# Microbial Co-Exclusion and Co-Occurrence:

# Making it All Add Up

Andrew Fernandes
Jean Macklaim
Gregor Reid
Gregory Gloor

Western

Schulich
MEDICINE & DENTISTRY

International Human Microbiome Congress
Hyatt Regency Vancouver
March 9-11,2011

Competitive
Exclusion
Principle

"Complete
Competitors
Can't Coexist"

They live all
Packed together

Or we can have
**mutualists** that
*love* to coexist!

# How do they interact?
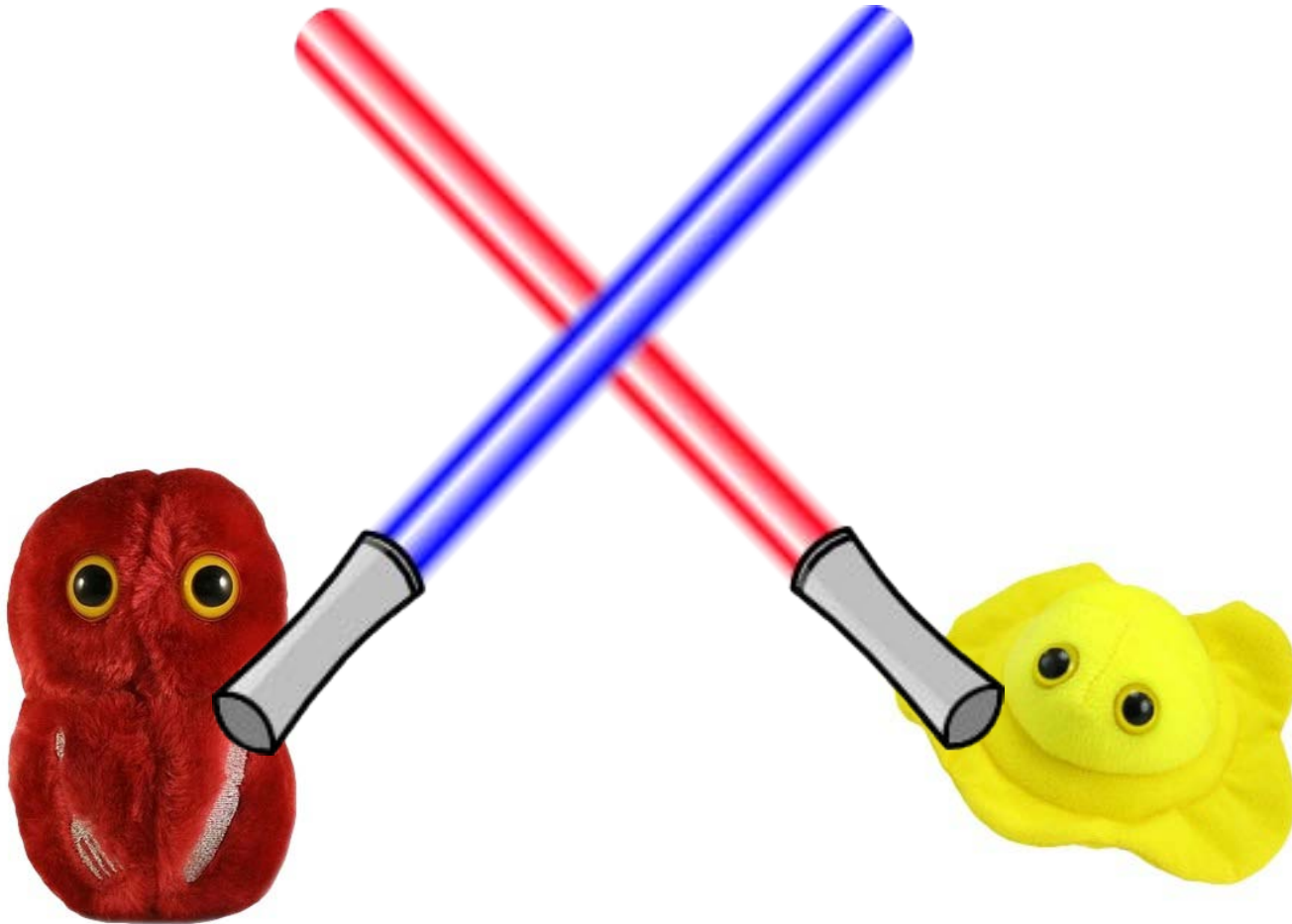
Consider just
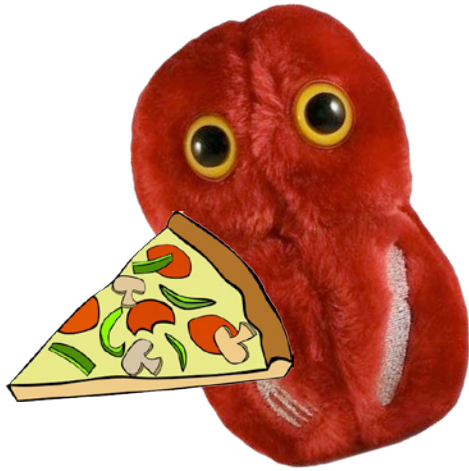**two** species...

... living together
in an environment.

# Co-Occurrence?

# Indifference?

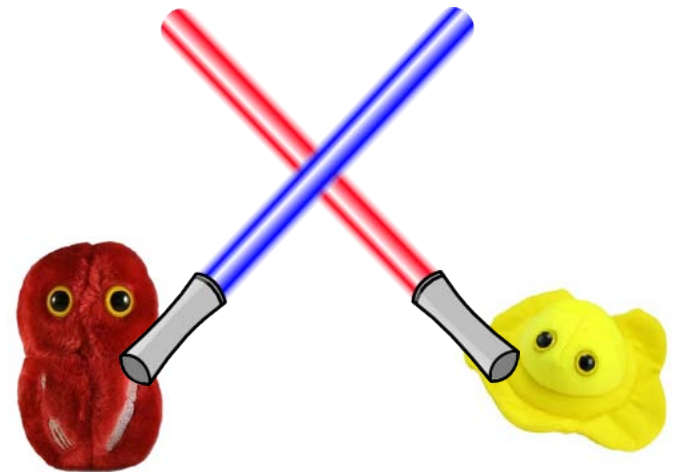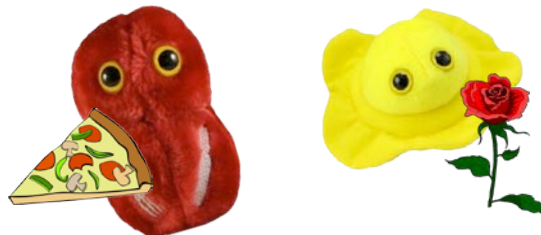# How can we tell *which* of these is occurring?

**Co-Occurrence**

**Co-Exclusion**
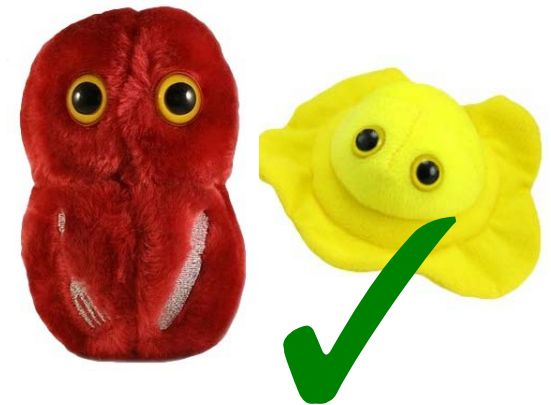
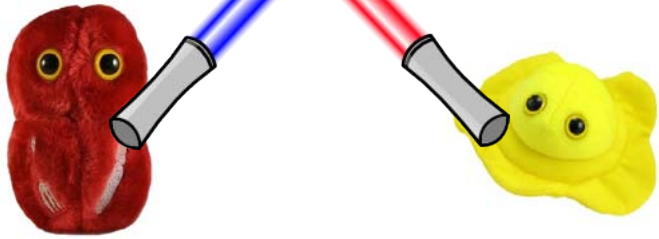**Indifference**

Co-Occurrence

Marginal Relative Abundance

Marginal Relative Abundance

**Co-Exclusion**
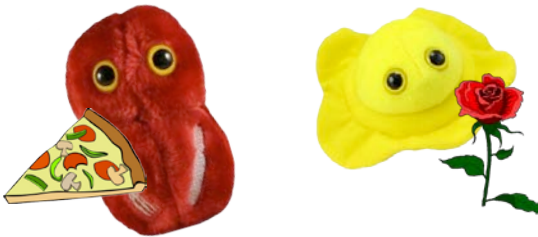
Marginal Relative Abundance
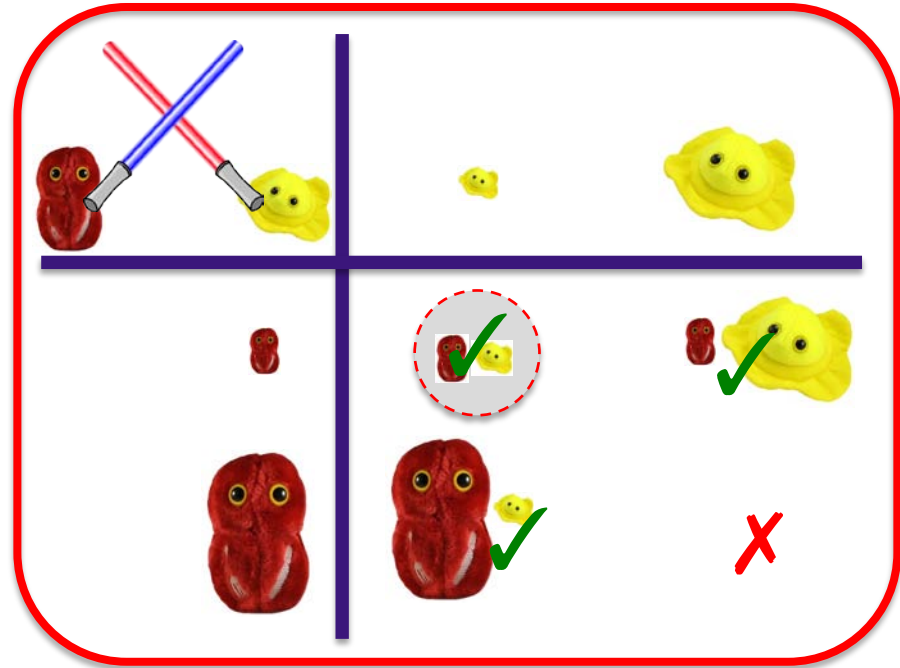
Marginal Relative Abundance
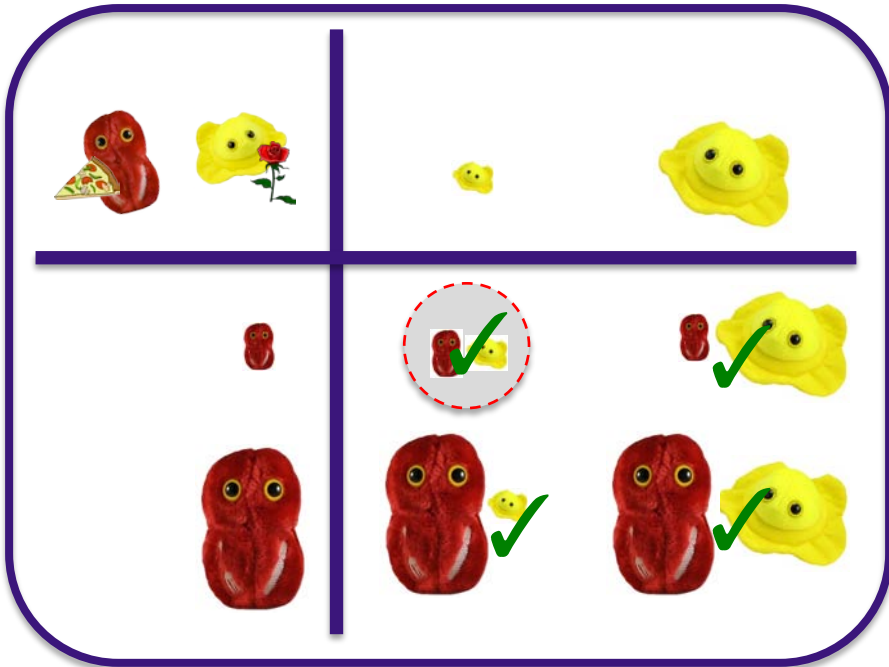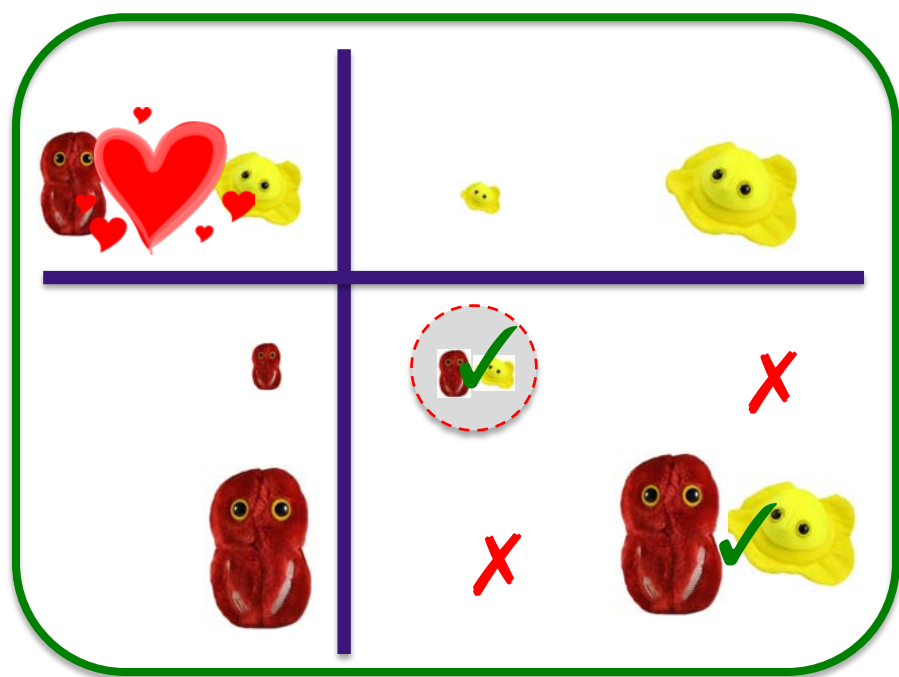
# Indifference

Marginal Relative Abundance

Marginal Relative Abundance

Let's recap…

These are **subtle** differences!

Human Vaginal Microbiome

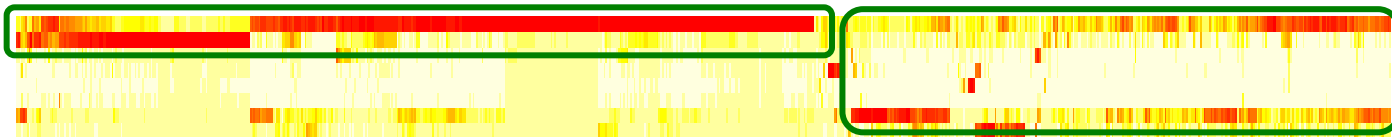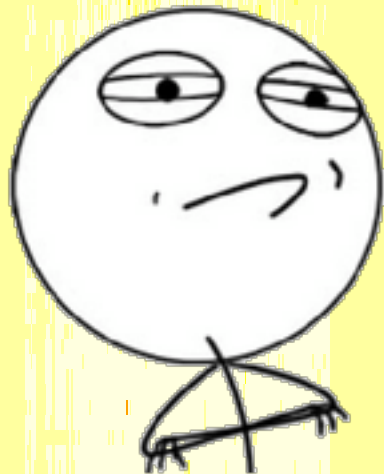"Hey! Look at these *great* results! They show **obvious** co-exclusion!" *

(normal *vs.* disease state)

Bacterium

seq_0_Lactobacillus_iners
seq_2_Lactobacillus_crispatus
seq_31_Lactobacillus_jensenii
actobacillus_gasseri.johnsonii

CHALLENGE ACCEPTED

"Sorry, boss. No they don't."

"What?! Sure they do!"

"Nope."

"Okay... **prove** it!"

Patients

- Gregory B. Gloor, Ruben Hummelen, Jean M. E. Macklaim, **Andrew D. Fernandes**, and Gregor Reid (2010, *Accepted*) Community Microbiome Profiling by Combinatorial Barcoding with Illumina Sequencing. *PLoS One* (Manuscript PONE-D-10-00044R1). Archived at ⟨http://arxiv.org/abs/1007.5075v1⟩.

- Ruben Hummelen, **Andrew D. Fernandes**, Jean M. E. Macklaim, Russell J. Dickson, John Changalucha, Gregory B. Gloor, and Gregor Reid (2010, *Accepted*) Deep Sequencing of the Vaginal Microbiota of Women With HIV. *PLoS One* (Manuscript 10-PONE-RA-19937).

**\*Dramatic Re-interpretation**
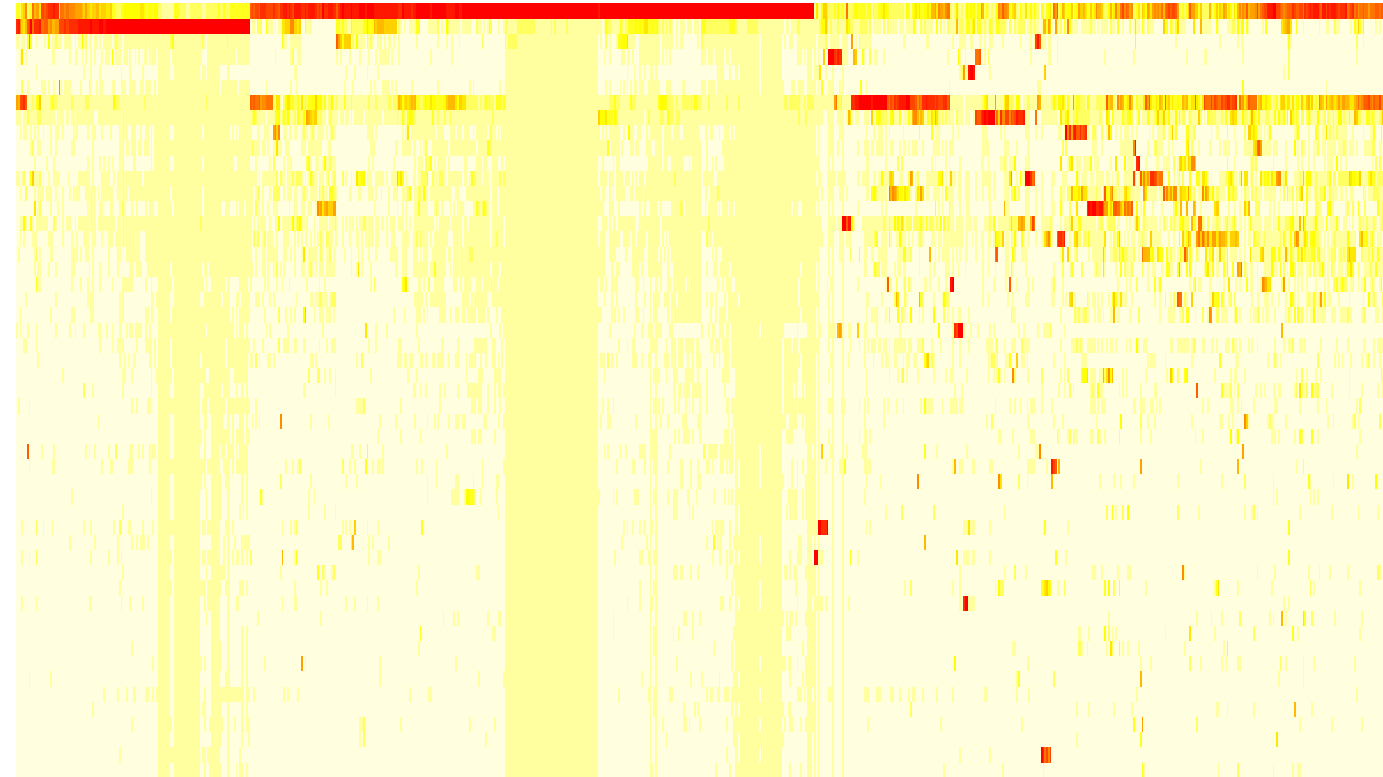
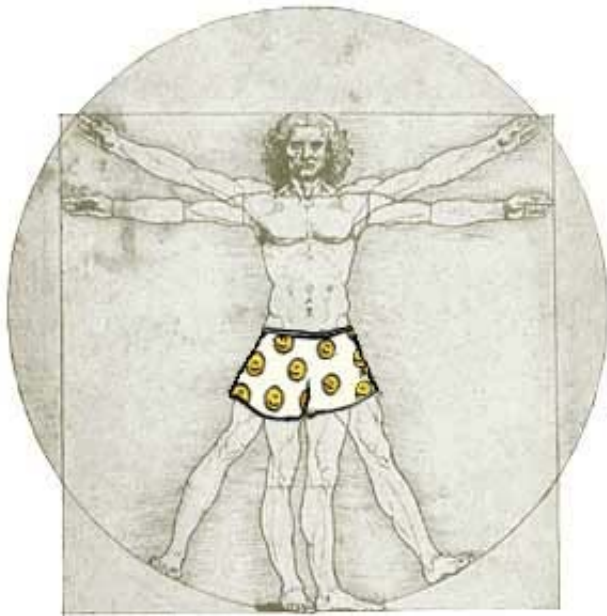# Where do results like this come from?

# Claim: Zeros, Logarithms, and PCA!



**Bacterium**

seq_0_Lactobacillus_iners
seq_2_Lactobacillus_crispatus
seq_31_Lactobacillus_jensenii
actobacillus_gasseri.johnsonii

**Patients**

# Zeros

Clinical Subject
(well, female...)

Swab Sample

PCR

Multiplex Next-Gen Sequence

| | *Lactobacillus iners* | *Lactobacillus crispatus* | *Lactobacillus jensenii* | *Lactobacillus gasseri* |
|---|---|---|---|---|
| Patient 1 | 49679 | 3177 | 21389 | 135 |
| Patient 2 | 7755 | 29752 | 989 | 368 |
| Patient 3 | 3286 | 5955 | 549 | 397 |
| Patient 4 | 2265 | 3263 | 13742 | 148 |
| Patient 5 | 10239 | 2926 | 226 | 100 |
| Patient 6 | 16376 | 20706 | 1037 | 79 |
| Patient 7 | 27313 | 4878 | 5320 | 92 |
| Patient 8 | 33006 | 1103 | 1186 | 176 |
| Patient 9 | 20504 | 1771 | 346 | 161 |

**V6 Read <u>Counts</u>**

ion torrent

(actually Illumina & Solid...)

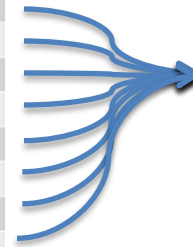# Counts *between samples (patients)* are **meaningless**!

## We are interested in sample **proportions**, only!

| | *Lactobacillus iners* | *Lactobacillus crispatus* | *Lactobacillus jensenii* | *Lactobacillus gasseri* | … … … |
|---|---|---|---|---|---|
| **Patient 1** | 49679 | 3177 | 21389 | 135 | … |
| **Patient 2** | 7755 | 29752 | 989 | 368 | … |
| **Patient 3** | 3286 | 5955 | 549 | 397 | … |
| **Patient 4** | 2265 | 3263 | 13742 | 148 | … |
| **Patient 5** | 10239 | 2926 | 226 | 100 | … |
| **Patient 6** | 16376 | 20706 | 1037 | 79 | … … |
| **Patient 7** | 27313 | 4878 | 5320 | 92 | … |
| **Patient 8** | 33006 | 1103 | 1186 | 176 | … |
| **Patient 9** | 20504 | 1771 | 346 | 161 | … … |

But proportions and counts are kind-of-almost the same thing, aren't they?

| | Lactobacillus iners | Lactobacillus crispatus | Lactobacillus jensenii | Lactobacillus gasseri |
|---|---|---|---|---|
| Patient 1 | 49679 | 3177 | 21389 | 135 |
| Patient 2 | 7755 | 29752 | 989 | 368 |
| Patient 3 | 3286 | 5955 | 549 | 397 |
| Patient 4 | 2265 | 3263 | 13742 | 148 |
| Patient 5 | 10239 | 2926 | 226 | 100 |
| Patient 6 | 16376 | 20706 | 1037 | 79 |
| Patient 7 | 27313 | 4878 | 5320 | 92 |
| Patient 8 | 33006 | 1103 | 1186 | 176 |
| Patient 9 | 20504 | 1771 | 346 | 161 |

$$p_i \approx \frac{n_i}{\sum_j n_j}$$

**Nope.**

**ONLY VALID IF $n_i$ IS NOT SMALL!**

# So where does $p_i \approx \frac{n_i}{\sum_j n_j}$ come from?

Consider a biome with only **two** species.

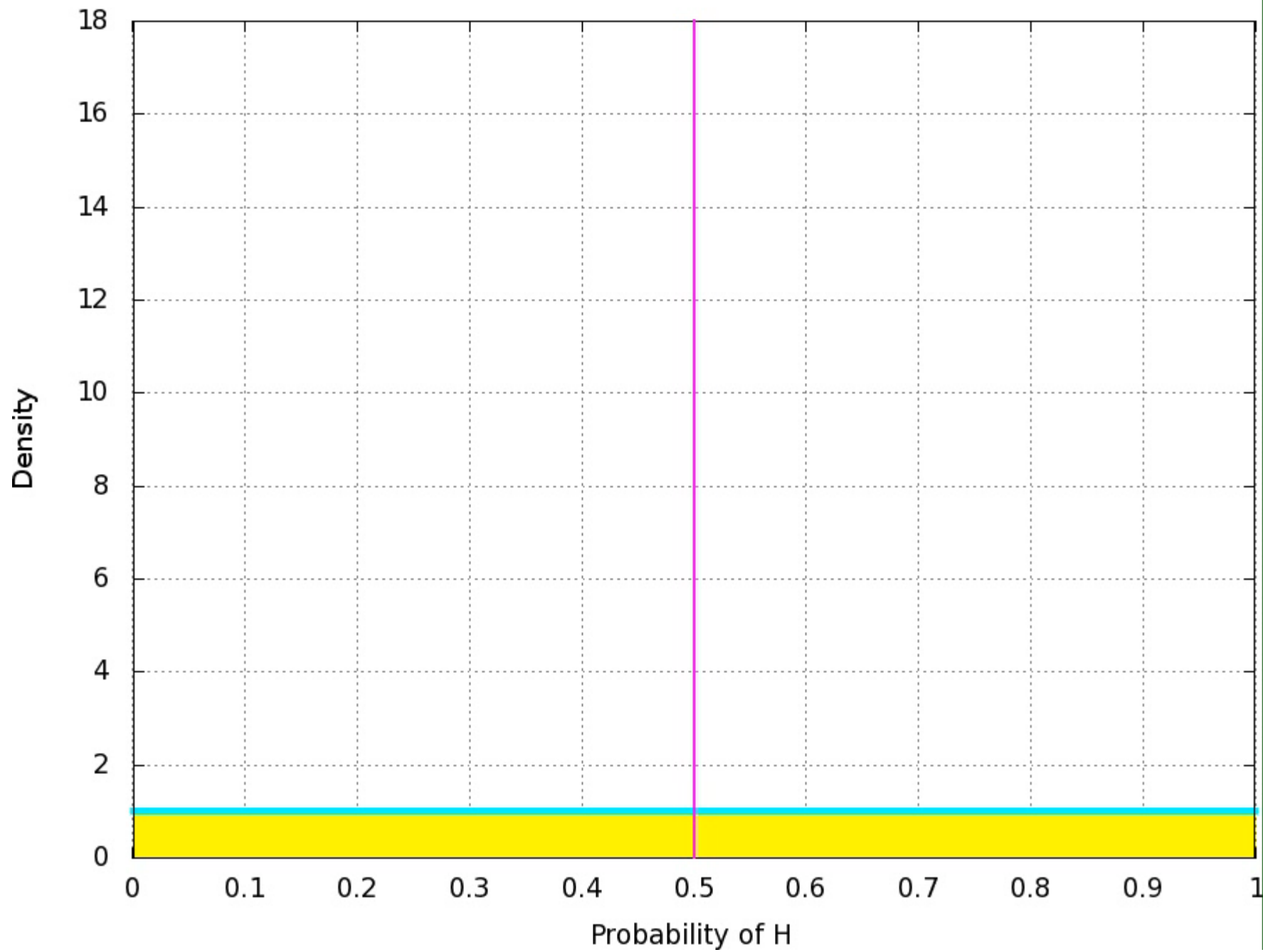$$\Pr\left(p_H, p_T \,|\, n_H, n_T\right) = \frac{(n_H + n_T)!}{n_H! \, n_T!} \, p_H^{n_H} \, p_T^{n_T}$$

$$
\begin{aligned}
\log\left[\Pr\left(p_H, p_T \,|\, n_H, n_T\right)\right] &= \log\left[\frac{(n_H + n_T)!}{n_H! \, n_T!}\right] + n_H \log\left(p_H\right) + n_T \log\left(p_T\right) \\
&= \log\left[\frac{(n_H + n_T)!}{n_H! \, n_T!}\right] + n_H \log\left(p_H\right) + (n - n_H) \log\left(1 - p_H\right)
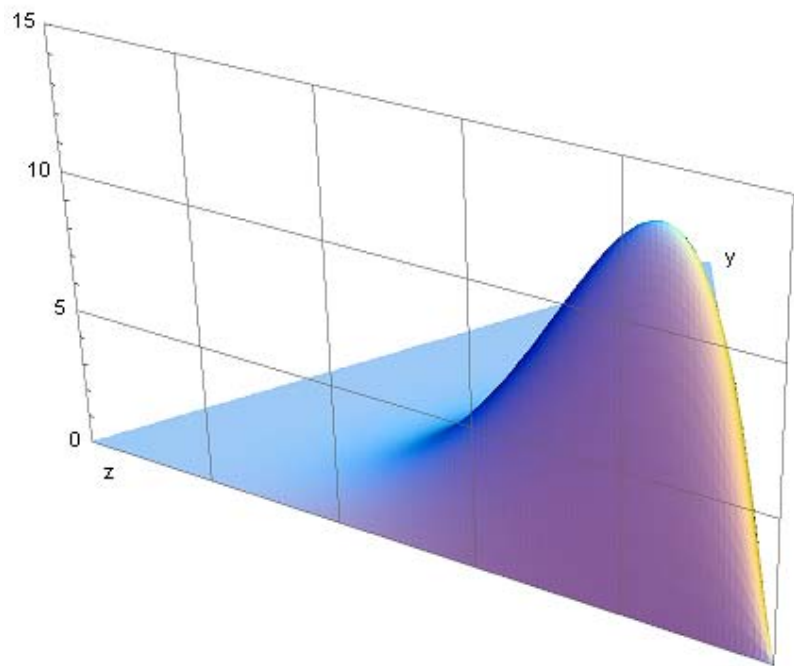\end{aligned}
$$

$$\partial/\partial p_H = 0 \qquad \Rightarrow \qquad \frac{n_H}{p_H} - \frac{n - n_H}{1 - p_H} = 0$$

$$\Rightarrow \qquad \Pr\left(p_H \,|\, n, n_H\right) \quad \text{is maximzed when} \quad p_H = \frac{n_H}{n}$$
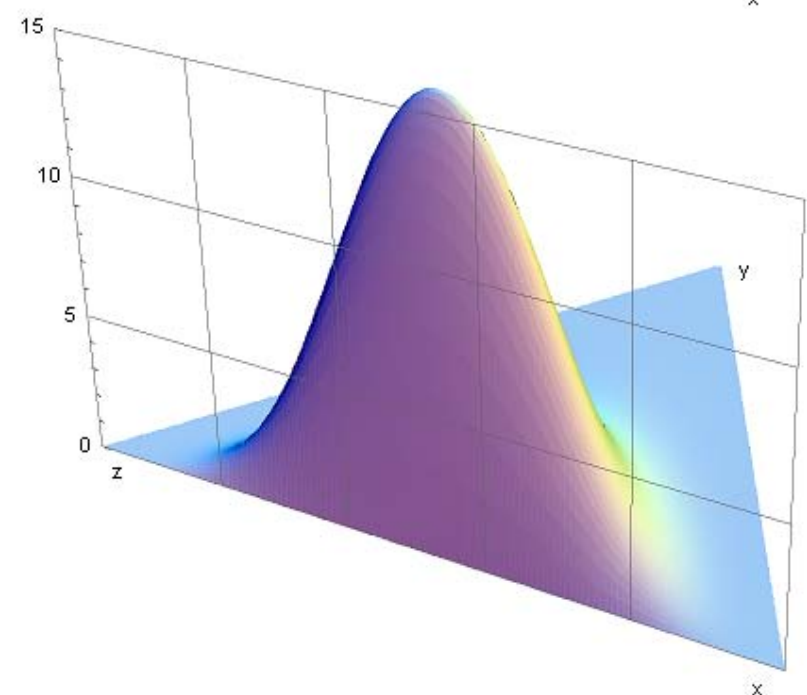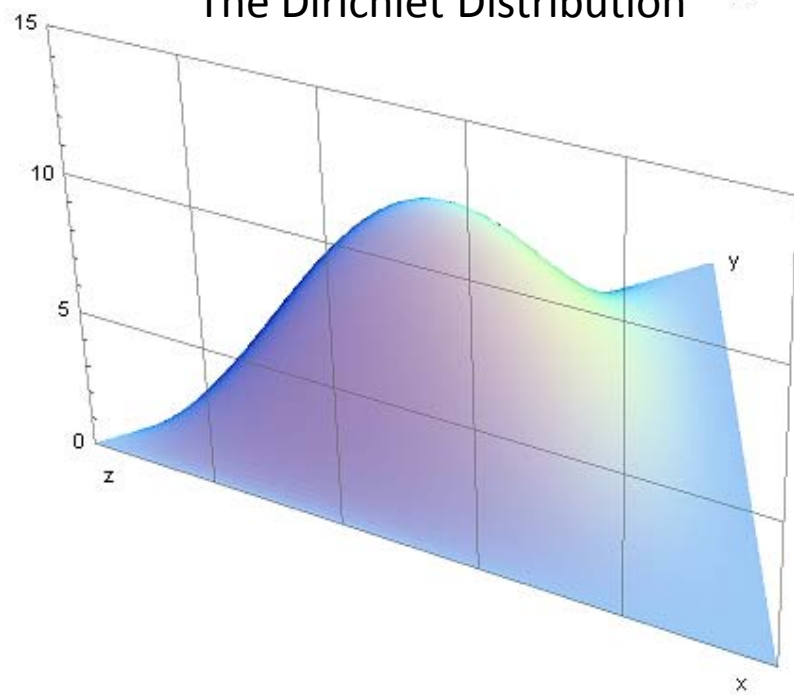
That's a **lot** of math for something that most people think is rather obvious…

N:000, H:000, T:000

The Dirichlet Distribution

The **point**… mathematically, even though

$$p_i \approx \frac{n_i}{\sum_j n_j} \qquad n_i = 0 \;\not\Rightarrow\; p_i = 0$$

# There is a **enormous** difference between a *biological* zero and a *mathematical* zero!

(Distinguishing "rare" *vs.* "impossible" is almost the entire
basis of Shannon's Theory of Information…)

$$\mathbb{E}\left[\log(p_i)\right] = \psi(\alpha_i) + \psi\left(\sum_j \alpha_j\right)$$

**(Never Zero!)**

$$\psi(z) = \frac{d}{dz}\ln\Gamma(z) \qquad \alpha_i = n_i + \frac{1}{2}$$

You can think of this as a "pseudo-count"

(You would be completely wrong, but at least you would feel comforted…)

# Logarithms

Why use $\log(p_i)$ ?

Both of these are **biologically** wrong!

"0.01% to 0.02% = 100% change!"

"1% abundance of a virulent pathogen **is** negligible!"

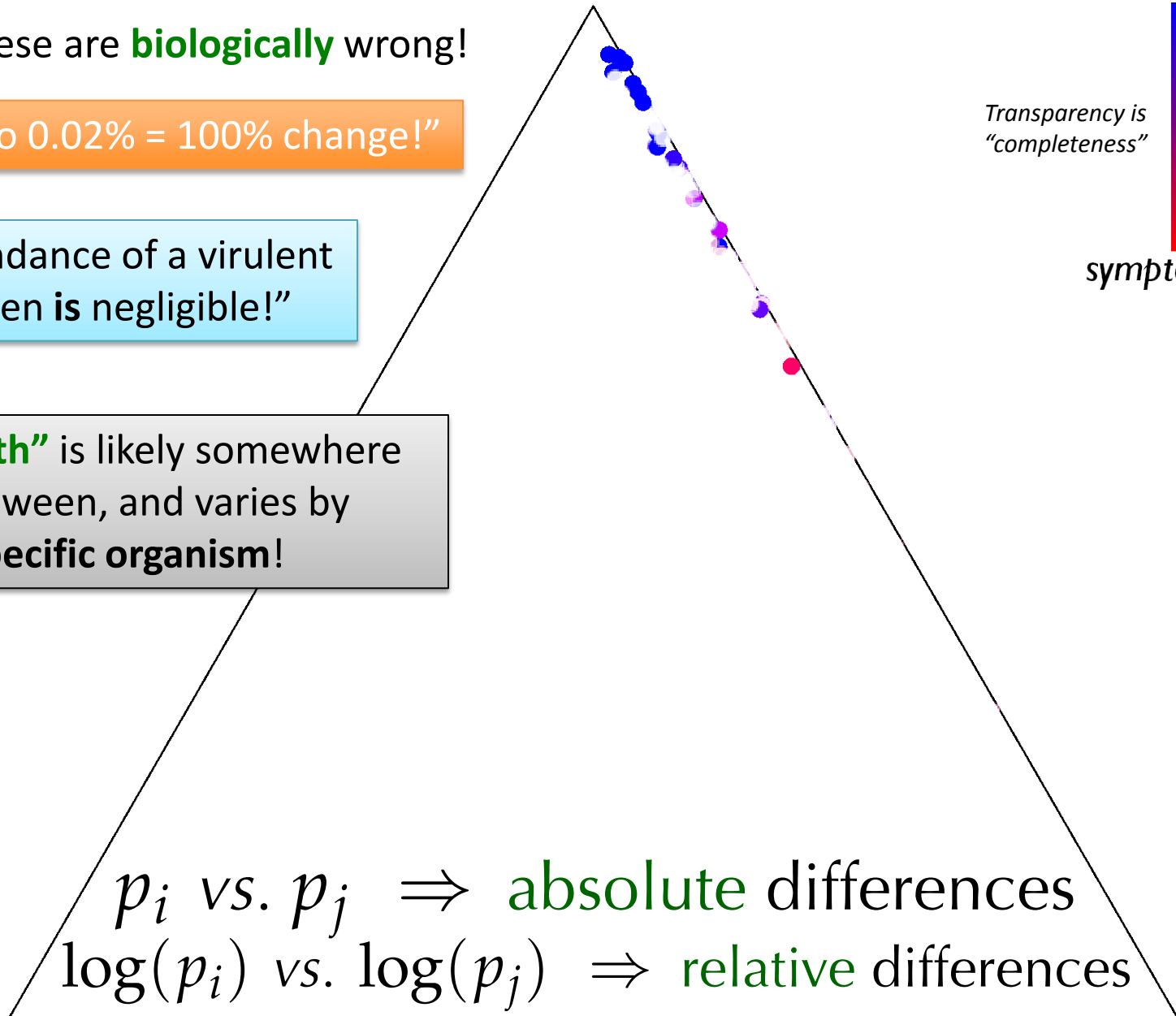The **"truth"** is likely somewhere in-between, and varies by **specific organism**!

Lactobacillus iners

*asymptomatic*

*Transparency is "completeness"*

*symptomatic*

$p_i$ *vs.* $p_j$ $\Rightarrow$ absolute differences
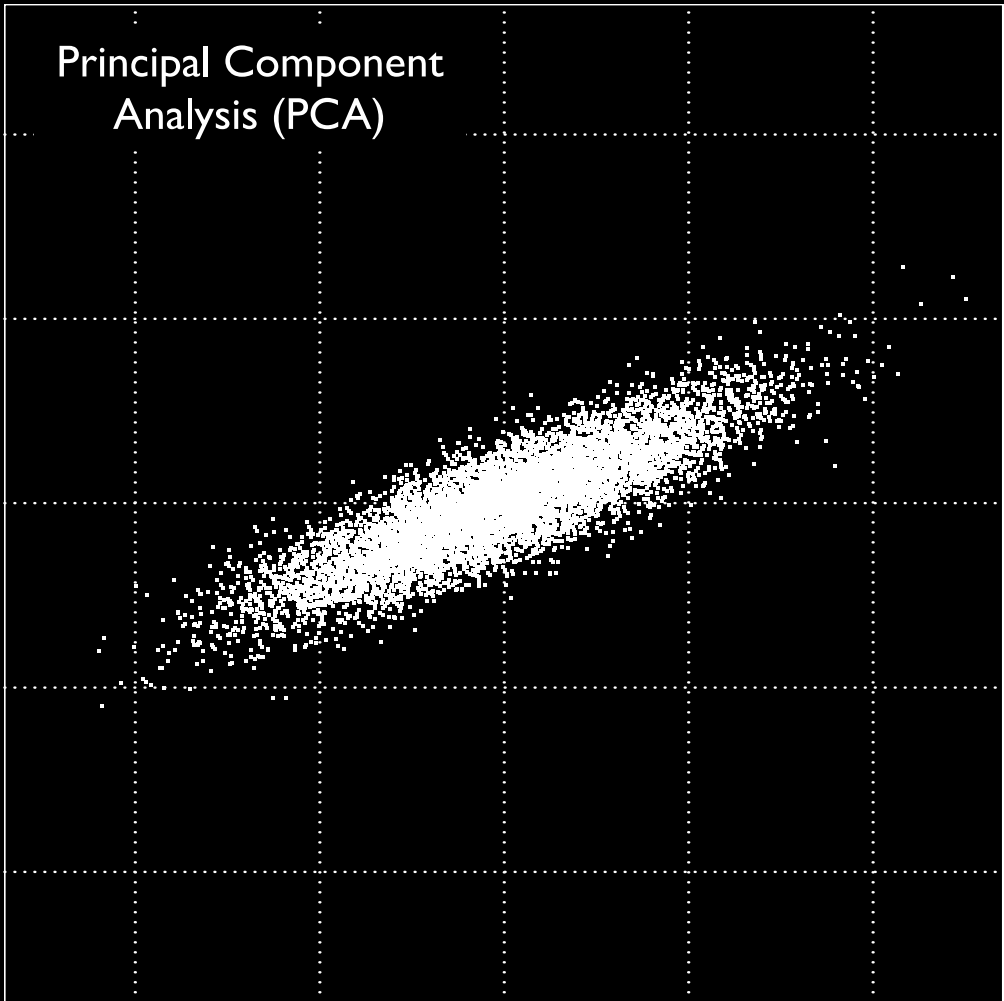$\log(p_i)$ *vs.* $\log(p_j)$ $\Rightarrow$ relative differences

Prevotella bivia

Gardnerella vaginalis

# Principal Component Analysis

Why use $\log(p_i)$ ?

Principal Component
Analysis (PCA)
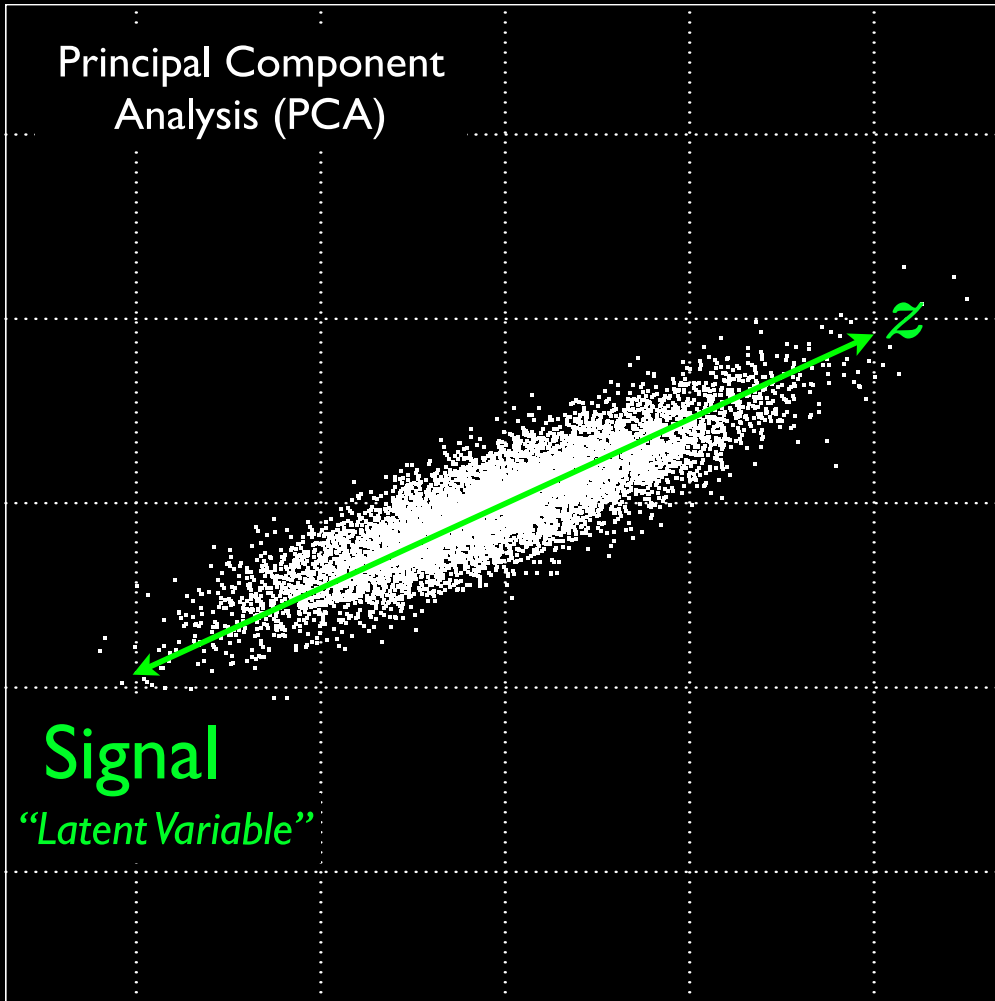
$x_2$

$x_1$

Principal Component
Analysis (PCA)

$x_2$

$z$

Signal

$x_1$

Principal Component Analysis (PCA)

$x_2$

$z$
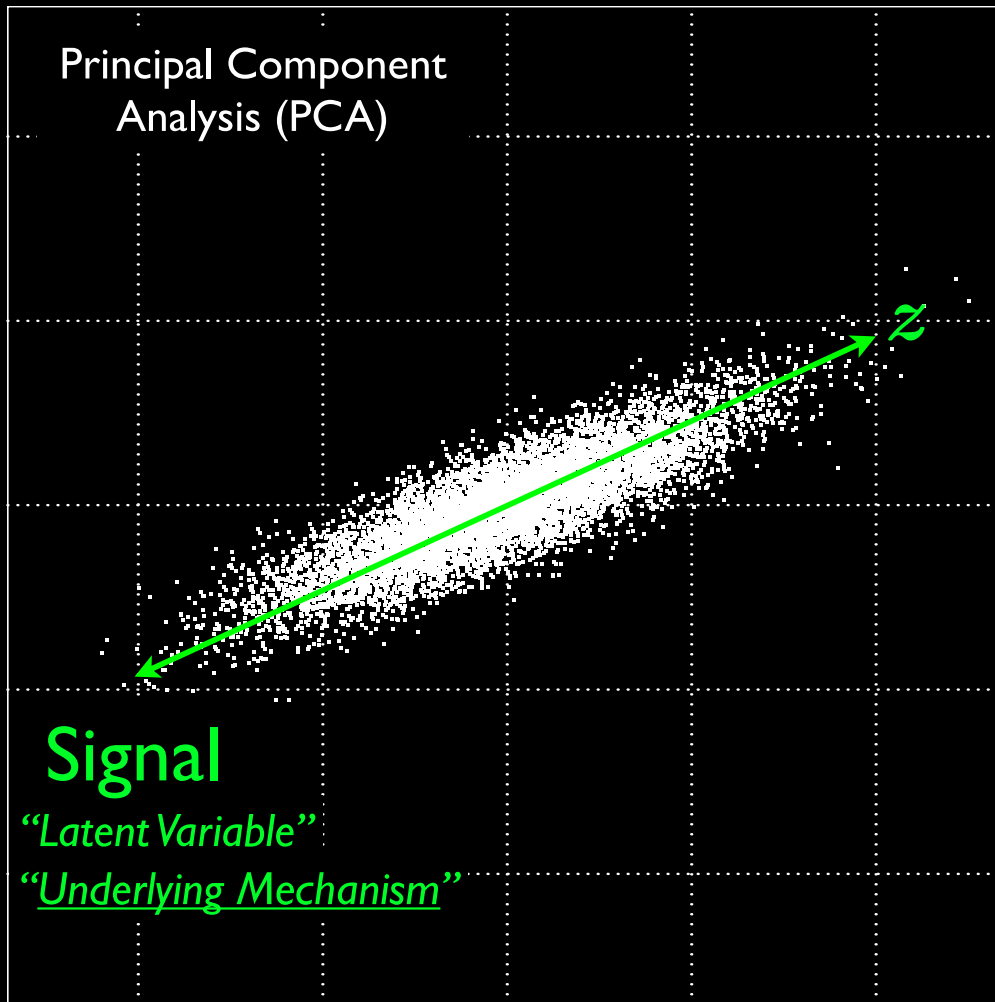
Signal
*"Latent Variable"*

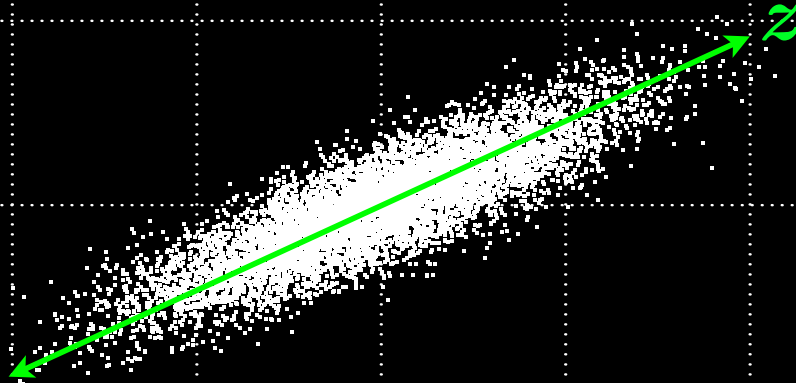$x_1$

Principal Component
Analysis (PCA)

$x_2$

$z$

Signal
*"Latent Variable"*
*"Underlying Mechanism"*
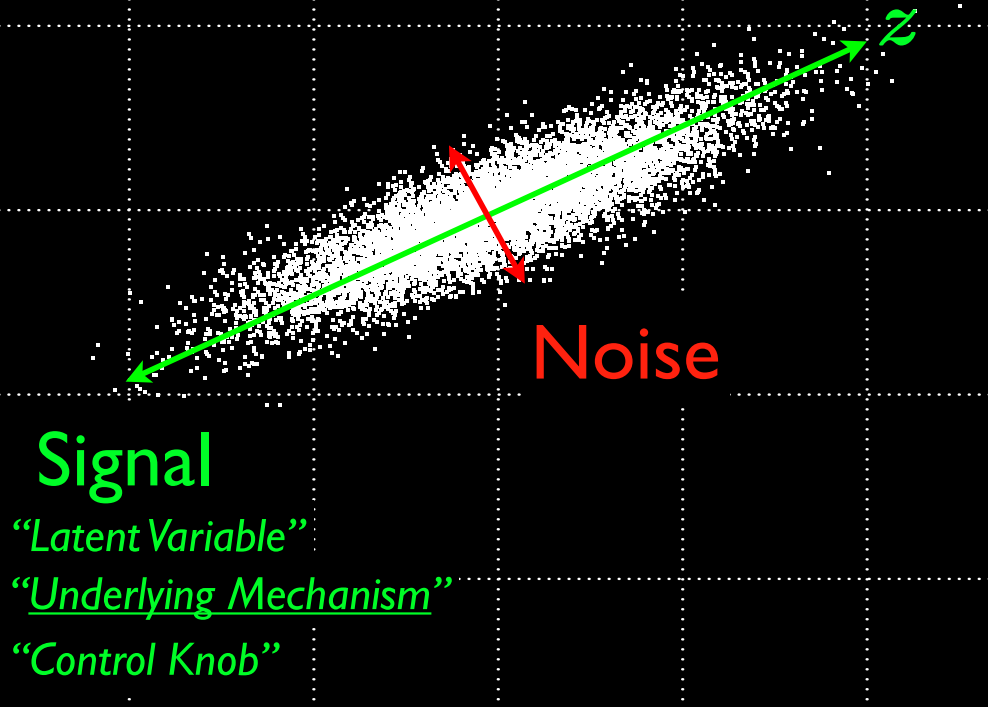
$x_1$

Principal Component
Analysis (PCA)

$z$

Signal
*"Latent Variable"*
*"Underlying Mechanism"*
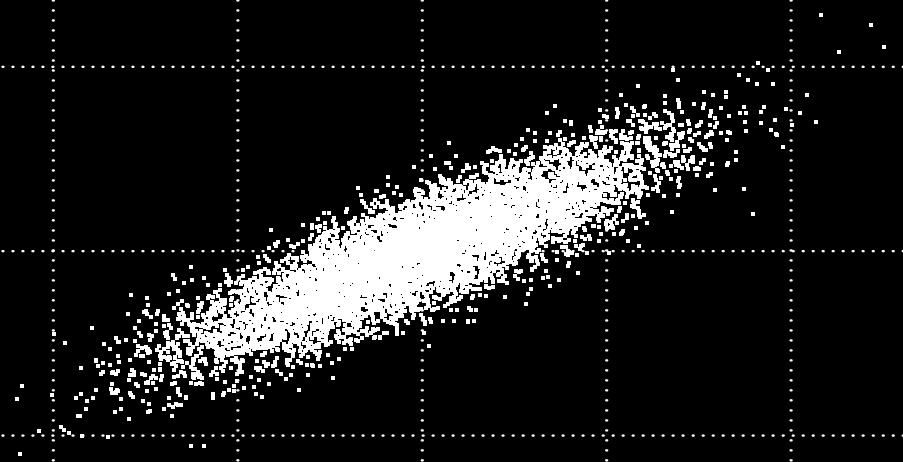*"Control Knob"*

Principal Component Analysis (PCA)

$z$

Noise

Signal
"Latent Variable"
"Underlying Mechanism"
"Control Knob"

onent
Analysis (PCA)

"Control Knob"

Component
Analysis (PCA)

*"Control Knob"*

Principal Component
Analysis (PCA)

$z$

Noise

Signal

*"Latent Variable"*

*"Underlying Mechanism"*

*"Control Knob"*

Principal Component Analysis (PCA)

The PCA variance decomposition has **strange** biological implications:

① Every fraction has the **same** proportion of stochastic variation

Noise

Signal

*"Latent Variable"*

*"Underlying Mechanism"*
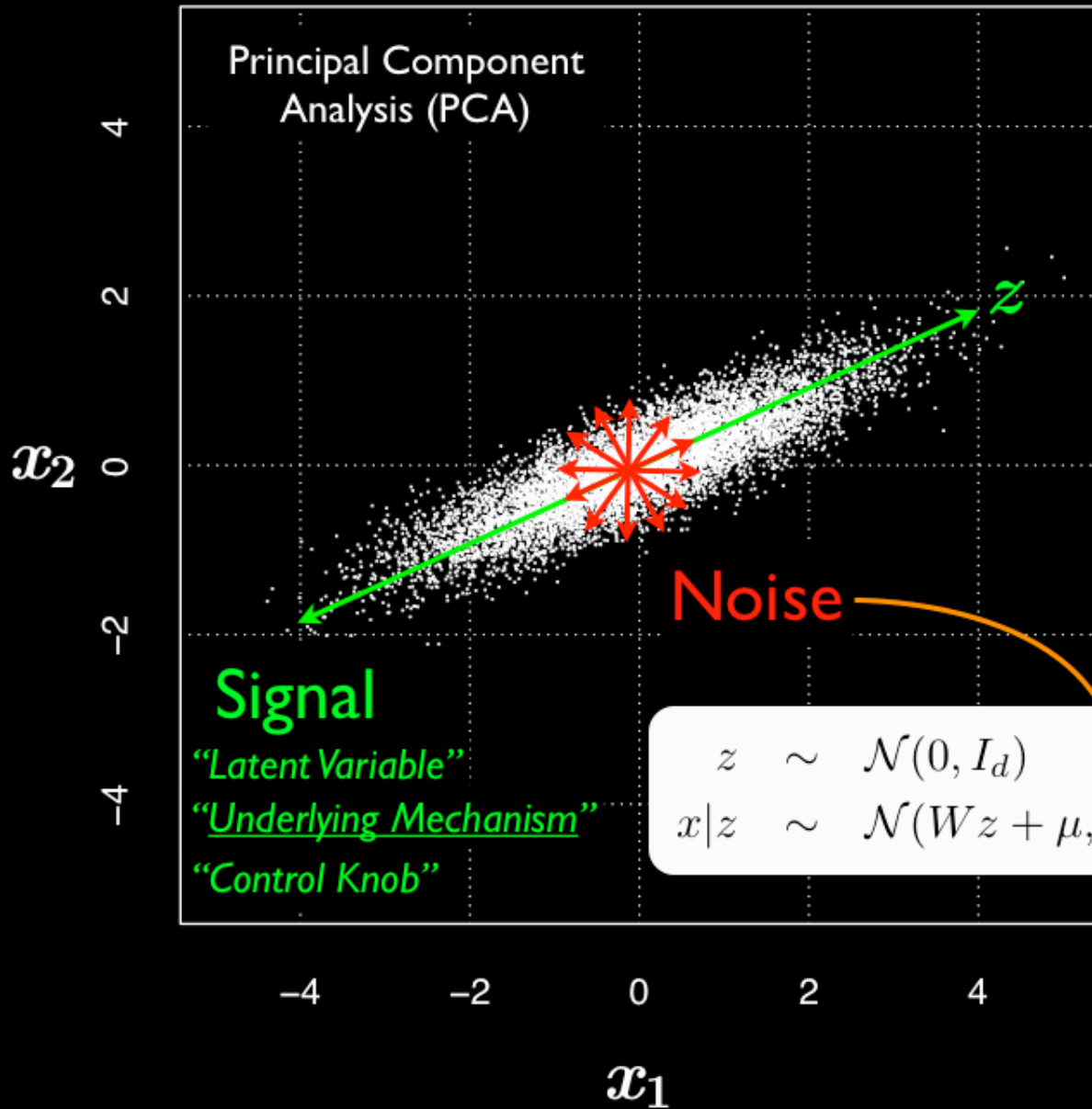
*"Control Knob"*

$$z \quad \sim \quad \mathcal{N}(0, I_d)$$
$$x|z \quad \sim \quad \mathcal{N}(Wz + \mu, \sigma^2 I_m), \quad \sigma > 0, \quad W \in \mathbb{R}^{md}$$

*Tipping and Bishop (1999)*

Principal Component Analysis (PCA)

$x_2$

$x_1$

$z$

Noise

Signal

*"Latent Variable"*

*"Underlying Mechanism"*

*"Control Knob"*

$$
\begin{aligned}
z &\sim \mathcal{N}(0, I_d) \\
x|z &\sim \mathcal{N}(Wz + \mu, \sigma^2 I_m), \quad \sigma > 0, \quad W \in \mathbb{R}^{md}
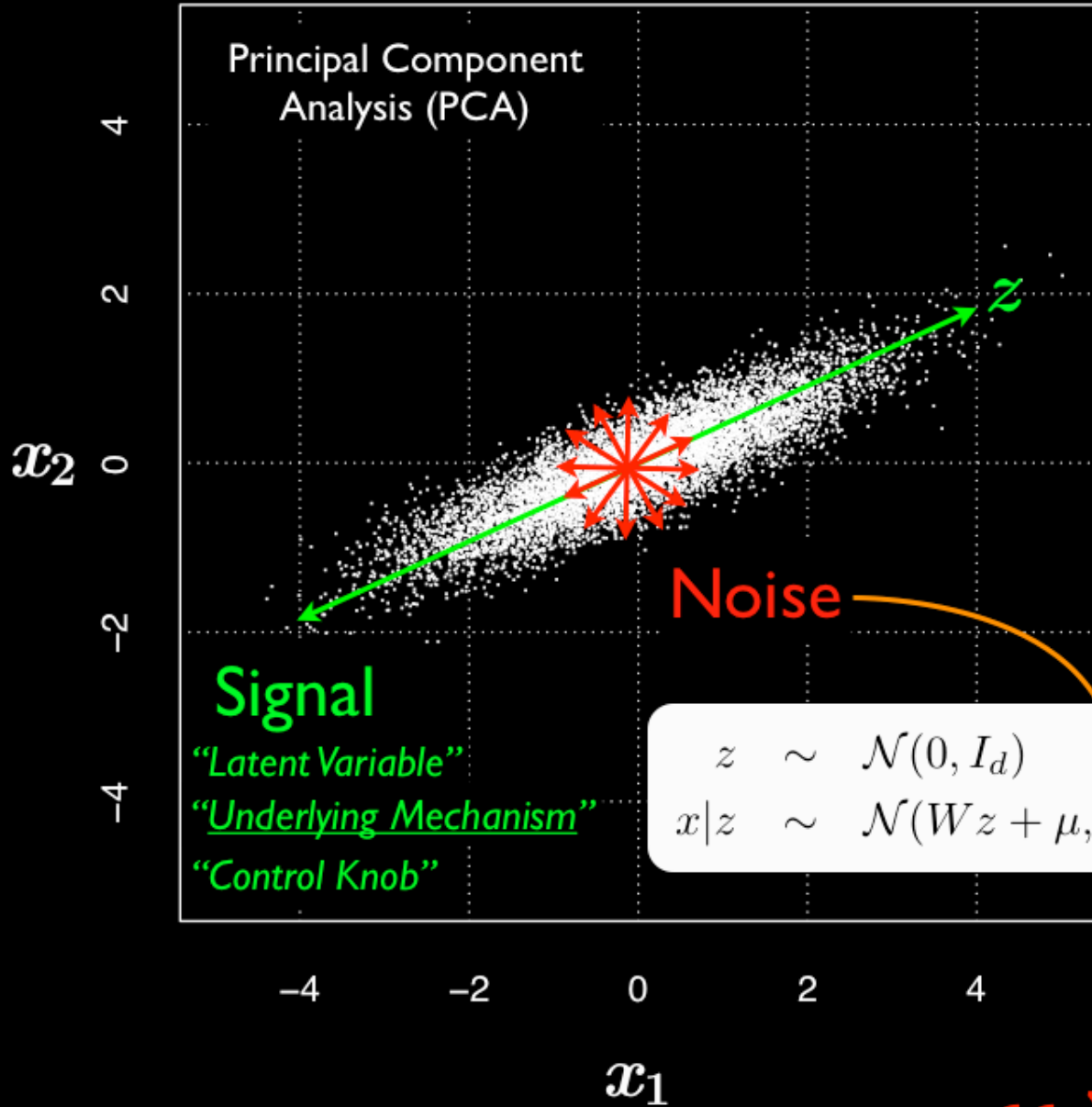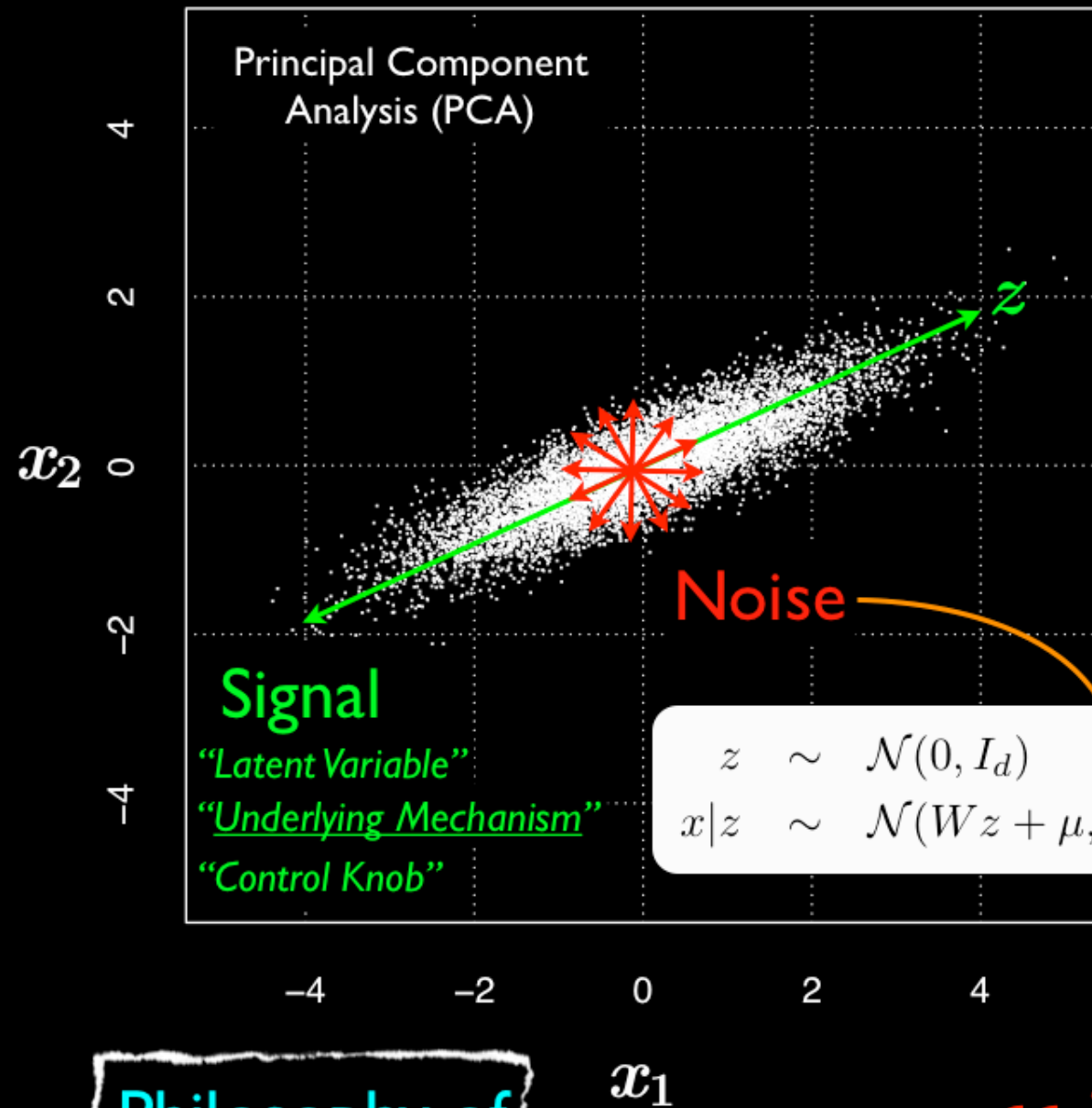\end{aligned}
$$

The PCA variance decomposition has **strange** biological implications:

① Every fraction has the **same** proportion of stochastic variation

② Results depend on the <u>units</u> of measurement

*Tipping and Bishop (1999)*

Principal Component Analysis (PCA)

$x_2$

Signal

*"Latent Variable"*
*"Underlying Mechanism"*
*"Control Knob"*

Noise

$x_1$

$$z \quad \sim \quad \mathcal{N}(0, I_d)$$
$$x|z \quad \sim \quad \mathcal{N}(Wz + \mu, \sigma^2 I_m), \quad \sigma > 0, \quad W \in \mathbb{R}^{md}$$

The PCA variance decomposition has **strange** biological implications:

① Every fraction has the **same** proportion of stochastic variation

② Results depend on the underlying of measurement

*Tipping and Bishop (1999)*

" $\sum$ or $\rho$ " ?

Principal Component Analysis (PCA)

$x_2$

$z$

Noise

Signal

"Latent Variable"
"Underlying Mechanism"
"Control Knob"

The PCA variance decomposition has **strange** biological implications:

① Every fraction has the **same** proportion of stochastic variation

② Results depend on the units of measurement

$$z \sim \mathcal{N}(0, I_d)$$
$$x|z \sim \mathcal{N}(Wz + \mu, \sigma^2 I_m), \quad \sigma > 0, \quad W \in \mathbb{R}^{md}$$
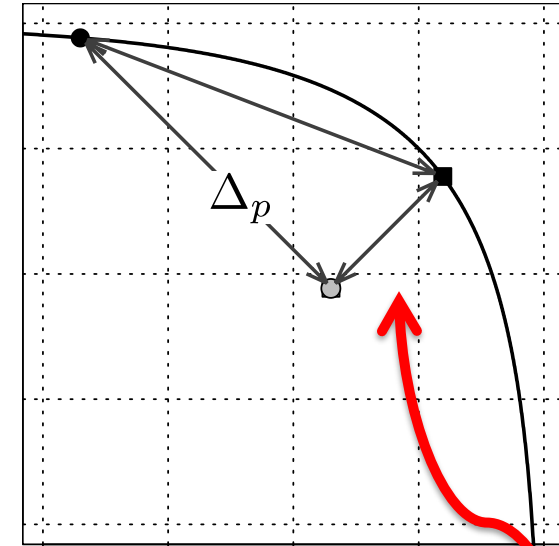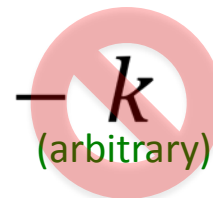
*Tipping and Bishop (1999)*

$x_1$

Philosophy of Modelling

"$\sum$ or $\rho$" ?

# PCA also introduces **<u>systematic distortions</u>** in the analysis!

$$\begin{bmatrix} p_i \\ q_i \\ r_i \end{bmatrix} \Rightarrow \begin{bmatrix} p_i/(p_i+q_i+r_i) \\ q_i/(p_i+q_i+r_i) \\ r_i/(p_i+q_i+r_i) \end{bmatrix}$$

$$\begin{bmatrix} \log(p_i/(p_i+q_i+r_i)) \\ \log(q_i/(p_i+q_i+r_i)) \\ \log(r_i/(p_i+q_i+r_i)) \end{bmatrix} \Rightarrow \begin{bmatrix} \log(p_i)-\log(p_i+q_i+r_i) \\ \log(q_i)-\log(p_i+q_i+r_i) \\ \log(r_i)-\log(p_i+q_i+r_i) \end{bmatrix}$$

$$\begin{bmatrix} \log(p_i)-\log(p_i+q_i+r_i) \\ \log(q_i)-\log(p_i+q_i+r_i) \\ \log(r_i)-\log(p_i+q_i+r_i) \end{bmatrix} \Rightarrow \begin{bmatrix} \log(p_i) \\ \log(q_i) \\ \log(r_i) \end{bmatrix} - \cancel{k}$$

(arbitrary)



$\Delta_p$

Distortions are
**routinely between
5% to 50% to 500%**
of datum distances!

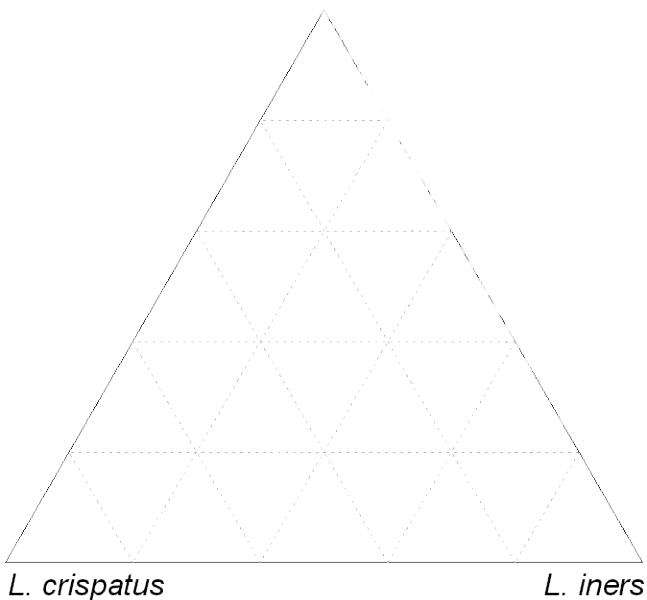# … taking into account …

- Zero-counts (not zero-proportions)
- Not-really-absolute and not-really-relative
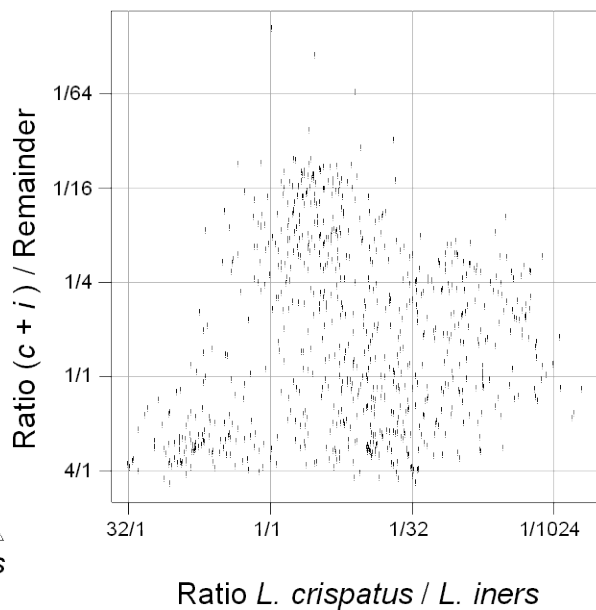- Removal of systematic distortions
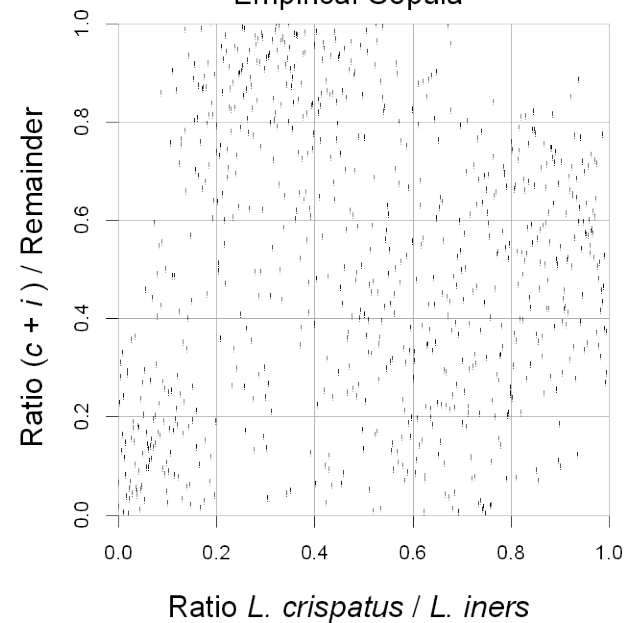- **Reparameterization invariance**

Co-Exclusion

?

L. iners
L. crispatus
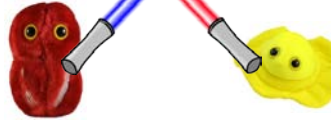
Remainder

L. crispatus          L. iners

Log Base 2

Ratio (*c* + *i*) / Remainder

1/64
1/16
1/4
1/1
4/1

32/1          1/1          1/32          1/1024

Ratio *L. crispatus* / *L. iners*

Empirical Copula

Ratio (*c* + *i*) / Remainder

1.0
0.8
0.6
0.4
0.2
0.0

0.0   0.2   0.4   0.6   0.8   1.0

Ratio *L. crispatus* / *L. iners*

No!

Co-Exclusion

?

L. crispatus
A. vaginae

**Remainder**

A. vaginae          L. crispatus

**Log Base 2**

Ratio ($v + c$) / Remainder

1/256
1/64
1/16
1/4
1/1
4/1

32/1          1/1          1/32          1/1024

Ratio *A. vaginae* / *L. crispatus*

**Empirical copula**

Ratio ($v + c$) / Remainder

1.0
0.8
0.6
0.4
0.2
0.0

0.0    0.2    0.4    0.6    0.8    1.0

Ratio *A. vaginae* / *L. crispatus*

## Yes!    (weakly)

## Quantifying "Uniformity"

Lactobacillus iners

Gardnerella vaginalis

*asymptomatic*

*Transparency is "completeness"*

*symptomatic*

0.00 weeks

Prevotella bivia

**Thank You!**

Leptotrichia amnionii