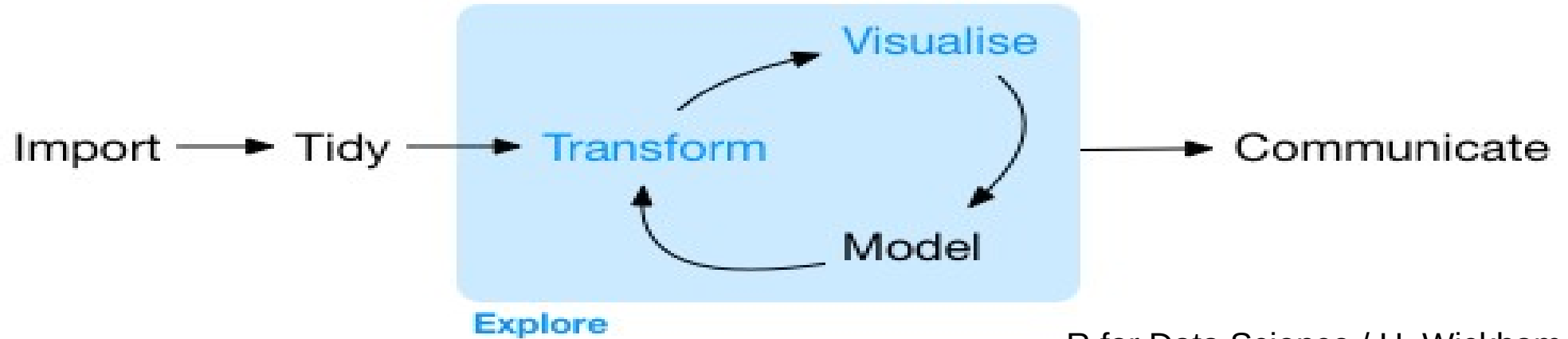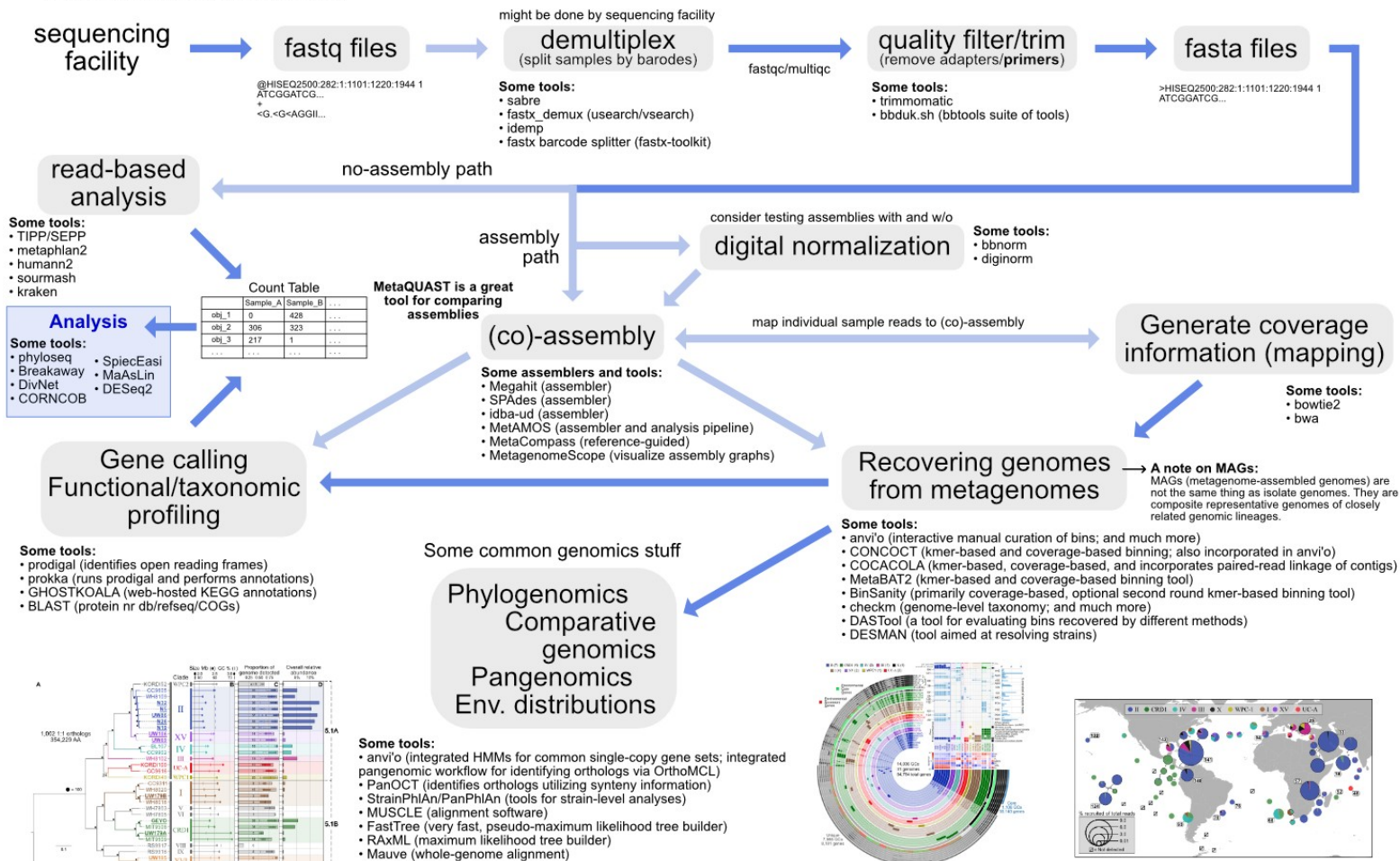# Data science workflow



R for Data Science / H. Wickham

# Overview of generic* metagenomics workflow

*This is generic; specific workflows can vary on the order of steps here and how they are done.

**When working with your own data you should never follow any pipeline blindly. There can be critical differences based on your data.**

sequencing facility → **fastq files**

@HISEQ2500:282:1:1101:1220:1944 1
ATCGGATCG...
+
<G.<G<AGGII...

might be done by sequencing facility

**demultiplex** (split samples by barodes)

**Some tools:**
• sabre
• fastx_demux (usearch/vsearch)
• idemp
• fastx barcode splitter (fastx-toolkit)

fastqc/multiqc

**quality filter/trim** (remove adapters/**primers**)

**Some tools:**
• trimmomatic
• bbduk.sh (bbtools suite of tools)

**fasta files**

>HISEQ2500:282:1:1101:1220:1944 1
ATCGGATCG...

**read-based analysis**

no-assembly path

**Some tools:**
• TIPP/SEPP
• metaphlan2
• humann2
• sourmash
• kraken

assembly path

consider testing assemblies with and w/o

**digital normalization**

**Some tools:**
• bbnorm
• diginorm

**Analysis**

**Some tools:**
• phyloseq         • SpiecEasi
• Breakaway      • MaAsLin
• DivNet           • DESeq2
• CORNCOB

**Count Table**

| | Sample_A | Sample_B | ... |
|---|---|---|---|
| obj_1 | 0 | 428 | ... |
| obj_2 | 306 | 323 | ... |
| obj_3 | 217 | 1 | ... |
| ... | ... | ... | ... |

**MetaQUAST is a great tool for comparing assemblies**

**(co)-assembly**

map individual sample reads to (co)-assembly

**Some assemblers and tools:**
• Megahit (assembler)
• SPAdes (assembler)
• idba-ud (assembler)
• MetAMOS (assembler and analysis pipeline)
• MetaCompass (reference-guided)
• MetagenomeScope (visualize assembly graphs)

**Generate coverage information (mapping)**

**Some tools:**
• bowtie2
• bwa

**Gene calling Functional/taxonomic profiling**

**Some tools:**
• prodigal (identifies open reading frames)
• prokka (runs prodigal and performs annotations)
• GHOSTKOALA (web-hosted KEGG annotations)
• BLAST (protein nr db/refseq/COGs)

**Recovering genomes from metagenomes**

→ **A note on MAGs:**
MAGs (metagenome-assembled genomes) are not the same thing as isolate genomes. They are composite representative genomes of closely related genomic lineages.

**Some tools:**
• anvi'o (interactive manual curation of bins; and much more)
• CONCOCT (kmer-based and coverage-based binning; also incorporated in anvi'o)
• COCACOLA (kmer-based, coverage-based, and incorporates paired-read linkage of contigs)
• MetaBAT2 (kmer-based and coverage-based binning tool)
• BinSanity (primarily coverage-based, optional second round kmer-based binning tool)
• checkm (genome-level taxonomy; and much more)
• DASTool (a tool for evaluating bins recovered by different methods)
• DESMAN (tool aimed at resolving strains)

Some common genomics stuff

**Phylogenomics Comparative genomics Pangenomics Env. distributions**

**Some tools:**
• anvi'o (integrated HMMs for common single-copy gene sets; integrated pangenomic workflow for identifying orthologs via OrthoMCL)
• PanOCT (identifies orthologs utilizing synteny information)
• StrainPhlAn/PanPhlAn (tools for strain-level analyses)
• MUSCLE (alignment software)
• FastTree (very fast, pseudo-maximum likelihood tree builder)
• RAxML (maximum likelihood tree builder)
• Mauve (whole-genome alignment)

astrobiomike.github.io

**Happy Belly Bioinformatics**

JOSE    10.21105/jose.00053

AstrobioMike

Orcid: 0000-0001-7750-9145

# Application:
# contig clustering for metagenome assembly



Repetitive element

DNA

Shotgun reads

Contigs

Collapsed contig

# Binning metagenomic contigs by coverage and composition

Johannes Alneberg, Brynjar Smári Bjarnason, Ino de Bruijn, Melanie Schirmer, Joshua Quick, Umer Z Ijaz, Leo Lahti, Nicholas J Loman, Anders F Andersson ✉ & Christopher Quince ✉

Cluster contigs by *co-occurrence* and *sequence composition*
in order to recover metagenomic species

CONCOCT

# Contig composition (4-bp strings)

Contig: TCGAAATTACGGTCGATTTTAAACTCGGTCTGGA...

k-mer content: "TCGA"

k=4: 136 tetramers (excl. palindromic sequences)

Each contig has a profile over the 136 tetramers

AAAA  AAAC  ....... TGAC  TGCA

Cluster contigs by *co-occurrence* and *sequence composition*
in order to recover metagenomic species

**Coverage**
53 samples

**Composition**
136 tetramers

| Coverage | Composition |
|---|---|

High-dimensional feature space (53 + 136)!

→ drop dimensionality with PCA (d~20; 90%)

→ multivariate gaussian clusters

→ infer the cluster number automatically

# Clustering contigs by coverage and composition (CONCOCT)

Variational Dirichlet Process multivariate mixture models work very well in practical applications with high-dimensional data



Alneberg et al. Nature Methods 2014

# Core microbiota variation (N = 5005)

Z-score across subjects: red – high abundance & blue – low abundance
Core microbiota shows remarkable variation across population.

Community states in vaginal microbiome of reproductive-age women

# Phylogenetic factorization of compositional data yields lineage-level associations in microbiome datasets

Alex D. Washburne ✉[1], Justin D. Silverman[2,3,4,5], Jonathan W. Leff[6], Dominic J. Bennett[7,8], John L. Darcy[9], Sayan Mukherjee[2,10], Noah Fierer[6], Lawrence A. David[2,4,5]

# Model

Θ

# Observations
(Data)

X

Posterior

Likelihood

Prior

$$P(\mathbf{\Theta}|data) \propto P(data|\mathbf{\Theta}) \times P(\mathbf{\Theta})$$

$$P(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Gaussian Likelihood

# N(x|μ,σ)

Likelihood

$$P(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$N(x|\mu,\sigma)\ N(\mu|\mu_0,\ \sigma_0)$$

Likelihood            Prior

$$N(\theta|x) \sim N(x|\mu,\sigma)\ N(\mu|\mu_0, \sigma_0)$$

Posterior          Likelihood          Prior

$\alpha$　$\beta$　$\mu_0$　$\sigma_0^2$

$\sigma^2$　$\mu$

$x$

$$x_i \mid \mu, \tau \sim \mathcal{N}(\mu, \tau) \quad \text{i.i.d.}$$

$$\mu \mid \tau \sim \mathcal{N}(\mu_0, n_0 \tau)$$

$$\tau \sim \text{Ga}(\alpha, \beta)$$

Likelihood   Prior

Data

x

y

μ

posterior

prior

likelihood

23

# Posterior for Gaussian mean

$$\mu \mid x, \tau \sim \mathcal{N}\left(\frac{n\tau}{n\tau + n_0\tau}\bar{x} + \frac{n_0\tau}{n\tau + n_0\tau}\mu_0 \quad, \quad n\tau + n_0\tau\right)$$

A Nasty Looking Likelihood

A Nasty Looking Likelihood

A Nasty Prior

http://doingbayesiandataanalysis.blogspot.com/2013/10/diagrams-for-hierarchical-models-we.html

# Estimating the posterior?

- Analytical solution..
- Gibbs sampling, Markov Chain Monte Carlo (MCMC)
- Variational Inference (VI)
- Hamiltonian Monte Carlo (HMC)
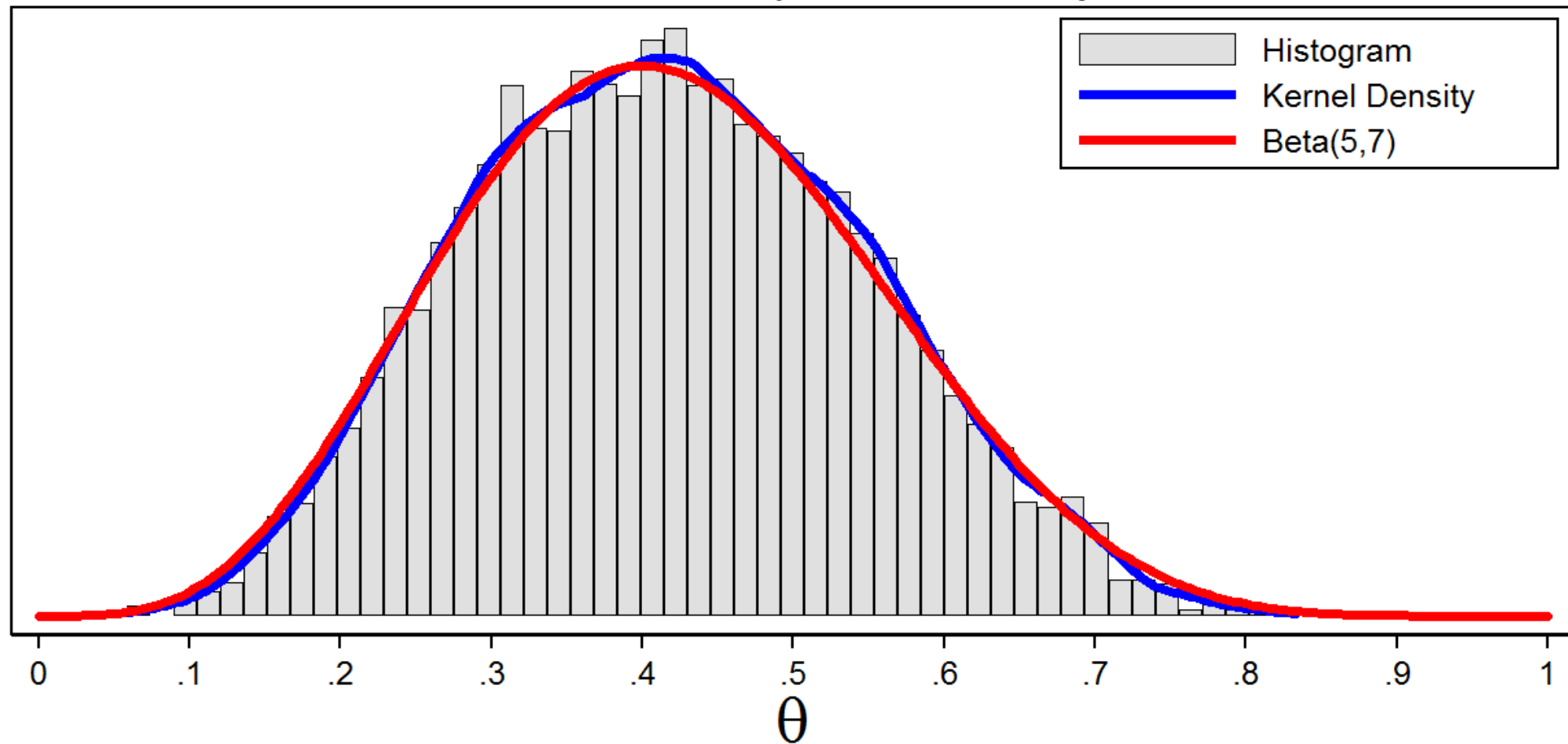- Approximate Bayesian Computation (ABC)
- ...

Our Best Guess of the Posterior Distribution

https://towardsdatascience.com/a-zero-math-introduction-to-markov-chain-monte-carlo-methods-dcba889e0c50

Our Best Guess of the Posterior Distribution

https://towardsdatascience.com/a-zero-math-introduction-to-markov-chain-monte-carlo-methods-dcba889e0c50

Our Best Guess of the Posterior Distribution

https://towardsdatascience.com/a-zero-math-introduction-to-markov-chain-monte-carlo-methods-dcba889e0c50

Our Best Guess of the Posterior Distribution

https://towardsdatascience.com/a-zero-math-introduction-to-markov-chain-monte-carlo-methods-dcba889e0c50

$n = 3000, \pi \approx 3.1133$

https://towardsdatascience.com/a-zero-math-introduction-to-markov-chain-monte-carlo-methods-dcba889e0c50

Draw $\theta_t \sim \text{Normal}(\theta_{t-1}, \sigma)$

$\text{Normal}(0.500, \sigma) = 0.497$

Comparison of the MCMC sample and
the theoretical posterior density

Legend:
- Histogram
- Kernel Density
- Beta(5,7)

$\theta$

# 15 year **prospective** view (Finland / FINRISK2002)

**2002** ~7000+ stool samples: omics and clinical measurements.

**2017** comprehensive health information from Finnish registers

15+ year follow up



N=7231

Area
- Helsinki
- Karelia
- Kuopio
- Lapland
- Oulu
- Turku



Deaths:
667

# Taxonomic profiles (N = 5005)

Z-score across subjects: red – high abundance & blue – low abundance - HITChip Atlas

# Cox proportional hazards

Hazard function

Baseline hazard

Data

Coefficients

$$\lambda(t) = \lambda_0(t) \exp(\mathbf{x}\beta).$$

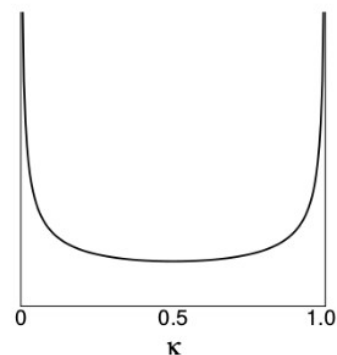$$\lambda(t) = \lambda_0(t) \exp(\mathbf{x}\beta).$$

Alternative priors: P(β)

Sparsity!

Horseshoe
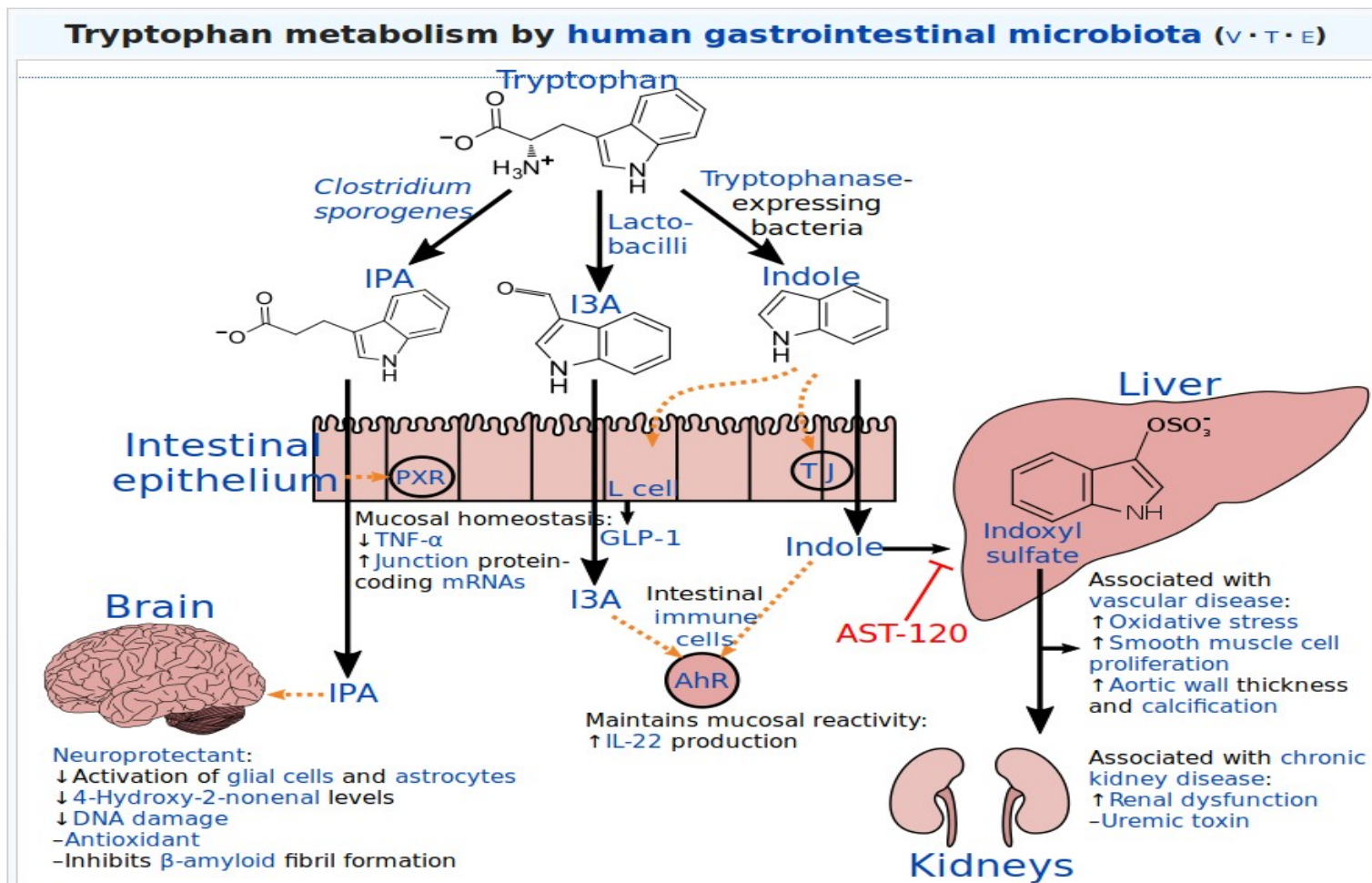
Gaussian prior

Horseshoe prior

Aki Vehtari

# Function, multi-omics, causality, mechanisms



Tryptophan metabolism by human gastrointestinal microbiota (v · t · e)
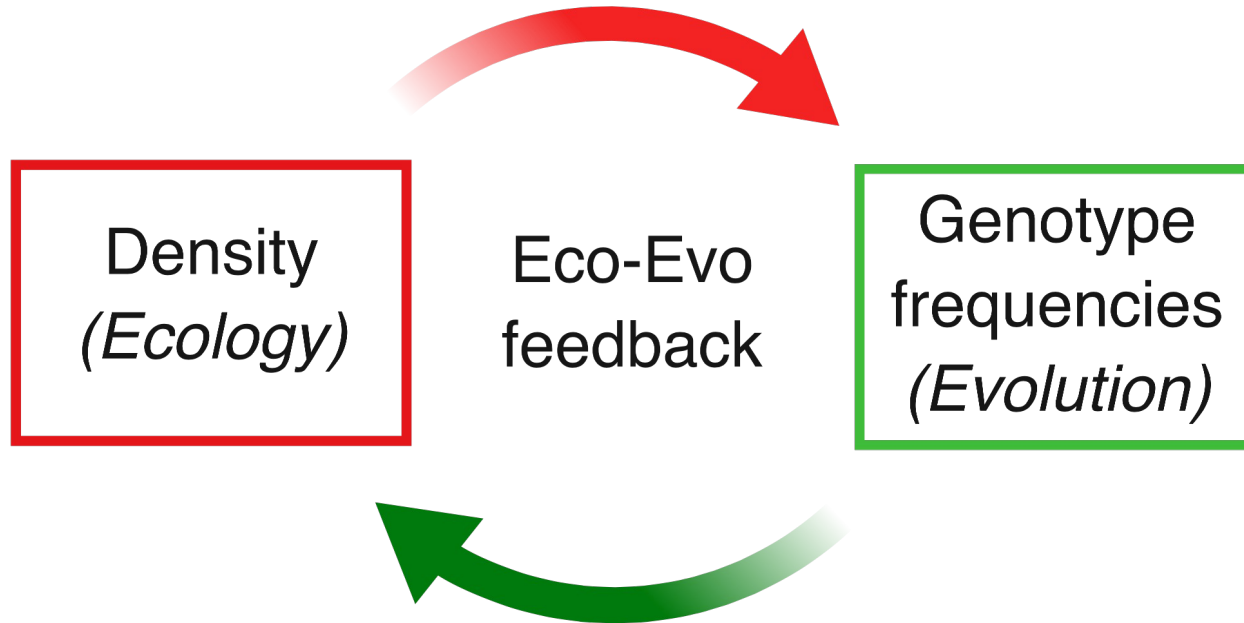
$$\begin{cases} X = W_x \mathbf{z} + \varepsilon_x \\ Y = W_y \mathbf{z} + \varepsilon_y \end{cases}$$

Multi-view learning

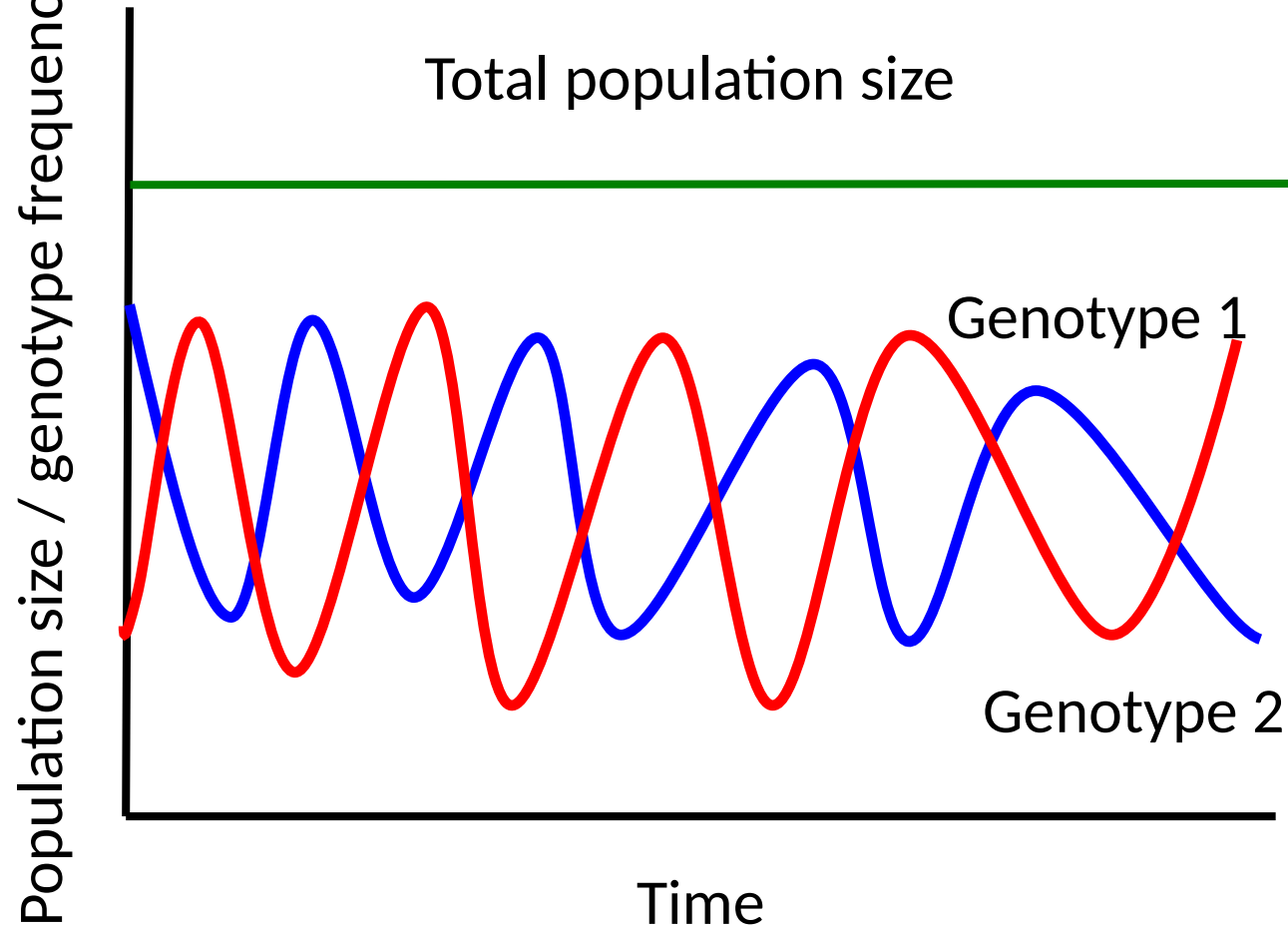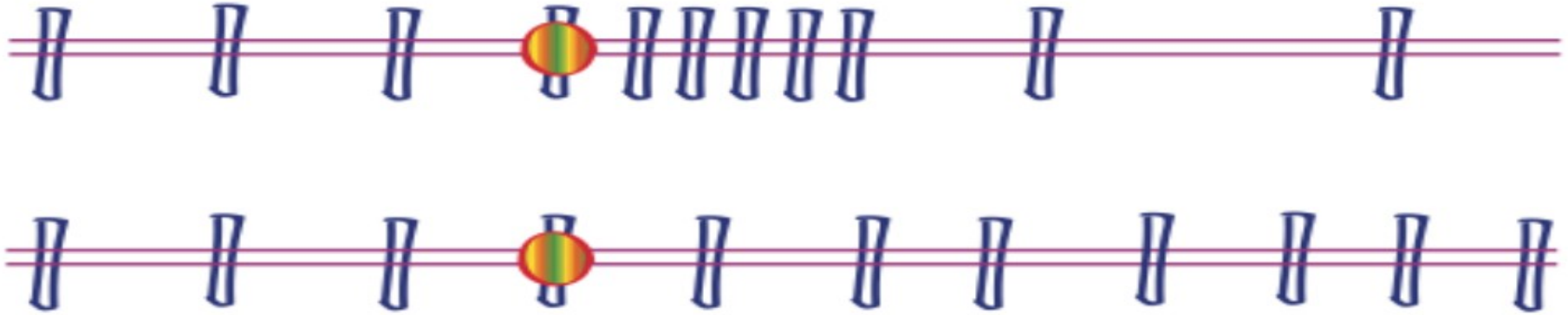# Ecology and evolution are connected processes

# We need information about the underlying genotype dynamics as well

# Longitudinal Data Analyses

# Eco-evo

# The genomic basis of Red Queen dynamics during rapid reciprocal host–pathogen coevolution

Andrei Papkou, (iD) Thiago Guzella, Wentao Yang, Svenja Koepper, Barbara Pees, Rebecca Schalkowski, Mike-Christoph Barg, Philip C. Rosenstiel, (iD) Henrique Teotónio, and Hinrich Schulenburg

Check for updates

# Reconciling taxon senescence with the Red Queen's hypothesis

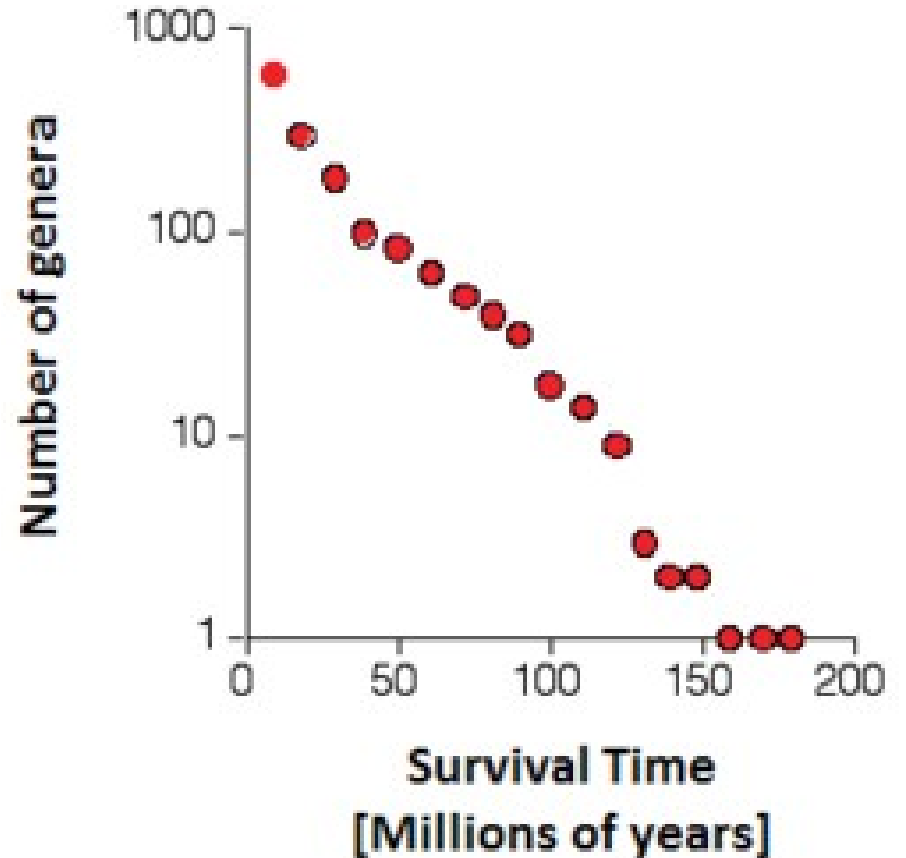Indrė Žliobaitė ✉, Mikael Fortelius & Nils C. Stenseth
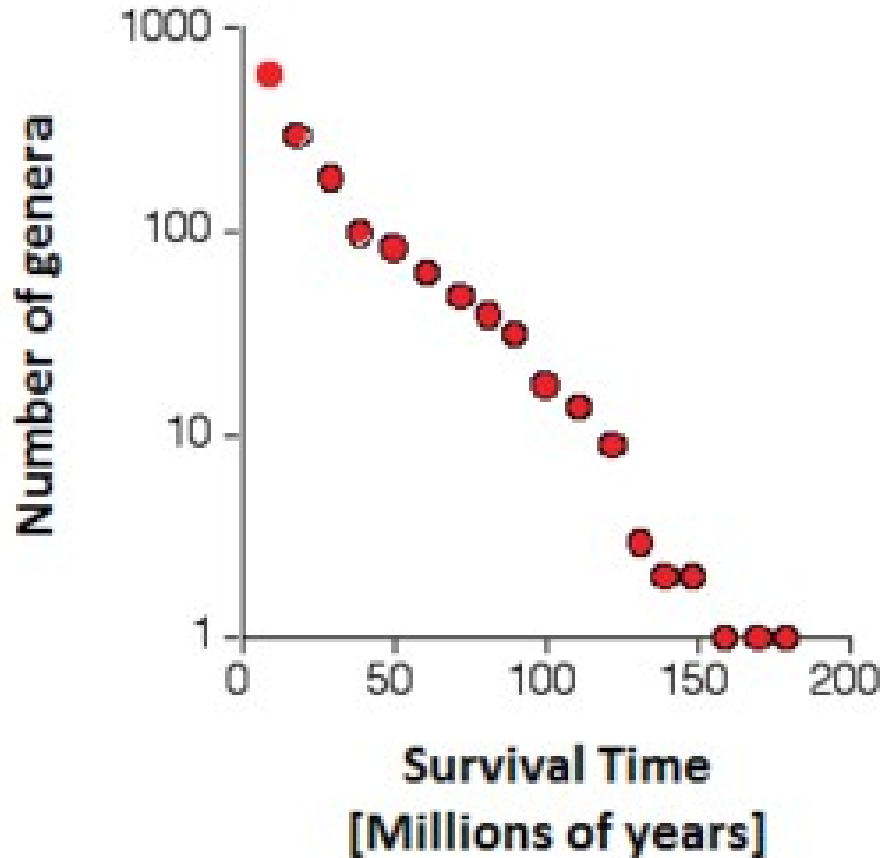
# Do species age?

**Hypothesis:** the longer a species has been around, the more likely it will go extinct.

**Rejection:** Fossil data did not support this.

Old and young species have a similar probability to go extinct at any point.

**(Leigh) Van Valen's "Law of Extinction":** Probability of extinction does not depend on the lifetime of the population / "Half-life of species"



A NEW EVOLUTIONARY LAW

Leigh Van Valen
Department of Biology
The University of Chicago
Chicago, Illinois  60637

ABSTRACT:

All groups for which data exist go extinct at a rate that is constant for a given group. When this is recast in ecological form (the effective environment of any homogeneous group of organisms deteriorates at a stochastically constant rate), no definite exceptions exist although a few are possible. Extinction rates are similar within some very broad categories and vary regularly with size of area inhabited. A new unit of rates for discrete phenomena, the macarthur, is introduced. Laws are appropriate in evolutionary biology. Truth needs more than correct predictions. The Law of Extinction is evidence for ecological significance and comparability of taxa. A non-Markovian hypothesis to explain the law invokes mutually incompatible optima within an adaptive zone. A self-perpetuating fluctuation results which can be stated in terms of an unstudied aspect of zero-sum game theory. The hypothesis can be derived from a view that momentary fitness is the amount of control of resources, which remain constant in total amount. The hypothesis implies that long-term fitness has only two components and that events of mutualism are rare. The hypothesis largely explains the observed pattern of molecular evolution.

Van Valen, Leigh (1973). "A new evolutionary law". Evolutionary Theory. 1: 1–30.

# Explanation?

**Red Queen Hypothesis:** organisms must constantly adapt, evolve, and proliferate in order to survive. A classic evolutionary theory (Van Valen, 1973).

Red Queen explains how her country differs from Alice's:

*"Now, here, you see, it takes all the running you can do, to keep in the same place. If you want to get somewhere else, you must run at least twice as fast as that!"*

A search for '**Red Queen**' on Google Scholar gives over a million hits, reflecting the enormous influence this idea has had and continues to have in a wide sector of science.
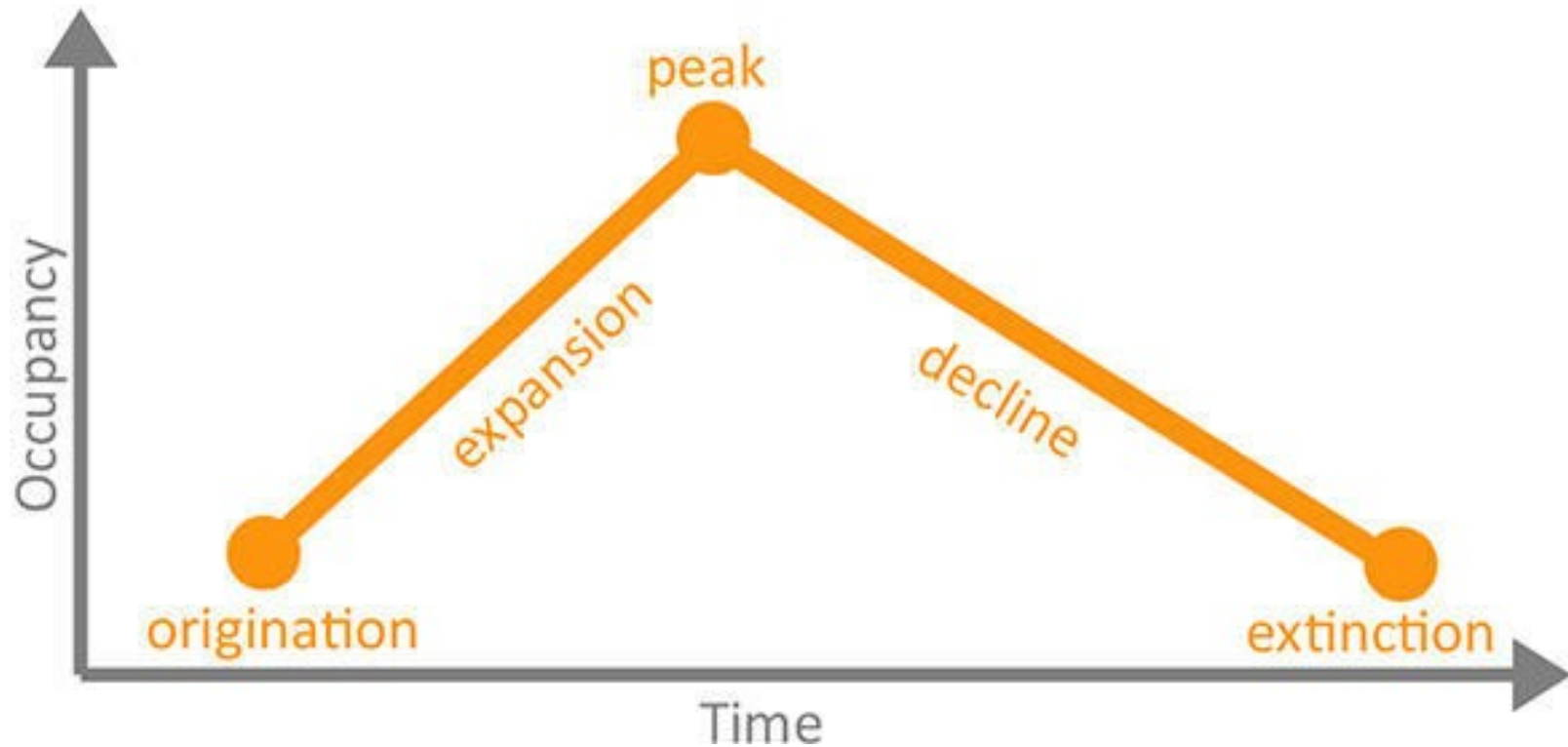


Image credit: Ika Österblad

Stenseth: "The Red Queen's Hypothesis has fascinated me from the very beginning since it, as an evolutionary hypothesis, explicitly brings in ecological interactions to explain large scale evolutionary patterns, such as rate of extinctions."


Image credit: Ika Österblad

# Hat pattern

**Contradiction!** Jernvall & Fortelius (2004). Do species age or not?
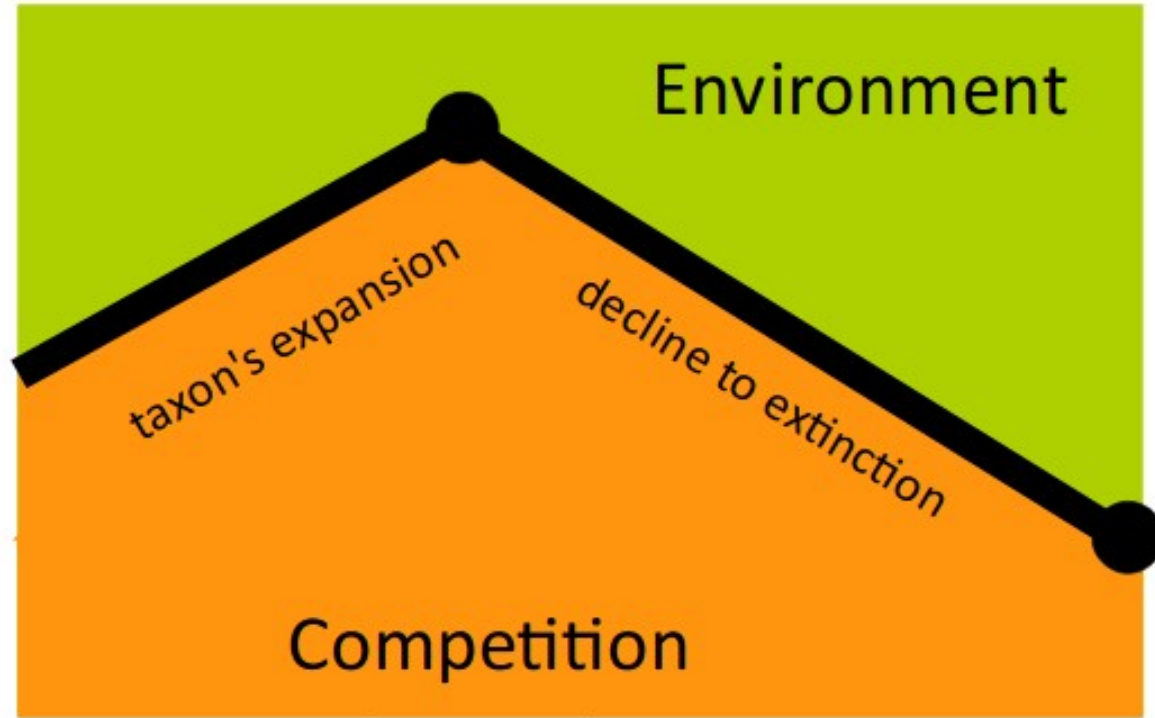
Originally reported in: Jernvall, J.&Fortelius, M.Maintenance of trophic structure in fossil mammal communities: site occupancy and taxon resilience. The American Naturalist164, 614-624 (2004).

**Expansion, stabilization, decline**

"*It takes all the running you can do, to keep in the same place*. **If you want to get somewhere else, you must run at least twice as fast as that!**"
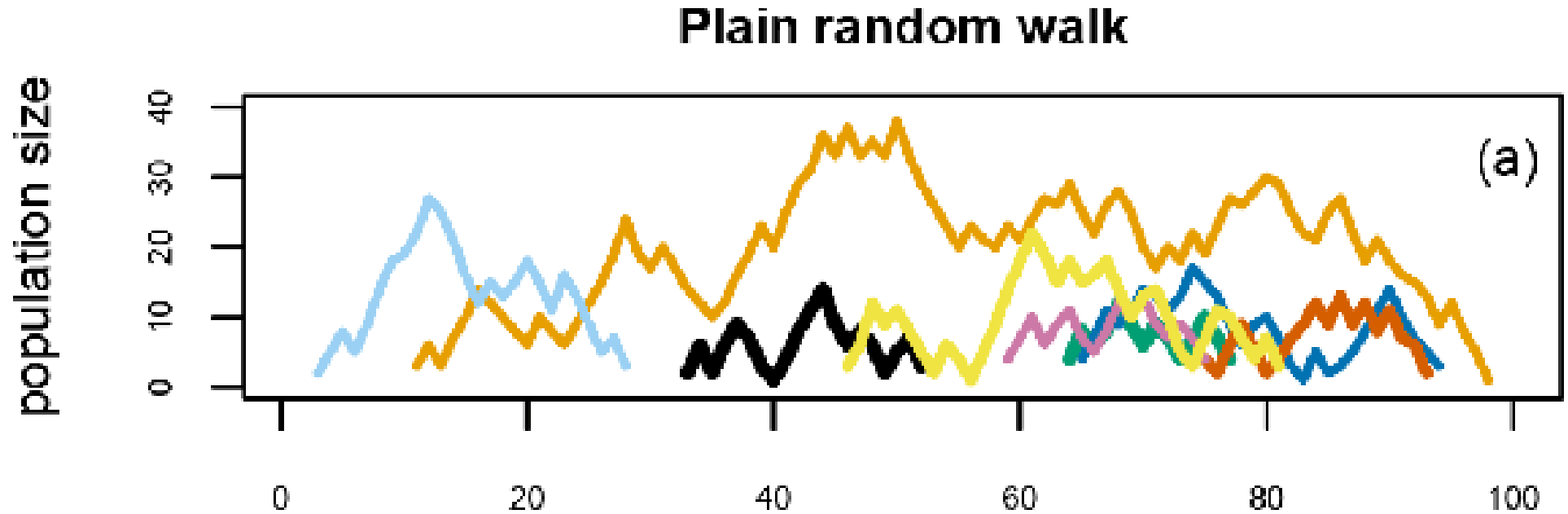
It is likely that a species will
get extinct by random chance
in the long run.

More than three peaks are very
seldom observed in real data.



**Plain random walk**

population size

(a)

**Plain random walk**

**Random walk with memory (correlated random walk)**

Increasing effect of Environmental change

Environment

taxon's expansion

decline to extinction

Expected contributions to evolutionary rates

Competition

Increasing effect of Competition

Memory effects emphasize the Red Queen effect.

**Random walk with memory and competition**

**Random walk with environmental change**

time units

**Resolution:** Hat Pattern *is* compatible with the Law of Constant Extinction, and predicted by the Red Queen's Hypothesis.

But the fact is that they didn't. In the end it took the fresh perspective of an outsider to realise that the two theories were connected and actually parts of the same puzzle. Zliobaite, with a background in computer science and credit analysis, who knew that a cessation of growth may signal the impending failure of businesses.



Through the Looking-Glass, and What Alice Found There (Lewis Carroll, 1871)

# How to choose a correct model?



Parametric assumptions:
(1) Independent samples
(2) Data normally distributed
(3) Equal variances

Type of data? — Discrete, categorical → Any counts < 5?

Continuous

Any counts < 5? — No → Chi-square tests, one and two sample

Any counts < 5? — Yes → Fisher's exact test

Type of question?

Relationships

Differences

Do you have dependent & independent variables?

Differences between what? — Means → One-sample t-test

Differences between what? — Variances → Fmax test or Bartlett's test

Multiple means Single variable
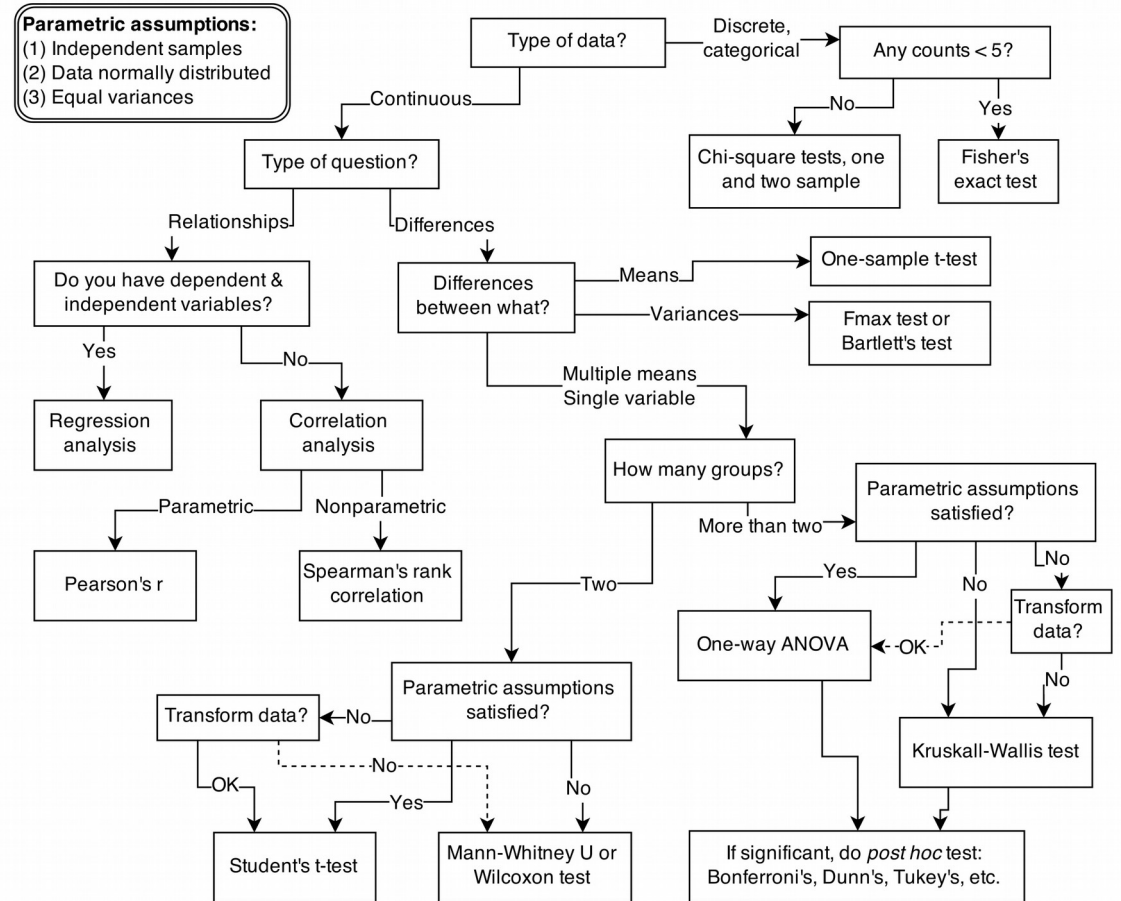
Yes → Regression analysis

No → Correlation analysis

How many groups?

Parametric assumptions satisfied?

More than two

Parametric — Pearson's r

Nonparametric — Spearman's rank correlation

Yes → One-way ANOVA ← OK

No → Transform data?

No → Transform data?

No → Kruskall-Wallis test

Two

Parametric assumptions satisfied?

Transform data? — No

OK

No

Yes

No

Student's t-test

Mann-Whitney U or Wilcoxon test

If significant, do post hoc test: Bonferroni's, Dunn's, Tukey's, etc.

# Statistical Rethinking

## A Bayesian Course with Examples in R and Stan

afternoon wait (mins)

morning wait (mins)

Richard McElreath

# Anvi'o in a nutshell

## Anvi'o: an advanced analysis and visualization platform for 'omics data

A. Murat Eren[1,2], Özcan C. Esen[1], Christopher Quince[3], Joseph H. Vineis[1], Hilary G. Morrison[1], Mitchell L. Sogin[1], Tom O. Delmont[1]

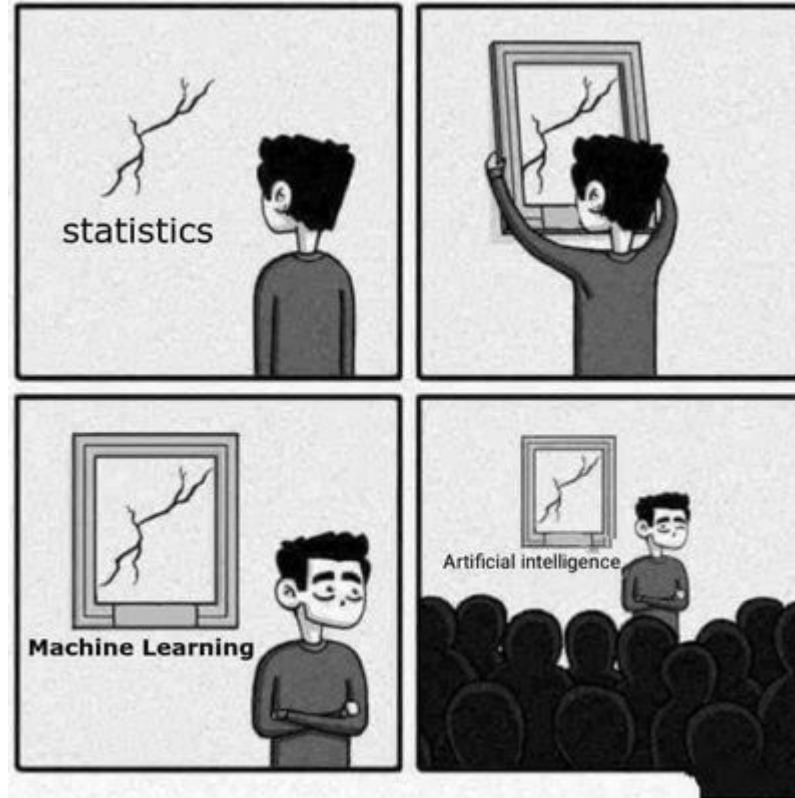Anvi'o is an open-source, community-driven analysis and visualization platform for 'omics data.

# Summary & Conclusions

## Statistical aspects

- supervised / unsupervised
- exploration / modeling
- mechanistic / non-mechanistic
- parametric / non-parametric
- deterministic / stochastic

# Learning goals

- Statistical thinking in microbiome studies

- The concept of open and reproducible research

- Familiarity with standard tools in amplicon profiling

- Looking at your own research problems in new ways

- Networking & collaboration!

WHEN YOU SEE A CLAIM THAT A COMMON DRUG OR VITAMIN "KILLS CANCER CELLS IN A PETRI DISH,"

KEEP IN MIND:



SO DOES A HANDGUN.