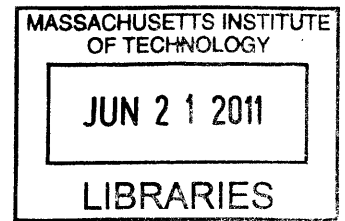


**Automation of Drosophila Gene Expression Pattern Image Annotation:
Development of Web-based Image Annotation Tool and Application of Machine
Learning Methods**

by

Anna Maria E. Ayuso

S.B., Electrical Engineering and Computer Science
Massachusetts Institute of Technology, 2007



ARCHIVES

Submitted to the Department of Electrical Engineering and Computer Science
in Partial Fulfillment of the Requirements for the Degree of

Master of Engineering in Electrical Engineering and Computer Science

at the
Massachusetts Institute of Technology

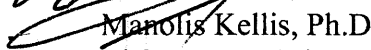
June 2011

©2011 Massachusetts Institute of Technology. All rights reserved.

Author: _____

Department of Electrical Engineering and Computer Science
May 20, 2011

Certified by: _____


Manolis Kellis, Ph.D
Associate Professor of Computer Science
Thesis Supervisor
May 20, 2011

Accepted by: _____

Dr. Christopher J. Terman Chairman, Masters of Engineering Thesis Committee

Automation of Drosophila Gene Expression Pattern Image Annotation: Development of Web-based Image Annotation Tool and Application of Machine Learning Methods

by

Anna Maria E. Ayuso

Submitted to the Department of Electrical Engineering and Computer Science
on May 20, 2011 in Partial Fulfillment of the Requirements for the Degree of
Master of Engineering in Electrical Engineering and Computer Science

ABSTRACT

Large-scale *in situ* hybridization screens are providing an abundance of spatio-temporal patterns of gene expression data that is valuable for understanding the mechanisms of gene regulation. *Drosophila* gene expression pattern images have been generated by the Berkeley *Drosophila* Genome Project (BDGP) for over 7,000 genes in over 90,000 digital images. These images are currently hand curated by field experts with developmental and anatomical terms based on the stained regions. These annotations enable the integration of spatial expression patterns with other genomic data sets that link regulators with their downstream targets. However, the manual curation has become a bottleneck in the process of analyzing the rapidly generated data therefore it is necessary to explore computational methods for the curation of gene expression pattern images. This thesis addresses improving the manual annotation process with a web-based image annotation tool and also enabling automation of the process using machine learning methods. First, a tool called LabelLife was developed to provide a systematic and flexible way of annotating images, groups of images, and shapes within images using terms from a controlled vocabulary. Second, machine learning methods for automatically predicting vocabulary terms for a given image based on image feature data were explored and implemented. The results of the applied machine learning methods are promising in terms of predictive ability, which has the potential to simplify and expedite the curation process hence increasing the rate that biologically significant data can be evaluated and new insights can be gained.

Thesis Supervisor: Manolis Kellis
Title: Assistant Professor of Computer Science

Acknowledgements

First and foremost, I profusely thank my thesis advisor **Professor Manolis Kellis**. It was around a year ago that I came out of left field from a job in the printing industry with no formal Biology background and Manolis agreed to take a chance on me and allow me the opportunity to work on such an exciting project. I am forever grateful. The Kellis lab has been a very fun and exciting place to work - full of many diverse people and interests. The melding of Computer Science and Biology is such a fascinating and quickly growing field and Manolis has been an enthusiastic and staunch advocate for its growth. There is no doubt many more revolutionary findings will come as a result of his enthusiasm, hard work, and incredible smarts. Manolis, thank you so much for the insights, enthusiasm, and guidance. I have been very privileged to work in your lab. Also, thank you for providing the financial support to make this project possible.

To **Chris Bristow** and **Charlie Frogner**, thank you both immensely for the guidance, support, and patience throughout this project. Chris, you were always there to help whenever I had any questions and I appreciate the energy you put into helping me understand background Biology concepts, past research in the field, Machine Learning methods, and so on. Charlie, you've also always been there for help and support. Your Machine Learning code was beautifully written and so helpful with this project and your exceptional project management skills really helped keep progress going and build out a LabelLife development team. I've learned so much from both of you over the past year. You both are very impressive people and I'm so appreciative that I had the honor to work with you. Thank you again.

To the **Berkeley Drosophila Genome Project (BDGP)**, this project wouldn't have been possible without data generated from your project.

To members of the **MIT Computational Biology Laboratory**, I feel honored to have been able to work in the same lab as all of you. You all are insanely smart people with many accomplishments under your belts. Thank you so much for all of the feedback on this project at the group lab meetings. All of it was so helpful in driving the direction of this project.

To the contributors of LabelLife and this project's motivation, **Stanislav Nikolov, Olivier Chatot, Bradley Wu, and Jennifer Wang**, it's been a pleasure working with all of you. Thanks for all of your effort and insights into further improving LabelLife. I've been so impressed with how quickly you all have been able to jump in and make immediate impact. It's been great working with all of you.

To the **LabelMe team**, specifically **Antonio Torralba** and **Bryan Russell**, thank you for writing an inspiring tool. I was privileged to have the opportunity to work with you both back in 2007 on LabelMe on an enjoyable Senior project. It's been incredibly fun the second time around working with the successful concept of a website that you all made and making it relevant for a Biological problem. Thanks for allowing us to extend your brain child for new and exciting purposes.

To my academic advisor, **Professor David Gifford**, thank you for helping me figure out my

Acknowledgements

career path. You met with me before I returned to school to offer guidance and support. I wasn't enrolled in school or an official advisee of yours, but you took time out of your busy schedule to help me and I am so grateful. I fondly remember the quick ChIP-sequencing lesson you gave me in your office before I went off to find a research group. Thank you again for everything.

To **Ed Keisling**, thank you for being my first real mentor. I count my lucky stars that I was able to start out my career at a place like Pegasystems with a boss like you. You were always encouraging and supportive, and you didn't seem to mind that I really didn't know what I was doing as a college sophomore. Thank you for taking the chance on me and encouraging my professional growth. I really look up to you and I'm so grateful for the opportunity I had to work with you and the stellar teams that you've run.

To my previous colleagues at **Vistaprint Inc**, thank you for developing me as a Software Engineer. I took the grueling Software Development classes at MIT, but nothing compared to immersion in the software systems at Vistaprint. (There's no better way to learn how to write robust code than to get calls at 3am from the production plants about software going horribly wrong.) It was truly an honor to work with such smart and creative people on fun problems. You all taught me so much about many different things and I can't thank you enough. You all not only helped me grow professionally, but also as an individual. It was sad to leave, but the encouragement and support was amazing and I not only miss you as colleagues, but as friends. I couldn't have asked for a better place to work right out of college. I'd like to thank you each as individuals, but that would take pages and pages. It suffices to say that you've all had a profound and lasting affect on me and I hold you all near and dear to my heart. Thank you again for an amazing experience and many learning opportunities.

To my former MIT field hockey and lacrosse coach, **Cheryl Silva**, thank you so much for everything. You helped develop me as an athlete and a person. I feel so privileged to have had a coach like you. Thank you for four great years of varsity athletics and thank you for allowing me the opportunity to help serve as an assistant coach this past year. Your words of wisdom "be part of something bigger than yourself" have stuck with me not only in athletics but in my career and everyday life. Thank you for everything. I look up to you in so many ways.

To my beloved cats, **Pooh** and **Moo Moo**. You both have turned me into a crazy cat lady. Thanks for offering your stress relief services and funny cat antics. I've also really enjoyed our long, late night conversations. (Kidding!) Thanks for holding on for me Pooh even though the odds have been against you with your brain tumors. You are the miracle cat (that is on a lot of miracle drugs). I couldn't have asked for better pets.

To **Karin**, who has been a second Mom to me. Thank you for the emotional support this past year. I relate to you on so many different levels and I appreciate how you've been there for me through the ups and downs of life. Just like you staked claim on Lorine, I told Brad that I staked claim on you. Thanks for everything.

To my boyfriend, life partner, little spoon, and man friend, **Brad**. I'm not so sure I would have had the guts to leave my cushy, engaging job in industry to pursue a personal interest if it weren't for your support and encouragement. It was the best decision ever and I'm so appreciative to you for

Acknowledgements

helping me get here. I feel very lucky to have you in my life. You've stood by my side in good and bad and I can't thank you enough. I love you very much. I can't wait for the next exciting chapter of life with you.

To my mother, **Kaye**. Thank you for everything Mom. You've always been there for me and I can't thank you enough. You're a great example of how to keep on fighting through adversity and I look up to you. Thank you for all of the sacrifices you made for me to get me where I am in life. I love you very much.

To my dearest little sister, **Kaity**. Even though you are my little sister, I truly look up to you in so many ways. You are an amazing (and sassy) woman and I am so lucky to have you for a sister and best friend. The way you carry yourself really inspires me. You are definitely the rock in our family. This has been a tremendous few years full of curve balls and you have been there for me every step of the way. Thank you so much for your emotional support and encouragement. You always know how to brighten my day and make me laugh. I really wouldn't be where I am without you. I love you so much Luscious.

Finally, to my courageous little brother, **Matt**, to you this thesis is dedicated. Matt, you are truly a very inspiring guy. Despite your physical and mental handicap, you go about life with such poise and optimism. Our family has truly been blessed that you are who you are and we've all learned so much from you. Over the years, you've taught me so much- compassion for others, resilience, facing adversity, and so on. You have inspired me to want to help others like you because you are such an amazing person. Your condition is why I've always been fascinated by genetics and it's because of you that I am where I am today. Thank you for teaching me so many life lessons little brother. You are a blessing. "Matt, I love you, call me."

Table of Contents

Table of Figures.....	10
Table of Tables.....	12
List of Acronyms/Definitions.....	13
1 Introduction.....	14
1.1 Motivation and Objective.....	14
1.2 Organization.....	16
2 Background and Previous Work.....	18
2.1 Drosophila melanogaster	18
2.2 The Berkeley Drosophila Genome Project.....	19
2.3 Automated Morphological Term Prediction For Gene Expression Image Pattern Data.....	21
2.4 LabelMe.....	30
2.5 Thesis Contribution	31
3 LabelLife: A Web-based Image Annotation Tool.....	32
3.1 Overview.....	32
3.2 Software Requirements.....	33
3.3 Design and Implementation.....	37
3.3.1 Data Model.....	39
3.3.2 System Architecture.....	41
3.3.3 User Interface.....	46
3.4 Use Cases.....	50
3.5 Ongoing Work.....	52
3.5.1 Integrating predictions into UI.....	52
3.5.2 UI Aesthetic Improvement.....	52
3.5.3 Data Download and Analysis Tools.....	53
3.5.4 Mechanical Turk.....	53
3.5.5 Search Functionality.....	54
3.5.6 Video.....	55
3.6 Other Uses For LabelLife.....	55
3.7 Summary.....	57
4 Machine Learning Methods For Morphological Term Prediction.....	58
4.1 Overview.....	58
4.2 Feature Extraction.....	59
4.3 Machine Learning Methods To Predict Morphological Terms For Gene Expression Pattern Images.....	61
4.3.1 Regularized Least Squares Regression (RLS).....	61
4.3.2 Support Vector Machine (SVM).....	63
4.3.3 Methods of Evaluating Results.....	64
4.4 Data Used For Prediction Problem	66
4.4.1 D. melanogaster Embryo Development Stage 4 To 6.....	66
4.4.2 D. melanogaster Embryo Development Stage 7 To 8.....	68
4.5 Experiments: Application of Machine Learning Methods To Predict Morphological Terms For Gene Expression Pattern Images	69

Table of Contents

4.5.1 D. melanogaster Embryo Development Stage 4 To 6.....	71
4.5.2 D. melanogaster Embryo Development Stage 7 To 8.....	72
4.6 Results.....	72
4.6.1 D. melanogaster Embryo Development Stage 4 To 6.....	73
4.6.2 D. melanogaster Embryo Development Stage 7 To 8.....	75
4.7 Discussion.....	77
4.7.1 Results.....	77
4.7.2 Future Work.....	78
5 Conclusion.....	80
6 Appendix.....	82
6.1 Controlled Vocabulary Terms.....	82
6.2 BDGP Data Breakdown.....	85
7 Bibliography	91

Table of Figures

Figure 1: Example data from the BDGP database for the <i>twi</i> gene in <i>Drosophila</i> [14].....	14
Figure 2: A <i>Drosophila</i> embryo stain pattern image generated from in situ hybridization.....	14
Figure 3: Adult <i>Drosophila melanogaster</i> [17].....	18
Figure 4: A <i>D.melanogaster</i> embryo expressing Gene CG13333 at different development stages.....	19
Figure 5: Development of midline glia through time [18].....	20
Figure 6: Common pipeline applied to the problem of predicting morphological terms for gene expression pattern images.....	21
Figure 7: Misclassified “Yolk Nuclei” gene expression pattern images.....	22
Figure 8: The way that the relationship of gene expression pattern image to controlled vocabulary term mapping is maintained in the BDGP process. Not all vocabulary terms apply to all images.....	23
Figure 9: LabelMe online annotation tool.....	31
Figure 10: LabelLife web-based annotation tool.....	33
Figure 11: Screenshot of LabelLife demonstrating the three different levels that controlled vocabulary terms can be applied to: ImageSets, Images, and ImagePolygons	35
Figure 12: Screenshot of LabelLife demonstrating the concept of ImageSets.....	35
Figure 13: Screenshot of LabelLife demonstrating the concept of Metadata.....	36
Figure 14: Screenshot of LabelLife demonstrating the concept of Controlled Vocabulary.	37
Figure 15: Screenshot of LabelLife labeling stage. One to many controlled vocabulary terms are required to complete the polygon.	37
Figure 16: Screenshot of LabelMe labeling stage. A free-form text label is required to complete the polygon.	37
Figure 17: LabelLife data model.....	40
Figure 18: Screenshot of LabelLife demonstrating where the data model objects apply in the user interface.	41
Figure 19: LabelLife Software Architecture	43
Figure 20: File system storage layout.....	46
Figure 21: LabelLife user interface.....	48
Figure 22: LabelLife UI Component: ImageSet Metadata Pane	48

Table of Figures

Figure 23: LabelLife UI Component: Image Metadata pane49

Figure 24: LabelLife UI Component: Main Image pane49

Figure 25: LabelLife UI Component: ImageSet pane50

Figure 26: LabelLife UI Component: VocabularyTerm pane50

Figure 27: Screenshot of LabelLife expanded ImageSet view. This is useful for viewing large ImageSets.51

Figure 28: Screenshot of LabelLife Image/Image Set Labeling View. The terms selected will be applied to the ImageSet as a whole. It can be seen that “dorsal ectoderm anlage in statu nascendi” and “ventral ectoderm anlage in statu nascendi” are already applied to this ImageSet.....52

Figure 29: Screenshot of LabelLife ImagePolygon Labeling View. The terms selected will be applied to the highlighted polygon in the Image.52

Figure 30: LabelLife Redesign. Image courtesy of Bradley Wu, an MIT undergraduate working on UI improvement.....53

Figure 31: A potential use for LabelLife could be labeling soft tissues in MRI brain scans.....56

Figure 32: A potential use for LabelLife could be labeling store fronts in GoogleMaps57

Figure 33: Feature Representation of gene expression pattern image. The embryo is converted to a binarized image indicating the presence/absence of a feature.....60

Figure 34: Example of an embryo gene expression pattern image run through image processing...61

Figure 35: Illustration of area under the curve.65

Figure 36: Zscore for all Stage 4 to 6 terms found in the data set sorted in descending order of frequency. Red color indicates positive feature detection and blue color indicates negative feature detection across all relevant embryo images for a term.....68

Figure 37: Effect of cleaning labels with LabelLife. The first row corresponds to the Zscore profile of the terms before cleaning. The second row corresponds to the Zscore profile of the same terms after cleaning. It can be seen that cleaning the data through LabelLife produces more striking features for classification.72

Table of Tables

Table 1: Stage 4 To 6 Morphological Term Breakdown.....	67
Table 2: Stage 7 To 8 Morphological Term Breakdown.....	69
Table 3: Change in number of positive labels before and after curating data in LabelLife.....	71
Table 4: AUC Per Gene, Stage 4-6, Pixel feature data, RLS classification w/Linear Kernel.....	73
Table 5: AUC Per Image, Stage 4-6, Pixel feature data, RLS classification w/Linear Kernel.....	74
Table 6: AUC Per Image, Stage 4-6, Pixel feature data, RLS classification w/Gaussian Kernel.....	75
Table 7: AUC Per Image, Stage 4-6, Pixel feature data before and after cleaning in LabelLife, RLS classification w/Linear Kernel.....	75
Table 8: AUC Per Image, Stage 4-6, Pixel feature data before and after cleaning in LabelLife, RLS classification w/Gaussian Kernel.....	75
Table 9: AUC Per Image, Stage 7-8, Pixel feature data, RLS classification w/Linear Kernel.....	76
Table 10: AUC Per Gene, Stage 7-8, Pixel feature data, RLS classification w/Linear Kernel.....	76
Table 11: AUC Per Image, Stage 7-8, Pixel feature data, RLS classification w/Gaussian Kernel.....	77
Table 12: AUC Per Gene, Stage 7-8, Pixel feature data, RLS classification w/Gaussian Kernel.....	77

List of Acronyms/Definitions

AUC	Area Under the Curve
BDGP	Berkeley Drosophila Genome Project
FPR	False Positive Rate
HTML	Hypertext Markup Language
ISH	<i>In situ</i> hybridization
LOOCV	Leave One Out Cross Validation
ML	Machine Learning
RLS	Regularized Least Squares
ROC	Receiver Operating Characteristic
SQL	Structured Query Language
SVM	Support Vector Machine
TPR	True Positive Rate
UI	User Interface
XML	Extensible Markup Language

1 Introduction

1.1 Motivation and Objective

Precise control of spatial and temporal patterns of gene expression is required for tissue and organ development. Genes expressed in similar patterns are often co-regulated and grouping genes based on shared expression patterns has been invaluable in determining gene regulatory mechanisms. Large-scale biological methods for evaluating gene expression are providing an abundance of valuable information to researchers. One technique for profiling spatial patterns of gene expression is RNA *in situ* hybridization (ISH) in which a labeled complementary nucleotide probe is hybridized to a fixed tissue sample and visualized using microscopy (Figure 1). This technique generates valuable spatial information and temporal staging can be inferred by the tissue morphology. Large-scale ISH screens have been carried out in a range of model organisms allowing for the profiling of thousands of genes. These screens generate tens of thousands of images making computational methods a necessity to quickly process and understand the results.

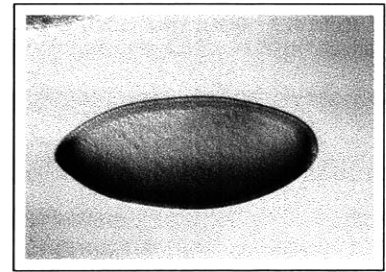


Figure 1: A *Drosophila* embryo stain pattern image generated from *in situ* hybridization

The Berkeley *Drosophila* Genome Project (BDGP) is an effort to provide gene expression data for all *Drosophila melanogaster* (fruit fly) genes across embryonic development to researchers. Currently the group has examined 7,153 genes that are documented with 97,882 digital images and the data is continuously growing [14]. Information is summarized using manual annotations completed by field experts that link images to controlled vocabularies of anatomical and developmental ontology terms (Figure 2) and this allows for the integration of spatial patterns with additional genomic datasets.

1 Introduction

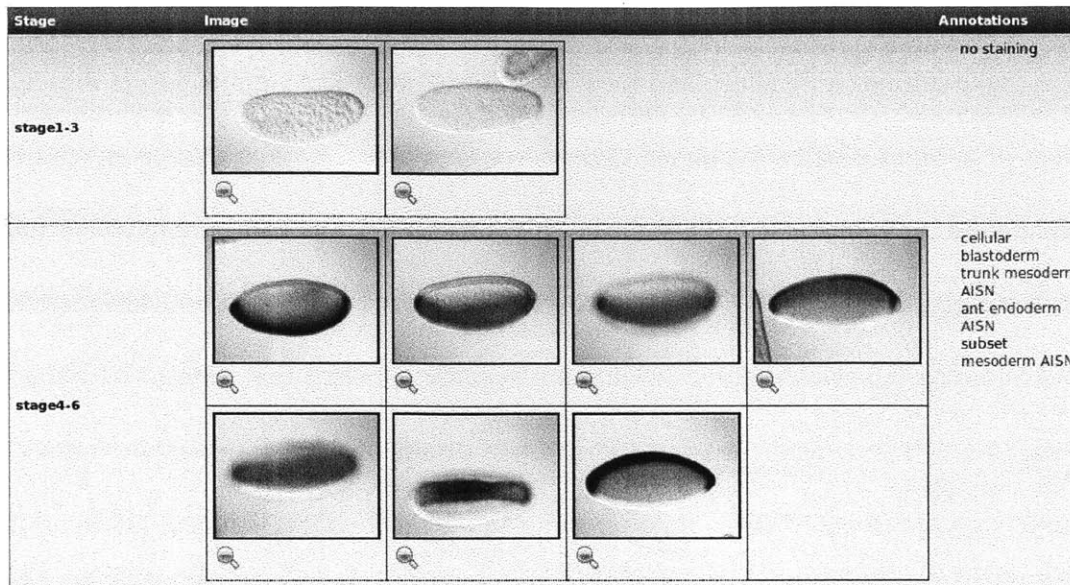


Figure 2: Example data from the BDGP database for the *twi* gene in *Drosophila* [14].

However manual curation is labor intensive and prone to errors. Several research groups have used this data and applied computational methods in an effort to predict anatomical and developmental terms based on the image data [1-13]. There are common challenges that each group has run into. One major challenge that affects the performance of term prediction algorithms is the fact that annotators currently apply multiple terms to images or sets of images instead of image features [2,4]. Term prediction algorithms have to compensate for loss of information in what feature in an image an annotator intended a term to apply to since terms are assigned in groups to groups of images. While some groups have had success in using image processing and machine learning methods to predict terms, even greater success is possible by enhancing the available data with a linkage between terms and their location within images. An existing web-based image annotation tool called LabelMe provides a foundation for accomplishing such goals.

LabelMe is a tool developed by the MIT Computer Science and Artificial Intelligence Laboratory (CSAIL) that is used by the computer vision community for annotating objects in real

1 Introduction

world photographs of various scenes [15,16]. The data can then be used in object recognition algorithms that take into account shape and location information. This application is directly relevant to the problem of assigning terms to stained gene expression patterns. The same concept of outlining stained regions with polygons and linking those polygons with terms can be applied to enhance the results of term prediction algorithms. However, the current system is not designed to manage large sets of biological images and there are domain specific concepts that need to be integrated. For example, in biological experiments, it is critical to track details of what reagents were used and when the experiment was carried out as there are often substantial batch effects that can confound the results. Furthermore, controlled vocabulary terms are essential to standardize the results and enable comparison across genes.

This thesis combines the previously mentioned concepts to address the meta-problem of analyzing a rapidly growing gene expression pattern data set to infer insights into gene regulation. The first objective of this thesis is the development of a web-based image annotation tool called LabelLife to provide a systematic and flexible software solution for annotating images, groups of images, and shapes within images using terms from a controlled vocabulary. The second objective of this thesis is the application of machine learning methods for automatically predicting vocabulary terms for a given image based on image feature data.

1.2 Organization

This thesis is organized into several sections. Chapter 2 provides background information on relevant concepts that are referred to in this thesis. This section also surveys existing literature on the problem of automatically predicting terms for gene expression pattern images and examines LabelMe, from which the web-based annotation tool written for this thesis draws much inspiration. Chapter 3

1 Introduction

introduces the first contribution of this thesis, LabelLife. The design and motivation behind the new web-based annotation tool and its relevance to the problem of annotation of gene expression pattern images is described in detail. Chapter 4 describes the machine learning methods used to predict controlled vocabulary terms for gene expression pattern image data. Background information on the methods is provided and the details of the experiments are laid out. Finally, results are presented and analyzed. Chapter 5 wraps up the two main contributions of this thesis with a discussion and conclusion.

2 Background and Previous Work

This section provides an overview on concepts touched upon by the work of this thesis. The following subsections provide more detail on *Drosophila melanogaster*, the Berkeley Drosophila Genome Project, previous work in automated annotation term prediction, and LabelMe, a web-based image annotation tool.

2.1 *Drosophila melanogaster*

The *Drosophila melanogaster* (Figure 3), a species of the fruit fly genus, is one of the most widely studied organisms in genetics. The focus on this species started in 1909 when Thomas Hunt Morgan of Columbia University began using the fruit fly to study the laws of Mendelian Inheritance [22]. The fruit fly was chosen for its fast reproduction rate and small size. Many insights were gained from Morgan's experiments about genetics by crossing flies of known phenotypes. It was not until 2000 that the whole genome of the *Drosophila melanogaster* was sequenced [28]. The 100 or so years

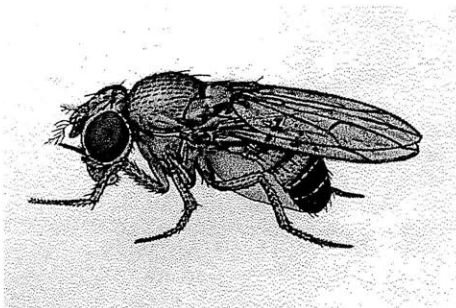


Figure 3: Adult *Drosophila melanogaster* [17]

in between have brought many discoveries as to the underlying biological mechanisms involved in genetics. Even though we know so much more today about the role that genetics plays in our lives, fruit flies are still widely studied because of the similarity of their DNA to humans [24] and also the speed at which experiments can be carried out.

Another aspect of fruit flies that is widely studied today is embryogenesis. Embryogenesis is the process of formation and development of an embryo. Tracking the development of the fly as an embryo allows for insights into the process of gene regulation. Modern science has conquered the

2 Background and Previous Work

hurdle of determining what DNA is present in a genome. Now, researchers are trying to better understand the complex signaling system that drives expression of this DNA. Applying *in situ* hybridization to known genes in fruit flies has allowed for the visualization of localized gene expression. Repeating this process at different developmental time points on a gene-by-gene basis yields spatio-temporal information that helps paint a more complete picture of gene expression as time progresses. Figure 4 illustrates this concept for a *D. melanogaster* embryo expressing the CG13333 gene across five development stages. The localization of the expression changes as embryogenesis progresses. One can also observe the detail introduced as the embryo develops.

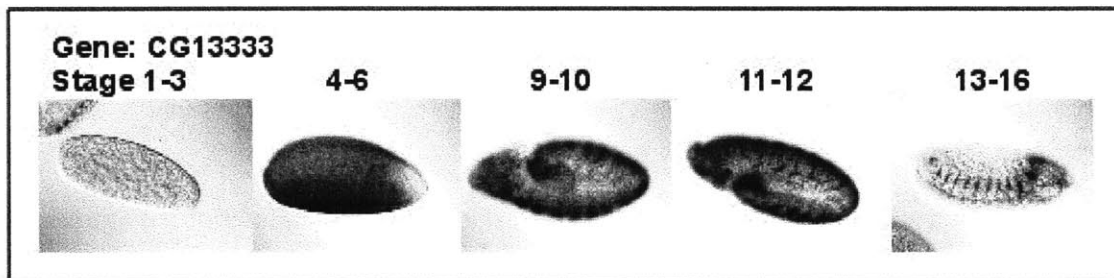


Figure 4: A *D.melanogaster* embryo expressing Gene CG13333 at different development stages.

2.2 The Berkeley Drosophila Genome Project

The Berkeley Drosophila Genome Project (BDGP) is an effort to provide a thorough and high quality collection of *Drosophila melanogaster* genomic data. One output of this effort is an abundance *in situ* hybridization gene expression pattern data for all *Drosophila melanogaster* (fruit fly) genes across embryonic development.

To generate this data, labeled RNA probes for each gene are hybridized to embryo tissue through *in situ* hybridization. *In situ* hybridization is the process of hybridizing a complementary

2 Background and Previous Work

nucleotide probe to a specific DNA or RNA sequence of interest in a tissue and visualizing the result with microscopy. The resultant stained embryo undergoes a quality control process. If the process goes smoothly, image acquisition is carried out to capture the result. Multiple high-resolution microscopic digital photographs are taken to capture several focal planes of the 3D embryo. The pictures are taken to capture the lateral, dorsal, and ventral views of the embryo with the anterior end aligned to the left and the posterior end aligned to the right. This process is repeated for all combinations of genes and development stages [13]. Embryogenesis is divided into set of development stages characterized by expected morphological markers. The BDGP subdivided these stages into 6 stage ranges (1-3, 4-6, 7-8, 9-10, 11-12, 13-16) that account for stage transitions with the largest noticeable affect on gene expression. It can be noted from the data that more intricate, complex stain patterns appear towards the later stages.

Once the images are settled upon for a gene/stage range combination, a human annotator manually assigns terms from a controlled vocabulary to indicate anatomical and developmental intermediaries present in the stain

pattern of the embryo. The controlled vocabulary was designed to account for the fact that embryonic structures might not be completely developed at a certain stage range with terms that indicate a development path. The terms “anlage in statu nascendi,”

anlage in statu nascendi ~ *anlage* ~ *primordium* ~ *differentiated organ*

Tracing back single-minded

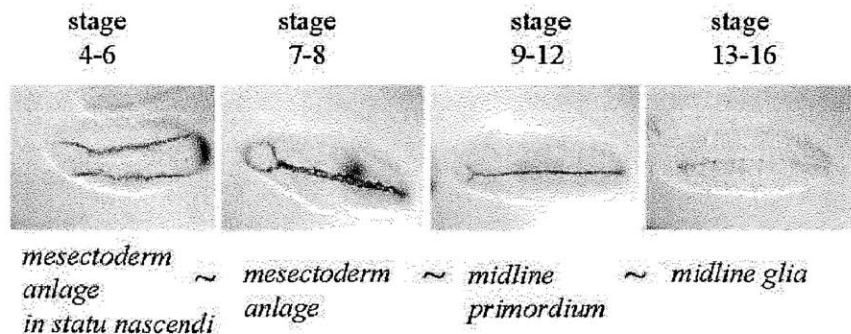


Figure 5: Development of midline glia through time [18]

“anlage,” “primordium,” and “organ” are appended to the fully developed term name to indicate

2 Background and Previous Work

development stage. Figure 5 demonstrates the development of the midline glia through time. One can observe the terms that indicate developmental intermediaries in the early stages. There is a hierarchal structure to the terms that form a directed acyclic graph. See Appendix 6.1 for the currently used terms. These terms are assigned as a group to the whole group of images by the expert curator.

Currently there are only a handful of expert human annotators working on the set. The project has processed over 7,000 of the known fly genes. The BDGP had two large spikes in the release of annotation data in 2007 and 2010. Since these spikes, the group aims to consistently deliver data for around 90 genes per month. The BDGP has implemented database and software changes to facilitate this flow of data to the public [14].

The BDGP makes all of its data publicly available through their website and they provide a web-based interface for searching existing image data and annotations. The website can be found at <http://www.fruitfly.org/index.html>.

2.3 Automated Morphological Term Prediction For Gene Expression Image Pattern Data

Various research groups have made attempts at automatically predicting annotation terms from gene expression pattern data generated by BDGP using computational methods. There is a common data analysis pipeline that is applied to this problem as illustrated in Figure 6. First, image processing is performed on the raw gene expression pattern

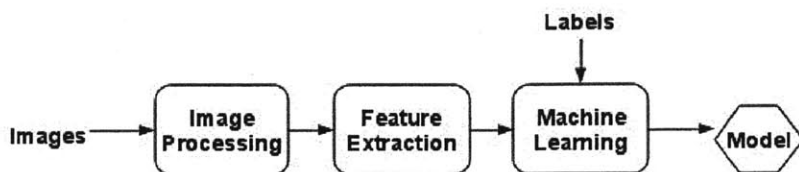


Figure 6: Common pipeline applied to the problem of predicting morphological terms for gene expression pattern images

2 Background and Previous Work

images to register embryos and threshold for stain intensity. Second, a method of feature extraction is used to generate data input for machine learning methods. Third, a machine learning method is performed on the input data with various parameters set to generate a predictive model. Fourth, unclassified images are run through the model and predictive ability is assessed using various metrics.

There are several common difficulties and considerations in carrying out this task that have been collectively noted by researchers who have worked on this problem before [1-13].

- **Groups of terms are assigned to groups of images**

This is the top most cited issue that has made automatic prediction of terms a difficult problem. Expert curators go through the images manually and combine their knowledge with what they see in the expression pattern images, microarray expression data, and published research to accurately assign one to many vocabulary terms to a set of images representing one development stage for a gene. One of those terms might have been meant for a pattern that shows up visibly well in one image out of that group but not the rest. The way the data is maintained, this information is lost and the term is associated to the whole set of images. Figure 7 illustrates this example for the Yolk Nuclei morphological term. The “good example” on the left shows how a Yolk Nuclei manifests in a gene expression pattern image. The “bad examples” on the right show other images which incorrectly receive the label of having a Yolk

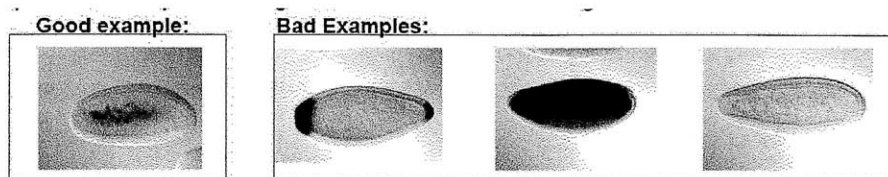


Figure 7: Misclassified “Yolk Nuclei” gene expression pattern images

Nuclei. When the terms are treated as ground truth

labels in machine learning methods, one has to be careful to determine which terms are actually

2 Background and Previous Work

relevant. Figure 8 further illustrates this data relationship.

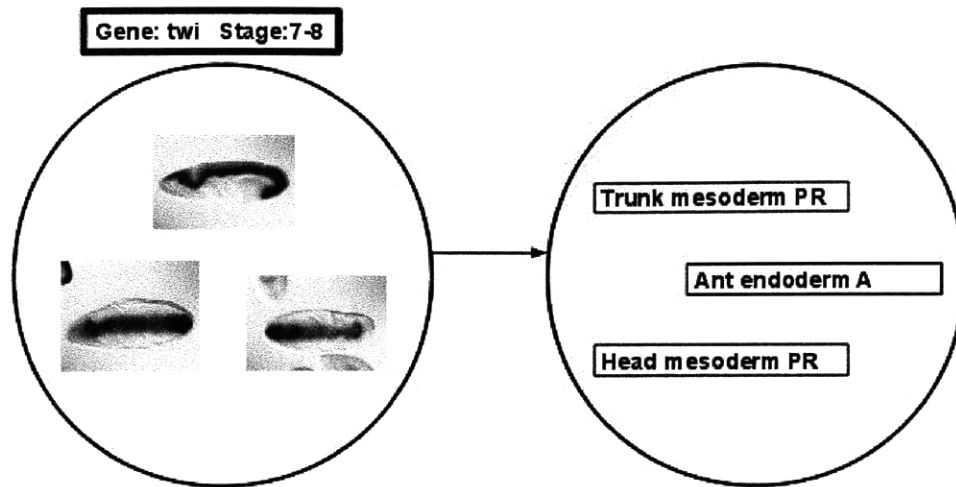


Figure 8: The way that the relationship of gene expression pattern image to controlled vocabulary term mapping is maintained in the BDGP process. Not all vocabulary terms apply to all images.

- **Multiple images per embryo/gene**

There are multiple images generated per gene/development stage. This is due to that fact that embryos are 3D objects and multiple pictures are necessary to capture enough focal planes to paint a complete picture. Lateral, dorsal, and ventral views are taken with the anterior end aligned to the left. An issue that arises from this aspect of the data is that some embryonic structures might be visible from one view but not another. Machine learning methods need to be able to account for these different views somehow and compare apples to apples.

- **Multiple developmental stages**

Similar to the problem of having multiple images per gene/development stage combination, there are multiple development stages. When looking at predicting terms for one gene, there will be several developmental stages to take into account. These need to be treated carefully because the terms in one stage generally do not show up in other stages. The controlled

2 Background and Previous Work

vocabulary was carefully crafted to account for developmental intermediaries of embryonic structures. Also, as the embryo develops, it becomes more defined and later stage embryos look vastly different than early stage embryos.

- **Embryos develop non-deterministically**

This is lesser of a problem but still worth mentioning. *Drosophila* embryogenesis is broken down into discrete development stage ranges where certain morphological markers are expected. While most embryos roughly adhere to this expected process, embryogenesis is a non-deterministic process and there is a chance that certain structures might appear sooner or later than expected given the stage range.

- **Image data has high dimensionality**

When carrying out machine learning to predict terms, it is necessary to generate feature data as an input to the algorithm. Using the raw image data itself yields feature data that has high dimensionality where the number of dimensions equals the number of pixels. Certain machine learning methods do not perform well with feature data that has this quality and so care must be taken to choose the right method for the input feature data. Attempts can be made to reduce the dimensionality of the data through a transformed representation, but care must be taken to ensure that classification still yields biologically significant results.

- **Term data is sparse**

Given all of the genes available in the database, controlled vocabulary terms are assigned at low frequencies as compared to the whole data set. This yields sparse term associations and data skew. Whether a term is assigned or not is treated as the label in machine learning methods. Having ground truth data that is sparse could yield less-than-ideal results if this fact is not taken

2 Background and Previous Work

into account. Often adding weights to the term classes to offset the skew in data is the most common way of remedying the problem.

Various groups have been working on developing automated pipelines for classifying gene expression pattern images since the early 2000s. There are a number of different feature representations as well as a number of different machine learning methods that could be applied to this problem. The groups that have investigated this topic have investigated a particular combination or a few combinations to evaluate predictive ability.

Choosing the best feature representation was a focus of the early research related to this problem. The focus of the problem evolved into predicting labels for gene expression pattern images after image search based on well-defined features was understood. Kumar *et al.* (2002) [7] began this effort in 2002. This group built a basic expression search tool called BEST to retrieve the best matching gene expression pattern images given an input gene expression pattern image. The group executed their software against 982 gene expression pattern images that were extracted from existing literature, since the BDGP was still in its infancy and these images were not readily available in 2002. The images were standardized by converting to a common dimension and maximizing the enclosed embryo to the boundaries. Stain patterns were extracted using an adaptive thresholding algorithm and the resulting images were stored as binary data where 1 indicated the presence of a feature and 0 indicated the absence of a feature. A distance metric was used to compare input gene expression patterns against ones stored in the database. The goal was to minimize the number of bit differences divided by the number of bits indicating features in one of the images. The result of a search displayed a list of images from the database and the percentage match achieved. This was the early work that inspired other groups to work towards automation and it laid out necessary preprocessing steps to

2 Background and Previous Work

achieve this goal. This work also provided researchers a way to visually search new expression patterns quickly.

Gurunathan *et al* (2004) [6] continued research on the best feature representation of gene expression pattern image data by comparing the Binary Feature Vector (BVF) and Invariant Moment Vector (IMV) feature representations. BVFs represent features as 1's and 0's where a 1 indicates that a feature is present. This is how images were represented in the BEST tool. IMV was tested as a viable substitute because this method was proven to work well in compensating for variations in scale, translation, and rotation. This group used 262 gene expression pattern images extracted from research papers. All of the embryos were between stages 1-8 and laterally aligned. The results indicated that while the IMV feature representation produced faster retrieval speeds, BVF still provided the most biologically meaningful matches.

The quest for the best feature representation was continued 5 years later in 2009 by Frise *et al* (2009) [9]. This group created a new feature representation of the *Drosophila* embryo in the form of an elliptical mesh grid. The group focused on stage range 4-6. They used this new representation to cluster images with similar expression patterns. The clusters that were learned were determined to be biologically relevant based on previous research.

Along the evolution of using the BDGP data for predicting morphological terms for gene expression pattern images, several groups used the data to learn different things, such as gene coexpression and development stage. In 2006, Ye *et al* [5] attempted to go beyond searching on pattern similarity and set up a system to classify which development stage range a gene expression pattern image belonged to. The group used Gabor filters and Histogram Equalization to extract features from pixel neighborhoods in the images. These features were then run through an implementation of Linear Discriminant Analysis (LDA) to extract the most relevant features and thus

2 Background and Previous Work

reduce the dimensionality of the data. Finally, these most relevant features were classified using both the Nearest Neighbor (NN) clustering algorithm and Support Vector Machines (SVM). The group used 2705 expression pattern images from the early stages. The classification was run varying the training sample size from 30 to 540. The results indicated that Nearest Neighbors and SVM were comparable in classifying developmental stage range for features generated from LDA. This group achieved accuracies that ranged from 76.55% to 88.91%.

Samsonova *et al* (2007) [8] focused on the problem of understanding gene coexpression by integrating non-homogeneous data sets. This group used both microarray and gene expression pattern image data in an attempt at predicting the spatial localization of gene expression. The group developed the concept of a functional unit and manually assigned genes that were associated with a functional unit of development. The problem then became predicting which genes should be associated with a functional unit. The group developed 15 functional groups and each group would have on average a few dozen genes associated. Using SVMs for classification, this group achieved accuracies as high as 80% for some functional units. This method required extensive biological insight into the data in order to transform the data to a new representation. Another group that investigated coexpression was Costa *et al* (2007) [12]. Costa *et al* investigated the problem of automatically predicting clusters of similarly expressed genes using semi-supervised machine learning methods. The raw gene expression pattern image data was used to generate pairwise constraints used in a Gaussian mixture model clustering scheme. The clusters found using the automated methods proposed by this group were determined to yield biologically meaningful clusters of related genes.

The following groups dedicated their research efforts to the problem that is laid out in this thesis. Zhou and Peng (2007) [4] were one of the early groups dedicated to the problem of automatically predicting morphological terms for gene expression pattern image. They developed a

2 Background and Previous Work

software system to carry out this recognition and classification task. First, the system would recognize which stage range the image belonged to and second, it would assign the appropriate morphological terms to that image based upon the stain pattern. Image features were extracted with a multi-resolution 2D discrete transform, followed by min-redundancy max-relevancy feature selection. Classification was carried up by implementing multiple parallel binary classification problems corresponding to each morphological term. Linear Discriminant Analysis (LDA) was used to implement the binary annotation classifier. To test the system, the group constructed two sets of features for each stage range, one containing two mutually exclusively assigned terms, and one containing two mutually inclusively assigned terms. This training setup was used to predict development stage range. The group then trained the parallel, binary classifiers on a subset of frequently occurring terms and generating probabilistic confidence scores. This approach demonstrated promising prediction results.

Peng *et al* (2007) [3] also contributed an automated pipeline for the problem of predicting morphological terms for gene expression pattern images. The group investigated the combination of three different feature representations; Gaussian Mixture Models (GMM), Principal Component Analysis (PCA), and wavelet functions. The goal of using all three representations was to adequately capture both local and global feature properties. The group also investigated three different classification methods; Linear Discriminant Analysis (LDA), Support Vector Machine (SVM), and Quadratic Discriminant Analysis (QDA). The results indicated that LDA consistently performed well across a small training class versus a large one. This was the most impressive result since a relatively good predictive model was achievable by using minimal data. SVM and QDA both performed well on a large training class, but not very well on a small training class. This result indicated overfitting of data.

2 Background and Previous Work

Shuiwang *et al* (2008) [1] also created a pipeline for the automated annotation of gene expression patterns using controlled vocabularies. The group extracted invariant features from individual images and constructed pyramid match kernels between sets of images to measure similarity in features. The specific features used were SIFT, shape context, PCA-SIFT, spin image, steerable filters, differential invariants, complex filters, moment invariants, and cross correlation. The kernels were combined using hypergraph spectral learning to learn which features were the most important and the output of this step was used in the classification problem. Testing data on this framework involved evaluating micro F1, macro F1, precision and recall against combinations of 1000, 1500, and 2000 sets versus 10, 20, 30, 40, 50, and 60 controlled vocabulary terms at a time. Multiple kernel combination methods were evaluated. As the number of terms being learned increased, the evaluation metrics decreased. In terms of macro/micro F1 scores, kernels integrated by either star or clique expansions achieved the highest performance up to around .60. In terms of precision and recall, SVM and Uniform achieved the highest performance up to around .67.

Shuiwang *et al* (2009) [2] used a bag-of-words approach for the problem of automating annotation of gene expression pattern images. This methodology was used to exploit the relationship between terms. To generate this representation, a visual codebook for each view – lateral, dorsal, and ventral - was generated. Each image was compared to this code book and the present features are quantized. For images in a set, the codebook-determined features were considered the bag-of-words. The classification stage made use of a binary classification scheme with a low-dimensional shared subspace to capture relationships between controlled vocabulary terms. For evaluation, this classification approach is compared against SVM and Pyramid Match Kernels (PMK) for classifying 10, 20, and 30 terms at a time. The AUC achieved by the shared-subspace multi-label formulation surpassed the other methods at .81 on stage range 4-6. This method also outperformed the other stage

2 Background and Previous Work

ranges. This work offered the most promising success rate in automatically predicting controlled vocabulary terms for gene expression pattern images.

2.4 LabelMe

LabelMe is a web-based image annotation tool developed in the MIT Computer Science and Artificial Intelligence Laboratory (CSAIL), specifically in the Computer Vision group. The tool provides an easy-to-use interface for hand-labeling objects in scenes. Figure 9 shows a screenshot of the website. The website is most commonly used to label real world scenes, but the tool could be applied to labeling any image.

The user simply points and clicks in their browser to generate polygon points around the object of interest. When the user completes the polygon, a prompt asks for a free form label to describe the polygon. This tool is widely accepted and used by the computer vision community with



Figure 9: LabelMe online annotation tool.

over 700,000 objects [15] labeled and counting in the main instance running at labelme.csail.mit.edu.

Full support for the tool has been built out in several ways. First, there is a Matlab toolbox to process and analyze the data. Second, the source code is made freely available with detailed installation instructions. Third, instructions are provided on how to integrate with Amazon's Mechanical Turk to outsource the labeling process to the public. Fourth, LabelMe video has been

2 Background and Previous Work

recently introduced to also support annotating video clips.

Data generated from LabelMe has been used in computer vision algorithms by the MIT Computer Vision group to understand scenes. Analysis has been performed to understand where in a scene an object class typically lies with a certain probability. The possibilities for analysis are endless and this tool can be used in many different ways. It is well known across the Computer Science community at the MIT and several research groups have used separate instances of LabelMe to carry out various research goals. LabelLife, the web-based image annotation tool developed as part of this thesis takes much inspiration from the vision of LabelMe.

2.5 Thesis Contribution

The overall goal of this thesis was to contribute to the problem improving annotation of gene expression pattern images for *D. melanogaster* embryos through tools and automated machine learning methods. The first contribution was the development of a web-based image annotation tool called LabelLife which builds on the inspiration of LabelMe to provide an intuitive and easy-to-use system for annotating gene expression pattern images. The second contribution was the application of machine learning methods to automatically predict morphological terms associated to gene expression pattern images. The first contribution, LabelLife, is discussed in detail in the next section.

3 LabelLife: A Web-based Image Annotation Tool

3.1 Overview

The first contribution of this thesis was the development of a web-based image annotation tool called LabelLife. Figure 10 shows a screenshot of the tool. This tool was built to address the current way that information is maintained for gene expression pattern term assignment and provide an efficient and easy-to-use interface for introducing new data linkages. As mentioned previously, the

BDGP has annotated thousands of fruit fly genes with controlled vocabulary terms that indicate anatomical and developmental intermediaries. These genes are labeled as groups of images

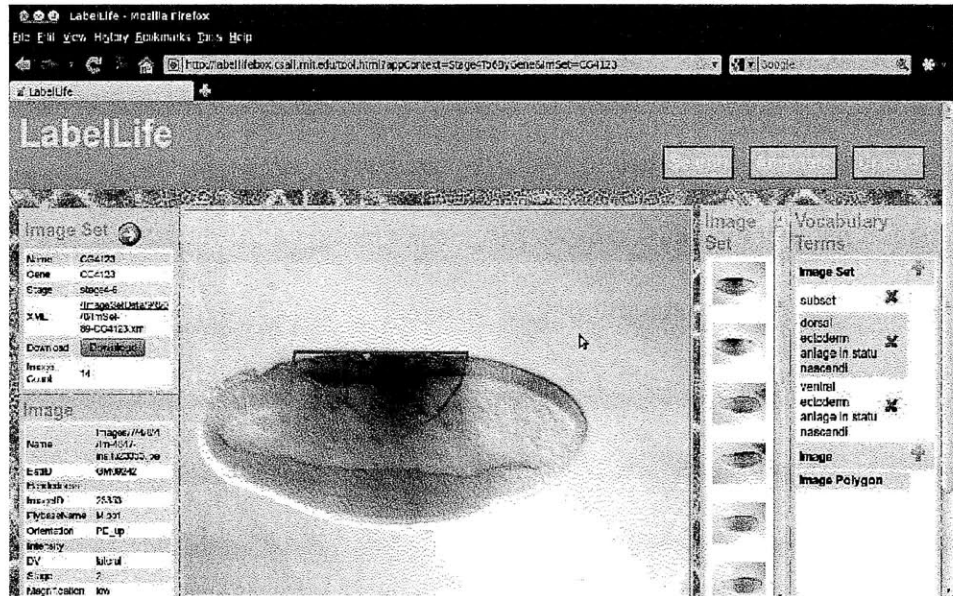


Figure 10: LabelLife web-based annotation tool

corresponding to the

gene and development stage range. One to many terms can be applied to this group of images. Granularity of detail is lost with this labeling scheme since a term that might only apply to a single image is forced to apply to the set of images. LabelLife was designed to allow for more precision in labeling and also improve the process by making it fast and efficient.

3.2 Software Requirements

There were several requirements that the LabelLife tool needed to satisfy:

- **Easy-to-use and intuitive interface**

There was already a tool for bulk annotating gene expression pattern images with controlled vocabulary terms that the BDGP uses. In order to create a tool that people would readily adopt, it was necessary to create something very easy-to-use and intuitive. The tool needed to be easy for a new user to install or access and obvious how to proceed with a task. All of the information needed to be displayed in a simple and uncluttered manner.

- **Open to the public**

LabelLife was meant to be embraced by the academic and research communities. In order to facilitate such adoption, it was necessary that the tool and data were both very open to all interested parties. The code needed to be open source and freely available. There also needed to be supporting tools for downloading and analyzing the data.

- **Ability to apply labels to shapes within images, images themselves, and group of images**

One major limitation of the current method of annotating gene expression pattern data is that detail is lost in what image or region within an image a term is meant to apply to. The tool needed to allow that detail to be captured somehow. The user needed to be able to outline a shape within an image. It was also important to capture relationships of labels at the image and image set levels as well. Figure 11 demonstrates the ability to apply controlled vocabulary labels to image polygons, images, and image sets in LabelLife.

3 LabelLife: A Web-based Image Annotation Tool

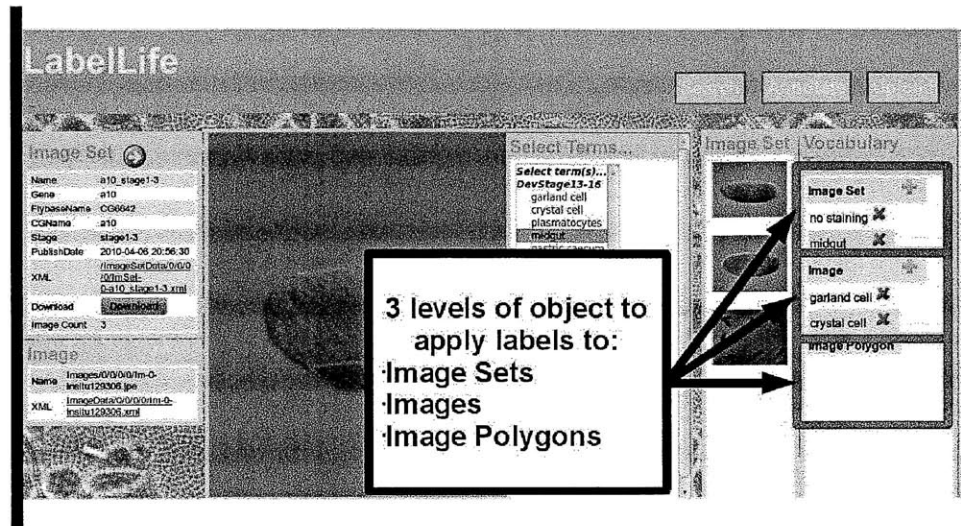


Figure 11: Screenshot of LabelLife demonstrating the three different levels that controlled vocabulary terms can be applied to: ImageSets, Images, and ImagePolygons

- **Image Sets**

For the problem of annotation gene expression pattern images, it was necessary to view images in sets grouped by gene and development stage. This is how they are currently labeled with terms so being able to visualize that data relationship was important. (Figure 12)

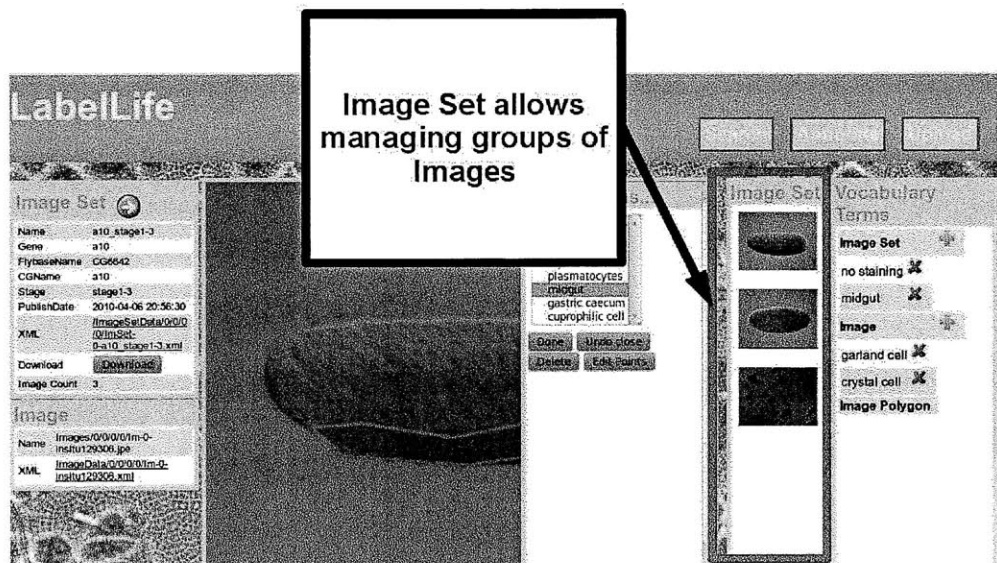


Figure 12: Screenshot of LabelLife demonstrating the concept of ImageSets

3 LabelLife: A Web-based Image Annotation Tool

- **Metadata**

There is a large amount of accompanying data that is important for manual curators to have easy access to when annotating gene expression pattern images. For example, it is important to know what focal plane of the embryo is being viewed – lateral, dorsal, or ventral. It might also be important to know when the experiment was done that generated the images or where on a specific chromosome the gene in question lies. There are many various pieces of information linked to the images that needed to be incorporated into the display. (Figure 13)

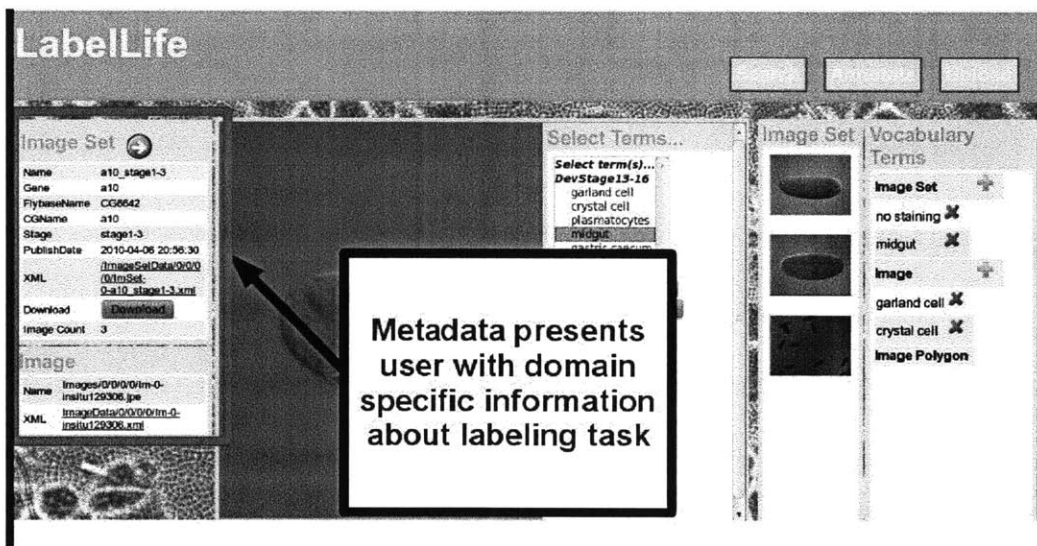


Figure 13: Screenshot of LabelLife demonstrating the concept of Metadata.

- **Controlled Vocabularies for labels**

When applying a label to a polygon, image, or set of images, the user needed to be restricted to a predefined set of vocabulary terms to ensure consistent labeling. Figure 14 demonstrates how users are restricted to a controlled vocabulary in LabelLife.

3 LabelLife: A Web-based Image Annotation Tool

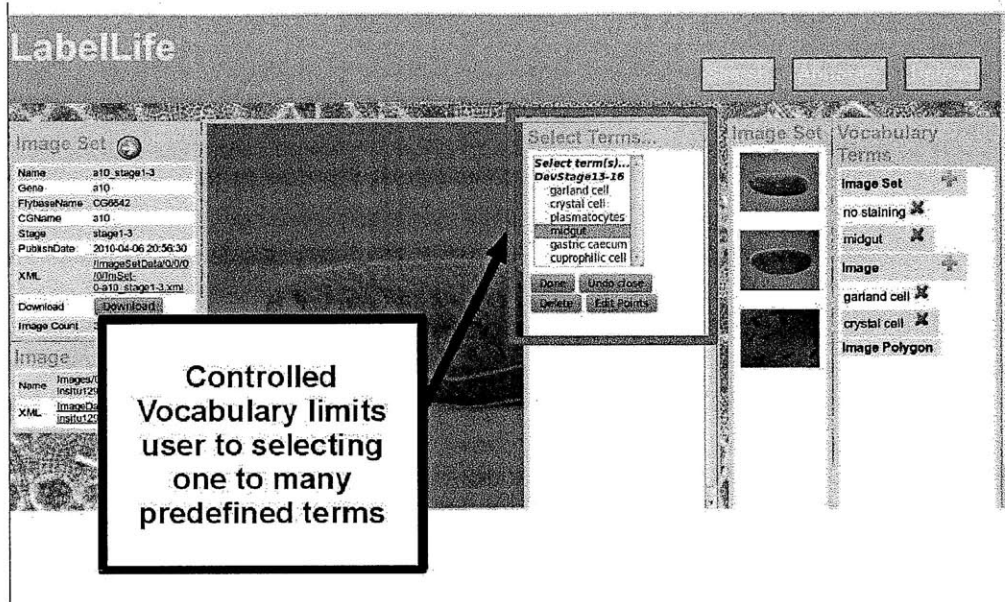


Figure 14: Screenshot of LabelLife demonstrating the concept of Controlled Vocabulary.

- **One-to-many mapping of polygons to labels**

The system needed to support the ability to apply multiple labels to a polygon, image, or set of images. Figure 15 demonstrates LabelMe's scheme for labeling polygons. A free-form text label is requested upon completion of the polygon. Only one label is accepted and for ambiguous objects, the variability in responses will most likely be high. For example, a picture of a dog might have labels such as “dog,” “canine,” “Lassie,” “golden doodle,” etc. Figure 16

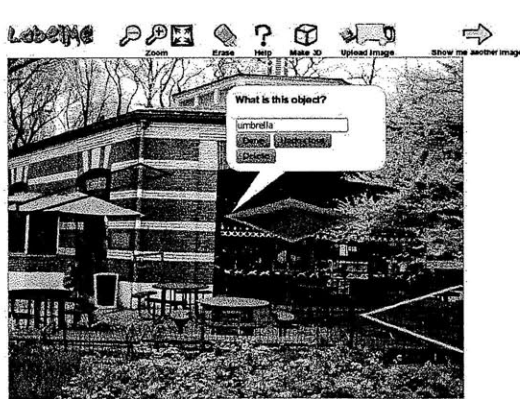


Figure 15: Screenshot of LabelMe labeling stage. A free-form text label is required to complete the polygon.

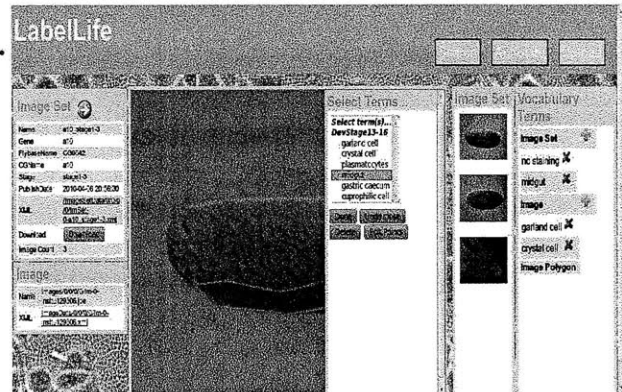


Figure 16: Screenshot of LabelLife labeling stage. One to many controlled vocabulary terms are required to complete the polygon.

3 LabelLife: A Web-based Image Annotation Tool

demonstrates LabelLife's scheme for labeling polygons. A label selection box prompts the user for one to many controlled vocabulary terms that were preloaded into the system. Given the task of labeling very specific embryonic structures, this structure in the system is important to reduce variability in labeling responses and increase annotation accuracy.

- **Integration of predictions into UI**

Once a database of polygon, image, and image set to label mappings were built, it would then be possible to generate label predictions based on the past data. The tool needed to incorporate these predictions in to the UI to offer suggestions to the user.

- **Generic Tool For Many Uses**

The new tool needed to be generic for many uses. The problem of assigning a vocabulary term to an image, set of images, or polygon within an image is not specific to the task of annotating gene expression pattern images. The implementation of the tool needed to be generic so that other groups could adopt it for other purposes.

3.3 Design and Implementation

LabelLife was designed to take into account all of the requirements outlined above. Since LabelLife takes so much inspiration from LabelMe, LabelMe was strongly considered for the task of labeling gene expression pattern images with controlled vocabulary terms. There were many desirable qualities that LabelMe had to accomplish this task. First, it is a very easy to use and intuitive tool to use. It is web-based and supported across multiple browsers. Users do not have to go through a complicated installation and configuration process. The interface is visual and guides the user through

3 LabelLife: A Web-based Image Annotation Tool

the process without wordy instructions. To outline shapes in images, the user only has to point and click and the interface guides the user to complete the polygon to proceed to the labeling step. Second, LabelMe is open to the public. The source code is freely and publicly available through the web site. The developers of LabelMe provide detailed instructions on how to get an instance running on both Windows and Linux. The database is also available for download as long as the user has an account and actively contributes to its growth. There is also a well developed Matlab toolkit for analyzing the data. Third, LabelMe provides the ability to label polygons in images for annotation. These qualities all made LabelMe a viable option for the task at hand, but there were several qualities lacking that necessitated the development of a new system.

LabelLife combined several attributes that worked well with LabelMe with several new features that were necessary for labeling gene expression pattern images. The design for LabelLife incorporated the requirements above that were not supported by LabelMe. First, LabelLife needed to support the labeling of not only polygons, but also images and image sets. The display also needed to handle this newly introduced relationship. Second, LabelLife needed to support a one-to-many mapping of polygons, images, and image sets to labels. Third, LabelLife needed to support controlled vocabularies. LabelMe currently allows free form text entry for assigning labels to polygons. The system was originally designed with this intent in mind to allow for learning from the user about what was an appropriate name for an object [26]. The development team at LabelLife has since integrated WordNet, an electronic dictionary, that generates semantic trees of labels for the analysis phase of LabelMe. [16]. For the purpose of labeling very specific anatomical terms in a developing fly embryo, the design decision was made to restrict users to a limited set of complex vocabulary terms. Fourth, LabelLife had to incorporate important domain specific metadata into the user interface to assist the labeler. Finally, LabelLife had to support presenting predictions to users about what labels are likely to

3 LabelLife: A Web-based Image Annotation Tool

be found in a particular image or stained region of an image.

To accomplish these specific design goals, much thought went into the design of the system.

The system is laid out in detail in the following sections.

3.3.1 Data Model

Figure 17 depicts LabelLife's data model. The following describes each module of the data model. Figure 18 illustrates a screenshot of LabelLife with the objects defined in the object model tagged on the display to communicate what parts of the UI these objects have an effect on.

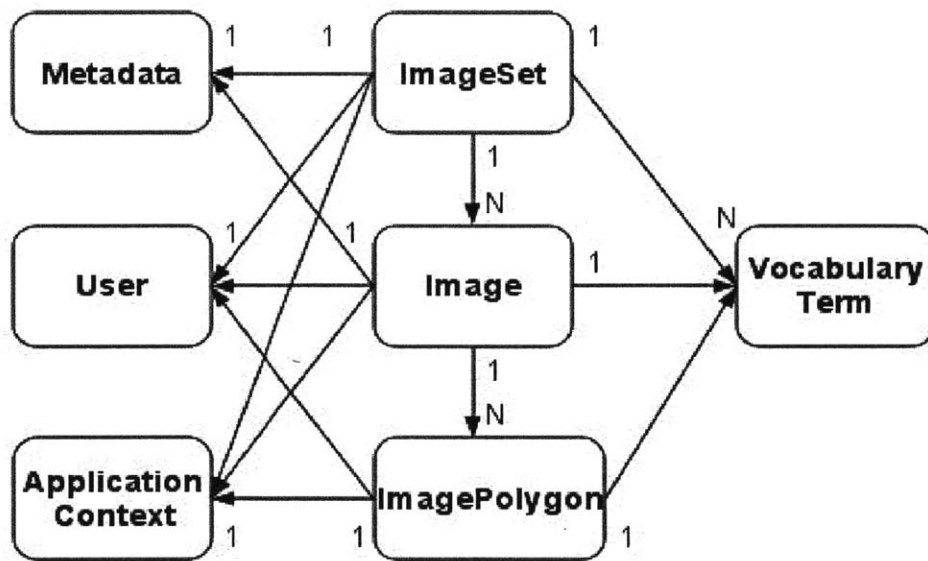


Figure 17: LabelLife data model.

3 LabelLife: A Web-based Image Annotation Tool

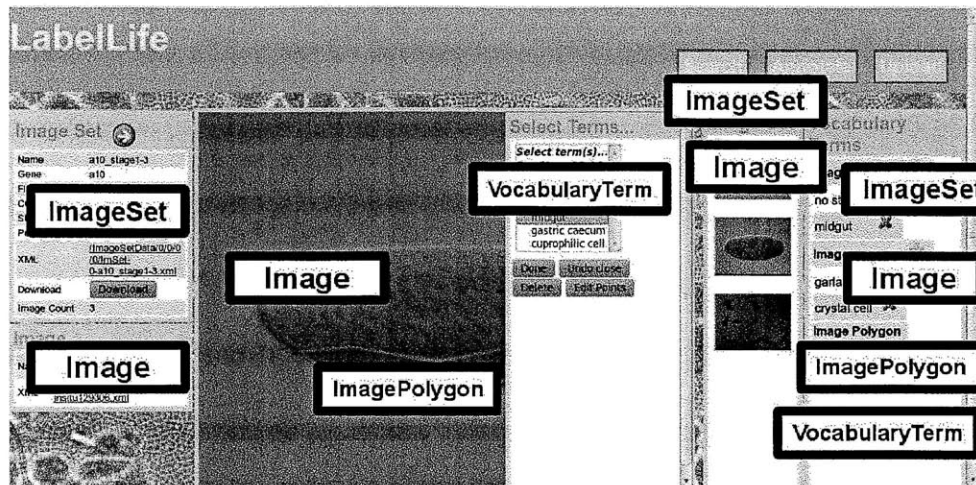


Figure 18: Screenshot of LabelLife demonstrating where the data model objects apply in the user interface.

- **Image**

Image objects are contained by ImageSets. Images contain ImagePolygons. Images also have associated Metadata and VocabularyGroup objects as well as User ownership information.

- **ImageSet**

ImageSets are containers for Images. They have associated Metadata and VocabularyGroup objects. Also, the User who created and modified the ImageSet is stored.

- **ImagePolygon**

ImagePolygons represent polygons that can be hand drawn by pointing and clicking on the website. Image objects contain one-to-many ImagePolygons. ImagePolygons themselves maintain a list of points, a User that created and modified them, and a VocabularyGroup of assigned controlled vocabulary terms.

- **ApplicationContext**

The ApplicationContext objects allows different data sets to stay separate in the LabelLife

3 LabelLife: A Web-based Image Annotation Tool

application. This object keeps track of allowed controlled vocabularies that can be used for applying terms to objects in the data set.

- **User**

The User object stores all necessary information about a user of the application. User information is stored to keep track of who creates and modifies objects. Also, all operations that can be performed require user identification.

- **VocabularyTerm**

A VocabularyTerm represents a controlled vocabulary term that can be applied to an object. The VocabularyGroup object ensures that one to many VocabularyTerms can be applied to an object.

- **VocabularyGroup**

A VocabularyGroup is a container for VocabularyTerms. VocabularyGroups can represent available controlled vocabulary terms for labeling. VocabularyGroups can also represent the actual terms that are assigned to Images, ImageSets, and ImagePolygons.

- **Metadata**

Metadata is a container of key/value mappings that provide relevant information for specific objects. In the context of the problem of labeling gene expression pattern images, examples of Metadata could be the Gene, orientation, staining intensity, and Chromosome location.

3.3.2 System Architecture

Figure 19 depicts LabelLife's system architecture. The system will be described in detail in the following sections.

3 LabelLife: A Web-based Image Annotation Tool

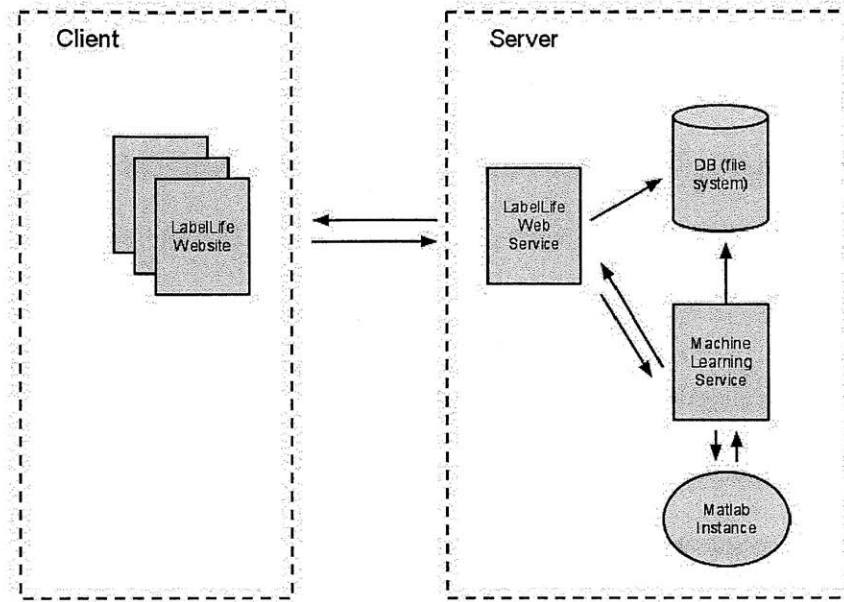


Figure 19: LabelLife Software Architecture

- **LabelLife Website**

Many instances of the LabelLife website can be run at a time. The code currently supports running the website in the Firefox browser and there are plans to support the other most popular browsers on the market today (i.e. Internet Explorer, Chrome, etc.). The website was written using HTML, Javascript, and CSS.

To communicate with the LabelLife Web Service, Javascript was written to execute AJAX calls to retrieve data. Data is passed between the website instances and the LabelLife Web Service in plaintext XML and both sides understand how to parse the result. There is a Javascript class to represent every object laid out in the Data Model. When the XML from a request is received, it is parsed and the information is used to generate its corresponding object.

3 LabelLife: A Web-based Image Annotation Tool

The LabelLife website was built to emulate the easy labeling properties that LabelMe possessed. Much of the code that allows point and click labeling was reused from LabelMe. To display shapes in Firefox, SVG (Support Vector Graphics) markup was used. To display shapes in Internet Explorer, VML (Vector Markup Language) markup was used. LabelMe also introduced a concept in the code of having layered canvases – Draw, Query, Rest, Select. There is a state machine that dictates the available actions and events allowable when each canvas is elevated to the top. This concept greatly simplified and organized the code in an easy-to-understand manner and so it was reused with modification in the LabelLife source code.

When the website first loads, a request is issued to the LabelLife web service to retrieve an ImageSet to display. The LabelLife Web Service returns an ImageSet shell with Image IDs that are contained in that set. The website then issues requests to retrieve every Image object from the LabelLife Web Service. The site is loaded with data as these requests are returned.

On the server side, the LabelLife website is hosted on an Apache2 HTTP Server instance running under Ubuntu Server 10.10. Apache2 was chosen to serve the website because of its positive reputation of stability, speed, and ease of configuration.

- **LabelLife Web Service**

One instance of the LabelLife Web Service runs on an Apache Tomcat Application Server under Ubuntu Server 10.10. The service itself is written in Java. It is run under a Tomcat application

3 LabelLife: A Web-based Image Annotation Tool

called Axis2. Axis2 is an Apache project that delivers a Web Services/SOAP/WSDL engine [21]. Using Axis2, one can quickly code a web service in Java and host it through Tomcat without worrying about the web service layer of the application. Serving and receiving requests, formatting packets, and queuing are all taken care of under this framework. Creating a service is simplified to creating a Java class where the method calls serve as calls into the web service.

This implementation scheme was chosen for LabelLife because of its ease of use. Most people working on the tool in the future will be undergraduate and graduate students. Not every Computer Science student is a web programming expert, but most CS students have been exposed to Java in some capacity, and so programming should be less of a barrier to finding interested applicants to contribute to the project. This scheme also allows for the quick implementation of new functionality. The framework to communicate with a central authority and carry out tasks on system objects is complete. Using the web service, it is straightforward how to go about inserting new functionality into the system.

Using Apache Tomcat, it was possible to serve the website under Tomcat instead of Apache2 HTTP Server. Most industry consensus across Internet forums suggested running websites on Apache2 HTTP Server and web services on Apache Tomcat Application Server because both applications have been optimized for those specific tasks. This advice was heeded when deciding how to run LabelLife's website and web service.

The LabelLife Web Service itself is responsible for serving as the point of contact for retrieving

3 LabelLife: A Web-based Image Annotation Tool

and updating data by the instances of the website. When the Web Service loads, it caches all data in memory for fast access and minimal disk I/O. When data updates are made through the Web Service, the service updates the objects in memory and also writes the changes out to disk.

- **Database**

The Database used for storing all necessary system information is a file system scheme. All objects are written out to disk in XML format. The file system storage layout is illustrated in Figure 20. When the LabelLife Web Service

```
Data/  
  Images/  
  ImageData/  
  ImageSetData/  
  ControlledVocabularies/  
  ApplicationContexts/  
  Users/
```

loads, the file system is read from the XML. The

Figure 20: File system storage layout

corresponding objects are created and loaded into memory. When data updates are made, these changes are written back to disk in XML. The decision was made to use flat, readable XML files on a file system as opposed to a SQL database for the ease of use. A new developer could come onto the project, open up an object's file, and understand the data contained. File systems are easy enough to back up with automated tasks and Java libraries for reading and writing to XML are well defined for this task. Using a SQL database has the added benefit of allowing the data to be queried, but there is added complexity for developers to manage this system especially without an extensive SQL background. Just like Java/Axis2 was chosen to implement the Web Service to foster ease of use and development, using flat files on a file system was also chosen for this same goal.

3 LabelLife: A Web-based Image Annotation Tool

- **Machine Learning Service**

This complete build out of the Machine Learning Service is planned for future work, but offering predictions to users is a main vision of the project and this concept has been already been worked into the system design of LabelLife. The Machine Learning Service is a back end service written in Java. It runs under a daemon process. The Machine Learning Service currently sits idly, but once complete, the service will periodically read the file system for data, run predefined machine learning tasks on the data to build a classification model, and save the model to disk. New data entered into the system will be run against this model to generate predictions and these predictions will be presented to the user. The heavy duty processing will not be carried out by the Machine Learning Service, but instead by an instance of Matlab also running on the same server. The Java Machine Learning Service will issue the appropriate commands to execute a predefined task. A third-party Java library called Matlab Control library is available and will be used to interact with Matlab from Java.

3.3.3 User Interface

Figure 21 illustrates LabelLife's user interface. Each component of the UI will be discussed in detail in the following subsections.

3 LabelLife: A Web-based Image Annotation Tool

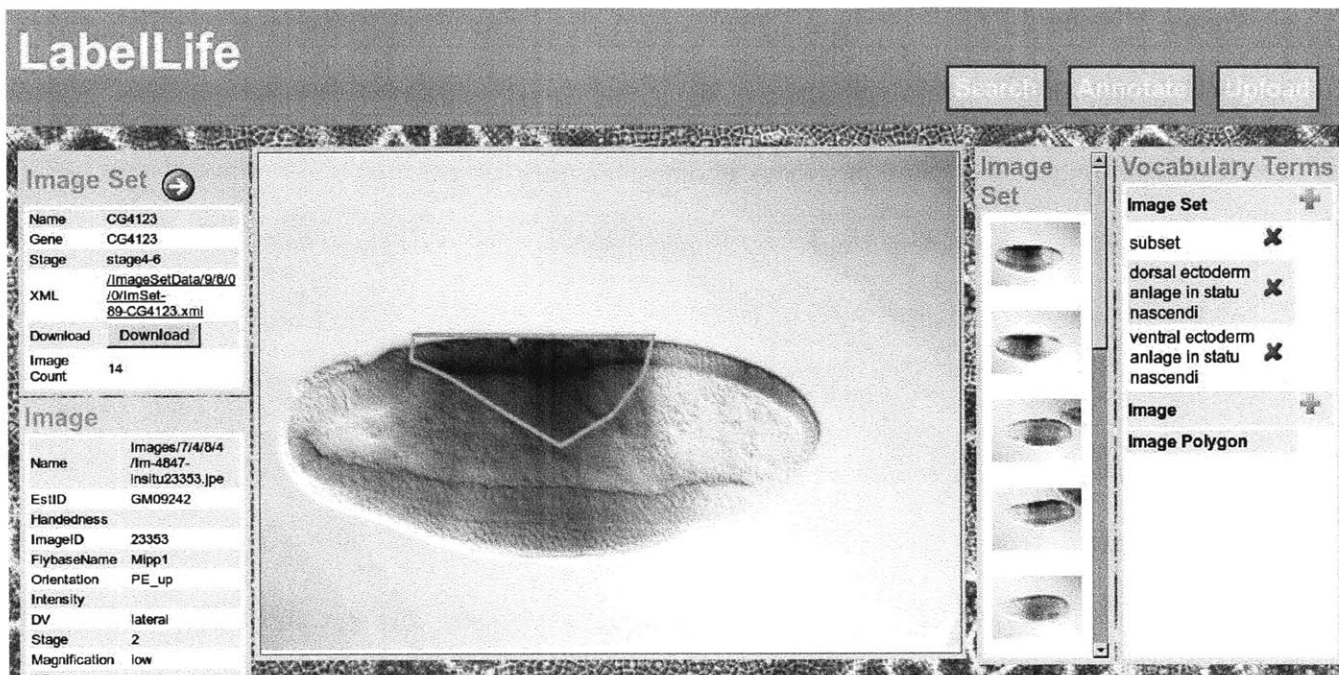


Figure 21: LabelLife user interface

- **ImageSet Metadata**

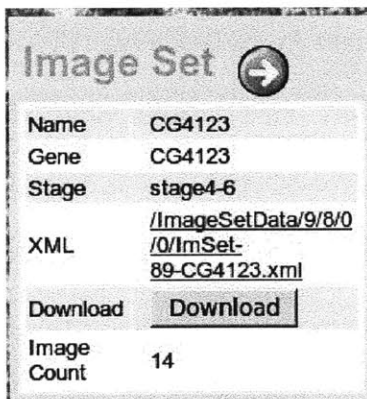


Figure 22: LabelLife UI Component: ImageSet Metadata Pane

This section of the interface displays metadata that is relevant to the whole ImageSet. The data contained here applies to all Images. The underlying storage of Metadata is simple Key/Value pair strings. Depending on how the data is loaded, the user could load anything here. It is flexible to support problem-specific data.

3 LabelLife: A Web-based Image Annotation Tool

- **Image Metadata**

This section of the interface displays metadata that is relevant to the displayed Image. For example, there is data indicating orientation, magnification, development stage, etc. Figure 23 shows an example of loaded Image Metadata.

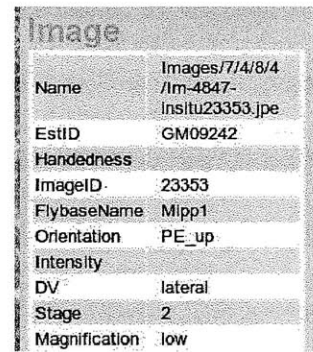
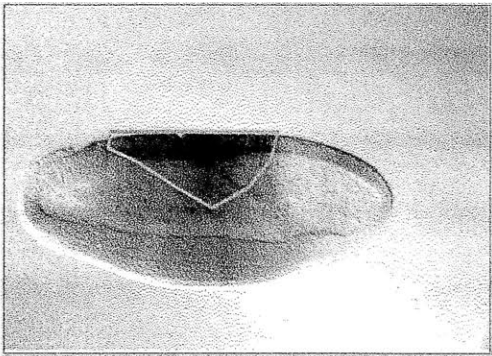


Image	
Name	Images/7/4/8/4 /Im-4847- Insitu23353.jpe
EstID	GM09242
Handedness	
ImageID	23353
FlybaseName	Mipp1
Orientation	PE_up
Intensity	
DV	lateral
Stage	2
Magnification	low

**Figure 23: LabelLife
UI Component: Image
Metadata pane**

- **Image pane**



**Figure 24: LabelLife UI Component:
Main Image pane
drawing is permitted.**

The Image pane displays the currently selected Image. When the Image loads, all of its metadata loads to the left in the Image Metadata pane. Also, all of the Image's assigned vocabulary terms are loaded to the right in the VocabularyTerm pane. This is the only area where point-and-click

3 LabelLife: A Web-based Image Annotation Tool

- **ImageSet pane**

The ImageSet pane provides a scrollable view of thumbnails of all Images contained in the loaded ImageSet. Clicking on one of the Images loads the information for that Image. Also, clicking on the ImageSet creates a larger display to show all Images contained in the set in a grid-like fashion. This is useful for large ImageSets. An example of this functionality can be seen in Figure 27.

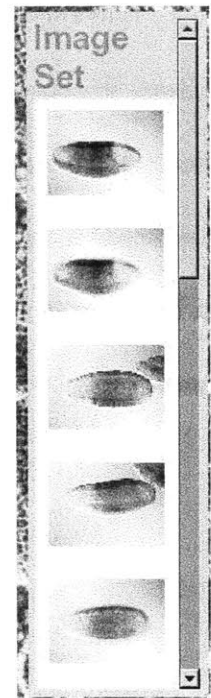


Figure 25:
LabelLife UI
Component:
ImageSet pane

- **VocabularyTerms**

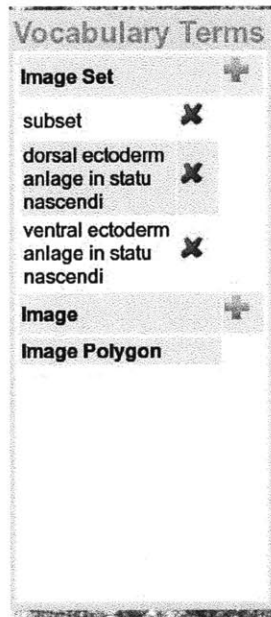


Figure 26: LabelLife
UI Component:
VocabularyTerm pane

The VocabularySet pane shows the terms that are applied at the three different objects levels: ImageSet, Image, and ImagePolygon. ImageSet and Image have add buttons after the heading. Selecting this button brings up the labeling prompt and terms can be added to these objects by multi-selecting terms and selecting 'Done.' There is no add button after the ImagePolygon heading. To add VocabularyTerms to an ImagePolygon, either a new ImagePolygon must be completed or an existing one must be selected to bring up the labeling prompt. Once this prompt is displayed, multiple terms can be selected and choosing 'Done' applies the terms to the shapes. Figure 28 shows the labeling prompt for

3 LabelLife: A Web-based Image Annotation Tool

Images and ImageSets. Figure 29 shows the labeling prompt for ImagePolygons.

3.4 Use Cases

The design and implementation of the LabelLife system ensures that the following actions are possible:

1. Adding and Removing ImagePolygons to Images
2. Adding and Removing VocabularyTerms to ImageSets, Images, and ImagePolygons.
3. Displaying Metadata.
4. Toggling through all Images in an ImageSet and updating Image-specific data.
5. Viewing different ImageSets.
6. Viewing different ApplicationContexts.
7. Pointing and clicking to add ImagePolygons to Images.

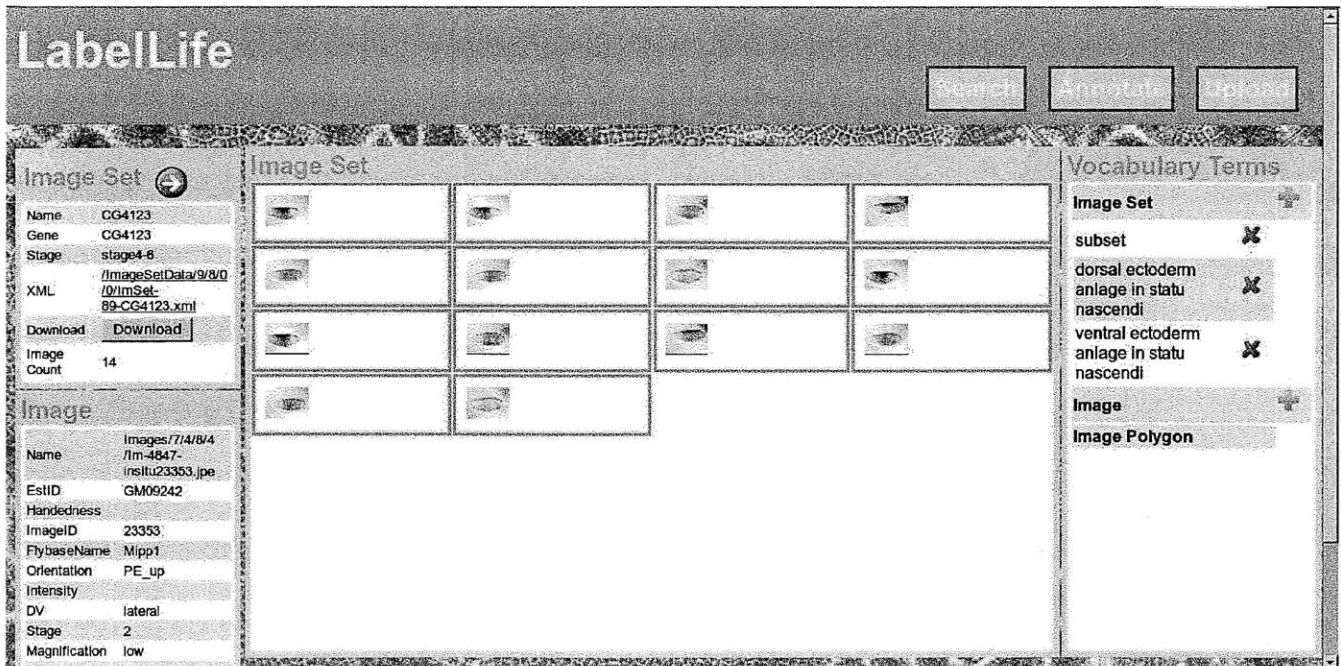


Figure 27: Screenshot of LabelLife expanded ImageSet view. This is useful for viewing large ImageSets.

3 LabelLife: A Web-based Image Annotation Tool

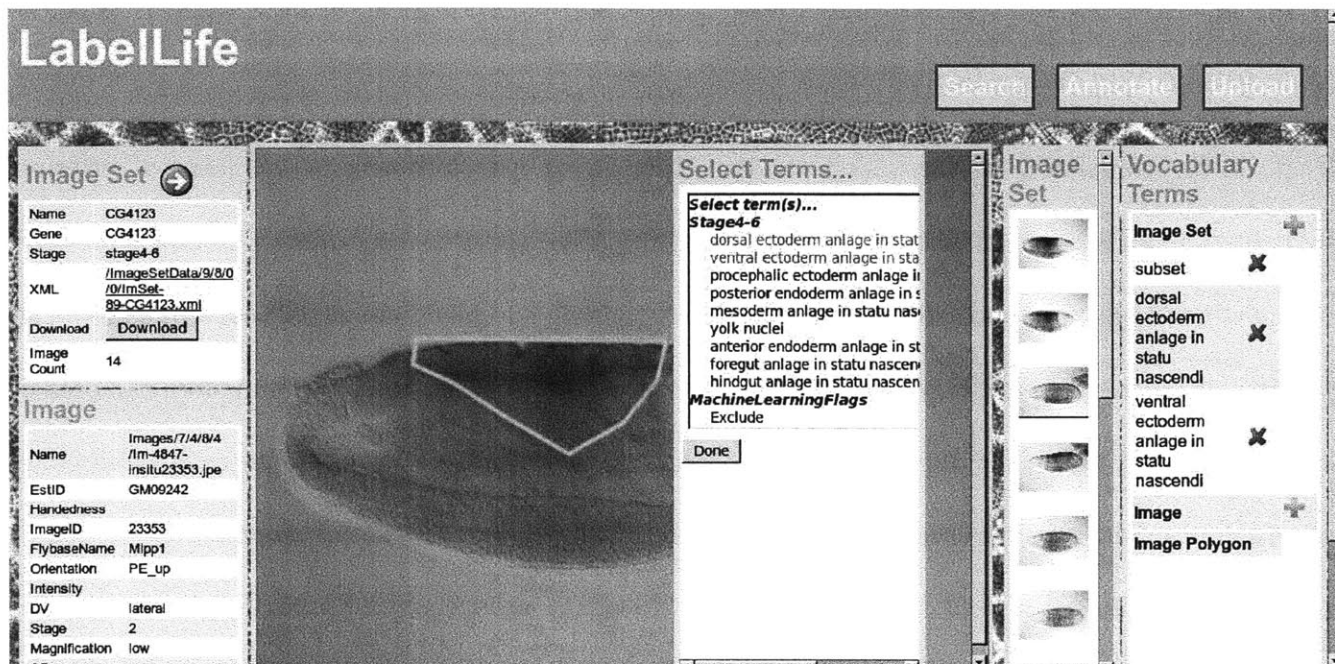


Figure 28: Screenshot of LabelLife Image/Image Set Labeling View. The terms selected will be applied to the ImageSet as a whole. It can be seen that “dorsal ectoderm anlage in statu nascendi” and “ventral ectoderm anlage in statu nascendi” are already applied to this ImageSet.

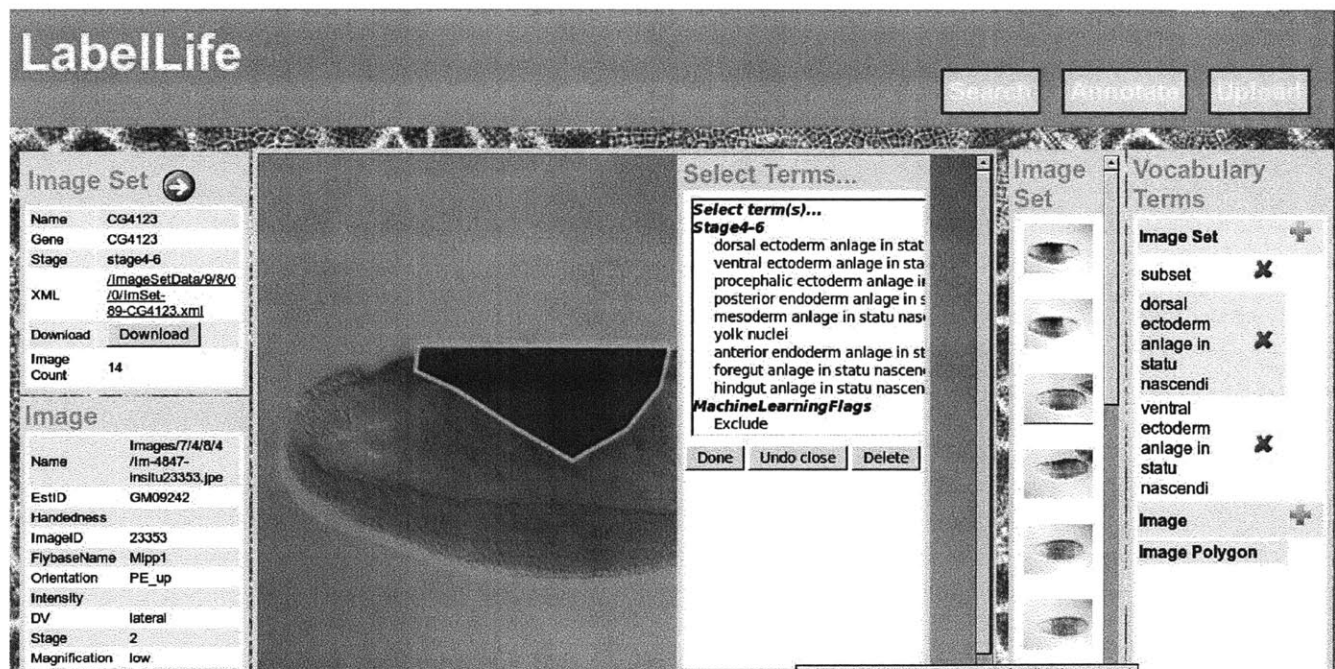


Figure 29: Screenshot of LabelLife ImagePolygon Labeling View. The terms selected will be applied to the highlighted polygon in the Image.

3 LabelLife: A Web-based Image Annotation Tool

3.5 Ongoing Work

3.5.1 Integrating predictions into UI

A big component of the vision for LabelLife is the user feedback in the form of offering predictions about the data. There are several types of predictions that can be offered. First, the tool can predict which ImagePolygons an Image has based on staining and pixel intensity. Second, the tool can predict which controlled vocabulary terms should be associated to an Image. Third, the tool can predict which controlled vocabulary terms should be associated to an ImagePolygon. The second two goals require the completion of the Machine Learning Service and for the methods to generate the prediction to achieve a certain level of accuracy.

3.5.2 UI Aesthetic Improvement

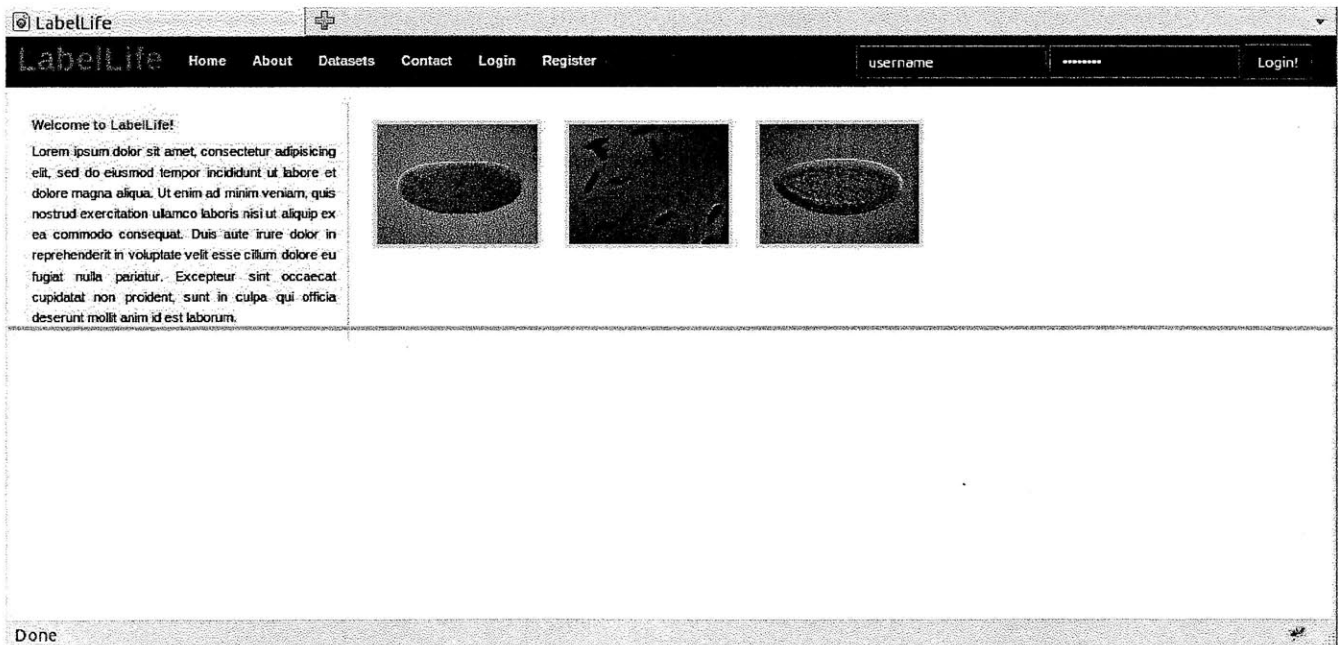


Figure 30: LabelLife Redesign. Image courtesy of Bradley Wu, an MIT undergraduate working on UI improvement

LabelLife is a tool that is meant to be adopted by the academic and research communities. For

3 LabelLife: A Web-based Image Annotation Tool

this to happen, the tool has to have a clean, polished, and professional UI to encourage adoption. There is room to improve what is already implemented to achieve these goals. Progress in this area has already been made by Bradley Wu, an undergraduate in MIT's Electrical Engineering and Computer Science department, as part of a senior project. His contribution to UI aesthetic improvement has involved changes to site formatting, adding user login, and creating a flow to progressing through the website (Figure 30). Another aspect of UI aesthetic improvements adding more tools for the outlining process. For example, the tool now currently supports pointing and clicking to designate points and connect them together. The tool could allow free from dragging to outline object or painting over stained regions. Jennifer Wang, an undergraduate in MIT's Electrical Engineering and Computer Science department, is currently working on this problem to find an efficient way to make drawing outlines faster and more intuitive.

3.5.3 Data Download and Analysis Tools

LabelLife is meant to be a completely open and available tool that researchers can use to perform analysis. The intention is there but work needs to be completed to make this a reality. Bradley Wu has also enabled downloading of data sets through his updates to the site. Currently, it is possible to download the data as an XML file. To make these downloads useful, a Matlab library needs to be built for processing the data.

3.5.4 Mechanical Turk

Currently, field experts are the ones assigning labels to groups of gene expression pattern images. When completing such a manual task, it is common to ask if the task could be farmed out to human labor in a way that produces accurate results. A vehicle for accessing cheap human labor is

3 LabelLife: A Web-based Image Annotation Tool

Amazon's Mechanical Turk [27]. People who need work done post tasks that are compensated a certain amount of money per hit (usually under \$1.00). Workers can peruse the available jobs to select something of interest and complete the work. Since labeling fruit fly embryos is such a specialized task, encouraging and teaching accuracy is a difficult problem, but a direction that LabelLife will attempt in the future. Stanislav Nikolov, an undergraduate in MIT's Electrical Engineering and Computer Science department, has configured LabelMe for the task of farming out the labeling of embryos to Mechanical Turk since LabelMe provides instructions and code. He has made significant strides by conducting training experiments on Internet users to accurately label different morphological features in the embryos. For LabelLife to catch up to this direction of interest, support must be built for Mechanical Turk and training must carefully be considered in order to achieve useful results. LabelMe developers went through the step of removing much of the normal LabelMe user interface for a very simple Mechanical Turk LabelMe user interface. These same steps must be carried out for LabelLife.

3.5.5 Search Functionality

LabelLife is a website that contains a large amount of data so there needs to be a fast and effective way to search for items of interest. There are many different things that could be searched. First, Image and ImageSet metadata could be searched to find Key/Value pairs of interest. Second, the controlled vocabulary terms could be searched to see which Images, ImageSets, and ImagePolygons have these terms assigned. Third, Images themselves could be searched for stains in a specific region. This would require image processing. Fourth, ImagePolygon shapes could be searched to see which Images contain similar shapes in a certain location of the image. This mode of search would require a drawing as an input. There are many ways for the data in LabelLife to be searched and these modes

3 LabelLife: A Web-based Image Annotation Tool

should be exposed to the user for fast access to the data.

3.5.6 Video

LabelLife is a tool developed specifically for the task of labeling gene expression pattern images with controlled vocabulary terms, but it was written in a generic way that it can be applied to other problems. Other groups have expressed interest in a tool that could also allow the labeling of video clips. LabelMe has realized this is an immediate need in the research community and they have introduced LabelMe video that uses Flash to display videos for annotation. Olivier Chatot, an undergraduate in MIT's Electrical Engineering and Computer Science department, is currently working on the problem of integrating video support into LabelLife using Flash.

3.6 Other Uses For LabelLife

One of the major requirements for the LabelLife tool was a generic implementation so that it could be adopted by other groups for other purposes. An example of a use for LabelLife outside of labeling gene expression pattern images is labeling mice in scenes and assigning behaviors to indicate what the mice are doing. This example came up in the Center for Biological and Computational Learning, a laboratory at MIT. This is applicable to many Biological experiments where mice are inflicted with a disease and researchers observe the mice to understand how the disease manifests. The mice could be outlined with a tool like LabelLife and altered behavioral states could be assigned with

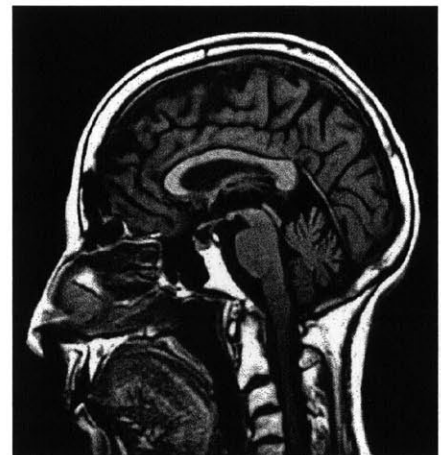


Figure 31: A potential use for LabelLife could be labeling soft tissues in MRI brain scans

3 LabelLife: A Web-based Image Annotation Tool

the controlled vocabulary terms. Another example of a use for LabelLife involves characterizing brain MRI images. (Figure 31) MRI technology is effective at displaying soft tissues and when used on the brain, one can differentiate between white matter, gray matter, and Cerebrospinal fluid (csf) [29]. Automated methods are sought after to separate these components. To train algorithms to do this, experts hand curate the voxels in MRI brain scans with white matter, gray matter, or csf. LabelLife would be a great tool to accomplish such a task with. The curator could load one brain as many slices in an ImageSet into LabelLife. The curator could then outline regions of the brain with polygons and assign the appropriate tissue label.

Another example of a use for LabelLife could be labeling storefronts in Google street view. (Figure 32) This tool could be used to view several vantage points of the same store and the store outline could be labeled with a polygon and assigned many terms such as “supermarket,” “organic,” “self check out”, etc. The feedback mechanism built into the UI could offer suggestions to users about

what image portions should be labeled, what stores could be present given other nearby stores, what other terms should be applied given the currently selected ones, and so on. These are just a few examples of real-world applications out of many that this

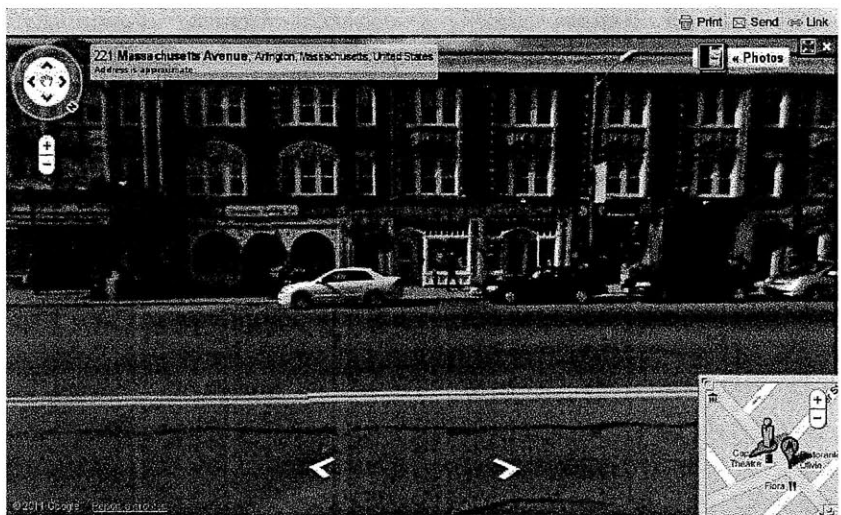


Figure 32: A potential use for LabelLife could be labeling store fronts in GoogleMaps

tool could serve. This tool not only allows for the easy, point-and-click annotation of shapes, images, and sets of images with multiple terms, it also allows for an interactive and adaptive user experience with the integration of predictions into the UI. LabelLife offers a powerful combination of features that

3 LabelLife: A Web-based Image Annotation Tool

are applicable to many problems outside of the Biological one that this thesis focuses on.

3.7 Summary

A tool called LabelLife was developed for the purpose of labeling gene expression image data with a set of controlled vocabulary terms. This application took great influence from LabelMe, a web-based image annotation system embraced by the Computer Vision community. LabelMe offered much inspiration, but new features were needed to support the specific problem at hand. A new web-based image annotation tool was built to incorporate the new functionality. Specific design decisions were made to make the tool easy to use and adoptable by many, and also to make the tool easy to develop for so new participants could easily contribute and make great impact. LabelLife was written in a generic way so that it could serve the purposes of many groups and their research problems. Many future improvements are planned to continually expand the possibilities of what LabelLife is capable of.

4 Machine Learning Methods For Morphological Term Prediction

4.1 Overview

The Berkeley Drosophila Genome Project (BDGP) has built up a large database of data that characterizes *Drosophila melanogaster* embryogenesis in an effort to understand gene expression. The resulting dataset comprises a large set of gene expression pattern images that depict a developing fruit fly embryo with stained regions that indicate where a gene is expressed. The images are grouped together by gene and development stage and these groups are assigned controlled vocabulary terms that indicate which anatomical developmental intermediary structures are present in the set.

The second contribution of this thesis was applying machine learning methods to the problem of automatically predicting which controlled vocabulary terms are associated with a gene expression pattern image. The rate at which this data is generated is far out-pacing the rate at which expert curators can manually assign labels to the data. Also, since labeling is a manual process, there is potential for human error. Incorporating machine learning predictions into this process could not only serve as feedback for the annotators to offer term suggestions or validate chosen terms, but it could also remove the annotator from the process altogether. Another potential benefit of incorporating machine learning to suggest terms is an increase in speed and efficiency in the labeling process. This is why methods that generate predictions with a high degree of accuracy are necessary.

The data relationship generated by the BDGP creates a problem that is well defined for machine learning methods. Machine learning encompasses algorithms that allow computational methods to learn from data and make predictions. Machine learning methods take a set of training examples, each represented by a vector of features, paired with labels, and learn to map new examples into the

4 Machine Learning Methods For Morphological Term Prediction

proper labels. The chosen method generates a model and can be reused to future training inputs to generate predicted labels.

In this case, the input features are the pixel (or some other transformation) of the gene expression pattern image data. The ground truth class labels are whether a specific term is applied to the image or not.

The following sections detail the steps in applying machine learning methods that were carried out.

4.2 Feature Extraction

Before attempting to learn from annotation terms from images, the image must be transformed into a format that will work well with machine learning algorithms. Feature extraction for this

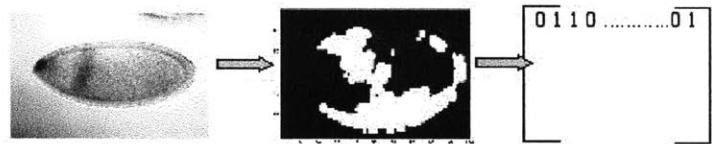


Figure 33: Feature Representation of gene expression pattern image. The embryo is converted to a binarized image indicating the presence/absence of a feature.

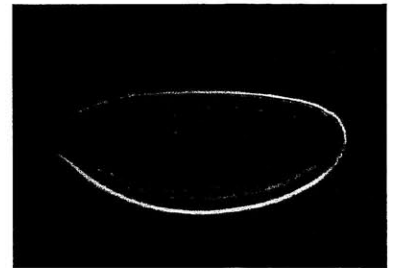
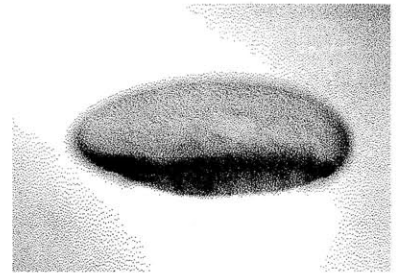
particular problem is typically broken down into two parts. First, image processing is performed to segment the embryo, align images, remove noise, and threshold pixel intensities. Second, the input images might be transformed into different feature representations for machine learning. One reason for transforming into a new feature representation is decreasing dimensionality in data.

The features used in the application of machine learning methods were supplied by Charlie Frogner, a graduate student in the Center for Biological and Computational Learning (CBCL) lab at MIT. Charlie and Tom Morgan, an MIT undergraduate, created an image processing pipeline for the fruit fly embryo images that carries out the following steps:

4 Machine Learning Methods For Morphological Term Prediction

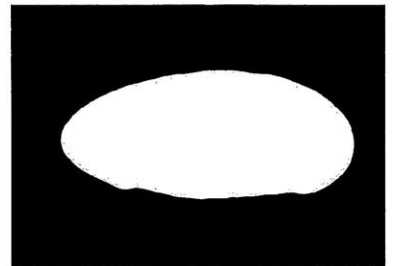
- **Embryo Segmentation**

First, the embryo must be distinguished from the background. A steerable filter is used to isolate embryo regions. This type of filter was chosen because it is robust to blurring outside of the microscope's focal plane. Also, variability in the process of generating the images sometimes produces images with multiple embryos that are touching. Searching for concavities is a tool used to extract the relevant embryo from the image.



- **Embryo Alignment**

Variation in the data generation process also yields embryos that are not ideally aligned – with the anterior region to the left and the posterior region to the right. Often times the embryo sits at a slight angle. To properly align the embryo, the major and minor axes of the embryo are determined using Principal Component Analysis (PCA). After the embryo is aligned, the image is cropped to the boundary of the embryo.



- **Expression Pattern Extraction**

The expression pattern was then extracted from the embryo by standardizing colors and luminance and then applying a learned stained filter to highlight features of interest.



Figure 34: Example of an embryo gene expression pattern image run through image processing.

4 Machine Learning Methods For Morphological Term Prediction

The result from this pipeline is a registered, aligned, thresholded, and downsampled embryo image of binary data where 1 indicates that a feature is present and 0 indicates that a feature is not present. These features correspond to the stain pattern in the gene expression pattern image. Figure 28 illustrates how the data is represented. The data was supplied in Matlab as a feature matrix. Each row corresponded to one image and each column corresponded to one pixel. (Figure 33) The final dimension of the image after being run through this pipeline was 60 x 100 pixels. This image was flattened to a row in the feature matrix and so the result has 6000 columns. Each pixel/column corresponds to a dimension in a machine learning problem. This feature matrix was used as the input to various machine learning Methods.

4.3 Machine Learning Methods To Predict Morphological Terms For Gene Expression Pattern Images

In order to generate an accurate prediction model, two supervised machine learning (ML) methods were explored: Regularized Least Squares Regression and Support Vector Machines. For both ML methods, generating predictions was treated as a set of separate and independent binary classification problems for each controlled vocabulary term present.

4.3.1 Regularized Least Squares Regression (RLS)

RLS is a machine learning approach that aims to minimize the sum of squared errors generated by the prediction. The problem can be set up as follows:

$$\text{Training Set: } S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \quad (1)$$

4 Machine Learning Methods For Morphological Term Prediction

$$\text{Inputs:} \quad \{x_1, x_2, \dots, x_n\} \quad (2)$$

$$\text{Labels:} \quad \{y_1, y_2, \dots, y_n\} \quad (3)$$

$$\text{Goal:} \quad \min \sum_{i=1}^n (f(X_i) - Y_i)^2 + \frac{\lambda}{2} \|f\|_k^2 \quad (4)$$

$$\text{Predictions:} \quad f(X) = K(K + \lambda I)^{-1} Y \quad (5)$$

Equation 1 shows the training set points. Each training point consists of the x input and the y label. The goal of the problem is to find a function $f(X)$ that minimizes the sum of the square loss and the RKHS norm (Equation 4). Equation 5 illustrates the resultant function that minimizes the square loss and generates a prediction based on an input point. This equation is written in term of matrices Y , the labels matrix, and K , the positive semidefinite kernel matrix representing the input training data. λ is the regularization parameter whose function will be described below [31,32].

The goal of RLS is to learn a coefficient matrix C that can be used to predict a label for new, unlabeled data. This coefficient matrix represents a decision boundary that can be applied to data points. Several parameter must be introduced to account for potential difficulties in creating a solid predictive model for the purpose of classifying gene expression pattern images with morphological terms.

First, offsets must be added. RLS is a linear classifier. An offset must be permitted otherwise the resulting decision boundary will be forced to traverse through the origin which might not be suitable for the data. This most certainly is not suitable for the problem of predicting terms for gene expression pattern images since the data has such high dimensionality, therefore offsets were included.

Second, weights must also added. Due to the nature of the gene expression pattern image data, the labeled data can have a preponderance of positive or negative examples of each class. This skew can introduce problems when generating decision boundaries in a classification problem. The resultant

4 Machine Learning Methods For Morphological Term Prediction

decision boundary can artificially ignore errors in the underrepresented class. If relatively few positive occurrences of a class are masked by many occurrences of the absence of that class and all data points are treated equally, the positive data points will add little contribution to the generated classifier. For this reason, it was necessary to introduce weights to add more weight to positive data points.

Third, a regularization parameter must be introduced. The purpose of a regularization parameter is to pull all of the coefficient parameters towards zero. The result is that parameters are set to 0 for linear dimensions that do not pertain to the problem. The goal of the regularization parameter is to impose smoothness and thereby ensure generalization of the solution to the RLS problem. The regularization parameter is another tool to help improve predictive ability.

Finally, due to the high dimensionality of the data, the problem can be converted to the dual format to make use of kernels. Kernels allow data to be mapped to a new feature space that might better suit the data and the prediction goals. For applying RLS to this problem, linear and Gaussian kernels were both tested. The Gaussian kernel was hypothesized to perform better because it performs better for data that has high dimensionality [20].

4.3.2 Support Vector Machine (SVM)

SVM is another supervised linear machine learning method for learning patterns from data. It is a maximum margin classifier just like RLS. This means that the goal of the problem is to maximize the distance between the decision boundary and all separable classes. The main difference is that the goal of this problem is to minimize the hinge loss as opposed to the squared error. Data points that lie directly on the maximum margin are referred to as support vectors. The optimal solution is guaranteed to be attainable by only considering the support vectors as input to the problem.

The problem is set up in a similar manner to RLS:

4 Machine Learning Methods For Morphological Term Prediction

$$\text{Training Set: } S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \quad (6)$$

$$\text{Inputs: } \{x_1, x_2, \dots, x_n\} \quad (7)$$

$$\text{Labels: } \{y_1, y_2, \dots, y_n\} \quad (8)$$

$$\text{Goal: } \min \frac{1}{n} \sum_{i=1}^n (1 - y_i f(x_i)) + \frac{\lambda}{2} \|f\|_H^2 \quad (9)$$

The difference can be seen in Equation 9 where the goal of the optimization problem is to find a function $f(X)$ that minimizes the hinge loss $(1 - y_i f(x_i))$ [33].

SVM is similar to RLS in that it needs extra parameters to run smoothly otherwise there is a chance that the classifier will not converge given certain data. For the reasons stated above, offsets and weights were included for running the SVM classifier. It is also possible with SVMs to convert to the dual with Kernels and use a new feature space to generate predictions. For the problem, SVM was run with weights, offsets, and a radial basis kernel.

4.3.3 Methods of Evaluating Results

Leave One Out Cross Validation (LOOCV)

Leave One Out Cross Validation (LOOCV) is an approach to evaluating the predictive ability of a classifier. LOOCV iterates over the training set of n data points n times and leaves one data point out each time. The left out data point is then classified using the determined decision boundary. Since the prediction and the ground truth are both available, this can be used to determine the error rate. The LOO error rate is an important metric because it indicates how well the classification boundary will do on unseen data. Low LOO error rate is desirable, but this metric alone does not illustrate the complete picture of predictive ability. For example, due to label skew there could be one positive ground truth

4 Machine Learning Methods For Morphological Term Prediction

label out of thousands of data points. If the generated model assigns the negative class to all points, the LOO error rate will be negligible since only one point was missed, but this is bad since the positive class cannot be effectively predicted. Low LOO error rate can sometimes be the result of data that has been overfit due to large skew in input features. Other metrics must be considered in conjunction with LOO error to obtain the complete picture.

Area Under the Curve (AUC)

Area under the curve (AUC) is another metric for evaluating classification performance. (Figure 35) Area under the curve refers to area under the Receiver Operating Characteristic (ROC) curve. This curve plots False Positive Rate (FPR) vs. True Positive Rate (TPR) varied over some parameter. The ideal case would be to have a True Positive Rate of 1 and a False Positive Rate of 0, but this is difficult to achieve in practice. The curve often resembled what is seen in Figure 34. The line $y = x$

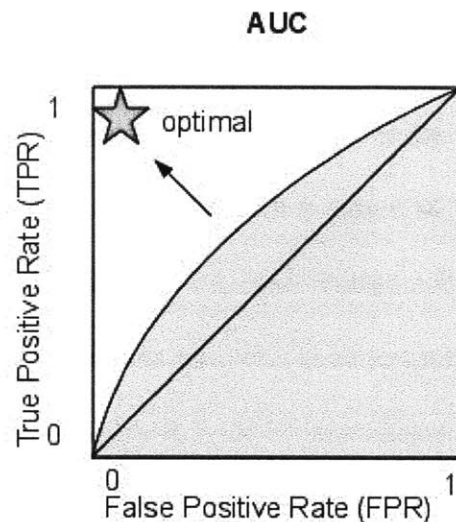


Figure 35: Illustration of area under the curve.

indicates that TPR always equals FPR and this means that the solution is no any better than guessing. The area under this curve would be 0.5. When assessing the performance of a classifier, the goal is to achieve a high AUC. The higher the AUC, the more occurrences of True Positives are seen than False Positives. A low AUC (lower than 0.5) can also be viewed as a positive. It means that there were more False Positives than True Positives and so the definition of the two classes could be flipped to achieve a high AUC. AUC was used as another metric to evaluate classification.

4.4 Data Used For Prediction Problem

4.4.1 *D. melanogaster* Embryo Development Stage 4 To 6

3,365 feature data points corresponding to images were used. All images were of the lateral orientation. These images spanned 1850 genes where each gene had on average 1.81 images per gene.

Term	Frequency	Freq %
ubiquitous	1048	31.14%
dorsal ectoderm anlage in statu nascendi	787	23.39%
procephalic ectoderm anlage in statu nascendi	666	19.79%
ventral ectoderm anlage in statu nascendi	665	19.76%
posterior endoderm anlage in statu nascendi	333	9.90%
faint ubiquitous	291	8.65%
mesoderm anlage in statu nascendi	278	8.26%
anterior endoderm anlage in statu nascendi	269	7.99%
yolk nuclei	264	7.85%
foregut anlage in statu nascendi	235	6.98%
hindgut anlage in statu nascendi	233	6.92%
amnioserosa anlage in statu nascendi	206	6.12%
visual anlage in statu nascendi	189	5.62%
trunk mesoderm anlage in statu nascendi	187	5.56%
pole cell	176	5.23%
head mesoderm anlage in statu nascendi	160	4.75%
ectoderm anlage in statu nascendi	124	3.68%
endoderm anlage in statu nascendi	120	3.57%
mesectoderm anlage in statu nascendi	102	3.03%
ventral ectoderm anlage	68	2.02%
dorsal ectoderm anlage	67	1.99%
clypeolabrum anlage in statu nascendi	50	1.49%
yolk	48	1.43%
trunk mesoderm anlage	39	1.16%
strong ubiquitous	32	0.95%
head epidermis dorsal anlage in statu nascendi	30	0.89%
head mesoderm anlage	30	0.89%
posterior endoderm anlage	14	0.42%
embryonic cortex	9	0.27%
head epidermis anlage in statu nascendi	4	0.12%
hypopharynx anlage in statu nascendi	2	0.06%
mesectoderm anlage	1	0.03%

Table 1: Stage 4 To 6 Morphological Term Breakdown

Each image had on average 1.99 terms per image. Table 1 shows that vocabulary terms that are present in this data set and the frequency they represent of the whole data pool.

Figure 36 shows Z score pixel intensity profiles for each term present in the data set with ubiquitous images excluded from the average of the data. It can be observed from this data set that classification is a

difficult problem since positive examples for frequently observed terms do not differ from the mean

4 Machine Learning Methods For Morphological Term Prediction

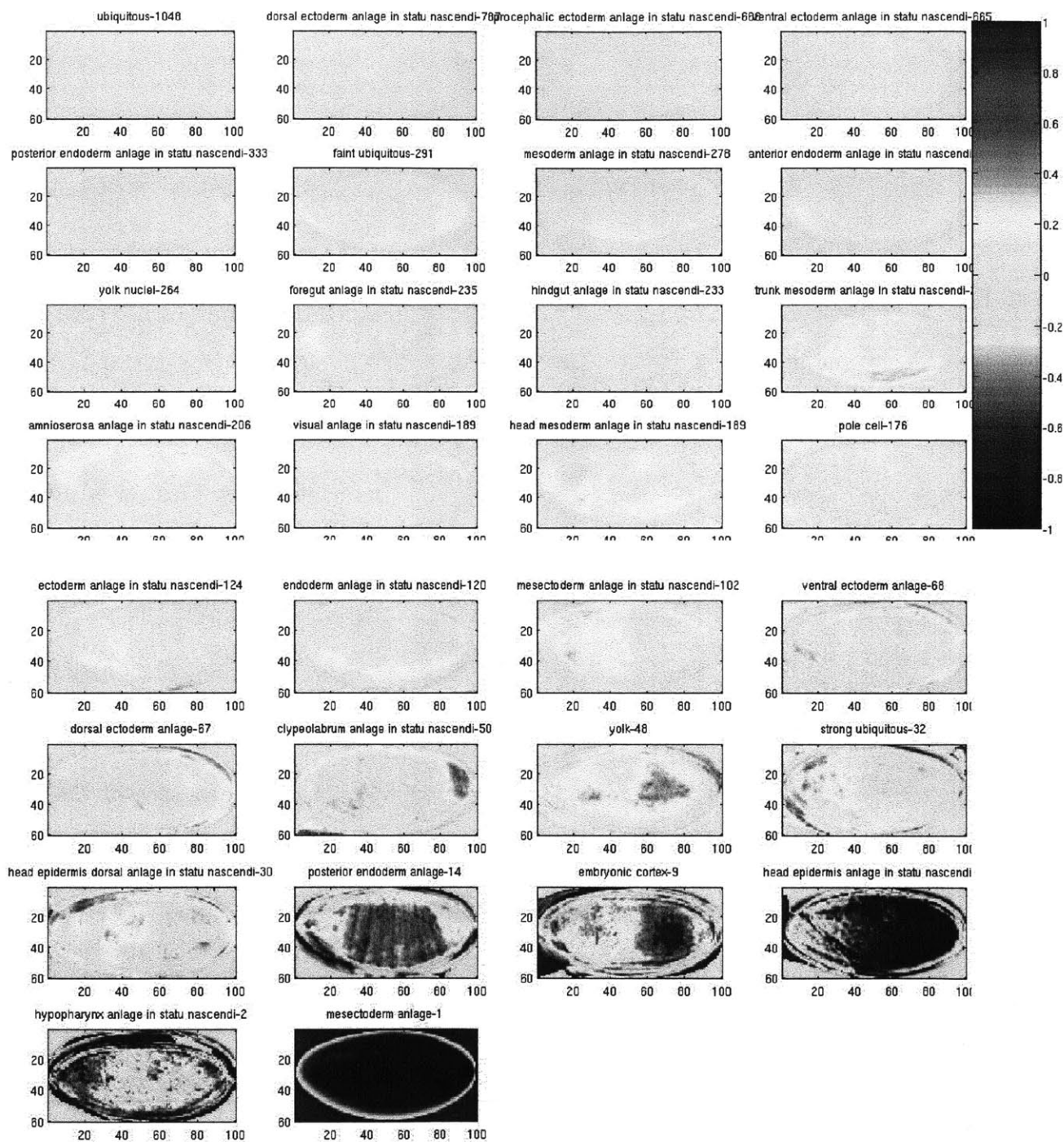


Figure 36: Zscore for all Stage 4 to 6 terms found in the data set sorted in descending order of frequency. Red color indicates positive feature detection and blue color indicates negative feature detection across all relevant embryo images for a term.

4 Machine Learning Methods For Morphological Term Prediction

for the majority of pixels.

4.4.2 *D. melanogaster* Embryo Development Stage 7 To 8

1,857 feature data points corresponding to images were used. All images were of the lateral

Term	Frequency	Freq %
ubiquitous	520	28.00%
posterior endoderm primordium P2	439	23.64%
ventral ectoderm primordium P2	417	22.46%
procephalic ectoderm anlage	368	19.82%
faint ubiquitous	268	14.43%
dorsal ectoderm primordium	238	12.82%
hindgut anlage	226	12.17%
pole cell	114	6.14%
foregut anlage	104	5.60%
yolk nuclei	94	5.06%
visual anlage	68	3.66%
amnioserosa anlage	65	3.50%
mesectoderm primordium	62	3.34%
clypeolabrum anlage	23	1.24%
strong ubiquitous	21	1.13%
central brain anlage	15	0.81%
head epidermis dorsal anlage	9	0.48%
hypopharynx anlage	9	0.48%
procephalic ectoderm primordium	7	0.38%
head epidermis lateral A	4	0.22%
amnioserosa	1	0.05%

Table 2: Stage 7 To 8 Morphological Term Breakdown

orientation. These images spanned 1,354 genes where each gene had on average 1.37 images per gene. Each image had on average 1.99 terms per image. Table 2 shows that vocabulary terms that are present in this data set and the frequency they represent of the whole data pool.

Figure 37 shows Z score pixel intensity profiles for each term present in the data set with ubiquitous images excluded from the average of

the data. This data set also indicates that the classification problem is a difficult one since the most frequently occurring terms have profiles that do not differ from the mean of the data much across all pixels.

4 Machine Learning Methods For Morphological Term Prediction

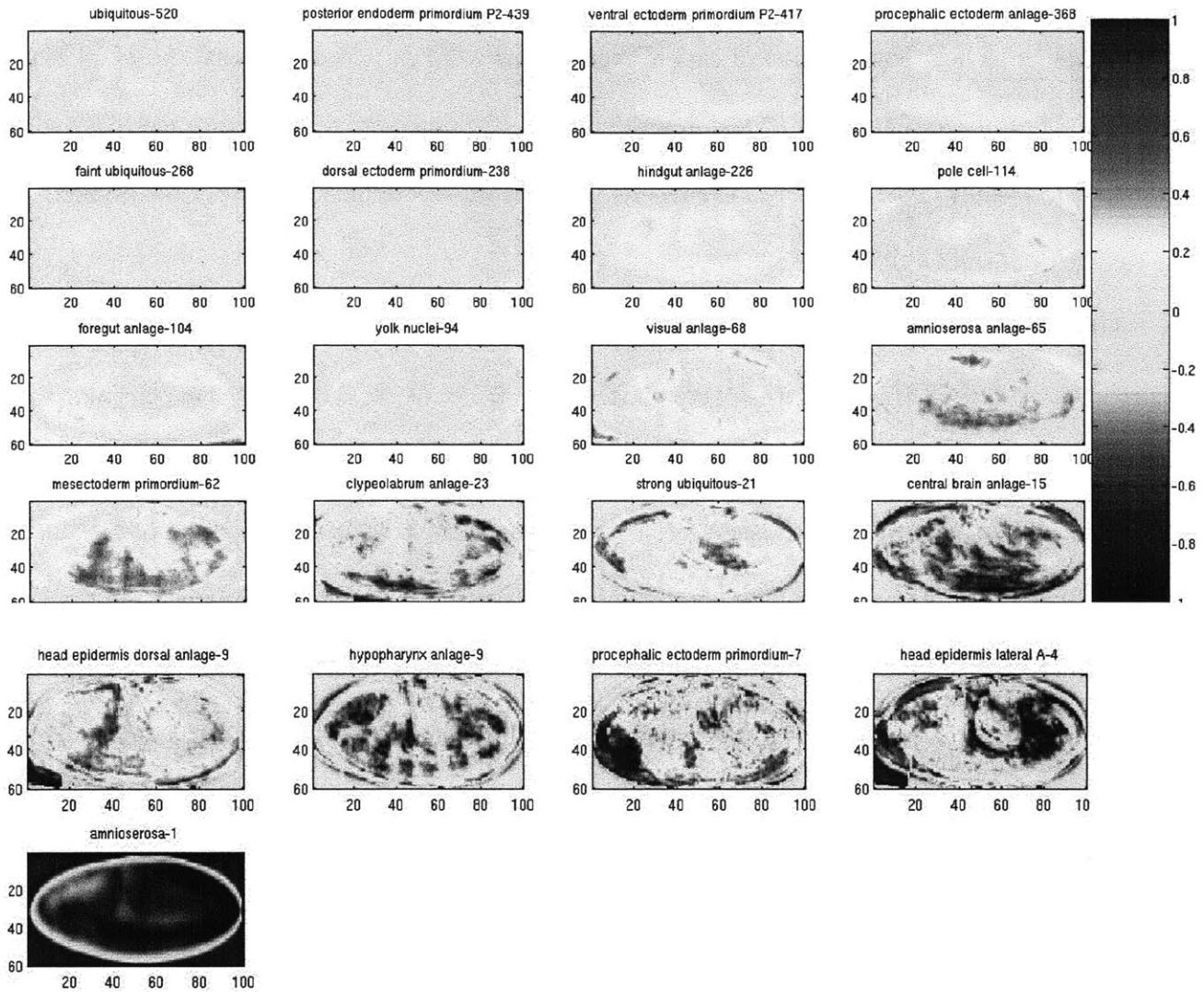


Figure 37: Zscore for all Stage 7 to 8 terms found in the data set sorted in descending order of frequency. Red color indicates positive feature detection and blue color indicates negative feature detection across all relevant embryo images for a term.

4.5 Experiments: Application of Machine Learning Methods To Predict Morphological Terms For Gene Expression Pattern Images

Various machine learning classification experiments were carried out against the pixel feature data input. Multiple combinations of features versus classification were tested across the Stage 4 to 6 and Stage 7 to 8 stage range data sets. The results were reported in terms of AUC per image, AUC per gene (Image Set), and LOO Error Rate.

For most experiments, the features used were the binarized pixel data features. To test out the efficacy of using LabelLife, four terms were chosen from the Stage 4 to 6 data set for cleaning. The terms were “yolk nuclei”, “dorsal ectoderm anlage in statu nascendi,” “ventral ectoderm anlage in statu nascendi,” and “mesoderm anlage in statu nascendi.” ImageSets were created for each of these terms in LabelLife and the data was manually curated to apply labels at the Image level where the term was actually appropriate. Table 3 shows the number of positive labels for these terms before and after cleaning. Figure 38 shows the difference introduced to the pixel intensity Z score by cleaning these labels.

Term	Positive Label Count	
	Before Curation	After Curation
yolk nuclei	264	91
dorsal ectoderm anlage in statu nascendi	787	95
ventral ectoderm anlage in statu nascendi	665	86
mesoderm anlage in statu nascendi	278	73

Table 3: Change in number of positive labels before and after curating data in LabelLife

4 Machine Learning Methods For Morphological Term Prediction

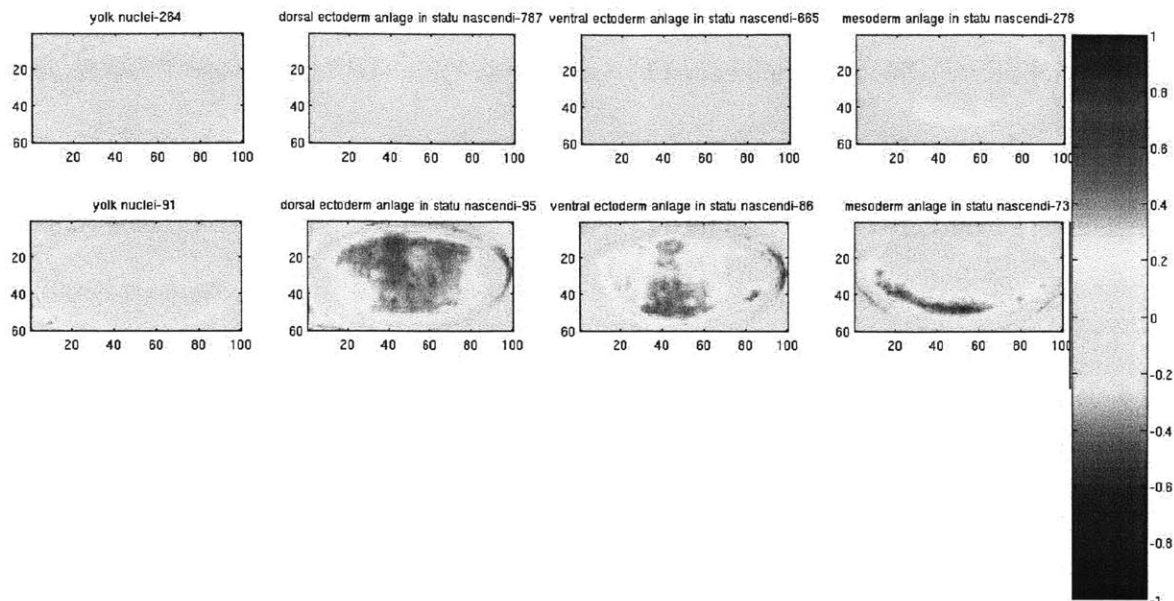


Figure 37: Effect of cleaning labels with LabelLife. The first row corresponds to the Zscore profile of the terms before cleaning. The second row corresponds to the Zscore profile of the same terms after cleaning. It can be seen that cleaning the data through LabelLife produces more striking features for classification.

Three different variations on machine learning methods were tested. First, regularized least squares classification was tested with weights, offsets, and a linear kernel. Second, regularized least squares classification was tested with weights, offsets, and a Gaussian kernel. Third, support vector machines were tested with weights, offsets, and a radial basis kernel.

The following lists the feature/classification method experiments that were carried out.

4.5.1 *D. melanogaster* Embryo Development Stage 4 To 6

- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Linear Kernel

4 Machine Learning Methods For Morphological Term Prediction

- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Gaussian Kernel
- Pixel Feature Image Data cleaned by LabelLife / Weighted Regularized Least Squares Classification With Linear Kernel
- Pixel Feature Image Data cleaned by LabelLife / Weighted Regularized Least Squares Classification With Gaussian Kernel

4.5.2 *D. melanogaster* Embryo Development Stage 7 To 8

- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Linear Kernel
- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Gaussian Kernel

4.6 Results

The results for each experiment are reported in terms of AUC per image, AUC per gene, and LOO error rate across the terms that were investigated. The top 10 frequently appearing terms from each data set were evaluated. AUC per gene was calculated in three ways using the coefficients generated from the classification problem applied to each image data. The coefficients were applied to the feature data to generate predicted labels. First, AUC per gene was calculated by averaging the predicted values across each gene (Image Set). Second, AUC per gene was calculated by taking the max of the predicted values across each gene. Third, AUC per gene was calculated by taking the max predicted value of the set if the ground truth label was one or the min predicted value of the set if the ground truth label was zero. This last method is cheating and could never be used in practice for a

4 Machine Learning Methods For Morphological Term Prediction

robust prediction scheme, but it was used as a benchmark to understand what AUC was attainable given the data.

The results of running SVM were excluded because the combination of parameters used did not produce meaningful results. The produced error rate was low, but AUC was also zero for every case which indicated a bug somewhere in the process. This is one area that requires future work.

4.6.1 *D. melanogaster* Embryo Development Stage 4 To 6

- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Linear Kernel

1. AUC Per Gene

Term	AUC – max	AUC – avg	AUC – cheating
ubiquitous	0.5520	0.5462	0.6888
dorsal ectoderm anlage in statu nascendi	0.7489	0.5790	0.8371
procephalic ectoderm anlage in statu nascendi	0.7322	0.5768	0.8272
ventral ectoderm anlage in statu nascendi	0.7316	0.5485	0.8243
posterior endoderm anlage in statu nascendi	0.7332	0.5846	0.8426
faint ubiquitous	0.5618	0.5540	0.7008
mesoderm anlage in statu nascendi	0.7124	0.5599	0.8274
anterior endoderm anlage in statu nascendi	0.6909	0.5667	0.8101
yolk nuclei	0.6992	0.5271	0.8274
foregut anlage in statu nascendi	0.6997	0.5915	0.8066

Table 4: AUC Per Gene, Stage 4-6, Pixel feature data, RLS classification w/Linear Kernel

2. AUC Per Image, LOO Error Rate Per Image

Term	AUC	LOO Error Rate
ubiquitous	0.5497	0.0056
dorsal ectoderm anlage in statu nascendi	0.5657	0.0051
procephalic ectoderm anlage in statu nascendi	0.5654	0.0048
ventral ectoderm anlage in statu nascendi	0.5294	0.0049
posterior endoderm anlage in statu nascendi	0.5238	0.0037
faint ubiquitous	0.5505	0.0034
mesoderm anlage in statu nascendi	0.5504	0.0033
anterior endoderm anlage in statu nascendi	0.5495	0.0033
yolk nuclei	0.4755	0.0033
foregut anlage in statu nascendi	0.5507	0.0031

Table 5: AUC Per Image, Stage 4-6, Pixel feature data, RLS classification w/Linear Kernel

4 Machine Learning Methods For Morphological Term Prediction

- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Gaussian Kernel

1. AUC Per Gene

Term	AUC – max	AUC – avg	AUC – cheating
ubiquitous	0.5500	0.5500	0.6796
dorsal ectoderm anlage in statu nascendi	0.6933	0.6933	0.7841
procephalic ectoderm anlage in statu nascendi	0.6757	0.6757	0.7775
ventral ectoderm anlage in statu nascendi	0.6734	0.6734	0.7689
posterior endoderm anlage in statu nascendi	0.6655	0.6655	0.7821
faint ubiquitous	0.5833	0.5833	0.7150
mesoderm anlage in statu nascendi	0.6441	0.6441	0.7563
anterior endoderm anlage in statu nascendi	0.6721	0.6721	0.7855
yolk nuclei	0.6374	0.6374	0.7583
foregut anlage in statu nascendi	0.6384	0.6384	0.7552

Table 6: AUC Per Gene, Stage 4-6, Pixel feature data, RLS classification w/Gaussian Kernel

2. AUC Per Image / LOO Error Rate Per Image

Term	AUC	LOO Error Rate
ubiquitous	0.5472	0.0079
dorsal ectoderm anlage in statu nascendi	0.5255	0.0072
procephalic ectoderm anlage in statu nascendi	0.5194	0.0068
ventral ectoderm anlage in statu nascendi	0.5335	0.0068
posterior endoderm anlage in statu nascendi	0.5301	0.0051
faint ubiquitous	0.5431	0.0048
mesoderm anlage in statu nascendi	0.5618	0.0047
anterior endoderm anlage in statu nascendi	0.5222	0.0047
yolk nuclei	0.5123	0.0046
foregut anlage in statu nascendi	0.5091	0.0044

Table 7: AUC Per Image, Stage 4-6, Pixel feature data, RLS classification w/Gaussian Kernel

- Pixel Feature Image Data cleaned by LabelLife / Weighted Regularized Least Squares Classification With Linear Kernel

4 Machine Learning Methods For Morphological Term Prediction

1. AUC Per Image / LOO Error Rate Per Image

Term	LabelLife Curated Data		Non-Curated Data	
	AUC	LOO Error Rate	AUC	LOO Error Rate
dorsal ectoderm anlage in statu nascendi	0.7225	0.0019	0.5657	0.0051
ventral ectoderm anlage in statu nascendi	0.7222	0.0018	0.5294	0.0049
mesoderm anlage in statu nascendi	0.6333	0.0017	0.5504	0.0033
yolk nuclei	0.4385	0.0020	0.4755	0.0033

Table 8: AUC Per Image, Stage 4-6, Pixel feature data before and after cleaning in LabelLife, RLS classification w/Linear Kernel

- Pixel Feature Image Data cleaned by LabelLife / Weighted Regularized Least Squares

Classification With Gaussian Kernel

1. AUC Per Image / LOO Error Rate Per Image

Term	LabelLife Curated Data		Non-Curated Data	
	AUC	LOO Error Rate	AUC	LOO Error Rate
dorsal ectoderm anlage in statu nascendi	0.5807	0.0029	0.5255	0.0072
ventral ectoderm anlage in statu nascendi	0.5494	0.0027	0.5355	0.0068
mesoderm anlage in statu nascendi	0.5399	0.0025	0.5618	0.0047
yolk nuclei	0.4952	0.0028	0.5123	0.0046

Table 9: AUC Per Image, Stage 4-6, Pixel feature data before and after cleaning in LabelLife, RLS classification w/Gaussian Kernel

4.6.2 *D. melanogaster* Embryo Development Stage 7 To 8

- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Linear Kernel

1. AUC Per Gene

4 Machine Learning Methods For Morphological Term Prediction

Term	AUC – max	AUC – avg	AUC – cheating
ubiquitous	0.5607	0.5693	0.6435
posterior endoderm primordium P2	0.5804	0.5384	0.6579
ventral ectoderm primordium P2	0.6419	0.5876	0.6972
procephalic ectoderm anlage	0.6543	0.6003	0.7102
faint ubiquitous	0.5327	0.5333	0.6278
dorsal ectoderm primordium	0.5889	0.5044	0.6598
hindgut anlage	0.6154	0.5838	0.6912
pole cell	0.5864	0.5929	0.6827
foregut anlage	0.6273	0.5503	0.6989
yolk nuclei	0.6690	0.5532	0.7439

Table 10: AUC Per Gene, Stage 7-8, Pixel feature data, RLS classification w/Linear Kernel

2. AUC Per Image / LOO Error Rate Per Image

Term	AUC	LOO Error Rate
ubiquitous	0.5626	0.0073
posterior endoderm primordium P2	0.5347	0.0070
ventral ectoderm primordium P2	0.5883	0.0068
procephalic ectoderm anlage	0.5964	0.0065
faint ubiquitous	0.5366	0.0058
dorsal ectoderm primordium	0.4928	0.0055
hindgut anlage	0.5878	0.0053
pole cell	0.5582	0.0039
foregut anlage	0.5583	0.0038
yolk nuclei	0.4649	0.0037

Table 11: AUC Per Image, Stage 7-8, Pixel feature data, RLS classification w/Linear Kernel

- Pixel Feature Image Data / Weighted Regularized Least Squares Classification With Gaussian Kernel

1. AUC Per Gene

Term	AUC – max	AUC – avg	AUC – cheating
ubiquitous	0.5495	0.5564	0.6387
posterior endoderm primordium P2	0.5737	0.5618	0.6413
ventral ectoderm primordium P2	0.5416	0.5219	0.5997
procephalic ectoderm anlage	0.5660	0.5464	0.6283
faint ubiquitous	0.5437	0.5456	0.6212
dorsal ectoderm primordium	0.5659	0.5384	0.6353
hindgut anlage	0.5573	0.5369	0.6318
pole cell	0.5574	0.5542	0.6402
foregut anlage	0.5435	0.5257	0.6068
yolk nuclei	0.5494	0.5273	0.6262

Table 12: AUC Per Gene, Stage 7-8, Pixel feature data, RLS classification w/Gaussian Kernel

4 Machine Learning Methods For Morphological Term Prediction

2. AUC Per Image / LOO Error Rate Per Image

Term	AUC	LOO Error Rate
ubiquitous	0.5755	0.0104
dorsal ectoderm anlage in statu nascendi	0.5573	0.0098
procephalic ectoderm anlage in statu nascendi	0.5077	0.0097
ventral ectoderm anlage in statu nascendi	0.5238	0.0092
posterior endoderm anlage in statu nascendi	0.5513	0.0081
faint ubiquitous	0.4985	0.0077
mesoderm anlage in statu nascendi	0.5078	0.0076
anterior endoderm anlage in statu nascendi	0.5515	0.0056
yolk nuclei	0.5162	0.0053
foregut anlage in statu nascendi	0.5004	0.0051

Table 13: AUC Per Image, Stage 7-8, Pixel feature data, RLS classification w/Gaussian Kernel

4.7 Discussion

4.7.1 Results

There are several observations that can be made from the results above. First, the linear Kernel when used with Regularized Least Squares (RLS) Linear Regression outperformed the Gaussian Kernel in all cases. This was not an expected result. Gaussian Kernels are known to handle high dimensionality in data well. One way to investigate this interesting outcome would be to vary specific parameters to the Gaussian Kernel to look for improvement. Second, it can be observed that predicting labels based on Genes (ImageSets) outperforms predicting labels based on Images in all cases. The results are comparable to the results that other groups working on this problem have achieved. The scheme of computing AUC that worked best involved taking the maximum predicted label out of each Image Set to represent the gene. The cheating scheme always performed the best but this could never be used in practice. Third, it can be observed from the data that using LabelLife to clean labels can offer significant improvement to predictive models. For the four terms chosen, three out of the four

4 Machine Learning Methods For Morphological Term Prediction

terms showed improvement in AUC when applying RLS with both linear and Gaussian kernels. RLS with a linear kernel showed the most striking improvement in AUC. This further solidifies the point that it is necessary to capture more detail in the data relationship between gene expression pattern images and terms. Grouping terms with ImageSets offers some predictive power, but grouping terms with Images increases that predictive power. It is theorized that going a step further and applying the labels to polygons within the images will offer the most predictive power.

4.7.2 Future Work

There are several things that should be investigated further as a follow up to this work. First, SVM was used for classification but meaningful results were not generated to report in this document. It is necessary to investigate the parameters used to not only achieve low error rate, but also a high AUC. The results produced a zero AUC indicating a bug in the process and this needs to be improved upon for the evaluation of this method. Second, to understand how the data sets used for this thesis would perform for other groups who reported high predictive ability, it is necessary to use the same data set here on their methods. The group that achieved the best predictive ability working on the same problem was Shiuwang *et al* (2008). They termed their method “multi-label learning.” Their method consisted of something similar to Regularized Least Squares Linear Regression with the addition of a shared subspace parameter of term relationships to optimize. The goal of this group was to better capture the relationships between terms. This group achieved high AUC and so this benchmark could be used to compare this data set against. This group's code is readily available and the data set used in this thesis should be tested against their term prediction pipeline. Third, a final point that should be investigated is using LabelLife to fine tune labels even more. An experiment was carried out to assign terms at the Image level for 4 terms to evaluate changes to predictive metrics. It was demonstrated that

4 Machine Learning Methods For Morphological Term Prediction

for three out of the four terms, using LabelLife to fine tune labels did result in predictive improvement. The next step of this experiment would be to draw polygons surrounding the stained regions in the gene expression pattern images and then assign terms to these regions. This would allow for the most detailed data relationship to be captures and it is theorized that this method of label assignment would continue to improve prediction metrics even more.

5 Conclusion

This thesis provided two contributions to the problem of automatically classifying *Drosophila melanogaster* gene expression pattern images with morphological controlled vocabulary terms. The first contribution of this thesis was a web-based annotation tool called LabelLife that was written to provide a means of accurately and efficiently annotating images. LabelLife was designed with a few key goals in mind. First, it was designed to provide a way of generating new data relationships not available in other tools. LabelLife allows labeling of sets of images, individual images, and shapes within images. LabelLife also allows a variable number of controlled vocabulary terms to be applied to each of these objects. Second, LabelLife was designed to be as intuitive and easy-to-use as possible with point-and-click drawing and a display full of relevant information. One is able to draw shapes with their mouse around regions of interest and apply labels. The user interface was built with simplicity and speed in mind. Third, it was designed to be an interactive experience that could eventually offer suggestions to users annotating images. Finally, it was designed to be generic and open to the public to promote adoption not only for the Biological community, but for all who have an image labeling task in mind.

The second contribution of this thesis was an analysis of the application of machine learning methods to the problem of classifying gene expression pattern images with terms. The results of this analysis provided a few key results. First, it was demonstrated that using Regularized Least Squares (RLS) linear regression, it was possible to achieve better predictions using a linear kernel than a Gaussian kernel. Second, it was shown that computing predictive metrics at the Gene (ImageSet) level as opposed to the individual Image level provided the best results. This is consistent with current relationship between the gene expression pattern image and label data. This is also consistent with the

5 Conclusion

characteristics of the data demonstrated by the Zscore profiles of the images. These images illustrated the inherent difficulty of this classification problem since many term profiles did not differ significantly from the mean of the dataset. Third, it was demonstrated that using LabelLife to clean the data and apply terms at the Image level to remove misclassified images greatly improved predictive ability.

The combination of the newly developed LabelLife tool and the analyzed machine learning methods significantly contributes to the problem of automatically predicting labels for gene expression pattern images. The tool provides a powerful, easy-to-use visual interface for quickly labeling data and capturing new data relationships. The initial analysis of machine learning methods provides promising ability to create predictive models from the data. These predictive models will be incorporated into the system to offer suggestions to the user to further enhance the labeling process. It was demonstrated that using LabelLife to hand curate data strengthens data relationships and further improves predictive ability. Using the tool in conjunction with machine learning methods as a feedback loop to continuously improve the methods will expedite the process of analyzing the large and growing *Drosophila* gene expression pattern data set. Simplifying and automating this process will increase the speed at which data can be processed, trends can be understood, and important discoveries concerning the complexities of gene regulation can be made.

6 Appendix

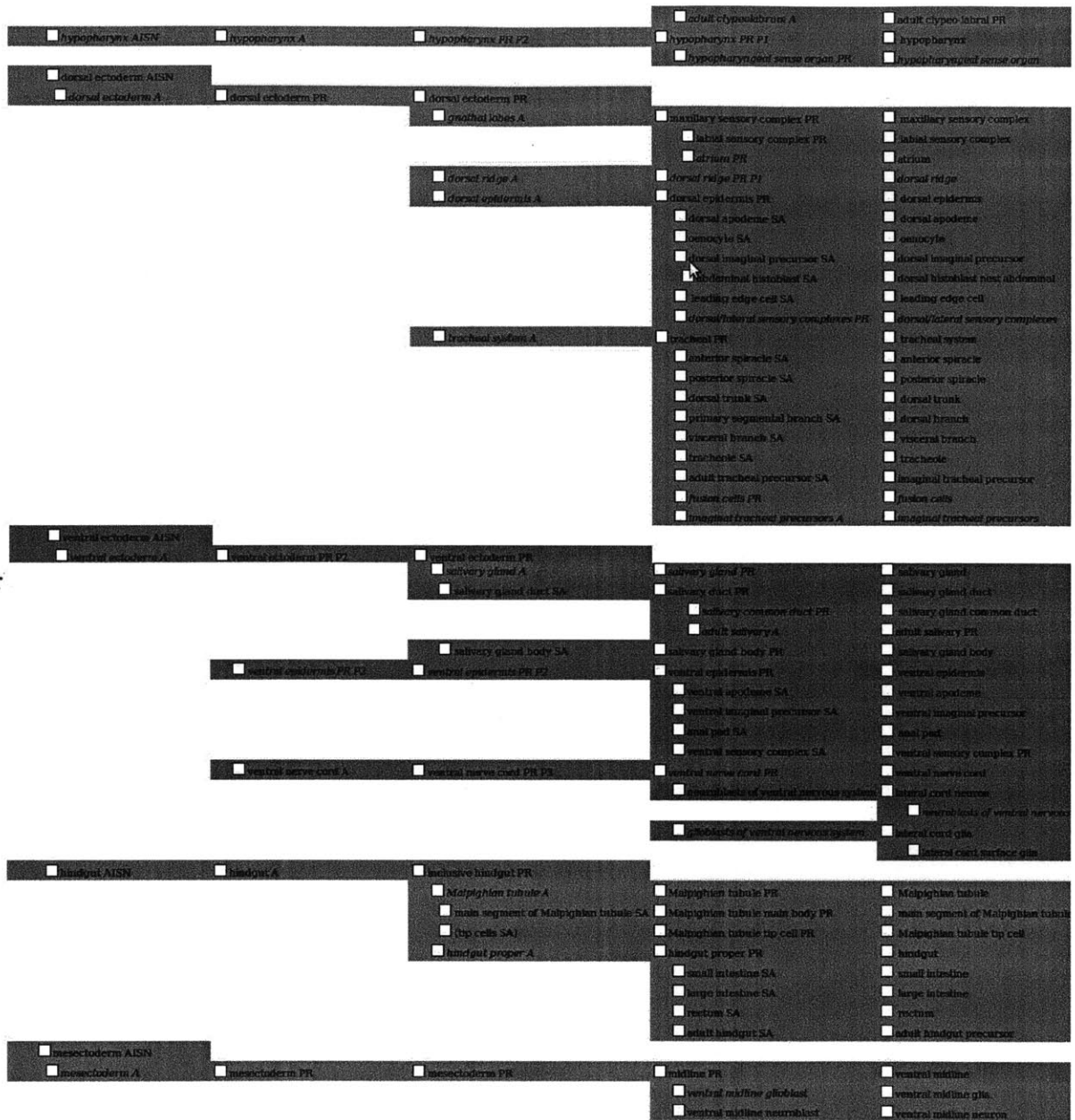
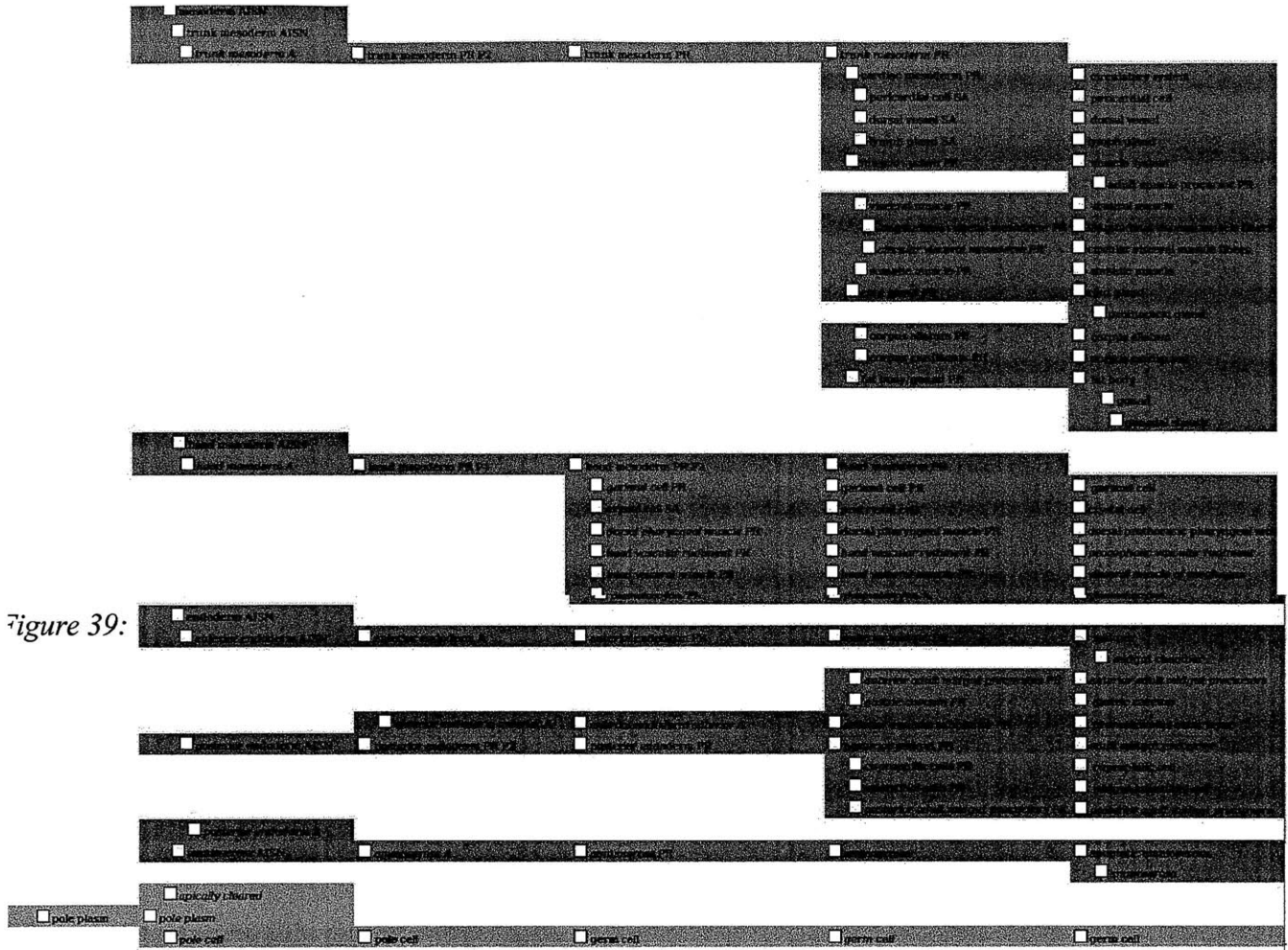


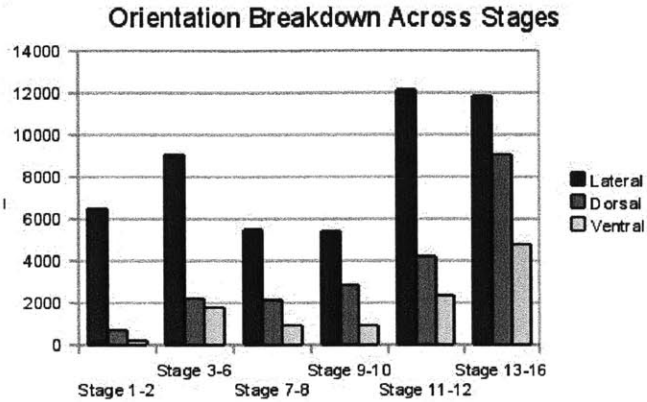
Figure 38:

6 Appendix



6 Appendix

6.2 BDGP Data Breakdown



Stage 7-8

Term	Frequency
ventral ectoderm primord	4252
procephalic ectoderm anl	3925
ubiquitous	3614
posterior endoderm primc	3248
dorsal ectoderm primordi	2654
no staining	2068
hindgut anlage	1814
faint ubiquitous	1160
foregut anlage	971
visual anlage	783
pole cell	734
amnioserosa anlage	670
yolk nuclei	596
mesectoderm primordium	595
clypeolabrum anlage	407
hypopharynx anlage	250
central brain anlage	219
strong ubiquitous	149
head epidermis dorsal an	134
head epidermis lateral A	114
antennal primordium anla	71
ventral epidermis primorc	46
procephalic ectoderm pri	28
amnioserosa	9
dorsal ectoderm	6
proventriculus anlage	4
anterior endoderm anterio	3
esophagus anlage	2

Stage 1-2

Term	Frequency
maternal	13718
no staining	6595
ubiquitous	945
pole plasm	587
strong ubiquitous	179
faint ubiquitous	56

Stage 3-6

Term	Frequency
cellular blastoderm	6672
ubiquitous	6003
dorsal ectoderm anlage in statu nascendi	5780
procephalic ectoderm anlage in statu nascendi	5277
ventral ectoderm anlage in statu nascendi	4874
no staining	2569
posterior endoderm anlage in statu nascendi	2173
mesoderm anlage in statu nascendi	2118
yolk nuclei	2018
hindgut anlage in statu nascendi	1933
anterior endoderm anlage in statu nascendi	1820
foregut anlage in statu nascendi	1678
amnioserosa anlage in statu nascendi	1630
visual anlage in statu nascendi	1611
faint ubiquitous	1556
ectoderm anlage in statu nascendi	1423
trunk mesoderm anlage in statu nascendi	1307
mesectoderm anlage in statu nascendi	1244
head mesoderm anlage in statu nascendi	1191
endoderm anlage in statu nascendi	1029
pole cell	669
clypeolabrum anlage in statu nascendi	603
ventral ectoderm anlage	597
dorsal ectoderm anlage	573
head epidermis dorsal anlage in statu nascendi	476
trunk mesoderm anlage	334
yolk	309
head epidermis anlage in statu nascendi	292
strong ubiquitous	223
hypopharynx anlage in statu nascendi	177
embryonic cortex	169
head mesoderm anlage	163
central brain anlage in statu nascendi	76
antennal primordium anlage in statu nascendi	53
mesectoderm anlage	47
posterior endoderm anlage	47

6 Appendix

Stage 9 to 10

Term	Frequency
procephalic ectoderm primordium	4489
trunk mesoderm primordium	3885
posterior endoderm primordium	3589
ventral nerve cord primordium P3	3462
head mesoderm primordium P2	3341
anterior endoderm primordium	3173
ubiquitous	3109
ventral ectoderm primordium	2951
inclusive hindgut primordium	2359
dorsal ectoderm primordium	2346
no staining	2274
foregut primordium	1392
faint ubiquitous	1247
visual primordium	1144
mesectoderm primordium	879
central brain primordium P3	701
yolk nuclei	672
amnioserosa primordium	562
amnioserosa	513
clypeolabrum primordium P2	470
germ cell	446
hypopharynx primordium P2	263
crystal cell specific anlage	255
dorsal epidermis anlage	226
gnathal lobes anlage	216
antennal primordium2	176
head epidermis lateral primordium P2	152
ventral epidermis primordium P2	148
strong ubiquitous	119
head epidermis dorsal primordium P4	118
salivary gland body specific anlage	113
segmentally repeated	110
tracheal system anlage	101
plasmatocytes primordium	77
proventriculus primordium	51
stomatogastric nervous system primordium P2	50
hindgut proper anlage	40
main segment of Malpighian tubule specific anlage	30
Malpighian tubule anlage	29
head epidermis ventral P2	28
labial complex specific anlage	24
maxillary organ specific anlage	24
esophagus primordium P2	22
proventriculus primordium P2	16
garland cell primordium	15
anterior endoderm anterior anlage	13
salivary gland anlage	12
salivary gland duct specific anlage	12
dorsal pharyngeal muscle primordium	9
dorsal ridge anlage	8
head visceral muscle primordium	6

6 Appendix

labral sensory complex primordium	543
procephalon primordium	537
embryonic central brain glia	468
embryonic leading edge cell specific anlage	447
lateral cord glia	434
hypopharynx primordium P1	415
lateral cord glia	398
lateral cord neuron	393
fat body specific anlage	378
ventral midline glia	372
pericardial cell specific anlage	347
salivary duct primordium	333
dorsal apodeme specific anlage	292
need new term(larval eye primordium)	280
germ cell	272
dorsal trunk specific anlage	263
head epidermis dorsal primordium P1	243
rectum specific anlage	231
visceral branch specific anlage	227
ventral apodeme specific anlage	225
head visceral muscle primordium	198
ventral midline neuroblast	190
lateral cord	182
Malpighian tubule main body primordium	169
segmentally repeated	163
esophagus primordium P1	141
lymph gland specific anlage	135
visual primordium2	132
epipharynx primordium P1	129
circular visceral mesoderm primordium	128
oenocyte specific anlage	128
embryonic inner optic lobe primordium	117
embryonic central brain	116
ventral imaginal precursor specific anlage	111
primary segmental branch specific anlage	109
glioblasts of ventral nervous system	106
lateral cord	104
hypopharyngeal sense organ primordium	98
strong ubiquitous	93
dorsal imaginal precursor specific anlage	90
corpus allatum primordium	75
dorsal ridge primordium P1	70
salivary common duct primordium	69
embryonic central brain surface glia	68
procephalic glioblasts	56
anterior spiracle specific anlage	56
head epidermis lateral primordium P1	55
lateral cord surface glia	55
anterior adult midgut precursors primordium	51
corpus cardiacum primordium	48

6 Appendix

head epidermis ventral P1	46
adult hindgut specific anlage	43
cuprophilic cells primordium	35
head vascular rudiment primordium	34
posterior adult midgut precursors primordium	34
adult clypeolabrum anlage	26
ring gland primordium	25
neuroblasts of central brain	24
interstitial cells primordium	16
tracheole specific anlage	15
Malpighian tubule tip cell primordium	13
dorsomedial neurosecretory cell	12
mushroom body primordium	11
embryonic central brain	10
esophageal ganglion primordium	8
embryonic dorsal pouch primordium	5
embryonic outer optic lobe primordium	5
gastric caecum primordium	3

6 Appendix

Stage 11-12

Term	Frequency
posterior midgut primordium	11064
anterior midgut primordium	9906
brain primordium	9413
dorsal epidermis primordium	7439
hindgut proper primordium	7177
ventral nerve cord primordium	6427
ventral epidermis primordium	5762
ubiquitous	5074
foregut primordium	4594
trunk mesoderm primordium	4189
tracheal primordium	4068
salivary gland body primordium	3926
clypeo-labral primordium	3021
somatic muscle primordium	2995
visceral muscle primordium	2783
faint ubiquitous	2610
amnioserosa	2603
plasmacytes anlage	2509
dorsal pharyngeal muscle primordium	2284
midline primordium	2225
no staining	2140
yolk nuclei	2088
fat body/gonad primordium	1960
head mesoderm primordium	1686
muscle system primordium	1602
posterior spiracle specific anlage	1498
procrystal cell	1412
stomatogastric nervous system primordium	1387
longitudinal visceral mesoderm primordium	1223
proventriculus primordium	1223
garland cell primordium	1167
atrium primordium	1067
gnathal primordium	1063
head epidermis primordium P1	1050
ventral sensory complex specific anlage	1035
sensory nervous system specific anlage	855
maxillary sensory complex primordium	838
cardiac mesoderm primordium	835
germ cell	792
labial sensory complex primordium	788
lateral cord neuron	773
salivary gland primordium	744
embryonic central brain neuron	734
pars intercerebralis primordium	724
dorsal/lateral sensory complexes primordium	711
antennal primordium 1	695
embryonic ganglion mother cell	624
anal pad specific anlage	553

6 Appendix

Stage 13-16					
Term	Frequency				
		embryonic salivary gland d	1376	embryonic main segment o	360
		embryonic antennal sense	1321	visceral muscle of esophag	360
ventral nerve cord	18507	embryonic maxillary senso	1299	dorsal imaginal precursor	328
embryonic dorsal epiderm	15923	embryonic salivary gland b	1189	embryonic Malpighian tubul	328
embryonic midgut	15608	embryonic rectum	1162	lateral cord neuron	288
embryonic brain	13852	embryonic/larval visceral b	1133	cat_term:common_C	270
embryonic ventral epiderm	13760	labral sensory complex	1095	embryonic central brain mu	270
embryonic hindgut	12435	ventral imaginal precursor	1050	ventral midline neuron	270
embryonic head epidermi	10013	embryonic labial sensory c	1039	adult foregut precursor	258
embryonic/larval muscle	8722	embryonic Bolwig's organ	1027	embryonic tracheole	248
embryonic foregut	8146	embryonic/larval circulatory	1010	adult muscle precursor prim	178
embryonic/larval tracheal	7480	embryonic optic lobe	990	dorsal ridge	138
embryonic central nervou	6845	embryonic salivary gland c	979	head epidermis lateral	128
embryonic proventriculus	6463	embryonic midgut interstitia	979	embryonic esophageal gang	118
embryonic hypopharynx	6057	embryonic/larval pericardial	960	embryonic hypocerebral ga	98
sensory system head	5982	ventral midline glia	913	dorsal metathoracic disc	98
embryonic epipharynx	5784	embryonic proventriculus ir	894	dorsal mesothoracic disc	98
embryonic/larval posterio	5399	cat_tem:endo_meso	880	head epidermis ventral	88
ventral midline	5343	embryonic foregut sensory	879	adult clypeo-labral primordi	88
embryonic/larval fat body	4871	no staining	867	lateral cord	88
embryonic anal pad	4852	embryonic leading edge ce	854	dorsal histoblast nest abdo	78
dorsal prothoracic pharyn	4828	adult hindgut precursor	828	embryonic frontal ganglion	78
embryonic/larval somatic	4651	lateral cord surface glia	810	embryonic inner optic lobe	78
embryonic/larval visceral	4469	cat_term:true_ubiquitous	717	macrophage	78
gonad	4242	embryonic corpus allatum	709	adult salivary primordium	68
atrium	4081	lateral cord glia	708	strong ubiquitous	68
ring gland	3840	embryonic small intestine	706	segmentally repeated	68
embryonic esophagus	3829	gonadal sheath	701	hypopharyngeal sense organ	58
ubiquitous	3795	embryonic central brain sur	689	embryonic outer optic lobe	58
ventral sensory complex	3713	embryonic central brain par	652	neuroblasts of central brain	58
sensory nervous system	3419	apoptotic amnioserosa	636	neuroblasts of ventral nerve	48
embryonic salivary gland	2672	embryonic/larval anterior sp	630	lateral cord	48
plasmacytes	2634	embryonic central brain glia	629	list_13-16	48
embryonic/larval garland	2593	amnioserosa	606	visual system	38
lymph gland	2524	longitudinal visceral muscle	594	procephalon	38
dorsal/lateral sensory cor	2520	embryonic proventriculus o	580	embryonic central brain	18
faint ubiquitous	2506	circular visceral muscle fib	574	procephalic vascular rudime	18
lateral cord glia	2379	embryonic proventriculus ir	503	posterior adult midgut precu	18
crystal cell	2309	clypeolabrum	496	embryonic proventricular ga	8
yolk nuclei	2241	adult midgut precursor	481	imaginal tracheal precursor	8
embryonic central brain n	2117	embryonic/larval garland ce	451		
head epidermis dorsal	2111	amnioserosa	450		
embryonic/larval oenocyt	2018	embryonic/larval dorsal bra	439		
lateral cord neuron	1922	embryonic corpus cardiacu	419		
embryonic central brain g	1892	embryonic cuprophilic cell	409		
embryonic/larval dorsal tr	1741	embryonic central brain neu	406		
germ cell	1690	adult eye primordium	396		
embryonic midgut chamb	1668	prothoracic gland	393		
embryonic gastric caecur	1418	embryonic central brain	383		
embryonic stomatogastric	1382	lateral cord	367		

7 Bibliography

7 Bibliography

- 1: Ji S, Sun L, Jin R, Kumar S, Ye J. (2008) Automated annotation of *Drosophila* gene expression patterns using a controlled vocabulary. *Bioinformatics*. 2008 Sep 1;24(17):1881-8. Epub 2008 Jul 16. PubMed PMID: 18632750; PubMed Central PMCID: PMC2519157.
- 2: Ji S, Li YX, Zhou ZH, Kumar S, Ye J. A bag-of-words approach for *Drosophila* gene expression pattern annotation. (2009) *BMC Bioinformatics*. 2009 Apr 21;10:119. PubMed PMID: 19383139; PubMed Central PMCID: PMC2680406.
- 3: Peng H, Long F, Zhou J, Leung G, Eisen MB, Myers EW. (2007) Automatic image analysis for gene expression patterns of fly embryos. *BMC Cell Biol*. 2007 Jul 10;8 Suppl 1:S7. PubMed PMID: 17634097; PubMed Central PMCID: PMC1924512.
- 4: Zhou J, Peng H. (2007) Automatic recognition and annotation of gene expression patterns of fly embryos. *Bioinformatics*. 2007 Mar 1;23(5):589-96. Epub 2007 Jan 19. Erratum in: *Bioinformatics*. 2007 Jun;23(12):1571-2.. PubMed PMID: 17237064.
- 5: Ye J, Chen J, Li Q, Kumar S. (2006) Classification of *Drosophila* embryonic developmental stage range based on gene expression pattern images. *Comput Syst Bioinformatics Conf*. 2006:293-8. PubMed PMID: 17369647.
- 6: Gurunathan R, Van Emden B, Panchanathan S, Kumar S. Identifying spatially similar gene expression patterns in early stage fruit fly embryo images: binary feature versus invariant moment digital representations. (2004) *BMC Bioinformatics*. 2004 Dec 16;5:202. PubMed PMID: 15603586; PubMed Central PMCID: PMC545963.
- 7: Kumar S, Jayaraman K, Panchanathan S, Gurunathan R, Marti-Subirana A, Newfeld SJ. (2002) BEST: a novel computational approach for comparing gene expression patterns from early stages of *Drosophila melanogaster* development. *Genetics*. 2002 Dec;162(4):2037-47. PubMed PMID: 12524369; PubMed Central PMCID: PMC1462359.
- 8: Samsonova AA, Niranjana M, Russell S, Brazma A. (2007) Prediction of gene expression in embryonic structures of *Drosophila melanogaster*. *PLoS Comput Biol*. 2007 Jul;3(7):e144. PubMed PMID: 17658945; PubMed Central PMCID: PMC1924873.
- 9: Frise E, Hammonds AS, Celniker SE. Systematic image-driven analysis of the spatial *Drosophila* embryonic expression landscape. (2010) *Mol Syst Biol*. 2010;6:345. Epub 2010 Jan 19. PubMed PMID: 20087342; PubMed Central PMCID: PMC2824522.
- 10: Mace DL, Varnado N, Zhang W, Frise E, Ohler U. (2009) Extraction and comparison of gene expression patterns from 2D RNA in situ hybridization images. *Bioinformatics*. 2010 Mar 15;26(6):761-9. Epub 2009 Nov 26. PubMed PMID: 19942587.
- 11: Tomancak P, Berman BP, Beaton A, Weizmann R, Kwan E, Hartenstein V, Celniker SE, Rubin GM. (2007) Global analysis of patterns of gene expression during *Drosophila* embryogenesis. *Genome Biol*. 2007;8(7):R145. PubMed PMID: 17645804; PubMed Central PMCID: PMC2323238.
- 12: Costa IG, Krause R, Opitz L, Schliep A. (2007) Semi-supervised learning for the identification of syn-expressed genes from fused microarray and in situ image data. *BMC Bioinformatics*. 2007;8 Suppl 10:S3. PubMed PMID: 18269697; PubMed Central PMCID: PMC2230504.
- 13: Tomancak P, Beaton A, Weizmann R, Kwan E, Shu S, Lewis SE, Richards S, Ashburner M, Hartenstein V, Celniker SE, Rubin GM. (2002) Systematic determination of patterns of gene expression during *Drosophila* embryogenesis. *Genome Biol*. 2002;3(12):RESEARCH0088. Epub 2002 Dec 23. PubMed PMID: 12537577; PubMed Central PMCID: PMC151190.

7 Bibliography

- 14: BDGP website. (2011) <http://insitu.fruitfly.org/cgi-bin/ex/insitu.pl>, May, 2011.
- 15: LabelMe website. (2011) <http://labelme.csail.mit.edu/>, May, 2011.
- 16: Russel B, Torralba A. (2008) LabelMe: a database and web-based tool for image annotation. *International Journal of Computer Vision*.
- 17: Adult *Drosophila melanogaster* image. <http://people.mokk.bme.hu/~daniel/Heni/regebbi/031%20%20Diptera/039%20%23%23%23%20ecetmuslica%20%28Drosophila%20melanogaster%29.jpg>, May 2011.
- 18: BDGP Controlled Vocabulary Image. http://www.fruitfly.org/ex/Annotation_files/image010.jpg. May 2011.
- 19: Apache HTTP Server Version 2.2 Documentation. <http://httpd.apache.org/docs/2.2/>. May 2011.
- 20: Apache Tomcat 6 Documentation. <http://tomcat.apache.org/tomcat-6.0-doc/index.html>. May 2011.
- 21: Apache Axis2 Documentation. <http://axis.apache.org/axis2/java/core/>. May 2011.
- 22: Purves, William K, et al. (2001) *Life: The Science of Biology*, 6th edition. Sinauer Associates, Inc., W. H. Freeman and Company, Sunderland, Massachusetts.
- 23: Duda, Richard O., Hart, Peter E., and Stork, David G. (2001) *Pattern Classification*, 2nd edition. John Wiley & Sons, Malden, Massachusetts.
- 24: Hartl, Daniel L. and Jones, Elizabeth W. (2002) *Essential Genetics, A Genomics Perspective*, 3rd edition. Jones and Bartlett Publishers, Sudbury, Massachusetts.
- 25: Brain MRI Image. http://3.bp.blogspot.com/_7JAArcms7_8/SKskm-HKgAI/AAAAAAAAAhc/o5MLUbfC-3I/s400/BrainMRI_Sagittal.jpg. May 2011.
- 26: Torralba, Antonio. Personal Communication. August 2010.
- 27: Amazon's Mechanical Turk. <https://www.mturk.com/mturk/welcome>. May 2011.
- 28: Adams MD, Celniker SE, Holt RA, *et al.* (2000). The genome sequence of *Drosophila melanogaster*. *Science* 287 (5461): 2185–95. doi:10.1126/science.287.5461.2185. PMID 1073113.
- 29: Greenberg, Julie. (2011) Laboratory Project 5, Post Processing Brain MRI Images. 6.555 Biomedical Signal and Image Processing Spring 2011.
- 30: Jaakkola, Tommi and Collins, Michael. (2010) 6.867 Machine Learning. Fall 2010.
- 31: Rifkin, Ryan M. and Lippert, Ross A. (2007) Notes on Regularized Least Squares. May 1, 2007.
- 32: Frogner, Charlie. (2011). Lecture 5 Notes on Regularized Least Squares. 9.520 Statistical Learning Theory and Applications Spring 2011.
- 33: Frogner, Charlie. (2011) Lecture 6 Notes on Support Vector Machines. 9.520 Statistical Learning Theory and Applications Spring 2011.