



SORBONNE
UNIVERSITÉ



STATISTICAL PHYSICS FOR
QUANTITATIVE BIOLOGY



SAPIENZA
UNIVERSITÀ DI ROMA

Emergent time scales of epistasis in protein evolution

Francesco Zamponi

Mathematical Physics Webinar, Rutgers University,
February 19, 2025

Bisardi, Rodriguez-Rivas, FZ, Weigt, Mol.Bio.Evo. (2022)

Di Bari, Bisardi, Cotogno, Weigt, FZ, PNAS (2024)

Rossi, Di Bari, Weigt, FZ, arXiv:2412.01969 (2024)



Martin
Weigt



Matteo
Bisardi



Sabrina
Cotogno



Juan
Rodriguez
-Rivas



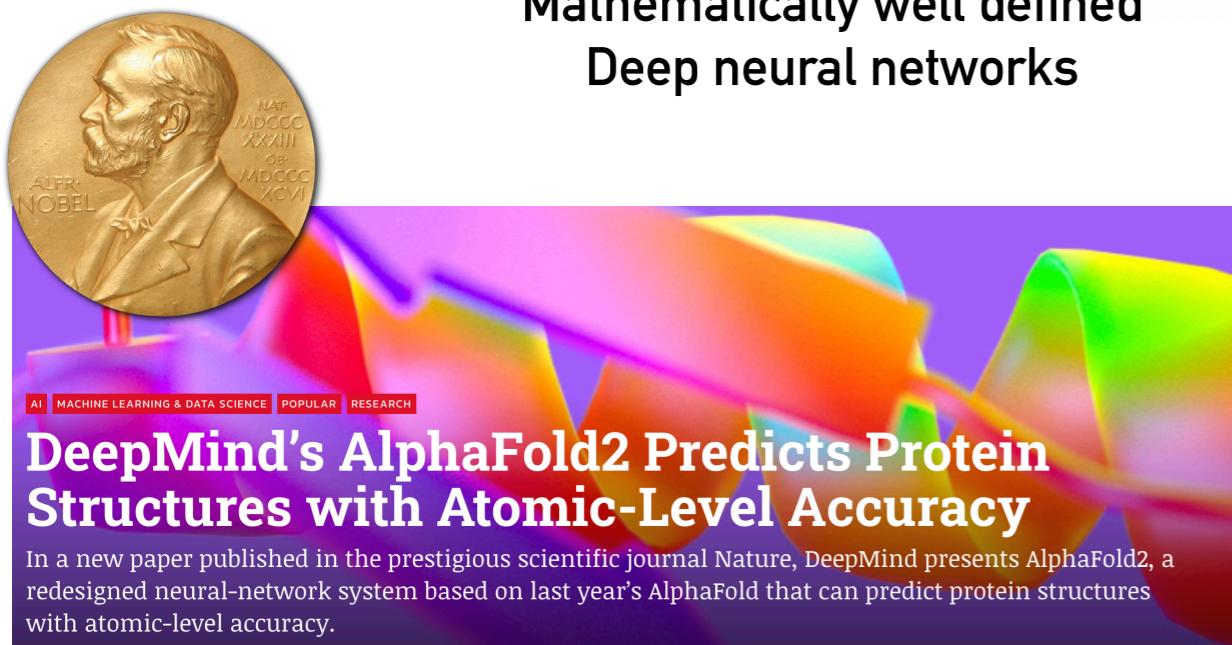
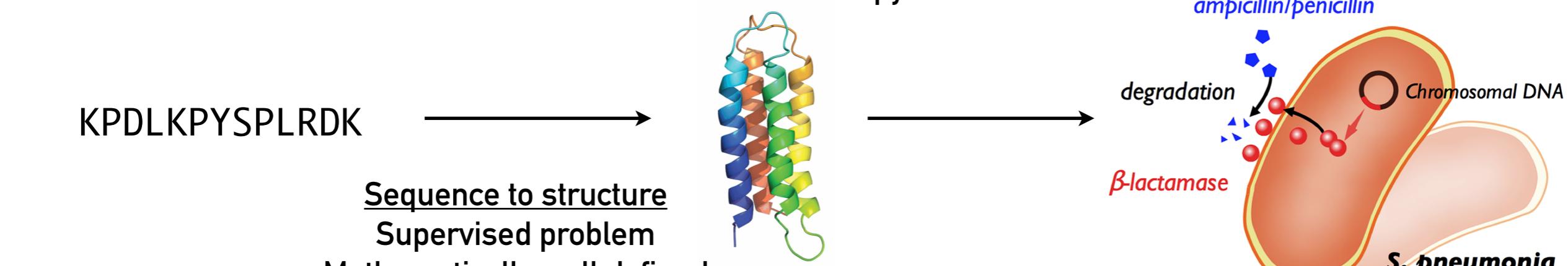
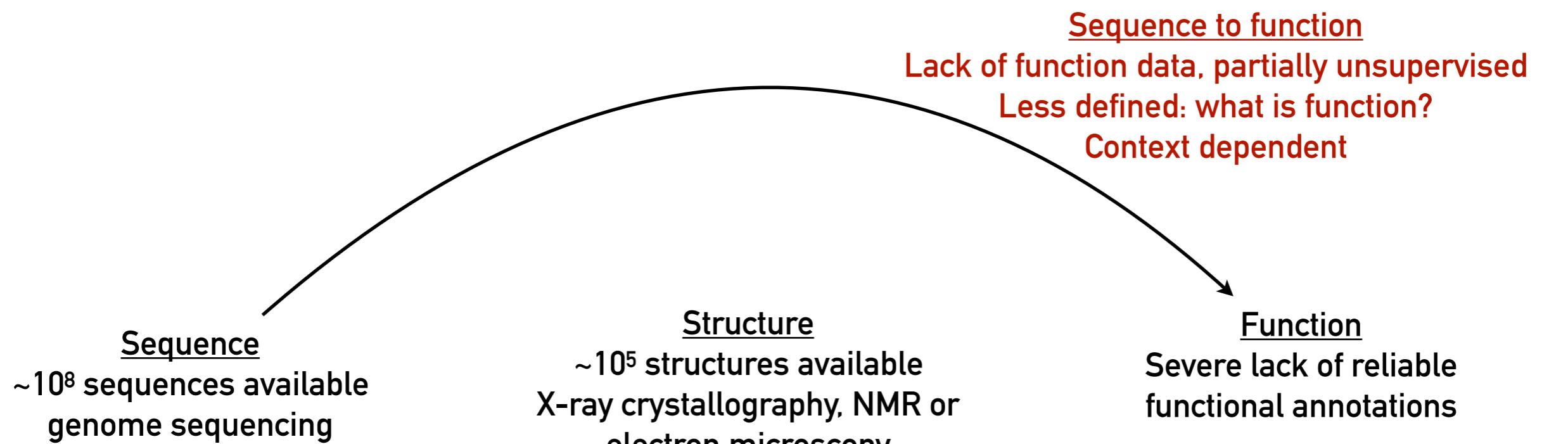
Leonardo
Di Bari



Saverio
Rossi

Goal

Sequence-to-structure-to-function paradigm



Long term goals

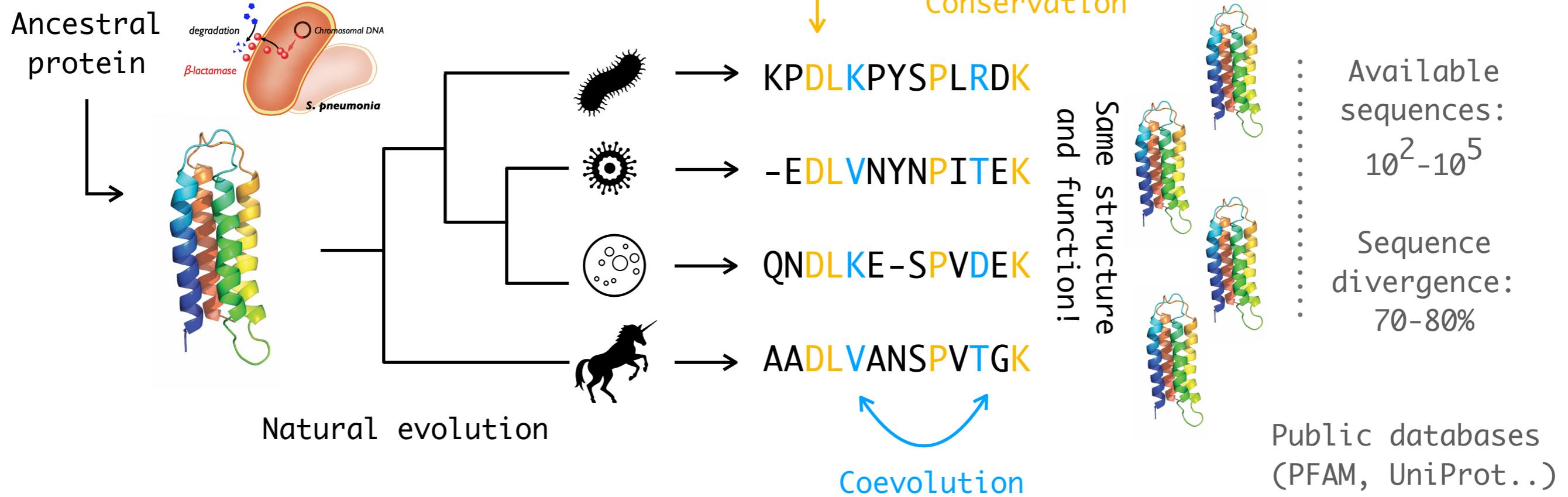
- Learn something about the sequence-to-function relationship
- Design artificial protein sequences with desired functions
- Understand how nature does it by evolution
- Understand and optimize protocols to do it in the lab (directed in vitro evolution)

Take home: we use a data-driven approach and statistical physics modelling

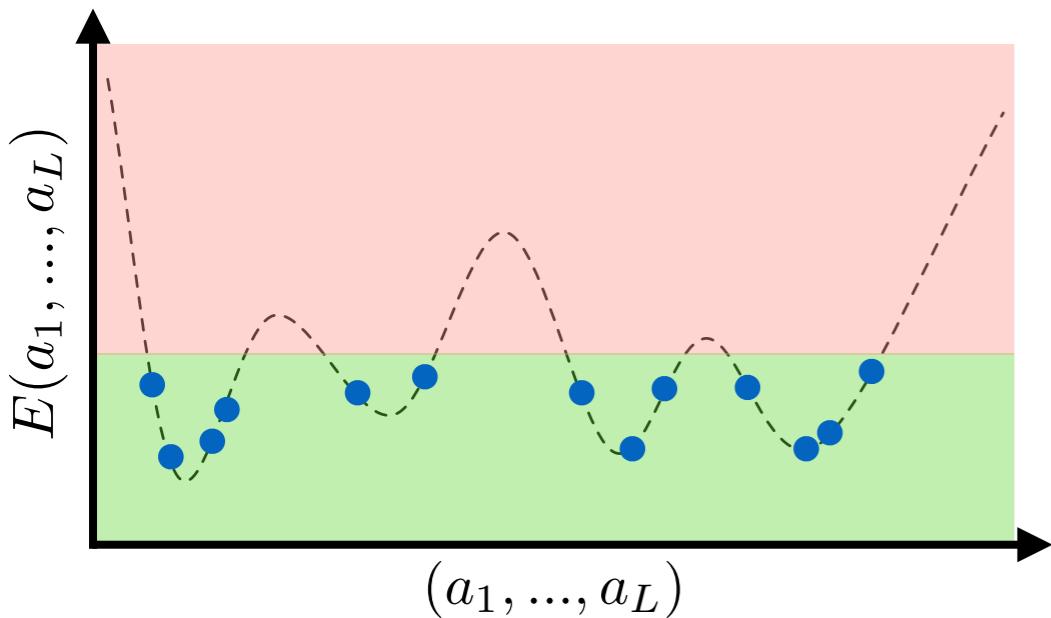


Data

Data #1: natural protein sequences



sequence landscape



sequence data

global sample

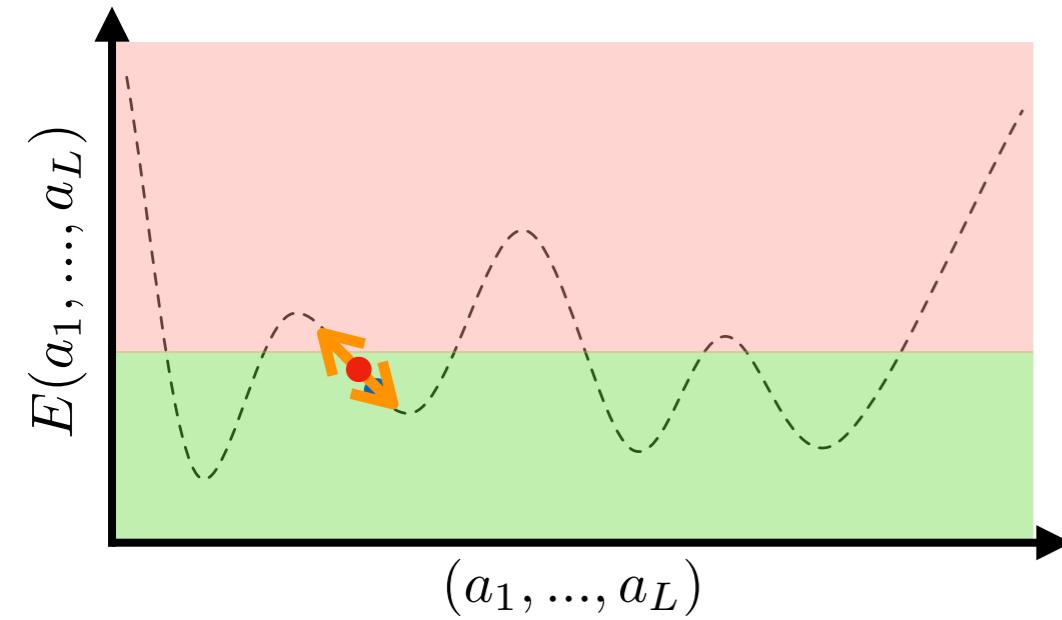
species 1
species 2
species 3
...

Data #2: deep mutational scans



Fowler and Fields, Nature Methods (2014)

sequence landscape



sequence data

reference sequence

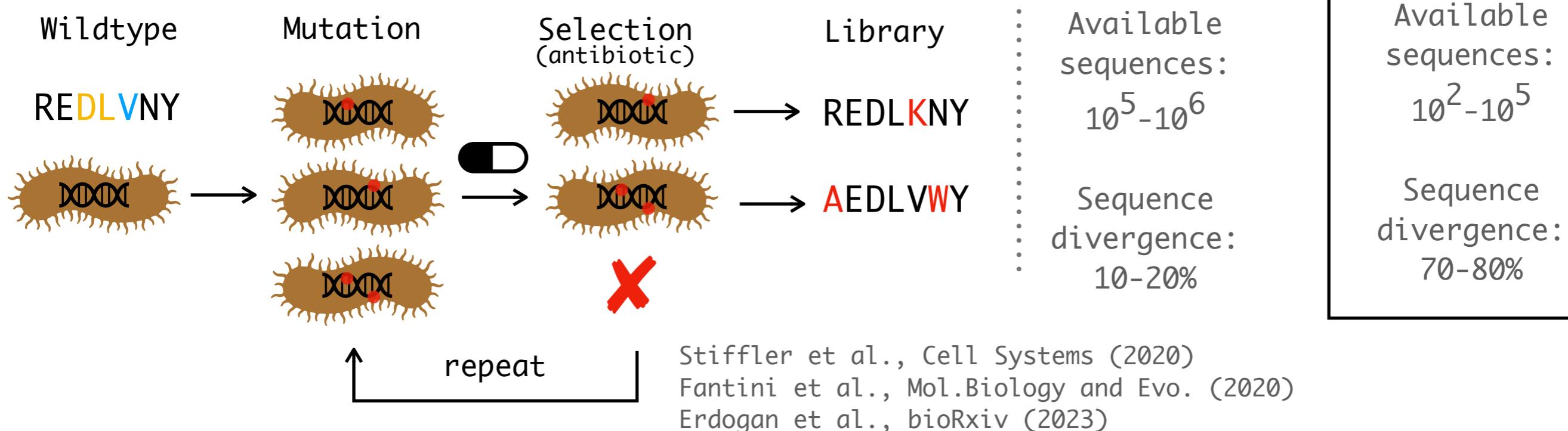
mutant sequences

- | | |
|-----|-------|
| • | E_1 |
| • | E_2 |
| • | E_3 |
| • | E_4 |
| • | E_5 |
| • | E_6 |
| ... | |
- phenotype

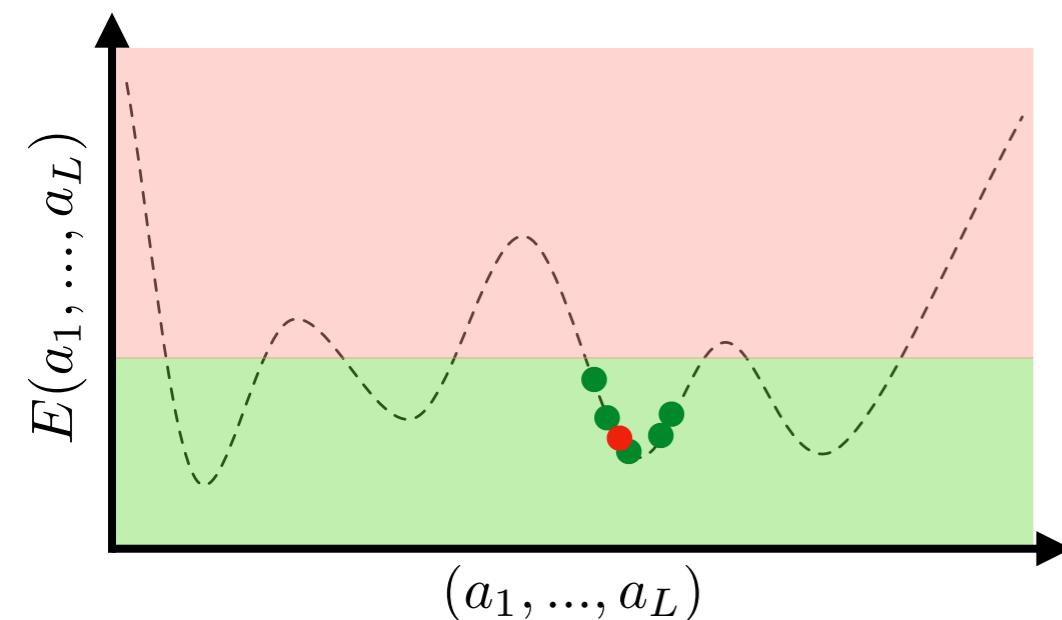
Data #3: in vitro evolution

(weak selection)

Natural evolution



sequence landscape



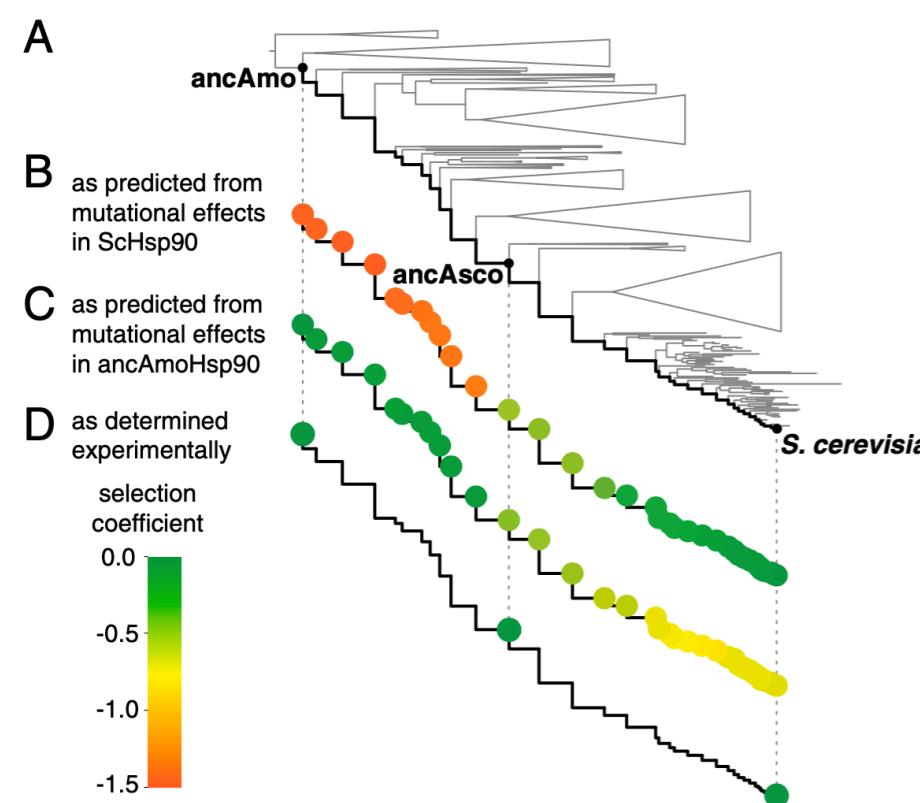
sequence data

reference sequence

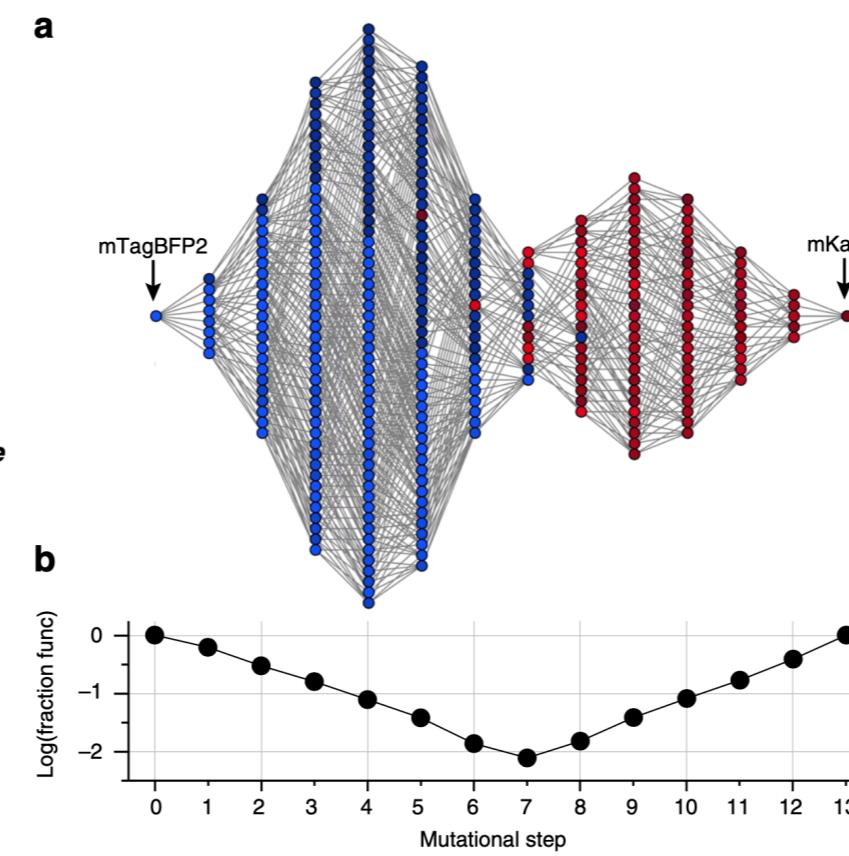
local sample

mutant 1
mutant 2
mutant 3
...

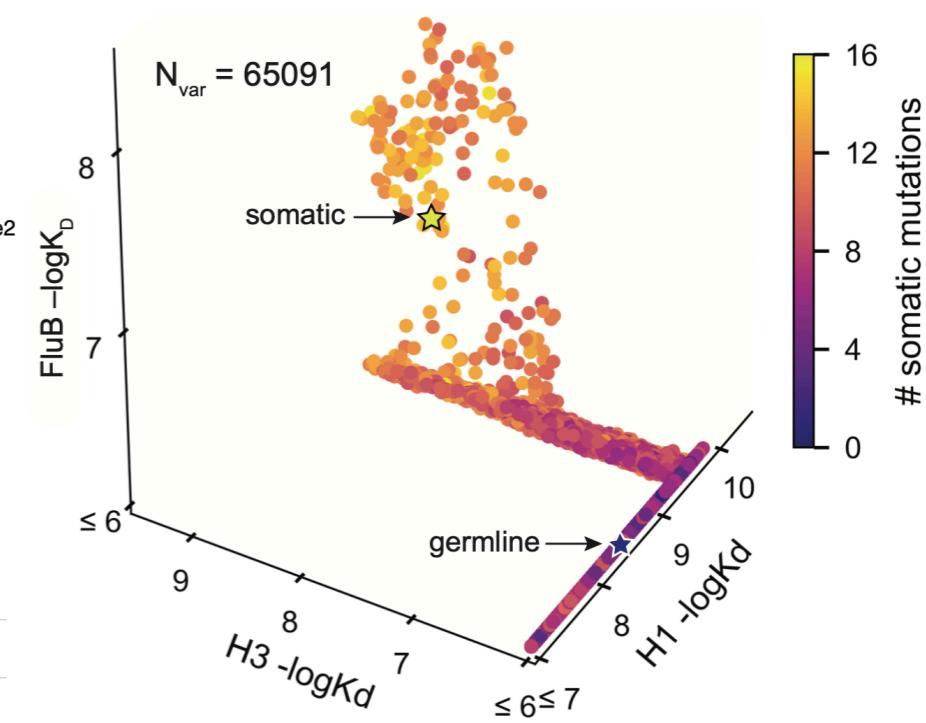
Data #4: path/phylogeny reconstruction



Starr et al. PNAS (2017)

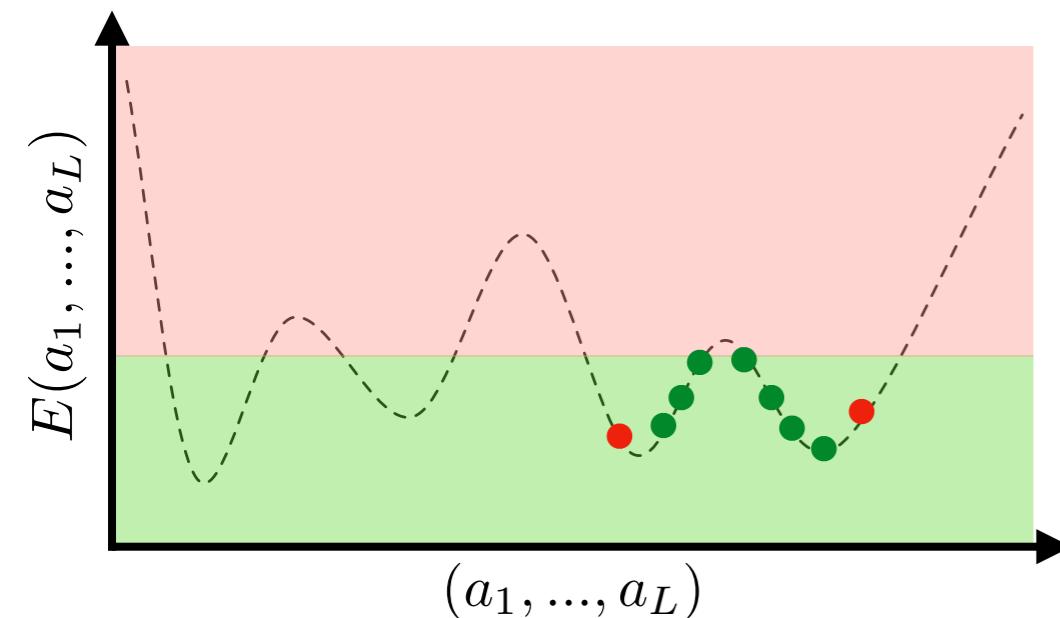


Poelwijk et al. Nat.Comm. (2019)



Phillips et al. eLife (2021)

sequence landscape



sequence data

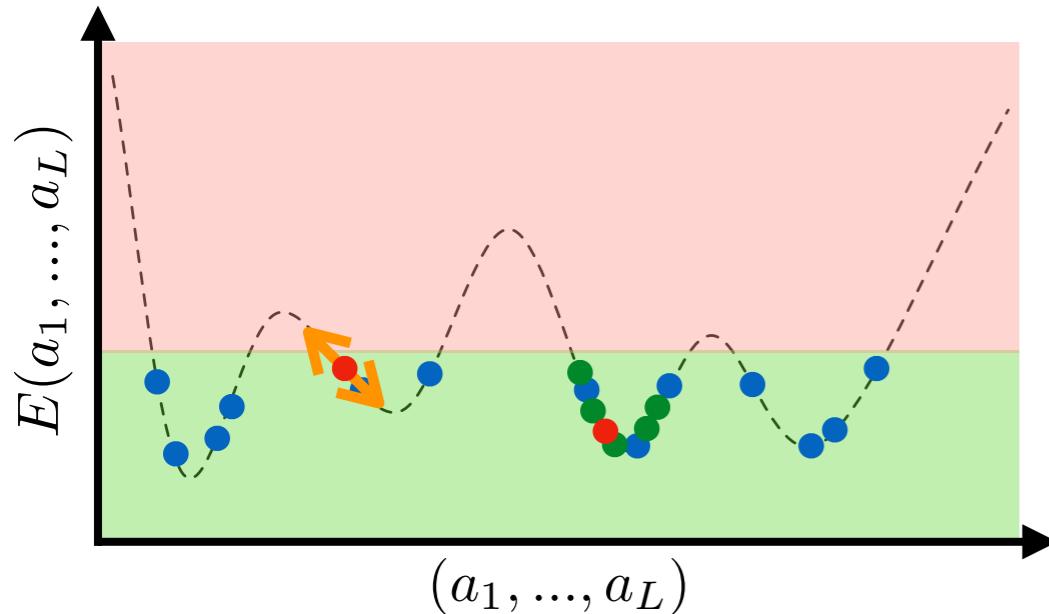
reference sequence

local sample

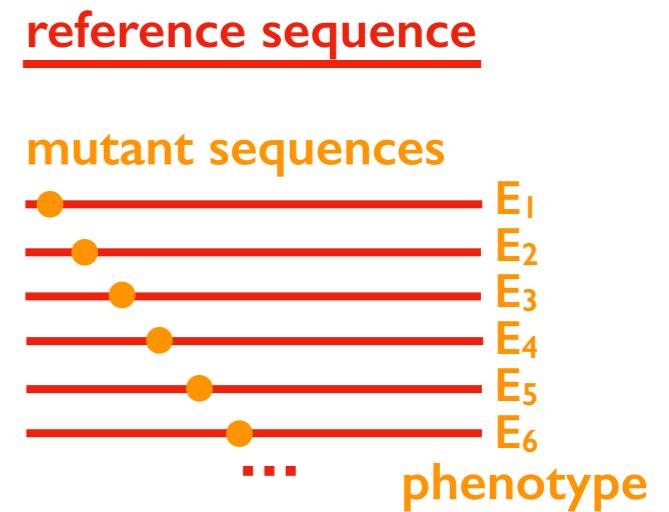
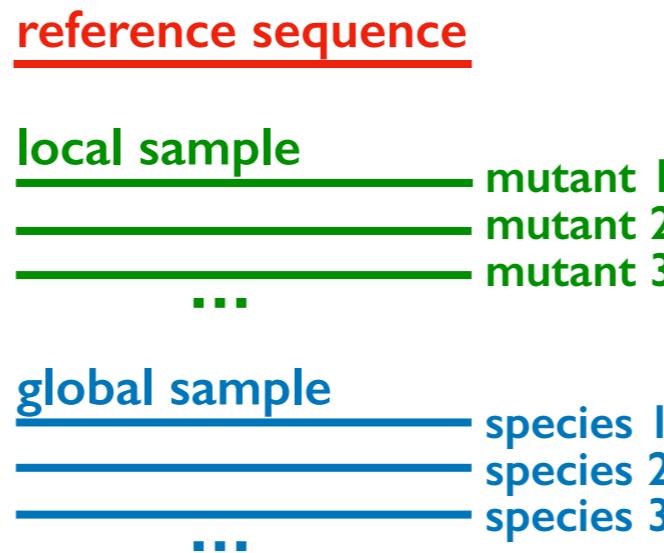
- path 1
- path 2
- ...
- path 3

Massive amount of data!

sequence landscape



sequence data



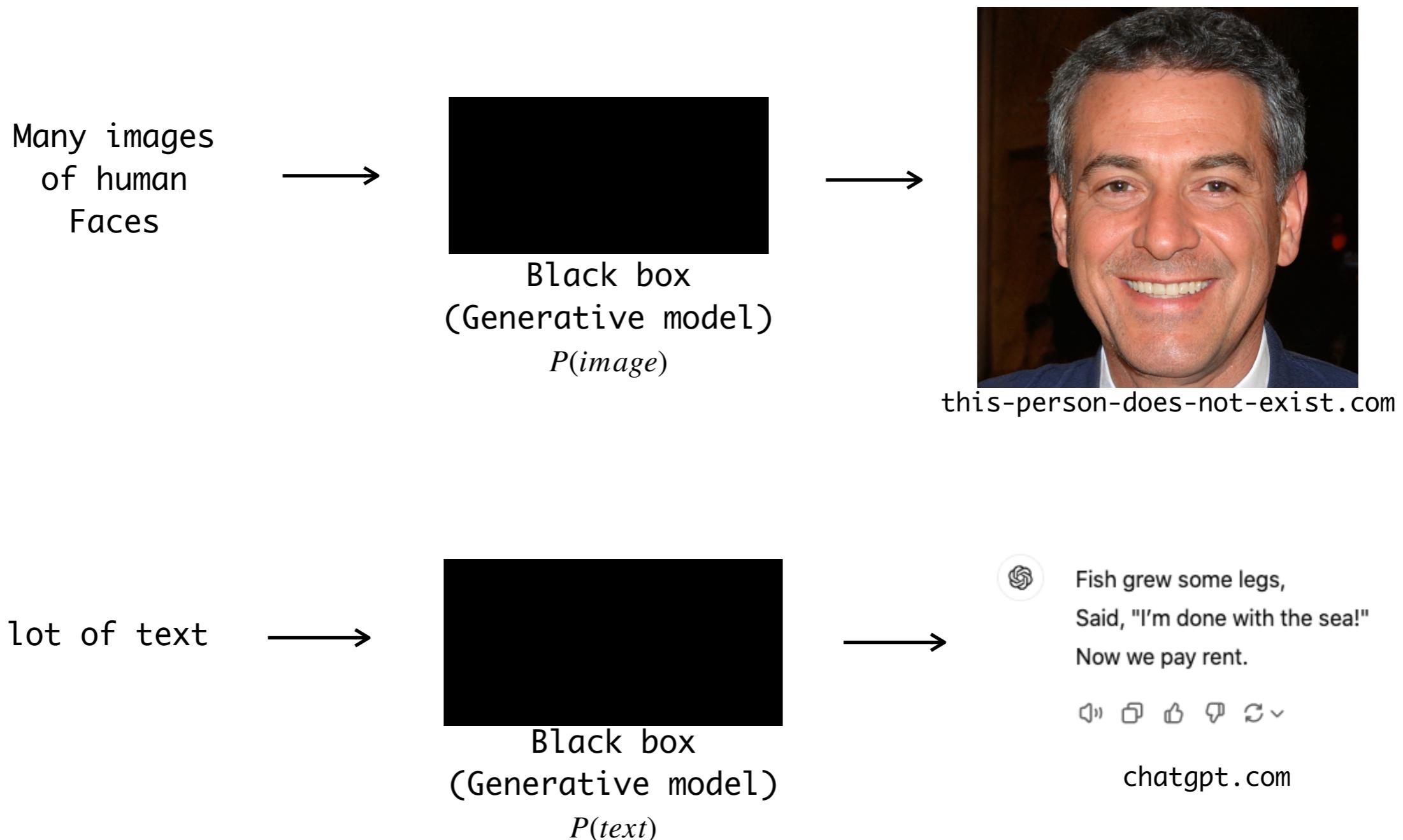
Characterize the dynamics in the sequence landscape:

- Model in vitro evolution (controlled environment)
- Model natural evolution (phylogeny, fluctuating environment)
- Optimize evolution protocols (antibiotic concentration, vaccination protocols...)
- Generate paths of artificial sequences connecting two natural ones

Take home: there are many data, and many more to come

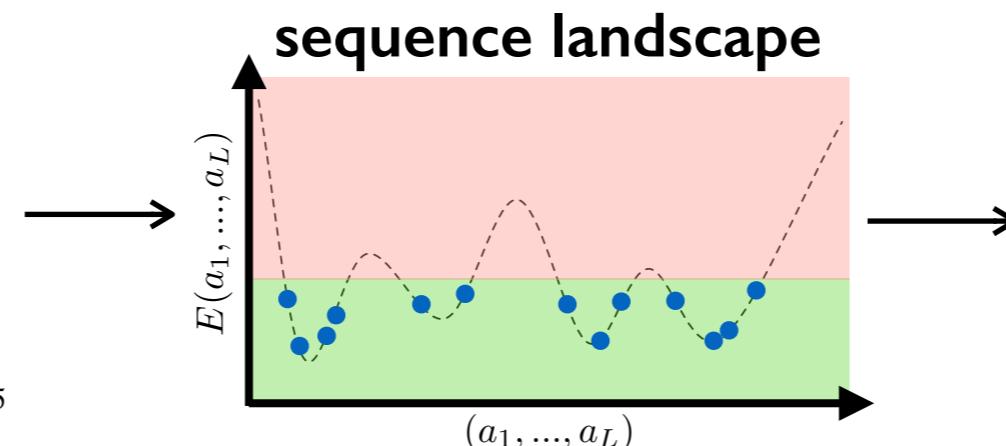
Generative modelling

Generative modelling (unsupervised)



Generative modelling of biological sequences

KPDLK**PYSPLRDK**
 -EDLV**NYNPITEK**
 QNDL**KE-SPVDEK**
 AADLV**ANSPVVTGK**
 Alignment of $10^2 - 10^5$ natural proteins



$$P(a_1, a_2, \dots, a_L) = \frac{1}{Z} \exp\left(\sum_i^{1,L} h_i(a_i) + \sum_{i < j}^{1,L} J_{ij}(a_i, a_j)\right) = \frac{1}{Z} e^{-E(a_1, \dots, a_L)}$$

Conservation Coevolution

Boltzmann learning

$h_i(a_i)$ $J_{ij}(a_i, a_j)$

Statistical physics approach: maximum likelihood, disordered Potts models

Coevolution is at the basis of all structure prediction methods (including AlphaFold)

Cocco et al., Reports on Progress in Physics (2018)

Simpler models

Because we have less data, but...

More interpretable

Energy efficient

Accessible to everyone

Trump Announces \$100 Billion A.I. Initiative

OpenAI, Oracle and SoftBank formed a new joint venture called Stargate to invest in data centers, building on major U.S. investments in the technology.

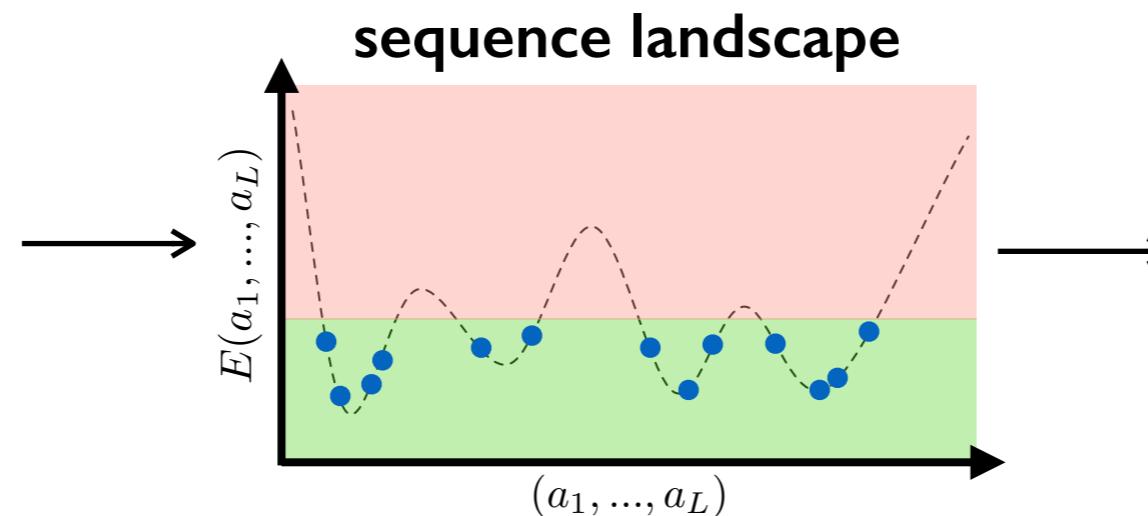
The New York Times



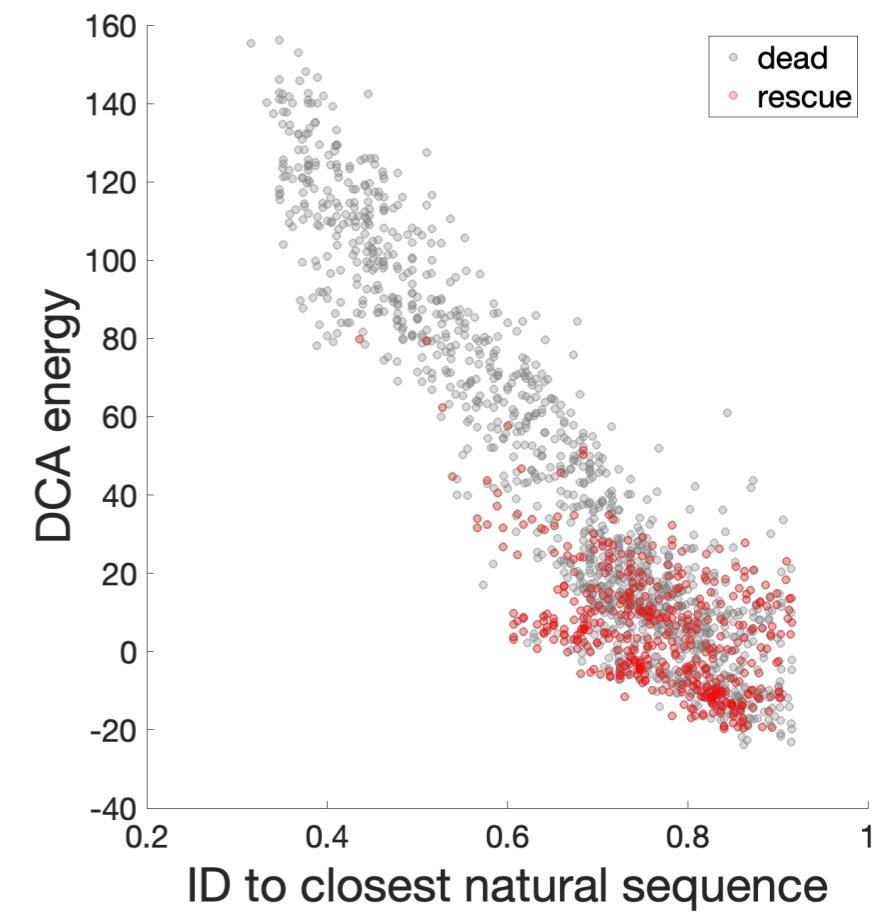
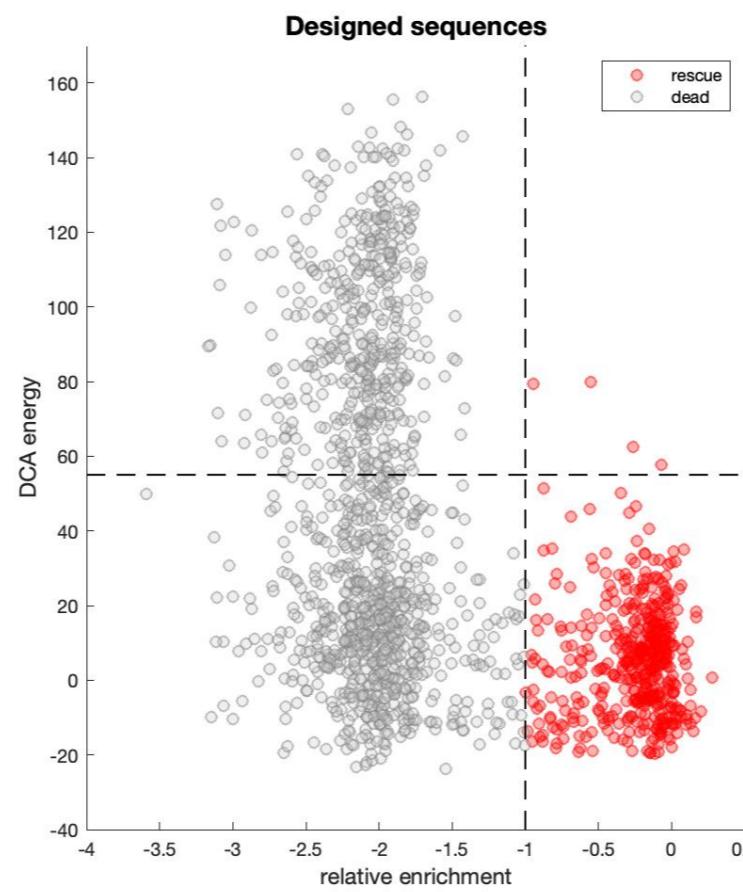
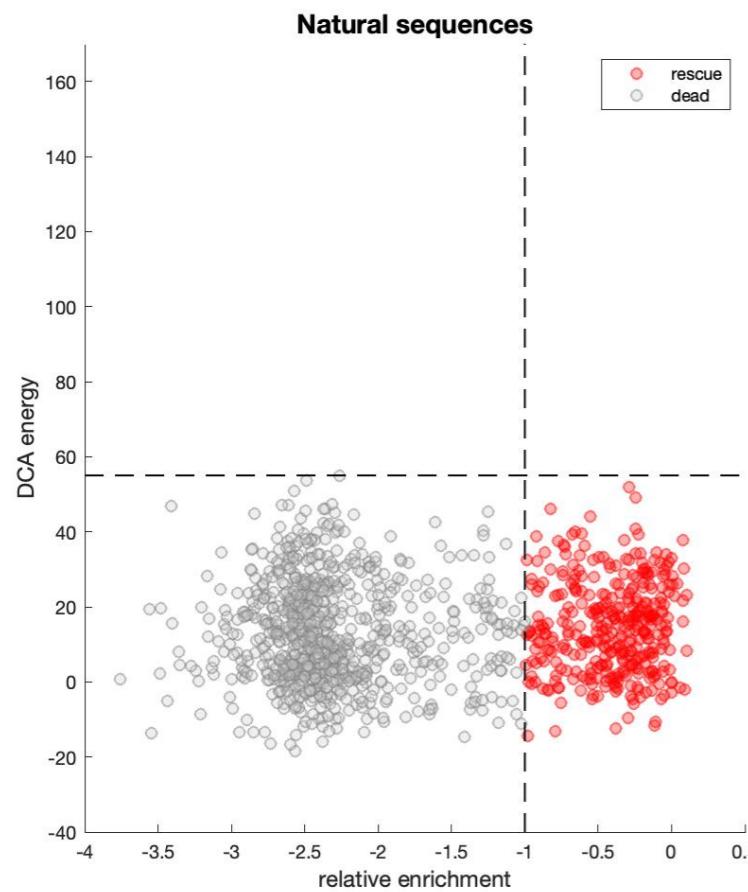
By Cecilia Kang and Cade Metz
 Reporting from Washington and S

Generative modelling of biological sequences

KPDLKPYSPPLRDK
 -EDLVNYNPIKEK
 QNDLKE-SPVDEK
 AADLVANSPTGK



CYDLVGVWEPTAK
 $P(a_1, a_2, \dots, a_L) \sim \exp[-E(a_1, \dots, a_L)]$



Typical protein with $L = 100$ (Chorismate Mutase expressed in E.Coli)

$20^L \sim 10^{130}$ possible sequences

$S[P] \sim 1.5 \sim \log(q)/2 \rightarrow 10^{65}$ functional sequences

Russ et al. Science (2020)

Generative modelling of biological sequences

Generative models have been experimentally validated!

Many other different architectures (GAN, VAE, deep or not)

Russ et al. Science (2020)

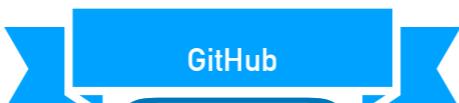
Repack et al. Nature Machine Intelligence (2021)

Hawkins-Hooker et al. PLoS Comp. Bio. (2021)

Two public packages that everyone can use: adabmDCA and arDCA

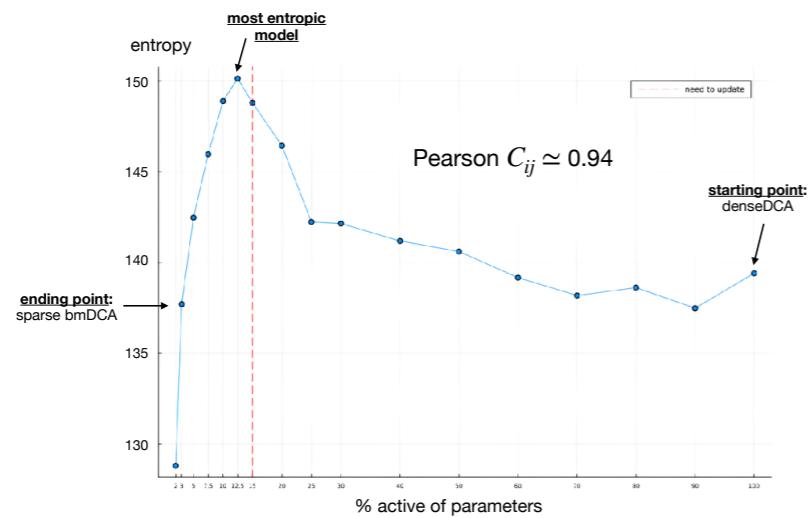


STATISTICAL PHYSICS FOR QUANTITATIVE BIOLOGY



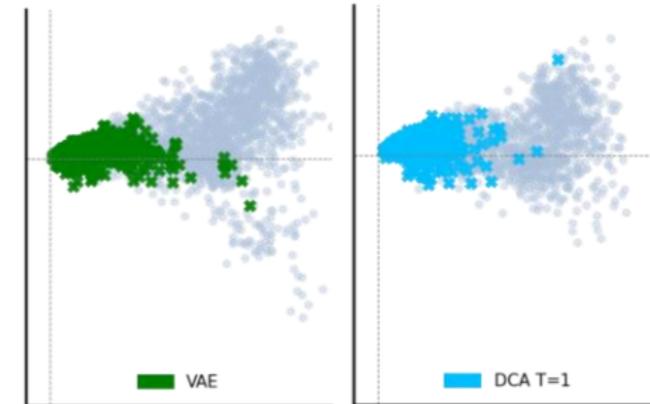
Trinquier et al. Nat. Comm. (2021)
Muntoni et al. BMC Bioinformatics (2021)
Rosset et al. Springer Protocol (to appear)

Information-based parameter reduction and entropy maximization



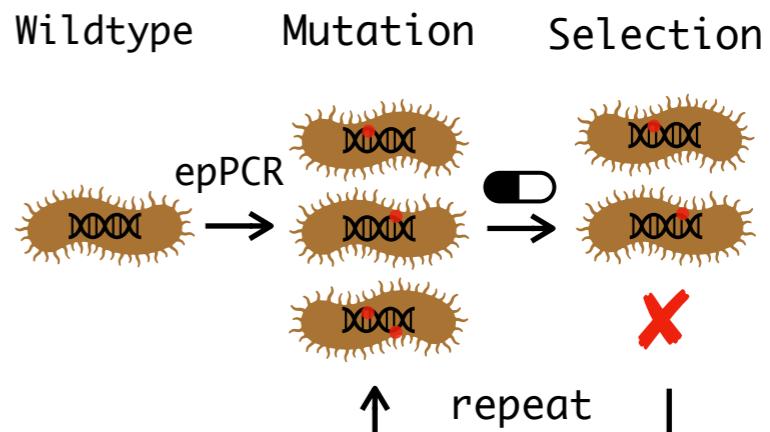
Barrat-Charlaix et al. PRE (2021)
Ongoing experimental test with Ranganathan's group

The model is simple and interpretable, and not generically outperformed by other methods



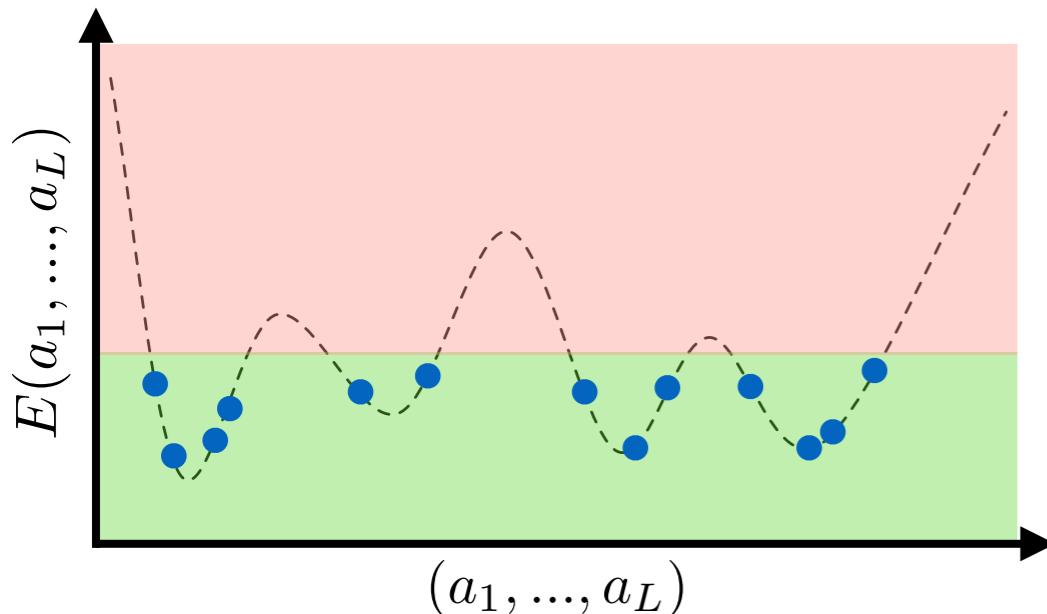
Calvanese et al. Nucleic Acids Research (2024)
Lambert, Opuu, Calvanese et al. bioRxiv (2024)

Model and direct evolution (today's main topic)



What comes next?

sequence landscape



Equilibrium sampling

$$P(a_1, a_2, \dots, a_L) \sim \exp[-E(a_1, \dots, a_L)]$$

Generate artificial sequences



Tune their properties



Take home: simple generative models can design artificial proteins (equilibrium sampling)

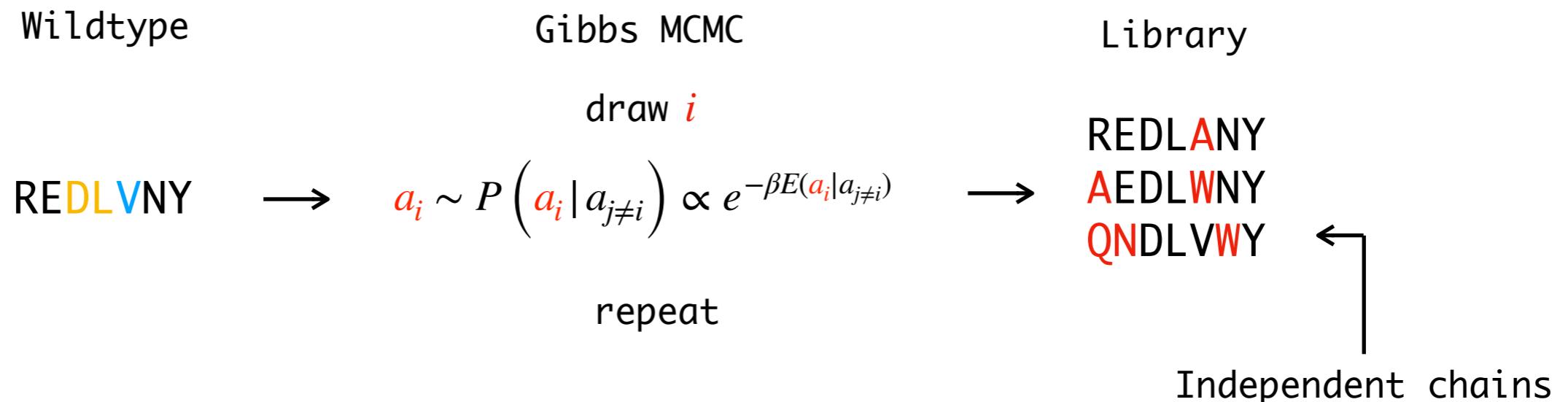
Next: characterize the shape of the sequence landscape and the dynamics

Dynamics

Local sampling ~ evolution?

$$P(a_1, a_2, \dots, a_L) = \frac{1}{Z} \exp\left(\sum_i^{1,L} h_i(a_i) + \sum_{i < j}^{1,L} J_{ij}(a_i, a_j) \right) = \frac{1}{Z} e^{-E(a_1, \dots, a_L)}$$

↓ Conservation
↑ Coevolution

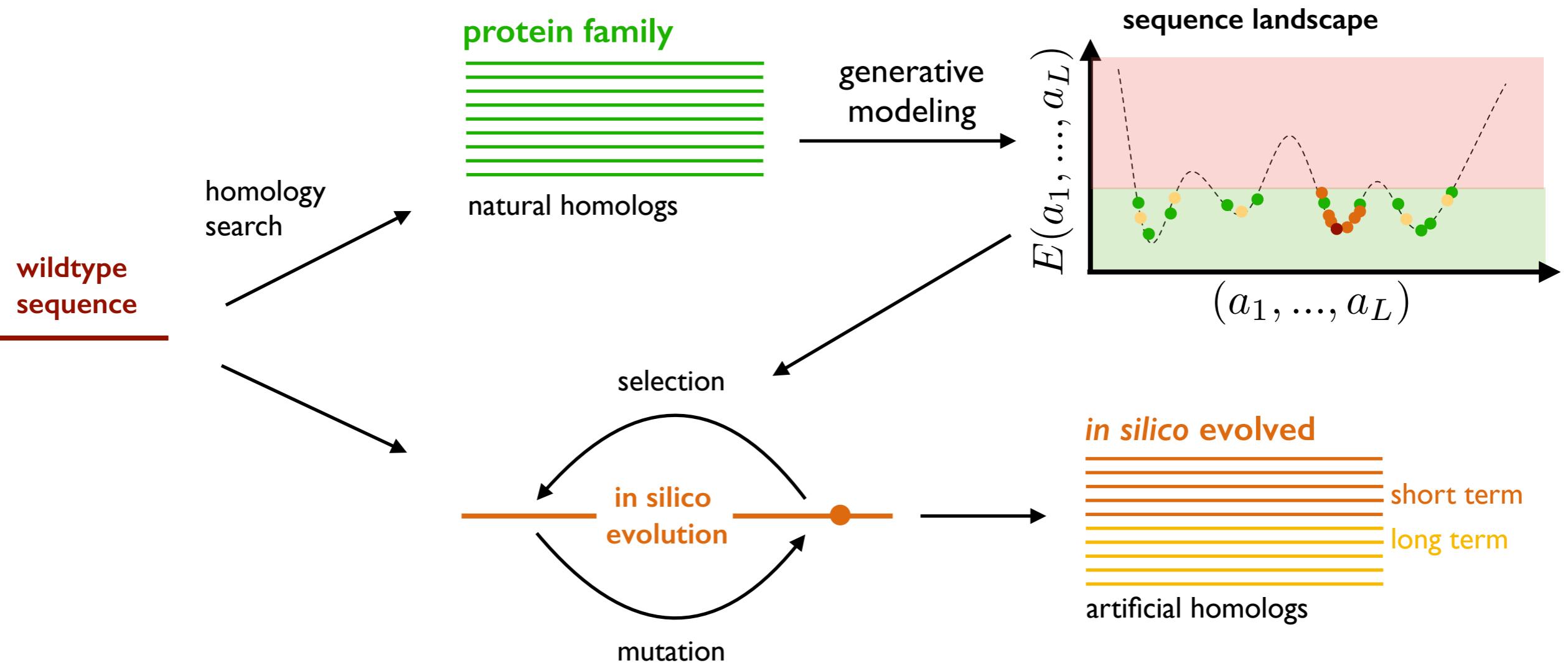


De la Paz et al., PNAS (2020)

a_i can be drawn:

- From the set of all 20 amino acids
- From the set of amino acids corresponding to codons with one mutation
- With or without gaps

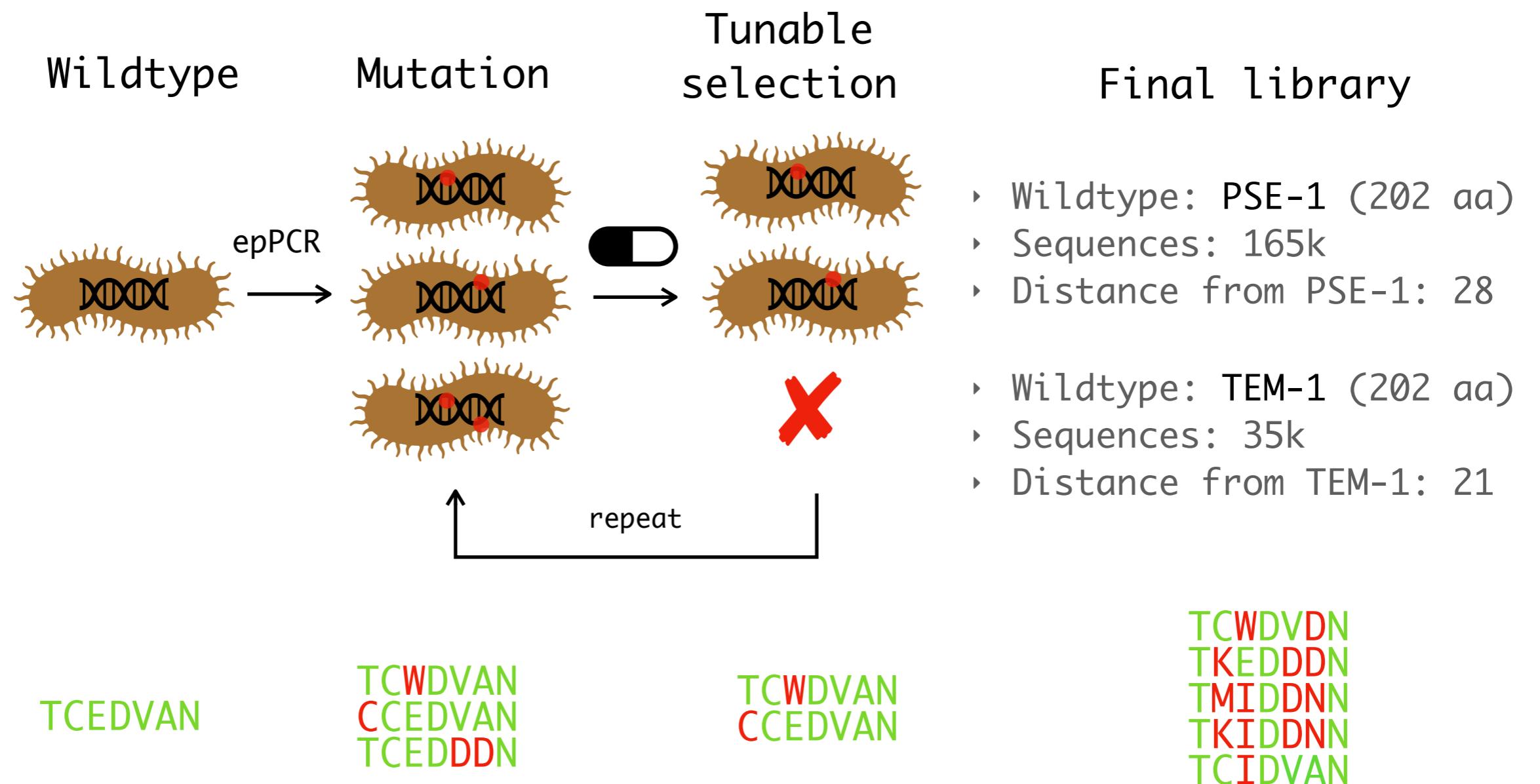
Local sampling ~ evolution?



Natural evolution \neq In-vitro evolution

Phylogenetic effects \neq Independent chains (star phylogeny)

Experimental data

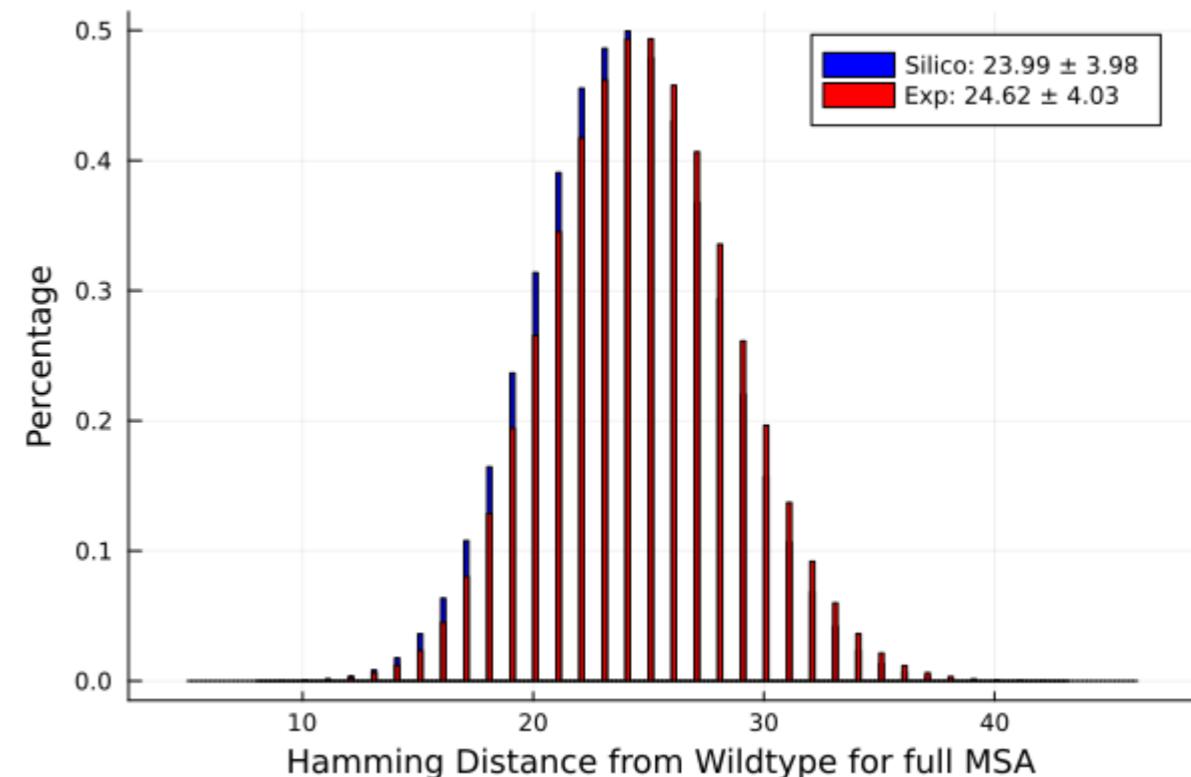
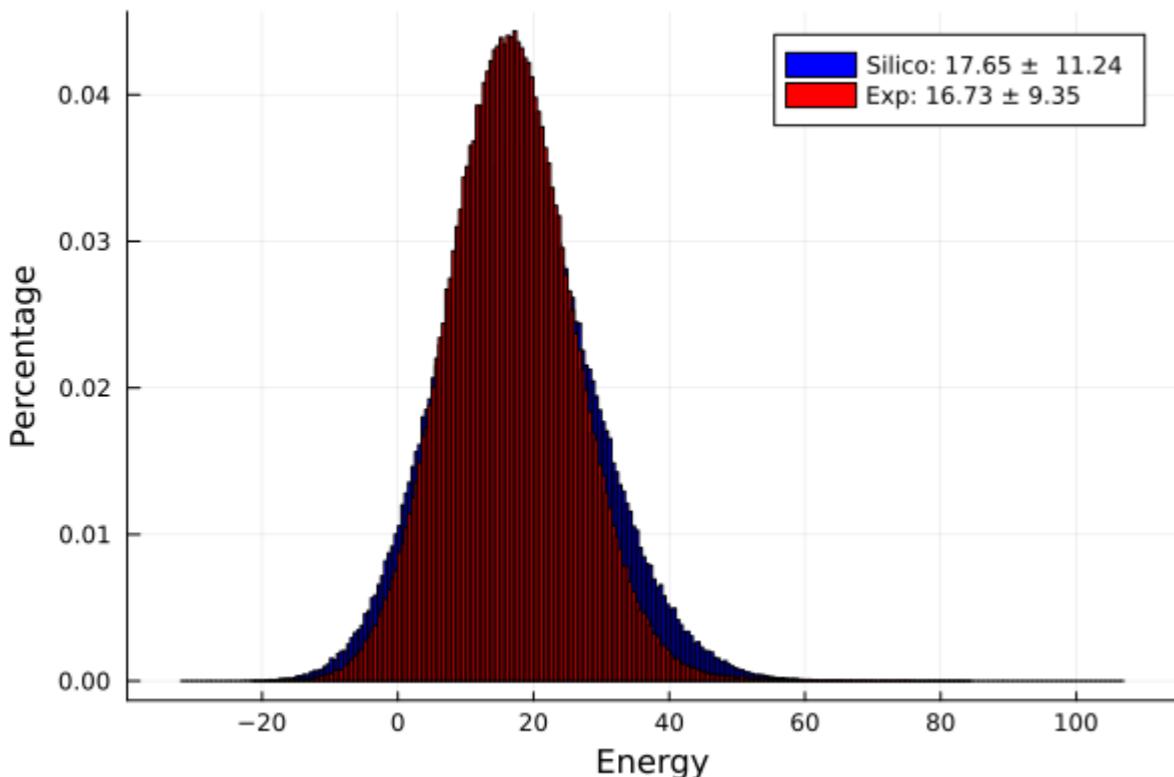


Stiffler et al., Cell Systems (2020)
 Fantini et al., Mol.Biology and Evo. (2020)
 New data (in progress): Erdogan et al., bioRxiv (2023)

Calibrating the model

Wildtype: PSE-1 Beta-lactamase, 20 generations of in vitro evolution

Stiffler et al., Cell Systems (2020)



Calibrate three parameters:

1. Artificial temperature that models selection strength
2. Number of MCMC steps from reference wild type
3. Number of sequences

Gibbs MCMC

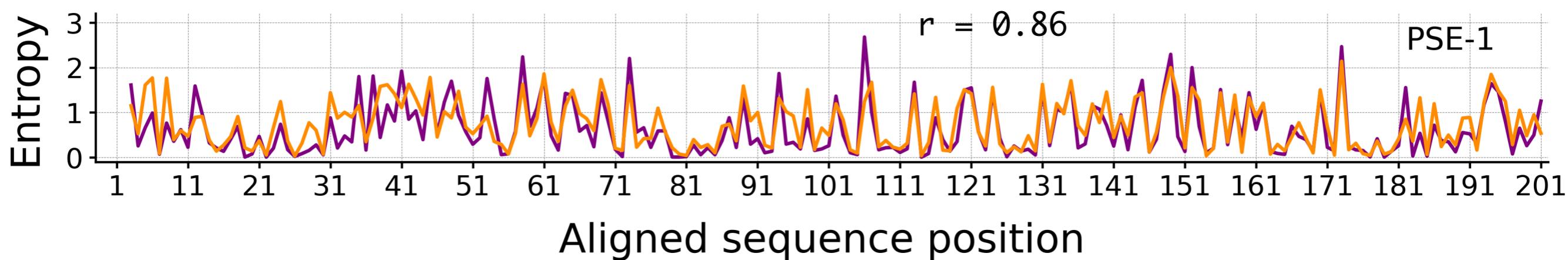
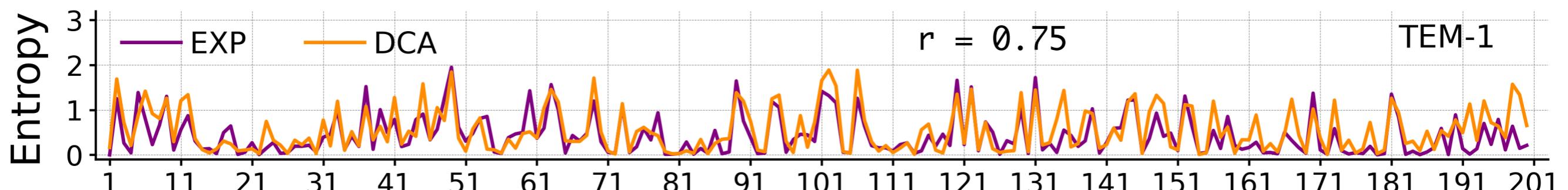
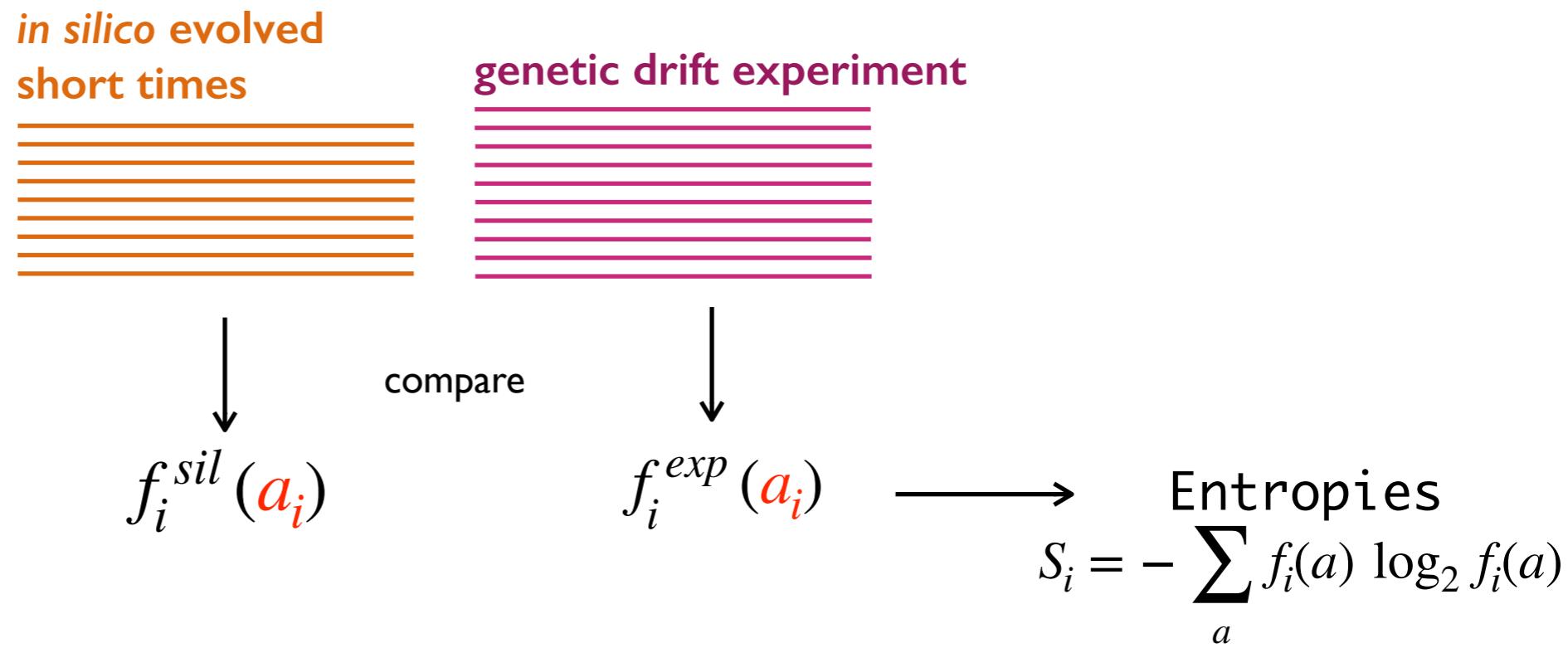
draw i

$$a_i \sim P(a_i | a_{j \neq i}) \propto e^{-\beta E(a_i | a_{j \neq i})}$$

repeat

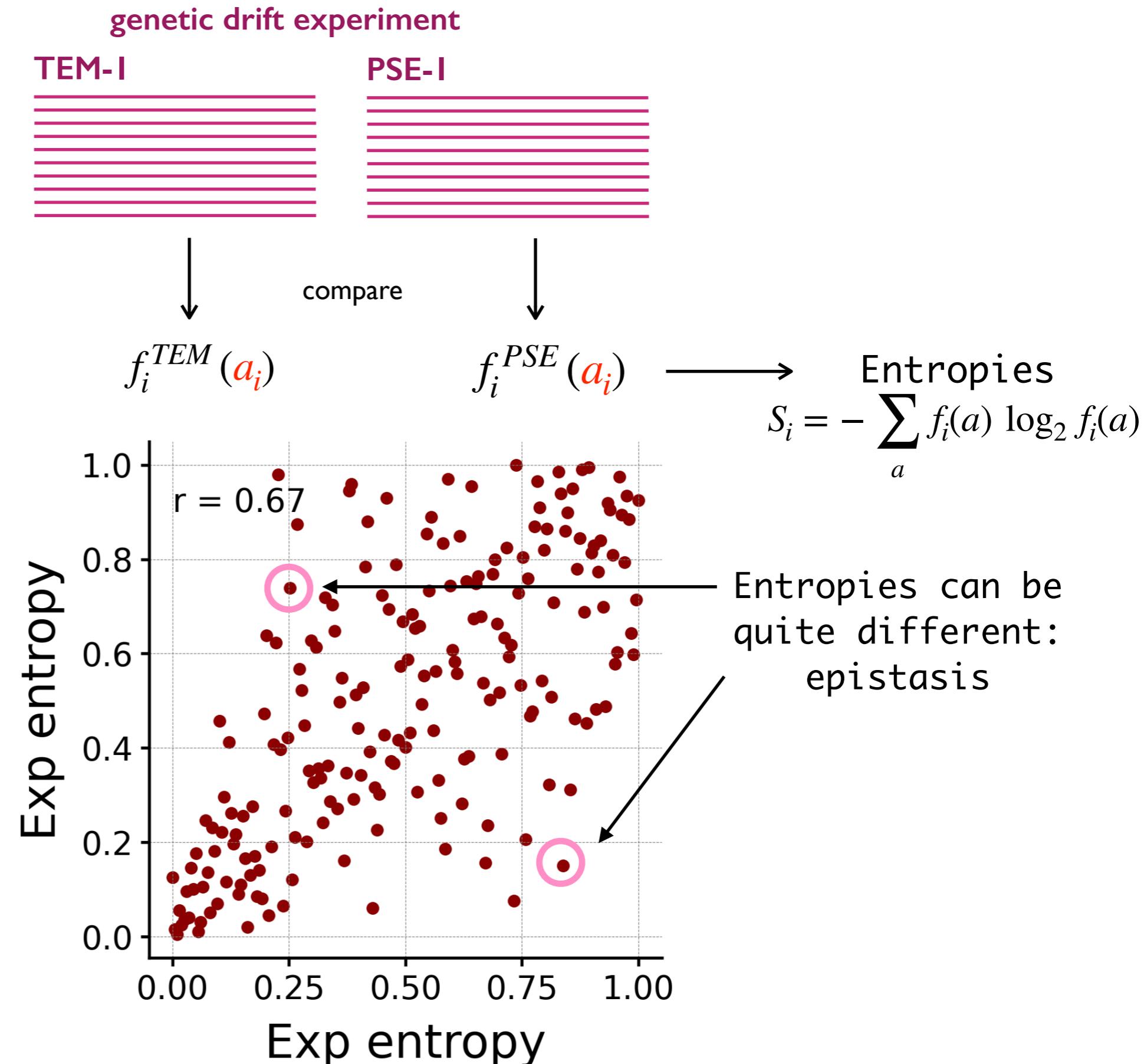
Short-term evolution

Calibrated:
 ► temperature
 ► steps
 ► number of sequences

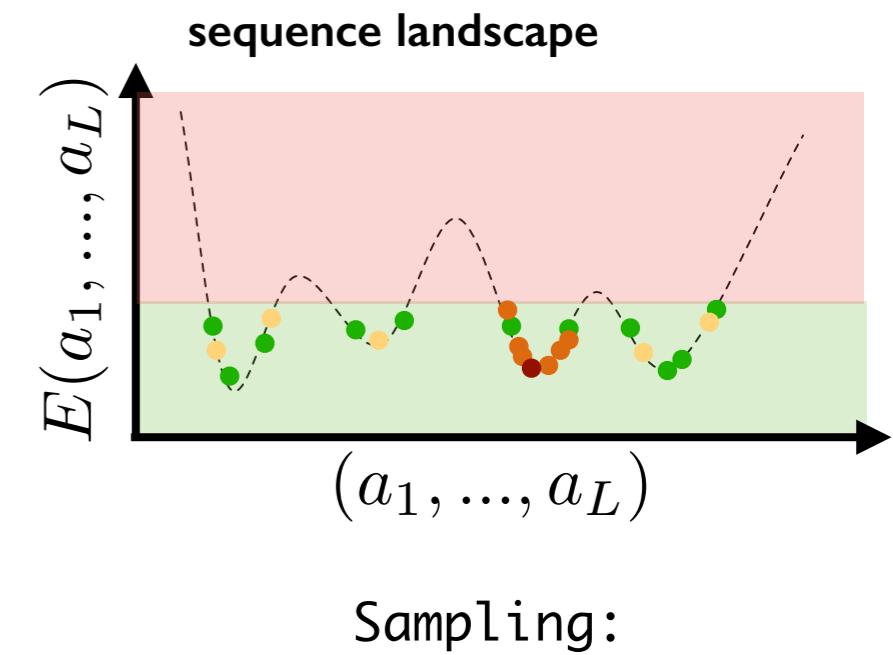
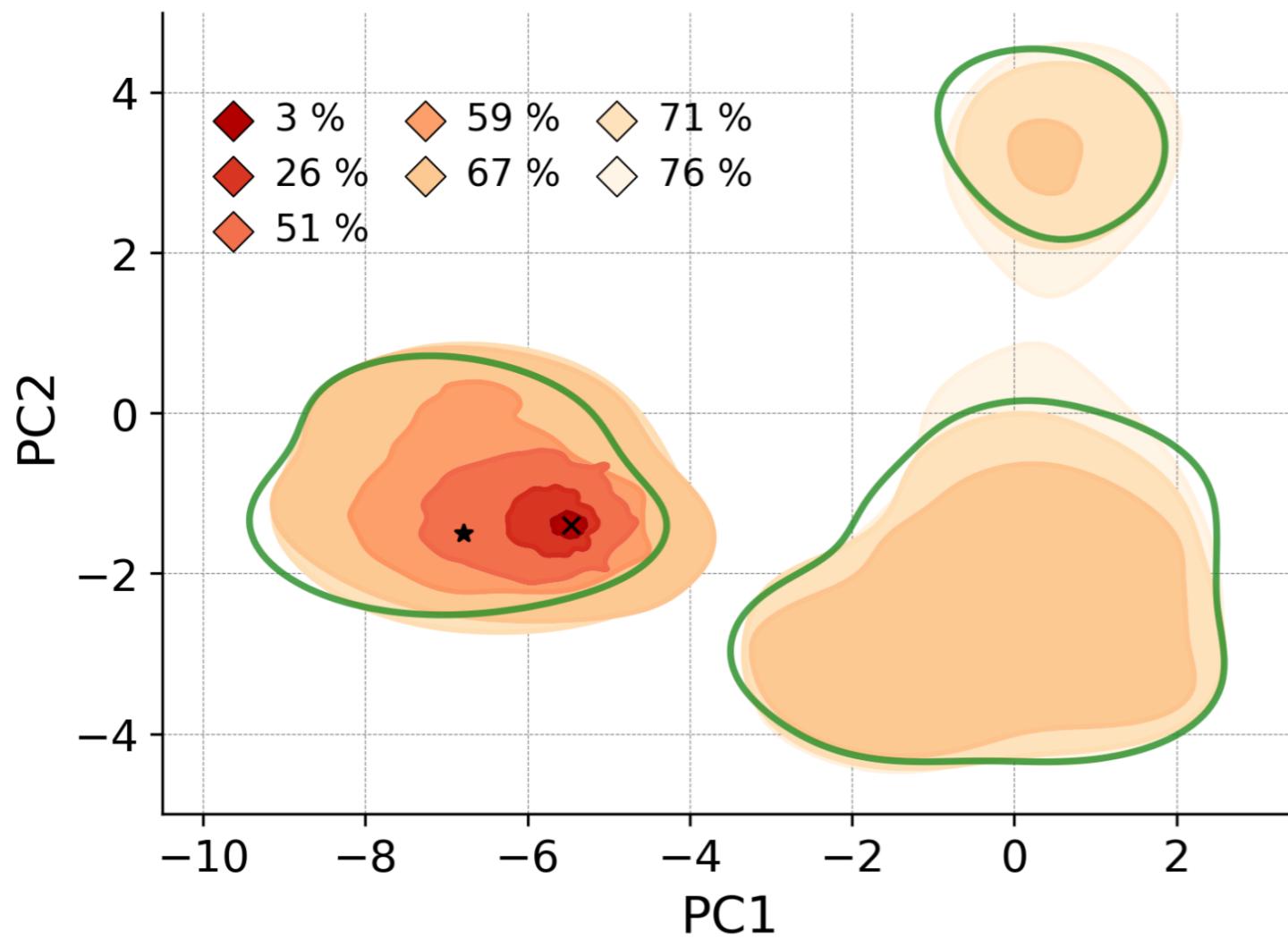
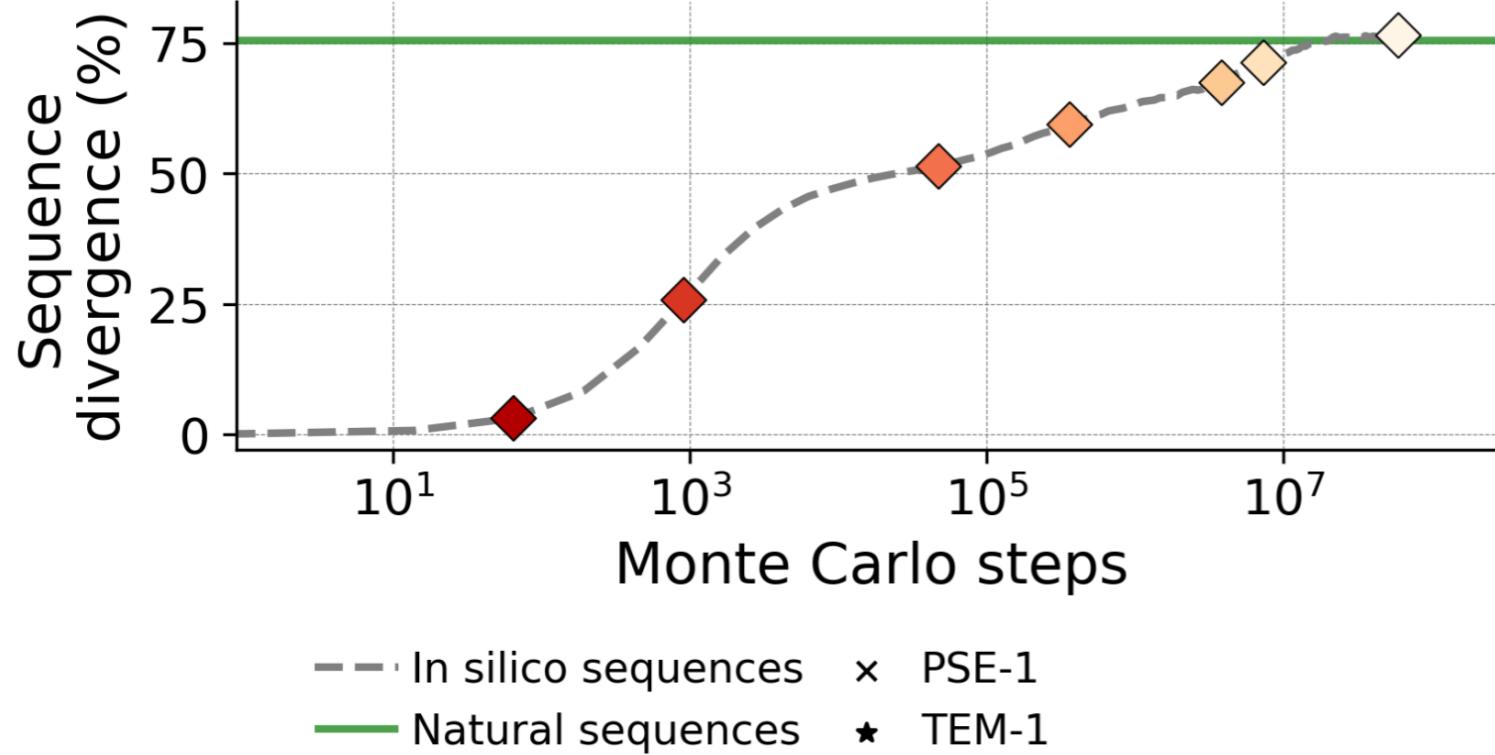


Predict ~400 numbers with 3 fitted parameters

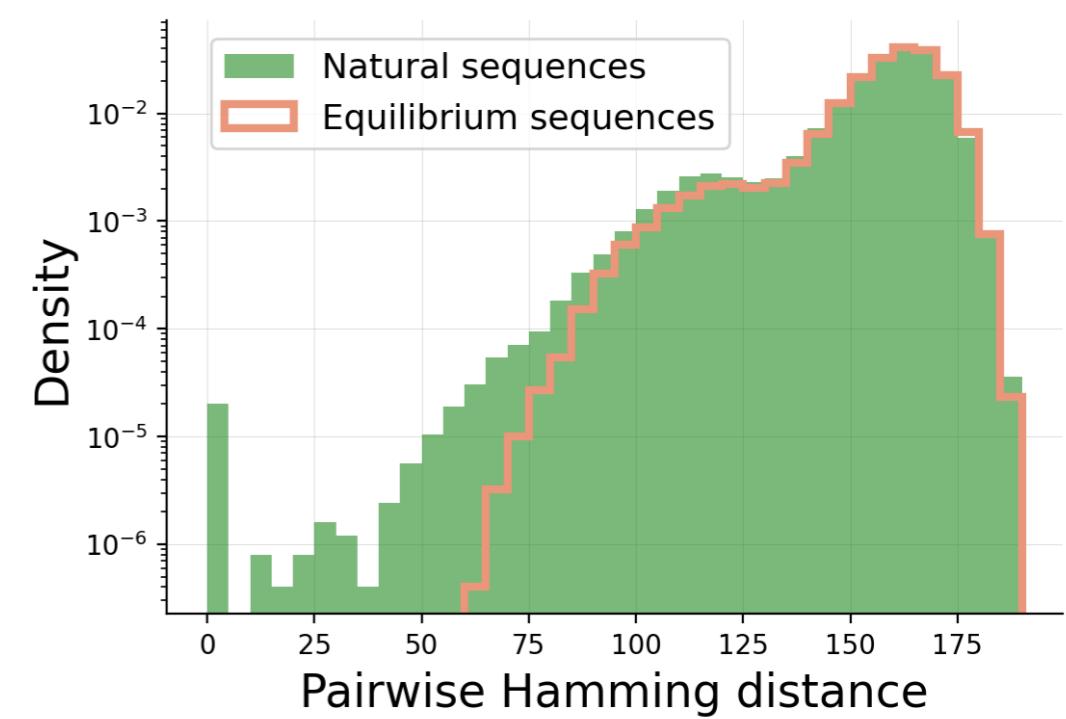
Coevolutionary couplings are needed



Long-term evolution



- Start from a wt sequence
- Sample 1000 trajectories
- Save the MSA at various times



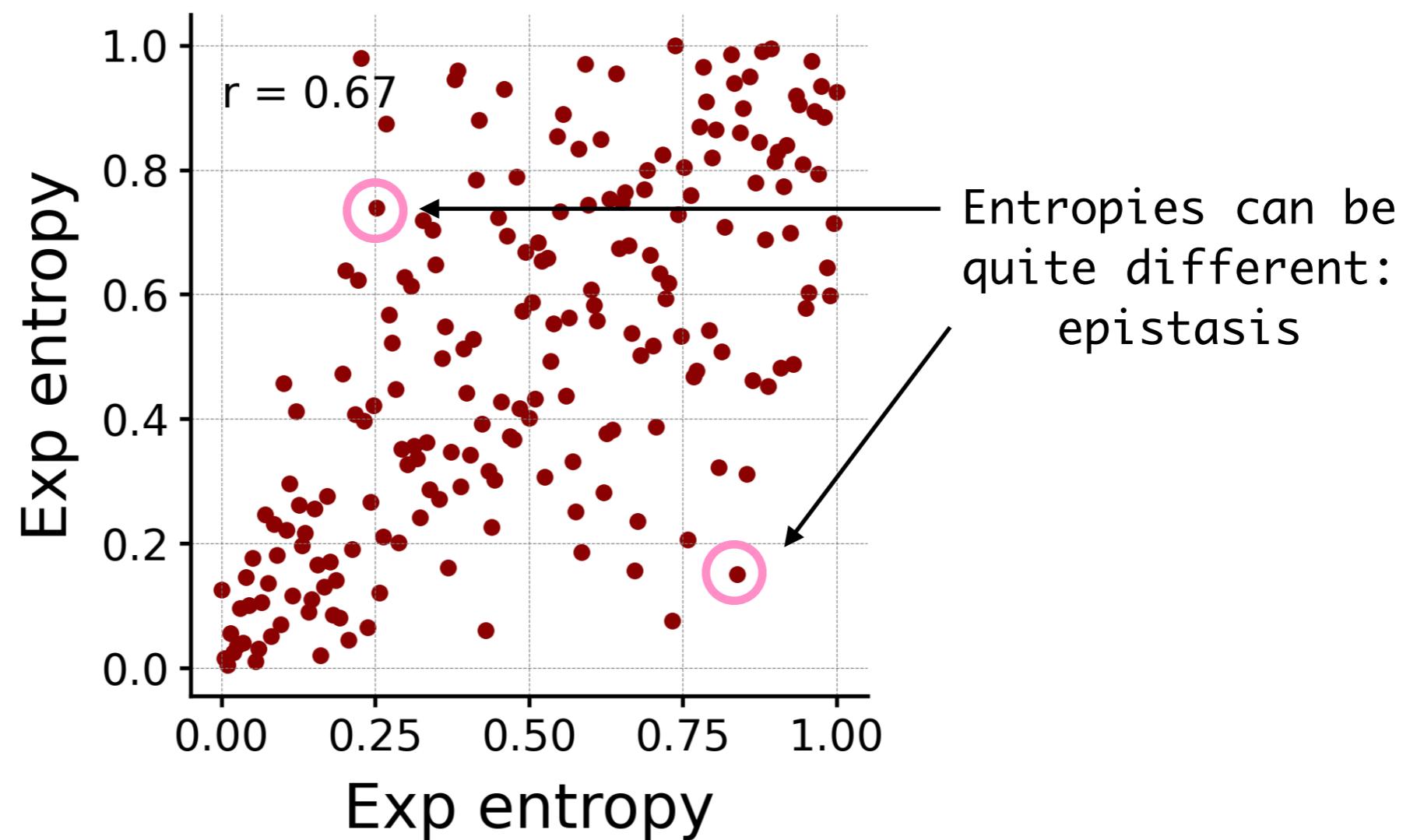
Take home:

Monte Carlo dynamics of generative models
can accurately describe in vitro evolution (short time, dynamics)

Epistasis

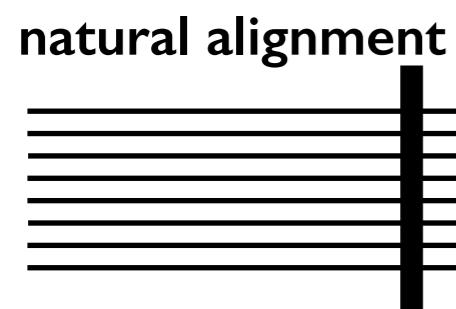
Epistasis

Context dependence of single mutations
(related to coevolution)



Quantifying epistasis

Mutability



CIE

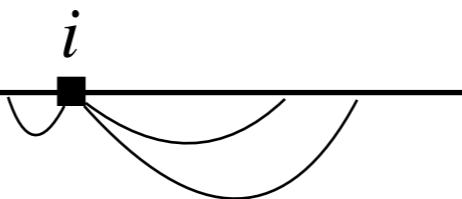
Context-INdependent



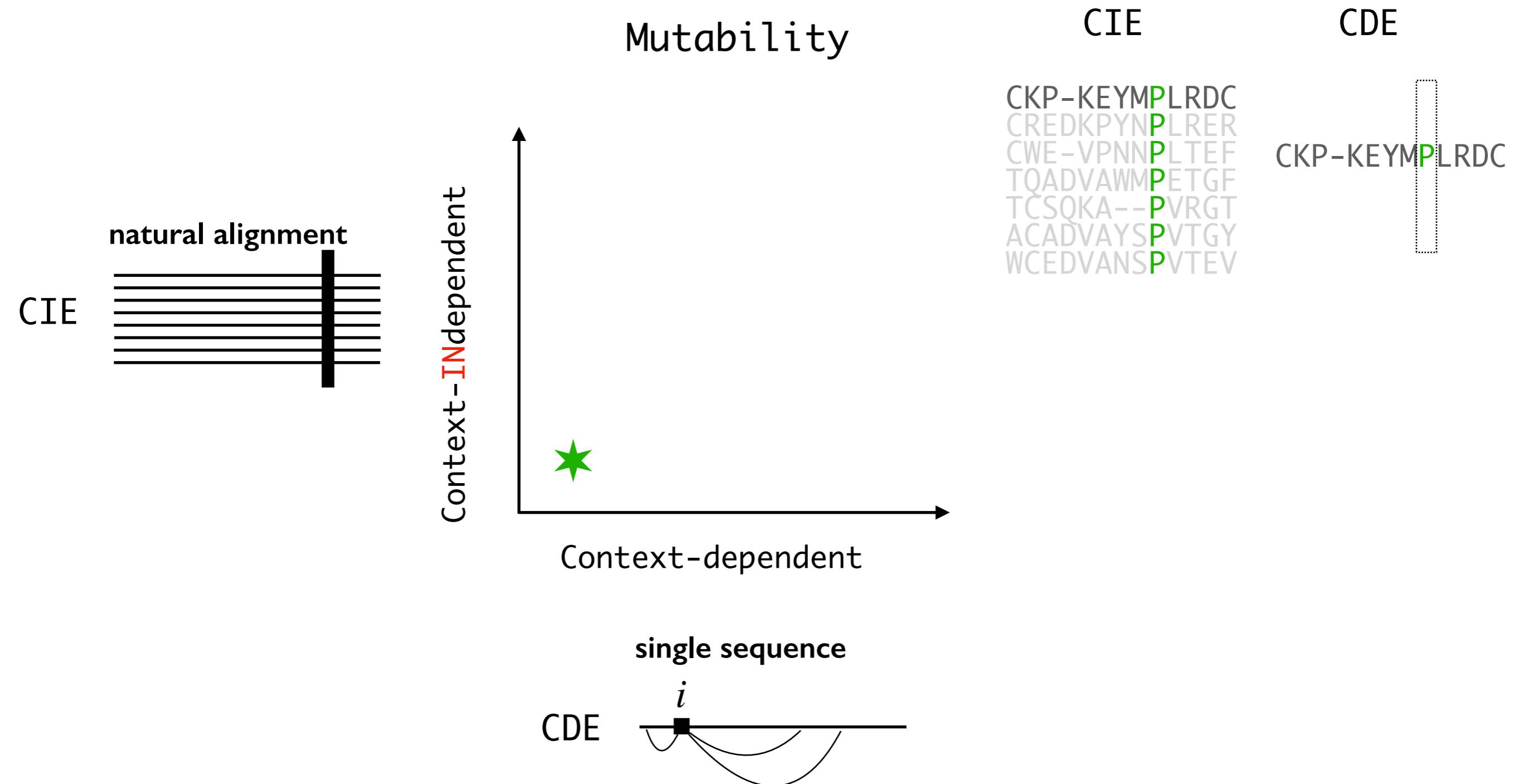
Context-dependent

single sequence

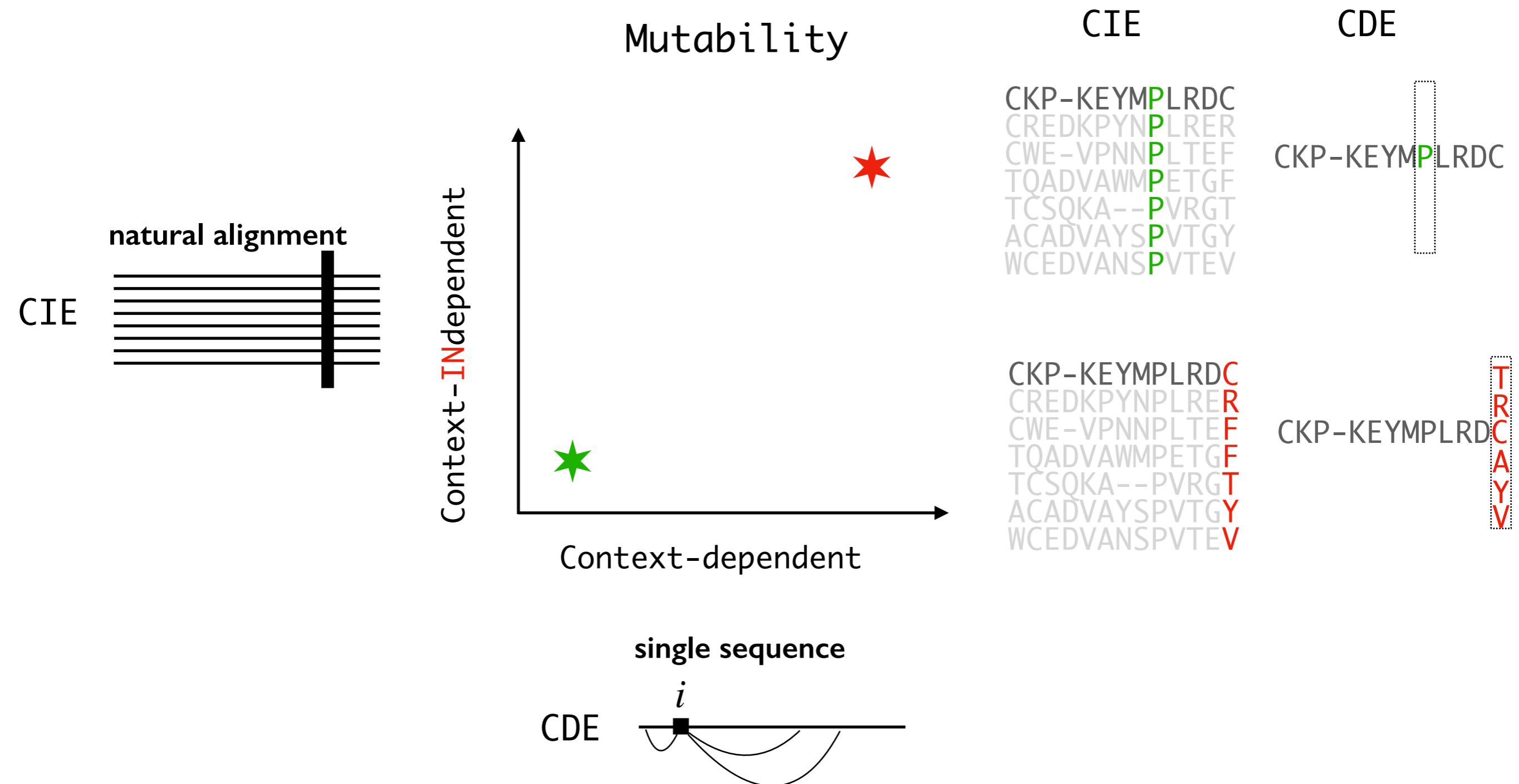
CDE



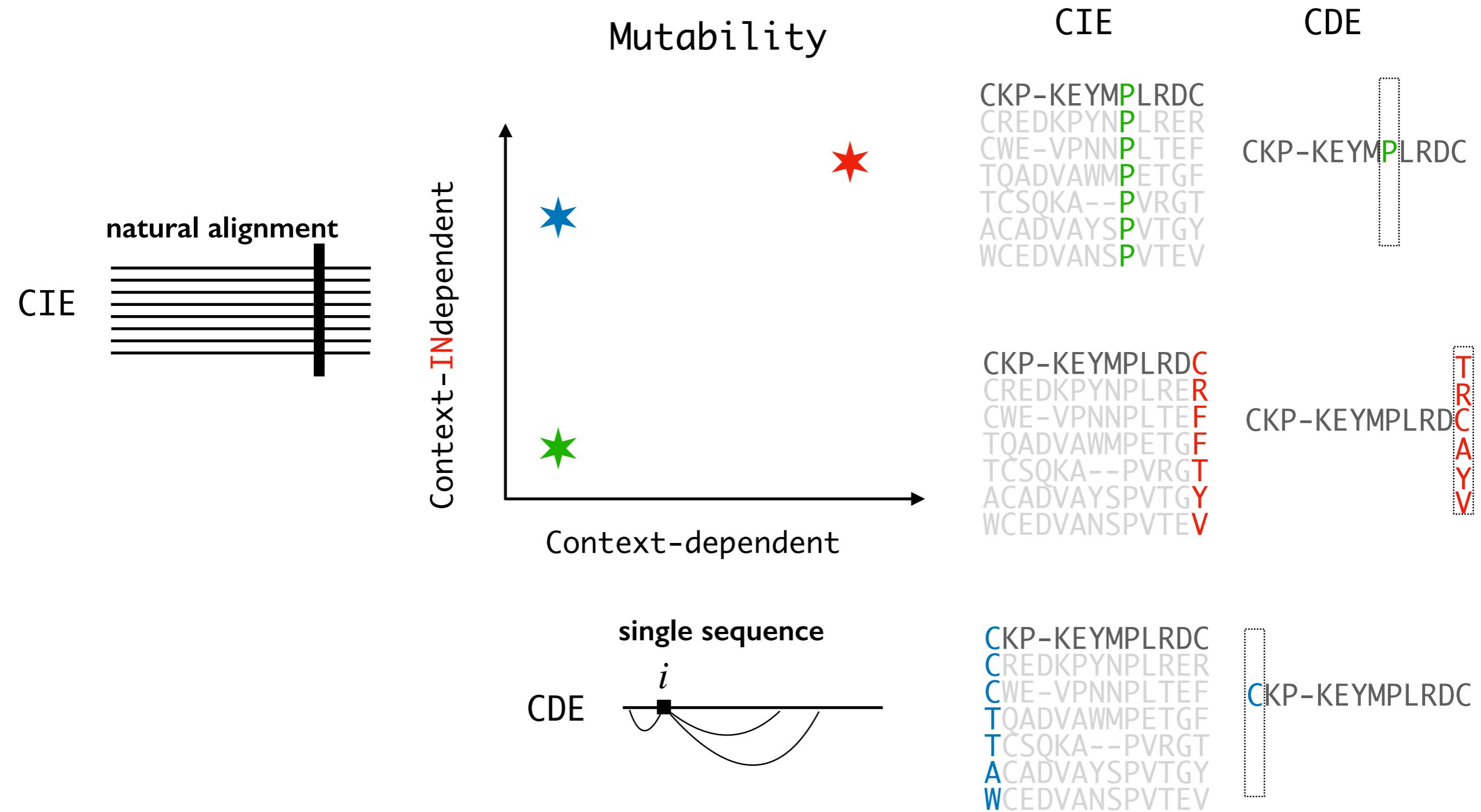
Quantifying epistasis



Quantifying epistasis



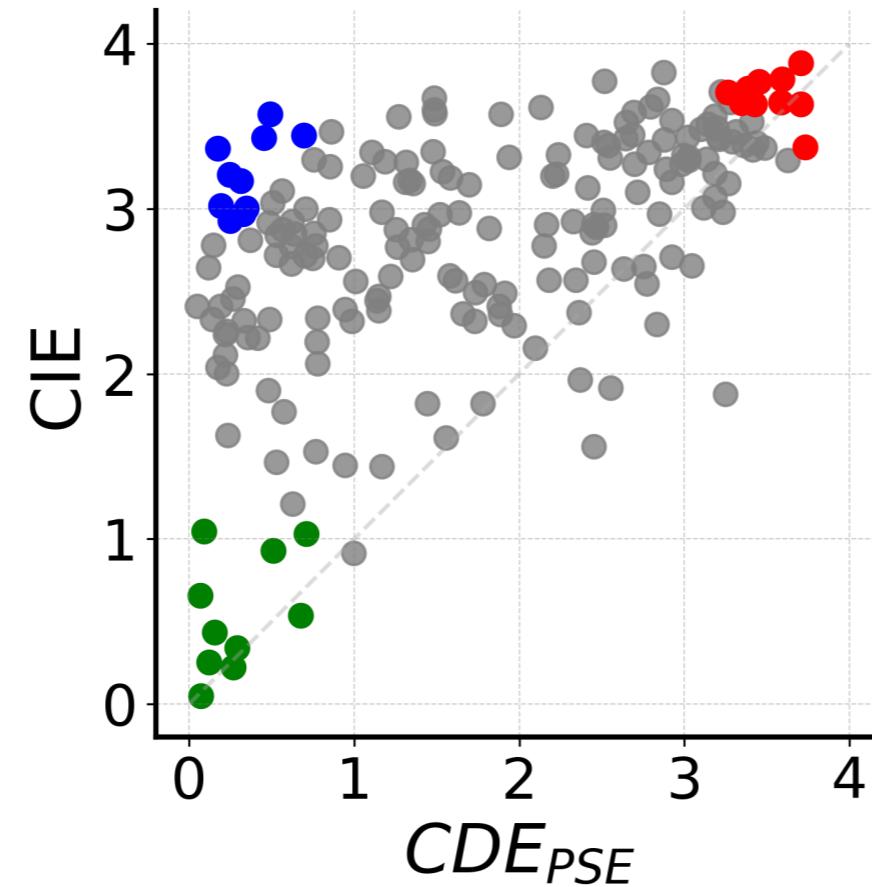
Quantifying epistasis



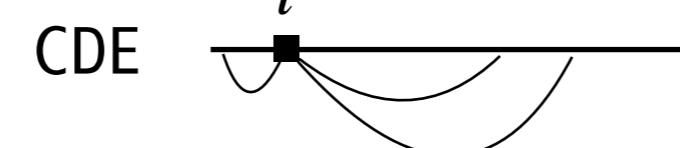
Quantifying epistasis

natural alignment

CIE



single sequence



CIE

CKP-KEYMPLRDC
CREDKPYNPLRER
CWE-VPNNPLTEF
TQADVAWMPETGF
TCSQKA--PVRGT
ACADVAYSPVTGY
WCEDVANSPVTEV

CDE

CKP-KEYMPLRDC



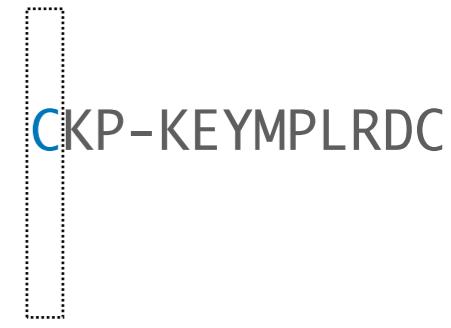
CKP-KEYMPLRDC
CREDKPYNPLRE
CWE-VPNNPLTE
TQADVAWMPETG
TCSQKA--PVRGT
ACADVAYSPVTGY
WCEDVANSPVTEV

CKP-KEYMPLRD

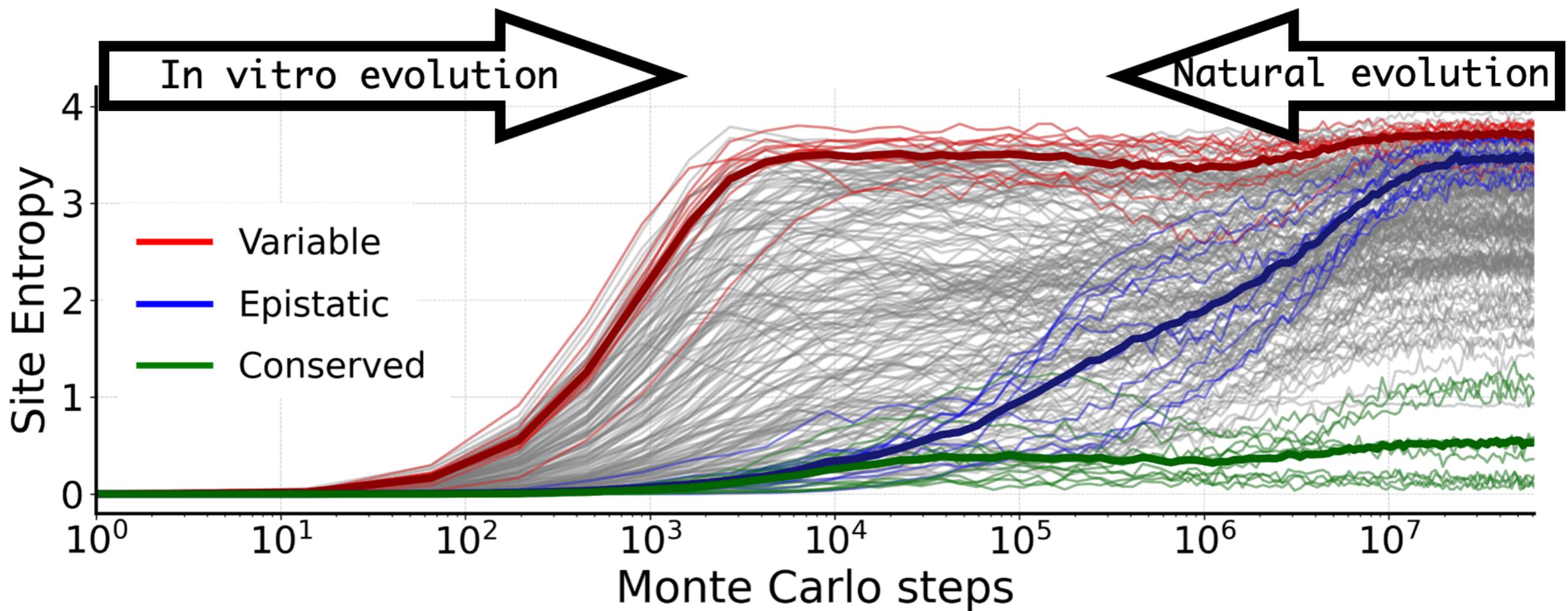
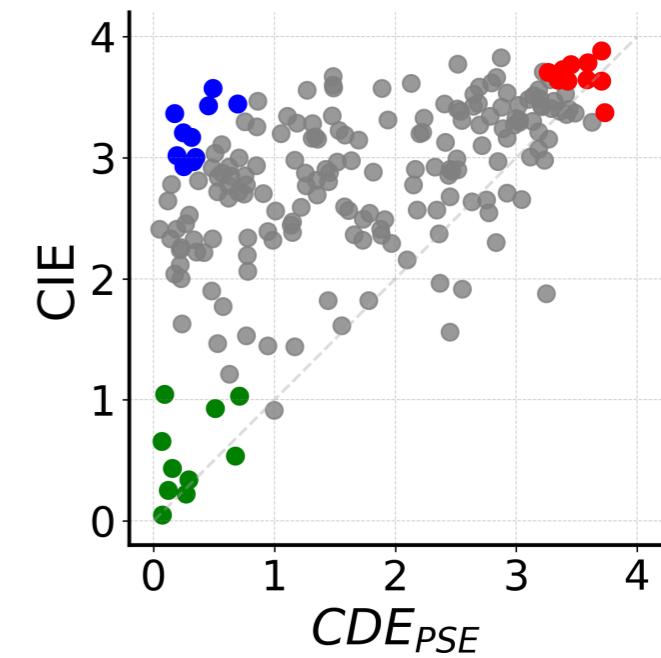
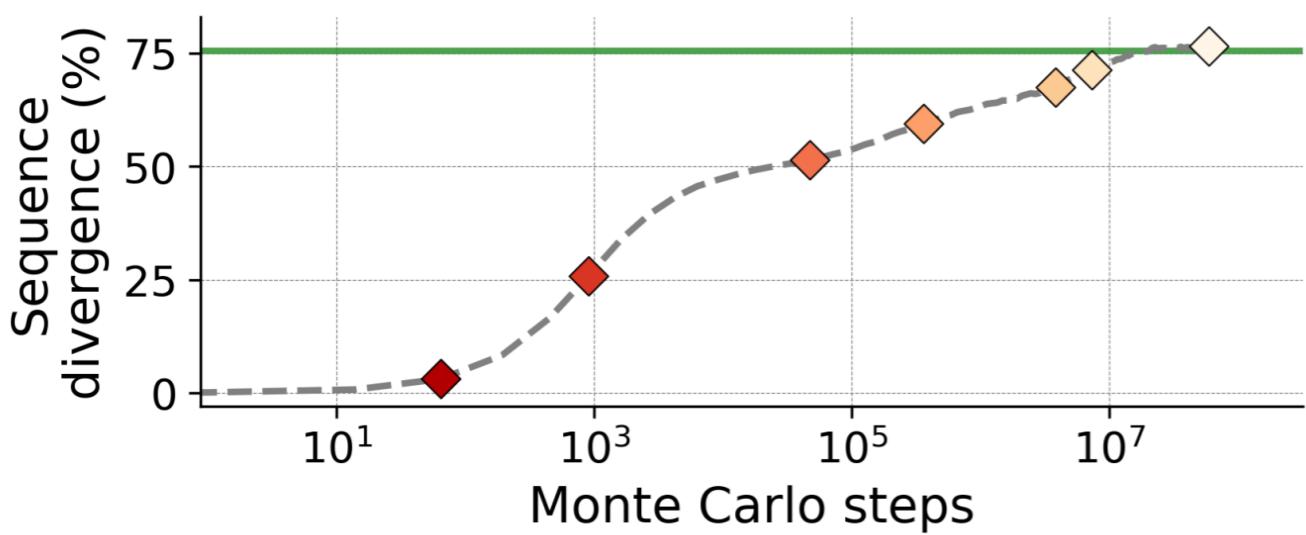


CKP-KEYMPLRDC
CREDKPYNPLRE
CWE-VPNNPLTE
TQADVAWMPETG
TCSQKA--PVRGT
ACADVAYSPVTGY
WCEDVANSPVTEV

CKP-KEYMPLRD



Emergence of epistasis



Take home:

Epistasis is very heterogeneous across sites

It emerges on intermediate time scales due to collective effects

Measures of heterogeneity



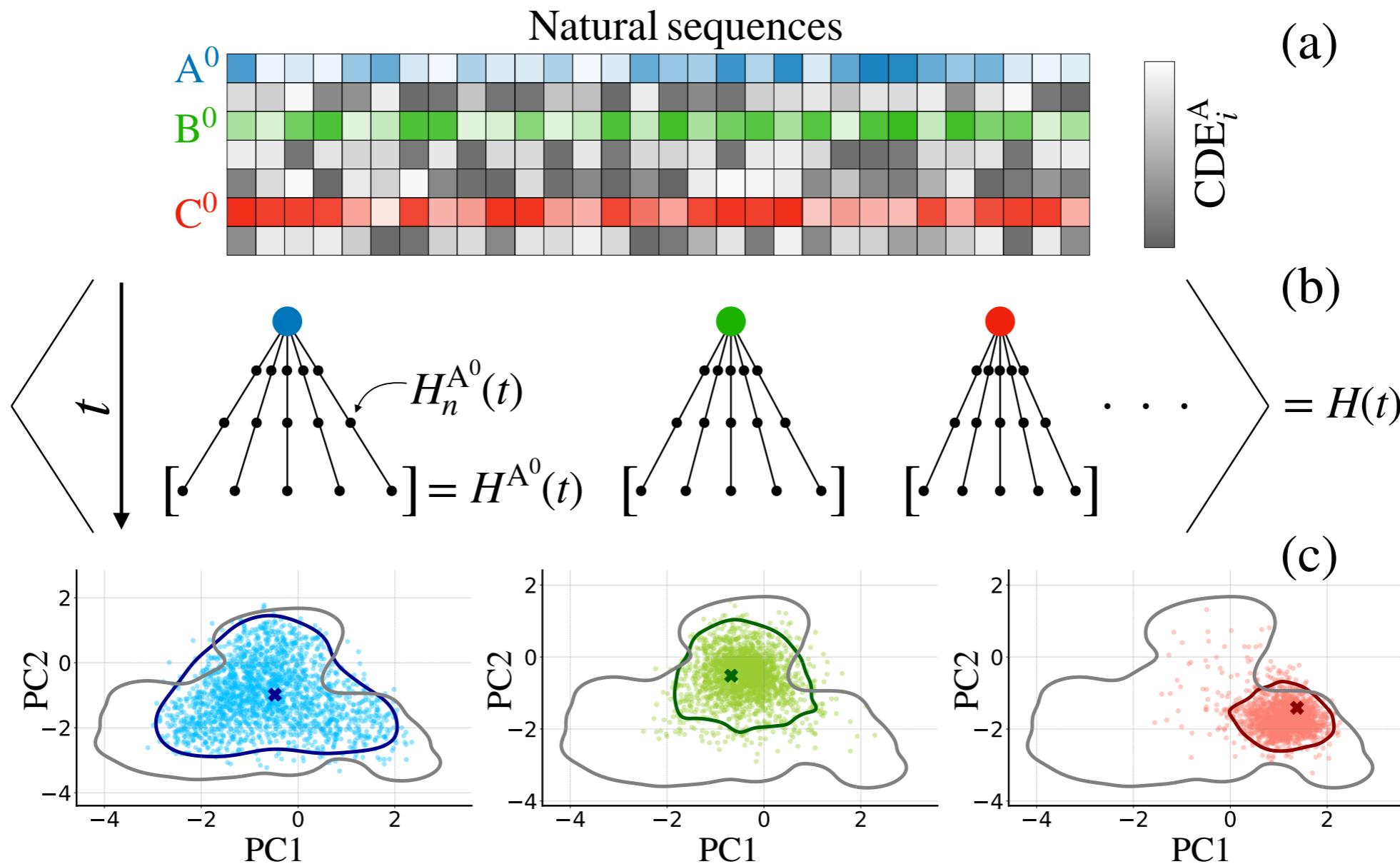
Saverio
Rossi



Leonardo
Di Bari

Rossi, Di Bari, Weigt, FZ, arXiv:2412.01969 (2024)

Different sources of fluctuations



H = number of mutations (Hamming dist.)

$$\chi_{dyn} = \langle \chi_{dyn}^A \rangle = \langle [H^2] - [H]^2 \rangle$$

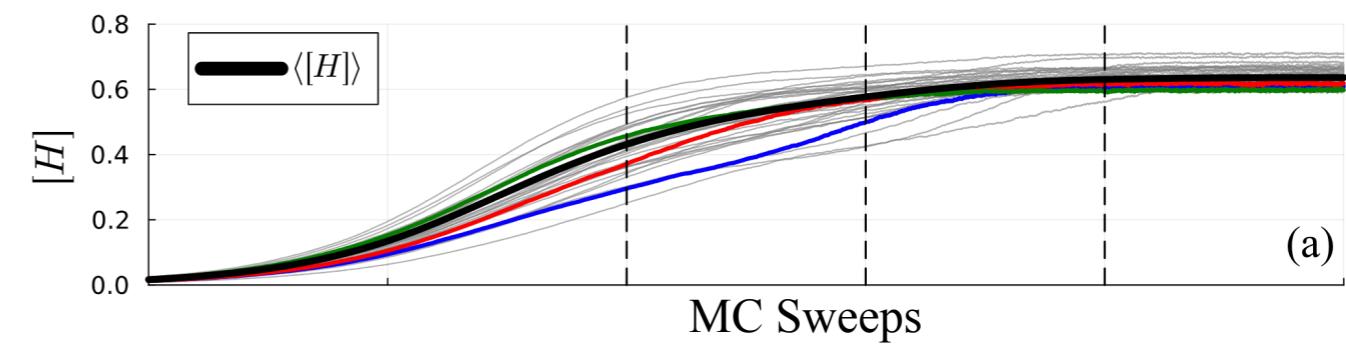
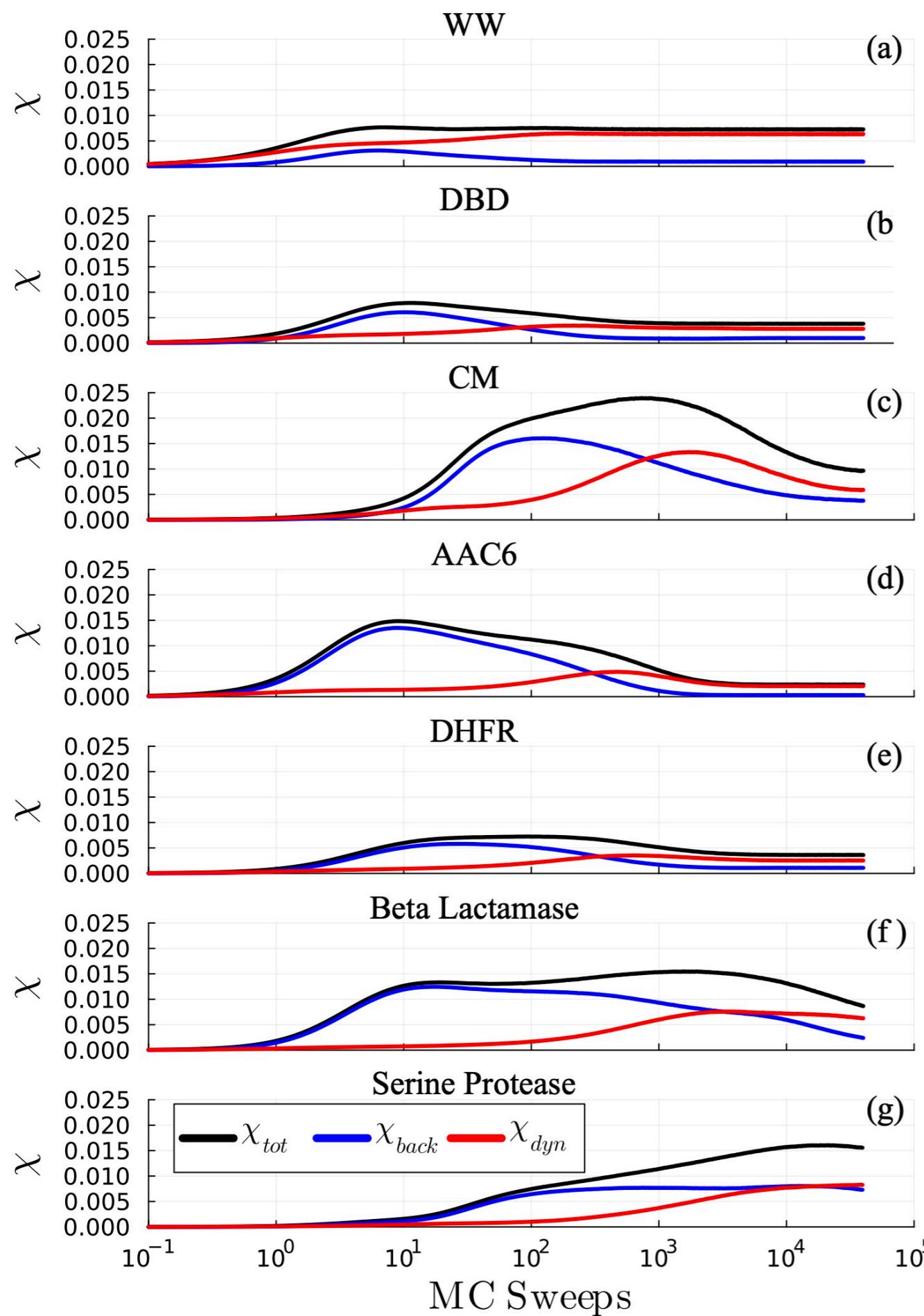
[•] over dynamical fluctuations

$$\chi_{back} = \langle [H]^2 \rangle - \langle [H] \rangle^2$$

$\langle \bullet \rangle$ over initial sequence A

$$\chi_{tot} = \langle [H^2] \rangle - \langle [H] \rangle^2 = \chi_{back} + \chi_{dyn}$$

Different sources of fluctuations



Dynamical fluctuations subdominant at short/intermediate times

They only take over at very long times

Evolution is partially predictable from initial condition

Different sources of fluctuations

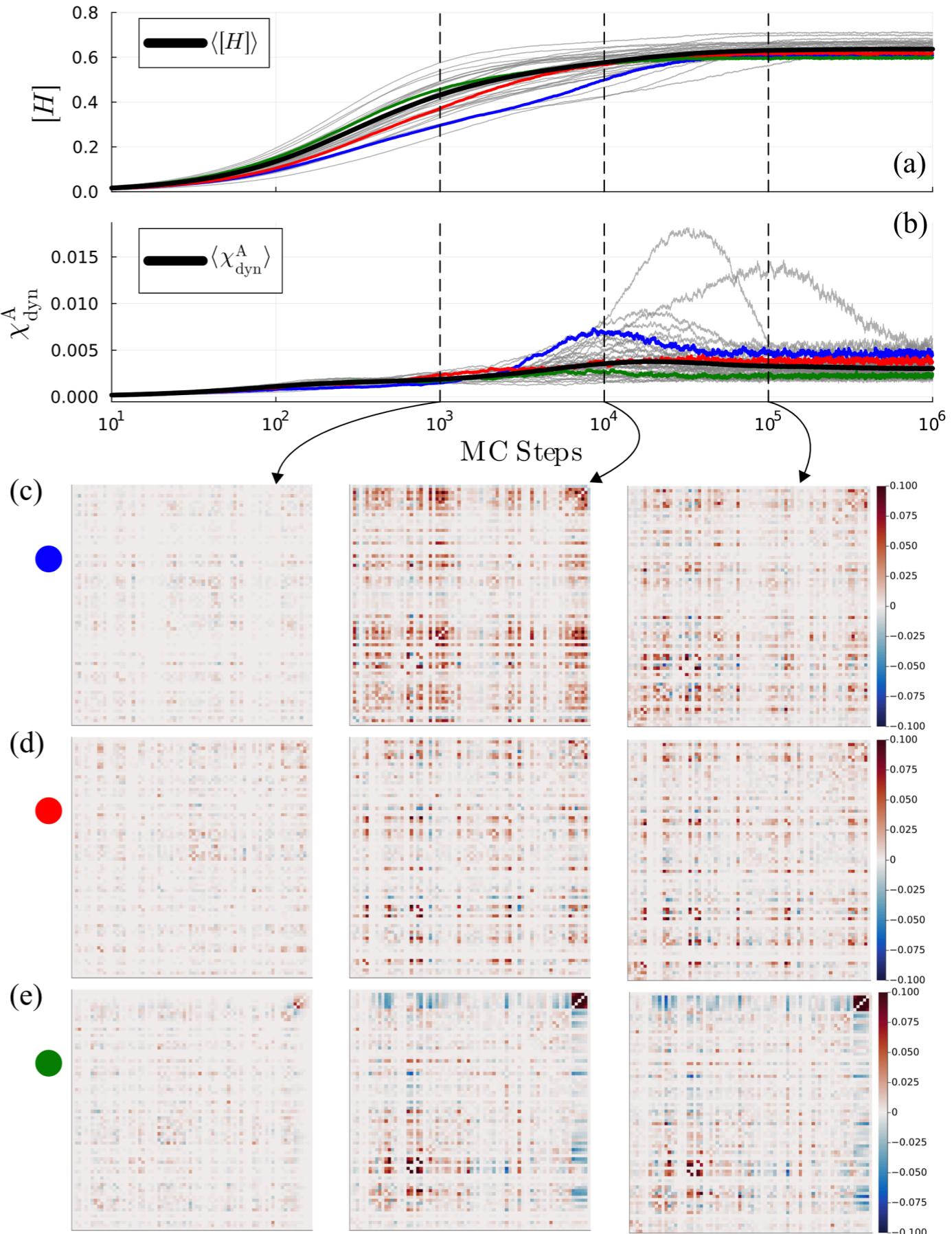
$$\chi_{dyn}^A(t,0) = \frac{1}{L^2} \sum_{ij} G_{ij,dyn}^A(t,0)$$

$$G_{ij,dyn}^A(t,0) = [\delta_{a_i^t, a_i^0} \delta_{a_j^t, a_j^0}] - [\delta_{a_i^t, a_i^0}] [\delta_{a_j^t, a_j^0}]$$

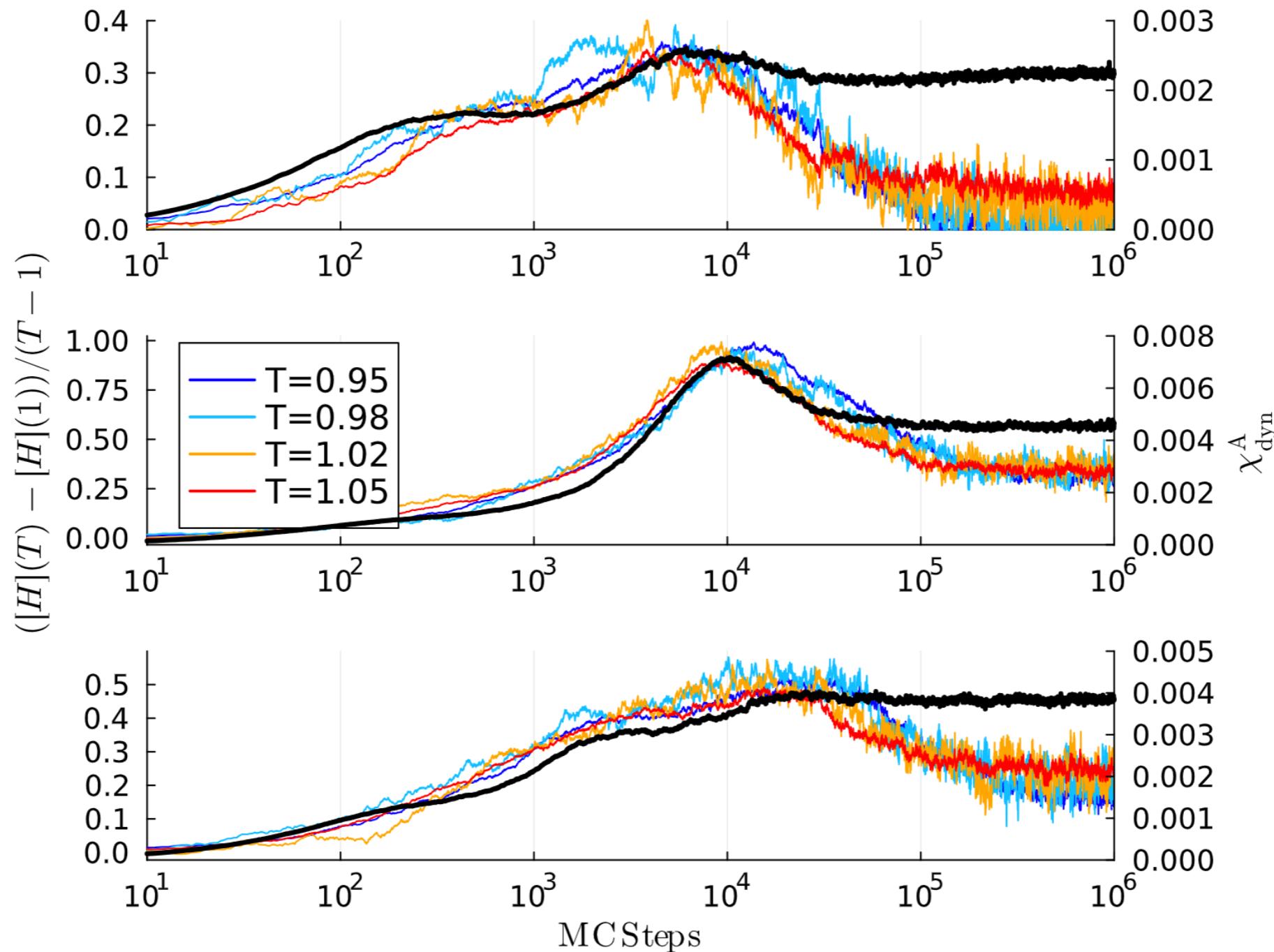
Clusters of strongly correlated sites

Depending on the initial sequence!

Sequence-specific dynamical ‘sectors’?



Fluctuation-dissipation relation



$$\chi_{dyn}^A(t,0) \approx \frac{d[H(t,0)]_A}{dT}$$

Berthier et al. Science (2005)

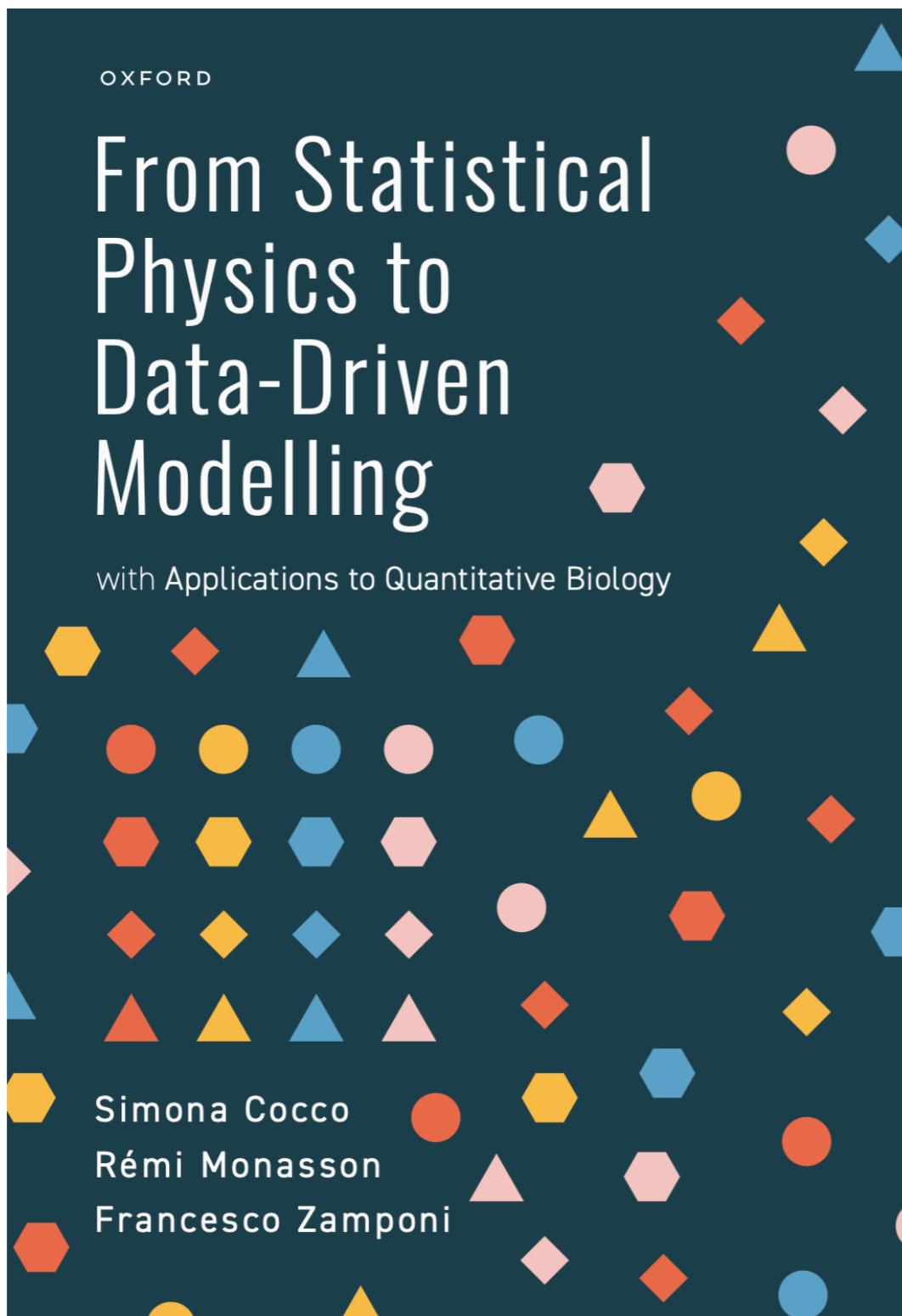
More heterogeneous sequences are more sensitive to perturbations
 (remember: T = selection strength)

Take home:

Epistasis is very heterogeneous across sequences

We can quantify it using dynamical correlations and identify cooperatively rearranging sites

Advertisement





SAPIENZA
UNIVERSITÀ DI ROMA

Thank you for your attention!

Bisardi, Rodriguez-Rivas, FZ, Weigt, Mol.Bio.Evo. (2022)

Di Bari, Bisardi, Cotogno, Weigt, FZ, PNAS (2024)

Rossi, Di Bari, Weigt, FZ, arXiv:2412.01969 (2024)

