

EPICViz: An interactive visualization of *C. elegans* embryogenesis and gene expression

Melissa Chiasson
chiasson@uw.edu

Timothy Durham
tdurham@uw.edu

Andrew Hill
ajh24@uw.edu

Ning Li
ningli30@uw.edu

ABSTRACT

Caenorhabditis elegans (*C. elegans*) is a widely used model for studying how a multicellular organism with many different tissue types develops from a single fertilized cell. Gene expression plays an important role in development, and the Expression Patterns In *C. elegans* (EPIC) study in the Waterston Lab at UW has generated measurements of the expression patterns for 227 genes in each cell of developing *C. elegans* embryos at about 1 min intervals for the first ~350 min of development. In addition, this dataset contains positional coordinates for each cell, its diameter, and its cell type and lineage. We developed a visualization tool to allow users in the *C. elegans* community to assess the gene expression patterns across development in particular tissues and cell lineages. We include a 3D plot of the embryo, 2D projections of the 3D embryo, a lineage tree, a principal components analysis (PCA) plot to cluster cells based on gene expression, and a heatmap showing expression patterns of specific genes. Users can manipulate the 3D plot to change views, highlight specific tissue types and lineages, and create gene expression enrichment reports for cells of interest. This visualization will make the EPIC data set accessible to a broad spectrum of *C. elegans* researchers and provide a platform for exploring development and facilitating hypothesis generation.

Author Keywords

Caenorhabditis elegans, *C. elegans*, development, embryogenesis, gene regulation, gene expression.

INTRODUCTION

C. elegans is a small roundworm used widely as a model organism in genetics and genomics. Its development has been well studied; each worm takes around 14 hours to grow from a single fertilized cell to a hatched larvae with 558 cells. This process of embryonic development progresses in a stereotyped pattern that follows an invariant cell lineage; the same branches in this tree always produce the same tissues in the hatched worm [1]

(Figure 1).

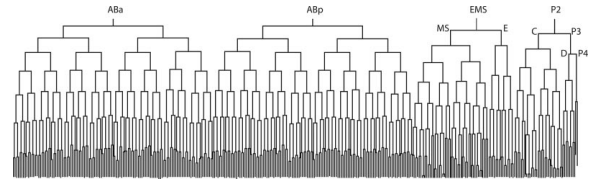


Figure 1. Lineage tree of the first 350 cells in *C. elegans* embryogenesis. Adapted from [3].

Development is the process by which cells derived from a single fertilized egg divide and differentiate to form the diverse tissues of the adult organism. In this process, cells start out very similar to each other, but over time express different sets of genes that specify the unique fate of each cell. Some of development is deterministic, encoded in the cells themselves; however each cell's fate is also influenced by the local environment and interactions with its neighbors. Thus, to understand development, we must consider the gene expression patterns of each cell within the broader context of the developing embryo.

Despite intense study of *C. elegans*, currently there does not exist a resource where one can interactively visualize *C. elegans* embryo development and interrogate how gene expression patterns change with time, lineage, or cell type. This limits exploration of data that could yield novel insights into the gene regulation of development.

RELATED WORK

C. elegans Gene Expression and Development Analysis and Visualizations

Because *C. elegans* development follows an invariant pattern, much work has focused on making computational tools that can track and annotate cells from microscopy images or movies [2], [3]. These tools are helpful for generating datasets like the one we used (described in **METHODS**) by identifying cells and tracking their positions over developmental time. Automated tracking and annotation of cells enables higher-level analysis of biological characteristics, like gene expression, that influence development.

Indeed, computational analysis of a subset of the data we use in our visualization yielded insights into *C. elegans* embryogenesis [4]. By analyzing expression of 127 genes,

the study found expression patterns associated with terminal tissue type and spatial positions in the embryo; it also found sequential patterns of gene induction that could be indicative of regulatory cascades. This kind of analysis is potentially of interest to a broad spectrum of the *C. elegans* research community, but not all of them have the computational expertise necessary to analyze the raw data. We hope that our visualization can help all in the community to see new patterns in this data set, while increasing the accessibility of the data to biologists without computational expertise.

In addition to developing embryogenesis cell tracking and gene expression software, 3D modeling of *C. elegans* anatomy has been explored. The OpenWorm project is an interactive 3D model of *C. elegans* anatomy and behavior, focusing on the muscular system and neural network [5]. Users can query specific cells and visualize the connections between them to better understand how neurons and muscles coordinate.

Work in Other Model Organisms

Gene expression visualizations have been explored in other model organisms. ViBE-Z is a software package that generates graphics of gene expression in the zebrafish larval brain [6]. This helps in identifying colocalization of specific genes in particular substructures of the brain, but does not allow for interactive exploration. In *Drosophila*, a tool called MULTEESUM was developed to compare spatial and temporal gene expression data at a specific stage in embryogenesis [7]. At this stage, the drosophila embryo is shaped like a hollow tube, so the authors could convert a 3D embryo to two dimensions. In the visualization, users can select specific cells in this 2D rendering and view small multiples that display expression patterns for genes of interest over time. Both of these tools include features, like small multiples and brushing and linking, we sought to include in our own visualization.

METHODS

Data

The Expression Patterns In *C. elegans* (EPIC – <http://epic2.gs.washington.edu/Epic2>) project has generated a dataset that describes the spatial orientation of every cell during the first ~350 minutes of *C. elegans* embryogenesis, its developmental lineage and cell fate, and expression measurements for a set of 227 genes. These values were derived from confocal microscopy movies of developing *C. elegans* embryos (Figure 2). In these movies, the nucleus of every cell is highlighted with green fluorescent protein (GFP), while the nucleus of cells expressing the gene of interest additionally light up with mCherry, a red fluorescent protein. Therefore, expression of a gene is indicated by the presence of red

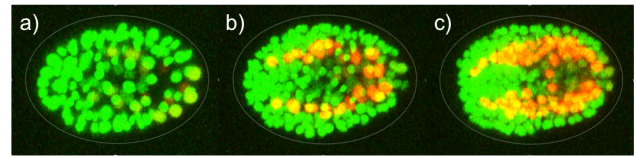


Figure 2. Expression of gene *hnd-1*, a muscle-specific transcription factor, at a) 0 minutes, b) 100 minutes, and c) 200 minutes. Notice the bilateral patterning of expression.

fluorescence, while the locations of surrounding nuclei are indicated by green fluorescence.

The EPIC data processing pipeline analyzes the movie for each gene to identify specific cells, record the positions of these cells, and quantify the fluorescence level as a measure of expression level. One challenge in working with these data is that the cell division timing (i.e. developmental rate) is not necessarily the same from movie to movie, and neither are the nucleus positions. To address the time synchronization of the movies and make gene expression patterns comparable, we used a version of the data in which each time point was mapped to a unified pseudotime scale defined based on the relative lifespans of individual identified cells in the different movies. To address the positional variability of the nuclei in the movies, we simply selected a single movie that was particularly well annotated, extracted and mean-centered the x, y, z coordinates, and used these as reference positions for mapping the gene expression values from all other movies. Another consideration in our data preparation was that, due to technical details of the way the fluorescence experiments are done, the levels of fluorescence are not comparable across movies. However, the timing of expression induction and the pattern of cells expressing a gene are comparable. Thus, we binarized the gene expression by setting a threshold for calling a gene expressed. We consider a gene expressed in a cell at a particular time if that cell measures greater than 2000 red fluorescence units and this fluorescence level is at least 10% of the maximum red fluorescence level in the entire movie. Another artifact of the data collection protocol is that the fluorescence data do not contain information about when genes turn off. Thus, once a gene is called expressed in a cell, it is also called expressed in all progeny from that cell for the rest of the time series.

3D Plot and Small Multiples

Most of our understanding of cell-cell connections and spatial relationships in development comes from two-dimensional representations, either views through a microscope or in other representations like the lineage tree. While these representations can be very effective, we sought to address an important limitation with our 3D approach. Embryogenesis is a biological process that takes place in three dimensions, and the orientations and connections among cells play an essential role in this process. Being able to identify cells and to watch them

undergo divisions and migrations in three dimensions, from any orientation, can greatly facilitate our understanding of which lineages are close together and which cell-cell connections might be important in forming different tissues.

To that end, we used the HTML library X3DOM to build a three-dimensional plot of all cells. Cells are represented as spheres that are sized and positioned at each time point based on the EPIC data.

Lineage Tree

We also include a plot showing the *C. elegans* lineage tree. The lineage tree is a powerful visualization for showing the relatedness of cells, and by highlighting cells in the lineage tree and the 3D plot the user can see how cells from different developmental branches are positioned relative to each other in the embryo. This graphic is also widely used in the *C. elegans* community, and provides a familiar view of cell origins and fate. One challenge is that the lineage tree is very broad. In order to show details of particular lineages while also allowing the user to see the broader context of the tree we implemented the D3 tree layout with a moveable Cartesian distortion (<https://github.com/d3/d3-plugins/tree/master/fisheye>).

Gene Expression Plots

We sought to give the user the ability to assess global expression changes to find interesting patterns, and then provide ways to access more detailed information about the genes involved in these patterns of interest. One challenge was dealing with the scale of the data set. At the final time point in the series there are almost 550 cells, and we must display expression information for 227 genes, which means we must show on the order of 120,000 data points. We used principle components analysis (PCA; performed in R) to cluster the cells by gene expression pattern and provide a high level summary of these trends. In order to allow users to further investigate a cluster of interest, we provide an additional view showing the gene expression as a matrix with genes on the rows and cells on the columns. Any gene that is expressed in a particular cell is represented by a filled-in rectangle in the matrix. Thus users can see which genes are driving the PCA clustering.

RESULTS

Overview

Our final visualization, EPICViz, consists of a user selection menu, 3D plot of the embryo, 2D small multiples of the embryo, a lineage tree, a PCA plot of gene expression, and a heatmap showing expression for all cells and genes (Figure 3). We describe each element in detail below.

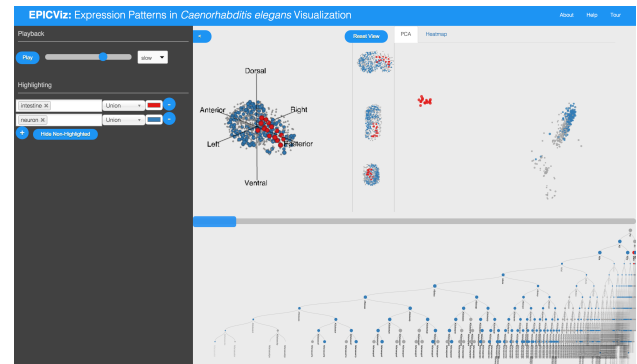


Figure 3. Overview of EPICViz.



Figure 4. User selection controls.

User Selection Menu

In a collapsible menu on the left hand side of our visualization, users will see options for playback and highlighting. Under playback, users can play and pause the time slider at specific time points, as well as select the speed at which the visualization will animate.

Under highlighting, users can select lineages, cell types, and tissue types of interest by typing or scrolling through a drop-down menu. Cells within these categories will then be highlighted in the color of the user's choice. Multiple selections can be made in the drop-down, and the user can specify whether the highlighting color should be applied to the union or the intersection of the cells specified by the drop-down selections. Cells that are not highlighted will appear as smaller grey spheres, but the user also has the choice to hide all non-highlighted cells.

As an example of how EPICViz can be used to understand gene expression and development, our figures will include highlights of intestine (in red) and neurons (in blue) (Figure 4).

3D Plot

The 3D plot shows cells as spheres, which divide and migrate as development progresses (Figure 5, panel A). The user can drag the 3D plot to see it from a variety of views. Mousing over a cell reveals the cell name and its coordinates in a pop-up. Clicking on a cell emphasizes it with a yellow outline and the representations of that cell in all of the other plots are similarly highlighted, allowing the user to quickly reference corresponding information in the different facets of the visualization.

In our example, intestinal cells are clustering towards the posterior pole of the embryo and are mostly internal. Neurons, on the other hand, are distributed more widely across the embryo and are concentrated on the exterior of the embryo.

2D Small Multiples

The 2D small multiples lie to the right of the 3D (Figure 5, panel B). These are orthogonal views along each axis, and they are linked to the 3D plot both in terms of time and in terms of brushing. The top view is from the left, the center is from the dorsal aspect looking straight down, and the bottom view is from the posterior looking along the frontal axis. Clicking on a small multiple moves the 3D plot to that view.

The small multiples show again the wider distribution of neurons on the exterior of the embryo in blue and the internal clustering of the intestinal cells towards the posterior in red. One can also easily see the left-right asymmetry of the positions of these tissue progenitors. Such asymmetries in development are important for the cells to construct the proper body plan [8].

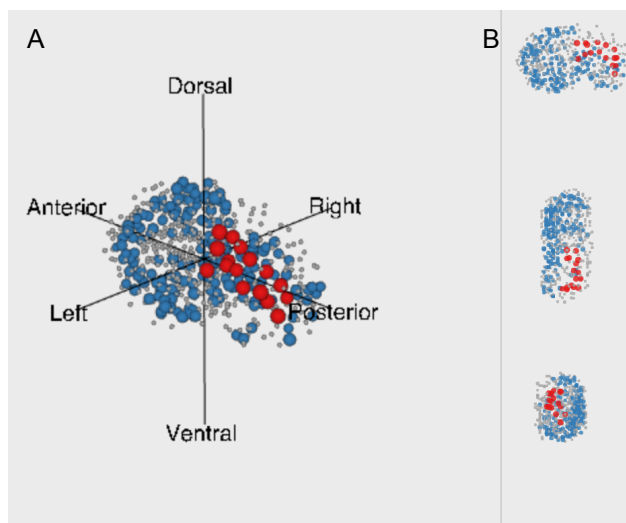


Figure 5. A) 3D plot and B) 2D small multiples. Intestine cells highlighted in red, neuron cells in blue.

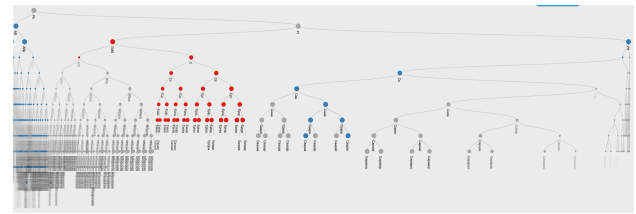


Figure 6. Lineage tree with intestinal cells in red, neurons in blue.

Lineage Tree

Cells within the lineage tree are populated as developmental time progresses. Highlights from the user selection menu are also linked to the tree, so cells of interest are easy to discriminate (Figure 6). Cells at the leaves of the lineage tree can be selected, resulting in a yellow highlight around the cell. As with the 3D view, these yellow highlights are propagated to all other plots to allow the user to quickly find their cell of interest throughout the visualization.

From the lineage tree, we can see how closely related intestinal cells and neurons are. Intestinal cells are descended from the EMS cell, which shares a parent with both P2-descended and AB-descended neuron cells.

PCA Plot of Gene Expression

Because our heatmap has values for approximately 500 cells and 227 genes at the end of the development time frame, we wanted a way to display gene expression information with reduced dimensionality. For this, we performed PCA for each cell and its 227 gene expression values.

Each dot within the PCA represents one cell, and users can watch as cells cluster into specific areas with developmental time (Figure 7). If a user clicks on a cell with a particular highlight (as defined in the user selection menu), a gene report pop-up is generated, which lists the gene name, WormBase ID (which links out to WormBase, a *C. elegans* database which has more detailed information on each gene), the fraction of cells in the selection expressing that gene, and a p-value to indicate how specific that gene's expression pattern is for the selection. P-values are calculated using a hypergeometric test on the number of selected cells expressing a particular gene versus the number of cells in the entire population expressing that gene and are corrected for multiple hypothesis testing using a Bonferroni correction.

Intestinal cells form a distinct cluster in the PCA relative to the more widely distributed neurons. Intestinal cells form part of the endoderm, which is thought to be the most evolutionarily ancient germ layer and initiates its differentiation program earlier than ectoderm or mesoderm [9]. The PCA plot clearly shows this trend as the intestinal cells form a distinct cluster quite early in the time series.

Heatmap of Gene Expression

On a tab next to the PCA plot, users can select to see a matrix of gene expression patterns (Figure 8). Every row in the heatmap is a gene, and every column is a cell. Like in the PCA plot, users can select colored highlights of interest, especially those exhibiting differential expression relative to other cells, and generate a gene report showing those genes that are most enriched among that cell selection.

Intestinal cells and neurons express many of the same genes, but there are some rows that are unique to each. This is seen by clicking on each tissue type’s highlights and comparing the resulting gene reports (Figure 9). In addition, the top-ranked genes in each selection include genes known to be important in intestinal and neural development, respectively. For example, *elt-2* is an essential gene in the intestine, and *cnd-1* is known to be an important regulator in motor neuron development.

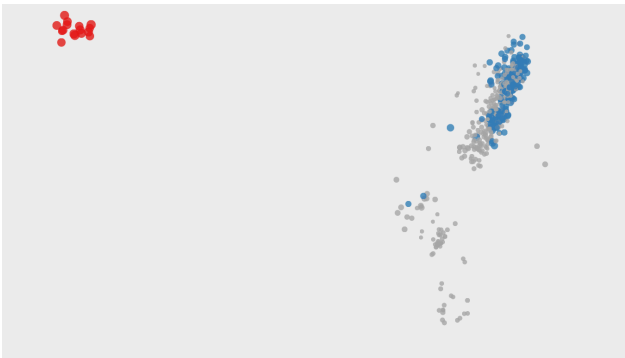


Figure 7. PCA plot of gene expression values. Intestinal cells (red) cluster differentially relative to neurons (blue).

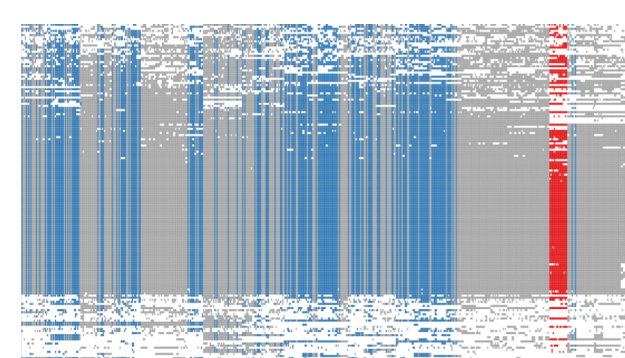


Figure 8. Heatmap of gene expression. Each row is a gene, each column is a cell. A filled square indicates expression. Intestines (red) and neurons (blue) highlighted.

INTESTINE					
Gene Name	Wormbase ID	Fraction Expressing	Fraction Expressing Selected	Gene Specificity p-val	
acp-5	WBGene00017427	0.02862862862862863	1	7.29030655689719e-29	
elt-2	WBGene00001260	0.02862862862862863	1	7.29030655689719e-29	
pdp-2	WBGene00003996	0.02862862862862863	1	7.29030655689719e-29	
F36A2.3	WBGene00009453	0.02862862862862863	1	7.29030655689719e-29	
T2B110.3	WBGene00012144	0.02862862862862863	1	7.29030655689719e-29	

NEURON					
Gene Name	Wormbase ID	Fraction Expressing	Fraction Expressing Selected	Gene Specificity p-val	
zif-11	WBGene00009839	0.7405405405405405	1	3.9496441754399846e-29	
cnd-1	WBGene00000561	0.21981981981981982	0.489247311827957	4.331506170125779e-24	
elt-3	WBGene00009899	0.818018018018018	1	1.2450491594926908e-18	
ZK185.1	WBGene00022681	0.7351351351351352	0.9354838709677419	5.087875361622902e-14	
isw-1	WBGene00002169	0.8306306306306306	0.989247311827957	9.401395967645181e-14	

Figure 9. Gene reports with top five enriched genes for intestine and neuron.

DISCUSSION

We envision EPICViz being used by the *C. elegans* community to explore gene expression patterns in cells, tissues, and lineages of interest. While computational analysis of a subset of the EPIC dataset has been performed (as in [4]), many potential users of this data lack the computational background to carry out such sophisticated analysis. In contrast, EPICViz can be used by anyone, regardless of computational experience. A user can select a cell, tissue, or lineage of interest, see where the selection lies in the developing embryo, and interrogate what genes are enriched in that population. This allows for straightforward identification of gene expression patterns, much simpler than watching the microscopy movie for each gene and trying to compare expression across samples.

In addition to identifying gene expression patterns among user selections, EPICViz can serve as a tool to generate hypotheses about genes and their functions. For those genes in the dataset that are poorly characterized, EPICViz can provide better data about where and when these genes are expressed during embryogenesis.

FUTURE WORK

We are currently working on ways to improve user selection, interface, and data management in EPICViz. With this version, users can select cell populations of interest by cell/tissue type and lineage to brush and link throughout all plots, but it is not possible to highlight based on a gene of interest. To facilitate gene-guided exploration of the data, we are planning to add gene names user selection menu drop-down so that users can specify a gene (or genes) of interest. In addition, at this point in time the screen is divided between all three major plots relatively equally; it is not possible for a user to resize a plot within the window. We would like to allow the user to expand or contract each element of the visualization so that he or she can focus on the plot(s) they find most informative. Finally, data loading and rendering performance is sub-optimal; we are working on building more efficient data structures that will quickly load and render to improve the performance of the visualization.

CONCLUSION

We have developed EPICViz, an interactive visualization of *C. elegans* embryogenesis that allows users to explore gene expression patterns in cells of interest. This tool facilitates exploration of development and gene regulation, including for users who lack the computational skills or resources to analyze confocal image data.

REFERENCES

- [1] J. E. Sulston, E. Schierenberg, J. G. White, and J. N. Thomson, "The embryonic cell lineage of the nematode *Caenorhabditis elegans*," *Dev. Biol.*, vol. 100, no. 1, pp. 64–119, Nov. 1983.
- [2] F. Long, H. Peng, X. Liu, S. K. Kim, and E. Myers, "A 3D digital atlas of *C. elegans* and its application to single-cell analyses," *Nat. Methods*, vol. 6, no. 9, pp. 667–672, Sep. 2009.
- [3] Z. Bao, J. I. Murray, T. Boyle, S. L. Ooi, M. J. Sandel, and R. H. Waterston, "Automated cell lineage tracing in *Caenorhabditis elegans*," *Proc. Natl. Acad. Sci. U. S. A.*, vol. 103, no. 8, pp. 2707–2712, Feb. 2006.
- [4] J. I. Murray, T. J. Boyle, E. Preston, D. Vafeados, B. Mericle, P. Weisdepp, Z. Zhao, Z. Bao, M. Boeck, and R. H. Waterston, "Multidimensional regulation of gene expression in the *C. elegans* embryo," *Genome Res.*, vol. 22, no. 7, pp. 1282–1294, Jul. 2012.
- [5] S. K. Andrey Palyanov, "Towards a virtual *C. elegans*: A framework for simulation and visualization of the neuromuscular system in a 3D physical environment," *In Silico Biol.*, vol. 11, pp. 137–147, 2012.
- [6] O. Ronneberger, K. Liu, M. Rath, D. Rueß, T. Mueller, H. Skibbe, B. Drayer, T. Schmidt, A. Filippi, R. Nitschke, T. Brox, H. Burkhardt, and W. Driever, "ViBE-Z: a framework for 3D virtual colocalization analysis in zebrafish larval brains," *Nat. Methods*, vol. 9, no. 7, pp. 735–742, Jul. 2012.
- [7] M. Meyer, T. Munzner, A. DePace, and H. Pfister, "MulteeSum: A Tool for Comparative Spatial and Temporal Gene Expression Data," *IEEE Trans. Vis. Comput. Graph.*, vol. 16, no. 6, pp. 908–917, Nov. 2010.
- [8] C. Pohl, "Left-right patterning in the *C. elegans* embryo," *Commun. Integr. Biol.*, vol. 4, no. 1, pp. 34–40, 2011.
- [9] T. Hashimshony, M. Feder, M. Levin, B. K. Hall, and I. Yanai, "Spatiotemporal transcriptomics reveals the evolutionary history of the endoderm germ layer," *Nature*, vol. 519, no. 7542, pp. 219–222, Mar. 2015.