



ROYAL INSTITUTE  
OF TECHNOLOGY

# Statistical genetics and direct coupling analysis

Erik Aurell

Unifying the epidemiology and evolutionary  
dynamics of pathogens

Stockholm May 29 – June 23, 2023

# Outline

*What is direct coupling analysis (DCA)?*

Physicists' jargon for what in statistics would be called inference in exponential families.

*What do I mean by statistical genetics (in this talk)?*

Mainly the phase of quasi-linkage equilibrium (QLE), at high rate of recombination.

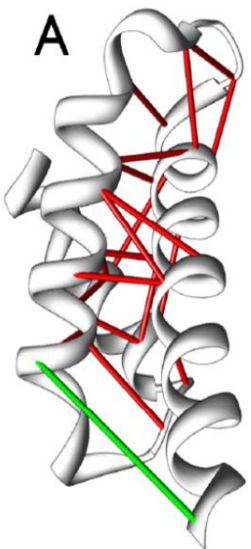
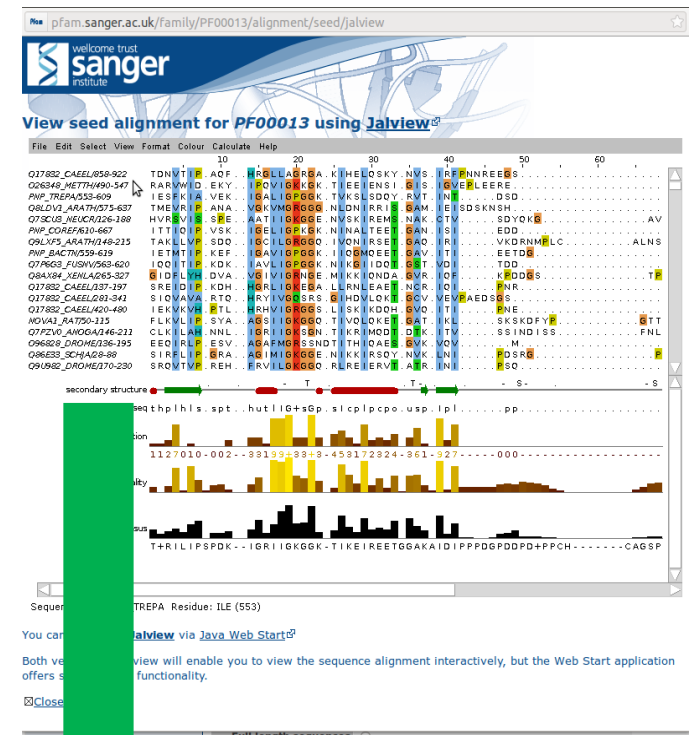
*What is the connection?*

You will see.

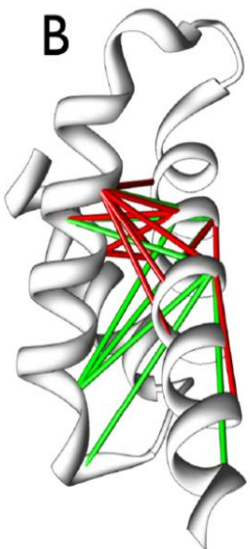
*References?*

Dichio, Zeng, Aurell, *Rep. Prog. Phys.* **86**  
052601 (2023) [and the original literature]

# DCA – flagship appl.



top DI pairs  
 direct coupling analysis (DCA)



top MI pairs  
 ranking by correlations

Weigt *et al*, PNAS 2009  
 ⋮  
 Morcos *et al*, PNAS 2011  
 ⋮  
 many others  
 ⋮  
[bit.ly/3Mr8351](http://bit.ly/3Mr8351)  
 ⋮  
 (courtesy S. Ovchinnikov lab)

$$P(\mathbf{x}) = \frac{1}{Z(h, J)} \exp \left( \sum_i h_i(x_i) + \sum_{ij} J_{ij}(x_i, x_j) \right)$$

“Considerable progress has recently been made by leveraging genetic information. It is possible to infer which amino acid residues are in contact by analysing covariation in homologous sequences, which aids in the prediction of protein structures”

Andrew W. Senior *et al Nature* **577**:706-710 (2020) [abstract]

LLLD**G**SSSLPESYFDMMKSF**A**KA**F**ISKANIGPHLTQVSVL**Q**YGSINTID  
 LLLD**G**SSSLPASYFEEMKSF**A**KA**F**ISKANIGPHHTQVSVL**Q**YGSITTID  
 LLLD**G**SSGFPASEFD**E**MKSF**A**KA**F**ISKANIGPQLTQVSVL**Q**YGSITTID  
 FVLD**G**SSSVRAS**Q**FEEMK**T**FVK**A**FIKKVNIGVGATQVSVL**Q**YGW**R**NILE  
 VLLD**G**STNIME**P**QFEEMK**T**FVK**E**L**I**KKVDIGNNGT**Q**ISV**V**QY**G**K**T**NTLE  
 FILD**T**SSSVGKDN**F**EK**I**RK**W**VADLVD**S**FDVSPDKTRVAV**V**L**Y**SDRPT**I**E  
 LAVD**T**S**Q**SME**I**QDLTVIKSVVDD**F**ISH**R**K-N---DRIGL**I**L**F**GT**Q**AY**L**Q  
 FLVD**T**SGSL**Q**KNGFDDEK**V**FVNSLLSHIRVSYK**S**T**Y**VSV**V**L**F**GT**S**ATID  
 LALD**T**SAT**T**GETILDHITRGA**Q**IGLAALS---DRSKVGV**V**L**Y**GEDHR**V**V  
 YVID**T**SGSMHGAK**I**EQ**T**RESMVA**I**LQDLH---EEDHF**G**IL**L**FERKIS**Y**W  
 FLID**T**SRSLGLRAY**Q**KE**L**Q**F**VERVLE**G**YE**I**GTNRTKVAV**I**T**F**SAGSR**L**E  
 ILLD**T**SSSIKIN**N**FDLIRK**F**VAN**I**IN**Q**F**E**VGRNGLMVG**M**AT**Y**S--RS**V**Q  
 FILD**T**SGSVGS**Y**N**F**EK**M**K**T**FVK**N**VVDF**F**NIGPK**G**THVAV**I**T**Y**ST**W**A--**Q**  
 FALD**T**ST**S**IG**S**Q**N**FEREK**Q**FVLA**F**VTD**M**DIGRSDV**Q**VSV**G**T**F**SDNARR**Y**  
 LLLD**T**SGSM**Q**GAA**I**EAL**L**SLK**D**EL-VK**N**S**I**AARR**V**E**I**A**I**V**T**FDSHIN**V**V  
 LLLD**T**SGSM**K**G**E**PLDAL**R**T**F**Q**Q**EL-DRDSLAK**R**VE**V**A**I**V**T**FNSD**V**E**I**V  
 LSVD**V**SL**S**MLARRLSAL**R**D**I**A**I**RFV**Q**K**R**K----NDRVGL**V**T**Y**S**G**EAL**R**  
 LAM**V**SGSM**Q**AN**R**LEAA**K**DVA**I**S**F**INN**R**NIG-----M**V**T**F**AGES**F**T**Q**  
 MSVD**V**SL**S**MLARR**L**TALK**N**IA**K**K**F**V**D**K**R**P----GDRIGL**V**T**Y**S**G**E**A**F**T**K  
 VLAD**V**SGSM**Q**G**E**P**I**AA-AA**F**TRY**L**-**Q**NE**V**-ASK**R**VE**V**AV**V**T**F**GT**V**AT**V**L

# DCA = model learning and parameter inference in biophysical applications

$$P(\mathbf{x}) = \frac{1}{Z(h, J)} \exp \left( \sum_i h_i(x_i) + \sum_{ij} J_{ij}(x_i, x_j) \right)$$

**Why not *correlation analysis* (which is a lot simpler)?**

Because DCA methods have empirically worked better, in particular for the flagship application of residue-residue contact prediction from tables of homologous protein sequences.

**Why not *maximum likelihood* (or *Bayesian estimates*)?**

Because a protein has maybe hundreds of amino acids. Inferring all these parameters from data using ML is slow.

**Why not just only use *deep learning*?** We'll get to that.

# 1<sup>st</sup> main method: elements of *inverse correlation matrix*

$$E(s) = \sum_i h_i S_i + \sum_{ij} J_{ij} S_i S_j \quad P^{\text{trial}}(s) = \prod_i P_i(S_i)$$

$$F^{nMF} = \sum_i H\left(\frac{1+m_i}{2}\right) + H\left(\frac{1-m_i}{2}\right) + \sum_i h_i m_i + \sum_{ij} J_{ij} m_i m_j \quad H(x) = -x \log x$$

$$\frac{\partial F^{nMF}}{\partial m_i} = 0 \quad \longrightarrow \quad m_i = \tanh\left(h_i^{nMF} + \sum_j J_{ij} m_j\right)$$

$$\chi_{ij} = \frac{\partial m_i}{\partial h_j} = c_{ij} \quad \text{An exact result, but used in } nMF \text{ approximation.}$$

$$\left(\chi^{nMF}\right)_{ij}^{-1} = \frac{\partial h_i^{nMF}}{\partial m_j} \approx \left(c^{-1}\right)_{ij} \quad \longrightarrow \quad \left(c^{-1}\right)_{ij} \approx \frac{1}{1-m_i^2} 1_{ij} - J_{ij}$$

**mean-field DCA:** Morcos et al *PNAS* (2011) [M Weigt] + many later contributions theory in Kappen & Spanjers *Phys. Rev. E* (2001) and in Nguyen, Berg & Zecchina (2017)

# 2<sup>nd</sup> main method: pseudo-likelihood maximization

Maximum likelihood  $P(\mathbf{S}) = \frac{1}{Z(\mathbf{h}, \mathbf{J})} \exp\left(\sum_i h_i S_i + \sum_{ij} J_{ij} S_i S_j\right)$

$$\Pr(\mathbf{S}^{(1)}, \dots, \mathbf{S}^{(n)}; \mathbf{h}, \mathbf{J}) = P(\mathbf{S}^{(1)}; \mathbf{h}, \mathbf{J}) \cdots P(\mathbf{S}^{(n)}; \mathbf{h}, \mathbf{J})$$

$$\mathbf{h}^*, \mathbf{J}^* \in \arg \max \left[ \sum_{ij} h_i \frac{1}{n} \sum_{s=1}^n x_i^{(s)} + \sum_{ij} J_{ij} \frac{1}{n} \sum_{s=1}^n x_i^{(s)} x_j^{(s)} - \log Z(\mathbf{h}, \mathbf{J}) \right]$$

Pseudo-maximum likelihood (avoids computing Z):

$$P(S_r | S_{\setminus r}) = \exp\left(h_r S_r + \sum_l J_{rl} S_r S_l\right) / \sum_y \exp\left(h_r y + \sum_l J_{rl} y S_l\right)$$

$$h_r^{plm}, J_{rl}^{plm} \in \arg \max \left[ \sum_{ij} h_i \frac{1}{n} \sum_{s=1}^n x_i^{(s)} + \sum_{ij} J_{ij} \frac{1}{n} \sum_{s=1}^n x_i^{(s)} x_j^{(s)} - f(h_r, J_{rl}, S_{\setminus r}) \right]$$

Julian Besag, *The Statistician* (1975); **plmDCA**, Ekeberg et al *Phys. Rev. E* (2013); **GREMLIN**, Kamisetty et al *PNAS* (2014); **CCMpred**, Seemayer et al *Bioinformatics* (2014)

# Flagship is now history



Google / DeepMind / AlphaFold

Andrew W. Senior *et al* "Improved protein structure prediction using potentials from deep learning", *Nature* **577**:706-710 (2020)



# Why DCA today?

You may not (yet) have a large number of labelled examples on which to train a more complex AI method. **Examples:** RNA, protein-protein interactions, fitness landscapes....

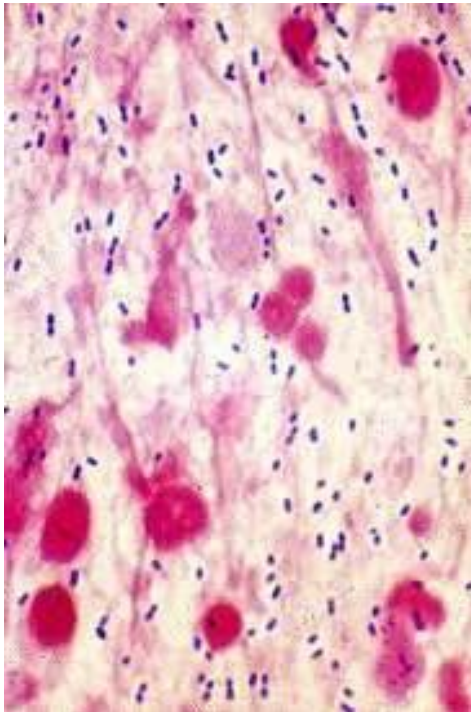
Your model might be too big for deep learning. **Example:** genome scale models.

You are actually using DCA as a family of inference methods to test something else. **This is the case today.**

But before I'll go there I'll give another motivational example.

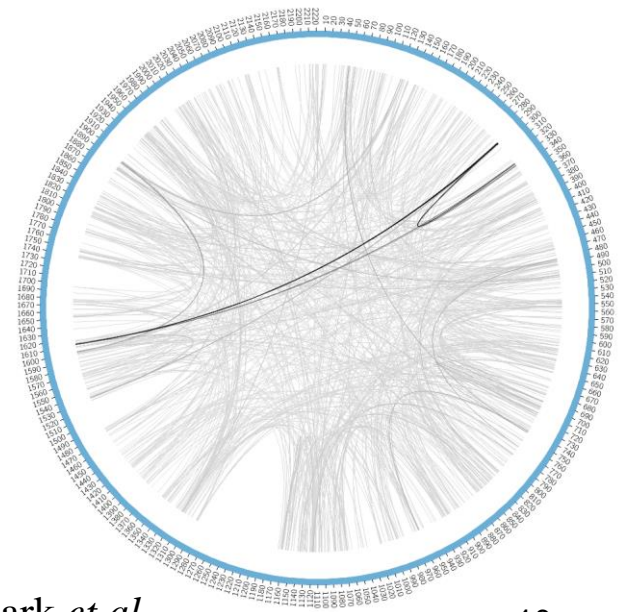
# Epistasis inference

ROYAL INSTITUTE OF TECHNOLOGY



The “Maela” data set: ~3,000 genomes of *Streptococcus pneumoniae*, retrieved from samples in the Maela refugee camp (Thailand / Myanmar).

The data had about 100,000 loci of variability, out of a genome 2.1Mbp (w/ some threshold).

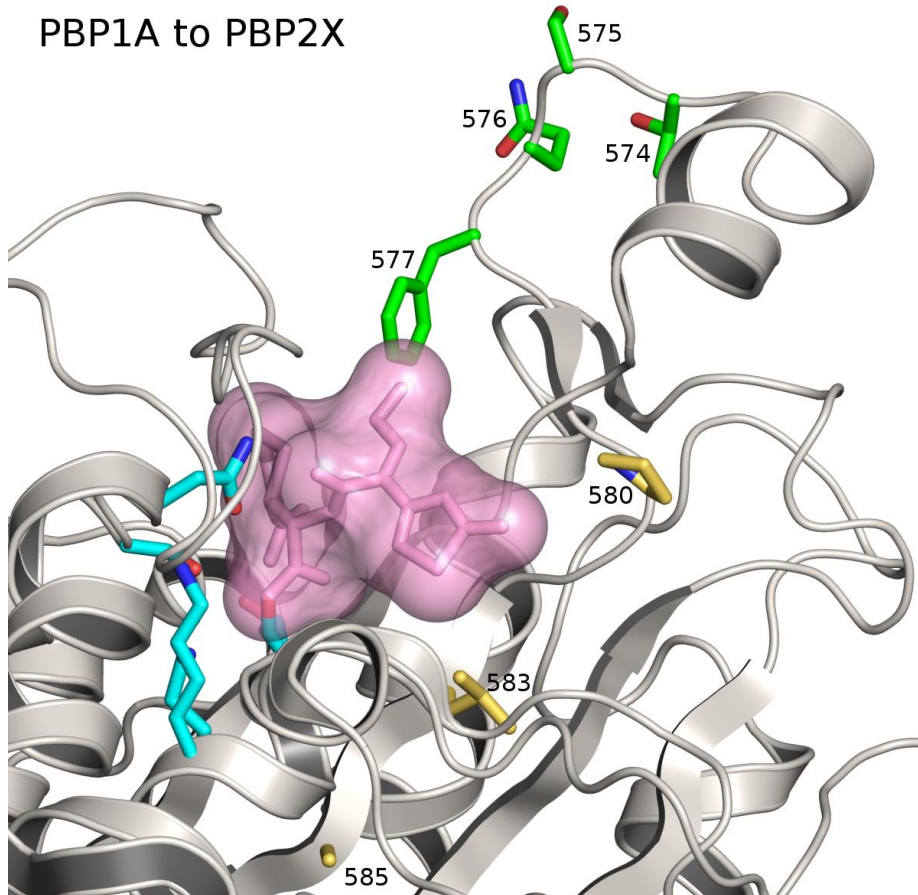


Nordita

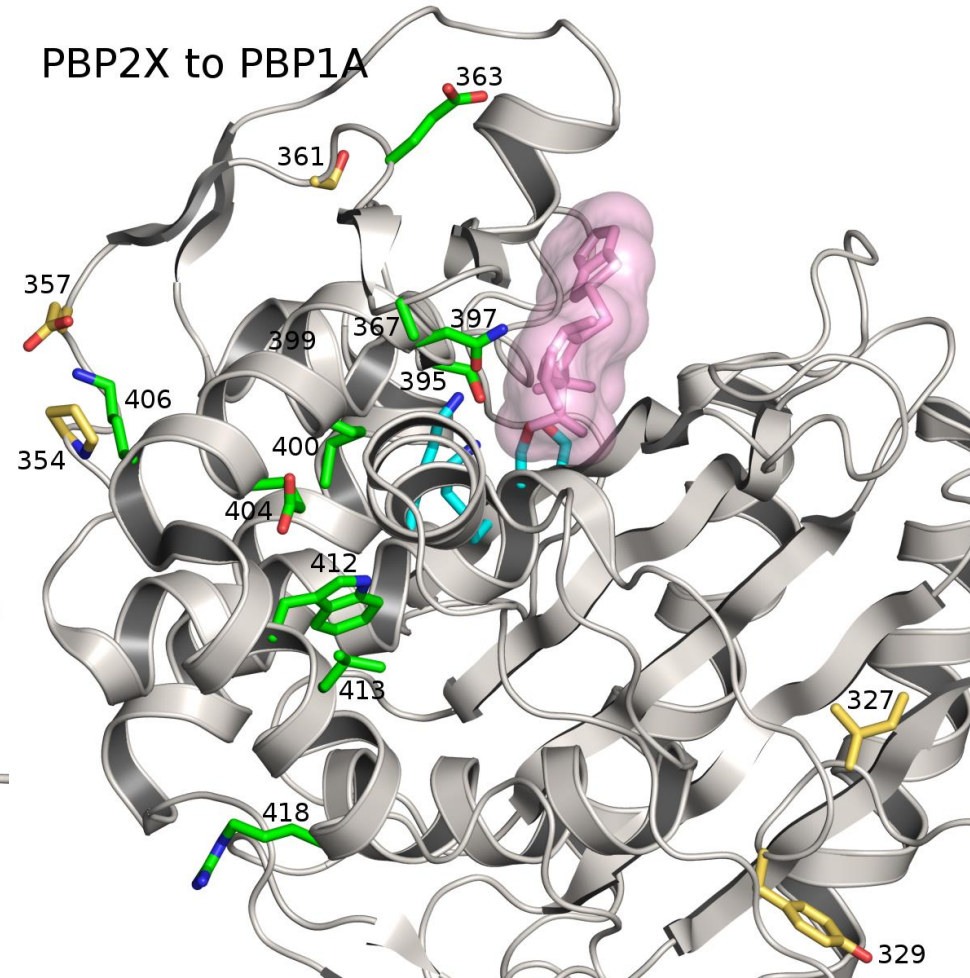
Skwark *et al*  
*PLoS Genetics* (2017)

# Epistatically coupled loci in proteins in the PBP family

PBP1A to PBP2X



PBP2X to PBP1A



[purple]  $\beta$ -lactam; [cyan] active site; [green and yellow] groups of predictions

# Some DCA on genome scale in bacteria and viruses

M. Skwark *et al*, "Interacting networks of resistance, virulence and core machinery genes identified by genome-wide epistasis analysis" *PLoS Genetics* 2017  
(*Streptococcus pneumoniae*, "Maela" data set) (*Streptococcus pyogenes M1*)

B. Schubert, R. Maddamsetti, J. Nyman, M. R. Farhat & D. S. Marks, *Nature Microbiology* 2019 (*Neisseria gonorrhoeae*)

Cui *et al*. [Daniel Falush] *eLife* 2020 (*Vibrio parahaemolyticus*)

C. Chewapreecha *et al* [Jukka Corander], *Molecular Biology and Evolution* 2022  
(*Burkholderia pseudomallei*, not quite DCA but by a similar method)

L Boeck *et al* [Julian Parkhill & R. Andres Floto], *Nature Microbiology* (2022)  
(*Mycobacterium abscessus*)

H-L Zeng *et al* [EA] *PNAS* 2020 (*SARS-CoV-2*)

E Cresswell-Clay & V Periwal, *Mathematical biosciences* 2021 (*SARS-CoV-2*)

J Rodriguez-Rivas *et al* [Martin Weigt] *PNAS* 2022 (*SARS-CoV-2*)

# Statistical genetics

A general understanding of population genetics in analogy with statistical physics. This has a long history starting with Hardy and Weinberg (1908), and Fisher and Wright in the 1920ies and 1930ies.

The objective today is the phase of *quasi-linkage equilibrium* first found by Kimura (1965), and how to use DCA to check it.

The *forces of evolution* taken into account in the analysis are *selection, mutations, genetic drift* and *recombination* (sex).

Obviously there are other forces of evolution (migration, etc.). But one must simplify somehow. It is complicated enough as is.

# Definitions

*Linkage equilibrium*: the distributions of alleles over loci are independent. Happens when recombination mix up genomes.

*Linkage disequilibrium (LD)*: distributions at alleles are not independent. Can be due to fitness or inheritance (or both).

*Formal*: A population is said to be in a *quasi-linkage equilibrium* (QLE) phase if (1) multi-genome distributions factorize and (2) single-genome distributions lie in an exponential family with no higher terms than in the fitness function. Which for quadratic fitness means

$$P(x) = \frac{1}{Z(h, J)} \exp\left(\sum_i h_i(x_i) + \sum_{ij} J_{ij}(x_i, x_j)\right)$$

Kimura *Genetics* **52**:875–890 (1965)  
 Neher & Shraiman *PNAS* **106**:6866 (2009); *Rev Mod Phys* **83**:1283 (2011)  
 formal definition in Dichio, Zeng, EA (2023)

# The Kimura-Neher-Shraiman theory (Neher-Shraiman version)

The distribution of genotypes in a population changes according to **selection, mutation, genetic drift** (finite- $N$ ) and **recombination**.

$\mathbf{g} = (s_1, s_1, \dots, s_L)$   $s_r = \pm 1$  “Ising genome”

$$P(\mathbf{g}, t + \Delta t) = \frac{e^{\Delta t F(\mathbf{g})}}{\langle e^{\Delta t F(\mathbf{g})} \rangle} P(\mathbf{g}, t) \quad F(\mathbf{g}) = \sum f_i s_i + \sum f_{ij} s_i s_j \quad \text{Fitness}$$

$$P(\mathbf{g}, t + \Delta t) = P(\mathbf{g}, t) + \Delta t \mu \sum_i [P(M_i \mathbf{g}, t) - P(\mathbf{g}, t)] \quad \text{Mutations}$$

$$P(\mathbf{g}, t + \Delta t) = (1 - r\Delta t)P(\mathbf{g}, t) + \Delta t r \sum_{\mathbf{g}_m, \mathbf{g}_f} C(\mathbf{g}, \mathbf{g}_m, \mathbf{g}_f) P(\mathbf{g}_m, t) P(\mathbf{g}_f, t)$$

Two haploid parents copy themselves, produce a child, and the rest of both genomes is discarded. Directly appropriate for some yeasts. One can modify the above to also cover bacterial recombination.

# Neher-Shraiman theory of QLE

Neher & Shraiman, *Rev Mod Phys* **83**:1283 (2011)

[for Potts not Ising] Gao, Cecconi, Vulpiani, Zhou, EA, *Phys. Biol.* **16** 026002 (2019)

Recombination is parametrized by a cross-over indicator variable  $\xi$

$$g^{(i)} = \xi_i g_m^{(i)} + (1 - \xi_i) g_f^{(i)} \quad C(\mathbf{g}, \mathbf{g}_m, \mathbf{g}_f) = C(\xi)$$

Recombination acts on pairwise dependencies through

$$c_{ij} = \sum_{\xi} C(\xi) [\xi_i(1 - \xi_j) + \xi_j(1 - \xi_i)]$$

Assume that  $P(\mathbf{g})$  is initially close to a Gibbs distribution of an Ising energy function  $(h_i, J_{ij})$  and recombination rate  $r$  is large:

$$\partial_t P(\mathbf{g}, t) = \dots \Rightarrow \dot{J}_{ij} = f_{ij} - r c_{ij} J_{ij} \Rightarrow J_{ij} = \frac{f_{ij}}{r c_{ij}}$$

In steady-state QLE the Ising parameters  $J_{ij}$  are proportional to pairwise fitness parameters  $f_{ij}$ , the proportionality being  $(r c_{ij})^{-1}$ .



# Kimura-Neher-Shraiman eq. for the DCA terms in QLE

$$J_{ij} = \frac{F_{ij}}{r \cdot c_{ij}}$$

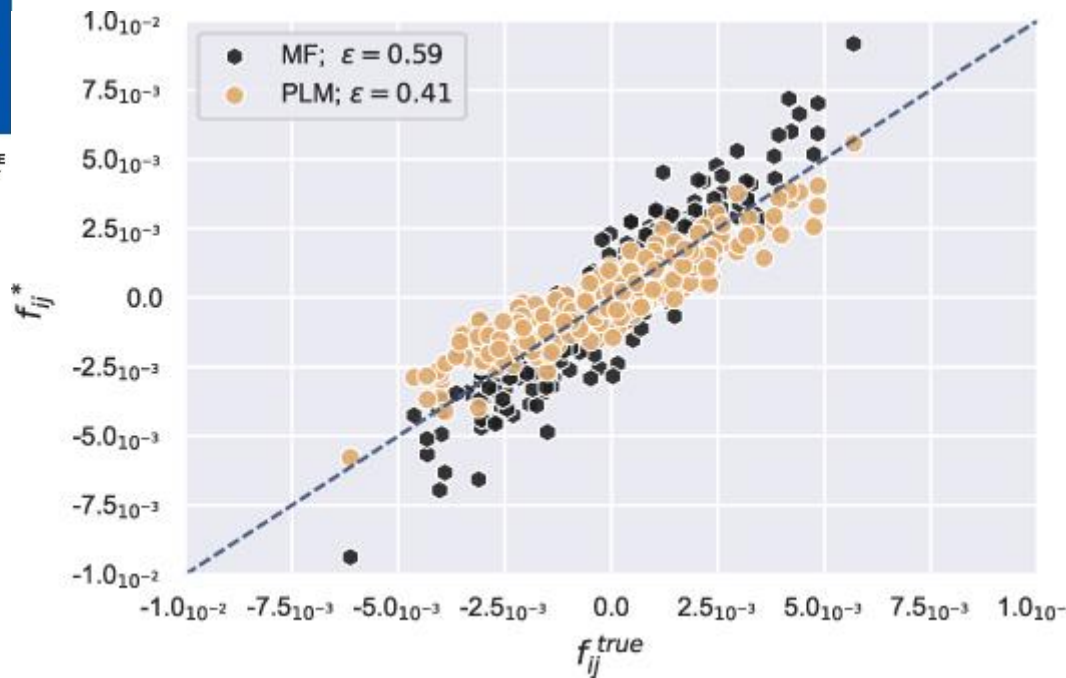
$r$  is an overall rate of recombination

$F_{ij}$  are the pairwise epistatic contribution to fitness

$J_{ij}$  are the Potts model parameters

$c_{ij}$  is the probability that alleles at loci  $i$  and  $j$   
are inherited from different parents

We can use DCA to test the above as a characteristic of QLE.



Example of a scatter plot for the reconstructed epistatic fitness components  $f_{ij}^*$  (y-axis) versus true underlying parameters  $f_{ij}^{\text{true}}$  (x-axis).

MF (mean-field) and PLM (pseudo-likelihood maximization) versions of DCA give similar reconstruction performance.

Simulation parameters of **FFPopsim** [Zanini and Neher *Bioinformatics* **28** 3332–3 (2012)]

	Value	Description
$N$	200	n. individuals
$L$	25	n. of loci
$T$	$2.5 \times 10^3$	n. of generations
$\omega$	0.5	crossover rate
$r$	[0.0:1.0]	rate of recombination
$\mu$	[0.005:0.1]	rate of mutation
$\sigma_e$	[0.001:0.02]	$f_{ij} \sim \mathcal{N}(0, \sigma_e)$

$$f_{ij}^* = r \cdot c_{ij} \cdot J_{ij}^* \quad c_{ij} \approx \frac{1}{2}$$

# Mauri-Zeng-Dichio-Aurell-Cocco-Monasson revised theor(ies)

Derived by a Gaussian closure on moments, but can also be done similarly to the Neher-Shraiman analysis. Several levels of inference formulae were found, out of which I will here only use the simplest (which NB bi-passes the need for DCA)

$$f_{ij}^* = \chi_{ij} \cdot \frac{4\mu + r c_{ij}}{(1-\chi_i^2)(1-\chi_j^2)} \quad \chi_i = \langle s_i \rangle \quad \chi_{ij} = \langle s_i s_j \rangle - \chi_i \chi_j$$

Note the presence of mutation rate  $\mu$ . The formula reduces to Kimura-Neher-Shraiman in the small-coupling regime and in the limit when  $\mu$  tends to zero.

Mauri, Cocco, Monasson, *Europhys Lett* **132** 56001 (2021)  
Zeng, Mauri, Dichio, Cocco, Monasson, *EA JSTAT* 2021 083501 (2021)

# KNS vs MZDACM

Regression of inferred epistasis ( $f_{ij}^*$ ) on underlying “true” epistasis ( $f_{ij}$ ).

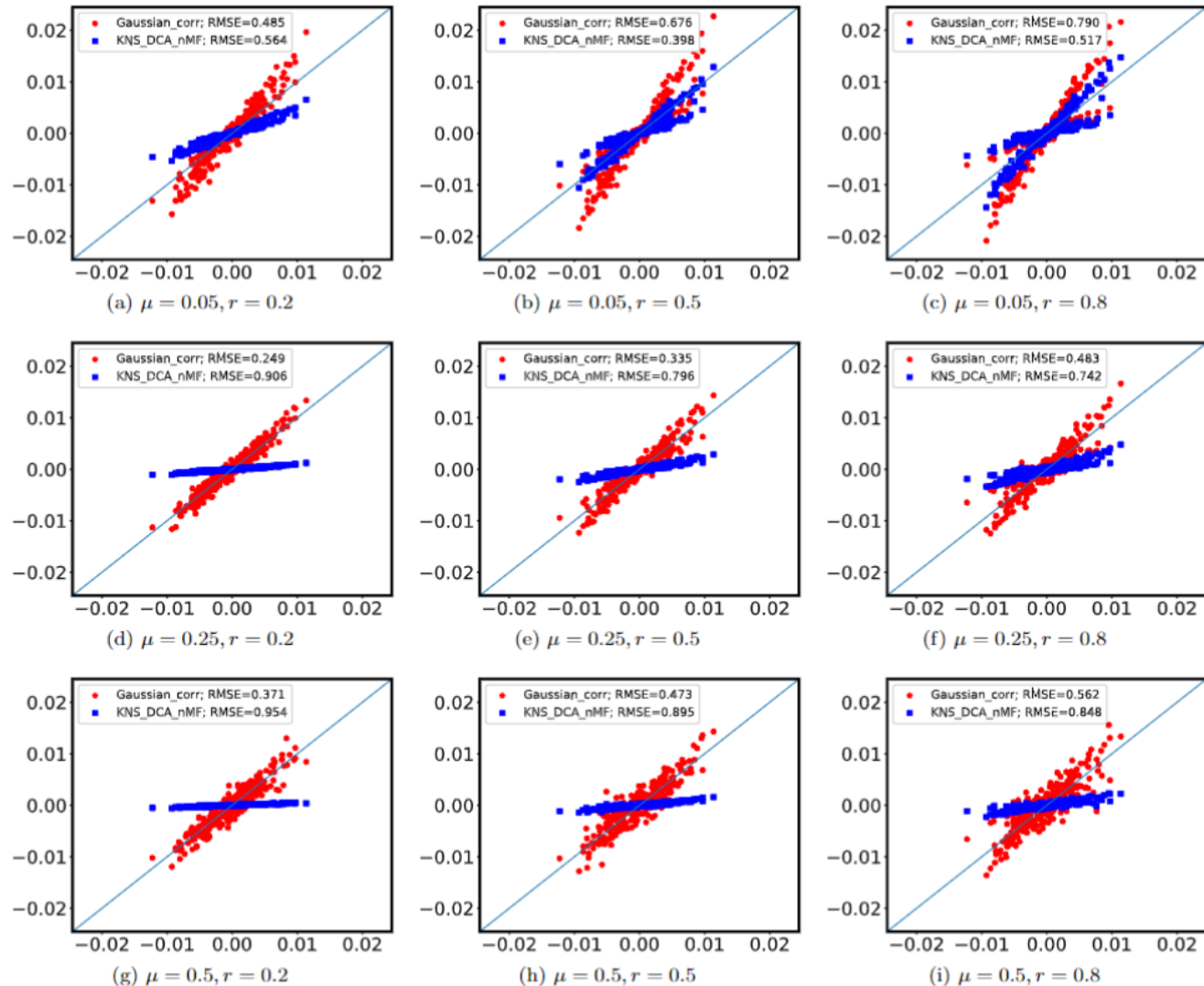
Comparison of the **KNS** formula:

$$f_{ij}^* = r \cdot c_{ij} \cdot J_{ij}^*$$

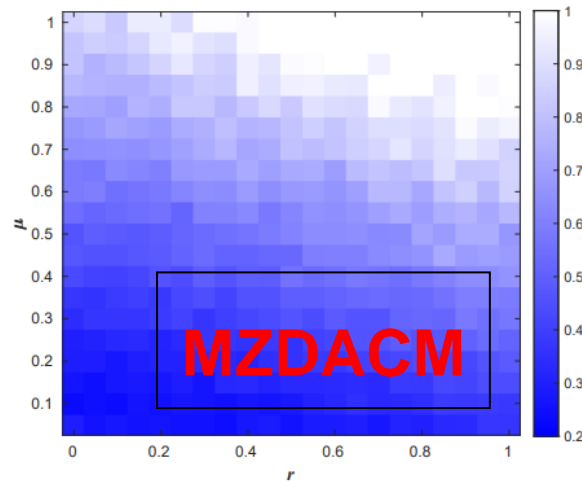
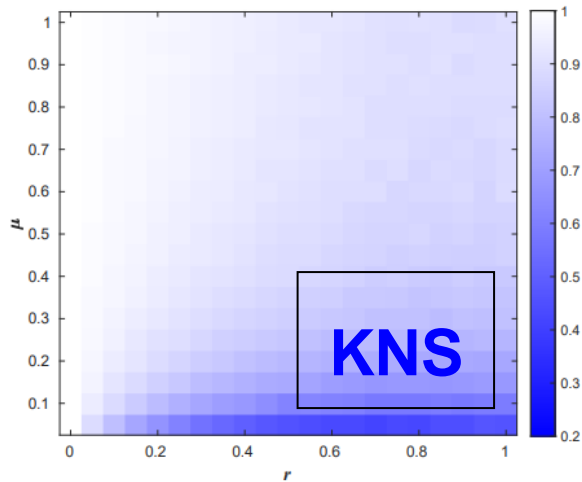
and the **MZDACM** formula;

$$f_{ij}^* = \frac{\chi_{ij} \cdot (4\mu + r c_{ij})}{(1 - \chi_i^2)(1 - \chi_j^2)}$$

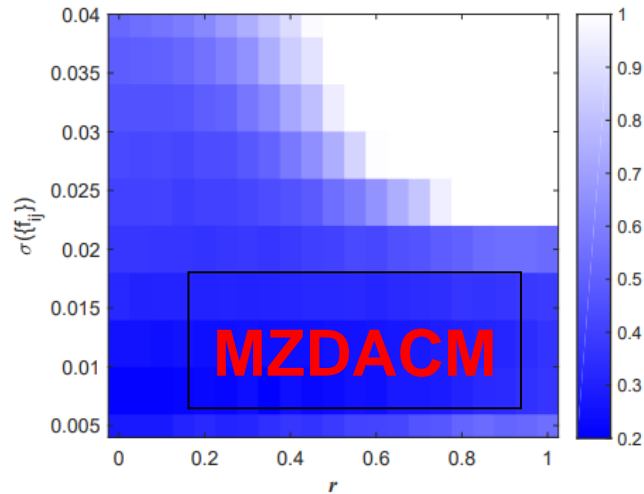
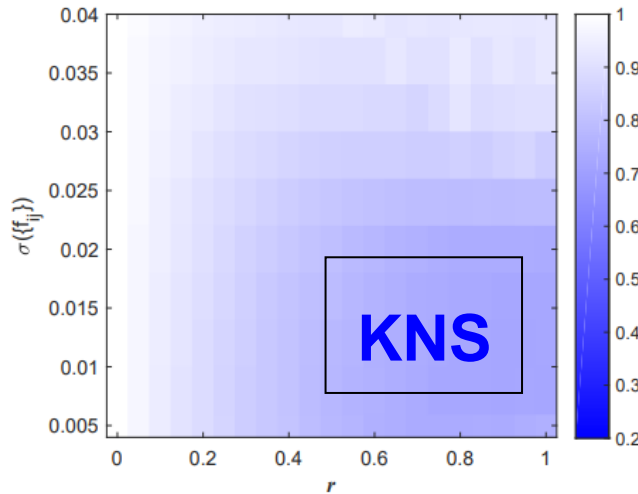
Zeng *et al* *JSTAT*  
083501 (2021)



# Performance phase diagrams



$\mu$  vs  $r$  at random additive fitness  $\sigma_a = 0.05$  and random epistatic fitness  $\sigma_e = 0.004$ . One realization for each parameter.



$\sigma_e$  vs  $r$  at mutation rate  $\mu = 0.2$ .

For other parameters, see paper.

Zeng *et al* JSTAT 083501 (2021)

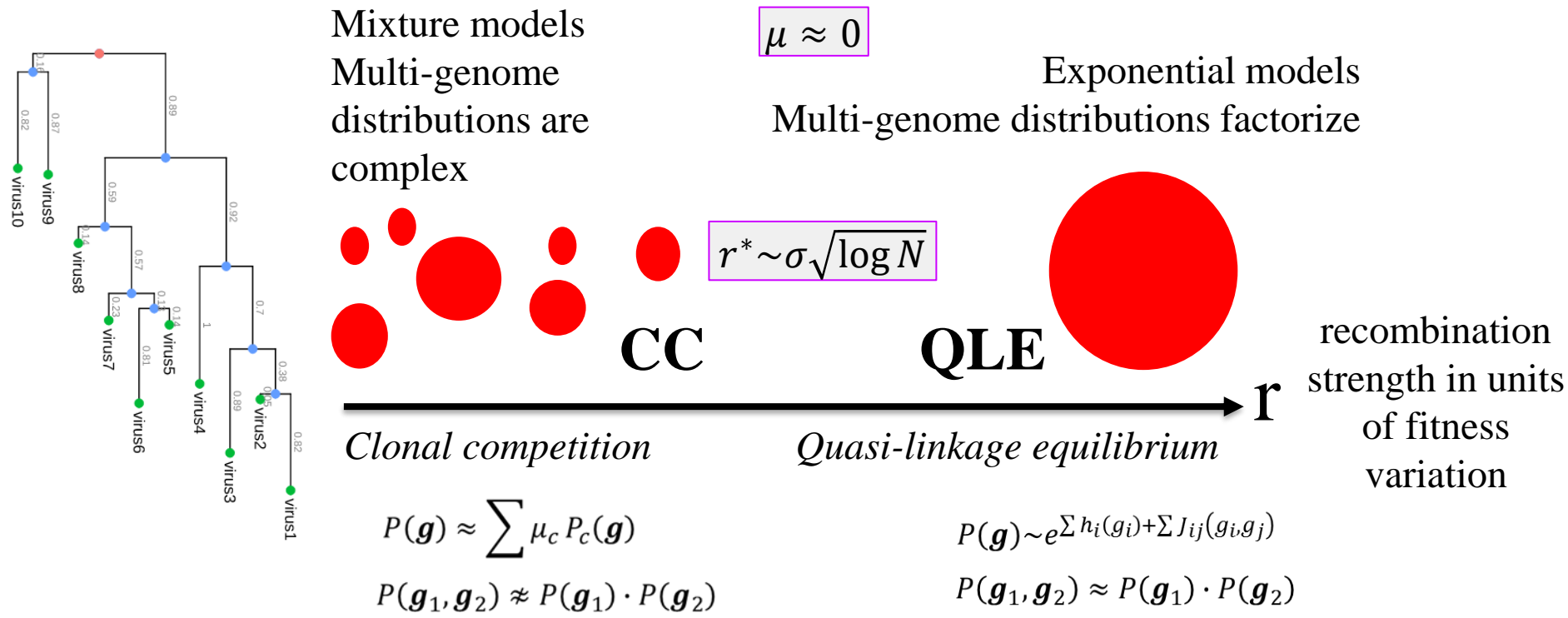
# Loss of QLE

*Rep. Prog. Phys.* **86** 052601 (2023) [arXiv:2105.01428]

and a brief review of earlier work

# QLE vs clonal competition

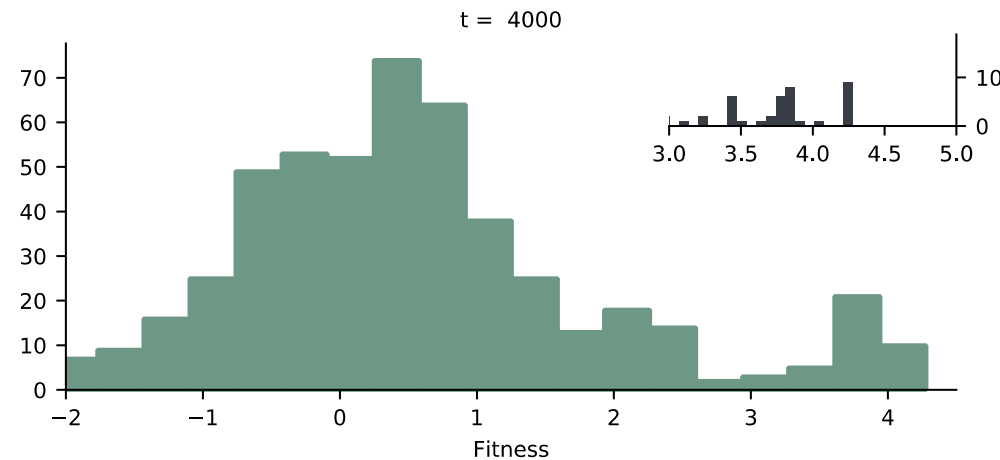
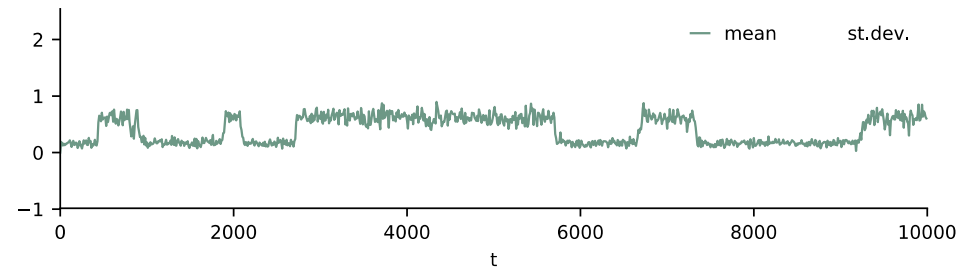
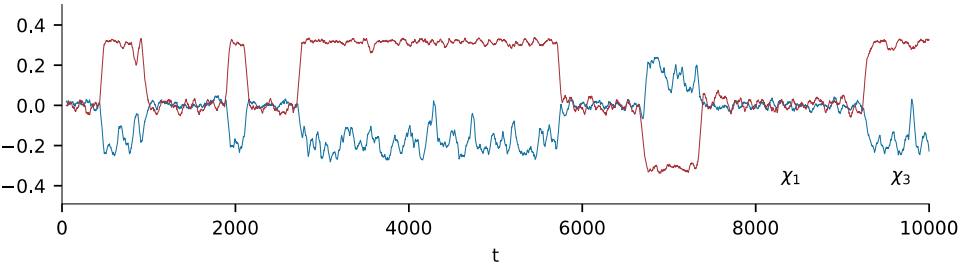
Neher & Shraiman *PNAS* **106**:6866 (2009); *Rev Mod Phys* **83**:1283 (2011);  
 Neher, Vucelja, Mézard, Shraiman *JSTAT* 01008 (2013)



At  $N = \infty$  there is no QLE! However,  $\sqrt{\log \mathcal{N}_{avo}} \approx 7,4\dots$

# Non-random coexistence

$$\mu \neq 0$$



At finite mutation rate the loss of QLE manifests itself differently. For finite populations appears an intermittent regime fluctuating between QLE and Non-Random Coexistence (NRC).

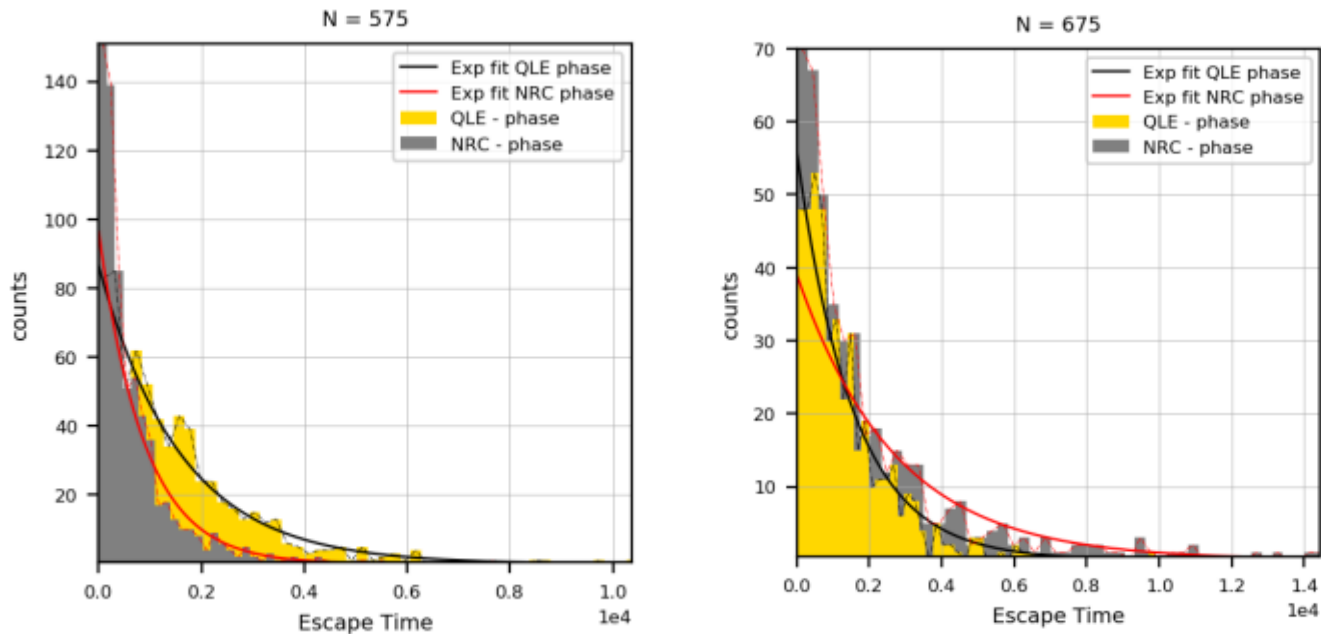
Total mean fitness in the population fluctuates, and is higher in NRC.

Snapshot of the fitness distribution at  $t = 4000$  in the above (NRC interval). Differently to QLE, the distribution is bimodal with a group of individuals at high fitness.

Similar to predictions in CC, though here no exact clones, due to mutations.



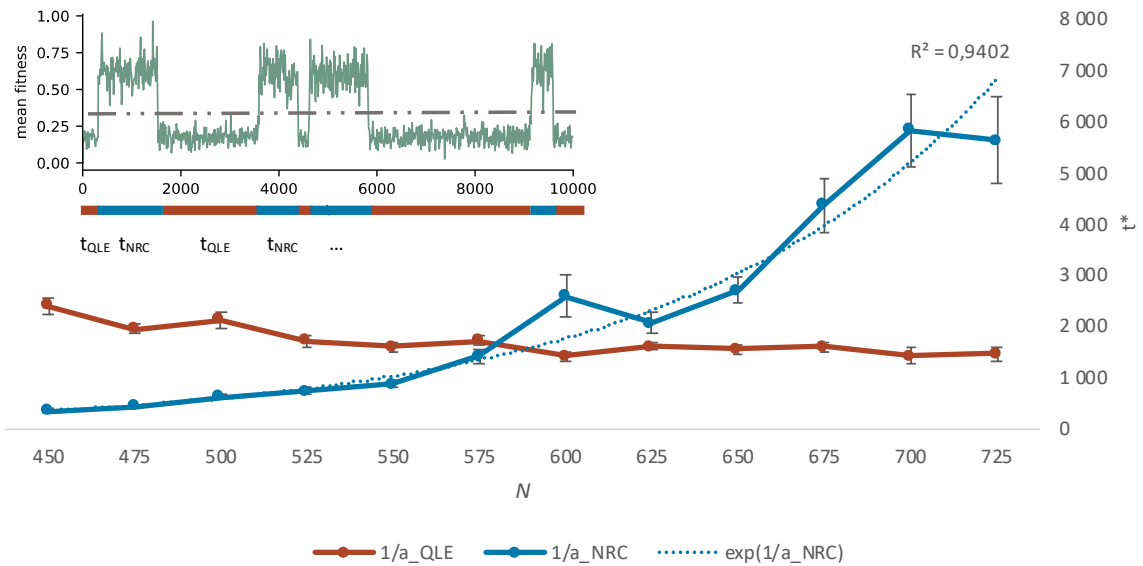
# Escape time distributions



Empirical distribution of escape times from respectively QLE and NRC. Simulations are run in a region of the parameter space (including  $N$ , here 575 and 675) where the systems dynamics visually jumps back and forth between QLE and NRC. Both distributions are well fitted as exponentials. The inverse rate is the mean escape time, in either direction. Other parameters:  $L = 25$ ,  $T = 1.5 \cdot 10^6$ ,  $\mu = r = \omega = 0.5$ ,  $\sigma_e = 0.029$ .

# Finite- $N$ dependence

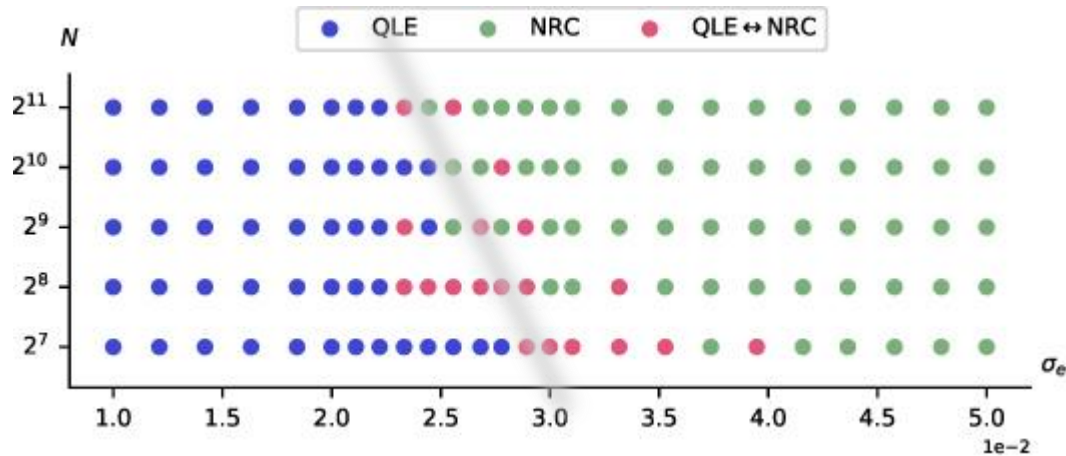
Estimated mean escape times from QLE and NRC.  
 Inset: The dynamics undergoes multiple transitions QLE  $\leftrightarrow$  NRC ( $T = 1.0 \cdot 10^4$ ).



The QLE  $\rightarrow$  NRC transition happens when an individual in a finite population finds a high-fitness state. Analogous to the biophysical problem of transcription factors finding a binding site. Expected waiting time  $N^{-1}$ .

The NRC  $\rightarrow$  QLE transition happens when a group of high-fitness individuals is lost from the population. Analogous to Muller ratchet. Expected waiting time exponential in  $N$ .

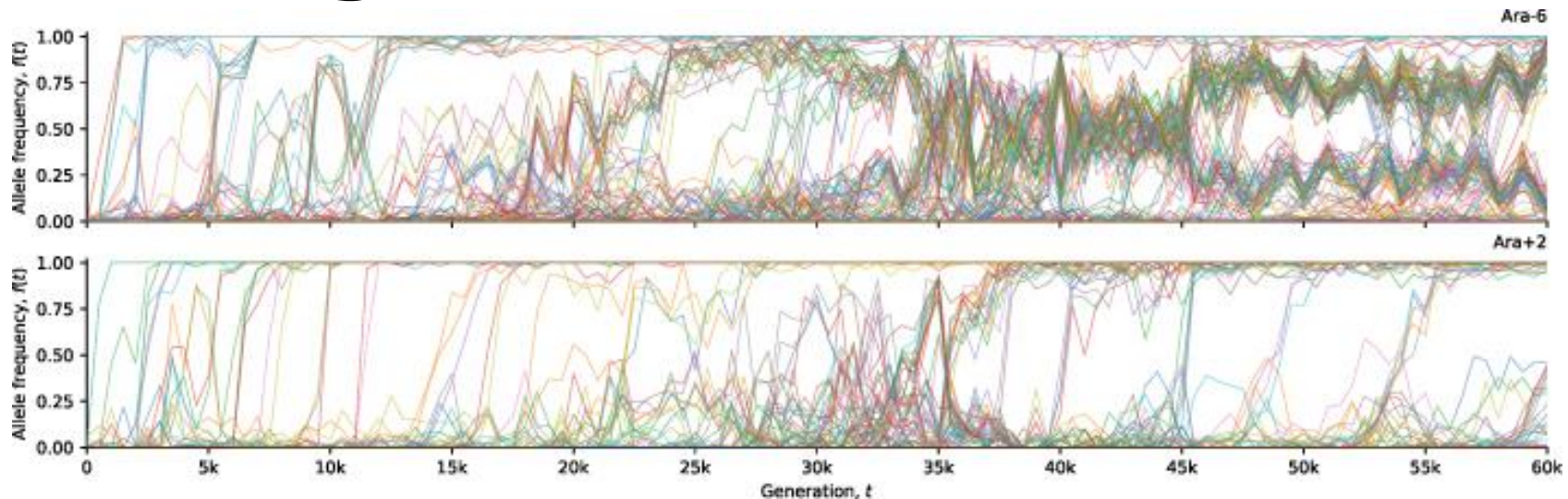
# “Phase diagram” in $(N, \sigma_e)$



A number of simulations are run for the same time ( $2.0 \cdot 10^4$ ). If the population remains in the QLE (NRC) the point is marked as **blue** (**green**). If at any point a transition QLE $\leftrightarrow$ NRC is observed, the corresponding point is marked as **red**.

The previous heuristic theory predicts that for high  $N$  we the population should *always* be in NRC (same as in the Clonal Competition loss channel). This seems to be in agreement with the simulations (provided there is at least one transition).

# Long-term evolution exps.



Allele frequency trajectories of all de novo mutations detected in 2 of the 12 LTEE populations, labelled respectively Ara-6 and Ara+2. Population Ara-6 (top row) shows quasi-stable coexistence of clades while Ara+2 (bottom row) shows mutations that fix rapidly. Quasi-stable coexistence was reported in 9 out of 12 LTEE populations [Good, McDonald, Barrick, Lenski, Desai 2017 *Nature* **551** 45–50 (2017)].

Figure previously unpublished, private communication from Profs B H Good and M M Desai, reproduced with permission.

# Outlook & loose ends

## SARS-CoV-2

H-L Zeng et al [EA], *PNAS* 2020 & *Phys Rev E* 2022

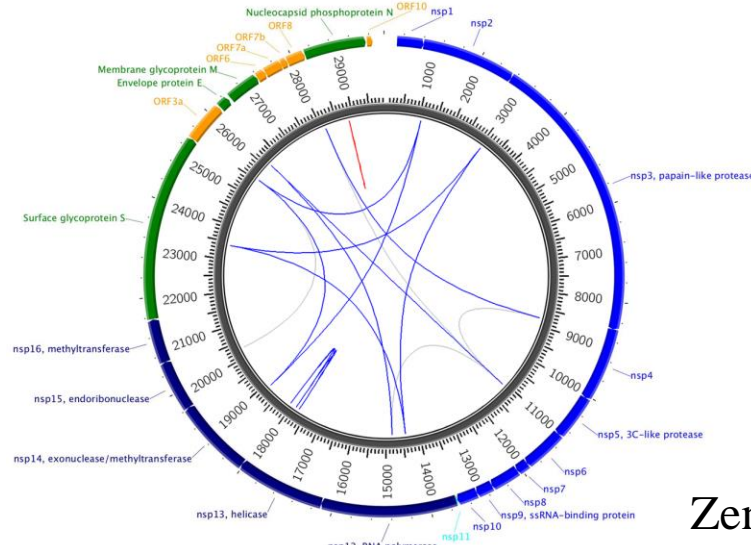
E Cresswell-Clay & V Periwal, *Mathematical biosciences* 2021

J Rodriguez-Rivas et al [Martin Weigt] *PNAS* 2022

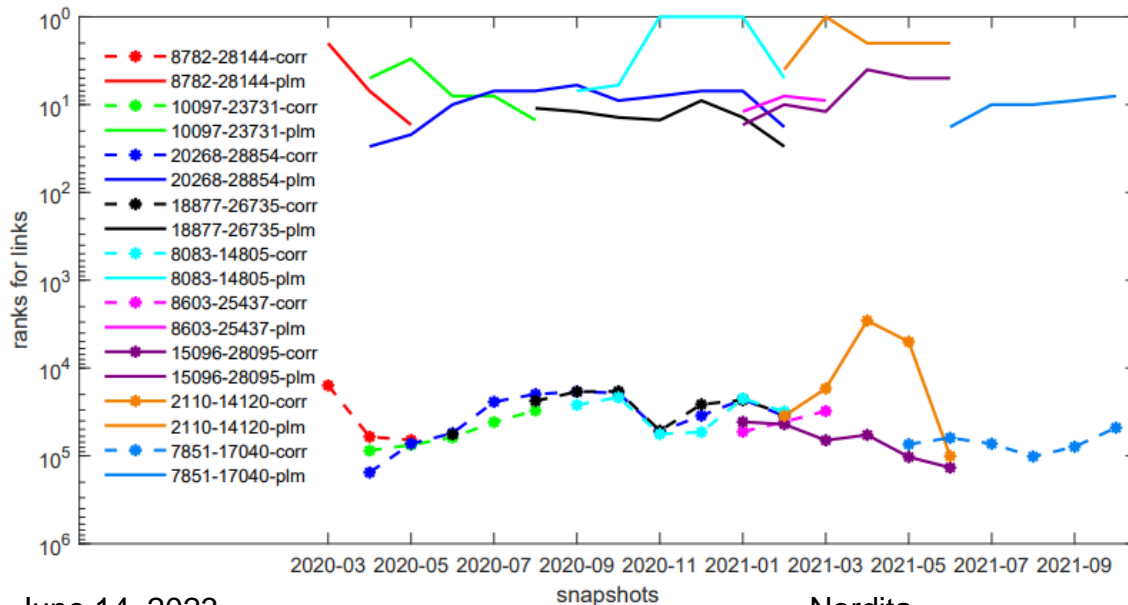
H-L Zeng & Y Liu (unpublished)

# SARS-CoV-2 genome-scale DCA

Initially we found an interesting set of predictions of epistatic interactions which were stable in GISAID data until August 2020. Many of these predictions were also found by Cresswell-Clay & Periwal, using a larger data set up to October 2020.



Zeng et al PNAS 2020 (Fig 1)



Later we found that predictions disappear when variability at one or both loci in a pair goes down to zero. Predictions are only stable on the time scales of months.

Phys Rev E 2022, Fig 5  
about 3,500,000 genomes

# Spike-spike

Zeng et al *PRE* 2022 Table I

August 2021					September 2021					October 2021				
rank	locus 1	AA-m.	locus 2	AA-m.	rank	locus 1	AA-m.	locus 2	AA-m.	rank	locus 1	AA-m.	locus 2	AA-m.
7	23284	D574D	25339	D1259D	7	23284	D574D	25339	D1259D	9	23284	D574D	25339	D1259D
16	21987	G142D	24410	D950N	15	21987	G142D	24410	D950N	11	21995	T145H	22227	A222V
67	22093	M177I	22104	G181V	45	21995	T145H	22227	A222V	15	21987	G142D	24410	D950N
70	22917	R452L	22995	K478T						135	21846	T95I	24208	I882I
71	22082	P174S	22093	M177I										
74	22081	Q173H	22093	M177I										
190	22082	P174S	22104	G181V										
195	22081	Q173H	22104	G181V										

TABLE I. Largest DCA terms with both terminals in Spike coding region, August-October 2021. Top-200 couplings computed as plmDCA scores are considered. For each of them in the three months displayed, there's the indication of the rank, the two loci involved and the corresponding amino acid (AA) mutations. Green color indicates that this mutation is found in delta variant. Red color indicates that this mutation is found in omicron variant. Couplings with one or both terminals colored green are attributed to a phylogenetic effect. The single pair with one terminal colored red is not attributed to a phylogenetic effect, the growth of omicron being later than October-2021.

# Spike-non-spike

Zeng et al *PRE* 2022 Table II

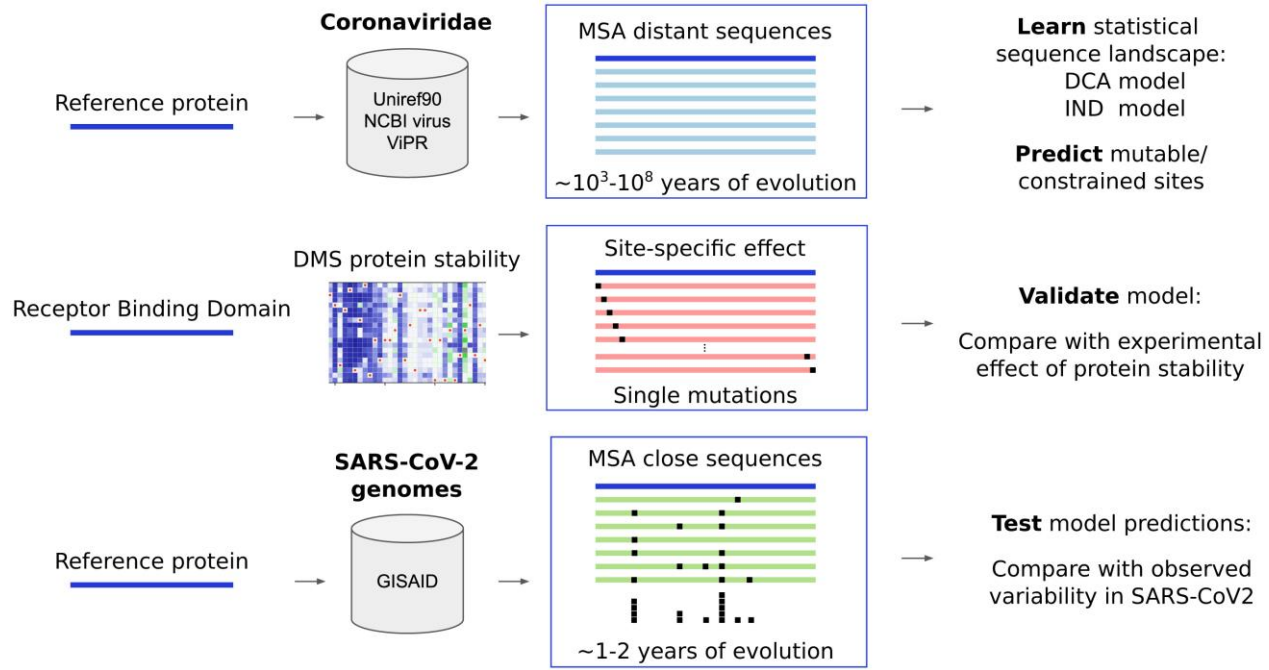
August 2021					September 2021					October 2021				
rank	Partner		Spike		rank	Partner		Spike		rank	Partner		Spike	
	locus	AA-m.	locus	AA-m.		locus	AA-m.	locus	AA-m.		locus	AA-m.	locus	AA-m.
1	17236	nsp13:I334V	24208	I882I	1	17236	nsp13:I334V	24208	I882I	1	17236	nsp13:I334V	24208	I882I
14	7851	nsp3:A1711V	21846	T95I	13	7851	nsp3:A1711V	21846	T95I	10	7851	nsp3:A1711V	21846	T95I
20	28461	N:G63D	24410	D950N	16	28461	N: D63G	24410	D950N	17	28461	N:D63G	24410	D950N
27	1048	nsp2:K81N	21846	T95I	36	1048	nsp2:K81N	21846	T95I	20	25614	ORF3a:S74S	21995	T145H
52	26107	ORF3a:E239Q	21897	S112L	52	25614	ORF3a:S74S	21995	T145H	21	25614	ORF3a: S74S	22227	A222V
57	27507	ORF7a:G38G	21897	S112L	57	26107	ORF3a:E239Q	21897	S112L	30	1048	nsp2:K81N	21846	T95I
62	18086	nsp14:T16I	22792	I410I	58	25614	ORF3a:S74S	22227	A222V	51	10977	nsp6:A2V	21846	T95I
76	27291	ORF6:D30D	24208	I882I	71	27507	ORF7a:G38G	21897	S112L	56	27291	ORF6:D30D	24208	I882I
79	1729	nsp2:V308V	22792	I410I	82	27291	ORF6:G30G	24208	I882I	60	26107	ORF3a:E239Q	21897	S112L
151	28007	ORF8:P38P	21846	T95I	83	11514	nsp6:T181I	22227	A222V	63	29253	N:S327L	21846	T95I
168	27604	ORF7a:V71I	21846	T95I	128	17236	nsp13:I334V	21846	T95I	64	18744	nsp14:T235T	24130	N856N
174	17236	nsp13:I334V	21846	T95I	151	18744	nsp14:T235T	24130	N856N	74	27507	ORF7a:G38G	21897	S112L
197	11514	nsp6:T181I	22227	A222V	190	5584	nsp3:T955T	22227	A222V	80	17236	nsp13:I334V	21846	T95I
					195	13019	nsp9:L112L	22227	A222V	124	15952	nsp12:S837S	21846	T95I
										153	26107	ORF3a:E239	21846	T95I
										163	28299	N:Q9L	21846	T95I
										190	27507	ORF7a:G38G	21846	T95I
										194	11562	nsp6:C197F	21897	S112L
										197	11514	nsp6:T181I	22227	A222V

TABLE II. Largest DCA terms with only one terminal in Spike coding region, August-October 2021. Top-200 couplings computed as plmDCA scores are considered. For each of them in the three months displayed, there's the indication of the rank, the locus in the Spike coding region and corresponding amino acid (AA) mutation, the locus in the partner coding region and corresponding amino acid (AA) mutation. Green color indicates that this mutation is found in delta variant. Red color indicates that this mutation is found in omicron variant. Pairs with one or both terminals colored green are attributed to a phylogenetic effect, while the several pairs with one terminal colored red are not, the growth of omicron being later than October-2021. Omicron mutations used here are taken from [67] on page 18, deletions not considered.

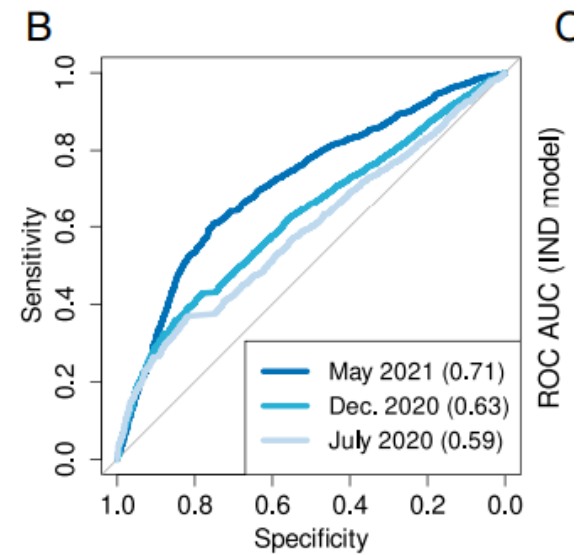


# Rodriguez-Rivas went deeper...

ROYAL INSTITUTE OF TECHNOLOGY



...used all coronavirus sequences to build an MSA, and then tested that on GISAID data.



Rodriguez-Rivas et al *PNAS* 2022, Fig 1

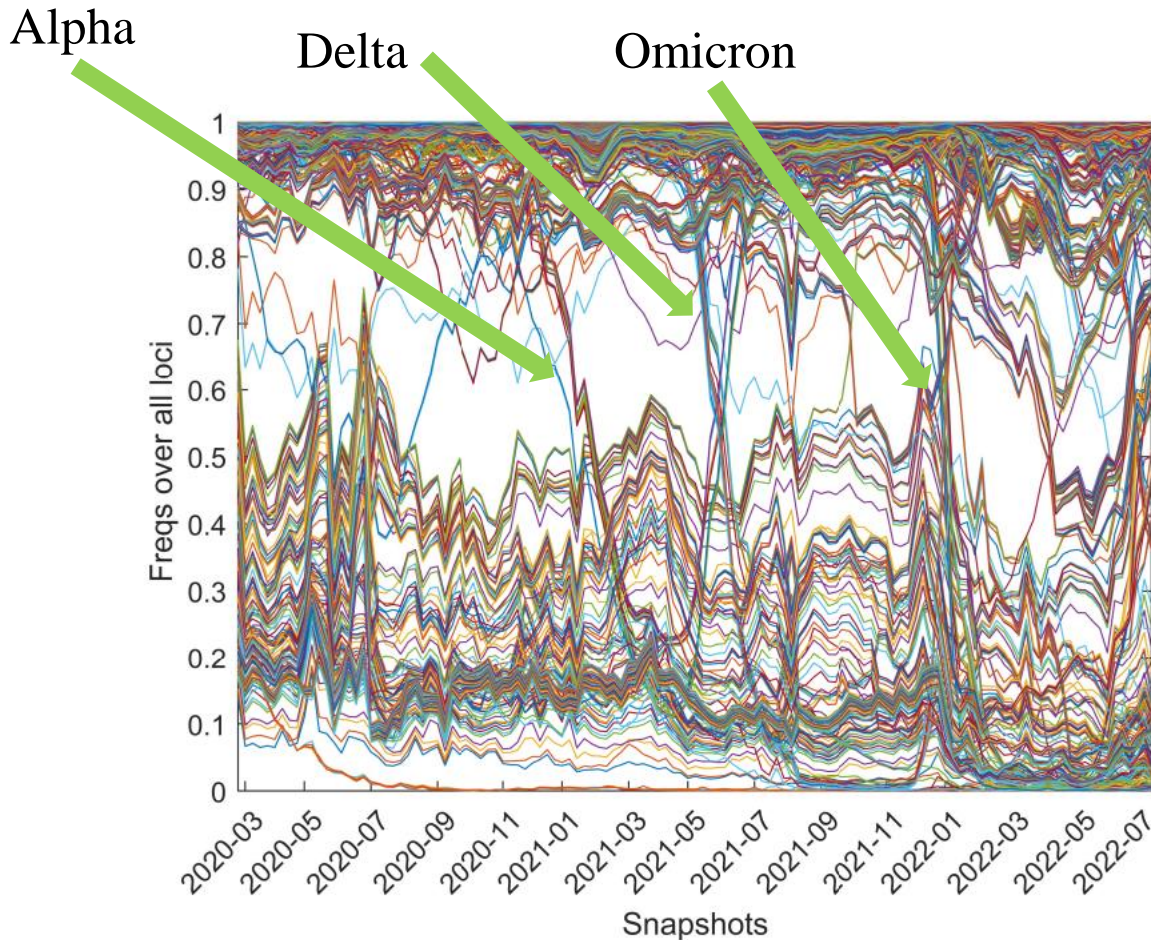
Rodriguez-Rivas et al *PNAS* 2022 Fig 4.B

$$\Delta E_{DCA}(i, b) = \log P_{DCA}(a_1, \dots, a_i, \dots, a_L) - \log P_{DCA}(a_1, \dots, b, \dots, a_L).$$

$$S_{IND/DCA}(i) = \frac{1}{q} \sum_{k=1}^q \Delta E_{IND/DCA}(i, b_k).$$

# SARS-CoV-2 perhaps also NRC?

Coronaviruses recombine. This has been observed in SARS-CoV-2, *in vivo*. Plots of allele frequencies at *all* loci show the well-known VoCs Alpha, Beta, Delta, Omicron...but also a bit more.



← frozen loci

Frequencies of all alleles on all positions per week from GISAID up to August 2022 [Zeng & Liu, unpublished] [ see also arXiv:2109.02962]

← An NRC phase? Most of these intermittently fluctuating loci lie in the 5' or 3' end of the SARS-CoV-2 genome.

# Thanks

Hong-Li Zeng

Vito Dichio

Yue Liu

Eugenio Mauri

Simona Cocco

Rémi Monasson



Vetenskapsrådet

Fabbio Cecconi

Chen-Yi Gao

Angelo Vulpiani

Hai-Jun Zhou

Boris Shraiman

Richard Neher

Benjamin Good

Michael Desai

National Natural Science Foundation of China (11705097), Natural Science Foundation of Nanjing University of Posts and Telecommunications (Grant Nos. 221101 and 222134), Swedish Research Council (Grant 2020-04980).