

# Ecological Genomics, pt. 2

you, your data, your perception and  
the hard realities

Christopher West Wheat



# Gordon Freeman



# Informatics and Biology

- We need to make sure we put the 'bio' into the bioinformatics
  - Do results pass 1<sup>st</sup> 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
  - Assessing assemblies, mapping, and expression
  - What is validation?
- Insights from candidate genes
  - Can Second Gen methods get us there?

# Core facilities and non-model species

Commonly heard statements that are not true:

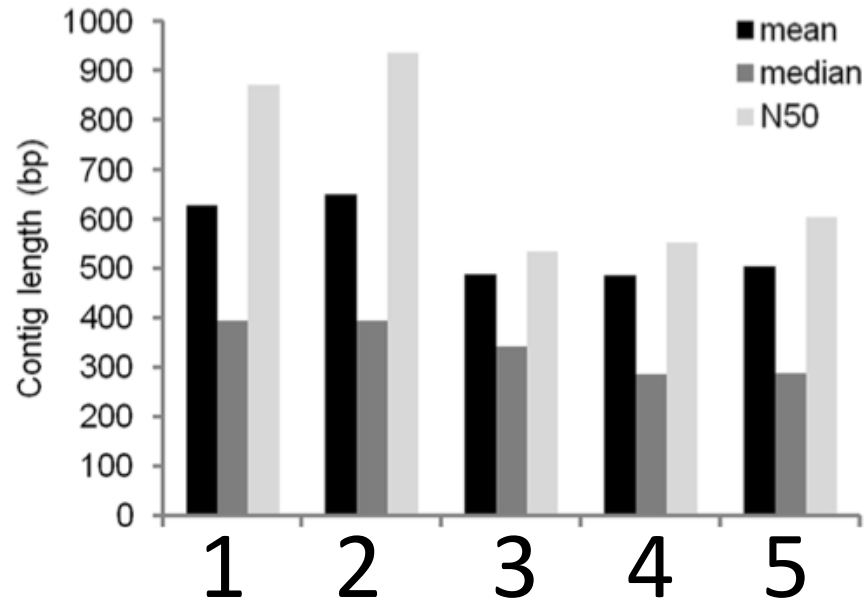
- You can't do RNA-Seq without a genome
- The best metric for Transcriptome Assembly assessment is N50 & # of contigs
- We'll have your data back in < 1 month

# Assessing transcriptome assembly

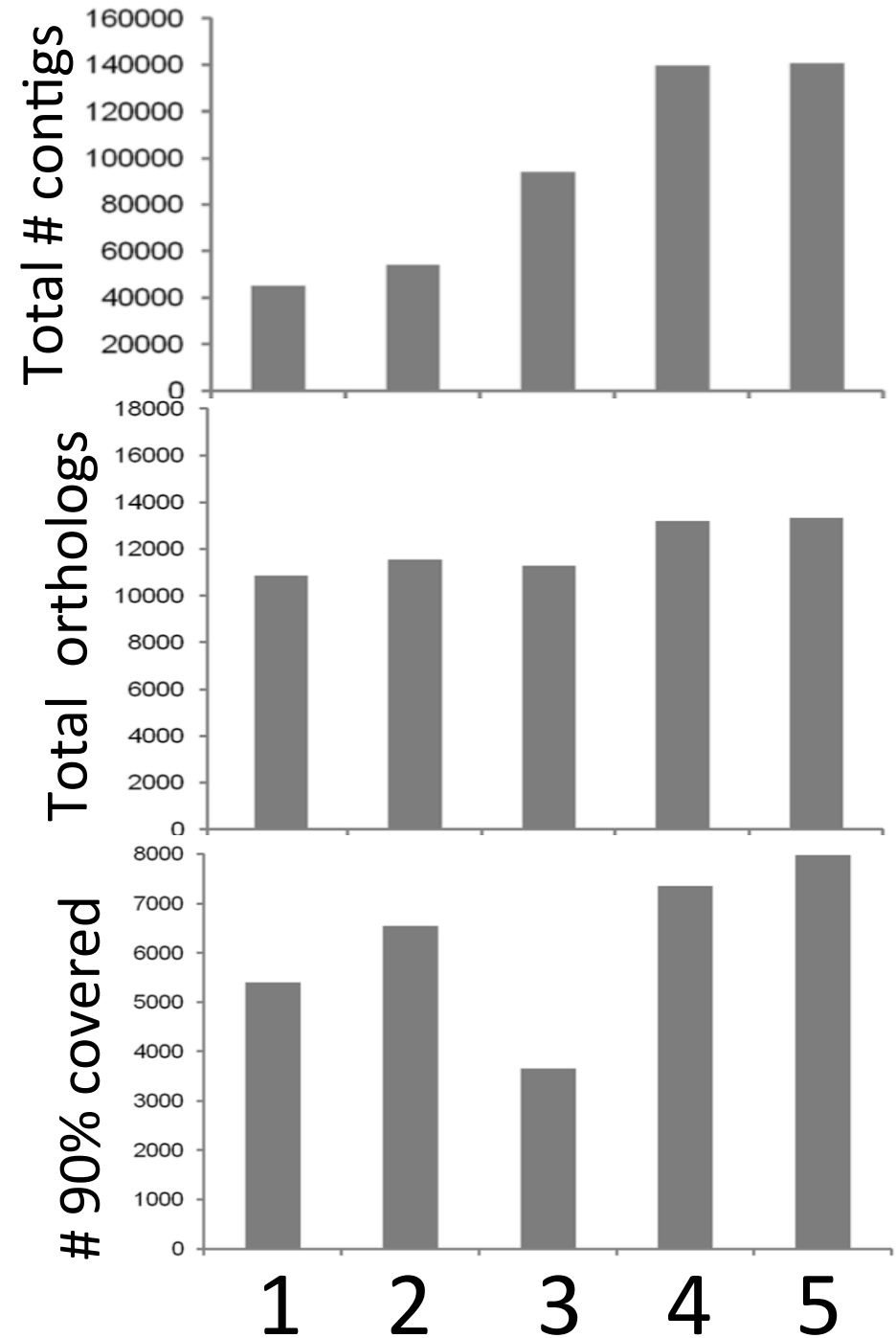
- Assessment metrics
  - Non-biological
    - N50, # of contigs
  - Biologically informative
    - # of orthologs identified
    - Ortholog hit ratio (OHR)

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

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

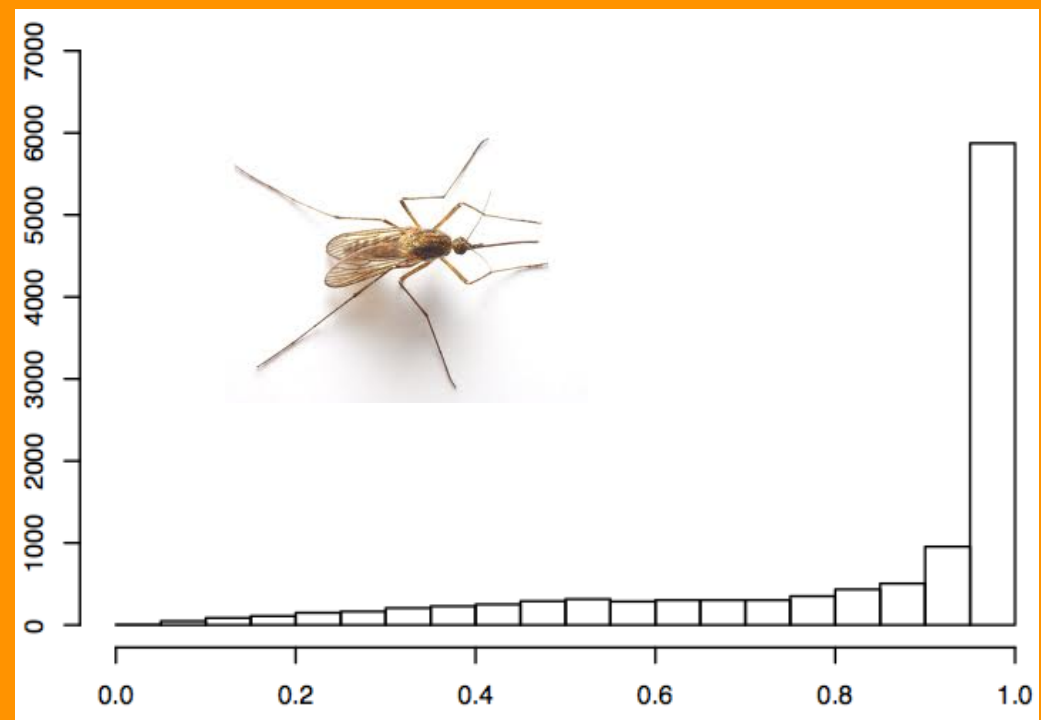
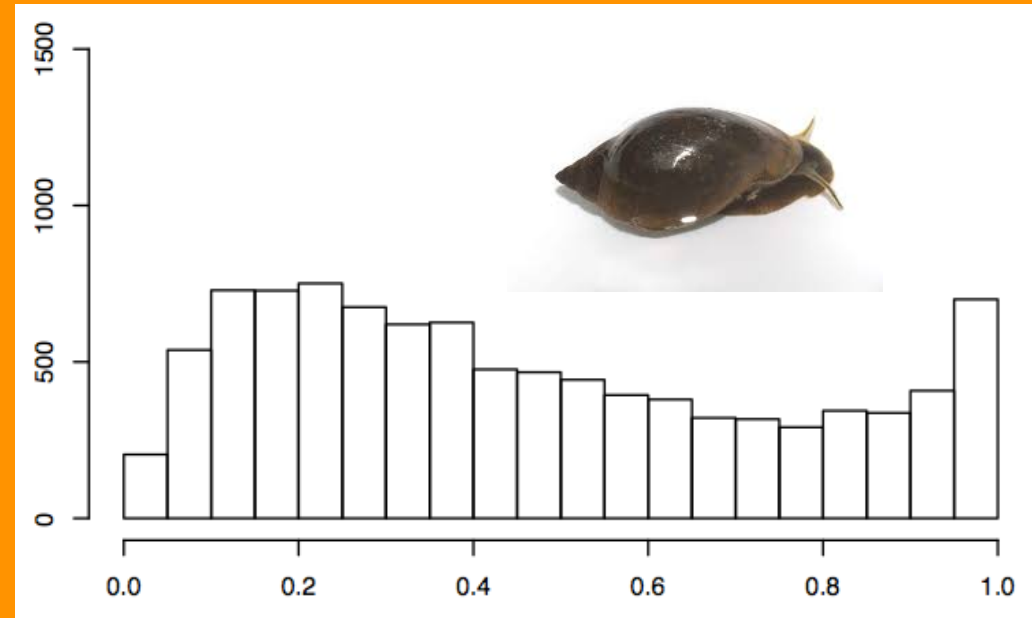


- 5 different TAs
- TA 2
  - Best N50, fewest contigs



# OHR graphs

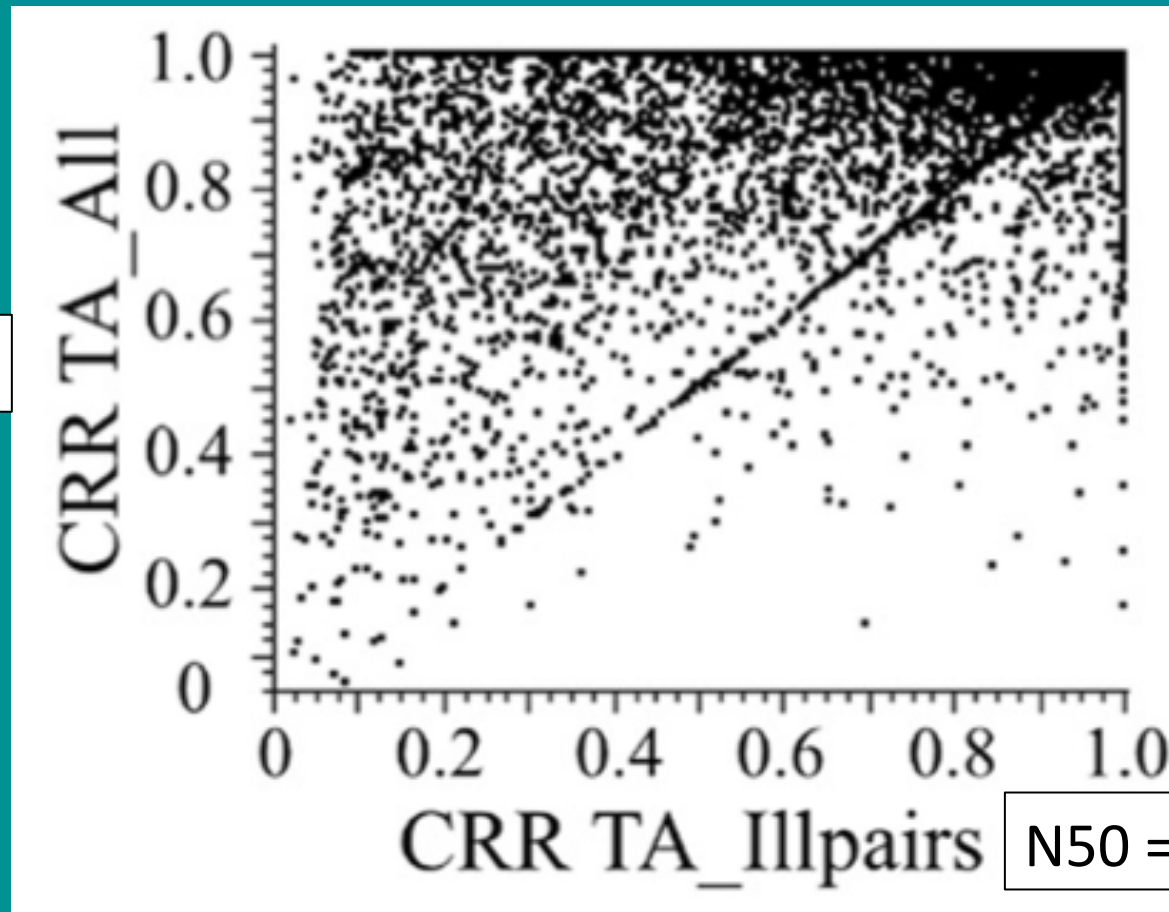
- Shows the number of unique orthologs hit
- Distribution of their reconstructed length





# Comparative OHR

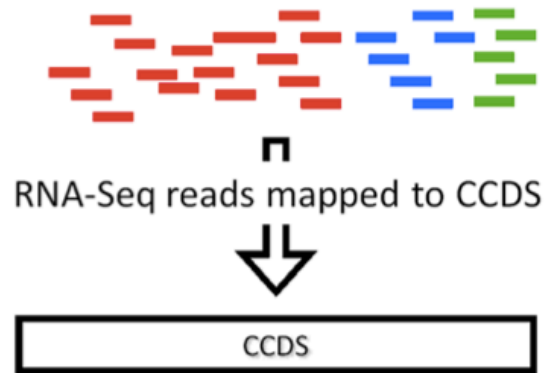
- Compare longest contig per ortholog for two assemblies
- Plot them against each other



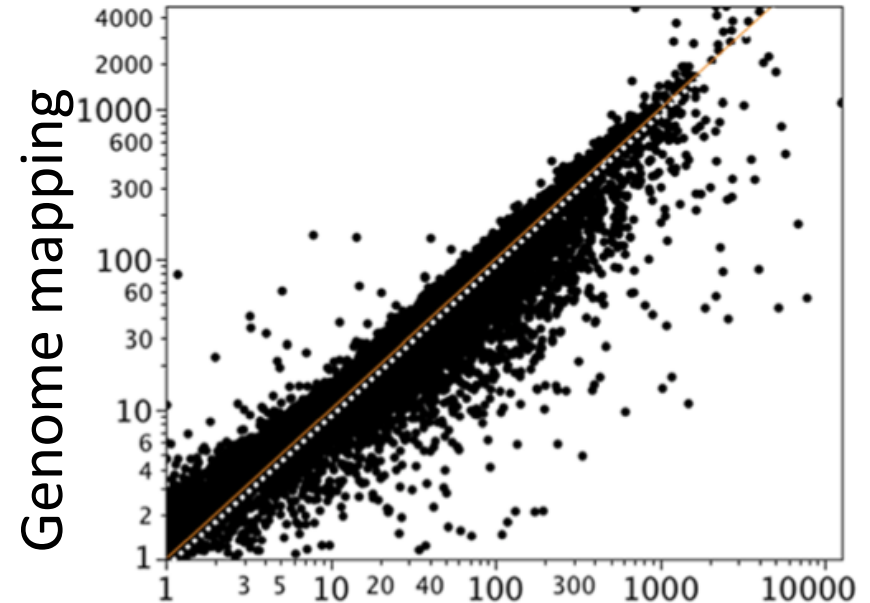
N50 = 610

N50 = 930

## Genome mapping

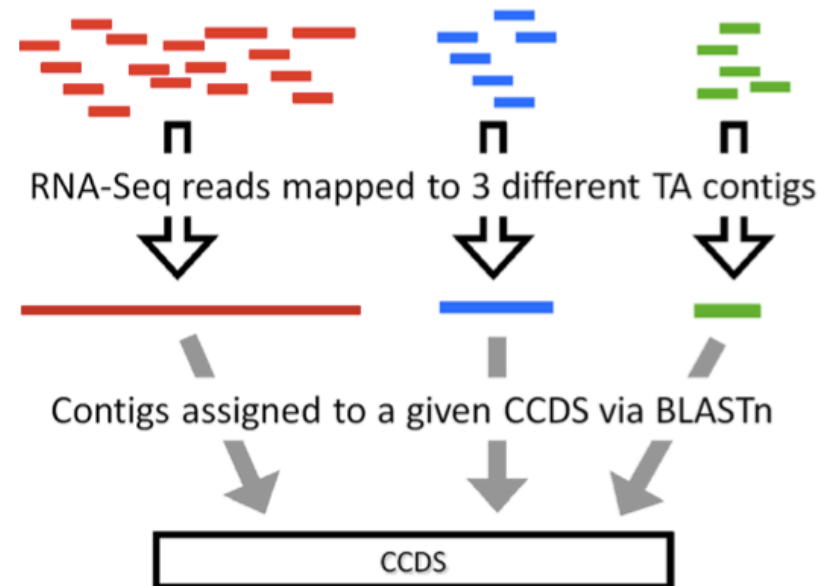


Spearman's  $\rho = 0.95$ ,  $P < 0.0001$



## RNA-Seq mapping: comparing genome vs. TA

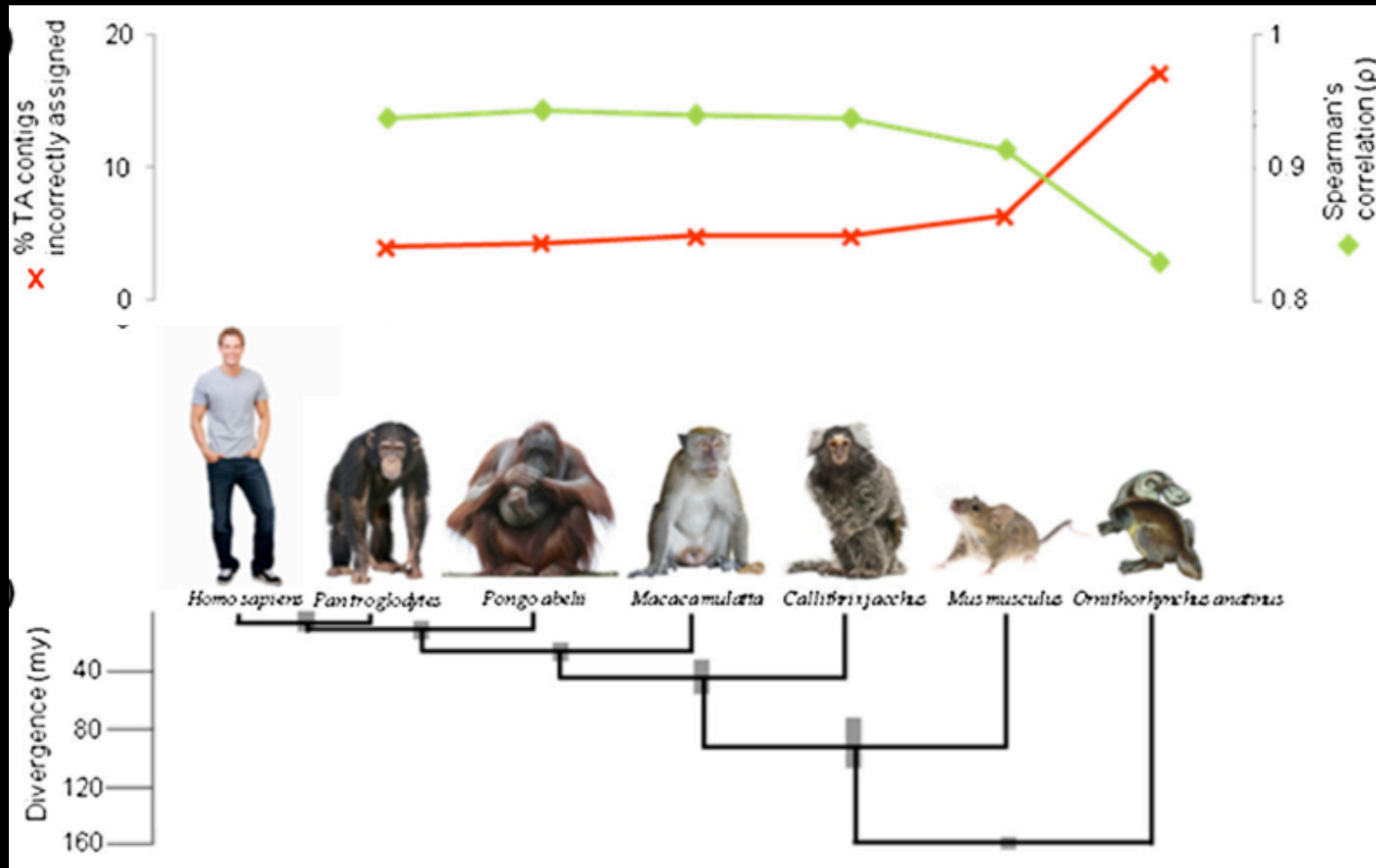
You can generate high quality data without a genome, for much of the transcriptome



# OHR can be calculated using predicted genes from divergent species



Emily  
Hornett



Hornett and Wheat 2012

# RNA-Seq



- Now we have a good assembly
- Ready for quantitative gene expression analysis
- 2 factor analysis with family effects

# *Bicyclus anynana*



**long** lifespan

**short**

**delayed** reproduction

**fast**

**inactive** behaviour

**active**

**high** fat reserves

**low**

**cryptic** wing pattern

**conspicuous**

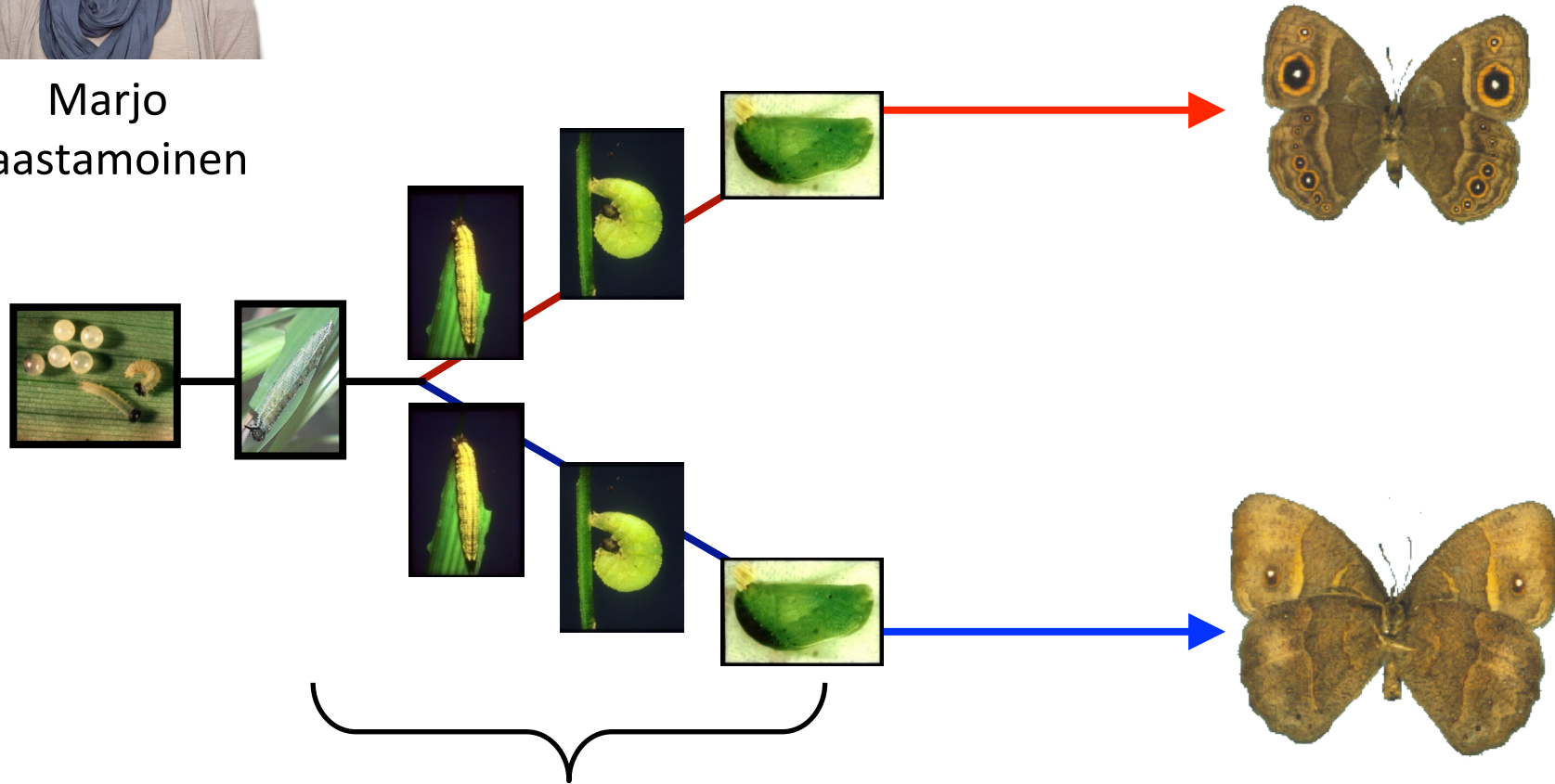
**Save  
energy,  
live long**

**Live  
fast,  
die  
young**



Marjo  
Saastamoinen

# *Bicyclus anynana*

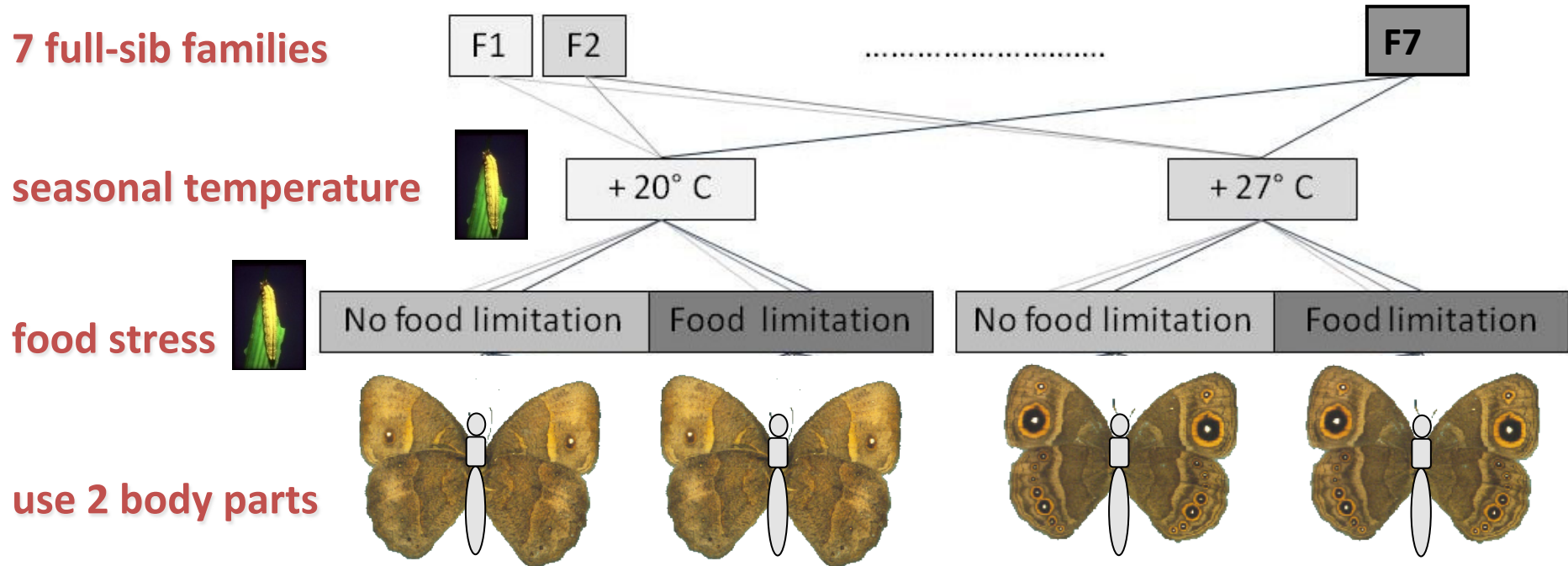


*sensitive period*

**environmental  
conditions**

**alternate  
phenotypes**

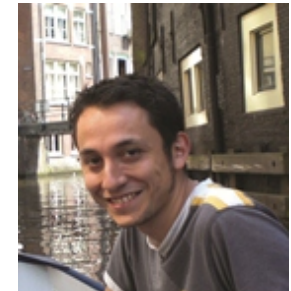
# Experimental design



- 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

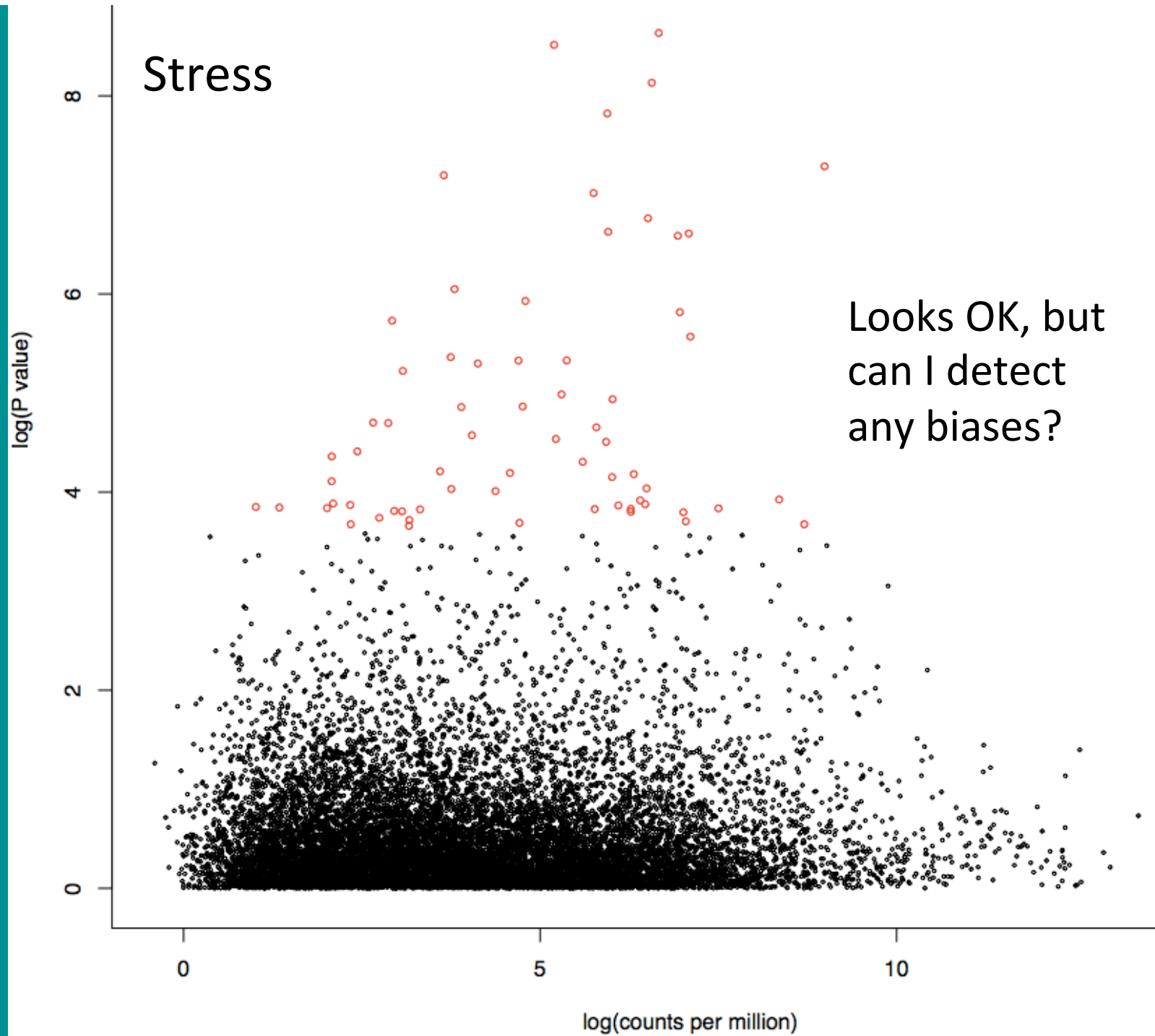
# edgeR



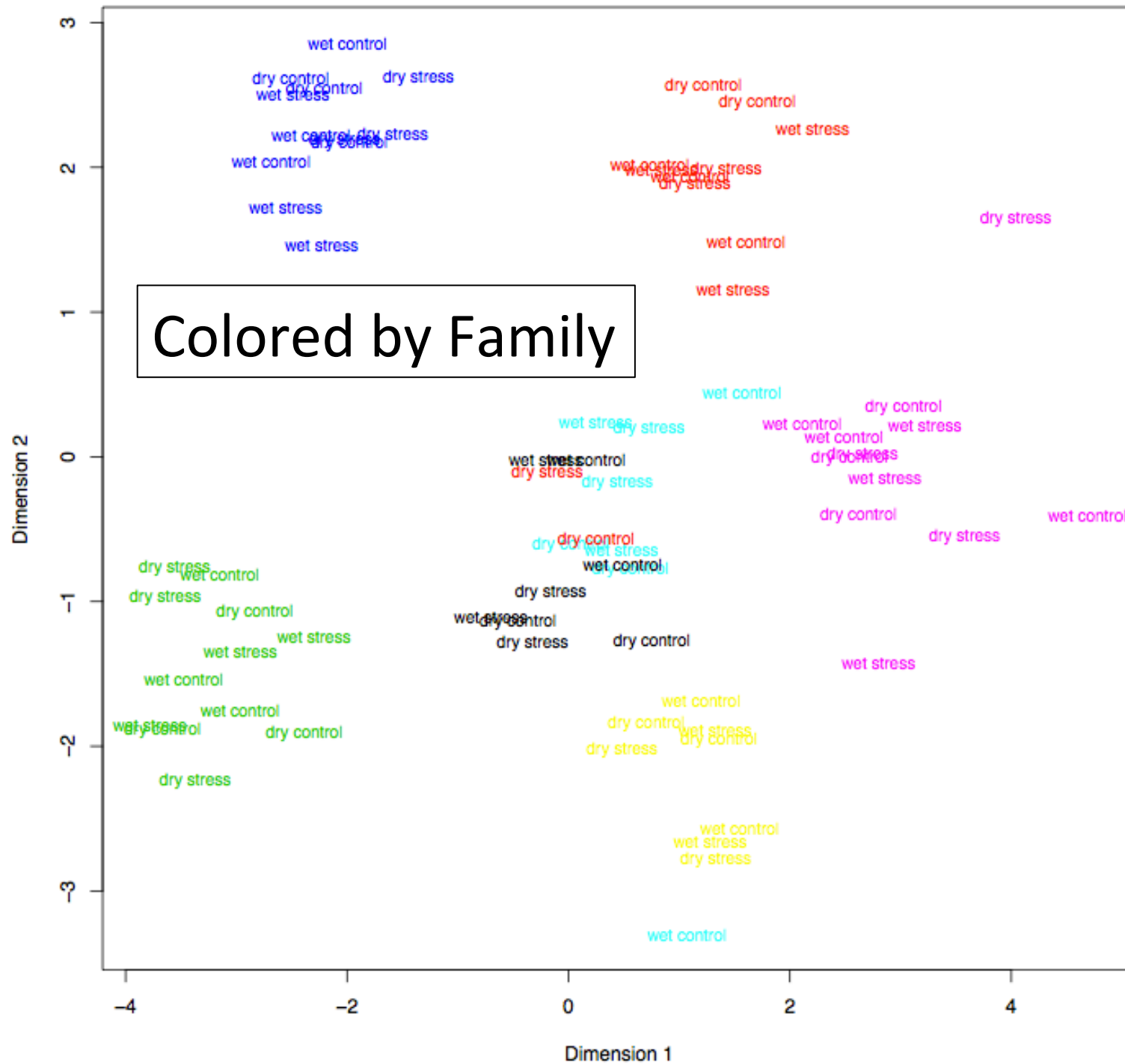
# reads ~ season + stress + family +  
 season\*stress + season\*family + stress\*family  
 season\*stress\*family

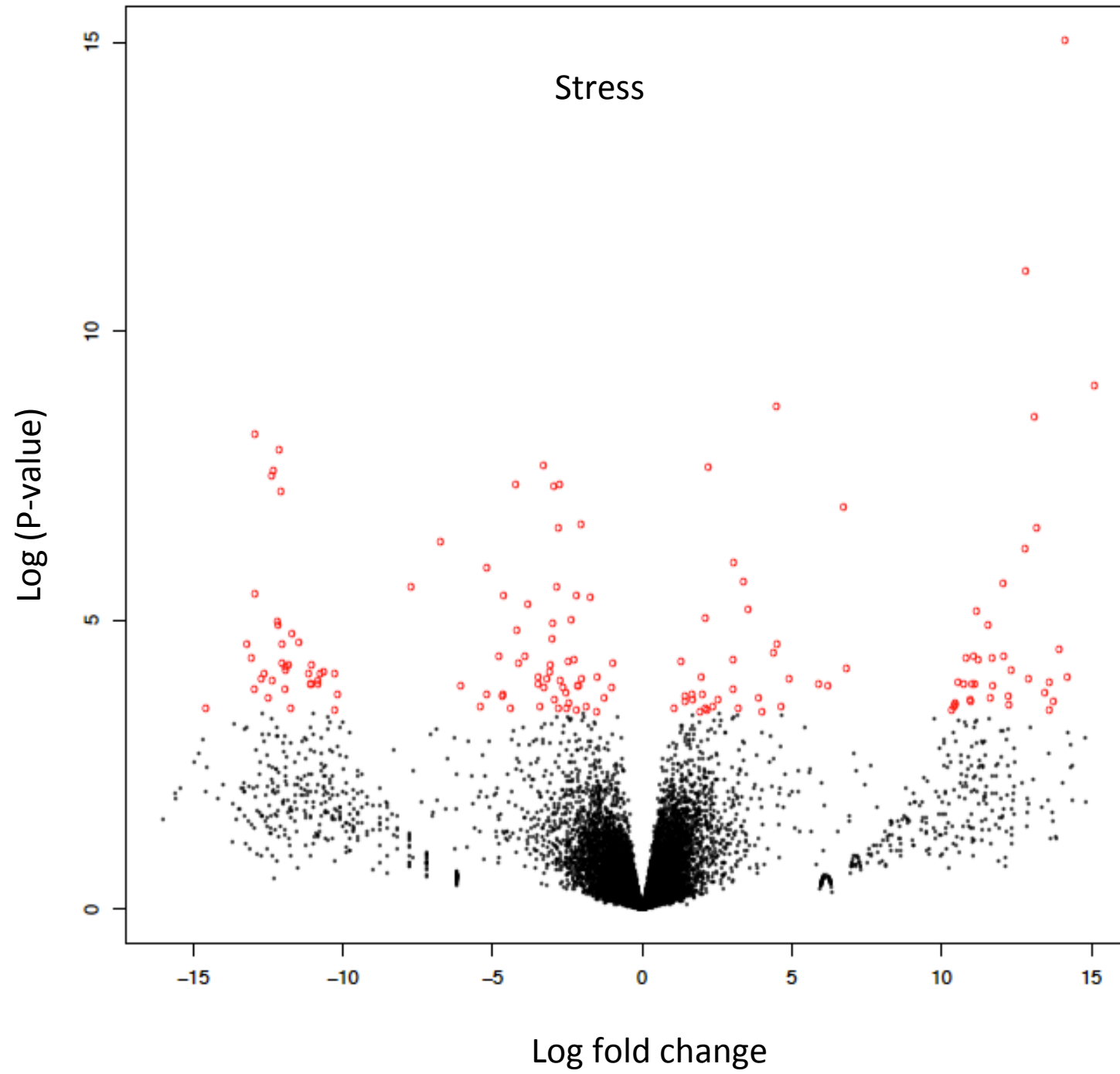


Stress



Looks OK, but  
can I detect  
any biases?



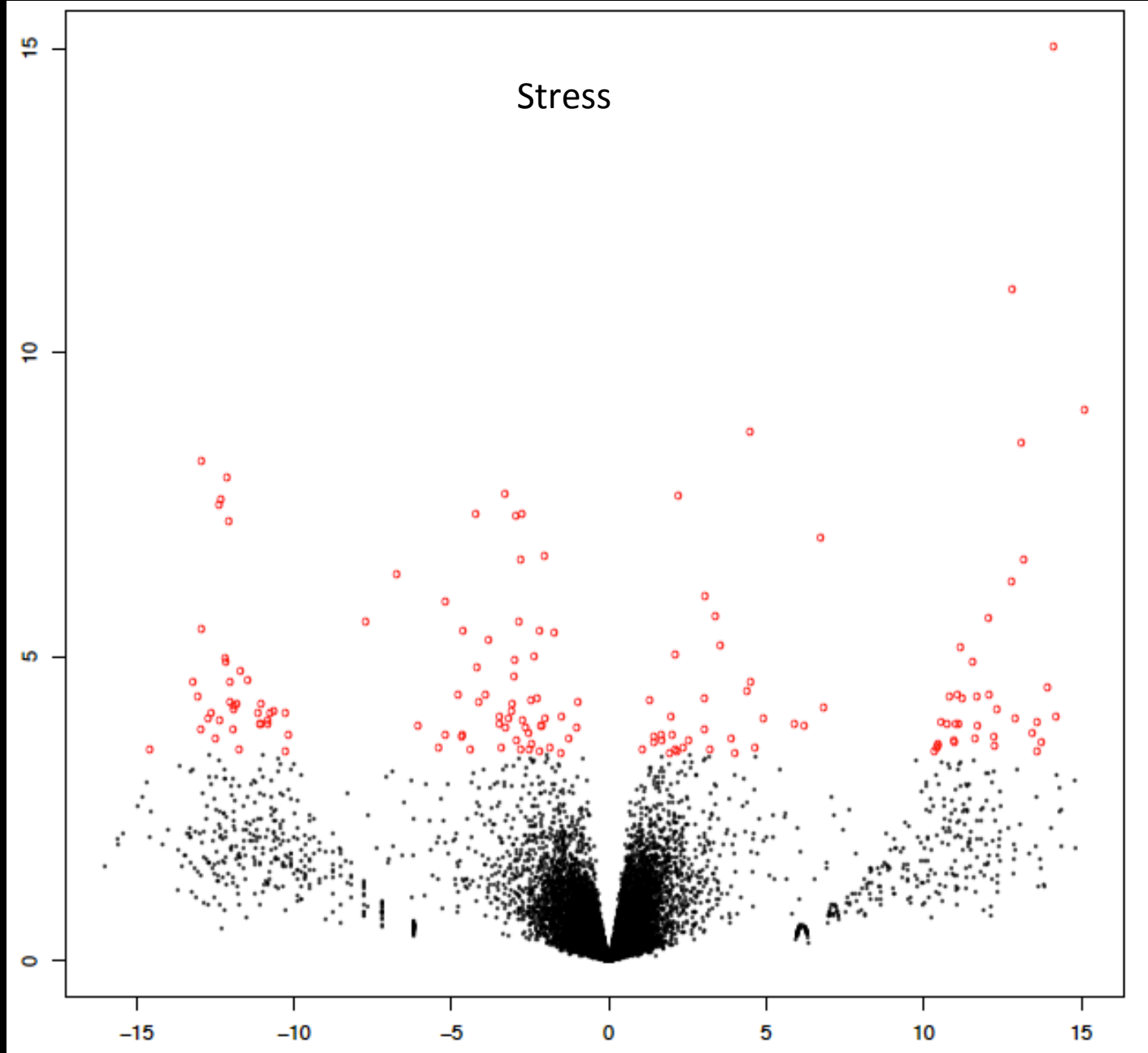


Log (P-value)



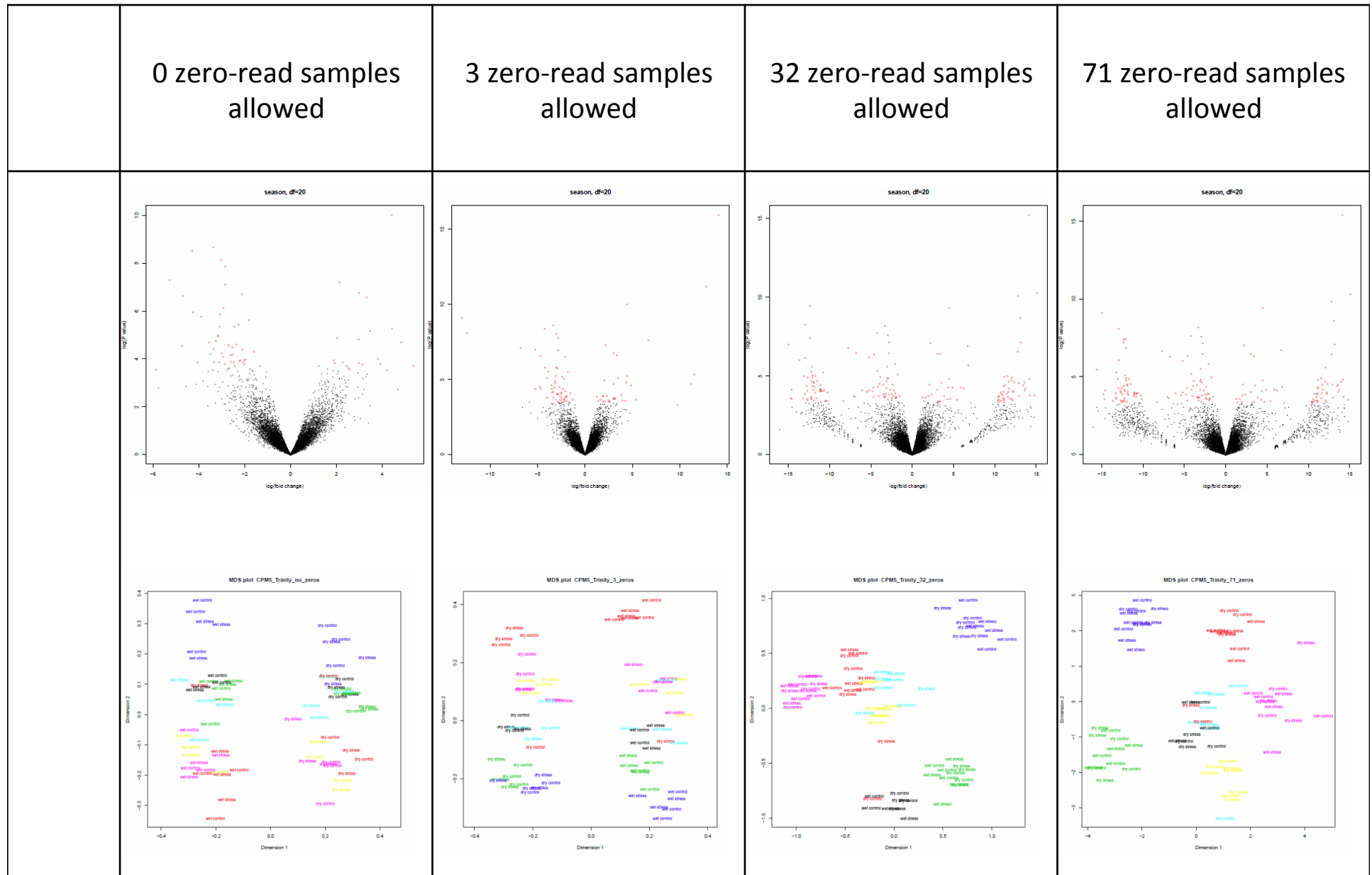
Log fold change

Log (P-value)



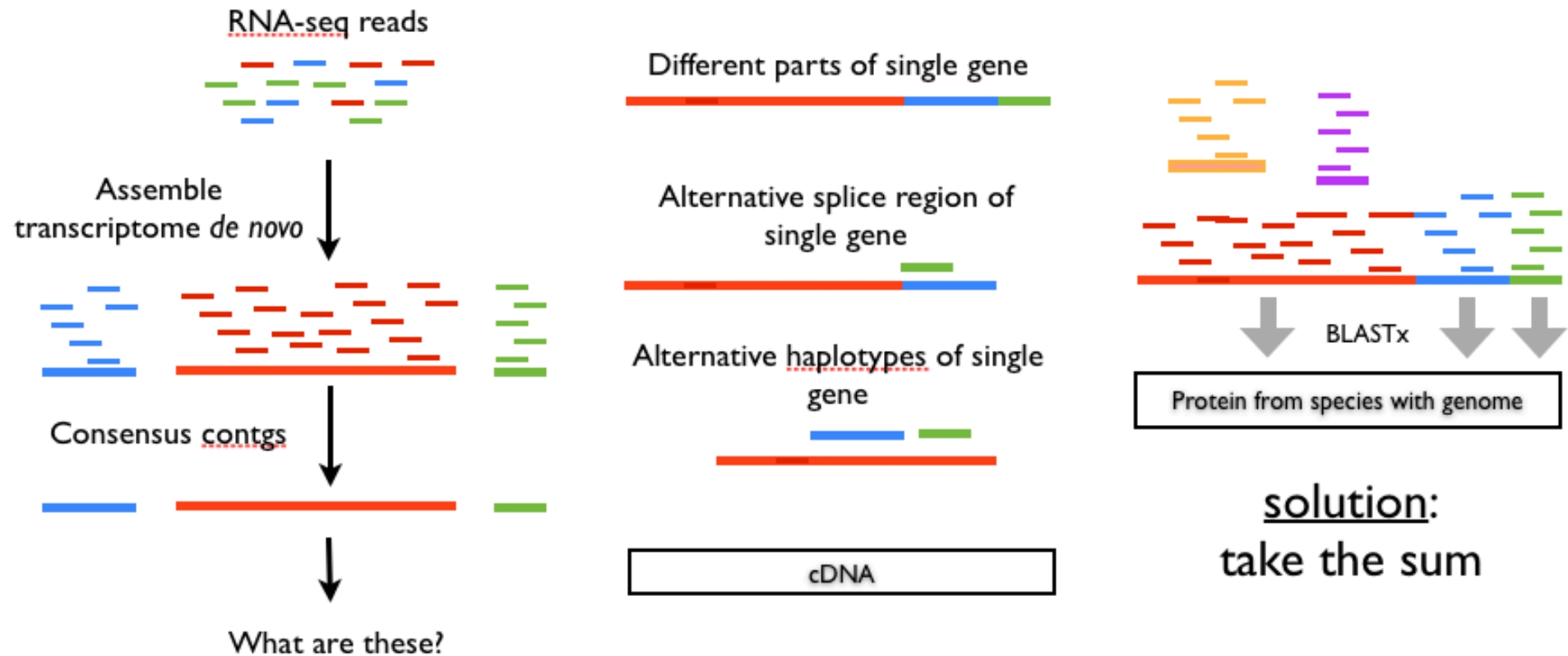
Log fold change

# Effect of filtering, mapping to Trinity contigs



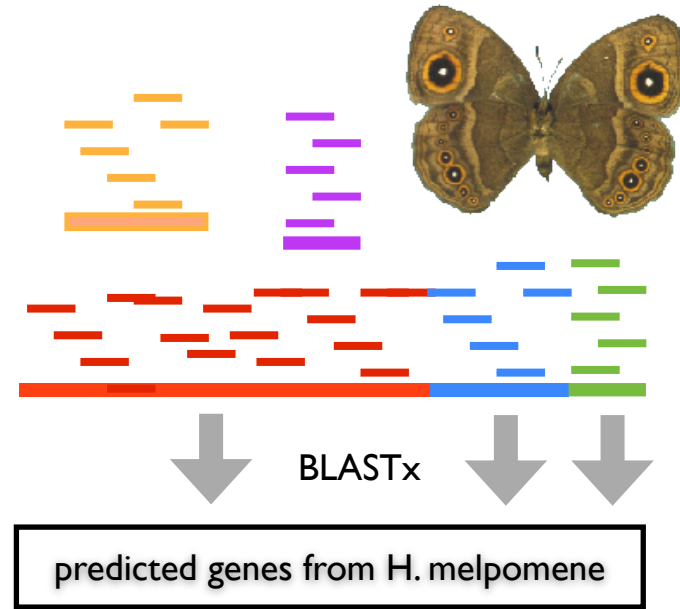


# What's happening?



Separate contigs made during assembly: SNPs X splicing  
Creates bias in expression pattern, with large family effect  
Summing by ortholog corrects this bias

# Effect of filtering when using sum method

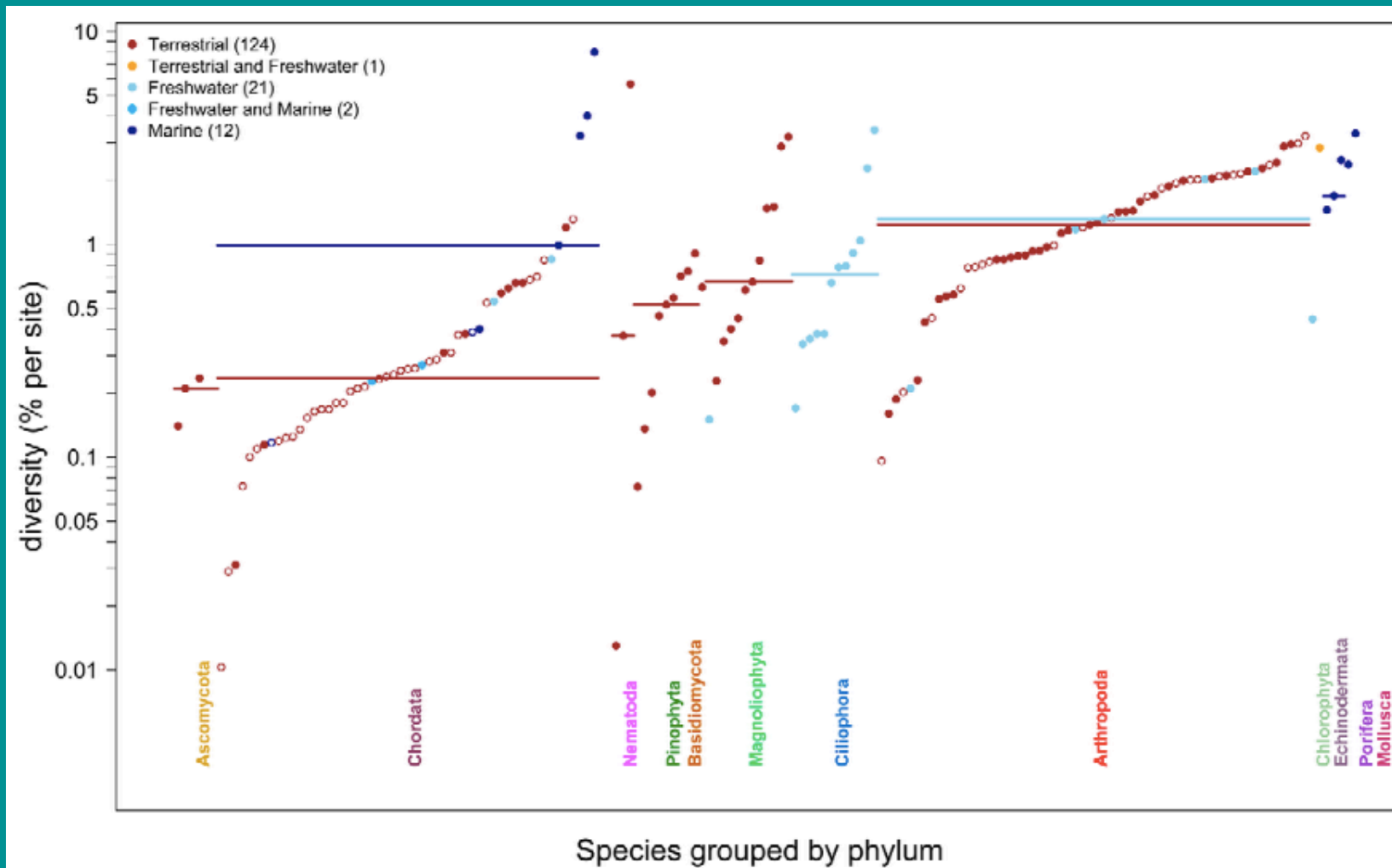


	0 zero-read samples allowed	3 zero-read samples allowed	32 zero-read samples allowed	71 zero-read samples allowed
CPM > 5	<p>MDS plot CPMS_Hmel_no_zeros</p>	<p>MDS plot CPMS_Hmel_3_zeros</p>	<p>MDS plot CPMS_Hmel_32_zeros</p>	<p>MDS plot CPMS_Hmel_71_zeros</p>

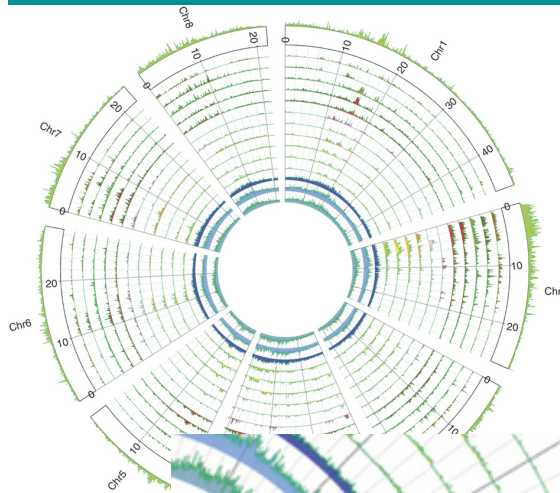


# Mapping reads in outbred species

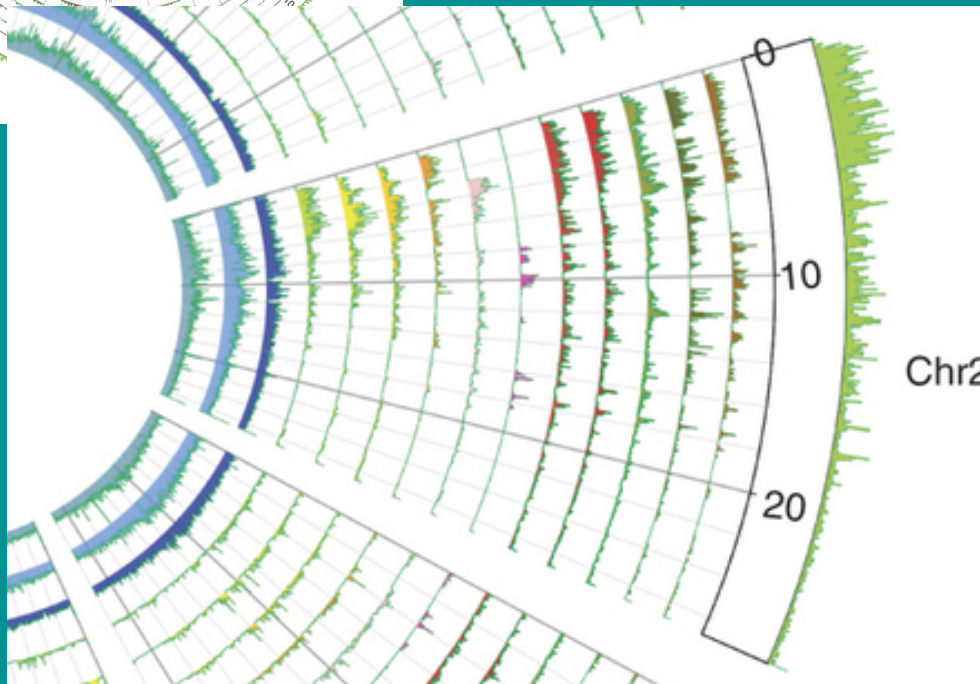
## Average genome polymorphism levels



# Is the mean where you want to look?



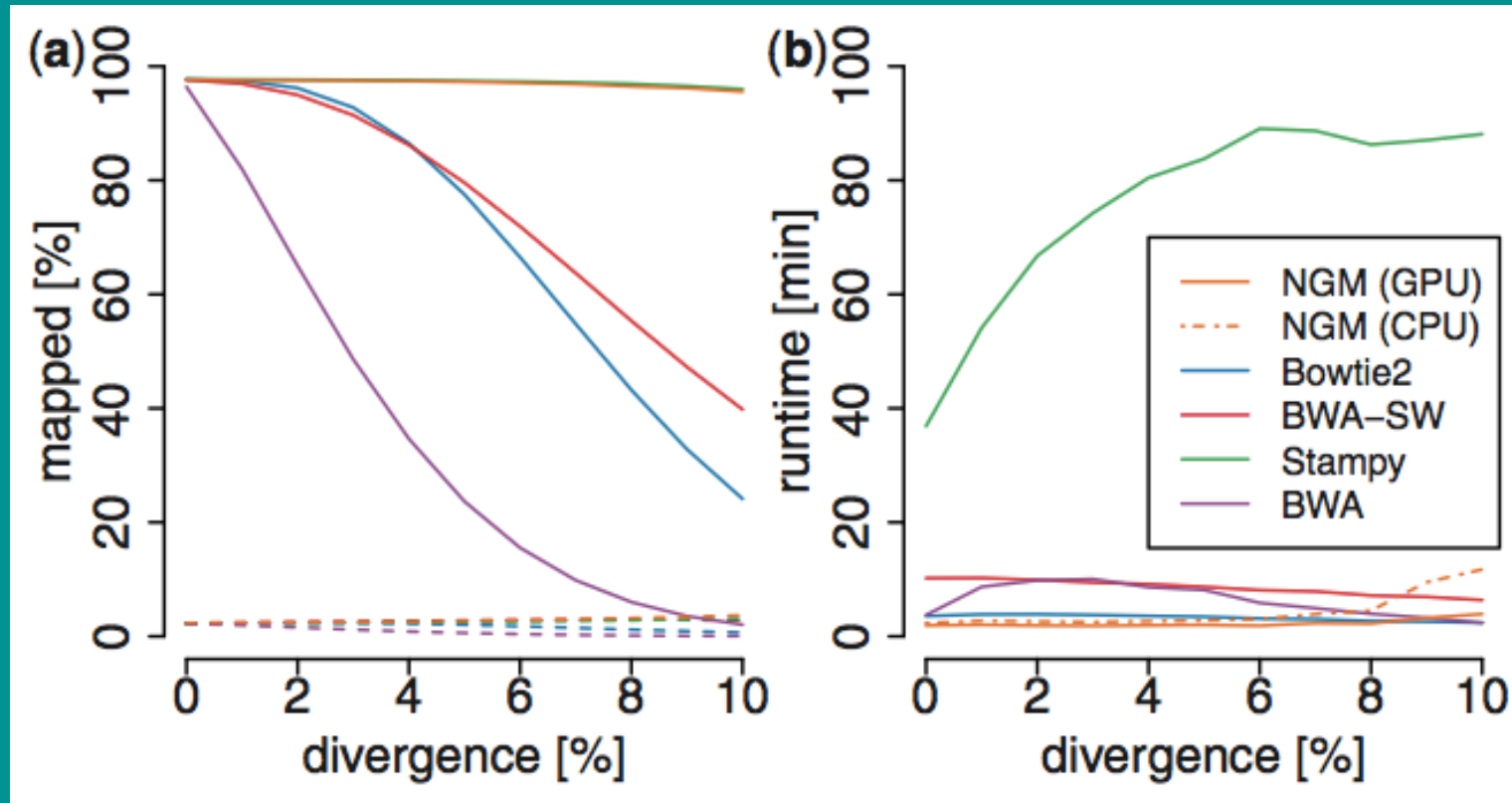
- Genes of interest are likely to have SNPs densities  $\gg$  genomic average
- These are not likely to get mapped
- Leads to biased expression values



- Resequencing data will be allele biased
- But perhaps only in small fraction of genome

50-kb nonoverlapping sliding windows estimated from a sample of 23 haploid Peach lines

# Allelic bias in read mapping



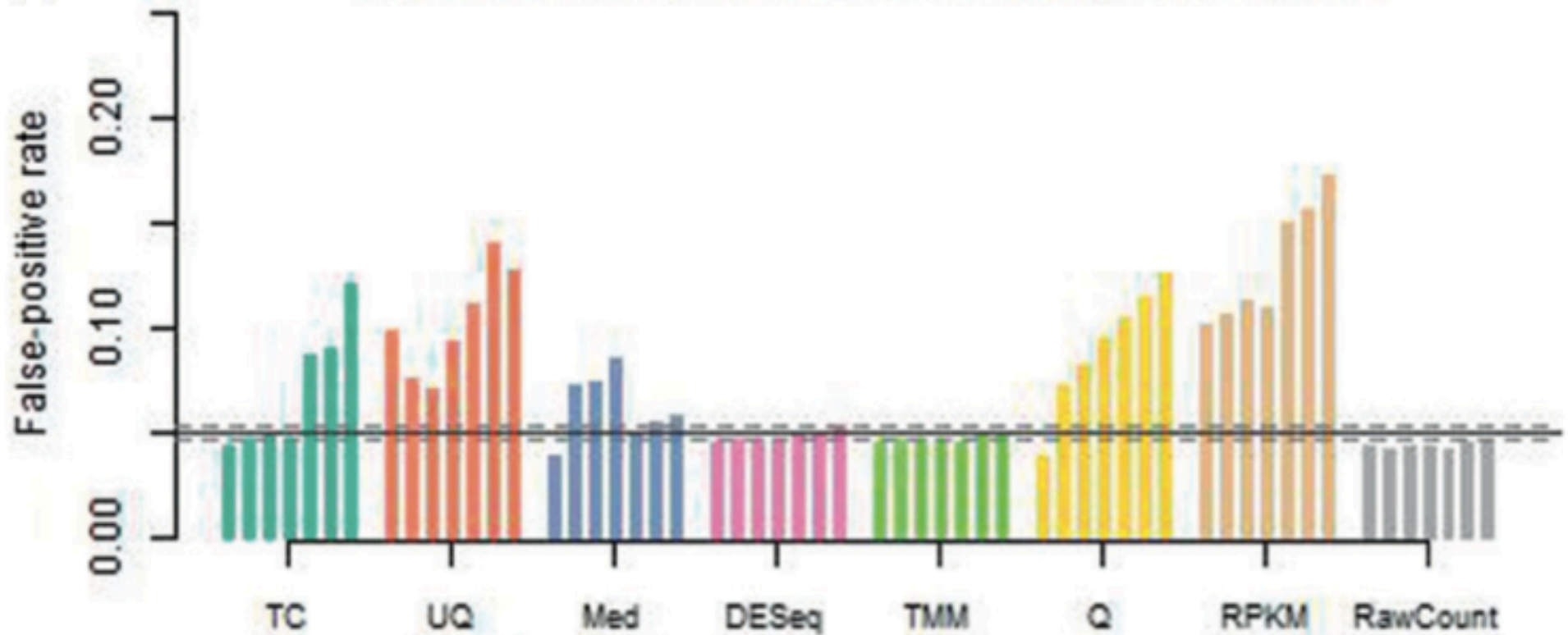
- Essentially identical to allele specific PCR bias ... but on a scale you can't detect unless you care to look
- Do your genes of interest have more than 3 SNPs / 100 bp?

# WILD WILD WEST

# Normalization

(a)

Equivalent library sizes / Presence of high count genes

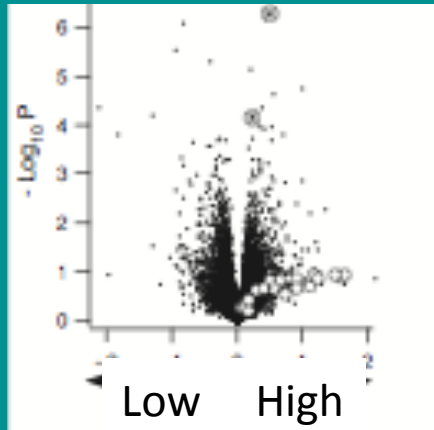


# Can genetic variation affect dispersal?

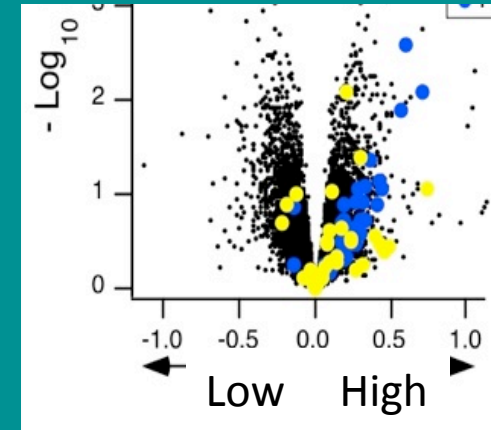
- Identifying such variation could help
  - Ecological & evolutionary study (theoretical models)
  - Conservation biology (captive breeding, predictions)



## Abdomen

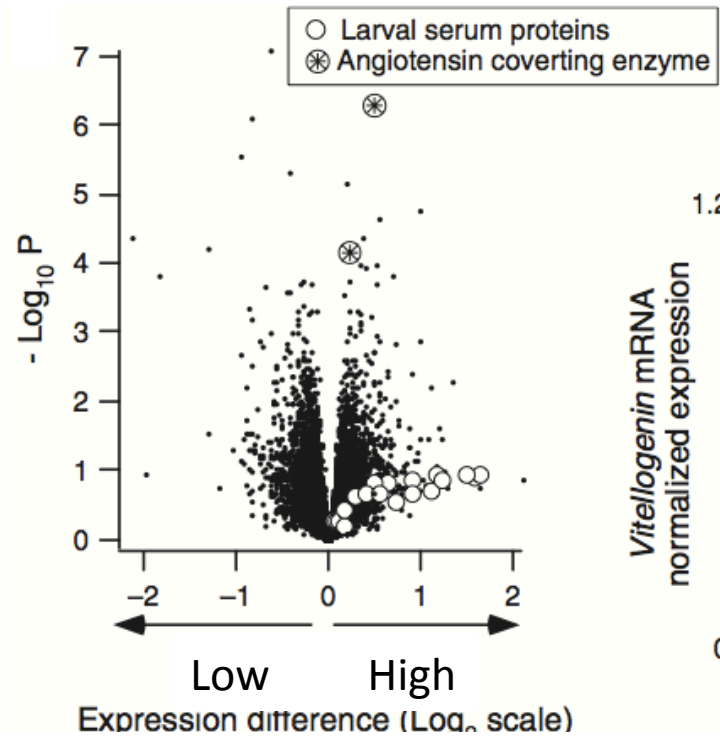


## Thorax



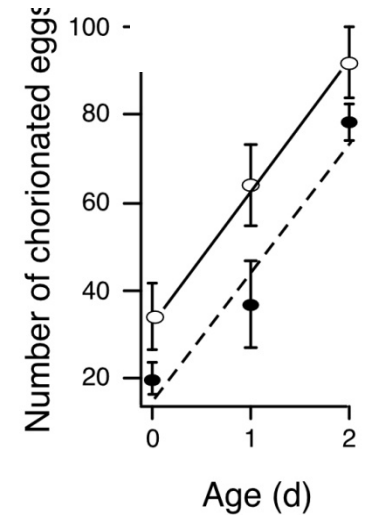
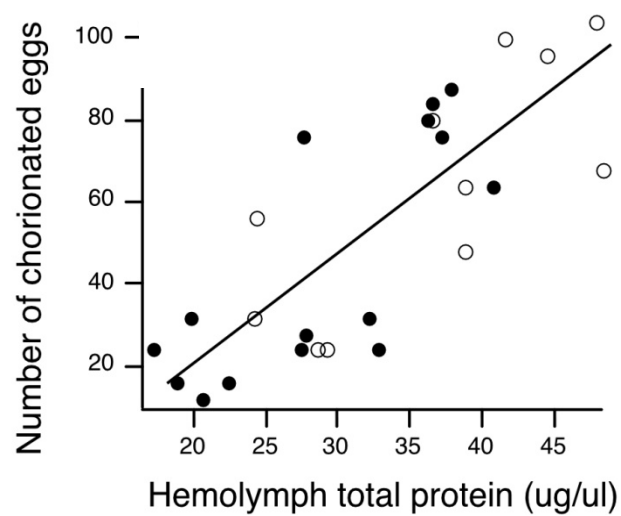
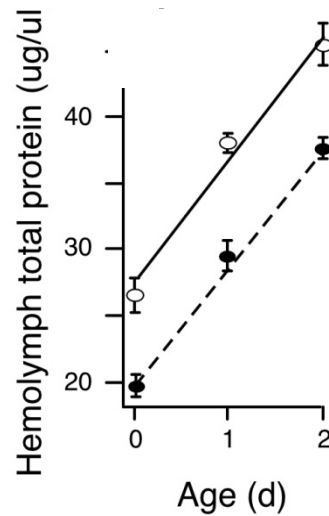
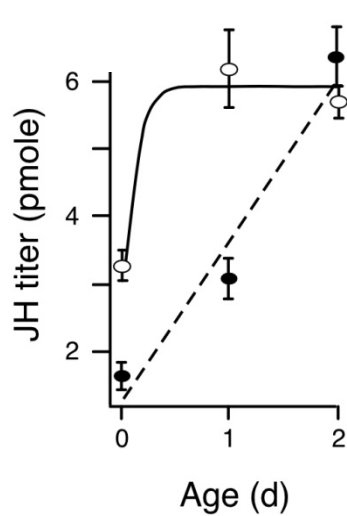
- Gene expression changes are primarily involved in protein allocation
  - Thorax - flight muscle performance
  - Abdomen - reproductive physiology
- A single gene could cause all these expression differences, but which gene?

# Butterfly dispersal genetics



● Low dispersal

○ High dispersal



# Most studies are annotation limited

- What is the biological meaning of the top Pvalue genes?
- Low Pvalue or expression genes are certainly important
- Gene set enrichments are key to insights
  - need network and regulatory insights relevant to the questions

Description	Uniprot	-log10P
Oxidoreductase.	Q9VMH9	7.087008
Hypothetical protein.		6.993626
SD27140p.		6.315473
	Q8SXX2	6.300667
SD01790p.	Q95TI3	5.316371
Electron-transfer-flavoprotein l	Q0KHZ6	5.1425
Pseudouridylate synthase.	Q9W282	4.784378
Hypothetical protein.	Q9VGX0	4.750469
CG14686-PA (RE68889p).	Q9VGX0	4.650051
Chromosome 11 SCAF14979, w	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





*Melitaea cinxia*  
454 sequence database

100 My



*Bombyx mori*  
Whole genome sequence,  
predicted gene set

320 My



*Drosophila melanogaster*  
Extensive genomic &  
functional resources

Assembly 2.0  
Contig\_57178  
Contig\_6821  
Contig\_1004  
Contig\_20226  
Contig\_27720  
Contig\_5260  
Contig\_27110  
Contig\_27390  
Contig\_26901  
Contig\_4713  
Contig\_20081  
Contig\_9982  
Contig\_15387  
Contig\_25362  
Contig\_36071



Bmori06 PepEd90  
BGIBMGA002704  
BGIBMGA003247  
BGIBMGA003248  
BGIBMGA003248  
BGIBMGA003248  
BGIBMGA003249  
BGIBMGA004806  
BGIBMGA004806  
BGIBMGA004865  
BGIBMGA004866  
BGIBMGA005329  
BGIBMGA006733  
BGIBMGA008859  
BGIBMGA008859  
BGIBMGA008859



Flybase gene ID  
CG33126  
CG6519  
CG6519  
CG6519  
CG6519  
CG6519  
CG33126  
CG33126  
CG33126  
CG33126  
CG3149  
CG6783  
CG4178  
CG4178  
CG4178

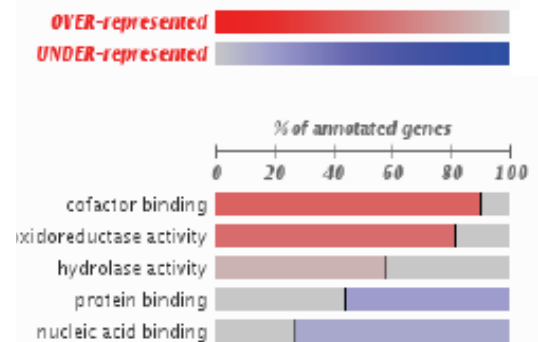
*D. melanogaster*  
lacks an orthologous  
reproductive  
physiology



### Gene Set Enrichment analysis using Gene Ontology database



#### Fatscan Analysis



# Life after your RNA-Seq experiment

- What are you likely to learn?
  - By measuring other aspects of the phenotype, we could at least validate and solidify our 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
  - Need data relevant to your phenotype and organism
    - Don't hesitate to make your own enrichment set

# 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

To address these we need to develop functional genomic insights in species with well understood ecologies that can be manipulated in the lab and in the field

## Widespread Cannibalism May Have Caused Prehistoric Prion Disease Epidemics, Science Study Suggests

Apr. 11, 2003 — Human flesh may have been a fairly regular menu item for our prehistoric ancestors, according to researchers. They say it's the most likely explanation for their discovery that genes protecting against prion diseases -- which can be spread by eating contaminated flesh -- have long been widespread throughout the world



Opinion

TRENDS in Genetics Vol.20 No.7 July 2004

# Balancing claims for balancing selection

Martin Kreitman<sup>1</sup> and Anna Di Rienzo<sup>2</sup>

Letter

Assessing the signatures of selection in *PRNP* from polymorphism data: results support Kreitman and Di Rienzo's opinion

Marta Soldevila<sup>1</sup>, Francesc Calafell<sup>1</sup>, Agnar Helgason<sup>2</sup>, Kári Stefánsson<sup>2</sup> and Jaume Bertranpetit<sup>1</sup>

Story time in genomics land .....

# Molecular spandrels:

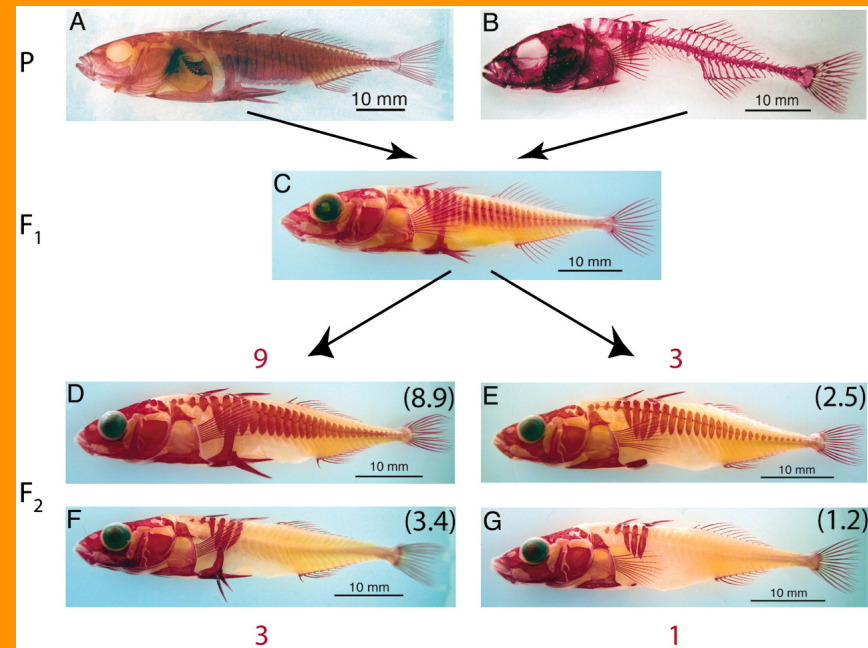
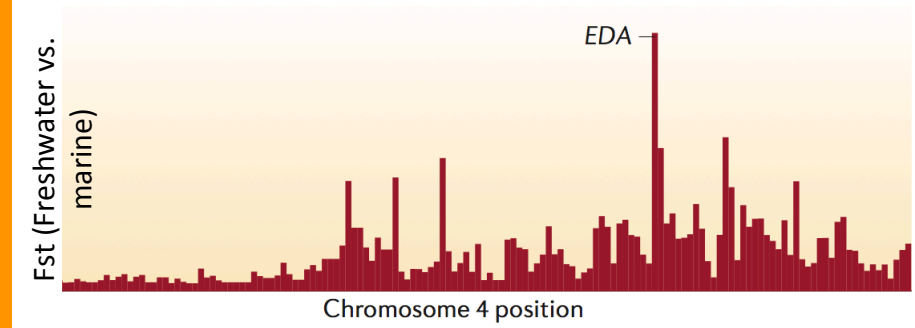
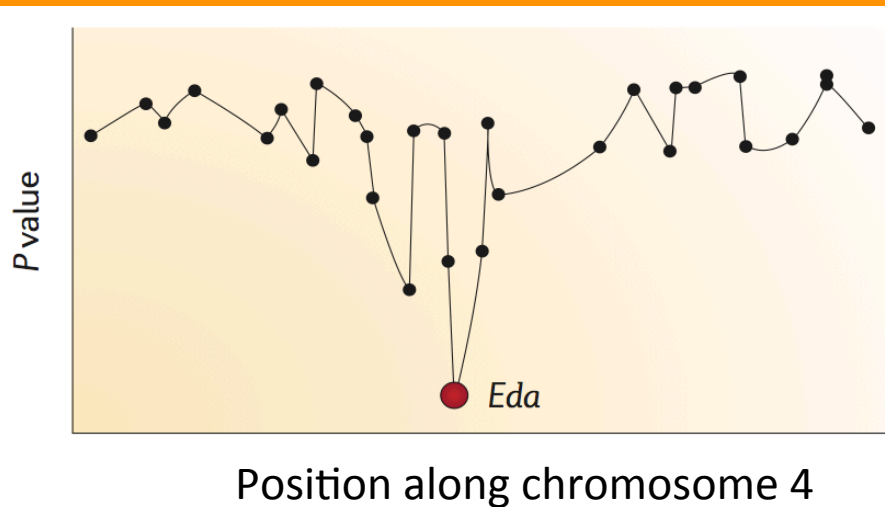
Story telling  
vs.  
causal understanding

Genomics of full of adaptive stories

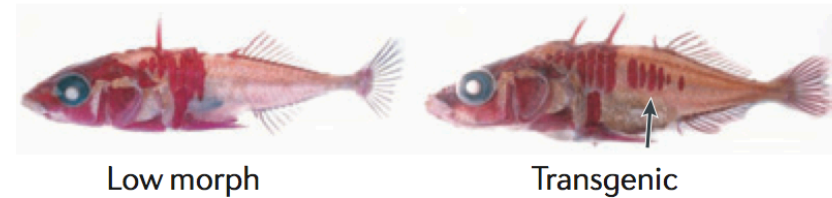
Functional and field validation of  
SNPs effects are needed to discern  
facts from fictions

# Model adaptation: the *Eda* gene

- Causes loss in body armor
  - Field association
  - QTL mapping



## Ac Gain-of-function

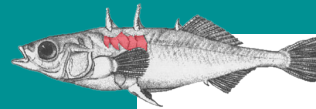


# Back to nature: do we know what we think we know?

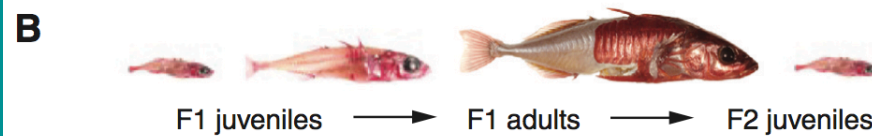
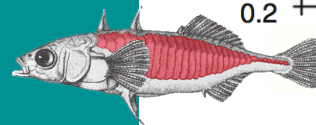
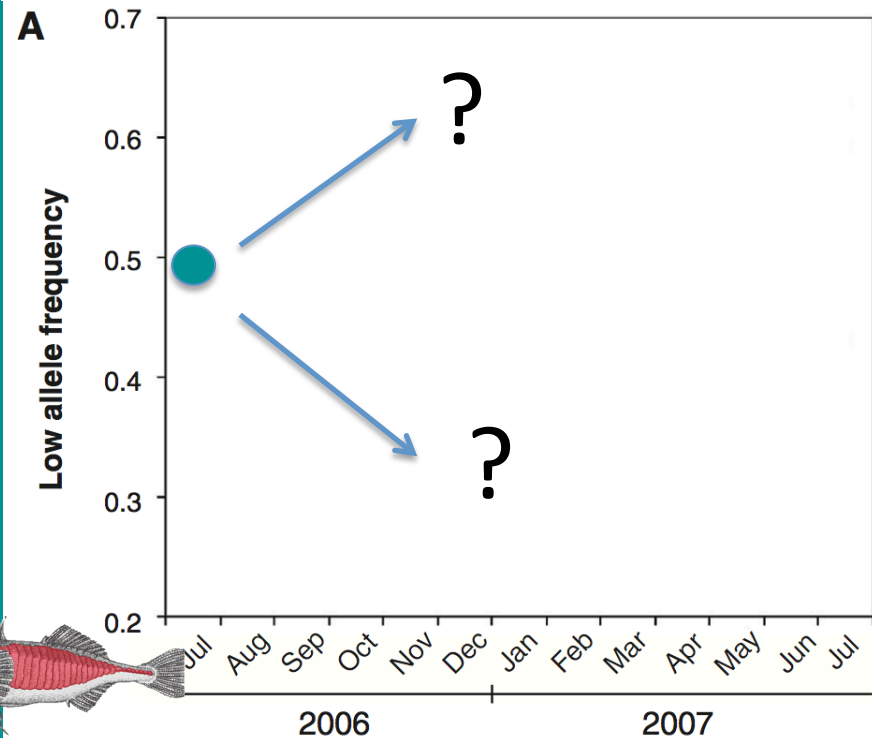
- Is low armor really adaptive in fresh water?
- Lets replay the selection event
  - Equal frequency *Eda* alleles in fresh water ponds

Studies in the field can uncover unexpected and complex selection dynamics

- Linked effect of other genes in the inversion on LG4?
- Is *Eda* even a target of selection?

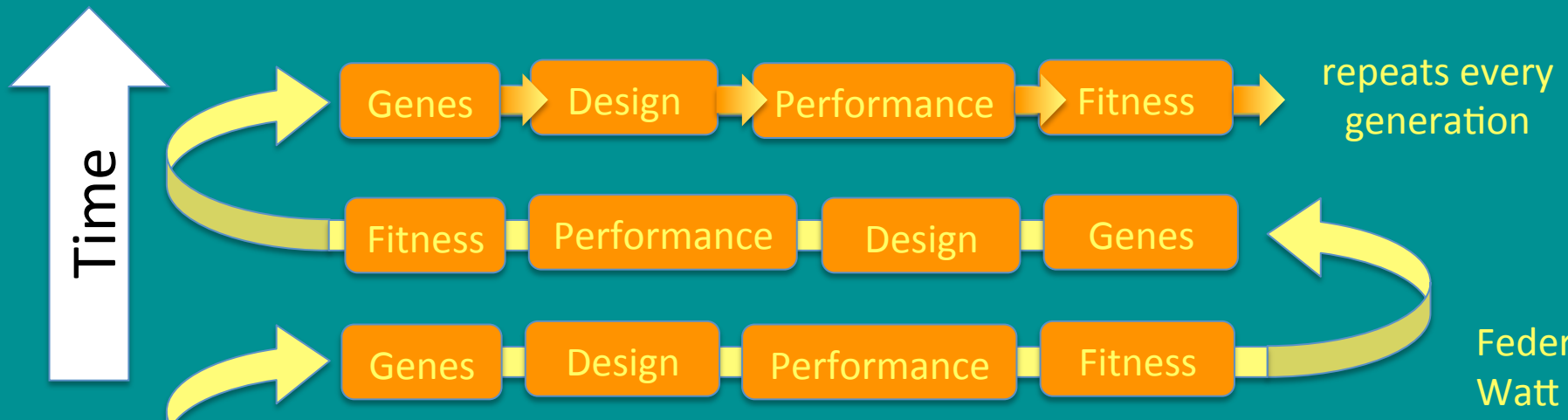


4 replicate freshwater ponds



Barrett et al. 2008 Science

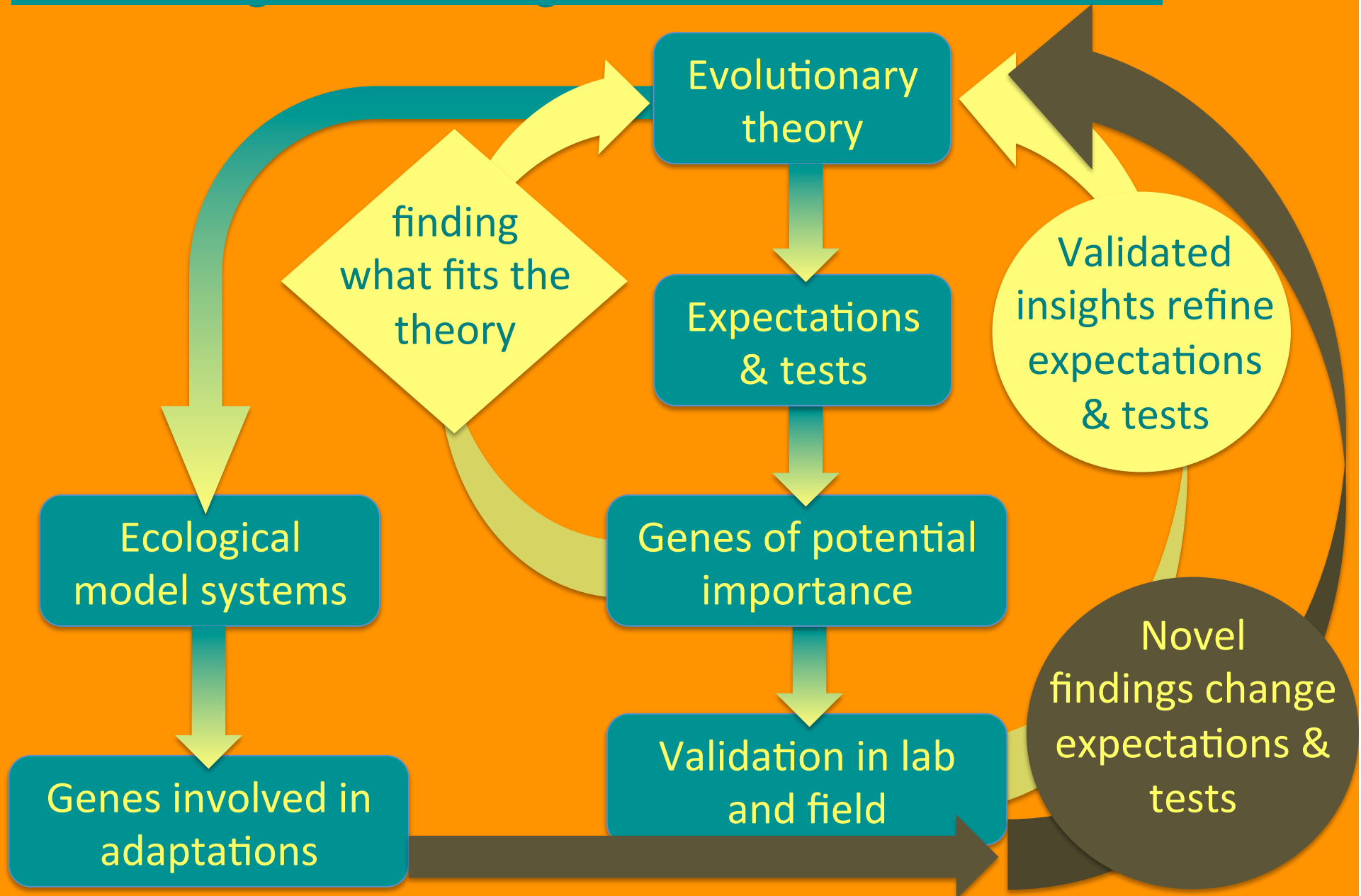
# Adaptation by natural selection



Feder & Watt 1992



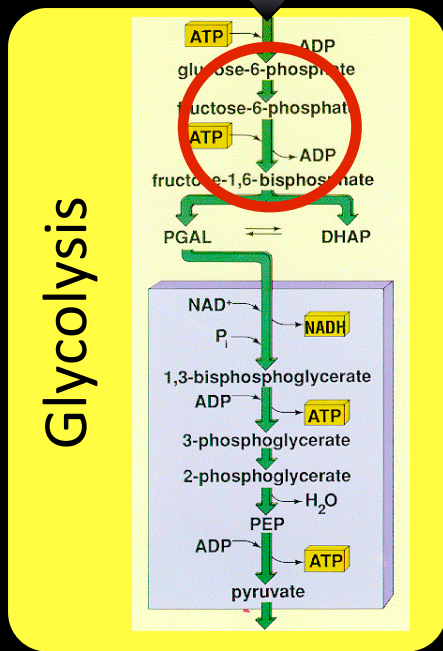
# Validating candidate genes moves us forward:



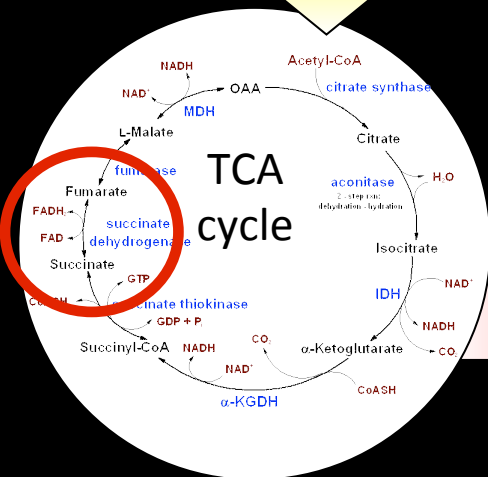
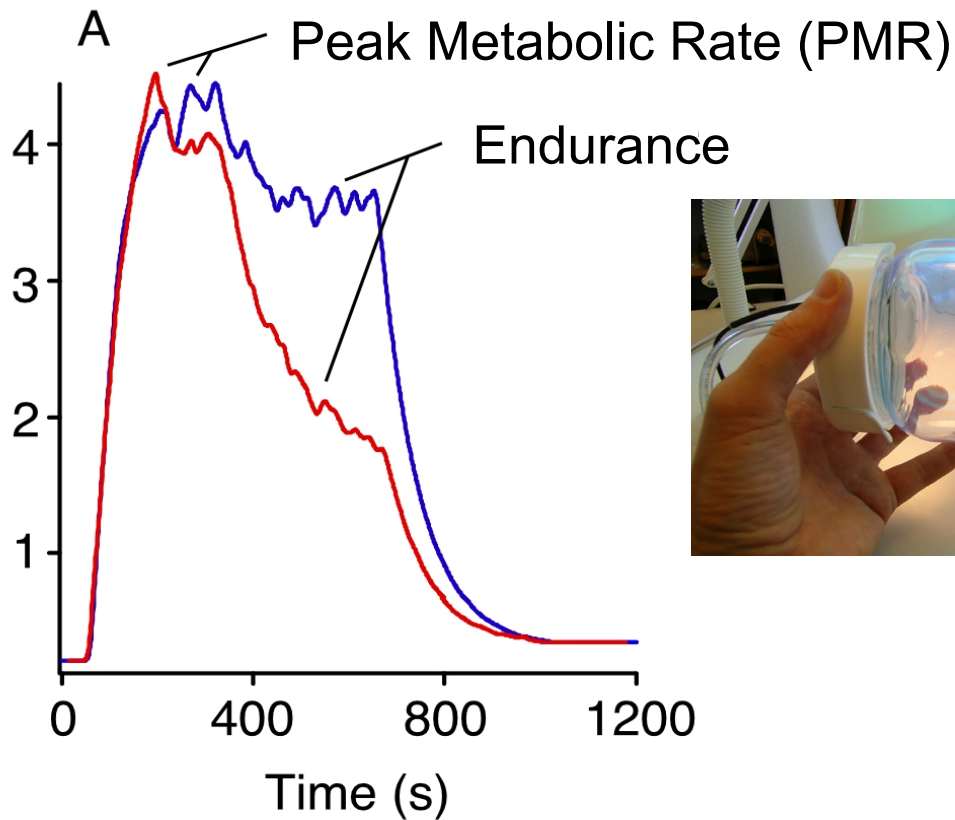
# Can you get there from here?

- Candidate genes associated with large effects on Darwinian fitness
  - Classic study systems in the wild
  - Validation process is still ongoing
- Can 'Second Generation' approaches find these same genes?
  - Sometimes not, or at least not easily
- Why?
  - Cause the modern tools aren't designed with such architectures in mind

Glucose

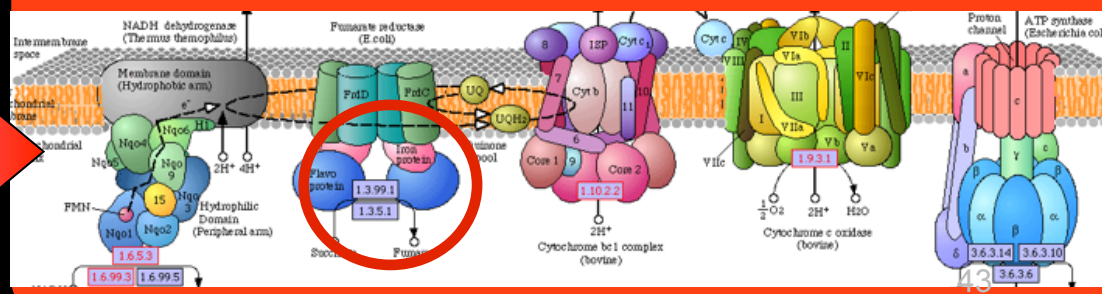


Flight metabolic rate (ml CO<sub>2</sub> hr<sup>-1</sup>)

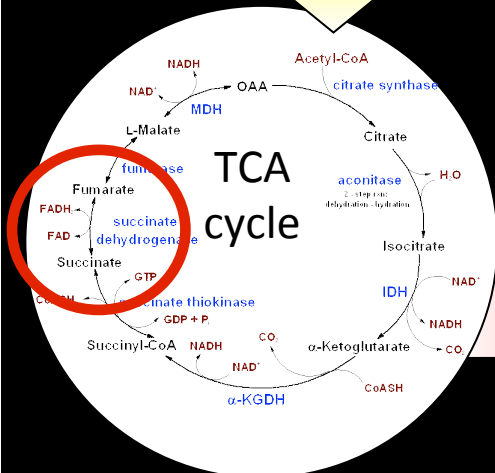
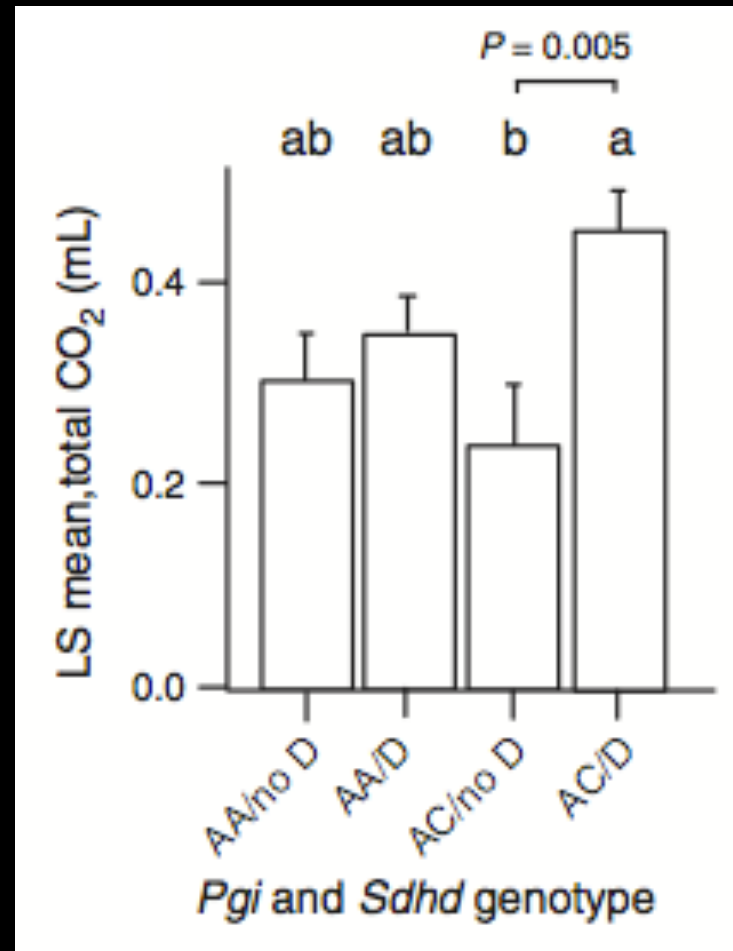
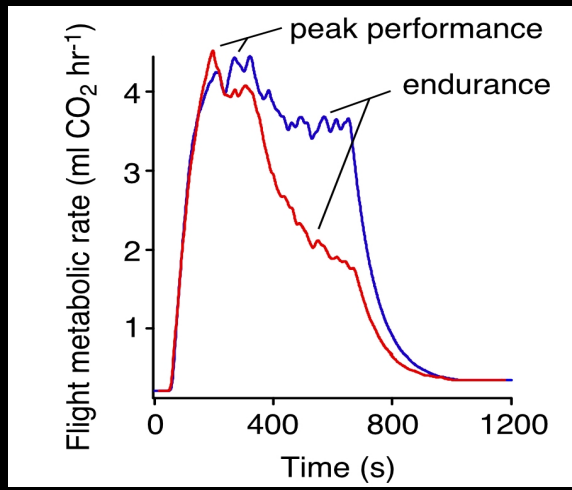
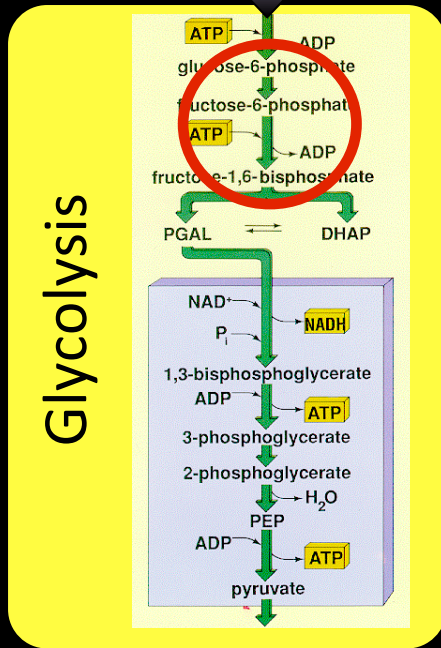


NADH

**Electron Transport**

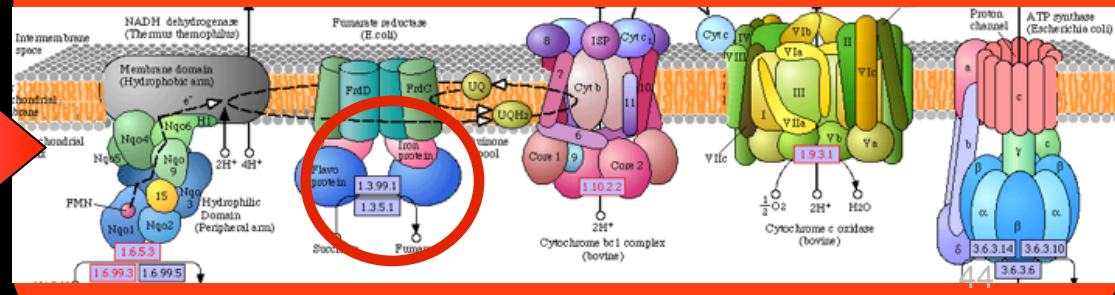


Glucose



NADH

Electron Transport



## Fitness effects in the wild

Year to year change in number of families living in 43 demes

*Pgi* & *Sdhd* SNPs are in linkage equilibrium



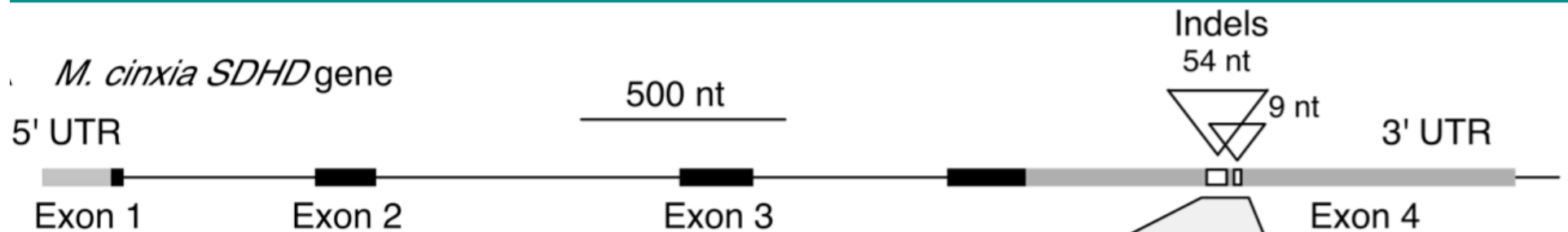
Source	(Full model $R^2 = 0.64$ )	d.f.	F ratio	P
Patch area		1	0.0001	0.99
Frequency <i>Pgi</i> F		1	10.9	0.002
Frequency <i>Pgi</i> F × Patch area		1	22.5	<0.0001
Frequency <i>Sdhd</i> M allele		1	19.1	0.0001
Frequency <i>Sdhd</i> M allele × Patch area		1	21.1	<0.0001

# Succinate dehydrogenase d (Sdhd)

3' UTR indel associated with performance and fitness in 5 studies across 3 populations



James H. Marden  
Penn State Univ.



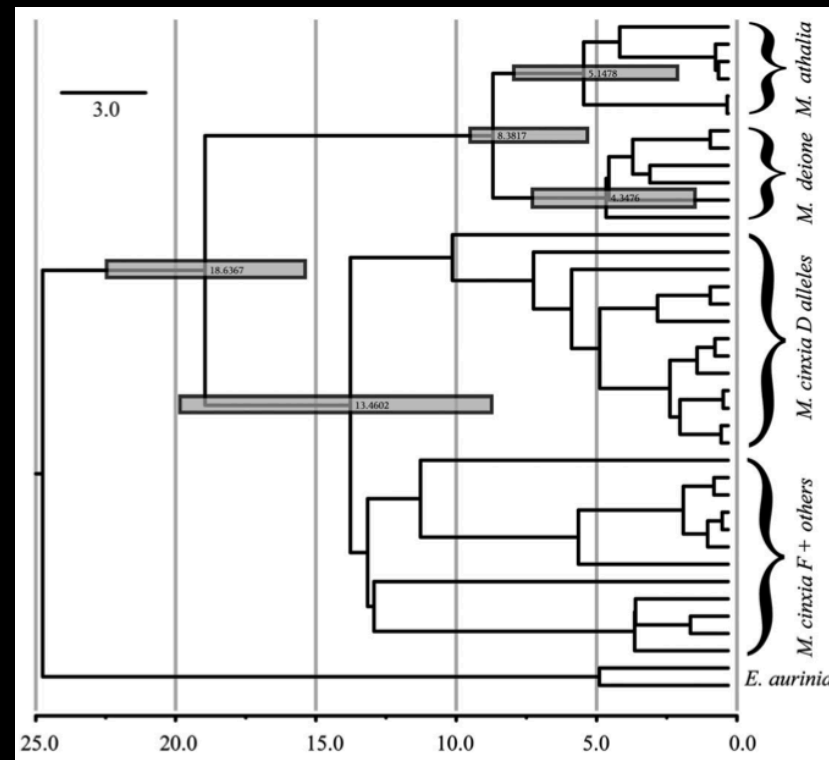
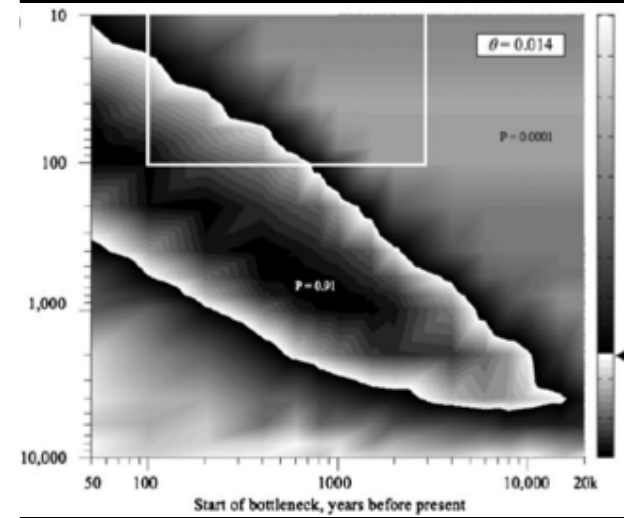
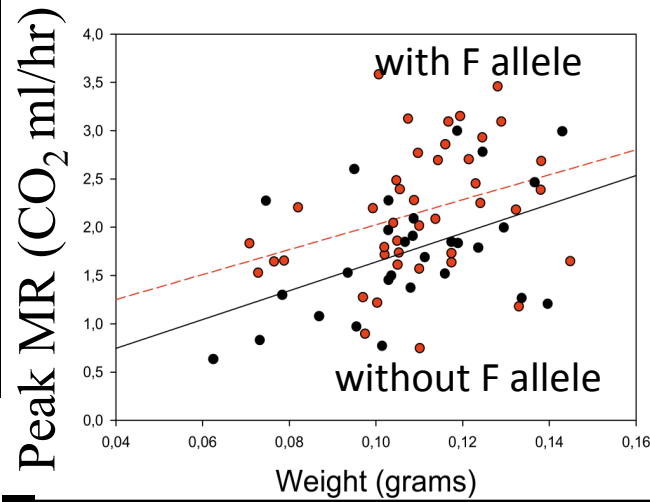
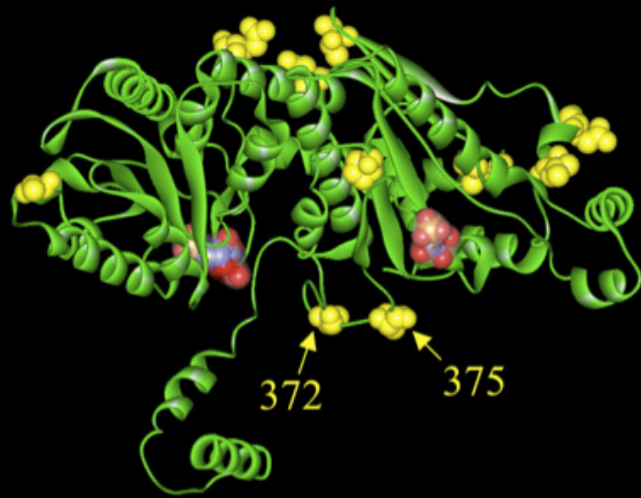
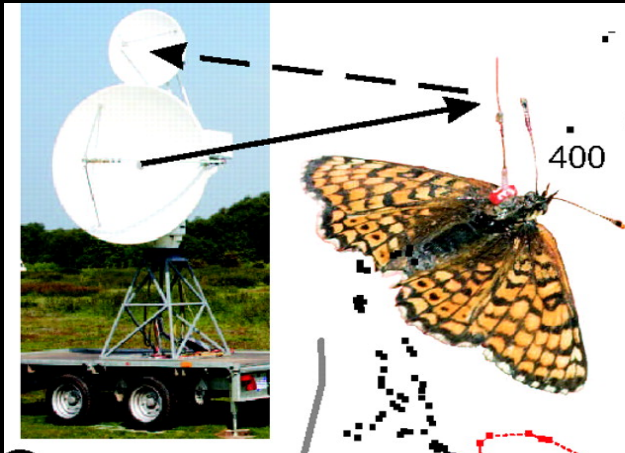
3' \_TAGAGTGATGGGTACAGAAAGT\_5' *M. cinxia* miR-71

I (Insertion)	TAAATGAGTTACCATATACTAAAATTTGTAAAAAAA-CTCATGTCTTTCTCTTAAATGTTAAAAACCTTTTATAATATATAAATAA
M (Mini Deletion)	TACTAC----ACCATATACTAAAATTTGTAAAAAAAACCTAAATGTCTTTCTCTTAAATGTTAAAAACCTTTTATAATATATAAATAA
D (Deletion)	TAAATG-----AAAAAAAAA-----AAAAAAAAACCTTTTATAATGTATAAATAA
E (Extra Deletion)	TAAATG-----AAAAAAAAA-----AAAAAAAAACCTTTT-----AAATAA

Assembly of this region would be biased to most common allele.  
Mapping would be allelic biased

Wheat et al. 2011; Marden et al. 2013

# Phosphoglucose isomerase (*Pgi*)



Niitepõld 2008; Wheat et al. 2010; Haag et al. 2005





# Conclusions:

What do we really know about .....

- molecular evolutionary dynamics?
- targets of selection in the wild?
- the age of species?
- transcriptomes?
- the performance of 2<sup>nd</sup> gen methods?

# Funding sources

Finnish Academy of Sciences (Finland)

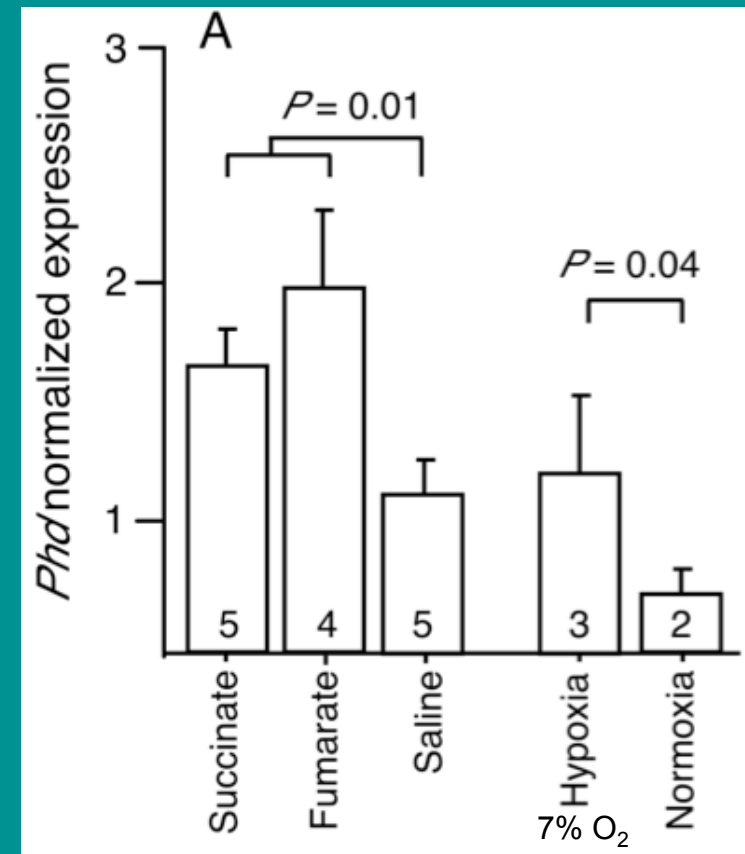
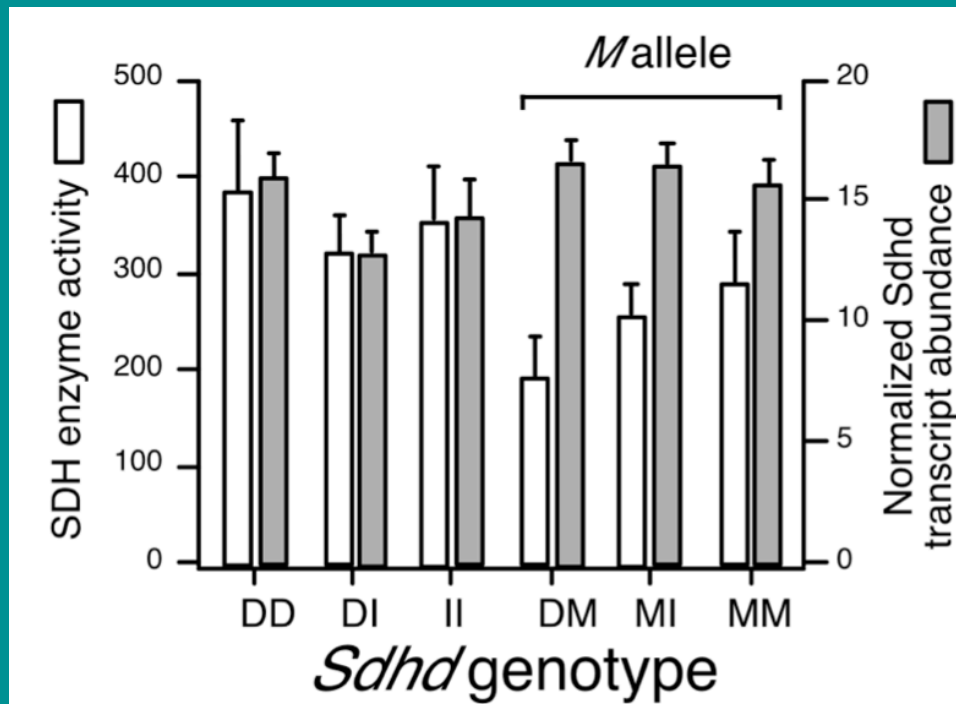
Vetenskapsrådet (Sweden)

Wallenberg Foundation (Sweden)



# How could *Sdhd* affect flight?

- Loss of function studies in humans result in constitutive activation of HIF pathway
- Increased flanking metabolites result in hypoxia signaling

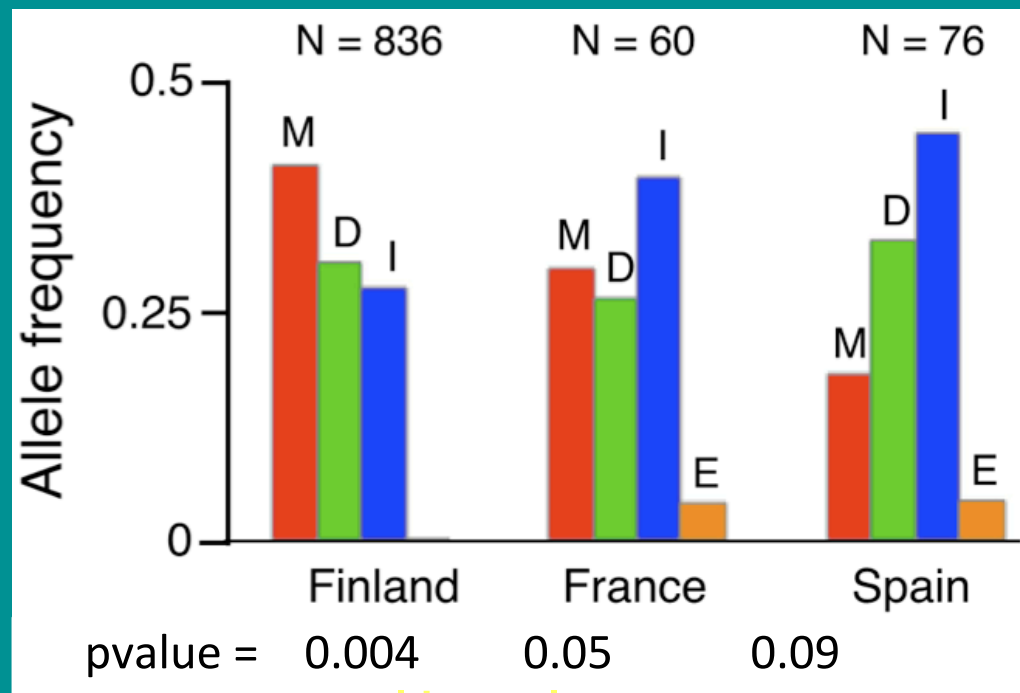


51 Suggests micro RNA down-regulation of SDH enzyme

Marden et al., accepted, Evolution

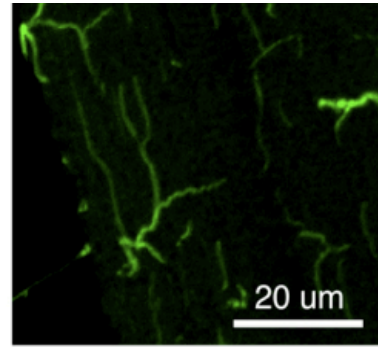
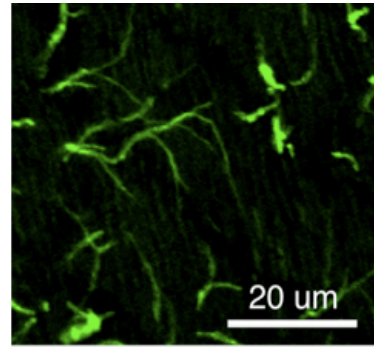
# Sdh<sup>d</sup> indel alleles

- Excess homozygosity / intermediate frequency within populations among alleles
  - Ewens Watterson tests suggestive of balancing selection

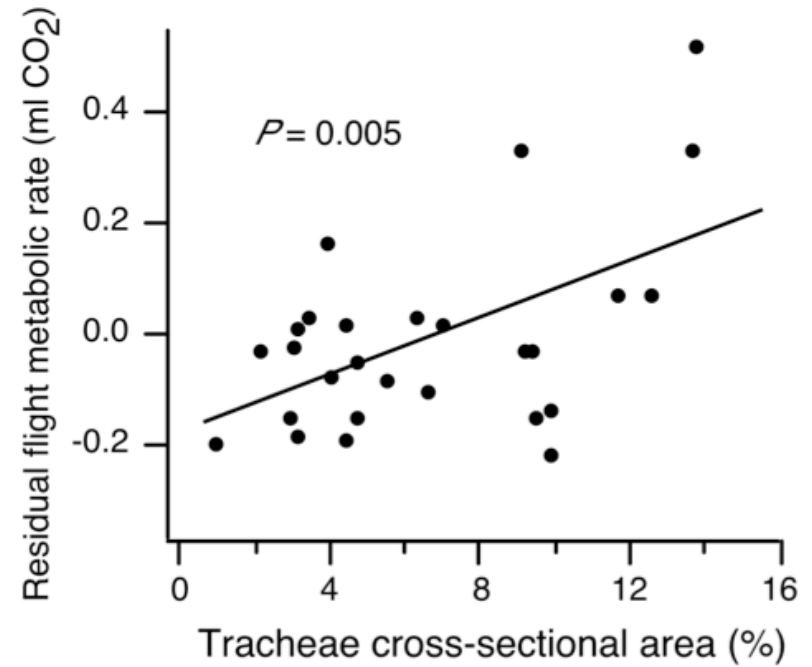
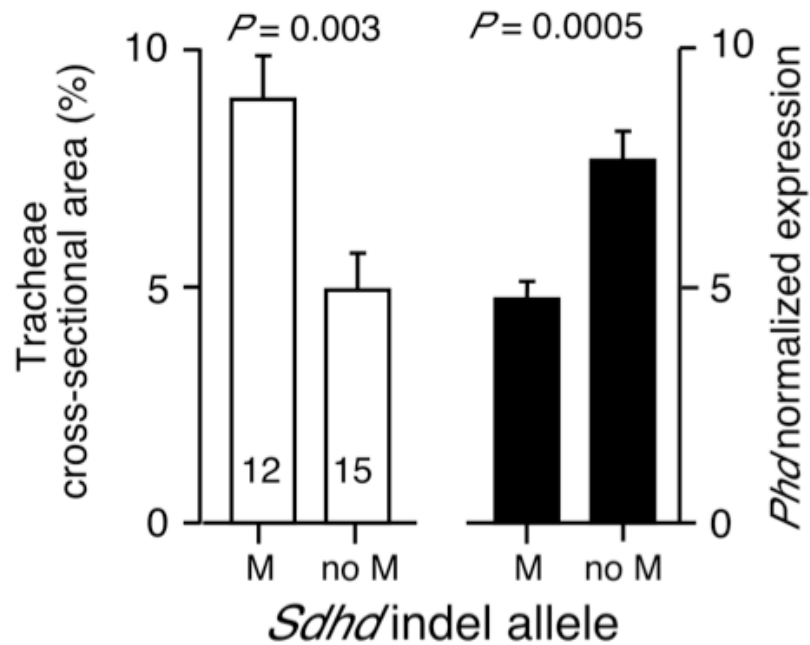


- M allele higher frequency in new vs. old populations
  - P = 0.04; N = 94 butterflies, 33 populations)

# Tracheael elaboration

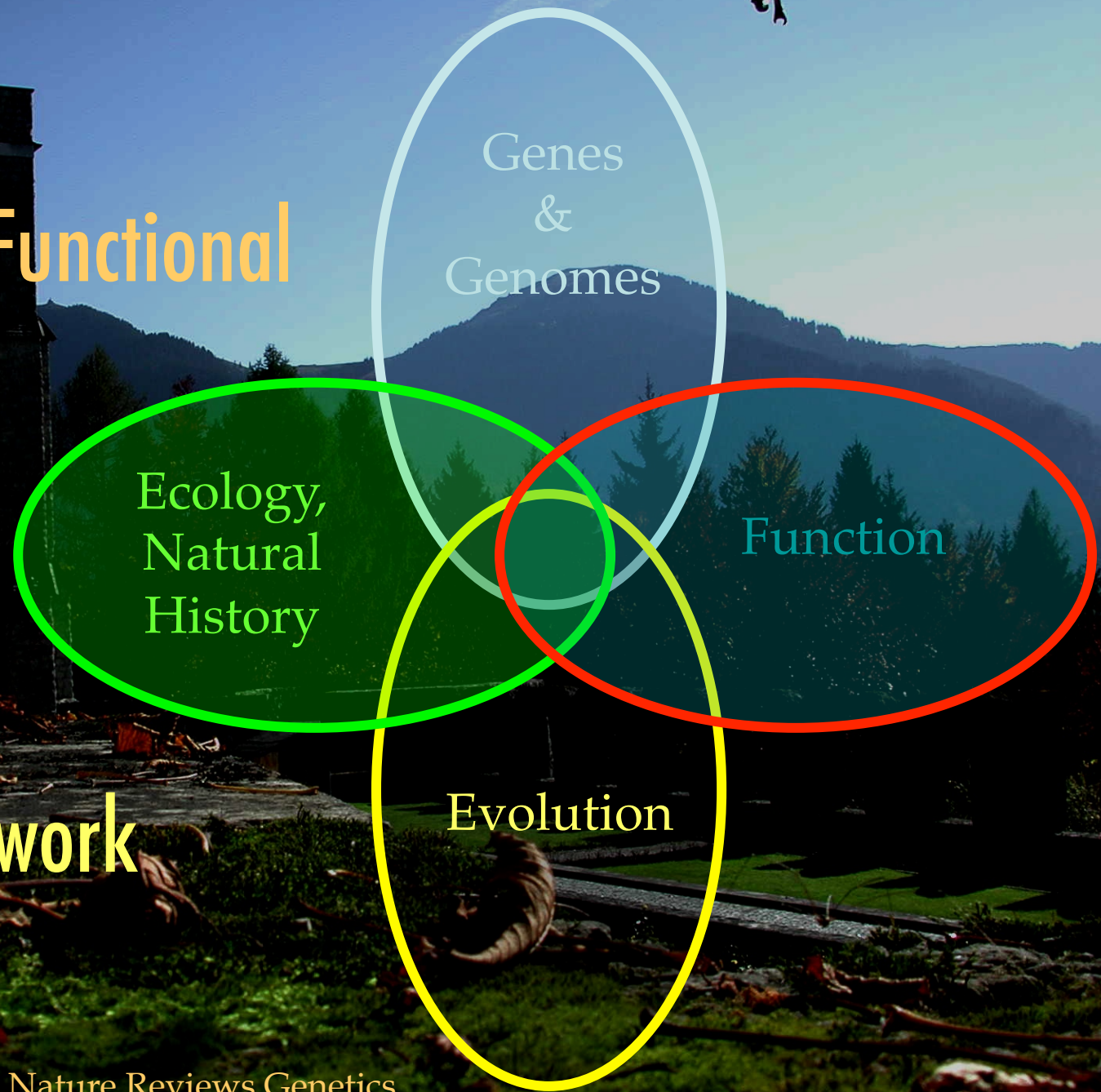


Two butterflies differing in tracheal elaboration



# Ecological & Evolutionary Functional Genomics

Integrative,  
collaborative work



Feder & Mitchell-Olds (2003) Nature Reviews Genetics  
Vasemägi & Primmer (2005) Molecular Ecology