# THE TROUBLE WITH
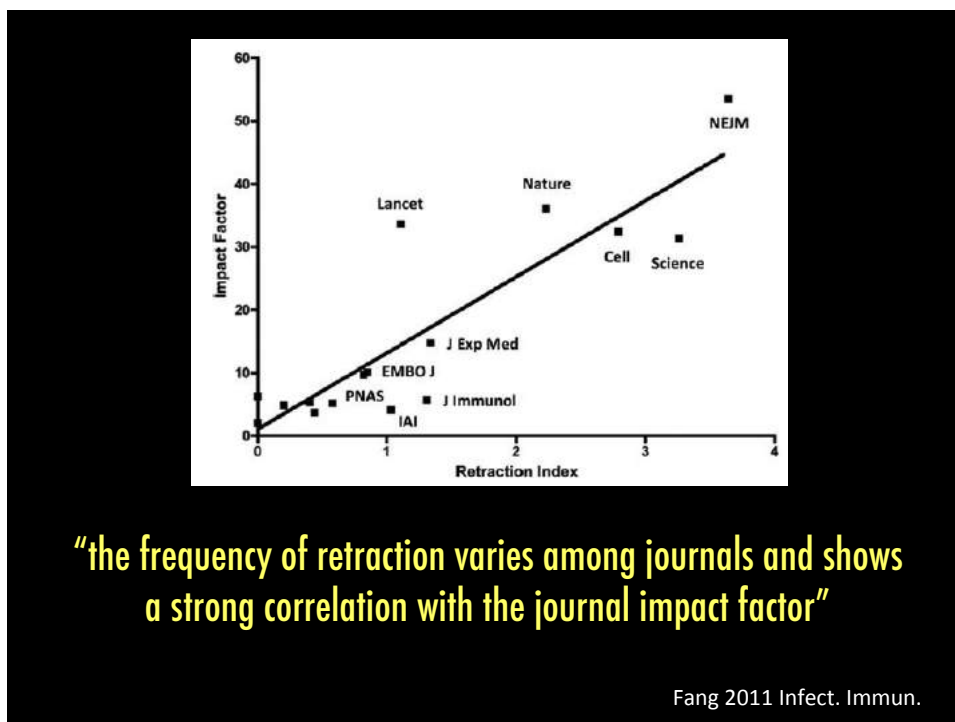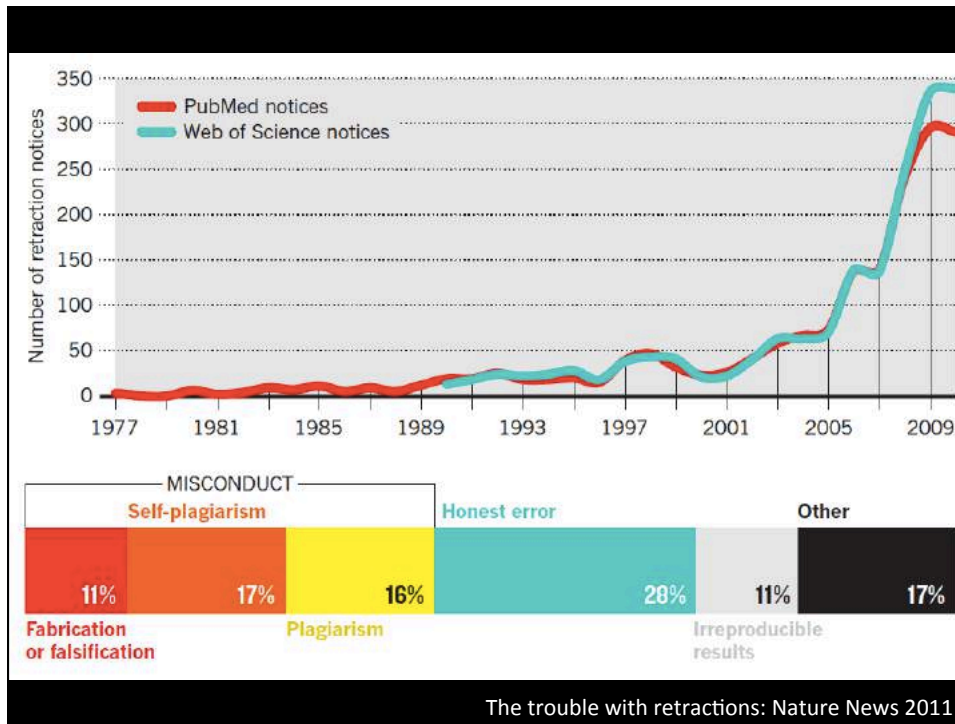
# RETRACTIONS

BY RICHARD VAN NOORDEN

*A surge in withdrawn papers is highlighting weaknesses in the system for handling them.*

The trouble with retractions: Nature News 2011



"the frequency of retraction varies among journals and shows a strong correlation with the journal impact factor"

Fang 2011 Infect. Immun.

## Publications with significant human error that have not been retracted

**PNAS**

### Comparison of the transcriptional landscapes between human and mouse tissues

"the expression for many sets of genes was found to be more similar in different tissues within the same species than between species"

**ARTICLE** 174 | NATURE | VOL 473 | 12 MAY 2011

doi:10.1038/nature09944

### Enterotypes of the human gut microbiome

we identify three robust clusters (referred to as enterotypes hereafter) that are not nation or continent specific … mostly driven by species composition

**LETTER** 228 | NATURE | VOL 502 | 10 OCTOBER 2013

doi:10.1038/nature12511

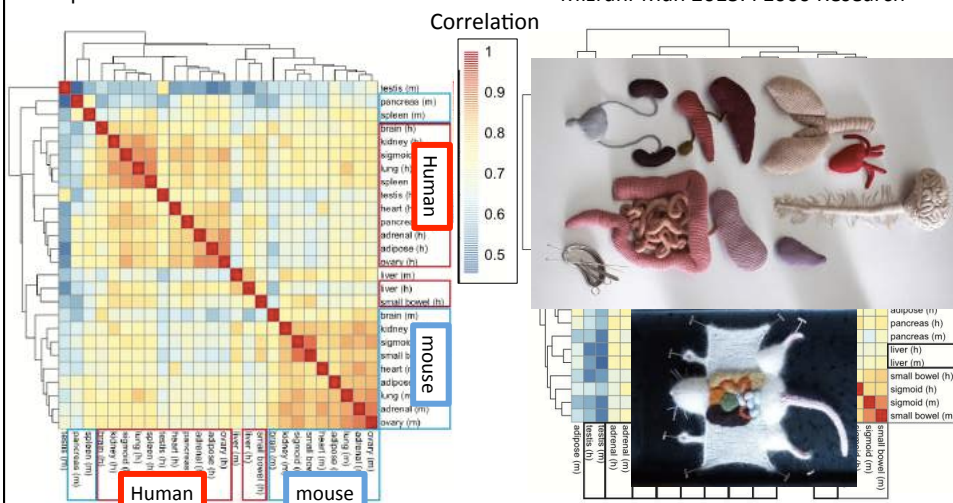### Genome-wide signatures of convergent evolution in echolocating mammals

**PNAS**

### More genes underwent positive selection in chimpanzee evolution than in human evolution

---

# Snyder mouse controversy

"the expression for many sets of genes was found to be more similar in different tissues within the same species than between species" Lin et al. 2014 PNAS

## Human – Mouse TMRCA ~ 90 MYA

## Brain – Kidney TMRCA?

"[after accounting] for the batch effect, …human and mouse tend to cluster by tissue, not by species" Gilad and Mizrahi-Man 2015. F1000 Research
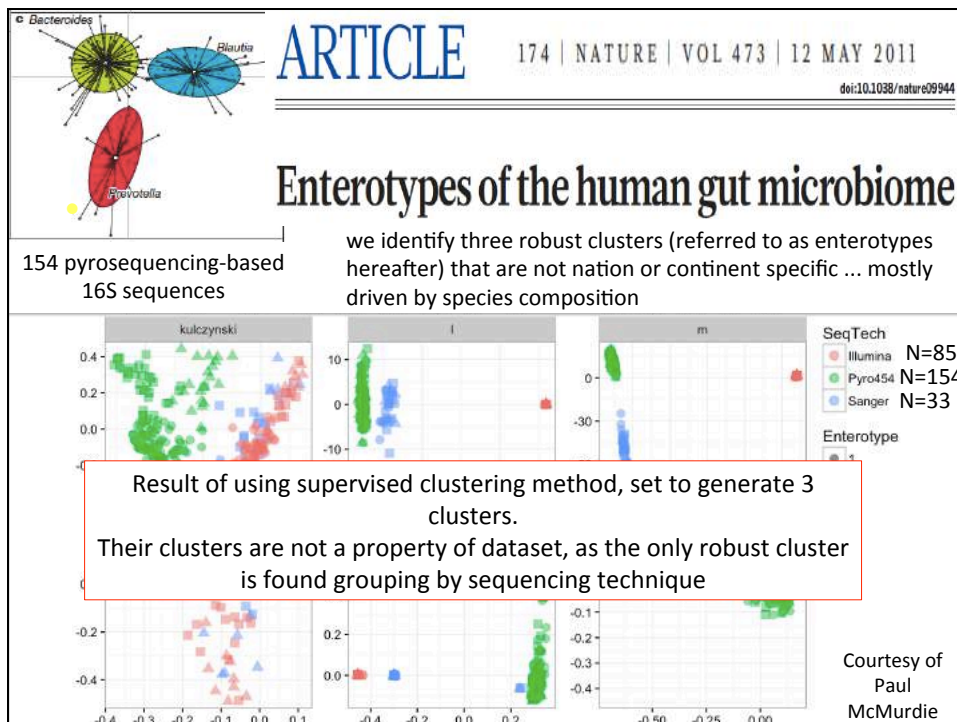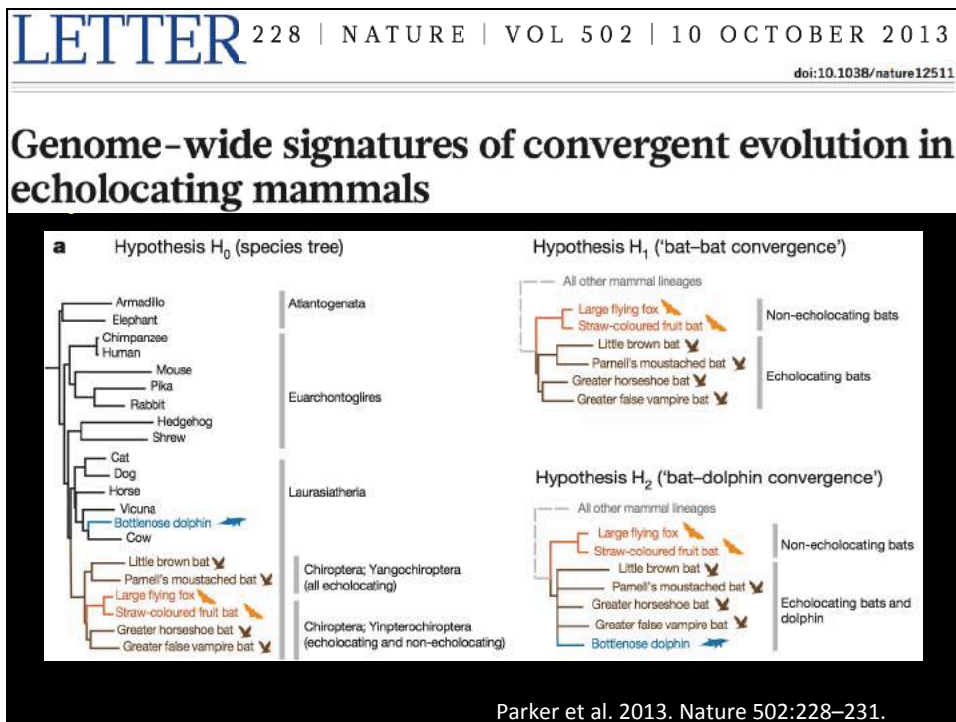


Correlation

## Batch effect: confounding sequencing grouping with biological grouping

| D87PMJN1 (run 253, flow cell D2GUAACXX, lane 7) | D87PMJN1 (run 253, flow cell D2GUAACXX, lane 8) | D4LHBFN1 (run 276, flow cell C2HKJACXX, lane 4) | MONK (run 312, flow cell C2GR3ACXX, lane 6) | HWI-ST373 (run 375, flow cell C3172ACXX, lane 7) |
|---|---|---|---|---|
| heart | adipose | adipose | heart | brain |
| kidney | adrenal | adrenal | kidney | pancreas |
| liver | sigmoid colon | sigmoid colon | liver | brain |
| small bowel | lung | lung | small bowel | spleen |
| spleen | ovary | ovary | testis | ● Human |
| testis | | pancreas | | ● Mouse |

**Solution = Keep technical effects orthogonal to biological**

- Mouse & Human in same lane, same tissues in same lane
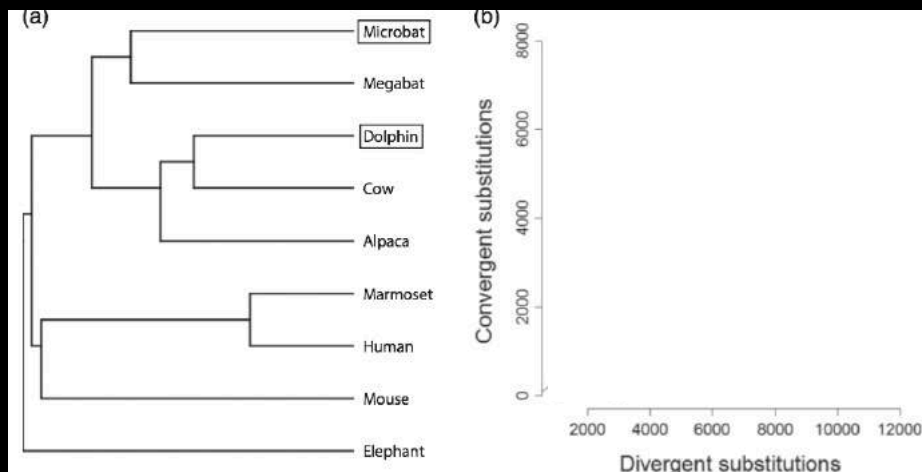  - Will your Core facility know to do this for you?

---



ARTICLE    174 | NATURE | VOL 473 | 12 MAY 2011
doi:10.1038/nature09944

# Enterotypes of the human gut microbiome

154 pyrosequencing-based 16S sequences

we identify three robust clusters (referred to as enterotypes hereafter) that are not nation or continent specific ... mostly driven by species composition

SeqTech
● Illumina  N=85
● Pyro454  N=154
● Sanger  N=33

Enterotype

Result of using supervised clustering method, set to generate 3 clusters.
Their clusters are not a property of dataset, as the only robust cluster is found grouping by sequencing technique

Courtesy of Paul McMurdie

## Genome-wide signatures of convergent evolution in echolocating mammals



Parker et al. 2013. Nature 502:228–231.

"Strong and significant support for convergence among bats and the bottlenose dolphin was seen in numerous genes linked to hearing or deafness, consistent with an involvement in echolocation."

- 2326 orthologous genes
- site-wise log-likelihood support (SSLS)
  - Negative values support convergence H1,H2
    - 824 mean support for H1
    - 329 mean support for H2

Hearing
Vision

# Parker et al. failed to conduct orthogonal 'test' of findings or estimate proper 'null' expectation



Thomas and Hahn 2015. Mol Biol Evol 32:1232–1236.

# What makes us difference from chimps?

## Is it really just 2%



---

More genes underwent positive selection in chimpanzee evolution than in human evolution

Margaret A. Bakewell, Peng Shi, and Jianzhi Zhang*

201 citations since 2007

Table 1. Genic positive s

| Comparison |
| --- |
| No. of genes analyzed |
| No. of PSGs |
| No. of PSGs |
| No. of PSGs |
| No. of PSGs |
| No. of synoi |
| No. of nons |
| Mean ω of a |
| Mean ω of |

Only 2 genes of original 59 were validated!!
(at bioinformatic level)

- Many chimpanzee-specific divergent sites are adjacent to indels
- removing nucleotides within five positions of indels abolishished most adaptive signals

# Evolutionary Inference = House of Cards?

The quality of our evolutionary inference

Is proportional to assumptions of orthology



# Orthologous genes ... can their phenotypic effects drift over evolutionary time?

- RNAi phenotypes assessed for 1,300 genes in two nematodes
  - TMRA ~ 24 MYA
  - 7% had divergent phenotypic effects (in lab, etc.)
  - Likely higher in nature



Verster et al. 2014. PLoS Genet

## So ... how many of you are sequencing a genome?

- What does that mean?

- What kind of genome are you generating?

- What is your question?
  - Short term vs. long term goals?
  - Are these in conflict?

# Is there a genome for humans?

## Genomes of 2,504 individuals from 26 populations



- You differ from references on average at:
  - 4 to 5 M SNPs
  - ~2k structural variants covering ~20 M bp

1000 Genomes Project Consortium (2015) Nature

# What does this mean

- Most species have lots of genomic polymorhism
  - SNPs are just the tip of the iceburg, lots of structural changes
  - Characterizing all the variation is very expensive

- But

- Very rarely will your questions require chromosomal level assembly
  - Thus you can get to your answers much faster and cheaper if you generate what you need rather than working for an ideal you don't need



Three years, ~300,000 Euros

Hill et al., in prep.

## Slide 1

Plutella xylostella
Bombyx mori
Manduca sexta
Operophtera brumata
Spodoptera frugiperda
Helicoverpa armigera
Plodia interpunctella
Chilo suppressalis
Papilio glaucus
Papilio machaon
Papilio polytes
Papilio xuthus
Lerema accius
Calycopis cecrops
Danaus plexippus
Bicyclus anynana
Melitaea cinxia
Heliconius melpomene
Heliconius erato demophoor
Leptidea sinapis
Phoebis sennae
Pieris rapae HiRise
Pieris napi
Pieris rapae

227                    729
Genome size (Mbp)

Published genomes vary dramatically in quality

**Which do you need NOW?**

Hill et al., in prep.

## Slide 2

# Depending on your question

Just sequence lots of genomes
Generate hypotheses
Test them

# Genomic signal of Diapause adaptation

Speckled Wood
(*Pararge aegeria*)

Genomic tools at start:
- mtDNA and microsat loci
- Extensive ecological studies > 10 years

Stockholm University

Karl Gotthard

Peter Pruisscher

Peter Pruisscher

---

Speckled Wood
(*Pararge aegeria*)

| Generations per year | % in diapause at 18 hours light |
|---|---|
| 1 | 100 % |
| 2 | 0 % |

What is the genetic basis of adaptation to day length?

MESPA: Mining Exons and Scaffolding on Poor Assemblies

Amino acid sequence:

Genomic contigs:

Find aligned regions:

Output:
NNNNNNN NNNNNNN
- scaffolds based upon exons
- cDNA of genes
- GFF files for the scaffolds (start, stop, exon boundaries)

Can use 1000s genes (much more than BUSCO):
• quantify # found in assembly and their length
• can scaffold these regions for better gene space coverage
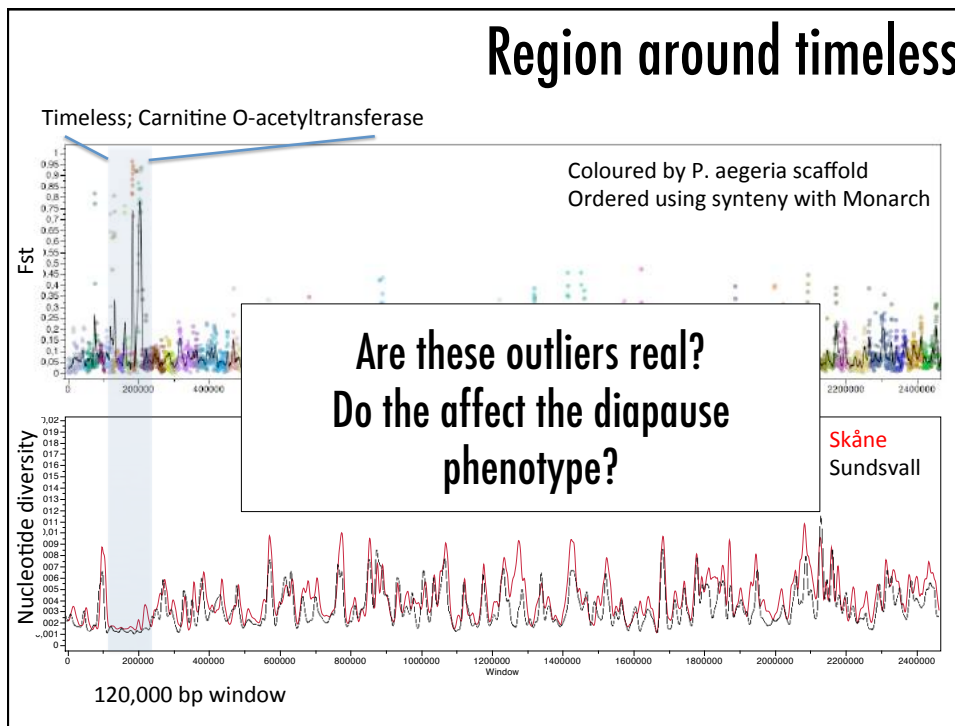• identify and work with these high quality scaffolds

Ram
Neethiraj

Neethiraj et al. 2017 ME

# Fst outlier analysis for candidates

EXON1 EXON2    EXON3

11,000 gene models & ~7 million SNPs

Quality Filtering

~ 114,000 SNPs of which 68,000 SNPs: FST >0.9

A/C

7 million SNPs

Filtering

~ 114.000 SNPs



# Fixed variation in genes

1. Intergenic regions contain+/- 67,604 Fixed SNPs

2. 67 gene models contain 209 fixed SNPs

3. Filter for SNPs in exons and introns



SNPs per gene model

| UniRef90_proteinnames | exon | gene | intergenic | Total | D.plex scaffold | Bmori_chr |
|---|---|---|---|---|---|---|
| Timeless | 2 | 0 | 0 | 2 | DPSC300014 | chr4 |
| Carnitine O-acetyltransferase | 3 | 25 | 1 | 29 | DPSC300014 | chr4 |
| Trypsin-like protein | 2 | 14 | 14 | 30 | DPSC300041 | chr5 |
| Vasa-like protein | 1 | 2 | 0 | 3 | DPSC300379 | chr19 |
| Period | 2 | 2 | 1 | 5 | DPSC30005 | chr1 |

Is there a foot-print of selection around these SNPs?

**Region around timeless**

Timeless; Carnitine O-acetyltransferase

Coloured by P. aegeria scaffold
Ordered using synteny with Monarch

Are these outliers real?
Do the affect the diapause phenotype?

Skåne
Sundsvall

120,000 bp window



**Region around timeless**

Timeless; Carnitine O-acetyltransferase

Coloured by P. aegeria scaffold
Ordered using synteny with Monarch

Rather than argue about the significance of this high Fst, I move quickly onto testing this hypothesis!

Skåne
Sundsvall

120,000 bp window

# Validating genomic hypothesis of Timeless

SNP genotyping in F2 cross

Clinal analysis



# Genomics to hypothesis to validation

- Use genomics to generate robust hypothesis
  - Orthogonal methods
  - Stong signals

- Validate upwards
  - Use independent biological samples
  - Higher level of biological organization
  - Simultaneously test hypothesis and its generality

# *Colias croceus*, the Clouded Yellow

Male          Female          Alba Female

Female limited alternative life history strategy (and/or reproductive strategy?)

Life History differences:

VS.

N=15          N=15

GWAS + genome + QTL mapping
(blood, sweat, tears)

-log 10 (FDR)

Contig 12   1                                        430,000

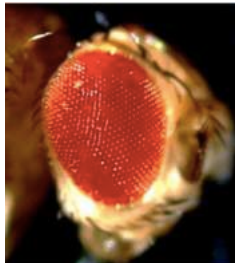(c)

log 10 (Read Depth)
orange

Alba

(d) Contig 12         220,000      230,000      240,000

*BarH-1*          DNA polymerase from
jockey-like TE

Alyssa Woronik, Phd

## Bar is functionally required in primary pigment cells of developing ommatidia
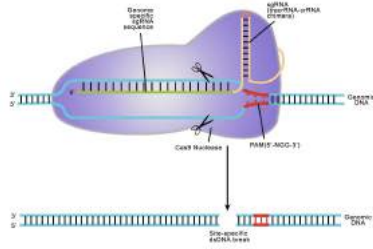
*Drosophila melanogaster*

Normal eye

Bar Knockout

Higashijima et al. 1992; Kang et al. 2013

- Hox gene transcription factor
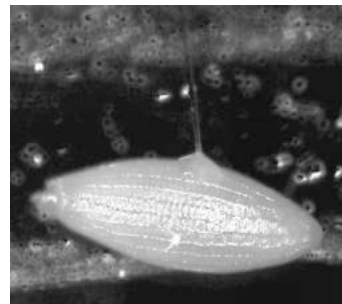- Repressor of other developmental genes

Should we spend money & time:

Validating mapping and GWAS?
or
Validate upward

Video

## Testing BarH1 hypothesis: CRISPR/Cas9 knockout of Bar



Allows us to remove function of BarH1
Up and working in 2 months
Masters project in lab

Injection of Cas9 protein + guideRNA into Colias croceus egg

## Developmental defects:

- Lack of pigment formation within ommatidia
- Equivalent to *Drosophila* phenotype

Phenotypes observed using 2 separate gRNA constructs, awaiting PCR validation

BarH1 knockout

John Hallman

# CRISPR/Cas9 results

| Individual | gRNA | Sex | Eye | Proboscis |
|---|---|---|---|---|
| CC58 | 3 | F | yes | yes |
| CC51 | 3 | M | yes | yes |
| CC31_2 | 3+4 | F | yes | yes |
| CC33 | 3+4 | F | yes | yes |
| CC31_1 | 3+4 | M | yes | yes |
| CC52 | bar5 | F | yes | yes |

- \>2000 eggs injected
  - Consistent developmental phenotype
- BarH1
  - Involved in development of eye, proboscis
  - Not involved in orange / white wing coloration
  - No sex specific effects

Woronik et al., in prep



When injected into Alba females color mosaic phenotype

Woronik et al., in prep

# 1001 ways for your pipeline to break

## An overview of genomic pipeline challenges

## Christopher West Wheat

---

# Informatics and Biology

- We need to make sure we put the 'bio' into the bioinformatics
  - Do results pass 1st principals tests
  - Always double check data from your core facility or service company
  - Use independent analyses as 'controls' on accuracy
    - What are your + and – controls?
    - Do independent methods converge?

- Need to re-assess our common metrics for potential bias in the genomic age
  - Bootstraps on genomic scale data
  - P-values, outlier analyses, demographic null models

# Outline

- Transcriptome analyses in non-model species
  - Walk through pipeline and highlight issues of concern
  - What is validation?

- Insights from candidate genes
  - Can Second Gen methods get us there?

# Pipeline Overview
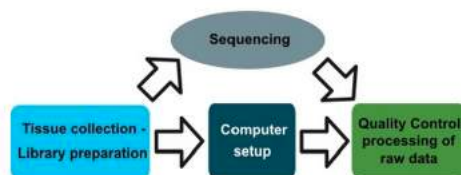
# Pipeline Overview

# Computer Infrastructure

RNAseq dataset:
4 conditions X 2 tissues X 3 families X 3 replicates = 72 X 10^6 reads

| | File Sizes (Gb) | CPUs | RAM (Gb) | Time |
|---|---|---|---|---|
| Raw files *gz | (1.5... | | | ~3 hours / file |
| Raw files expanded | | | | |
| TA assembly | | | | weeks |
| Mapping (BAM) | | | | hours / file |
| Annotation | 10... | | | ~6 – 12 days |
| Analysis | < 20 Mb | 4 | 4 | ~< 1 hour |
| Visualization | BAM files | ≥ 4 | ≥ 8 | |

Get ready for your data by downloading similar sized dataset from the Short Read Archive. Do not wait till it arrives

# Pipeline Overview

Sequencing

Tissue collection - Library preparation → Computer setup → Quality Control processing of raw data

## Core facilities and non-model species

Statements from core facilities that are not true:

- Here is your data

- You can't do RNA-Seq without a genome

- We'll have your data back in < 1 month

# Pipeline Overview

# Gene Ontology: order in the chaos

- Addresses the need for consistent descriptions of gene products in different databases in a species-independent manner

- GO project has developed three structured controlled vocabularies (ontologies) that describe gene products in terms of their associated
  - biological processes
  - cellular components
  - molecular functions



http://www.geneontology.org/

# Comparisons among annotation tools





Radivojac et al.: **A large-scale evaluation of computational protein function prediction**. *Nat Meth* 2013, **10**:221–227.
Falda et al. **Argot2: a large scale function prediction tool relying on semantic similarity of weighted Gene Ontology terms**. *BMC Bioinformatics* 2012, **13**:S14.

# Batch processing for GO terms

# Pipeline Overview



# Template mismatch effects: excellent yeast study



Nookaew et al 2012

# Does alignment software matter?

# Mappers don't appear to matter

**Wrong**

- Genomic scale data can hide widespread biases that unless you specifically look, are hard to find

- Mapping programs differ in their settings and design
  - DNA to DNA vs. RNA to DNA
  - Are usually compared using species without much genetic variation
  - Indels, splicing, SNPs all affect mapper performance

# SNP effects can be large



Nookaew et al. **A comprehensive comparison of RNA-Seq-based transcriptome analysis from reads to differential gene expression and cross-comparison with microarrays: a case study in Saccharomyces cerevisiae.** *Nucleic Acids Research* 2012, **40**:10084–10097.

# Insertions & deletions (indels) have large effects



Nookaew et al. **A comprehensive comparison of RNA-Seq-based transcriptome analysis from reads to differential gene expression and cross-comparison with microarrays: a case study in Saccharomyces cerevisiae.** *Nucleic Acids Research* 2012, **40**:10084–10097.
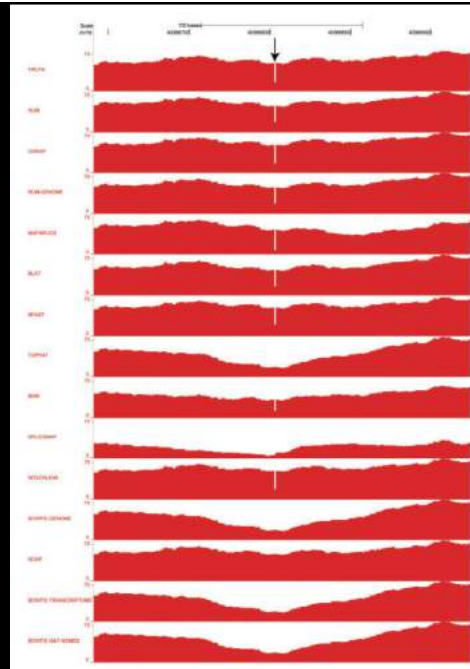
# 15 mapping results

**Dramatic differences in ability to handle a 2 bp insertion in reference compared to reads**
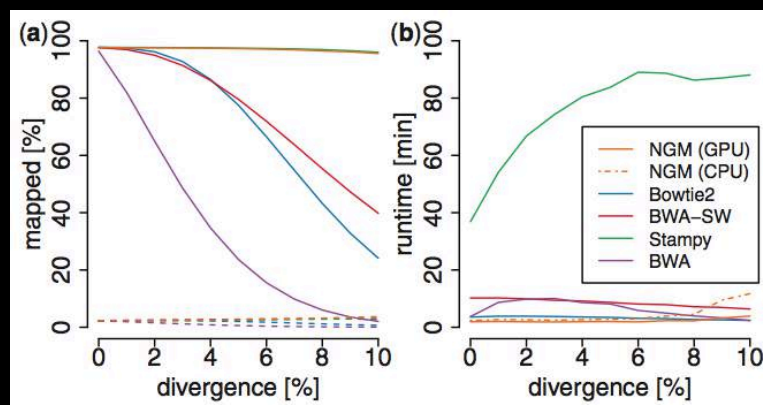
**TopHat, SpliceMap, Bowtie and Soap**

- – do not identify indels
- – they fail to accurately align reads to these regions



Grant GR, Farkas MH, Pizarro A, Lahens N, Schug J, Brunk B, Stoeckert CJ, Hogenesch JB, Pierce EA: **Comparative Analysis of RNA-Seq Alignment Algorithms and the RNA-Seq Unified Mapper (RUM)**. *Bioinformatics* 2011, doi:10.1093/bioinformatics/btr427.

# Allelic bias in read mapping



- **Essentially identical to allele specific PCR bias ... but on a scale you can't detect unless you care to look**
- **Do your genes of interest have more than 3 SNPs / 100 bp?**
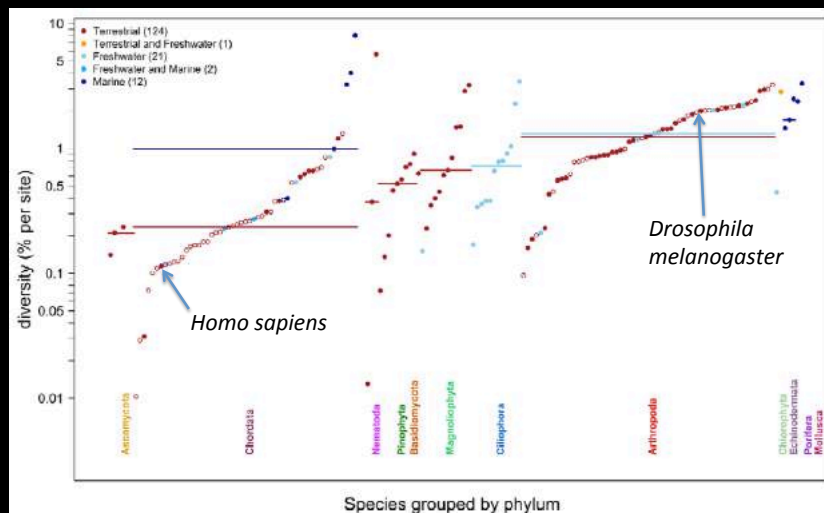
Sedlazeck et al. 2013 *Bioinformatics*

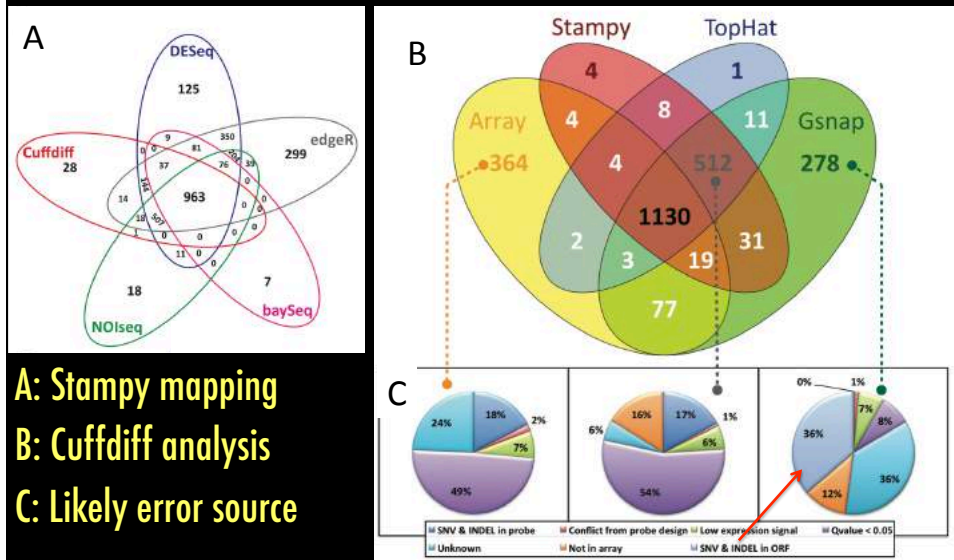# 100 bp window with 4 – 5 SNPs differing from reference



# Mapping reads in outbred species

## Average genome polymorphism levels (ignores indels)



*Drosophila melanogaster*

*Homo sapiens*

Leffler *et al.* 2012 *Plos Biol*

# Sig. expression differences by method

A



B



C

A: Stampy mapping
B: Cuffdiff analysis
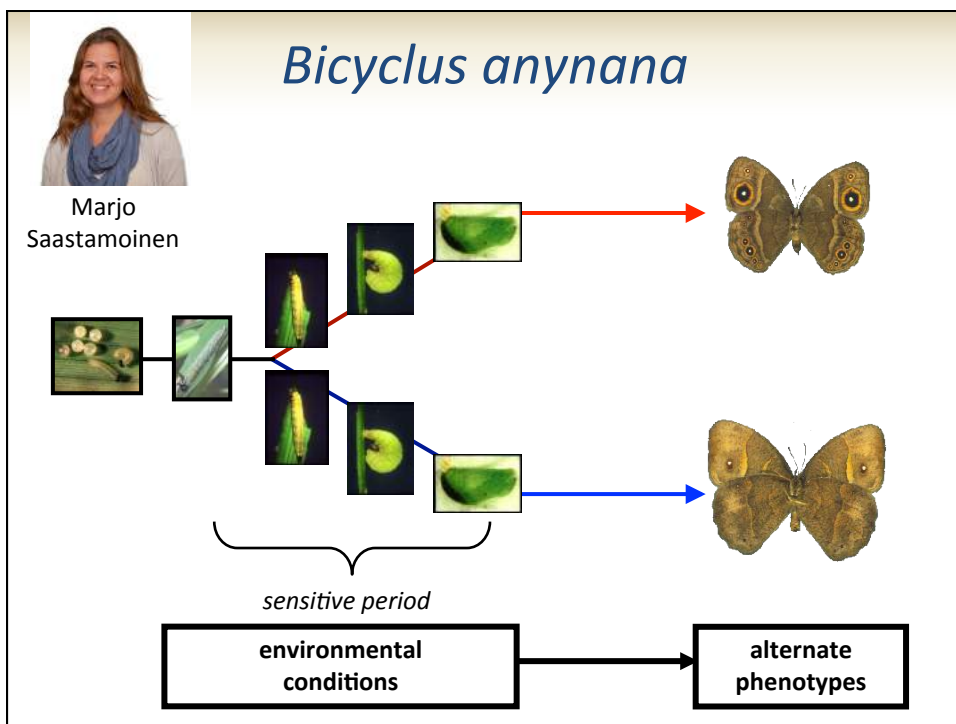C: Likely error source

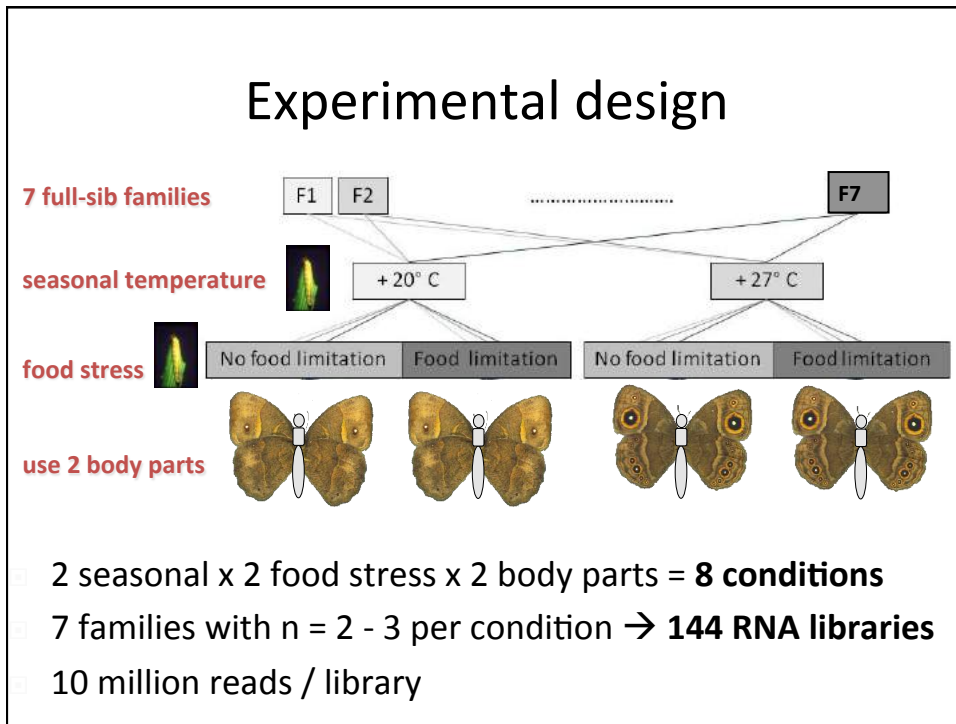# RNA-Seq

Real world example

2 factor analysis with family effects

**Bicyclus anynana**

*Save energy, live long*

*Live fast, die young*

| long | lifespan | short |
| delayed | reproduction | fast |
| inactive | behaviour | active |
| high | fat reserves | low |
| cryptic | wing pattern | conspicuous |



**Bicyclus anynana**

Marjo Saastamoinen

*sensitive period*

| environmental conditions | → | alternate phenotypes |

# Experimental design

**7 full-sib families**   F1  F2  ................  F7

**seasonal temperature**   + 20° C        + 27° C

**food stress**   No food limitation | Food limitation    No food limitation | Food limitation

**use 2 body parts**

- 2 seasonal x 2 food stress x 2 body parts = **8 conditions**
- 7 families with n = 2 - 3 per condition → **144 RNA libraries**
- 10 million reads / library

---

Vicencio Oostra

| body part | # libraries | # clean reads (per library) | # nucleotides (per library) | GC content |
|---|---|---|---|---|
| abdomen | 72 | 15,261,019 | 3,052,203,767 | 45% |
| thorax | 72 | 15,633,416 | 3,126,683,150 | 46% |
| total | 144 | 2,224,399,290 | 444,879,858,000 | 45% |

14 samples: one from each family, thorax and abdomen          69,075 contigs
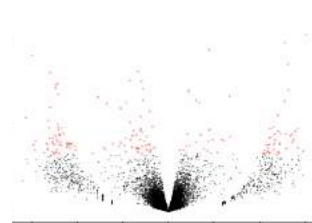
# edgeR          Bioconductor

```
# reads ~    season + stress + family +
            season*stress + season*family + stress*family
            season*stress*family
```

Season

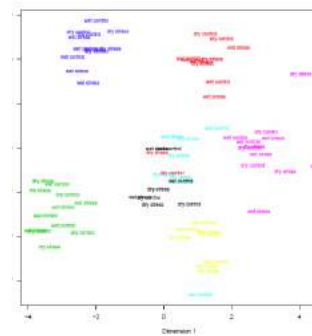What should I be looking at first?



Colored by Family

76

Stress

## Effect of filtering the mapping to Trinity contigs

71 zero-read samples
allowed

## GLM results

- Plastic responses:
  - Effects without any interaction with Family

season x treatment x family
**116**

**22**          **27**

**23**

**115**          **43**

seasonal x family

**15**

stress x family

- Genetic response:
  - Effects that have an interaction with family
  - Potential targets of natural selection

```
reads ~    season + stress + family + season*stress  +
           season*family + stress*family + season*stress*family
```

## Most studies are annotation limited

- **What is the biological meaning of the top P-value genes?**
- **Low P-value or expression genes are certainly important**
- **Gene set enrichments are key to insights**
  - Thus, annotation is very important

| Description | Uniprot | -log10P |
|---|---|---|
| Oxidoreductase. | Q9VMH9 | 7.087008 |
| Hypothetical protein. | | 6.993626 |
| SD27140p. | | 6.315473 |
| | Q8SXX2 | 6.300667 |
| SD01790p. | Q95TI3 | 5.316371 |
| Electron-transfer-flavoprotein I | Q0KHZ6 | 5.1425 |
| Pseudouridylate synthase. | Q9W282 | 4.784378 |
| Hypothetical protein. | Q9VGX0 | 4.750469 |
| CG14686-PA (RE68889p). | Q9VGX0 | 4.650051 |
| Chromosome 11 SCAF14979, wh | Q8T058 | 4.506043 |
| | | 4.470413 |
| , complete genome. (EC 1.6.5.5 | | 4.445501 |
| RNA-binding protein. | | 4.374033 |
| Hypothetical protein. | Q9VPL4 | 4.369727 |
| Peptidoglycan recognition-like | | 4.206247 |
| Angiotensin-converting-related | Q8SXX2 | 4.172776 |
| Lachesin, putative. | Q9I7H7 | 4.056174 |
| Secretory component. | Q9VVK5 | 3.981175 |
| Putative adenosine deaminase | Q9VVK5 | 3.980728 |
| | | 3.95787 |

7 of 20 (35%) no Uniprot ID

# Sources of error

Transcriptome assembly can be huge source of bias:
- Fragmentation creates multiple contigs of same gene
- SNPs and alternative splicing generates more contigs
- 1 locus = frag. X SNPs X alt. splicing = many contigs

We can observe effects in expression analyses:
  – Family effect mapping bias
  – Pseudo-inflation in Gene Set Enrichment Analyses

# Put the BIO in your informatics!!

### Use independent analyses as 'controls' on accuracy
#### — What are your + and - controls?

|  | Analysis # 1 | Analysis # 2 | Analysis # 3 |
|---|---|---|---|
| Mapper | TopHat2 | STAR | ? |
| Normalization | none | TMM | TMM |
| Analysis | PCA | RSEM | EDGER |

## Should independent methods converge?

# Interrogate your results

- "you need to be in charge of the analysis" – B. Cresko

- This will give you confidence
  - Bring freedom to your findings (no waterboarding)

- Graph your results – visualize the patterns
  - PCA or MDS plot
  - P-value distributions

- Assess gene copy number in gene set enrichment analyses (GSEA)
  - Do these levels fit to 1$^{st}$ principals expectations?
  - Do you have extra copies due to your Transcriptome assembly?

# A major challenge for Ecological Genomics

- What causes natural selection in the wild?
  - How does genetic variation at one region of the genome interact with its environment (genomic, abiotic, and biotic)

- DNA alone can't tell us about selection dynamics in the wild
  - Molecular tests are very weak and uninformative about selection dynamics

- Research community is demanding actual demonstration of natural selection when making claims of adaptive role
  - Triangulate!!!!



Molecular spandrels:

Story telling
vs.
Causal understanding

Genomics is full of adaptive stories

Functional and field validation of SNPs effects are needed to discern facts from fiction

Storz & Wheat 2010 *Evolution*          Barrett & Hoekstra 2011 *Nat Rev Genet*

Team Alba

Constantí Stefanescu

Alyssa Woronik

Mike Perry

Maria Celorio

John Hallmén

Philipp Lehmann

Kalle Tunström

Erik Philip-Sörensens stiftelse
FÖR FRÄMJANDET AV GENETISK OCH HUMANISTISK VETENSKAPLIG FORSKNING

Stockholm University

Knut och Alice Wallenbergs Stiftelse



# Thanks!
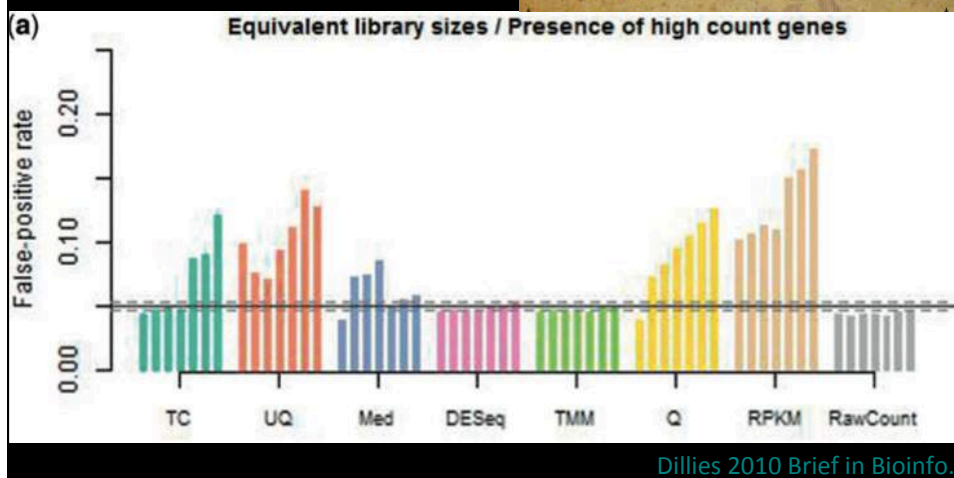
Stockholm University

ACADEMY OF FINLAND

VETENSKAPSRÅDET

Knut och Alice Wallenbergs Stiftelse

# Common mistakes

- Blindly trusting bioinformaticians: look at your data!!!
- Mapping reads to a very divergent genome
  - Only most conserved genes map: bias due to divergence and mapping thresholds
- Not accurately assessing a TA
  - Your template determines quality of results
- Not enough reads, replication, or statistical power
  - Large amounts of data to not change fundamental statistics (never pool unless necessary)
- Not assessing likely biases in analyses
  - Try different mapping thresholds & analysis methods to assess convergence of biological signal
  - Assess alternative splicing and duplication potential in findings
- Data size and computational power are demanding
  - Download data and work with it before your real data comes.

# Normalization matters

WILD WILD WEST



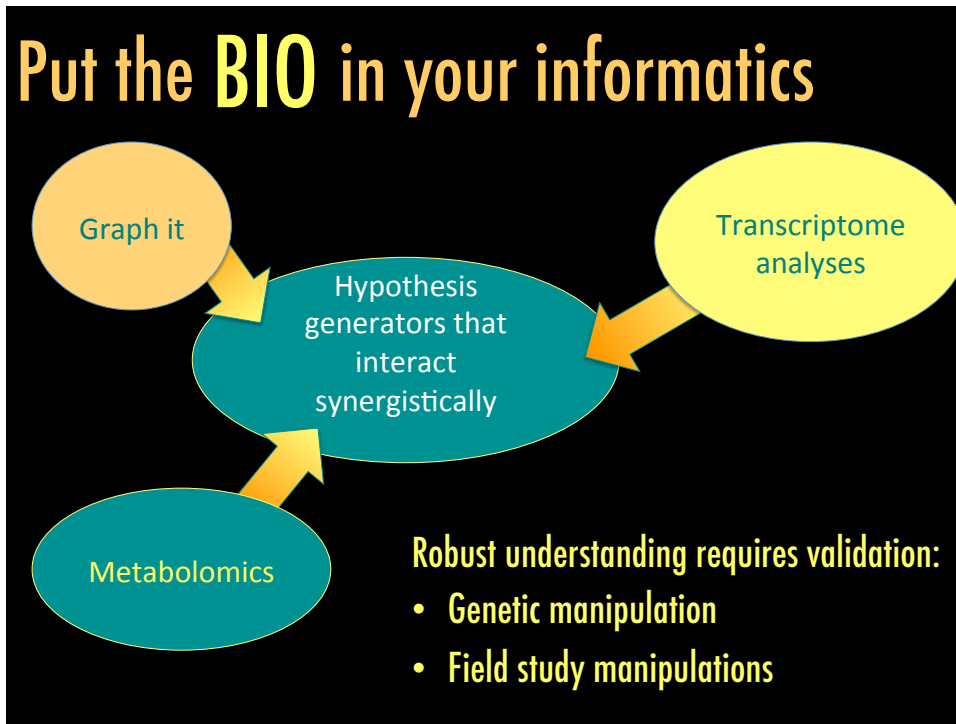Dillies 2010 Brief in Bioinfo.

# Life after your RNA-Seq experiment

— What are you likely to learn?
- By measuring other aspects of the phenotype, you can validate and solidify your transcriptome insights

— What may limit your insights?
- Single gene analyses can be restrictive
  - Statistically: FDR is very conservative
  - Biologically: genes work in networks varying in expression and direction across pathways

— Possible solutions
- Gene set enrichment analysis: harness the functional network
- Collect additional data relevant to your phenotype and organism
  - Don't hesitate to make your own enrichment set, measure hormones and metabolites.

# RNAseq Resources

- Papers
  - Oshlack A, Robinson MD, Young MD: **From RNA-seq reads to differential expression results**. *Genome Biol* 2010, **11**:1-10.
  - Haas BJ, Zody MC: **Advancing RNA-Seq analysis**. *Nat Biotechnol* 2010, **28**:421-423.
  - Grant GR, Farkas MH, Pizarro A, Lahens N, Schug J, Brunk B, Stoeckert CJ, Hogenesch JB, Pierce EA: **Comparative Analysis of RNA-Seq Alignment Algorithms and the RNA-Seq Unified Mapper (RUM)**. *Bioinformatics* 2011, doi:10.1093/bioinformatics/btr427.
  - Wolf JBW: **Principles of transcriptome analysis and gene expression quantification: an RNA-seq tutorial**. *Molecular Ecology Resources* 2013, doi:10.1111/1755-0998.12109.
  - Nookaew I, Papini M, Pornputtapong N, Scalcinati G, Fagerberg L, Uhlen M, Nielsen J: **A comprehensive comparison of RNA-Seq-based transcriptome analysis from reads to differential gene expression and cross-comparison with microarrays: a case study in Saccharomyces cerevisiae**. *Nucleic Acids Research* 2012, **40**:10084-10097.
  - De Wit P, Pespeni MH, Ladner JT, Barshis DJ, Seneca F, Jaris H, Therkildsen NO, Morikawa M, Palumbi SR: **The simple fool's guide to population genomics via RNA-Seq: an introduction to high-throughput sequencing data analysis**. *Molecular Ecology Resources* 2012, **12**:1058-1067.
- Websites
  - http://www.rna-seqblog.com/
  - Google anything that comes to mind
- Workshops
  - http://evomics.org/
  - EBI online
    - http://www.ebi.ac.uk/training/online/course/ebi-next-generation-sequencing-practical-course/rna-sequencing/rna-seq-analysis-transcriptome
- Colleagues
  - Email colleagues and ask questions early, rather than late.

http://sfg.stanford.edu/guide.html

# Validating candiate genes moves us forward:

# Put the **BIO** in your informatics

Graph it

Transcriptome analyses

Hypothesis generators that interact synergistically

Metabolomics

**Robust understanding requires validation:**
- **Genetic manipulation**
- **Field study manipulations**

A great place to start, but not stop

# Model adaptation: the *Eda* gene

- Causes loss in body armor
  - Field association
  - QTL mapping
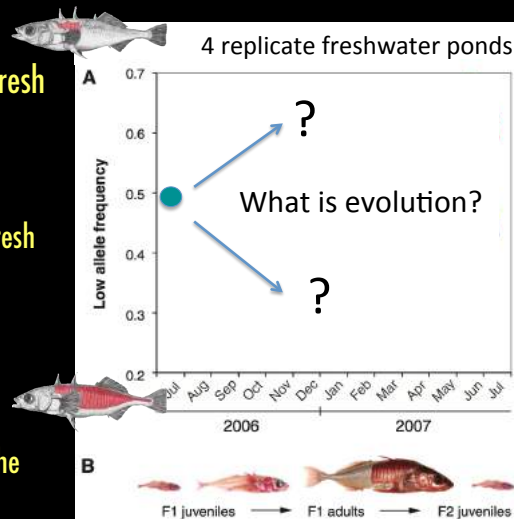  - Gain-of-function assay



Position along chromosome 4



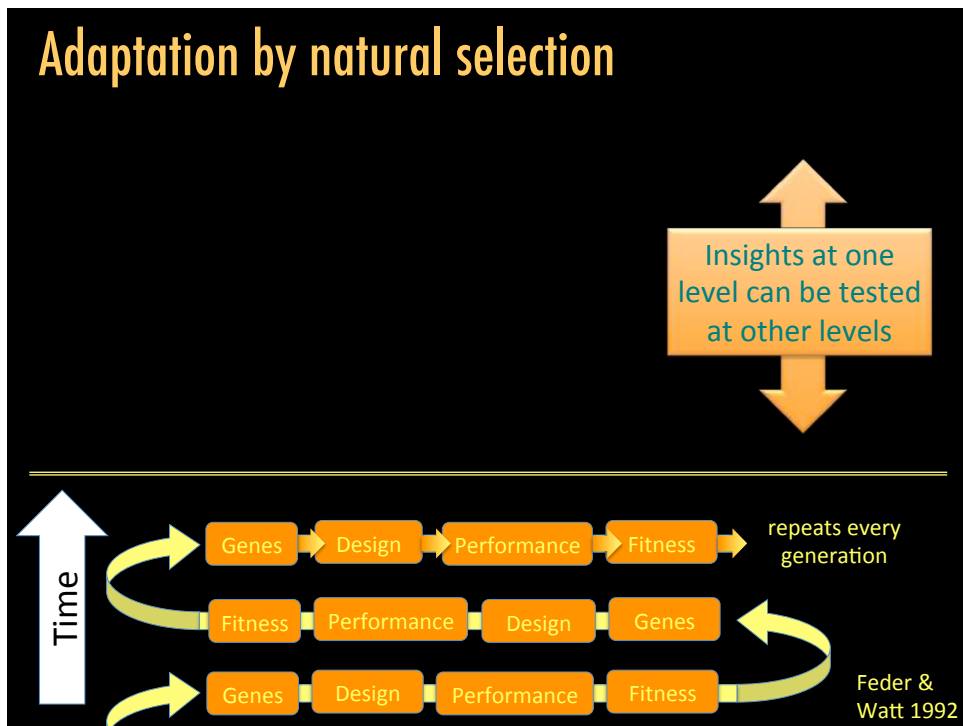# Back to nature: do we know what we think we know?

- Is low armor really adaptive in fresh water?

- Lets replay the selection event
  - Equal frequency *Eda* alleles in fresh water ponds

Studies in the field can uncover unexpected and complex selection dynamics
- Linked effect of other genes in the inversion on LG4?
- Is Eda the target of selection?



4 replicate freshwater ponds

?

What is evolution?

?

Barrett et al. 2008 Science

# Adaptation by natural selection

Insights at one level can be tested at other levels

Time

| Genes | Design | Performance | Fitness | repeats every generation |

| Fitness | Performance | Design | Genes |

| Genes | Design | Performance | Fitness |

Feder & Watt 1992

---

## Assessing transcriptome assembly

- **Assessment metrics**
  - Non-biological
    - N50, # of contigs
  - Biologically informative
    - # of orthologs identified
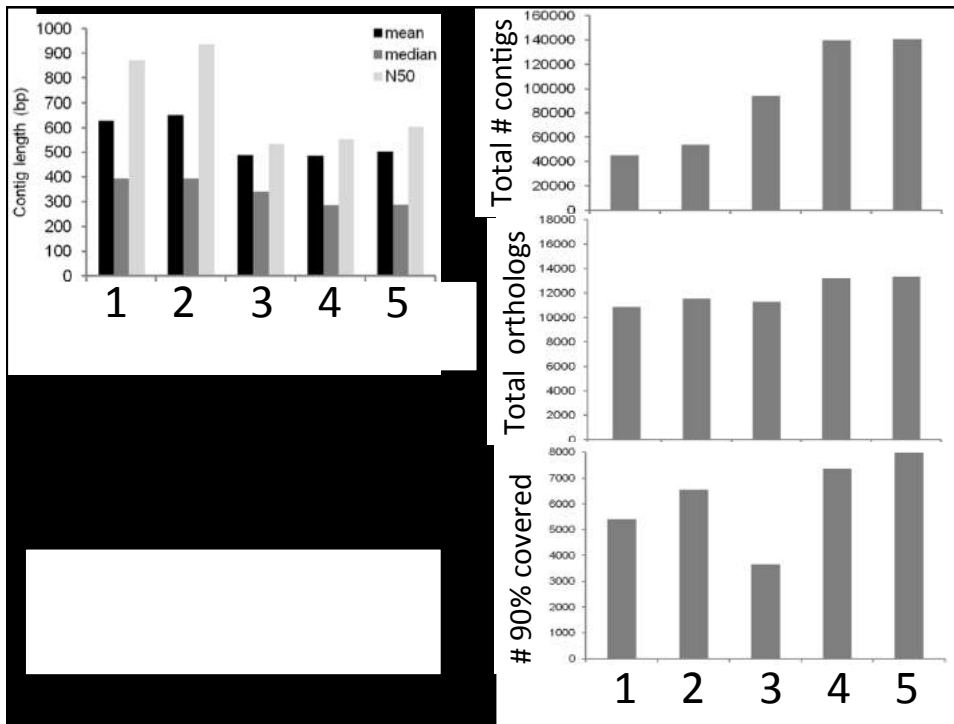    - Ortholog hit ratio (OHR)

$\alpha / \beta$ :
1 = complete
< 1 = % covered

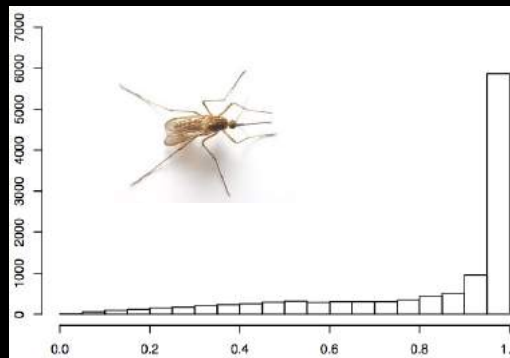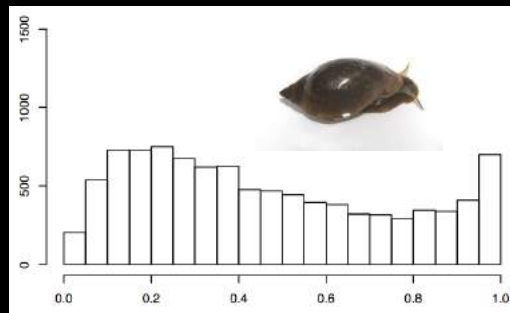$$\alpha / \beta = \frac{\text{TA contig} \quad \text{Length} = \alpha}{\text{Ortholog} \quad \text{Length} = \beta}$$

Hornett & Wheat 2012; O'neil & Emrich 2013 *BMC Genomics*
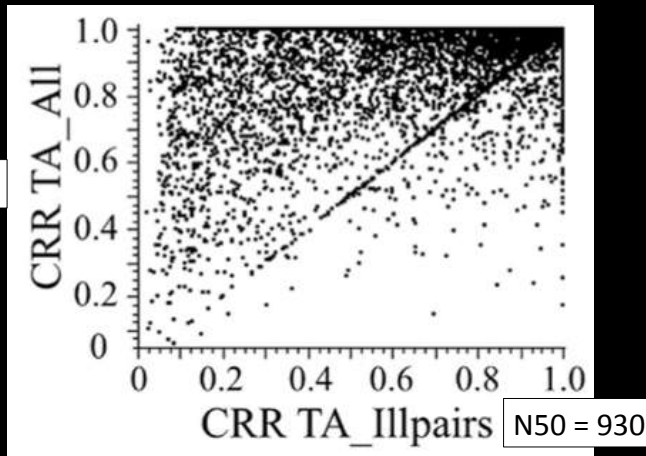
## OHR graphs

- Shows the number of unique orthologs hit
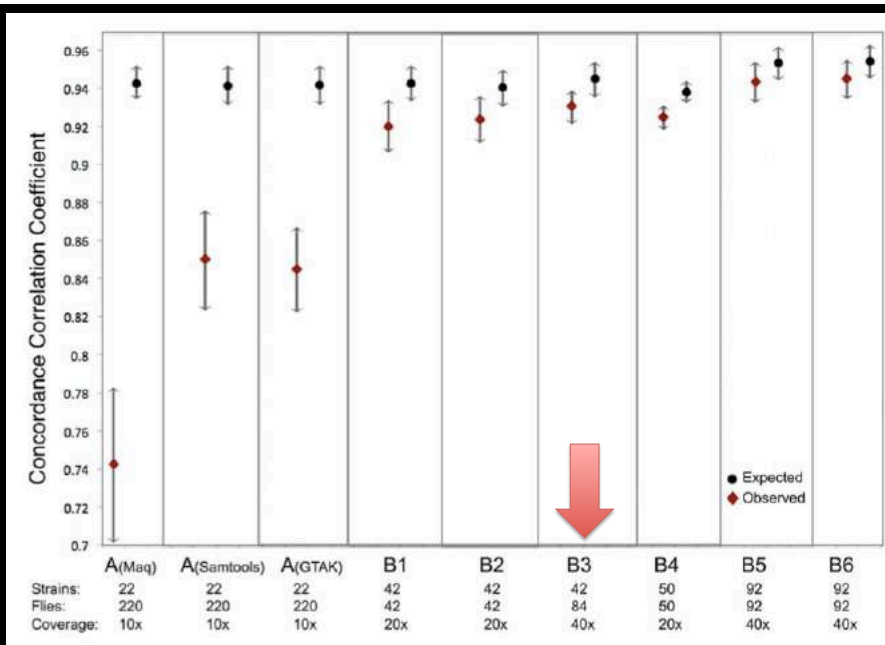
- Distribution of their reconstructed length

**Comparative OHR**

- Compare longest contig per ortholog for two assemblies
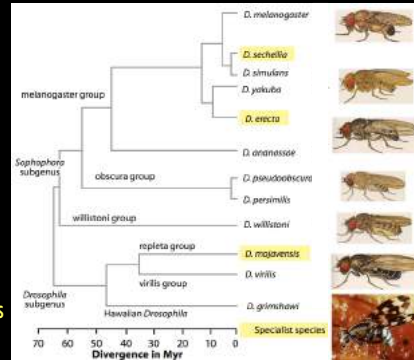- Plot them against each other

N50 = 610
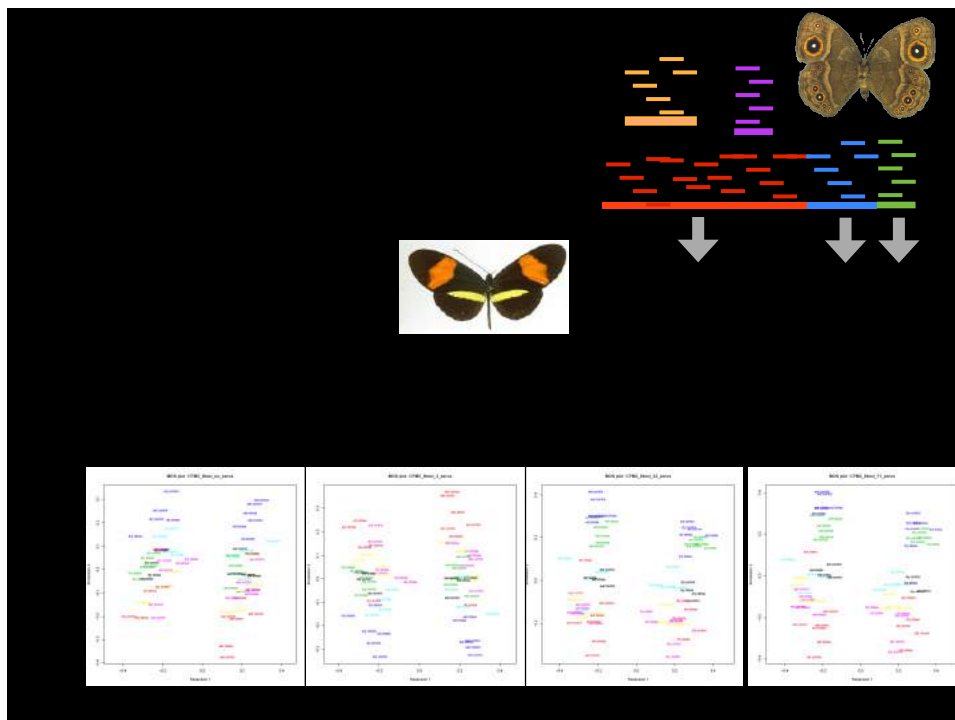
N50 = 930

Hornett & Wheat 2012 BMC Genomics

## Assessing MESPA accuracy



- Input
  - *D. virilis* AA sequences
  - *D. melanogaster* DAS
    - Pool-Seq data (n=50 individuals one population)
    - CLC assembly (kmer = 63, bubble = 2000)
    - N50 = 11,000 (but can work with smaller N50)
- Output
  - Gene models of *D. melanogaster* for putative orthologs

- Assessment:
  - *D. virilis* protein sequences & *D. melanogaster* genome assembled from pooled n=50

With AA set > 60 My divergent from poor genome, MESPA can accurately scaffold > 80% of the length for > 80% of the genes with > 95% accuracy

**Group 1**

**Group 2**

## What's the genetic difference?

In 2015, how should we answer this?
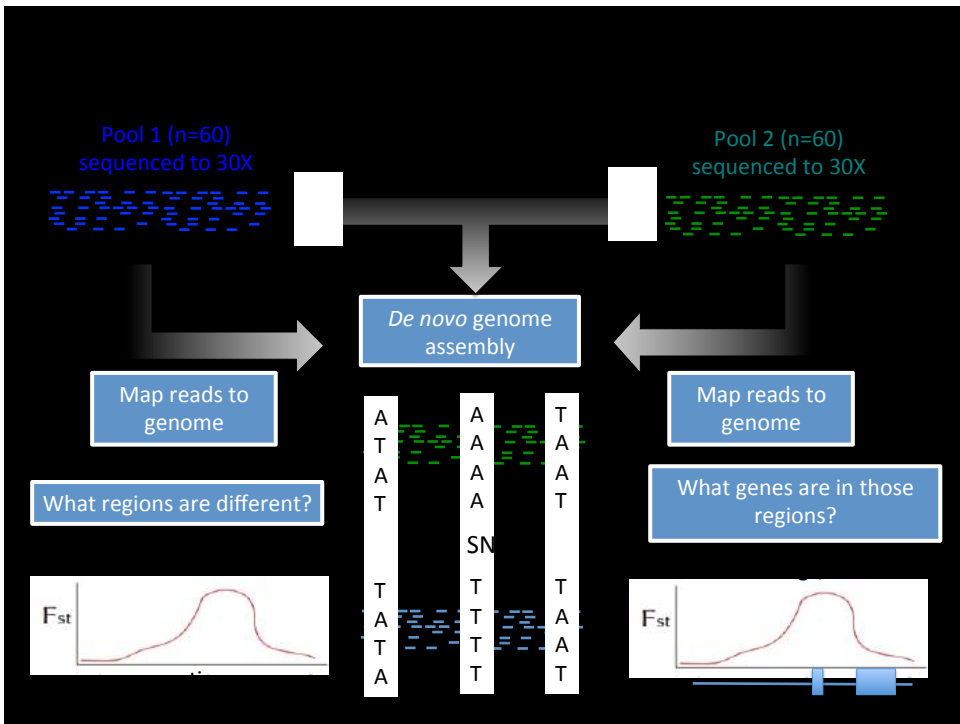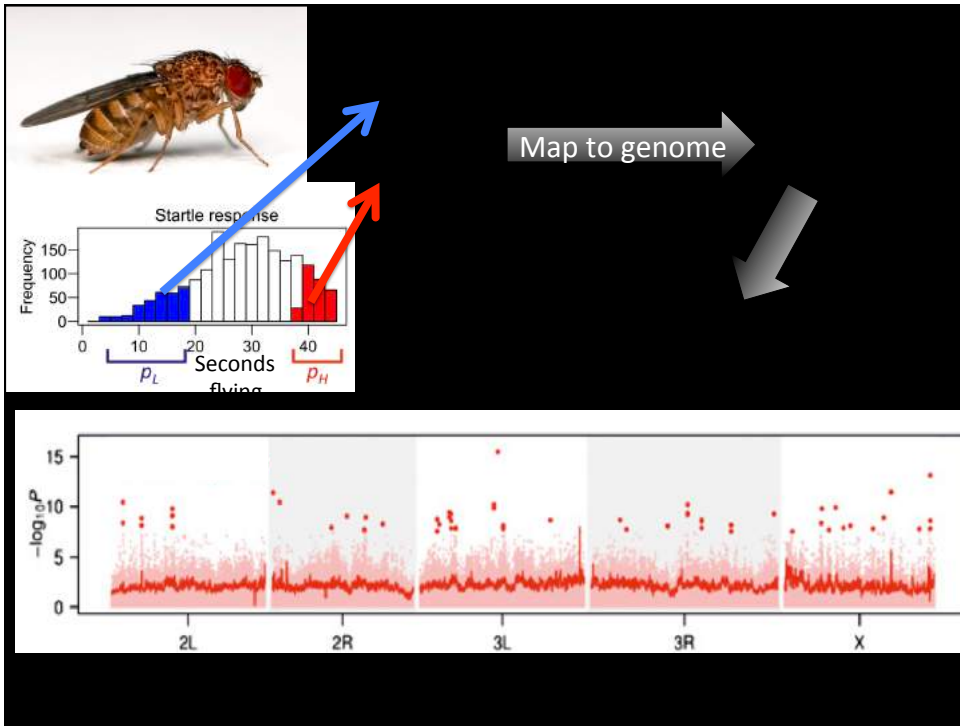
# Just sequence it!

---

**Group 1**

**Group 2**

## What's the genetic difference?

What's the cheapest/easiest experimental design?

- Sequence the be-jesus out of each group
  - >25 X genomic coverage of >50 haploid genomes per group
- Make a simple genome & map this data to it!
- Use good stats to ask what regions are different
- Figure out what those regions are
  - Invest your resources in these regions and their functional role

Map to genome



Pool 1 (n=60) sequenced to 30X

Pool 2 (n=60) sequenced to 30X

De novo genome assembly

Map reads to genome

Map reads to genome

What regions are different?

What genes are in those regions?

$F_{st}$

$F_{st}$

# Can this really work?

# Case study # 1

## LETTER

### *doublesex* is a mimicry supergene

K. Kunte[1]*, W. Zhang[2]*, A. Tenger-Trolander[2], D. H. Palmer[3], A. Martin[4], R. D. Reed[4], S. P. Mullen[5] & M. R. Kronforst[2,3]

One of the most striking examples of sexual dimorphism is sex-limited mimicry in butterflies, a phenomenon in which one sex—usually the female—mimics a toxic model species, whereas the other sex displays a different wing pattern[1]. Sex-limited mimicry is phylogenetically widespread in the swallowtail butterfly genus *Papilio*, in which it is often associated with female mimetic polymorphism[1–3]. In multiple polymorphic species, the entire wing pattern phenotype is controlled by a single Mendelian 'supergene'[4]. Although theoretical work has explored the evolutionary dynamics of supergene mimicry[5–9], there are almost no empirical data that address

pattern. However, Clarke and Sheppard found virtually no evidence for recombination in *P. polytes*[13], although they did recover apparently recombinant phenotypes in other species, such as *P. memnon*[14]. Over the past few decades, supergene mimicry has received considerable theoretical attention[5–9], but there are almost no empirical data that address the molecular basis of a supergene. One example from *Heliconius* butterflies, which involves supergene mimicry but not sexual dimorphism, suggests that supergenes may be the result of chromosomal inversions that lock multiple adjacent genes into a single, non-recombining unit[15].

Kunte et al. 2014 Nature

## Polymorphic, sex-limited mimicry



K Kunte et al. Nature **000**, 1-4 (2014) doi:10.1038/nature13112

Non-mimetic female forms

Mimetic female forms

*cyrus*

*polytes*

N=15

N=15

Mapped reads to *de novo* genome



N=30
60 X coverage

N=30
60 X coverage
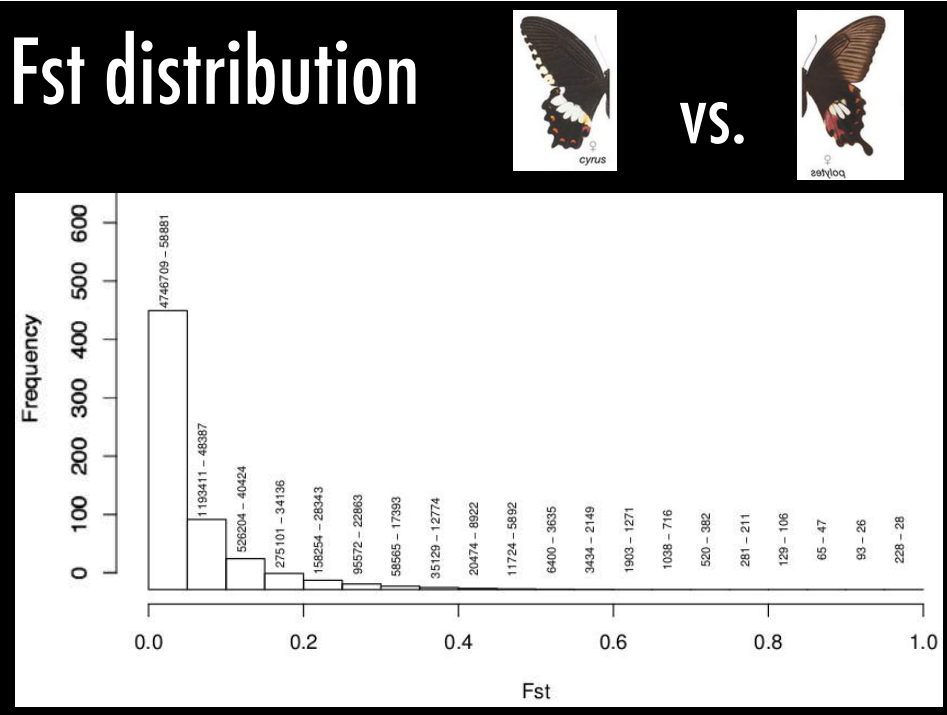
*De novo* genome assembly

Map reads to genome

Map reads to genome

What regions are different?

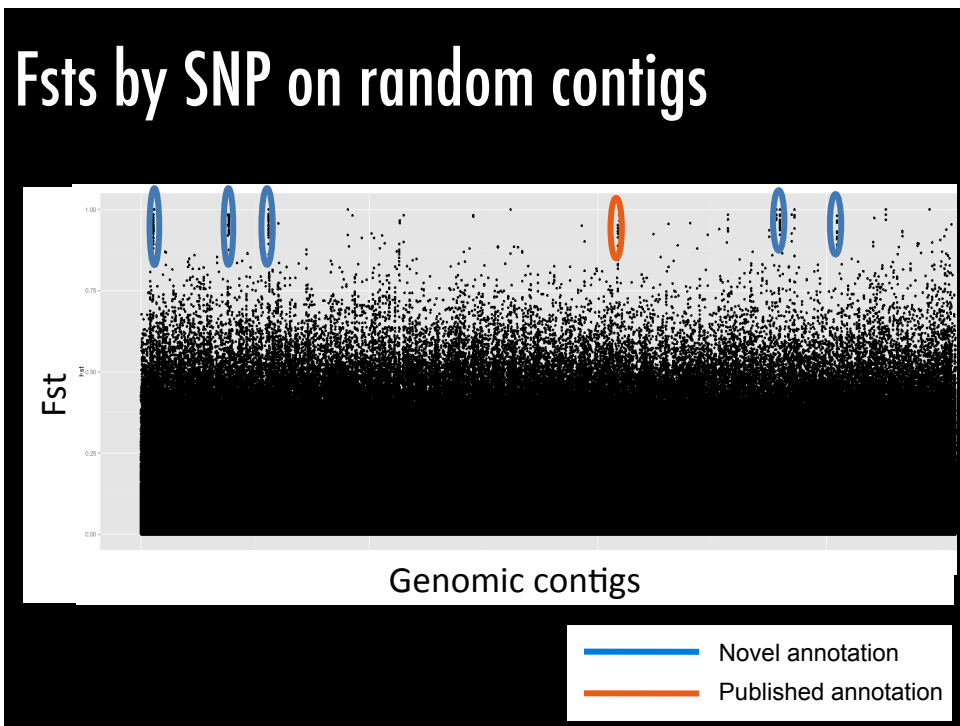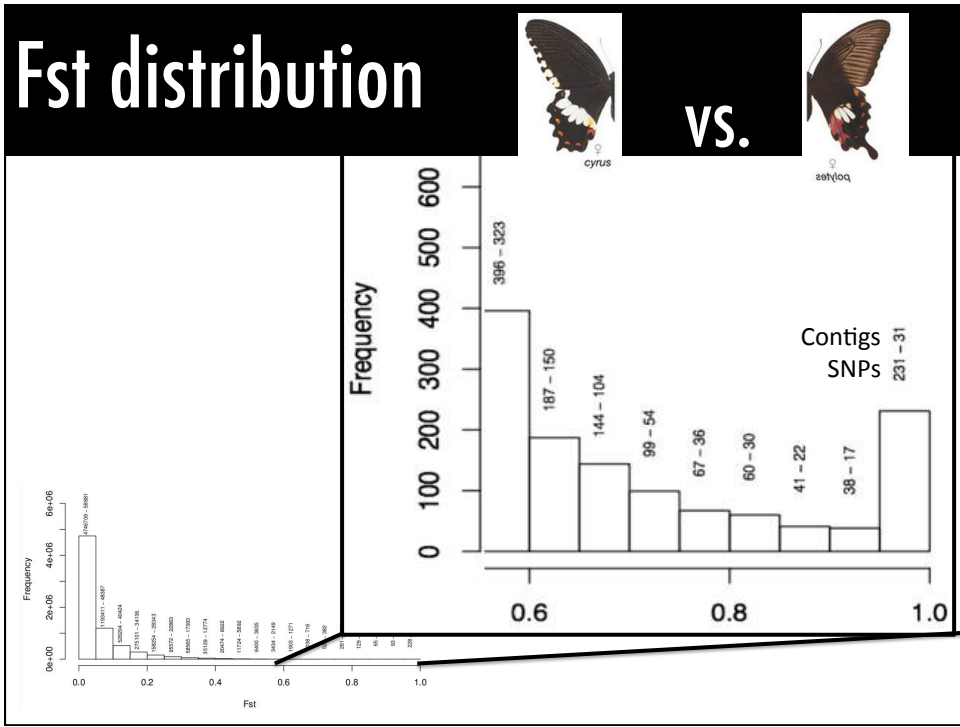What genes are in those regions?

$F_{st}$

$F_{st}$

A T A T

A A A A SN T T T T

T A A T

T A T A

T T T T

T A A T

# Can we find the same genomic regions?



# Fst distribution    vs.

# Fst distribution

vs.



# Fsts by SNP on random contigs

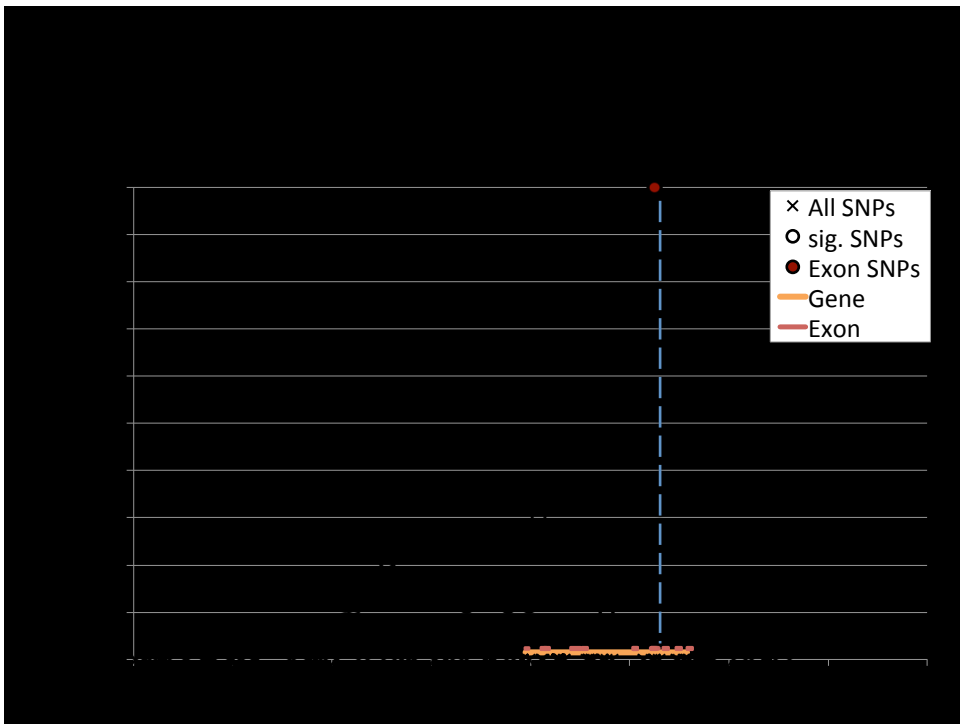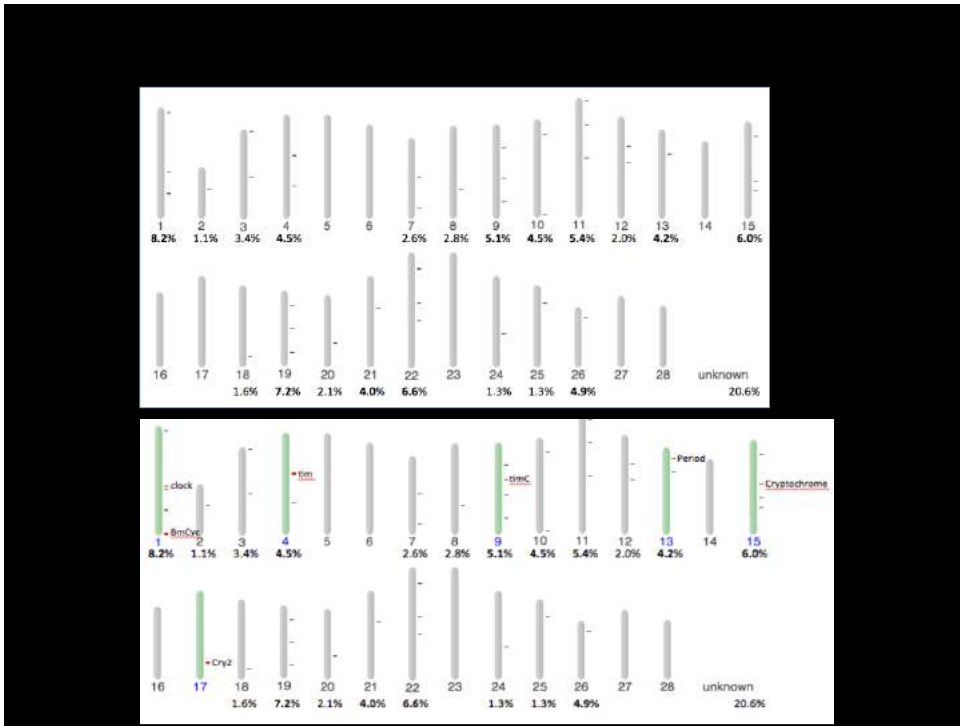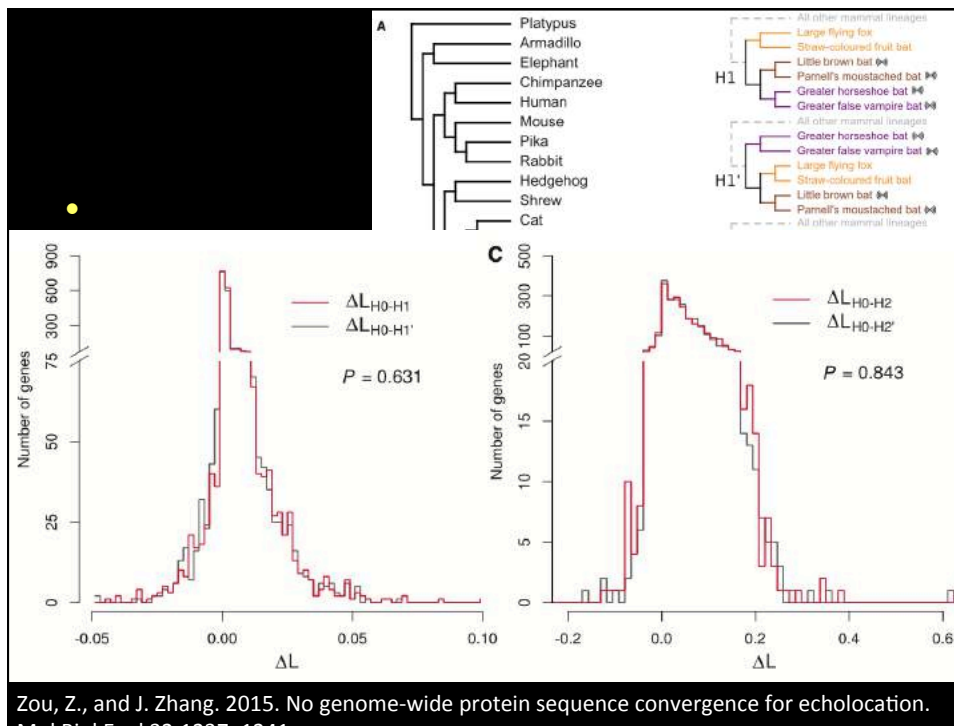Genomic contigs

Novel annotation
Published annotation

# Results

- Identified *doublesex* immediately

- Found new genes missed in pubication
  - p270
  - RNA directed DNA polymerase
  - Arginine/Serine rich coiled-coil protein 2

- Now searching to see if these are all near each other

Zou, Z., and J. Zhang. 2015. No genome-wide protein sequence convergence for echolocation.

---

# *De novo* RNA-Seq: Do you need a genome?

**No, but there are important biases & limitations**
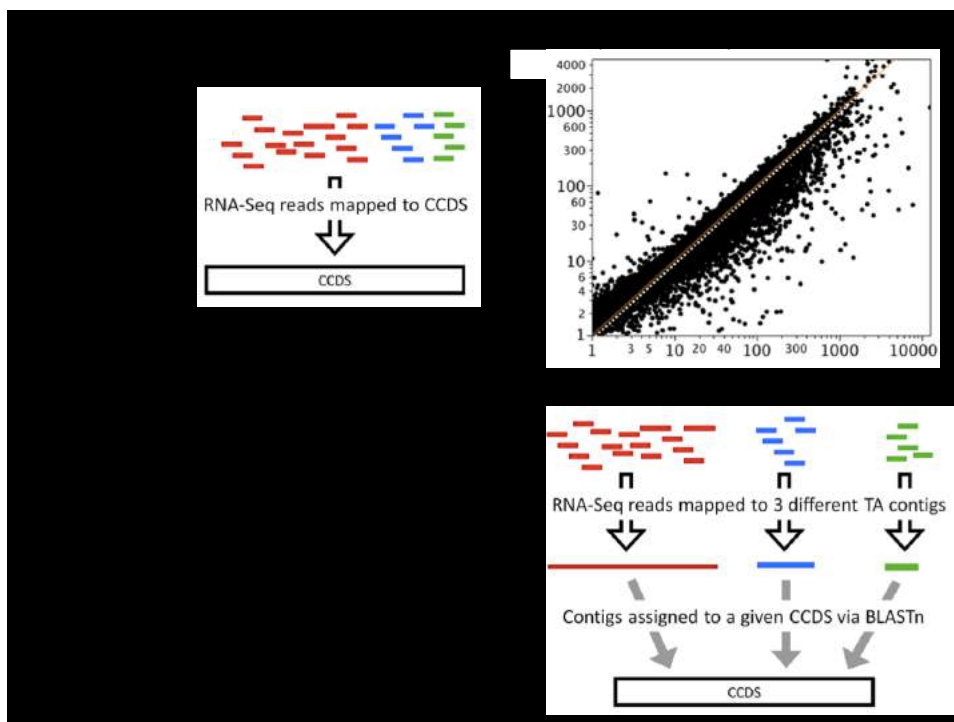
- **TA mapping limitations**
  - No exon level resolution but this will change soon
  - No coding information on identified SNPs unless you build gene feature files on contigs
- **TA mapping biases unique to it**
  - Spicing may cause mapping problems if locus is collapsed, but generally OK to not assume a gene model
- **TA mapping biases shared with genomic mapping**
  - SNP and indel effects
  - gene duplication (are reads mapping to the right place)

# Map to TA vs. Genome:
## which is better?

Template effects:

- Mismatch :
  - SNPs (single nucleotide polymorphisms)
  - Indels (insertion or deletion polymorphisms)

- Pseudo-inflation
  - An increase in the copy number of a gene that arise from genome assembly errors or TA errors

- Gene model errors
  - If the models in your genome are bad, this will affect results
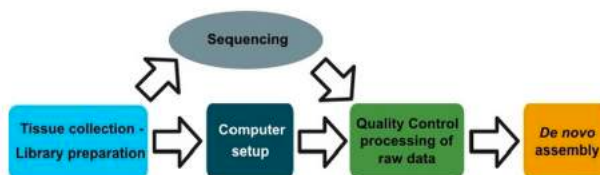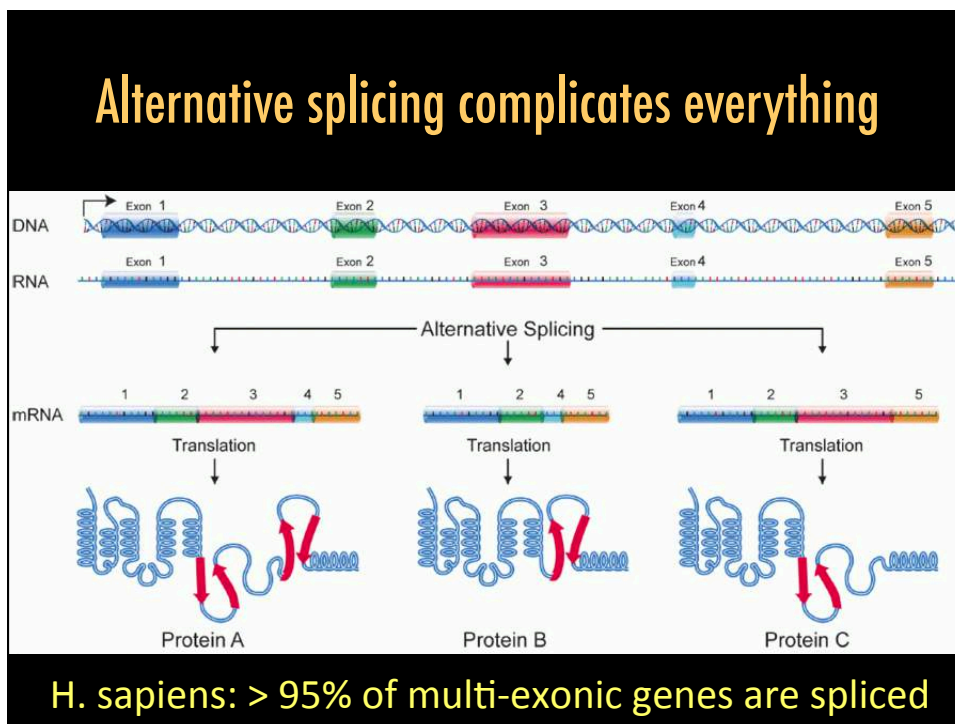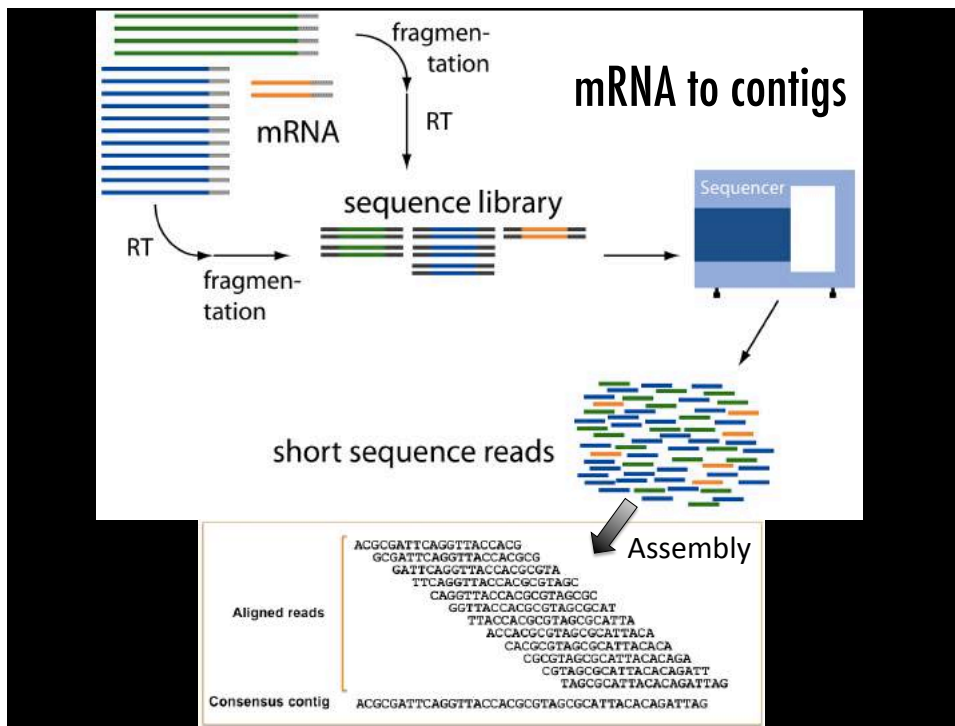
## Duplication levels in RNA-Seq data

- Common in transcriptome work

- Starting with lots of high quality RNA increases
  - mRNA amount for sequencing
  - Decreases need of core facility to PCR your sample

- Moderate amounts of PCR duplication are OK
  - ~ 20% expected
  - > 50% perhaps problematic if correlated with experimental design
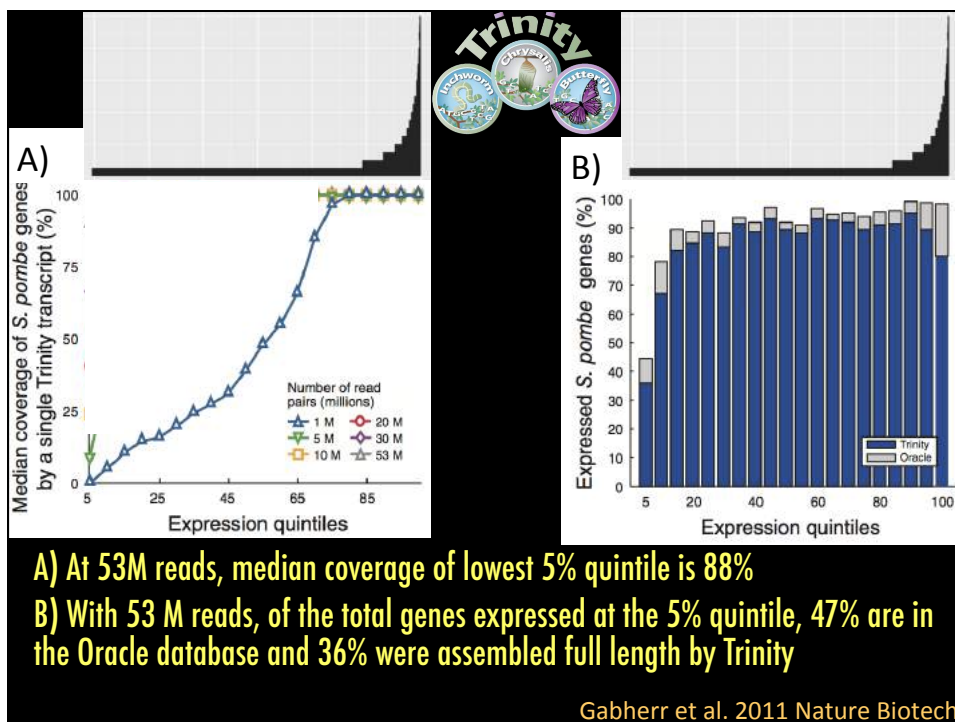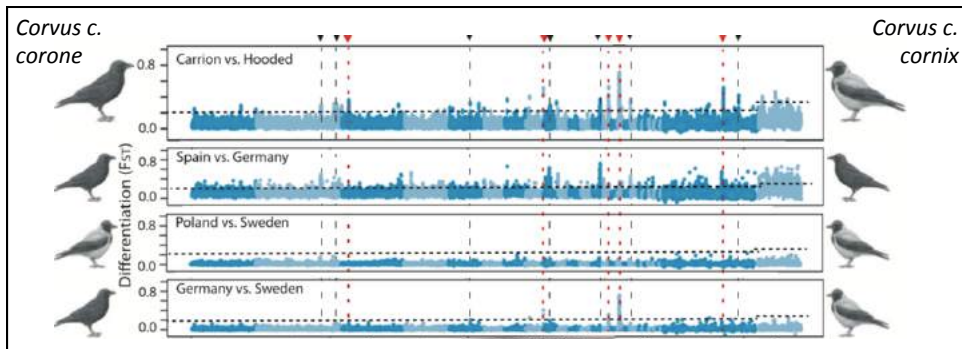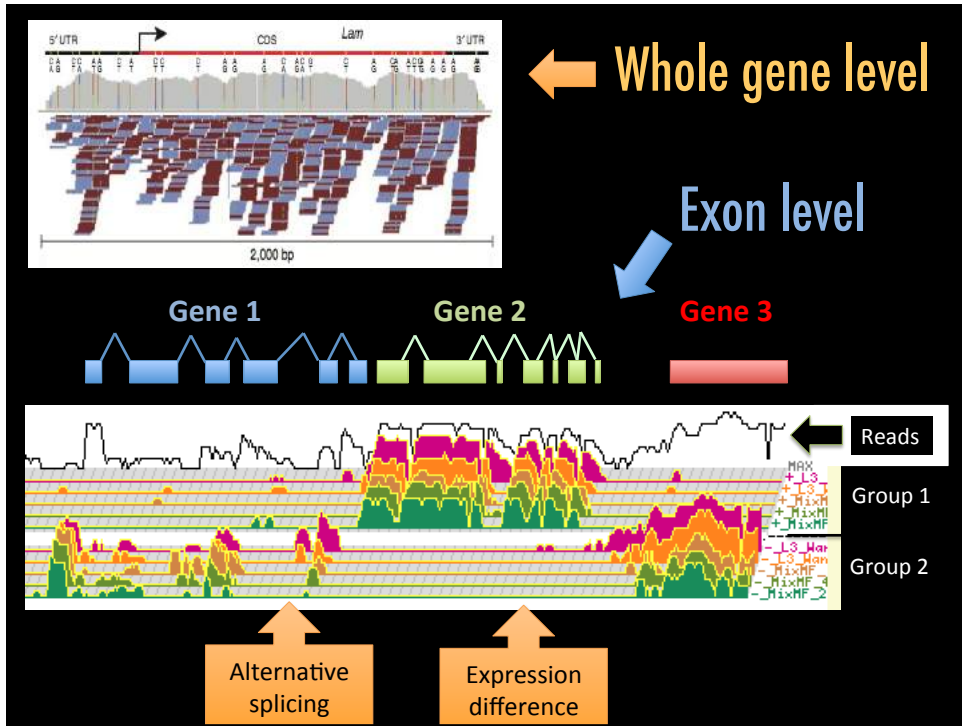  - Clone_filter program in STACKS is excellent assessing this

# Pipeline Overview

mRNA to contigs



Alternative splicing complicates everything

H. sapiens: > 95% of multi-exonic genes are spliced

*De novo* transcritpome assembly

Reconstructs splice isoforms using PE Illumina data



A) At 53M reads, median coverage of lowest 5% quintile is 88%
B) With 53 M reads, of the total genes expressed at the 5% quintile, 47% are in the Oracle database and 36% were assembled full length by Trinity

Gabherr et al. 2011 Nature Biotech

# Islands of speciation or background selection?



$D_{xy}$:
An absolute measure of differentiation, increase due to mutations

Fst:
A relative measure of differentiation, increases due to freq. change

The absence of high Dxy in regions of high Fst suggest a role of background selection driving these patterns rather than genomic 'islands' driving speciation.

Cruickshank and Hahn. 2014. Molecular Ecology.



Tree Height = Tan $a$ X Distance

$a$

Distance

Is it an adaptation?



*Colias croceus*, the Clouded Yellow

But what causes Alba?

GWAS + genome + QTL mapping (blood, sweat, tears)

BarH1 gene