Ecological Genomics, pt. 1

you, your data, your perception and the hard realities

Christopher West Wheat



Goal of this lecture

• Present a non-typical view of ecological genomics

• Make you uncomfortable by sharing my nightmares

• Encourage you to critically assess your results in light of publication biases

Disclaimer

I'm a positive person

I like my job and the work we all do

I'm just sharing food for thought

What if

50% of your favorite studies were just wrong?

How would that affect your expectations?

Publication replication failures

- Biomedical studies
 - Of 49 most cited clincal studies, 45 showed intervention was effective
 - Most were randomized control studies
 - Of the 34 that were later replicated, 41% were directly contradicted or had much lower effect sizes.
- Mouse cocaine effect replicates in three cities

 Highly standardized study
 Average movement was 600 cm, 701 cm, and > 5000 cm in the cities

Lehrer 2010 Ioannidis 2005 JAMA

Can publication bias increase effect size?



Decreasing effect size with increasing sample size

correlations between fluctuating asymmetry and individual attractiveness in various studies



Palmer 2000 Ann. Rev. Eco. Sys.

Decreasing effect size with increasing replication means what?



Palmer 2000 Ann. Rev. Eco. Sys.

Why Most Published Research Findings Are False

A research finding is less likely to be true when:

ne studies conducted in a field have a small sample size when effect sizes are small when there is a greater number and lesser pre-selection of tested relationships where there is greater flexibility in designs, definitions, outcomes, and analytical modes when there is greater financial and other interest and prejudice when more teams are involved in a scientific field, all chasing after statistical significance by using different tests

Ioannidis 2005 Plos Med.

There are lies, damn lies, and

Are datasets too big to fail?

What do follow-up studies reveal?

How can we gain confidence in our work?

Outline

- What is the genomic architecture of phenotypes?
- What is the power of molecular tests of selection?
- What does dissection of a classic comparative genomics study reveal?

Non – adaptive



disease, aging, height, etc.

Adaptive



salinity, color, resistance, etc.

generally ...



1000's of loci, each of small effect size

One or several loci of large effect

Is this a publication bias?

Will your trait have 1000's of small effect genes, or a few genes of large effect?

Sear (2010) ... Is bigger always better?

Rockman (2011) ... All that's gold does not glitter



What is the architecture of a causal variant?



How	predictable are	
C	daptations?	

	Plants	Animals
Coding ¹	71	163
Cis-regulatory	26	48
Other ²	16	7
Total	113	218
Null ³	67	32



	Morphology	Physiology	Behavior
Coding ³	62	170	2
Cis-regulatory	43	29	2
Other ⁴	3	20	0
Total	108	219	4
Null ⁵	41	58	0

Stern & Orgogozo 2008 Evolution

How do we find the genes that matter?

• Molecular tests of selection are popular, but ... —What are their assumptions and power?

- What are these tests detecting? —What is a footprint of selection?
 - How are they formed?
 - How large are they?
 - How long do the last?



What power do we have to detect balancing selection?

What is statistical power?

Power is the probability that the test will reject the null hypothesis when the alternative hypothesis is TRUE

Using ANOVA, you want power > 90% at reasonable sample size, right?

What power do we have to detect balancing selection?

		Width of window (bp)							
ρ	25	50	100	200	1000				
1	85.6	90.2	92.8	93.5	83.8				
3	80.8	85.3	86.3	83.5	44.7				
10	69.0	69.9	64.5	51.0	4.1				
30	48.1	42.5	31.0	15.7	0.1				
100	20.5	15.6	8.9	2.4	0.0				
Tajima's D									
% finding selection of 5000 simulations									

 For *Drosophila melanogaster*, power = 50% with window size of 200 bp, using 24 diploid individuals.

• For species with larger population size, power likely lower

• Recombination and gene conversion destroy 'footprint' rather quickly

Nordborg and Innan 2003 Genetics

Directional selection: an example of the expectations of hard selection

- Population genomics has been dominated by developing methods to detect hard sweeps for past two decades
 - But a 'null model' has been elusive

ATGGTAGGTCATATTGATCAGGGTGAATGTGCTAGAACATA ATGCTAGATCAAAGTGATCATGGTGAATGTGCTAGAACATA ATGGTAGATCAAATTGATCATGGTGCATGTGCTAGATCATA ATGCTAGATCATATTGATGATGGTGAATGTGCTAGATCATA ATGCTAGATCATATTGATCATGGTGAATGTGCTTGAACATA ATGCTAGGTCATATTGATCATGCTGAAAGTGGTAGATCATA







Parallel adaptation in fresh water lakes via hard sweeps

Marine population



Predominant form

Individual genome sequencing: powerful insights



2-5 X per individual, sliding 2500 bp window, 500 bp step

Jones et al. 2012 Nature

What type genomic regions are selected upon?



How common are such hard selective sweeps?



 Does your favorite test for selection rely upon such events?
 MK-test needs repeated events
 Fst outlier, EHH, Tajima's D, etc.



Storz 2005 Molecular Ecology



Hard vs. soft or incomplete sweeps in populations



.

What do soft sweeps look like?



How common were hard sweeps in our history?



- "classic sweeps were not a dominant mode of human adaptation over the past 250,000 years"
- "much local adaptation has occurred by selection acting on existing variation rather than new mutation"

1000 Genomes PC 2010 Science Hernandez et al. 2011 Science

How common are soft sweeps in your species?

Thought experiment:

Do most species respond to selection in the lab? Yes Why? Because they have existing variation in population If populations have variation, is selection likely to act on it? Yes What does this tell us about frequency of soft selection in wild?



Garud, Messer, Buzbas and Petrov 2013 ArchivX

Age and type of selection matters

- Novel mutation, large mutation, hard sweep selected to fixation

 High probability of detection
- Old mutation, polygenetic, soft sweep of incomplete fixation
 - Low probability of detection
- Finding the causal mechanism
 - Coding > expression
 - SNPs > more complex mutations (indel, TE, CNV)
 - Ongoing gene flow, grouping by phenotype across replicate populations helps a lot
- What is the relative frequency of these?
 - What will be the architecture of your phenotype?
 - What does your method have the highest power to detect?



Get ready, here come the 1000ⁿ genomes

- Roughly 20 arthropods sequenced to date — plans to sequence 5,000 more
- Many other large scale projects coming online







- Unprecedented data for studying:
 - Phylogenetic relationships
 - Genome evolution
 - Functional insights into genes and genomic features (e.g. regulation and inheritance)

Classic study: Evolution of genes and genomes on the *Drosophila* phylogeny



Drosophila 12 Genomes Consortium 2007 Nature

Tempo and mode of chromosome evolution



 > 20 My, chromosomal order completely reshuffled in Diptera Drosophila 12 Genomes Consortium 2007 Nature



Single-copy orthologues Conserved homologues Patchy homologues (with mel.) Patchy homologues (no mel.) Lineage specific

Selection dynamics across functional categories



• 33.1% of single-copy orthologues have experienced positive selection on at least a subset of codons.

Drosophila 12 Genomes Consortium 2007 Nature

Gene Family Evolution across 12 Drosophila Genomes

- One fixed gene gain/ loss across the genome every 60,000 yr
- 17 genes are estimated to be duplicated and fixed in a genome every million years



Drosophila 12 Genomes Consortium 2007 Nature Hahn et al. 2007 Plos Genetics

Comparative Genomics : a house of cards?

- Data scale is too large to thoroughly assess errors ...
 Its likely 50% of what you think you know is wrong (it's true for me)
 What is reality?
- All conclusions, at some stage, rest upon
 - Simple bioinformatics



- Exploring two pillars of this paper, their error and repercussions
 - Gene alignments in detecting positive selection
 - Calibrations in temporal analysis

Established studies allow ...

Follow up studies to reveal limitations Robust findings to emerge with age

Inferring selection dynamics:



33.1% of single-copy orthologues have experienced positive selection on at least a subset of codons.

How robust are these conclusions?

Codon based tests of selection

Positive selection

f.ex. effector genes

 d_N

Neutral evolution f.ex. pseudogenes

d

Purifying selection f.ex. housekeeping genes



> 1 positive sel.
= 1 neutral
< 1 purifying sel.</pre>

IMPRS workshop, Comparative Genomics

Classic study: Evolution of genes and genomes on the *Drosophila* phylogeny



Drosophila 12 Genomes Consortium 2007 Nature

dN/dS estimates by aligner 6690 orthologs

 5 alignment methods

 Little agreement of the different dN/dS estimates



Markova-Raina & Petrov 2011 Genome Biology

Comparing results across methods is responsible bioinformatics!!!!!

Since we can't look at our data, we need approaches that allow 1st principal assessments



Markova-Raina & Petrov 2011 Genome Biology

Alignment has larger effect than biology

 Number of significant genes in common across 1, 2, 3, 4, or all 5 of the alignment methods

		420)				43	2	[]	1		440
KDF	RN	DQD	DE	EE	DE	E.		AE	SS	EI	IED	DDG
KDF	RNI	DQD	DE	EE	DE	E.		AE	SS	EN	IED	DDG
KDF	RN	DQD	DE	EE	DE	E.		PE	SS	EI	IED	DDG
KDF	RN	DQD	AE	Ε·	DE	E.	·	AE	SS	EC	ED	DDG
KDF	RN	DQD	DE	Ε·	E	EN	E	SS	EN	EC	ED	DDG
QDF	RT	DQD			ED	E.		G S	SS	DI	ED	EEA
		420					43	0				440
KDF	RN	DQD	DE	EE	DE	EA	E	SS	Е·	1	ED	DDG
KDF	RNI	DQD	DE	ΕE	DE	EA	E	SS	E۰	- N	ED	DDG
KDF	RNI	DQD	DE	ΕE	DE	EP	E	SS	E۰	- N	ED	DDG
KDF	RNI	DQD	AE	۰E	DE	EA	E	SS	Ε·	- 0	ED	DDG
KDF	RNI	DQD	DE	۰E	I E	EV	E	SS	EN	EC	ED	DDG
QDF	RTI	DQD	ED			EG	S	SS	D -	- 0	ED	EEA



- Two alignment results
 - Top (Tcoffee) with 3 site sel. sites
 - Bottom (ProbCons) indicates
 region has a 60% alignment
 probability

Markova-Raina & Petrov 2011 Genome Biology

Temporal inference:

fact or fiction?



Timing of divergence



- Directly affects rate estimates
- Deriving unbiased dates from molecular data

 Large field of software development



- Bayesian methods, while potentially informative and unbiased
 - Can be easily, and are routinely, abused



Classic study: Evolution of genes and genomes on the *Drosophila* phylogeny



Drosophila 12 Genomes Consortium 2007 Nature



- Drosophila clade
 - Schizophora
 constrained to
 maximum of 70 Ma
 - Without constraint, goes to 115 Ma

What's reality?



Episodic radiations in the fly tree of life (Wiegmann et al. 2011 PNAS)

Determining objective priors is challenging





Obbard et al. 2012 Mol. Biol. Evol.

Priors directly influence posteriors



Prior distributions matter

- Integrative science is challenging
- Discuss or collaborate with experts to evaluate your approach.



Wheat and Wahlberg 2013 Trends Ecology & Evolution

How do we gain dating confidence when we are in the dark?

- Fossils and DNA are likely to rarely agree
- How can we assess the temporal signal in the DNA in a robust manner?
 - Reducing prior biases and using lots of DNA data, while modeling likely violations of analysis models



Wheat and Wahlberg 2013 Trends Ecology & Evolution



Microevolution effects

Previous examples were at deep evolutionary time scales

Surely such problems don't exist at the within genera level Right?

Recombination violates dN/dS tests

Codeml inferred selection:

False positives can increase to over 30%



- 13% of sites simulated at omega = 2.5
- Sample size = 30 sequences

Anisimova 2003 Genetics

Posterior distribution estimates of substitution rates from mitochondrial control region from Beringian bison





Ho et al. 2007 Systematic Biology

Time dependent rates of molecular evolution

Significant implications for phylogeographic studies that use fixed rates to assess demographic with environmental change



Post-genomics challenge

"What we can measure is by definition uninteresting and what we are interested in is by definition unmeasureable" - Lewontin 1974

"What we can assemble in the genome may, by definition, be uninteresting and what we are interested in is by definition very difficult to sequence and assemble and annotate and estimate"

- indels & inversions
- gene family dynamics
- demographic and selection dynamics
- temporal estimates



Goal of this lecture

- Present a non-typical view of ecological genomics
 - So you have a more complete view of the field
- Make you uncomfortable
 - Provide a context for understanding your results
- Encourage you to rethink the reality presented by publication biases
 - Overcoming this bias is a continual challenge

