

# WORMBASE – NEMATODE BIOLOGY AND GENOMES

Paul Davis, Wellcome Trust Genome Campus.

WormBase Consortium, (PI List) Lincoln Stein – OICR, Paul Sternberg – CALTECH,  
John Spieth - GSC, Richard Durbin – WTSI

WormBase [www.wormbase.org](http://www.wormbase.org)

WormBase is the major public online database resource for the *Caenorhabditis* research community. The database was developed primarily for the nematode *C. elegans* but expanded to host genomes and biological data from other closely related nematode species including *C. briggsae*, *C. remanei*, *C. brenneri*, *C. japonica* and *Pristionchus pacificus*. WormBase has developed tools to mine the data held within the database and compare the hosted species. Over the years we have developed a variety of curation pipelines which often begin in a "first-pass" literature curation step. This involves a brief overview of the literature before directing it to specialised data curators who extract all relevant information. Curators focus on particular data types or experimental techniques such as gene structure changes (see the Sequence curation poster), variations, phenotypes or RNAi and their expertise in these fields make curation efficient. WormBase works with many other groups and consortiums to validate, process and integrate both large and small scale data resources. WormBase also provides data that will be of interest to the wider biomedical and bioinformatics communities allowing researchers to utilise the information and techniques offered by nematodes to study wider aspects including medicine and disease.

## Automated First Pass Paper Curation

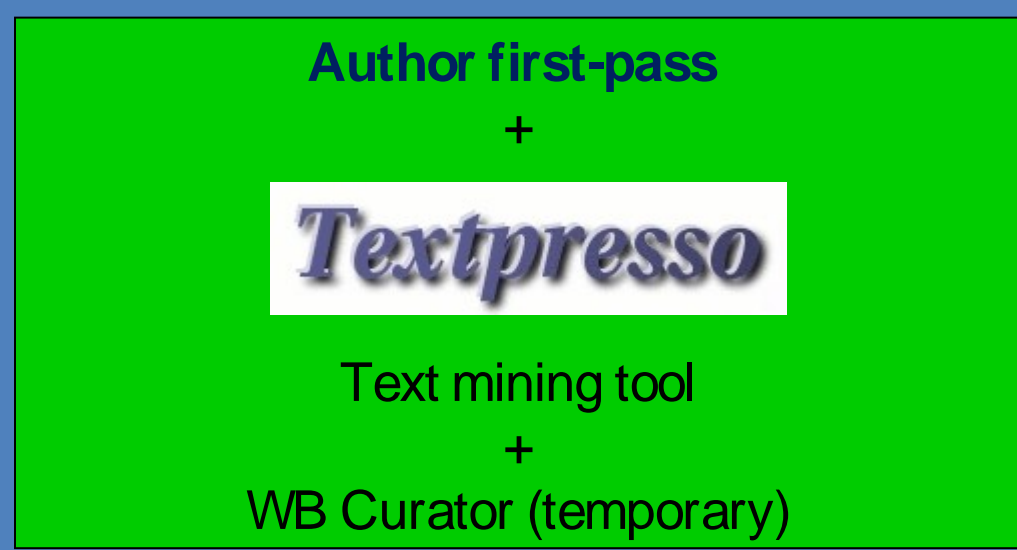
### Paper and Data Type Identification:

There is a project underway to move from a manual approach to a semi automated pipeline with author input. Currently we are in a transitional phase moving from WB Curators towards Authors and Text Mining.



searched  
using keyword  
'elegans',  
Manual selection  
of papers

PDFs Download  
Automatically and  
stored in a database.



27 Data types extracted (2009)

Data curator      Data curator      Data curator      Data curator

### Author First-Pass Form

Please click the box next to the type of data your publication includes.

If this is not a primary research article, please click here. You may ignore the fields below. Thank you.  
Click the "X" to find out more about the data type.  
This paper has already had data submitted, using a new

#### Species :

- ☐ *C. elegans* ? Add information.collapse
- ☐ *C. elegans* other than *N2* (wildtype) ? Add information.collapse
- ☐ Nematode species other than *C. elegans* ? Add information.collapse
- ☐ Non-nematode species ? Add information.collapse

#### Gene Identification and Mapping :

- ☒ Genes studied in this paper ? Add information.generated
- ☐ Newly cloned gene ? Add information.generated
- ☐ Newly created alleles ? Add information.generated
- ☐ Genetic mapping data ? Add information.generated

#### Gene Function :

- ☒ Mutant, RNAi, Overexpression, or Chemical-based Phenotypes. Please specify your data type.
  - ☐ Small-scale RNAi (less than 100 individual experiments) ? Add information.generated
  - ☐ Large-scale RNAi (greater than 100 individual experiments) ? Add information.generated
  - ☐ Overexpression phenotype ? Add information.generated
  - ☐ Chemicals ? Add information.generated
- ☐ Microarray analysis ? Add information.generated
- ☐ Tissue or cell site of action ? Add information.generated
- ☐ Time of action ? Add information.generated
- ☐ Molecular function of a gene product ? Add information.generated
- ☐ Homology of a human disease-associated gene ? Add information.generated

#### Interactions :

- ☐ Genetic interactions ? Add information.generated
- ☐ Functional complementation ? Add information.generated
- ☐ Gene product interaction ? Add information.generated

#### Regulation of Gene Expression :

- ☐ New expression pattern for a gene ? Add information.generated
- ☐ Microarray ? Add information.generated
- ☐ Alterations in gene expression by genetic or other treatment ? Add information.generated
- ☐ Regulatory sequence features ? Add information.generated
- ☐ Position frequency matrix (PFM) or position weight matrix (PWM) ? Add information.generated

#### Reagents. :

- ☒ *C. elegans* antibodies ? Add information.generated
- ☒ Integrated transgenes ? Add information.generated
- ☐ Transgenes used as tissue markers ? Add information.generated

#### Protein Function and Structure :

- ☐ Protein analysis in vivo ? Add information.generated
- ☐ Analysis of protein domains ? Add information.generated
- ☐ Covalent modification ? Add information.generated
- ☐ Structural information ? Add information.generated
- ☐ Mass spectrometry ? Add information.generated

#### Genome Sequence Data :

- ☐ Gene structure correction ? Add information.generated
- ☐ Sequencing repeat alleles ? Add information.generated
- ☒ New SNPs, not already in WormBase ? Add information.generated

#### Cell Data :

- ☐ Ablation data ? Add information.generated
- ☐ Cell function ? Add information.generated

#### In Silico Data :

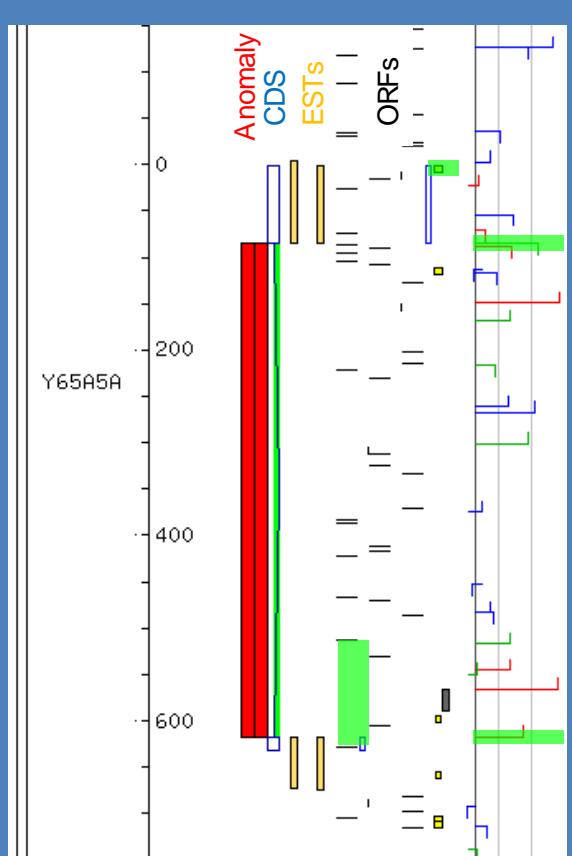
- ☐ Phylogenetic analysis ? Add information.generated
- ☐ Other bioinformatics analysis ? Add information.generated

#### Other :

Authors will be automatically contacted once their paper is downloaded, at this point the paper will not be visible to the curators. Once a set period of time has elapsed the text mining will be conducted and work distributed within the consortium.

## Curation Anomaly display

(see sequence curation poster for more details)

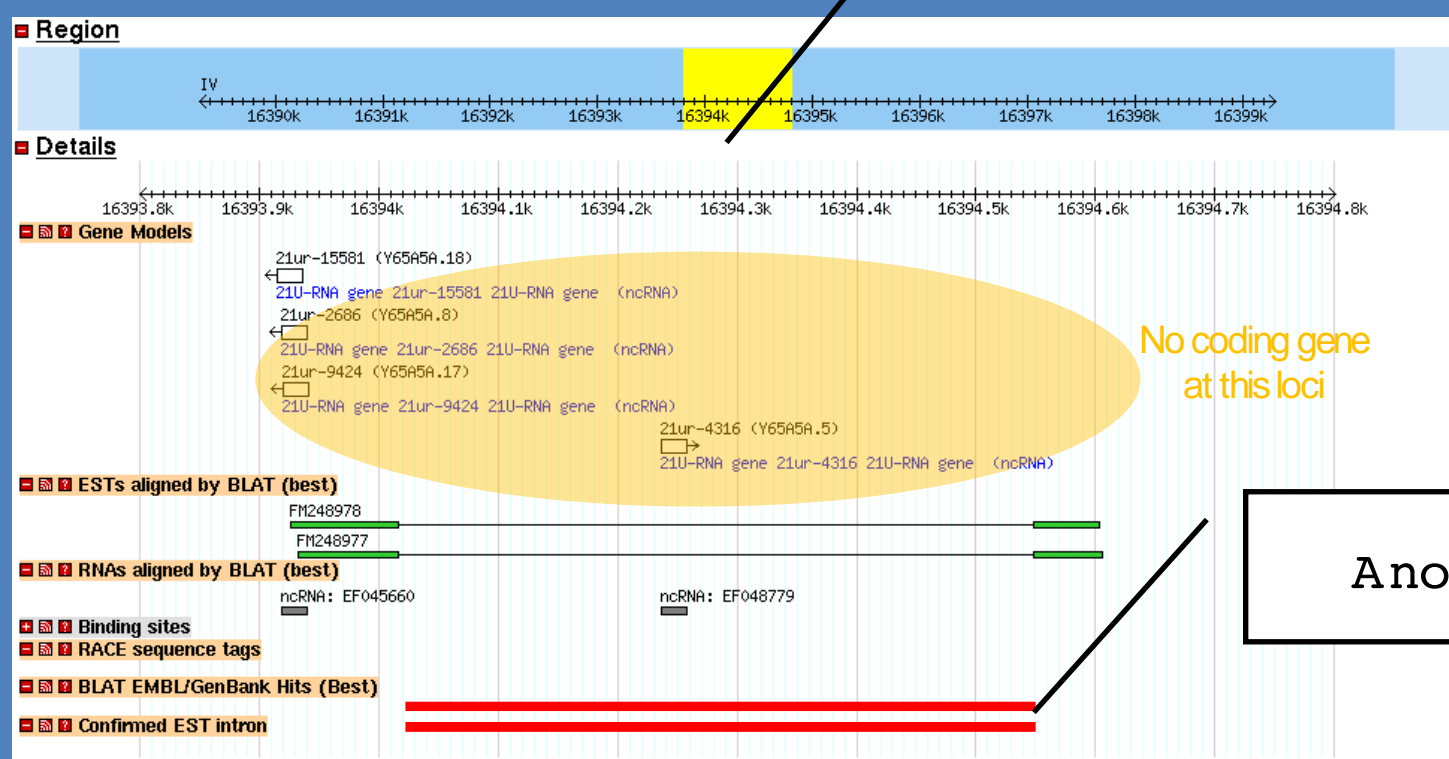


Curators use the Fmap within ACeDB to add/modify genes that go into the WormBase dataset. Curators use pre-computed anomalies to identify gene models that require attention as well as missing genes. Here we have an Fmap display with a Confirmed intron not in gene model anomaly. The curator created this two exon gene model based on the 2 EST sequences which contain a single intron.

Users can choose to see these anomalies/possible errors in the genome browser by selecting from these tracks.

- ☐ Misc. Curation Anomalies (All on 1st pass)
- ☐ Confirmed intron
- ☐ Confirmed EST intron
- ☐ Repeat coverage even
- ☐ Unmatched protein
- ☐ Merge genes by protein
- ☐ Mismatched EST
- ☐ Unmatched TSL
- ☐ Unmatched gene predictions
- ☐ Split gene by protein groups
- ☐ Unmatched TSL
- ☐ Unmatched WGA
- ☐ Weak intron splice

Genome browser prior to new gene annotation.



Anom a lies

## Gene Summary Page: Concise descriptions

Curators are working to produce a manually generated concise description for all *C. elegans* genes (extending to include tier II nematodes). The aim of this is to produce an abstract like summary of the gene and it's function so that WormBase users get a good understanding of the gene, with a minimum amount of effort.

### Gene Summary for unc-13

Specify a gene using a gene name (unc-20), a predicted gene ID (R1345.9), or a protein ID (CE0271.1) [unc-13]

Identification IDs:	Main name	Sequence name	WB Gene ID
[Identification] [Location] [Function] [Expression] [Gene ontology] [Genetics] [Phenotype] [Regulation] [Bibliography]	unc-13 (UNC-13) via Person evidence: Jonathan Hodgkin	ZK524.2	WBGene00006752
Concise Description:	unc-13 encodes at least five protein isoforms that regulate neurotransmitter release by altering the conformation of syntaxin. UNC-13 proteins are required for normal pharyngeal pumping and thrashing in liquid, normally short lifespan, normally large brood sizes, and full adult body size. UNC-13 proteins have orthologs in vertebrates and <i>Drosophila</i> . UNC-13 proteins are complex, with multiple C2, prothol ester-binding, and DUF1341 domains. UNC-13 protein form is localized to most or all synapses; many of the unc-13 mutant alleles with viable phenotypes are transcript-specific, while homozygotes with an unc-13 null (deletion) allele die as paralyzed first-stage larvae. [details]		

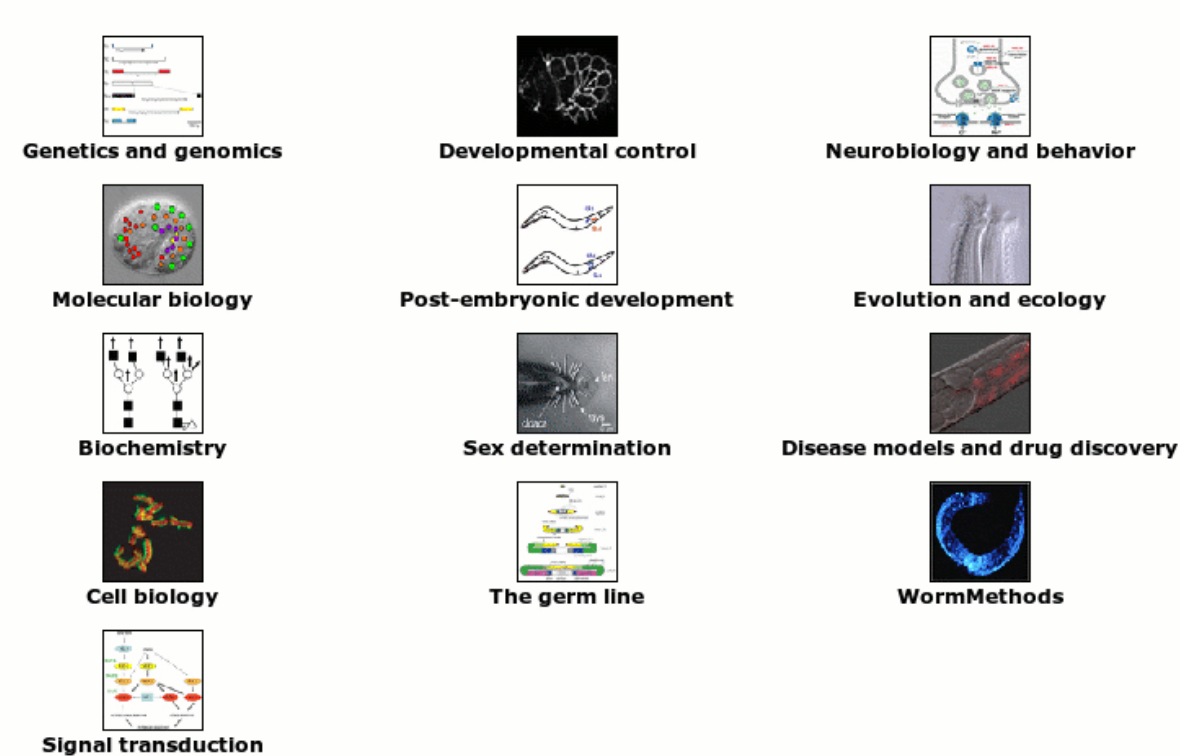
## WormBook

WormBook is the online text companion to WormBase, the *C. elegans* model organism database. WormBook contains original reviews on all aspects of *C. elegans* biology and up-to-date descriptions of technical procedures used to study this animal.



WormBook is a comprehensive, open-access collection of original, peer-reviewed chapters covering topics related to the biology of *Caenorhabditis elegans* (*C. elegans*). WormBook also includes WormMethods, an up-to-date collection of methods and protocols for *C. elegans* researchers.

#### WormBook Sections



Complete Chapter Listings  
By Section | By Publication Date

Photo Credits

## Phenotype ontology

Our Phenotype Ontology has been modified to curate nematode species other than *C. elegans* "N2" strain

### Phenotype Ontology

A hierarchy-based ontology  
1823 terms, 66% defined,  
55% associated with a variation

- Classes
  - Variant
    - behavior\_variant
    - development\_variant
    - morphology\_variant
      - cell\_morphology\_variant
      - organ\_system\_morphology\_variant
      - organism\_morphology\_variant
        - body\_region\_morphology\_variant
        - developmental\_morphology\_variant
        - adult\_body\_morphology\_variant
        - clonal\_body\_morphology\_variant
        - larval\_body\_morphology\_variant
        - lumpy
        - organism\_morphology\_variable
        - organism\_segment\_morphology\_variant
        - body\_morphology\_variant
        - head\_morphology\_variant
        - tail\_morphology\_variant
        - sexually\_dimorphic\_morphology\_variant
        - pericellular\_component\_morphology\_variant
        - physiology\_variant
        - pigmentation\_variant
        - unclassified
- Relationships
  - Obsolete

### Modifications required

Changes to Term Names  
from " \_abnormal" to " \_variant" .

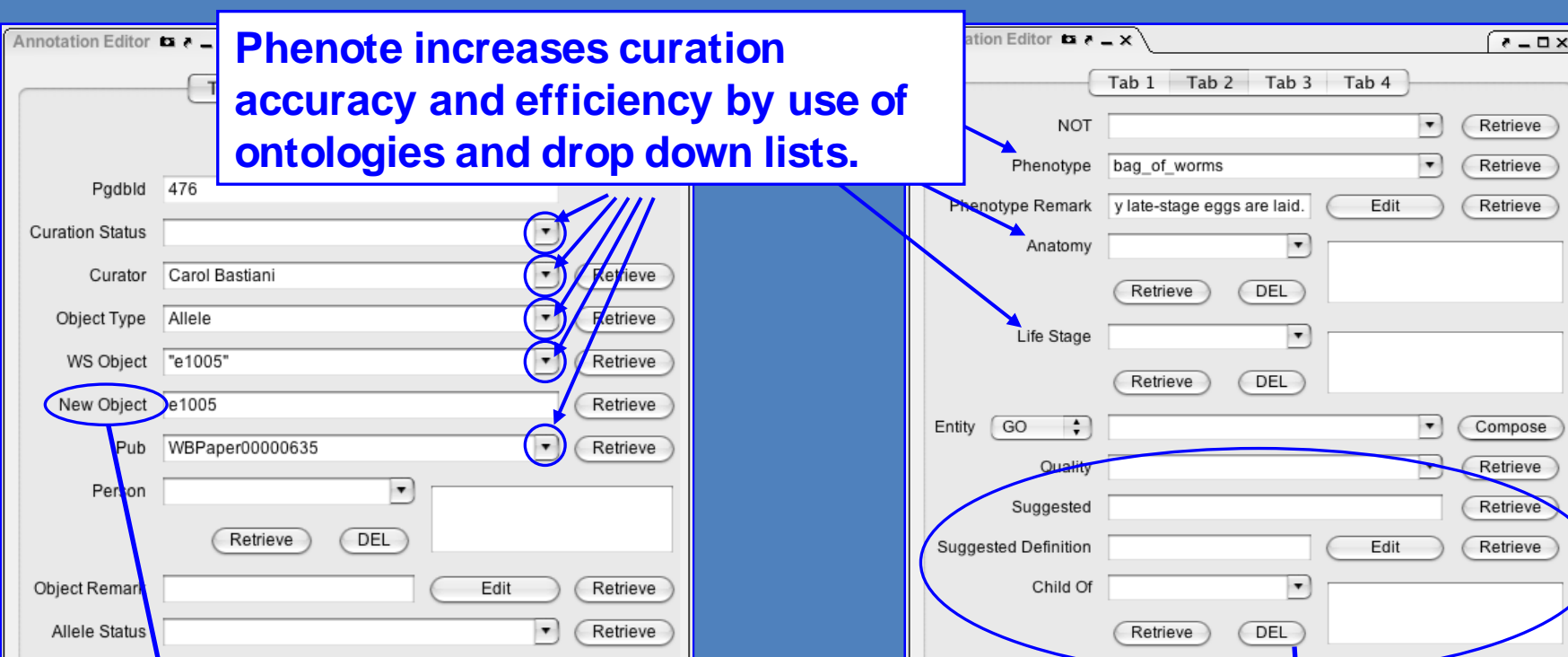
Changes to Definitions  
Use of " control animals" rather than " wild-type" or " N2" (*C. elegans* strain).

*C. elegans*-specific terminology, e.g., " hermaphrodite" , were removed from definitions when possible.

Example:  
WBGene0000037: **egg\_morphology\_abnormal**  
Def: "Any deviation in the overall structure or appearance of fertilized oocytes that are deposited by adult hermaphrodites."

Changed to:  
WBGene0000037: **egg\_morphology\_variant**  
Def: "Any variation in the overall structure or appearance of fertilized oocytes that are laid compared to those laid by control animals." "synonym": "egg\_morphology\_abnormal" (The " \_abnormal" version of the term is kept as a synonym so people used to these terms will still be able to find them.)

Phenotype curation captures multiple attributes reported by authors and requires the efforts of many data curators



COORDINATED WITH OBJECT CURATORS:  
If object (allele or transgene) does not exist in the latest release of the database, an e-mail is automatically sent to the curator responsible for creating those objects.

COORDINATED WITH ONTOLOGY CURATOR:  
Phenotype curators can request a term, send a suggested definition and hierarchy placement through the Phenote interface. New terms are automatically assigned to the record when they are approved.

### OTHER ATTRIBUTES CAPTURED INCLUDE:

Genotype, Treatment, Nature of allele (recessive, semi-dominant, dominant), Penetrance (incomplete, low, high, complete), Maternal effect (strictly maternal, with maternal effect), Paternal effect, Temperature sensitivity, Haploinsufficiency, Allele type (amorph, hypomorph, etc.).

Phenotypes are linked to genes through allele or RNAi curation

Gene	Summary	Phenotype	Supporting Evidence
unc-13	Summary	Phenotypes reported as observed	Supporting Evidence
unc-13	Summary	Phenotypes reported as NOT observed	Supporting Evidence

3626 / 23709\* genes with alleles were annotated with phenotype data (includes NOT annotations) as of WS200

	May 2008 WS188	March 2009 WS200
Allele-phenotype connections	9771	15951
Alleles Curated (total # alleles)	28% (15326)	34% (17448)
Papers curated (total papers flagged)	NA	23% (+125 unflagged papers)