

Christopher West Wheat





- 1995 2001 PhD California
- 2002 2005 Postdoc Germany
- 2005 2008 Postdoc Finland
- 2009 unemployed 4 month, spent all savings
 > 50 job applications, 1 grant application
- 2009 visiting scientist Germany
 - 1 job offer UK
 - 1 grant Finland
- 2012 started tenure track Sweden

What was important?

- Being able to move
- Chasing the money & skills
- Learning how to:
 - Write publications, grants
 - Believe in my ideas/skills

Needed to put science first, while having lots of fun along the way





What do you study?

Rough totals:

- Invertebrates: 7
- Fish: 8
- Mammals: 5
- Microbes: 16
- Plants: 4
- Humans: 4

What are your goals?

- Finding and study genomic regions that matter
- Investigating ecological
 processes
 - metagenomics
- Investigating physiology

 RNAseq

Goal of this lecture

- Present a critical view of things genomic
- Make you uncomfortable by sharing my nightmares
- Encourage you to critically assess findings and expectations in light of easy errors and publication biases

Disclaimer

I'm a positive person

I love my job and the work we all do

I'm just sharing scrumptious food for thought



TABLE 1 Variable	le nucleotides from the coding region of the Adh	locus in D. melanogaster, D. simulans and D. yakuba							
005 781 0 789 7 808 Å	D. metanggaster D. smulas a b o d + f g h 1 j k 1 a b o d + f T T T T T T T T T T T T T T T	0 yakuto a b c d + f g h i j k 1 0 C C C C C C C C C C C Syn. Fixed 0 C C C C C C C C C C Syn. Fixed							
816 G 834 T 859 C 807 C 870 C		Adaptive protein evolution at						ion at 1	the
950 G 974 G 963 T 2019 C	т.ттт Т.тттт	Adh loous in Droconhile							
1044 T 1043 C 1049 C	TT	Adn locus in <i>Drosophila</i>							
1901 G 1527 T 1531 C 1960 T		18888888888888888888888888888888888888							
1178 C 1584 C 1590 C 1196 C		G G G G G G G G G G G G G G G G G G G	John H. McDonald & Martin Kreitman						
1199 C 1202 T 1203 C 1229 T 1250 T		C C C C C C C C C Dyn. Poly. C C C C C C C C C C Dyn. Pised 	Department of Ecology and Evolutionary Biology Princeton University						
1235 C 1244 C 1265 C 1271 A		29n. Poly. 0 C C C C C C C C C C C C C C C C C C C	Princeton, New Jersey 08544, USA Nature 1991						
1277 T 1285 C 1296 C 1304 C	A A	C C C C C C C C C C C C Dyn. Pland 						Natar	0 1551
1425 C 1431 T 1443 C 1452 C	**************************************	CCCCCCCCCCCCC bys. Poly. Poly. bys. Poly. bys. Poly.							
1490 A 1504 C 1508 C 1504 T		Bys. Full 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0							
		ter	D. simulans D. yakuba						
	Con.	abcdefgh	i j k l	ab	cdef	abc	defghijk	1	
781	G	ттттттт	тттт					~ Repl.	Fixed
789	Т					CCC		C Syn.	Fixed
808	A		T	T T	 т т т т	GGG		G Repl.	Poly
834	т			C C	C			- Syn.	Poly.
859	C					GGG	GGGGGGGGG	G Repl.	Fixed
867	C					GGG	GGAGGGGG	G Syn.	2 Poly.
870	C	TTTTTTT	TTTT					- Syn.	Fixed
		UUU 🔪 _{Phe}	UCU		UAU 👔	Tur	UGU Cvs		
		UUC 🖌 ' '''ë	UCC (Ser	UAC 🕽	i yi			
		UUA	UCA (Ser	UAA 5	<u></u>	UGA Stop		
		UUG 👌 Leu	UCG		UAG }	Stop	UGG Trp		
		CUUN	CCUN		CAUN		CGU		
			CCC	-	CAC	His	CGC .		
		CUA Leu	CCA	Pro	CAA		CGA Arg		
		CUG			CAG	Gln	CGG		
		0007	0007		, ond j				







was used to test the null hypothesis, that the proportion of replacement substitutions is independent of whether the substitutions are fixed or polymorphic. G=7.43, P=0.006.



Adaptive protein evolution at the *Adh* locus in *Drosophila*

John H. McDonald & Martin Kreitman

Department of Ecology and Evolutionary Biology, Princeton University, Princeton, New Jersey 08544, USA

From DNA to Fitness Differences: Sequences and Structures of Adaptive Variants of *Colias* Phosphoglucose Isomerase (PGI)

Christopher W. Wheat,*†¹ *Ward B. Watt*,*† *David D. Pollock*,*†² *and Patricia M. Schulte**†³ *Department of Biological Sciences, Stanford University and †Rocky Mountain Biological Laboratory, Crested Butte, Colorado





If the biomedical science has the most money and oversight, then

Their findings should be robust:

- Repeatable effect sizes
- The same across different labs
- The same across years

Publication replication failures

• Biomedical studies

- Of 49 most cited clincal studies, 45 showed intervention was effective
- Most were randomized control studies (robust design)
- Mouse cocaine effect study, replicated in three cities — Highly standardized study

Ioannidis 2005 JAMA; Lehrer 2010











But surely, this doesn't apply to genomics

Or does it?

Outline

- Are these biases inherent in genomic studies?
- Why is this happening?
- How can we try and overcome these problems?



There are lies, damn lies, and

But wait, is that fair?

Are these really lies?

Where does this bias come from?

- Population heterogeneity
 - Space and time
- Publication culture
 - Large & significant effects publish fast and with high impact
 - Small & non-significant effects publish slow with low impact



Apophenia

The tendency to seek and see patterns in random information and view this as important





Story telling of Type 1 errors Celebration of the false positives

Genomics is too big to fail

- Making errors is extremely easy
- Results will very likely be significant, and sometimes dramatically so
- In non-model systems, rarely have replication studies
- You must always question your bioinformatics before falling in love with your results

When results are better than you could have dreamed,





Why? this was a technical artifact called a batch effect. confounded 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?								



FORENSIC BIOINFORMATICS AND REPRODUCIBLE **RESEARCH IN HIGH-THROUGHPUT BIOLOGY**

"Data processing, however, is often not described well enough to allow for exact reproduction of the results,

Thanks: Malachi Griffith

Baggerly and Coombes 2009

Cell Lines

















But there are lots of errors out there ...

In most instances, this is scientific progress ...

But, you must navigate these to calibrate your expections and approaches

Bioinformatics: get it right!

Can happen using the most basic tools / steps in genomics:

- Clustering of groups
- Mapping of reads against genome
- Comparative sequence alignment

















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

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



Aligner has a larger effect than	Number of significant genes			
biological signal		12 genomes, M7/8		
	Aligner	95% (a)	99% (b)	
	АМАР	817	213	
	MUSCLE	1043	306	
	ProbCons	1013	281	
	T-Coffee	1290	479	
	ClustalW	902	261	
	PRANK	468	49	
		.00	.,	
	Markova-Raina	& Petrov 2011	Genome Biolog	







How do we avoid Apophenia?

- Double check your tables and analyses
 - Plot your data, look at it, does it make sense?
- Test your hypotheses in an independent way
 - Test your findings using separate data and a different analysis
 - Functional Validation

Published studies allow ...

You to practice your bioinformatics

Assess their repeatability

Papers need enough details for replication







On the importance of negative results

JOURNAL OF NEGATIVE RESULTS

- There is a great need, and little incentive to publish negative results
- How can we change this?
 - Free publication charges
 - Change the name from negative to ?
 - **????**

