# Non-model arthropod assembly: past, present and future

Scott Emrich
On behalf of VectorBase and i5K

Eck Institute for Global Health
University of Notre Dame

## VectorBase
Bioinformatics Resource for Invertebrate Vectors of Human Pathogens

# Brief overview

- Brief self introduction

- The VectorBase BRC

- "10 simple rules for a successful genome project" – ode to PlOS series

# DIRTY ROTTEN BUGS?

## Arthropods Unite to Tell Their Side of the Story

NOT DIRTY!
NOT ROTTEN!
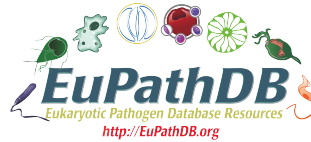Get to Know us!

Written and illustrated by **GILLES BONOTAUX**

DON'T LET
THIS BAD BUG
BITE YOU

# A brief introduction to BRCs



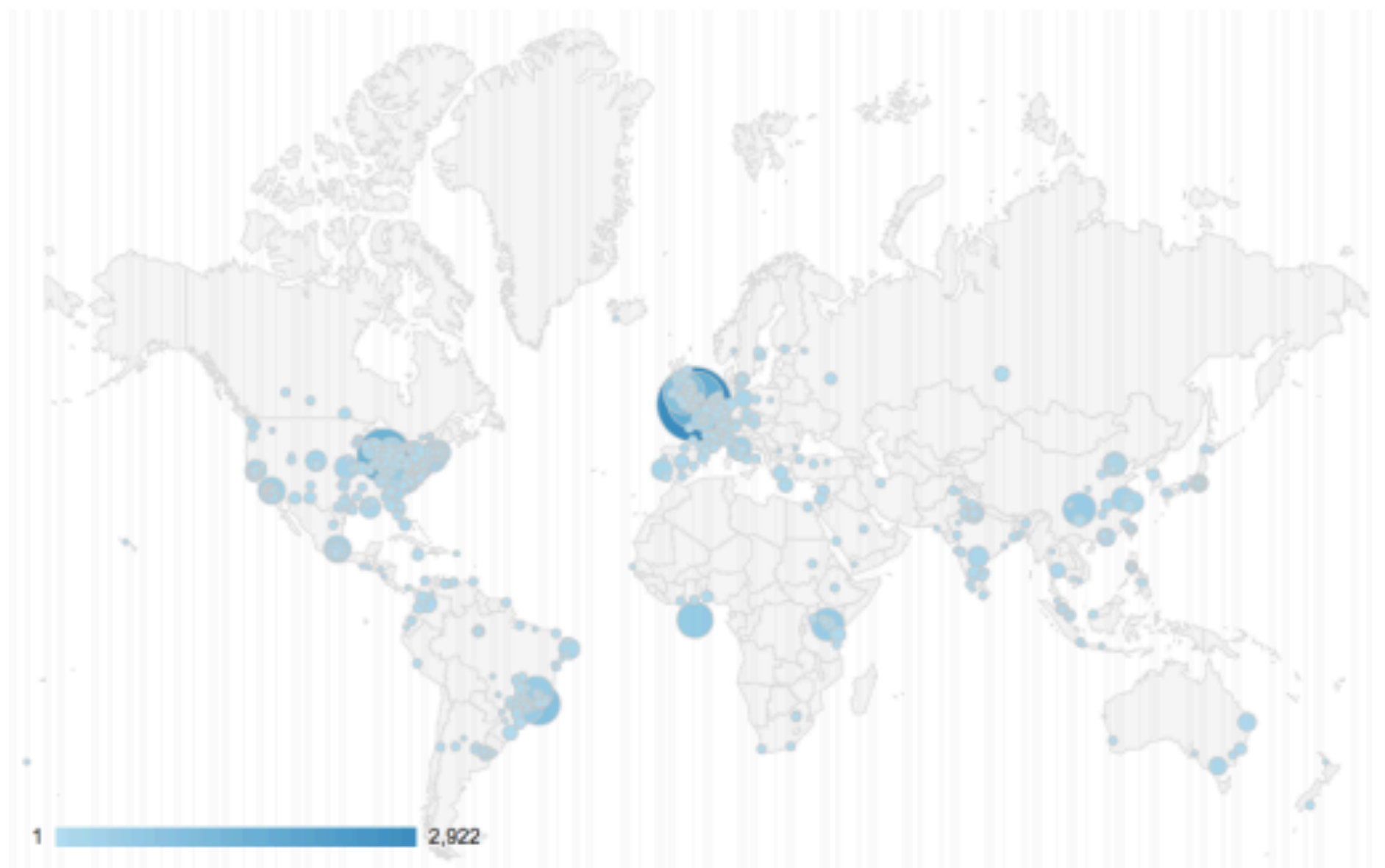- VectorBase is a genome resource for invertebrate vectors of human pathogens

- Funded by NIH-NIAID as part of a wider group of NIAID BRCs (see above) for biodefense and emerging and re-emerging infectious diseases

- Third contract started Fall 2014 (for up to 5 more years)

# VectorBase Scientific User Community by Loci



VectorBase Scientific User Community by Loci

# Tools

## Biomart
Use for (small and big scale) data mining queries that are not as easy or even possible to do using VectorBase Search

## BLAST
Finds regions of local similarity between sequences. Available data sets include contigs, scaffolds, chromosomes ESTs, RNAseq, transcripts and peptides.

## ClustalW
Can be used to generate input files for HMMER. After running a job just click on the link "Send to HMMER".

## Expression Browser and Map
Currently hosting microarrays mostly from *An. gambiae* and *Ae. aegypti*. Data from different publications is processed through the same pipeline so that results can be compared side-by-side.

## Galaxy
Galaxy is an open, web-based platform for data intensive biomedical research.

## Genome Browser
Makes genomic data accessible. Data is not only the genome sequence itself, but also other features such as comparisons between species including *in silico* and experimental data.

## HMMER
It looks for homolog genes, but unlike BLAST it aims to be more accurate and better to detect remote homologs. Input file is a protein multiple sequence alignment (MSA) from ClustalW.

## Ontology Browser
Ontologies are the structural framework for organizing information and are used in the Expression Browser and PopBio. You can also use the Ontology Browser for your research.
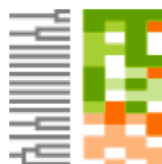
## Population Biology (PopBio)
Is part of our ongoing efforts to integrate genomic, phenotypic (including insecticide resistance) and population data (including SNPs).

## Web Apollo
Is an instantaneous, collaborative, genome annotation editor. Web Apollo is designed to support geographically dispersed researchers.

# Community gene annotation using WebApollo

Pre Web
Apollo

Early 2014
Web Apollo 1

December 2015
Web Apollo 2



User Downloads and edits
GFF with text editor or
WebApollo.

Inconsistent GFF with
many errors.

Common application for all users. ✓

Need one instance per species.
Hard to access meta data.

One Instance for
all species ✓
Reports on demand ✓

36 species
>100 users

# Summary of Community Annotation



Splits (15%) — 725

New genes (15%) — 715

Merges (6%) — 263

Modifications (64%) — 3036

Total Events 4739
29 Species

- 1/3 of annotations are major gene-model changes (splits, merges)
- 2/3 of annotations are minor modifications
- 2/3 of modifications are to the CDS
- ½ of CDS modifications gain sequence
- Accrued relevant InterPro Domains.

# Searchable track Meta Data

Publications ID, External Database ID, Description ....



Several hundred RNASeq tracks with associated metadata have been automatically added to VectorBase

# RNAseq tracking and analysis

# Population Biology

Welcome to VectorBase's Population Biology (PopBio) resource: a database and associated tools for visualisation, search and analysis of a wide range of population data, including genotypes, insecticide resistance and other phenotypes, and field collection metadata. You can interactively query all geo-tagged data (>99% of samples) using the map interface. Text-based search and browse is also available.

| Collection sites map | Insecticide resistance map | Search | Browse projects |
|---|---|---|---|

Phenotypes, genotypes and assays can be searched, browsed, and viewed in special views (see above) as we collect enough data

Large external data in resource:
- 30,000 observations from the Malaria Atlas Project
- >200 individuals with high-throughput genotyping and metadata
- 5500 insecticide resistance assays, including data from the President's Malaria Initiative

Sample summary data

Anopheles funestus  Anopheles arabiensis
Anopheles gambiae  Anopheles merus

**SEARCH**

Auto-completing available for suggestions

Taxonomy and ontology-aware (search with higher level concepts, such as "aquatic environment catch" or

Simple logic – does what you expect it to do when searching for several species or insecticides

**VIEW MODES**

Samples view: basic collection metadata for all population biology data in VectorBase

IR phenotypes view: one data point for each measured insecticide resistance phenotype

More view modes are planned, including genotypes and other phenotypes.

An. messeae
An. funestus
An. gambiae
An. arabiensis
An. atroparvus
An. darlingi
An. albitarsis
An. stephensi
An. fluviatilis
An. nuneztovari
An. sacharovi
An. punctulatus
An. subpictus
An. minimus
An. lesteri
Ae. aegypti
An. culicifacies
An. quadrimaculatus
An. superpictus
Others

Species key (color coded)

VectorBase's PopBio resource contains insecticide resistance data from a range of assay protocols and reported in a variety of measures and units, such as percent mortality, lethal concentration (e.g. LC50) and lethal time (e.g. LT95).

To aid the user in discovering geographical regions of resistance we have rescaled all comparable data, and inverting value ranges where appropriate. These rescaled values are used to color the map markers (from blue to red).

# Fitness trade-offs in different habitats

## M-form (colluzzi)



- More permanent
- Available year-round
- Allows slower development
- Predator-rich

## S-form (gambiae)



- Ephemeral
- rainy-season dependent
- Requires rapid development
- Largely predator-free

# Ecological adaptation via reverse ecology

- Beneficial mutations should carry signatures

——— Divergence          ——— Diversity



Target of Positive Selection

# Some genomic regions display footprint of strong, recent selection



Lawniczak, Emrich et al. 2010 Science

TEP1rB is (almost) exclusive to M from West Africa

White et al. PNAS (2011)

# Integrating variation and metadata across taxa

Display of variant data in genomic context

# "10 simple rules"
# for a genome sequencing effort

- The i5K folks were most interested with these "nuggets;" note these are from my experience and therefore are my opinions

- Experience from my efforts is really the best teacher – luckily we have a great, very open community of researchers to lean on

# Rule #1:
## "Bioinformaticians are not 'alchemists'"

- The quality of the manuscript is usually directly associated with the quality of
  - The data
  - The underlying questions


- Arthropod assembly still degrades with the complexity of the sample/genome


- Long reads are better but not a panacea
  - We failed to assembly *An. gambiae*'s Y and probably will continue to fail with current technology

# Mostly the same methods, different results



A. freeborni        A. minimus        A. albimanus

| Species | Genome (bp) | N50 (bp) | # Scaffolds |
|---|---|---|---|
| *A. gambiae* PEST | 273,093,681 | 49,364,325 | 7* |
| *A. minimus* | 201,793,324 | 10,313,149 | 678 |
| *A. arabiensis* | 246,567,867 | 5,604,218 | 1,214 |
| ***A. funestus*** | **225,223,604** | **671,960** | **1,392** |
| *A. stephensi* | 212,639,700 | 4,319 | 125,197 |

# Rule #2:
# "A good genome is a useful genome"

- No matter what the assembly method/sequencing tech, our results and others indicates quality directly correlates with heterogygosity

- In Anopheles at least, simpler libraries and DISCOVAR *de novo* can assemble the gene regions OK (Love et al., 2016)

- The corollary is hybrid approaches don't work unless you bound variation
    - Our "protocol" from gambiae Y is to sequence full sibs

- Bandage plot of the Aedes cell line genome most recently added using PacBio (Mark Kunitomi; UCSF)

# Rule #3:
# "Don't be afraid to ask for help"

- Every genomics expert was at one time a novice

- We have a great community, and multiple venues exist to network/share. One coming up soon is the Arthropod Genomics Symposium to be held at Notre Dame

- Send emails as needed to your bioinformatics resource; we'd love to help!

# Rule #4:
## "Don't put all of your eggs in one basket"

- The cost is usually both time and resources, but the best "flagship" manuscripts look at multiple cool biological features

- Set up working groups in your community, with a respected leader and have them write 1-2 pages with at least one figure

- You can be surprised what you find!

# Rule #5:
## "Use deadlines to 'herd the cats'"

- One of the great aspects for me as a genome project veteran is how engaged people are on their favorite organism(s)

- Academics, though, are overcommitted and realized enthusiasm is variable

- Set a schedule and have the working groups present at some regular frequency; this ensures both results and a feedback loop

# Rule #6:
## "Teamwork can be everything"

- Some (most?) of my most successful projects involved a good friend/collaborator or two

- Examples
  - Drs. Mara Lawnziack and Alisha Holloway (M/S)
  - Dr. Matt Hahn and others (gambiae complex)
  - Drs. Scott Egan and Greg Ragland (Rhagaletis)
  - Drs. Igor Sharakov, Jake Tu, Adam Phillippy, others (gambiae Y)

# Rule #7: (KG)
# "Share everything"

- Everyone knows of a horror story, but the value of open data / sharing is much higher than hoarding

- Make your genomes publically available as soon as they are frozen for the community and try and combine forces

- Archive these in national resources (e.g., GenBank)

- Your publication quality per effort expended together is much better than your group alone

CURRENT ISSUE // ARCHIVE // NEWS & MULTIMEDIA // AUTHORS // ABOUT // COLLECTED ARTICLES // BROWSE BY TOPIC // EARLY EDITION // FRONT MATTER

> Early Edition > Andrew Brantley Hall, doi: 10.1073/pnas.1525164113

CrossMark
click for updates

# Radical remodeling of the Y chromosome in a recent radiation of malaria mosquitoes

Andrew Brantley Hall[a,1], Philippos-Aris Papathanos[b,c,1], Atashi Sharma[d,1], Changde Cheng[e,f,1,2], Omar S. Akbari[g], Lauren Assour[h], Nicholas H. Bergman[i], Alessia Cagnetti[b], Andrea Crisanti[b,c], Tania Dottorini[c], Elisa Fiorentini[c], Roberto Galizi[c], Jonathan Hnath[i], Xiaofang Jiang[a], Sergey Koren[j], Tony Nolan[c], Diane Radune[i], Maria V. Sharakhova[d,k], Aaron Steele[h], Vladimir A. Timoshevskiy[d], Nikolai Windbichler[c], Simo Zhang[l], Matthew W. Hahn[l,m], Adam M. Phillippy[j], Scott J. Emrich[e,h], Igor V. Sharakhov[a,d,k,3], Zhijian Jake Tu[a,n,3], and Nora J. Besansky[e,f,3]

## Don't Miss

Customize your free PNAS alerts to receive a notification as new content becomes available.

## Article Tools

- Article Alerts ▶
- Export Citation ▶
- Save for Later ▶

Magda D. Lugon[x], David Majerowicz[c,mm], Paula L. Marcet[kk], Marco M. ... gy[k], Ana C. A. Melo[a,b], Fanis Missirlis[nn], Theo Mota[oo], Fernando G. Norieg..., Raquel L. L. Oliveira[a], Gilbert Oliveira-Silveira[c], Sheila Ons[e], Lucia Pag... ual[e], Marcio G. Pavan[g], Nicolás Pedrini[v], Alexandre A. Peixoto[b,g], Marcos ..., Francisco Prosdocimi[c], Rodrigo Ribeiro-Rodrigues[qq], Hugh M. Roberts..., Didac Santesmasses[ff,gg], Renata Schama[b,g], Eloy S. Seabra-Junior[ss], L..., Matheus Souza-Gomes[r], Marcos Sterkel[c], Mabel L. Taracena[c], Marta ..., Raul Ursic-Bedoya[d], Thiago M. Venancio[b,x], Ana Beatriz Walter-Nuno..., Richard K. Wilson[h], Erwin Huebner[ww], Ellen M. Dotson[kk,2,4], and Ped...

*An evolving threat*
How gene flow sped the evolution of the malarial mosquito pp. 28, 42, & 43

# Assembly improvement using synteny (in Diptera)



An. gambiae 2R (element 2)

- Leverage synteny to improve assemblies of phylogenetic clusters (Anopheles, Glossina)
- Test data already available (16 genomes)

# Rule #8: (KG)
## "Listen to the teacher(s)"

- No matter how much you share, someone will have to do the work of organizing and usually writing the paper

- This is a labor of love, but can be made easier with good junior (hungry) people leading the analysis

- There is plenty of time/opportunity for followup efforts once the first effort is done

# Rule #9: (KG-ish)
## "Warm food and cold beer are good for you"

- It costs resources, but the smoothest projects are ones that start over dinner and a beverage

- Try to fold in an organization meeting or two into your "big" meeting: ESA, Crete meeting, ASTMH, etc.

- An hour spent talking over beer makes much of the earlier rules much easier

# Rule #10:
## "Bigger is usually better"

- Using older PacBio chemistries, only local assembly continuity (contigs) is better than long-range library assisted Illumina

- The DNA demands have been large, but our best "single shot" results have been newer PacBio on size selected libraries
  - Residual haplotypes need to be removed by hand or software, esp. in inversions

- We also have had relatively good results with HiC scaffolding, which we have applied to all our proposed references (as have others)
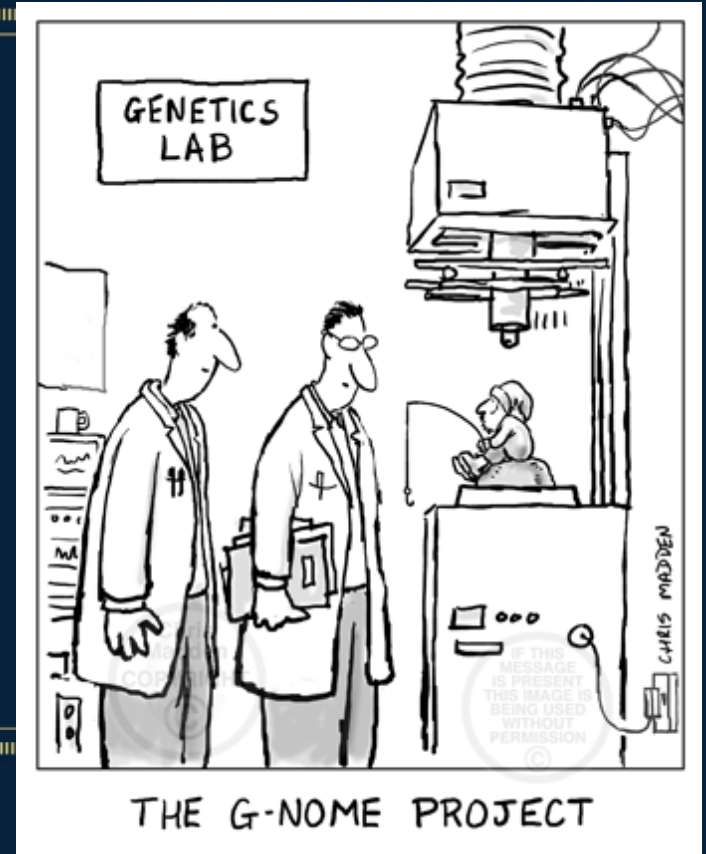
# Followup and funding

Notre Dame Bioinformatics
(@NDBioinformatics)

VectorBase team (EBI, ND, Imperial)

NIH/NIAID for funding

semrich@nd.edu
@ScottEmrich



THE G-NOME PROJECT



ANOPHELES ALIAS MALARIA MOSQUITO

WANTED ... DEAD OR ALIVE

## Questions?