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [11](#) [>>](#) ]

<b>Consensus name</b>	<a href="#">TAl1 1 1 1</a>
<b>Containing</b>	8797 ESTs
<b>Size rank</b>	1
<b>Best protein matches</b>	No significant hits

<b>Consensus name</b>	<a href="#">TAl1 7 300 1</a>																														
<b>Containing</b>	1209 ESTs																														
<b>Size rank</b>	2																														
<b>Best protein matches</b>	<table border="1"> <thead> <tr> <th>Alignment</th> <th>Ext DB Entry</th> <th>Description</th> <th>Score</th> <th>P value</th> </tr> </thead> <tbody> <tr> <td><a href="#">08AVK9</a></td> <td><a href="#">08AVK9</a></td> <td>Similar to nuclease sensitive element binding prot</td> <td>279</td> <td>4.70e-23</td> </tr> <tr> <td><a href="#">P21573</a></td> <td><a href="#">P21573</a></td> <td>Nuclease sensitive element binding protein 1 (Y bo</td> <td>270</td> <td>4.20e-22</td> </tr> <tr> <td><a href="#">08AVY9</a></td> <td><a href="#">08AVY9</a></td> <td>Similar to nuclease sensitive element binding prot</td> <td>270</td> <td>4.20e-22</td> </tr> <tr> <td><a href="#">090376</a></td> <td><a href="#">090376</a></td> <td>Y-box binding protein</td> <td>264</td> <td>1.90e-21</td> </tr> <tr> <td><a href="#">006066</a></td> <td><a href="#">006066</a></td> <td>Nuclease sensitive element binding protein 1 (Y bo</td> <td>260</td> <td>4.60e-21</td> </tr> </tbody> </table>	Alignment	Ext DB Entry	Description	Score	P value	<a href="#">08AVK9</a>	<a href="#">08AVK9</a>	Similar to nuclease sensitive element binding prot	279	4.70e-23	<a href="#">P21573</a>	<a href="#">P21573</a>	Nuclease sensitive element binding protein 1 (Y bo	270	4.20e-22	<a href="#">08AVY9</a>	<a href="#">08AVY9</a>	Similar to nuclease sensitive element binding prot	270	4.20e-22	<a href="#">090376</a>	<a href="#">090376</a>	Y-box binding protein	264	1.90e-21	<a href="#">006066</a>	<a href="#">006066</a>	Nuclease sensitive element binding protein 1 (Y bo	260	4.60e-21
Alignment	Ext DB Entry	Description	Score	P value																											
<a href="#">08AVK9</a>	<a href="#">08AVK9</a>	Similar to nuclease sensitive element binding prot	279	4.70e-23																											
<a href="#">P21573</a>	<a href="#">P21573</a>	Nuclease sensitive element binding protein 1 (Y bo	270	4.20e-22																											
<a href="#">08AVY9</a>	<a href="#">08AVY9</a>	Similar to nuclease sensitive element binding prot	270	4.20e-22																											
<a href="#">090376</a>	<a href="#">090376</a>	Y-box binding protein	264	1.90e-21																											
<a href="#">006066</a>	<a href="#">006066</a>	Nuclease sensitive element binding protein 1 (Y bo	260	4.60e-21																											

<b>Consensus name</b>	<a href="#">TAl1 20 176 4</a>																														
<b>Containing</b>	789 ESTs																														
<b>Size rank</b>	3																														
<b>Best protein matches</b>	<table border="1"> <thead> <tr> <th>Alignment</th> <th>Ext DB Entry</th> <th>Description</th> <th>Score</th> <th>P value</th> </tr> </thead> <tbody> <tr> <td><a href="#">P16878</a></td> <td><a href="#">P16878</a></td> <td>Keratin, type II cytoskeletal (XENCK55(5/6))</td> <td>1481</td> <td>1.30e-150</td> </tr> <tr> <td><a href="#">09I9P5</a></td> <td><a href="#">09I9P5</a></td> <td>Inner-ear cyokeratin</td> <td>1387</td> <td>1.10e-140</td> </tr> <tr> <td><a href="#">093532</a></td> <td><a href="#">093532</a></td> <td>Keratin, type II cytoskeletal cochleal (Cytokerati</td> <td>1292</td> <td>1.40e-130</td> </tr> <tr> <td><a href="#">08HZR5</a></td> <td><a href="#">08HZR5</a></td> <td>Keratin 7</td> <td>1269</td> <td>3.90e-128</td> </tr> <tr> <td><a href="#">096CL4</a></td> <td><a href="#">096CL4</a></td> <td>Similar to keratin 6A</td> <td>1204</td> <td>4.20e-127</td> </tr> </tbody> </table>	Alignment	Ext DB Entry	Description	Score	P value	<a href="#">P16878</a>	<a href="#">P16878</a>	Keratin, type II cytoskeletal (XENCK55(5/6))	1481	1.30e-150	<a href="#">09I9P5</a>	<a href="#">09I9P5</a>	Inner-ear cyokeratin	1387	1.10e-140	<a href="#">093532</a>	<a href="#">093532</a>	Keratin, type II cytoskeletal cochleal (Cytokerati	1292	1.40e-130	<a href="#">08HZR5</a>	<a href="#">08HZR5</a>	Keratin 7	1269	3.90e-128	<a href="#">096CL4</a>	<a href="#">096CL4</a>	Similar to keratin 6A	1204	4.20e-127
Alignment	Ext DB Entry	Description	Score	P value																											
<a href="#">P16878</a>	<a href="#">P16878</a>	Keratin, type II cytoskeletal (XENCK55(5/6))	1481	1.30e-150																											
<a href="#">09I9P5</a>	<a href="#">09I9P5</a>	Inner-ear cyokeratin	1387	1.10e-140																											
<a href="#">093532</a>	<a href="#">093532</a>	Keratin, type II cytoskeletal cochleal (Cytokerati	1292	1.40e-130																											
<a href="#">08HZR5</a>	<a href="#">08HZR5</a>	Keratin 7	1269	3.90e-128																											
<a href="#">096CL4</a>	<a href="#">096CL4</a>	Similar to keratin 6A	1204	4.20e-127																											

Figure 4. Browsing abundant transcripts. Shown (in this case) are the three most abundant transcripts, as assessed from the EST clustering analysis. If a translated BLAST (WuBLASTX) search of the consensus sequence produced significant protein sequence hits, the five highest scoring ones are displayed. In this example a transcript similar to that for y box binding protein-1 is highly expressed (second most abundant) in this library of *X. tropicalis* ESTs.