

# CSE 6242 Project: Improving the Vectorbase Bioinformatic Data

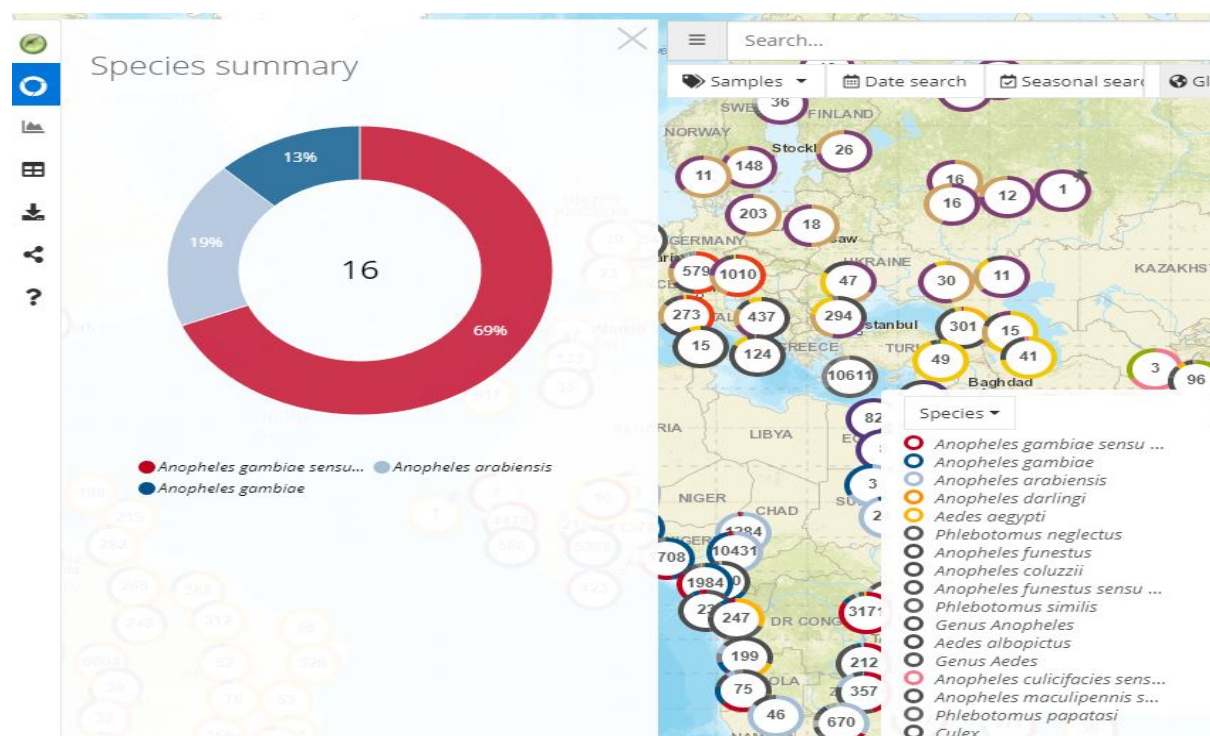
Jon Gerhart, Will Hutwagner, Bridget Neary, Mani Jain, Joey Gonzales

## Introduction - Motivation

New molecular analysis technologies, along with the wide availability of increasing quantities of biological data have allowed researchers an unprecedented view of vector-pathogen-host interactions involved in the transmission of many infectious diseases<sup>[1]</sup>. VectorBase is a resource for studying vectors of human pathogens, with genomes and other content for over 35 organisms.<sup>[2]</sup> VectorBase has the stated goal that its valuable resource is one of several Bioinformatics Resource Centers supported by the NIH, underpinned by community contributions and is crucial in the progress of this field. We hope to collaborate with VectorBase to further their stated goals of offering users open access and interfaces enabling meaningful application and exploration of this data<sup>[4]</sup>.

## Problem definition

We aimed to streamline the user interface and data curation process and to develop interactive visualization tools to help with hypothesis formulation in medical entomology. VectorBase uses Esri Leaflet to visualize ArcGIS data collected during the normal course of experimental and control efforts. VectorBase uses a visually busy navigation system for this service, in large part due to the high dimensionality and sparseness of the data.





**Figure 1:** Example of current VectorBase population browser UI

## Survey

1. This navigation service would benefit from streamlining and organization into a unified system. In addition, VectorBase uses violin plots to visualize large sets of phenotypic data, which obfuscates some of the key desired dimensions of the data request and merits redaction of some data.
2. Currently, many of the gene annotation links in VectorBase are deprecated which reduces its usability as a unified reference to other resources. The NCBI refseq database will be updated to retrieve updated gene annotations with the help of NCBI's Entrez system, by using Entrez programming utilities (E-Utilities and Entrez Direct) and bulk transfer via FTP. Further, we plan to fully automate this process and provide it to VectorBase<sup>[7]</sup>.
3. The development of a visualization/clustering tool for quick meta-analysis of large sets of data, such as large sets of genes. We can use statistical algorithms to arrange these data into similar expressions and view them via a heat map<sup>[8]</sup>. We plan to take advantage of machine learning aiming to produce subgroups of related points. These subgroups can be organized into dendrograms for analysis (neighbor-joining tree)<sup>[9]</sup>. We can find new genes for a specific phenotype by developing a text mining system to look for the relevant information in Gene Ontology, PUBMED, and other databases<sup>[10]</sup>.

We will need general software, algorithmic, and visualization-related components to achieve the three major aims of the project. For this, we can draw inspiration from BIRDS, an international project to advance the fields of bioinformatics and information retrieval by developing new data structures and algorithms. The scope of our project is much broader than the state-of-the-art techniques developed by the BIRDS project may inspire us to develop VectorBase using similar techniques<sup>[11]</sup>. We can also explore a set visualization tool similar to the OnSet system by Sadana et al. to analyse the degree of similarity between clusters<sup>[12]</sup>. Libbrecht and Noble also provide a useful overview of machine learning techniques as used in genetics and genomics that we think will support our aims well<sup>[13]</sup>.

## Proposed method

Intuition (*Why should it be better than the state of the art?*)

We anticipate that this to be a critical improvement for the field of vector biology. VectorBase is a powerful tool for biologists and the most widely used resource on the internet. Its user interface requires improvements to make the most of the existing infrastructure. By streamlining the navigation and visualization of the data, we will be able to provide a more efficient and powerful tool for biologists.

expanding the capabilities of the infrastructure will allow the resource to continue to be up-to-date. Harnessing these two improvements we will be able to implement a new tool that allows users to find biologically relevant insights to their data based on text annotations from various bioinformatics resources and clustering methods that improve on the current data presentation tool that has static clustering per species. By automating data curation and streamlining presentation we improve the impact of the service. Rather than all VectorBase sequences being organized based on user requested subsets of submitted annotations, the ma

and clustering algorithms we plan to implement can expose unexpected correlations and make excellent targets for further research.

### **Aim 1:** Improve existing visualization in vector base

Our updated UI for VectorBase required modifying the existing JavaScript implementation, which uses the Esri Leaflet API for the interactive map and D3.js for data visualization. Currently, when a filtered data selection is made, VectorBase displays it as a scatterplot, which switches to a violin plot if the selection is larger than a certain threshold. VectorBase reports that their users are confused by the violin plots, and would like to implement a single visualization style that is easy to read for datasets of any size.

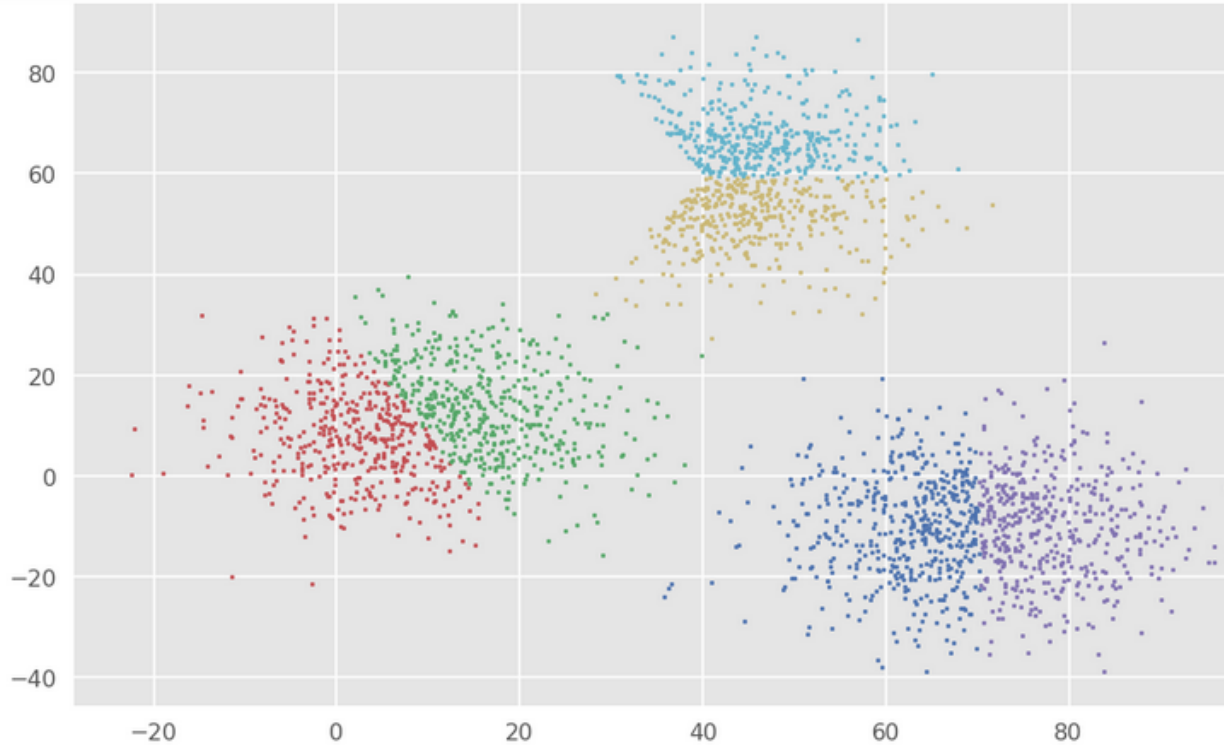
### **Aim 2:** Automate gene data curation in vectorbase

VectorBase users reported trouble accessing links to genes recently updated because VectorBase does not have an automated way to update HTML links when they are changed by their counterparts. We used lxml and a regex search system to retrieve all nucore hyperlinks and scraped for any hyperlinks in the body that redirected to the new location and output a master list containing any newer links. This implementation is functionally completed.

### **Aim 3:** Visualizing the relationship between expressed genes

We extracted expressed genes from Vector Base database and used k-means to cluster genes based on their common features that we extracted from NCBI, EMBL, and GenBank (curated bioinformatics databases). We choose Gene ontologies and molecular pathways as the number of clusters was the hyperparameter that we tuned and experimented with. We used a dataset of 2500 genes. The optimal number of clusters that gave us good results is 10 clusters.





**Figure 2:** Clustering of expressed genes chosen in the pilot dataset with 6 number of clusters

With the improved population biology map, User selects from a list of co species and enters a list of comma separated Ensembl gene names. We use the B obtain gene ontologies for each gene. We use these ontologies to query the same hierarchy of ontologies that range from specific to broad. This hierarchy is then t and relatedness between the list of genes is determined as the number of edges th leaves. This relatedness is visualized as a force directed graph where each node i each edge represents the existence a common ontology term, weighted for the dis lowest common ancestry. A python based user interface web-server was made fo component using Flask. The interface contains a scroll-bar option for species dat text bar for list of genes.

After retrieving data relevant to the gene from the three databases and their class plan to evaluate each of these classes using machine learning algorithms based o distance such as FuzzyWuzzy<sup>[17]</sup> and using similarity computed from n-gram app we will prepare our own separate small reference set with words having biologic selective terminology. We intended to implement a fuzzy text matching algorithm appearance of each word based on frequencies and matches with the reference se enable us to deliver biologically significant and meaningful results to the user.

enable us to deliver biologically significant and meaningful results to the user.

We then implement a scraper to ensure that the texts are indeed related to the gene of interest. We filter out texts with low word count, include beneficial language with biochemical and cellular terminology, and remove non-beneficial language such as “unannotated” and “hypothetically.” One last quality analysis is performed on the list of texts ensuring that the texts are acceptable and possibly improving the accuracy of the machine learning model to find keys or patterns to improve upon the mining and



Finally, we present our gene data through an interactive visualization implemented in JavaScript/D3.

## **Approaches**

We divided the undertaking into three aims, which will be evaluated and separately:

### **Aim 1:** Improve existing visualization in vector base

VectorBase described issues with visualizing large numbers of samples from their browser, when there are multiple species being visualized at the same time. In the vectorbase code, all the data points are collapsed in a single violin plot, which is simple, but loses information relevant to the user's query. We intend to refactor the code to include a collection of stackable, color coded overlapping fields. Based on our experience, we believe this will be more interpretable.

The map supports different views of the data, from which the user may select over open secondary visualizations. VectorBase reported that their component for viewing resistance statistics for selected samples did not clearly communicate data on growth phenotypes and metrics at the same time, so we developed a mockup component that better supports simultaneous visualization. The example of that D3 code is visualized



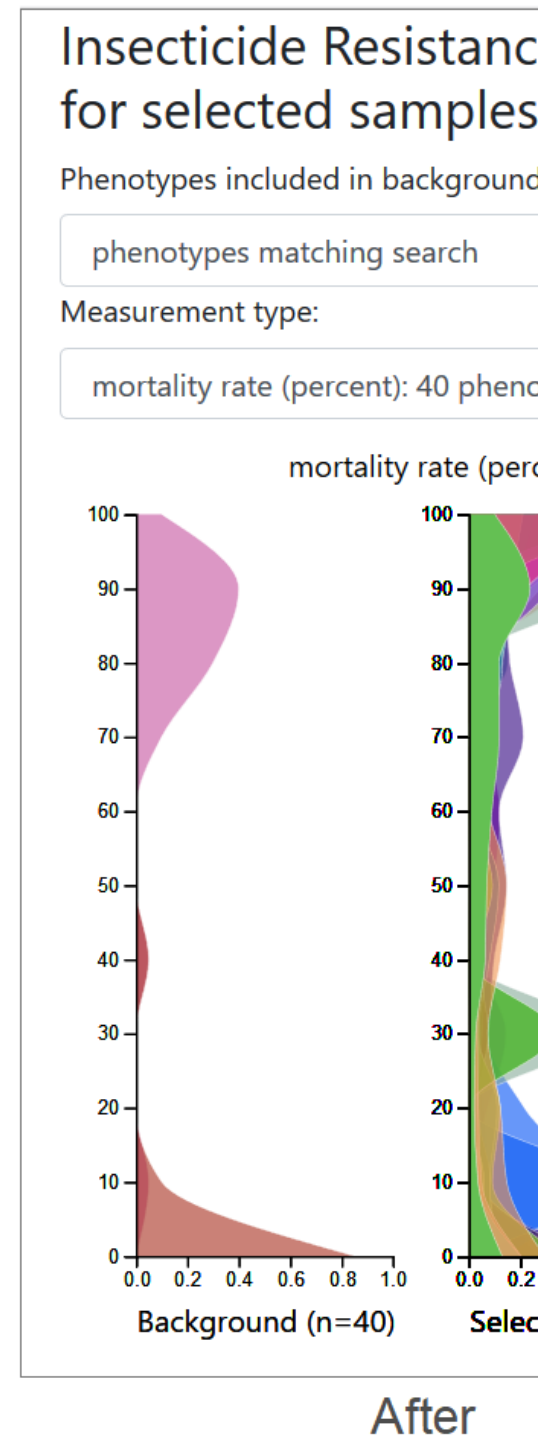
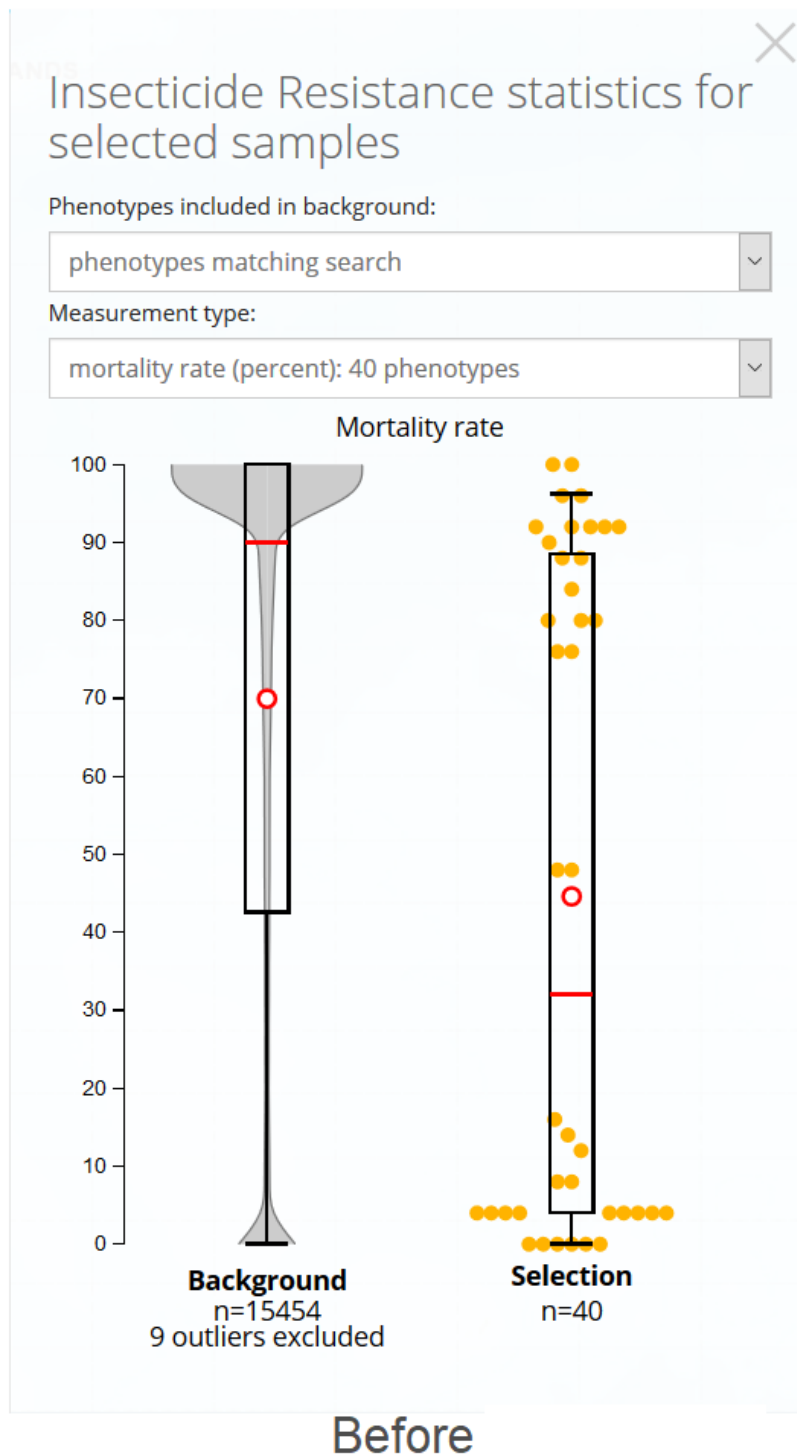


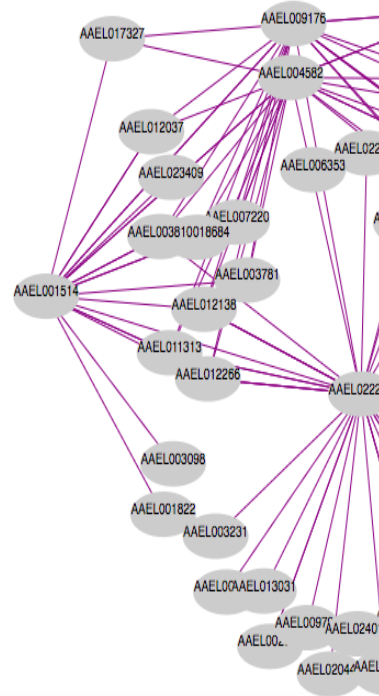
Figure 3: Visualization before and after implementation in D3.

## Aim 2: Automate gene data curation in VectorBase

VectorBase users reported trouble accessing links to genes recently updated by NCBI. VectorBase does not have an automated way to update HTML links when

changed by their counterparts. We coded a simple script that checks and updates the links in VectorBase pages when run.

**Aim 3:** Combine lessons learned from aims 1 and 2 into a working applet (NetGene) that provides summaries of gene relations from NGS data.

[home](#)

## Experiments/ Evaluation

*Description of testbed:*

Aim 2 outcomes are self-evident.

Aim 3 will be evaluated via surveys of user experiences.

### Details of the experiments:

Aim 1 was not completed until April 25th due to difficulties with implementation of Aim 1 proved to be more challenging than anticipated, so only a static data selection was completed. Simulations of Aim 1 received favorable responses from entomologists interviewed. Aim 2 has been trialed successfully, but Vectorbase was delayed in implementation due to a recent database update. Aim 3 is not yet feature complete due to experimentation, due to delays in Aim 1 information and talent transfer, but multi-

experimentation in visualization techniques, edge tracing and culling, and interaction led to a solid concept that could be further refined into a competitive data parsing

Experiments in aim 3 were in 3 phases. The first phase was identifying the measure distance between genes based on GO terms. Originally, we wanted to use the ontology hierarchy according to relationships between the ontology terms, but this was unavailable, so each edge was given a weight of 1. We also experimented on finding distances between genes in this hierarchy. Initially we used the Floyd-Warshall algorithm for finding shortest paths between all genes. This was too computationally expensive to implement in real time, so we devised a method of comparing lists of parent terms and identifying the lowest one for each gene pair, if present. This cut the computation

The second phase of experimentation for aim 3 was defining an edge weight that was informative for the user to define the appropriate cutoff for their gene set in terms of as an edge. This phase is still in development but the gene distances will likely be as the distribution tends to be log-normal. Ongoing experimentation also involves clustering algorithms. We validated our graph clustering using networkx clustering

### *Conclusions and discussion:*

Due to personnel issues, we did not complete aim 1 within a reasonable time frame for the full functionality. In addition, the aim 1 delays cascaded into the development of the front-end and interactivity portions. Further, aim 2 testing was delayed to a point where we went through a blanket update of all hyperlinks instead of automating the process. This work is purely academic. Aim 3 represents an interesting approach to organizing the presentation of gene IDs by ontology, but is not currently feature complete, missing interactivity, interface and clustering components. The NetGO team regrets this

### **Distribution of team member effort**

Jonathan Gerhart coordinated with the VectorBase team, developed project reports, vectored and designed poster, designed, tested and wrote aim 2, managed and contributed to modules 5 and 6 design phase.

Bridget Neary coordinated team efforts, wrote key elements of aim 3, coordinated phases of development of aim 3, assisted in experimental design, and proofread reports

reports.

Will Hutwagner drafted, developed and debugged the web server and API components of aim 3, coordinated aim 3 development and provided logistical support for the poster presentation.



Mani Jain drafted, tested and reviewed visualization techniques for aim 3 all phases of debugging, contributed significant insight to attempts to implement proofread reports and communications with VectorBase.

Joey Gonzales-Dones implemented the VectorBase PopBio insecticide re for aim 1 and tried to provide general assistance with other tasks where possible. only team member that is not a bioinformatics student, and thus found it more ch overwhelming to contribute to the project, so he wasn't able to contribute as much as the other team members.

