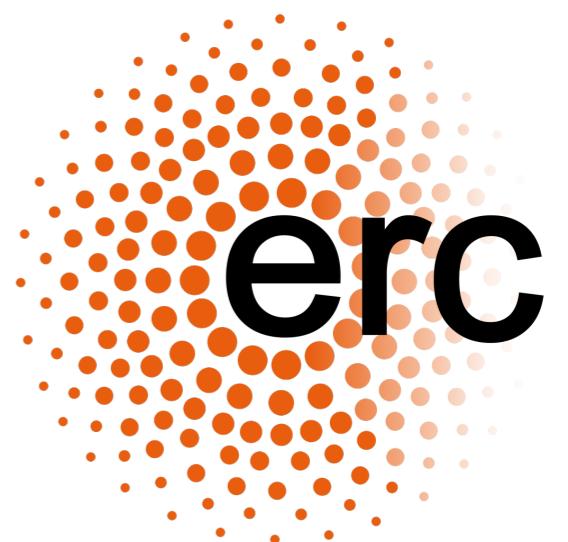


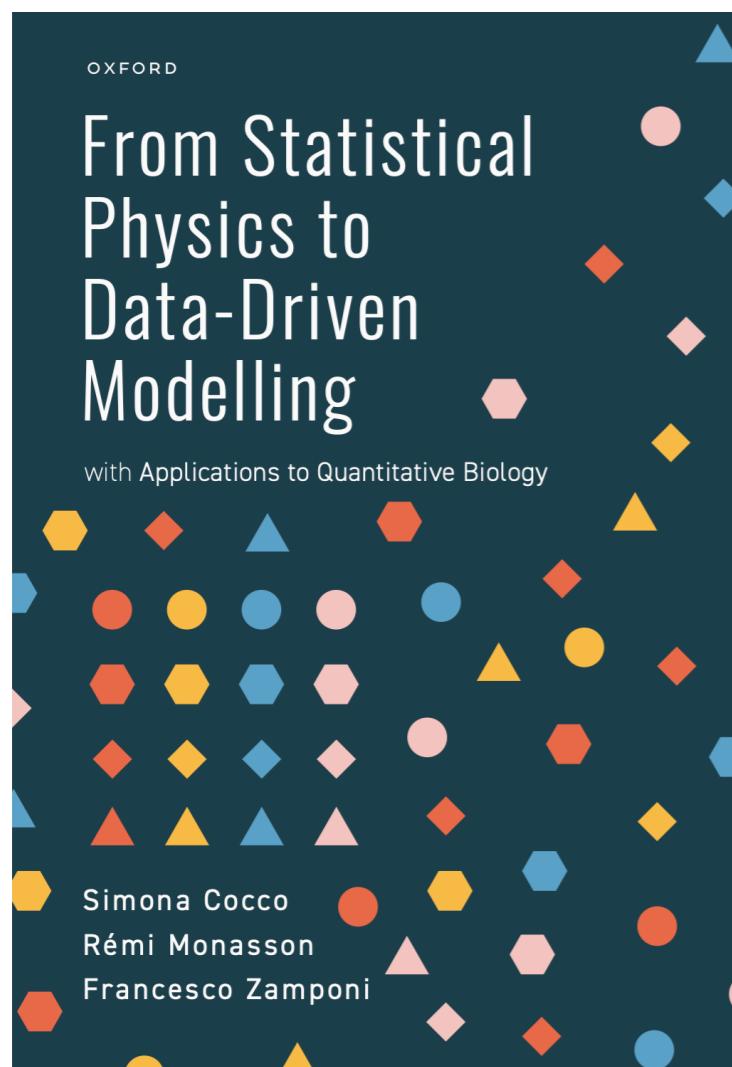
Rough fitness landscapes: from protein evolution to protein design

Francesco Zamponi

Gulliver, ESPCI, October 17, 2022

SIMONS FOUNDATION





Background

2016-2021: teaching with S.Cocco and R.Monasson

- Bayesian and high-dimensional inference
- Regularization
- Graphical models
- Supervised and unsupervised learning
- Time series analysis

With coding tutorials using real biological data

Oxford University Press, 2022

Promotion code: ASPROMP8



SORBONNE
UNIVERSITÉ



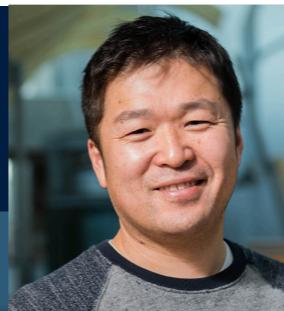
Since 2017: collaboration with Martin Weigt's group

Since 2021: collaboration with Nobuhiko Tokuriki's lab

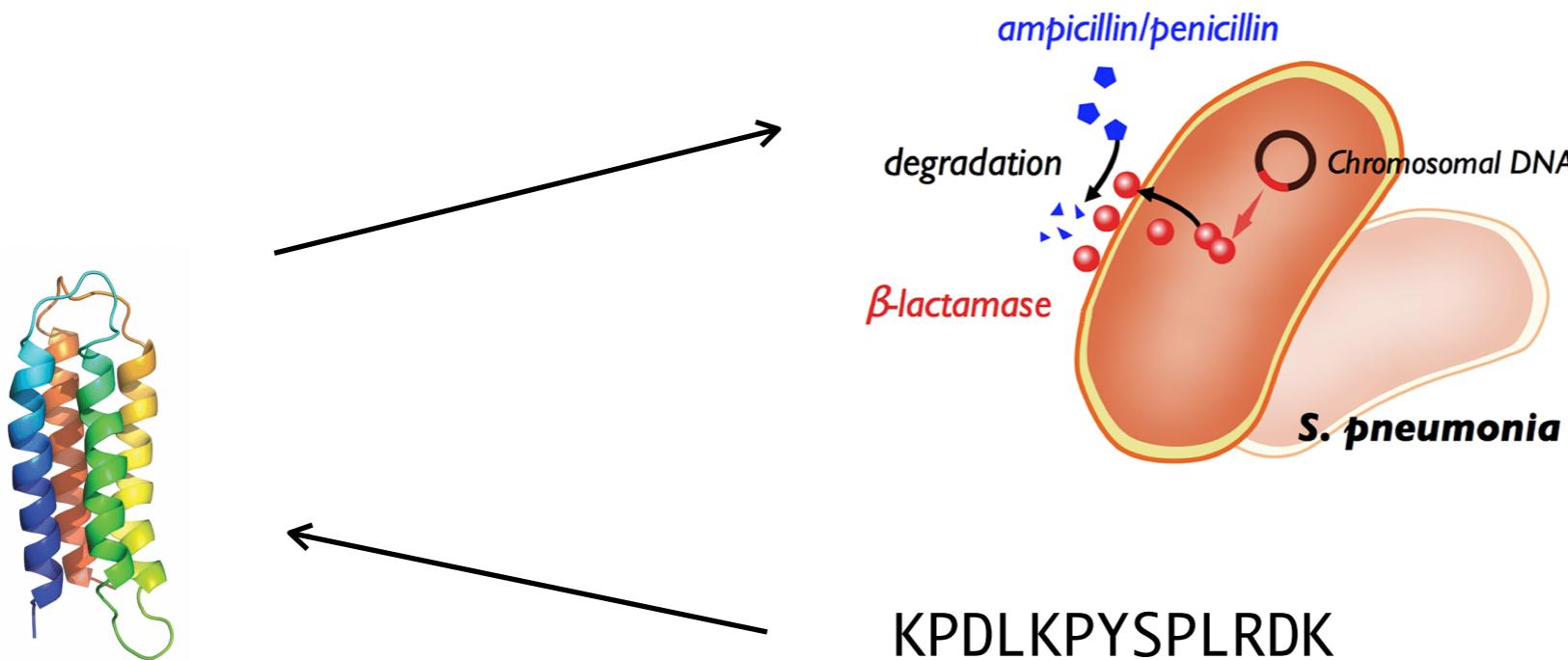


THE UNIVERSITY OF BRITISH COLUMBIA

Tokuriki Lab - Michael Smith Laboratories



Sequence-to-function paradigm

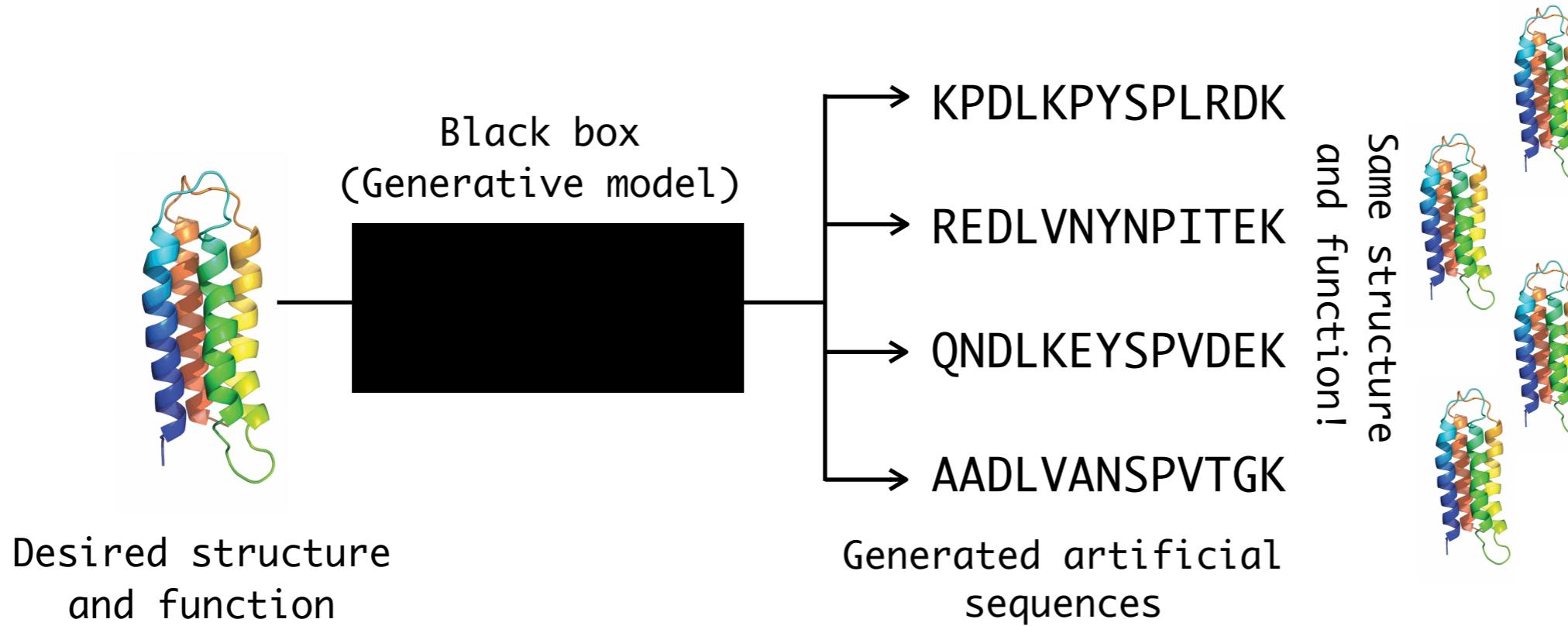


AI | MACHINE LEARNING & DATA SCIENCE | POPULAR | RESEARCH

DeepMind's AlphaFold2 Predicts Protein Structures with Atomic-Level Accuracy

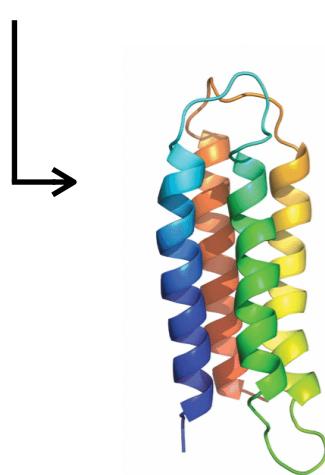
In a new paper published in the prestigious scientific journal Nature, DeepMind presents AlphaFold2, a redesigned neural-network system based on last year's AlphaFold that can predict protein structures with atomic-level accuracy.

Inverse problem: protein design

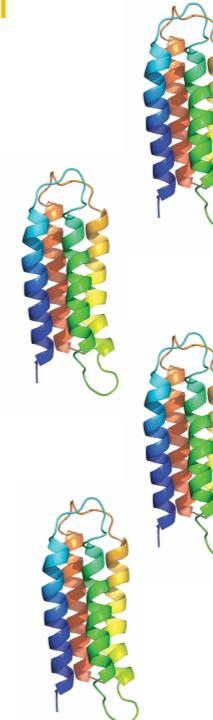
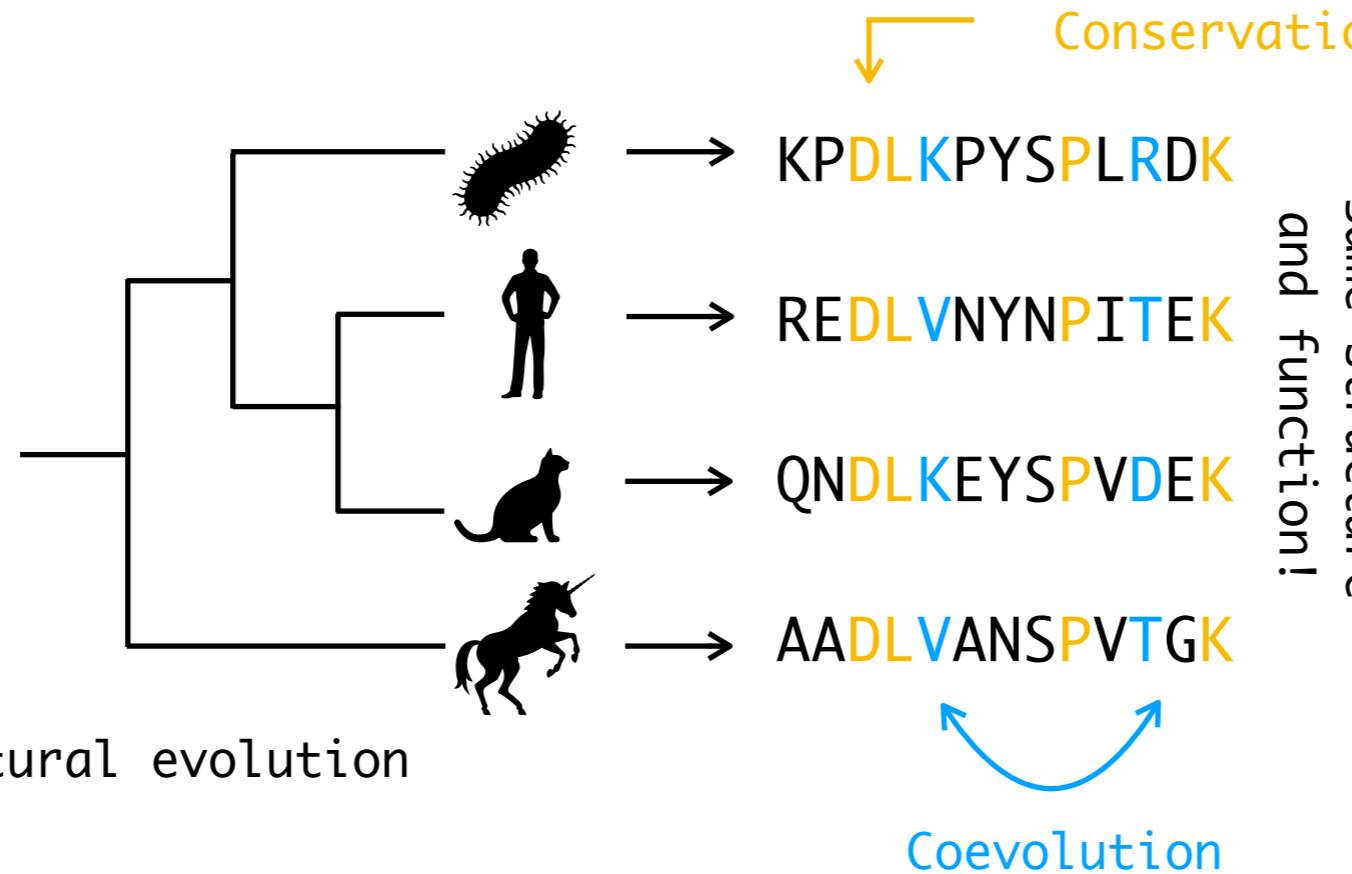


Data #1: natural protein sequences

Ancestral protein



Natural evolution

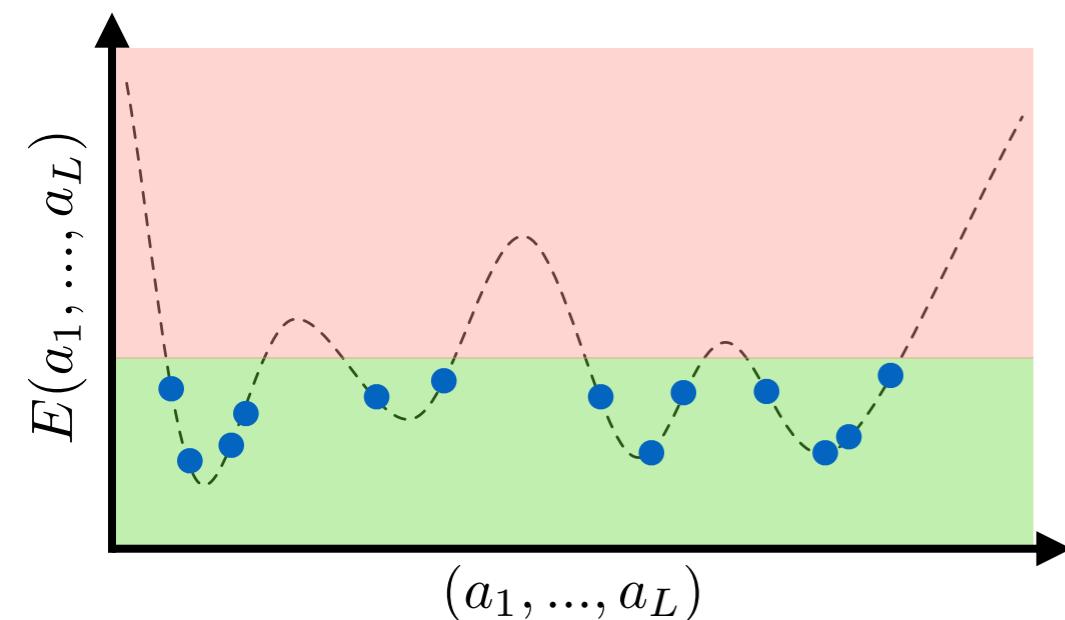


Available sequences:
 $10^2 - 10^5$

Sequence divergence:
70-80%

Public databases
(PFAM, UniProt...)

sequence landscape

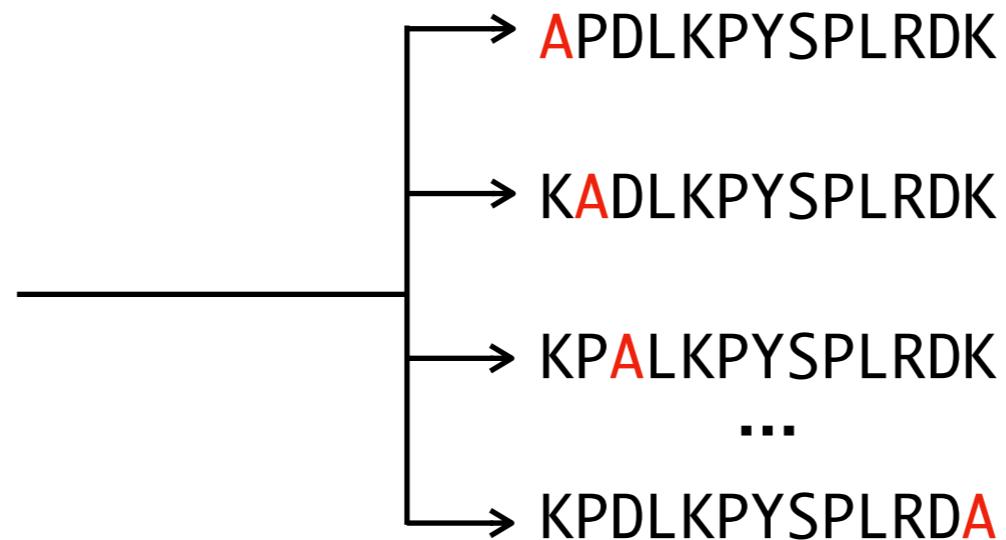
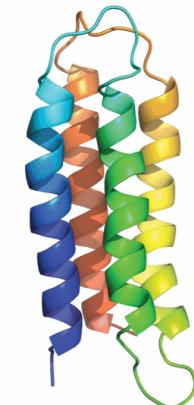


sequence data

global sample

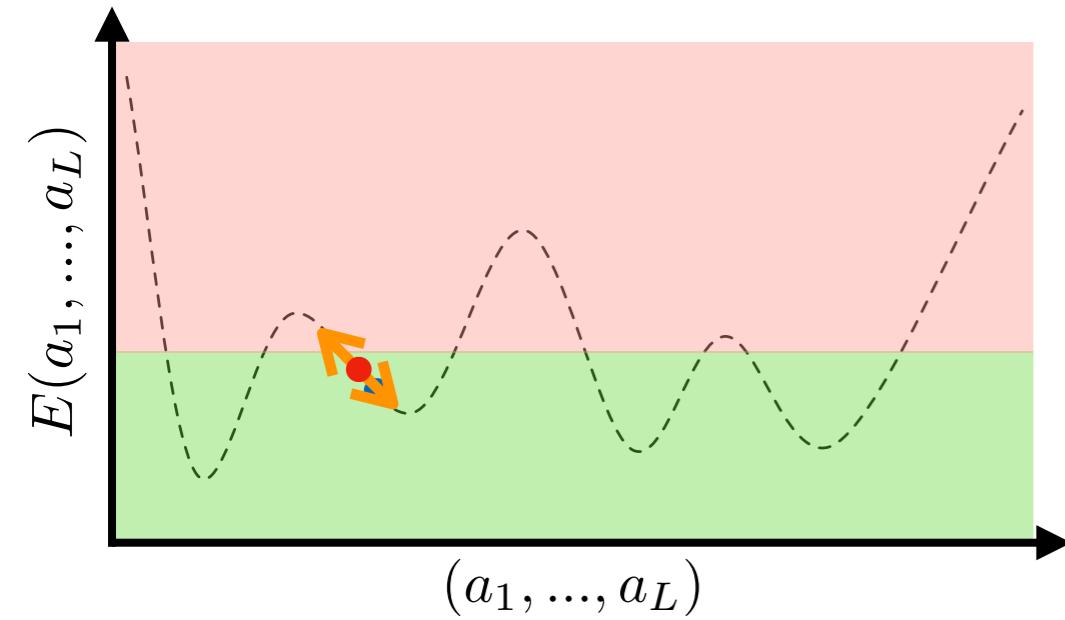
species 1
species 2
species 3
...

Data #2: deep mutational scan



Available sequences:
 $\sim 10^4$
 Sequence divergence:
 1

sequence landscape



sequence data

reference sequence

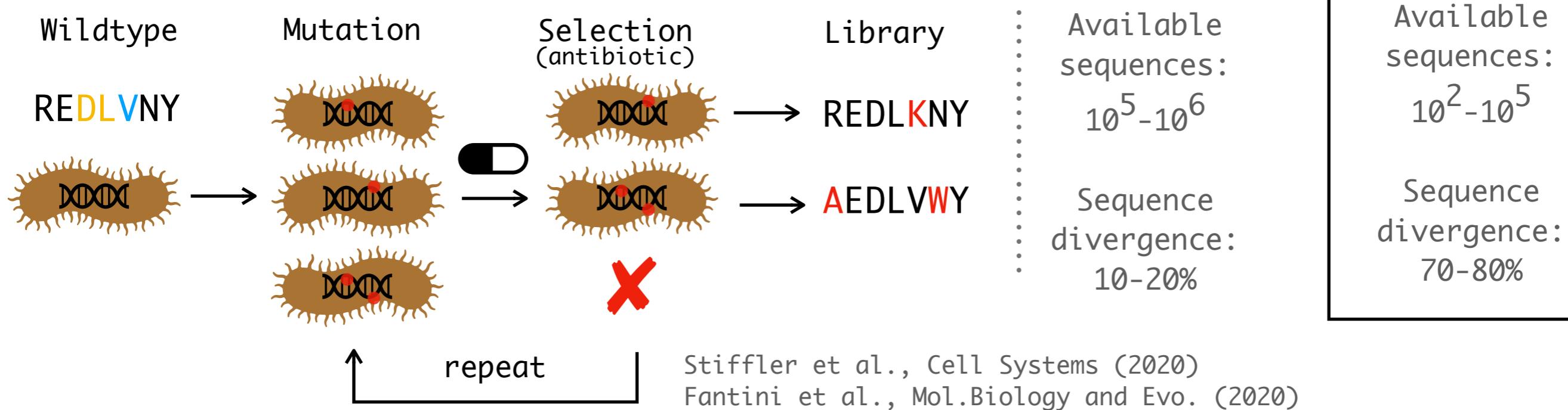
mutant sequences

- E₁
- E₂
- E₃
- E₄
- E₅
- E₆

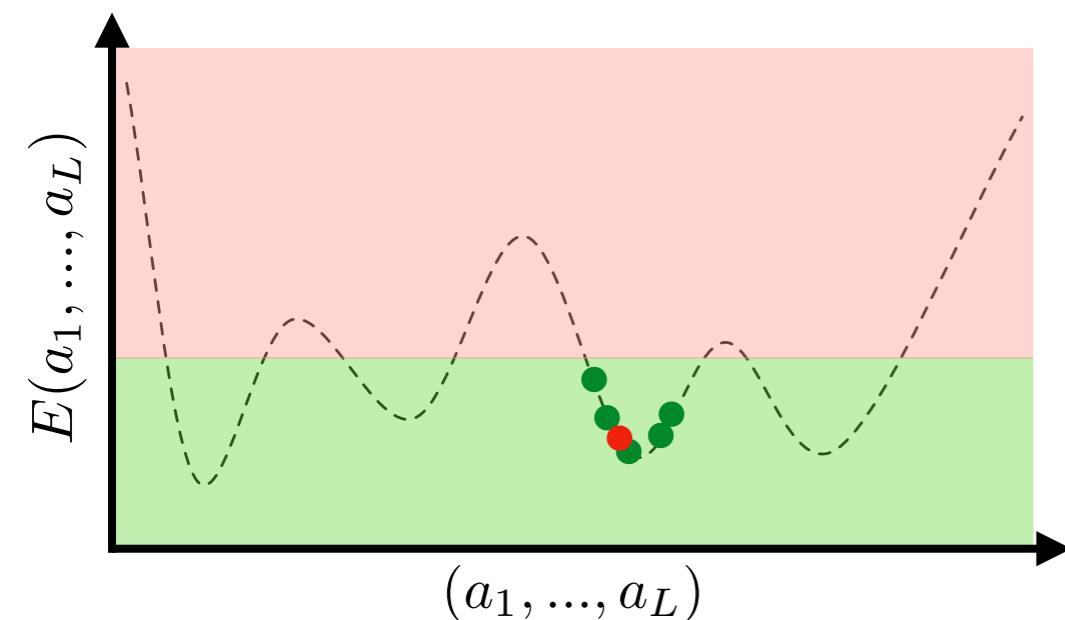
... phenotype

Data #3: in vitro evolution

Natural evolution



sequence landscape



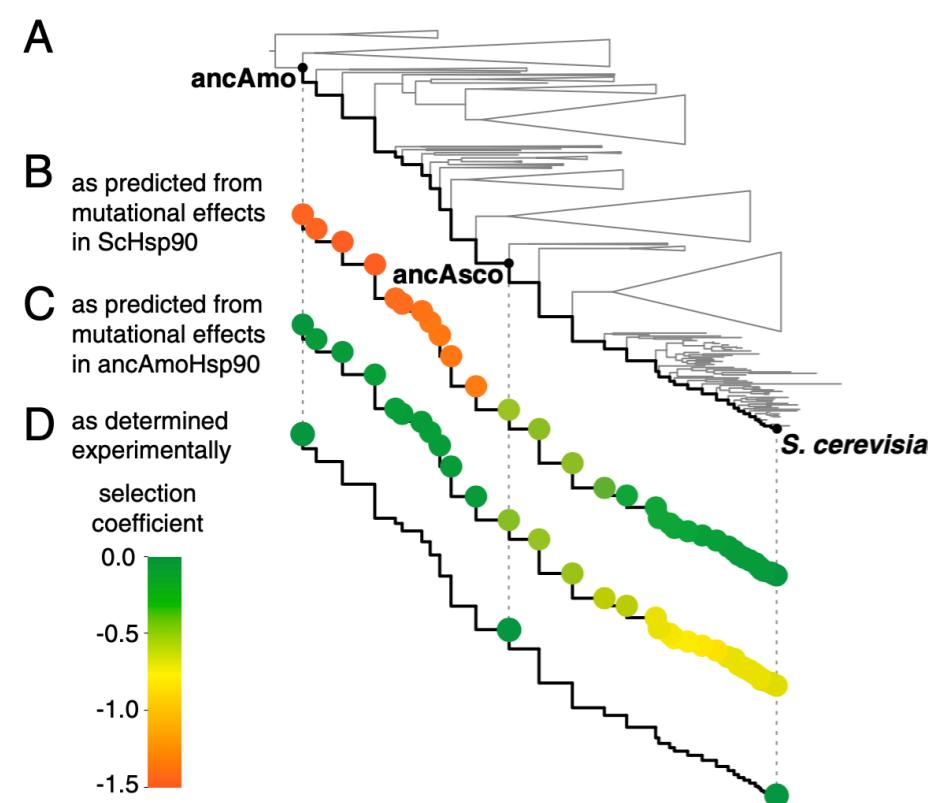
sequence data

reference sequence

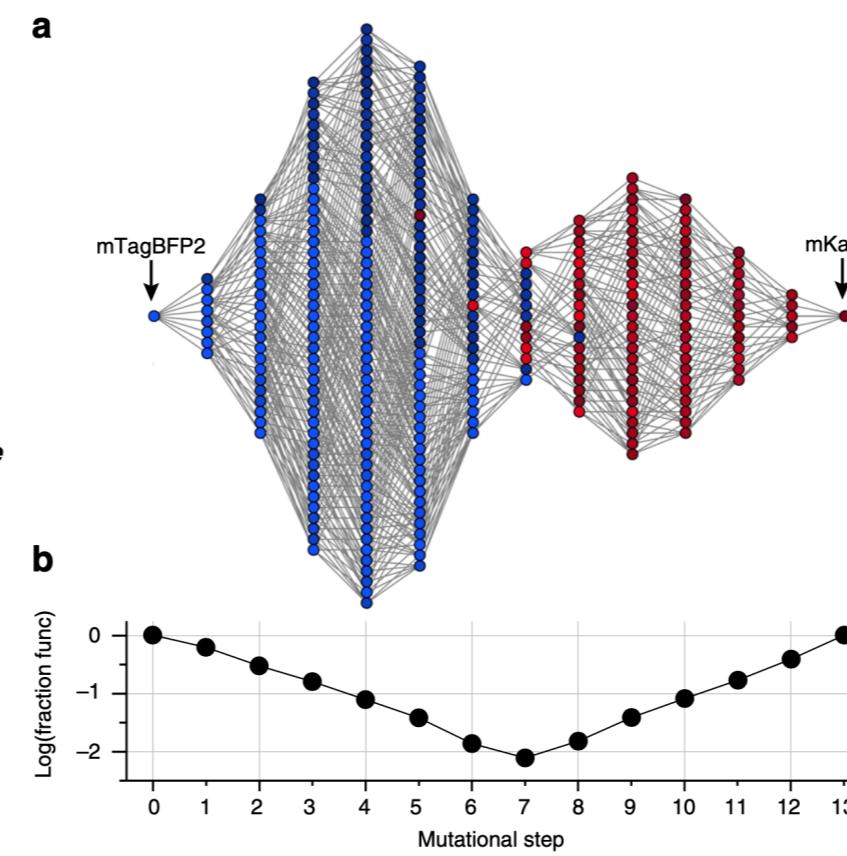
local sample

- mutant 1
- mutant 2
- mutant 3
- ...

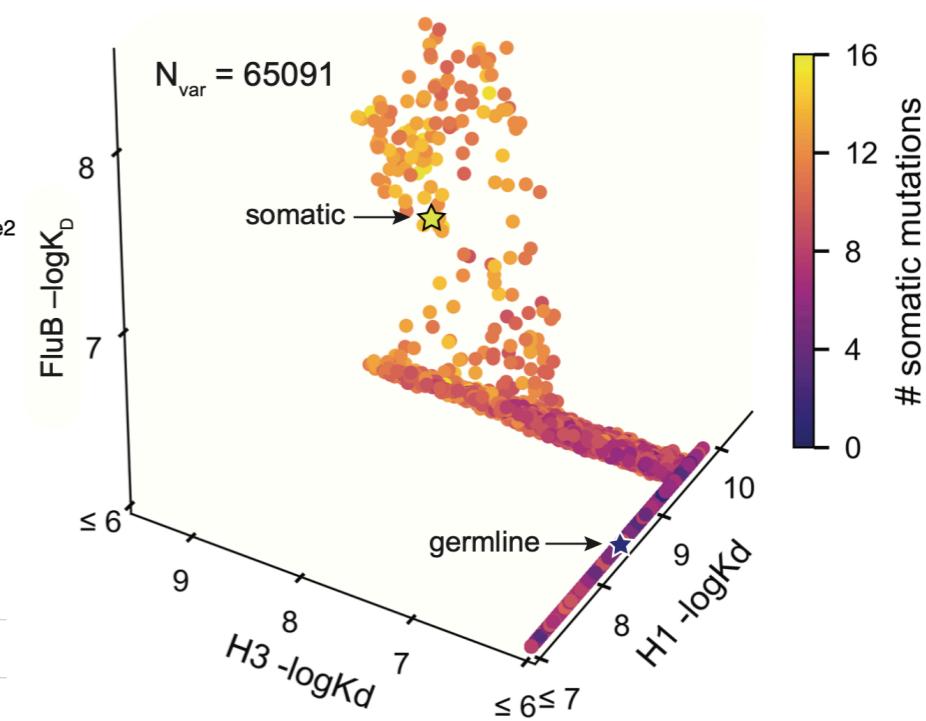
Data #4: path/phylogeny reconstruction



Starr et al. PNAS (2017)

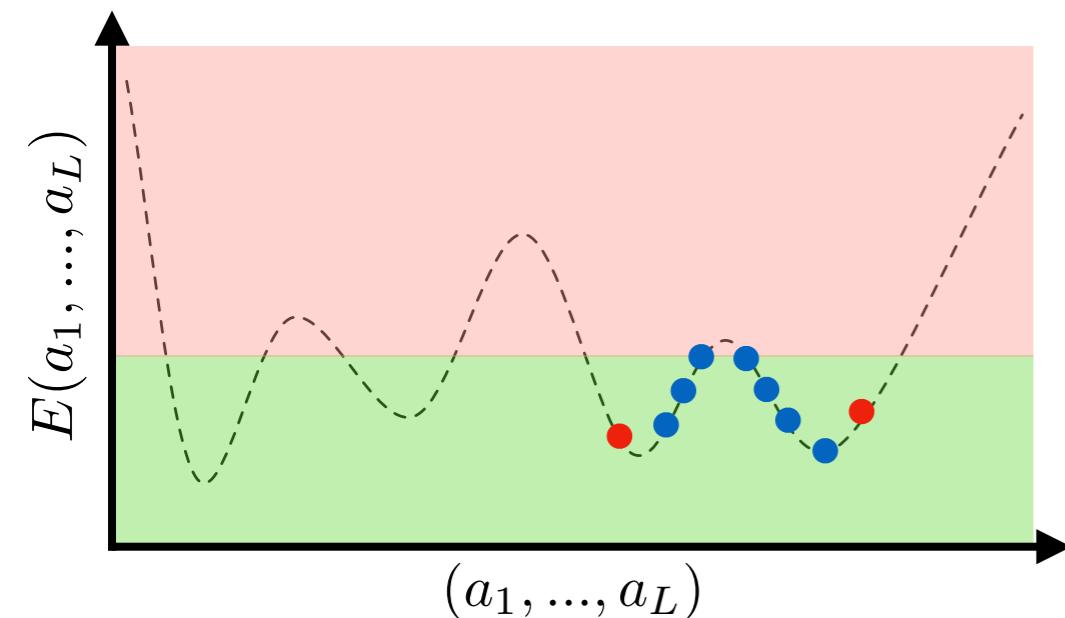


Poelwijk et al. Nat.Comm. (2019)



Phillips et al. eLife (2021)

sequence landscape



sequence data

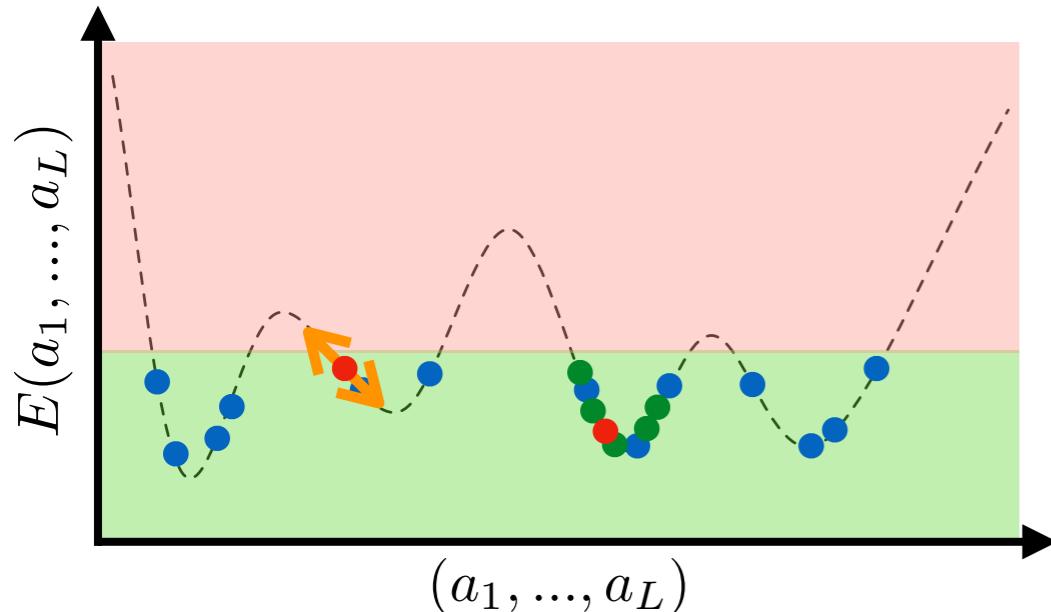
reference sequence

local sample

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

Massive amount of data!

sequence landscape



sequence data

reference sequence

local sample

mutant 1

mutant 2

mutant 3

...

global sample

species 1

species 2

species 3

...

reference sequence

mutant sequences

E₁

E₂

E₃

E₄

E₅

E₆

...

phenotype

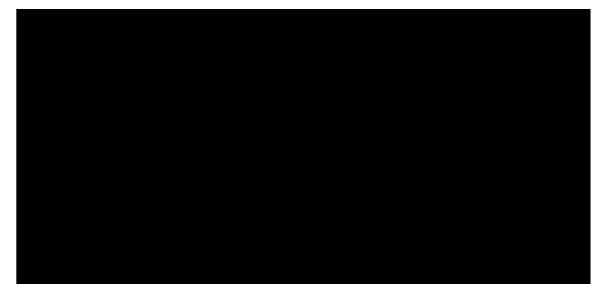
Big data-driven questions:

- Generate artificial sequences
- Generate artificial sequences with desired properties
(binding affinity, fluorescence, antibiotic resistance, catalytic activity...)
- Model in vitro evolution (controlled environment)
- Optimize evolution protocols (antibiotic concentration, vaccination protocols...)
- Understand natural evolution (phylogeny, fluctuating environment)
- Generate paths of artificial sequences connecting two natural ones



Generative modeling

Many images
of human
Faces



Black box
(Generative model)
 $P(\text{image})$



this-person-does-not-exist.com

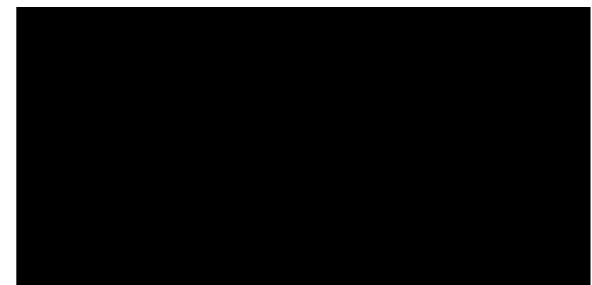
KPDLKPYSPLRDK

REDLVNYNPITEK

QNDLKEYSPVDEK

AADLVANSPTVGK

Alignment of $10^2 - 10^5$
natural proteins



$P(a_1, a_2, \dots, a_L)$



CYDLVGWEPATAK

This protein does not
exist in nature,
but is functional!

Russ et al. Science (2020)

Repack et al. Nature Machine Intelligence (2021)

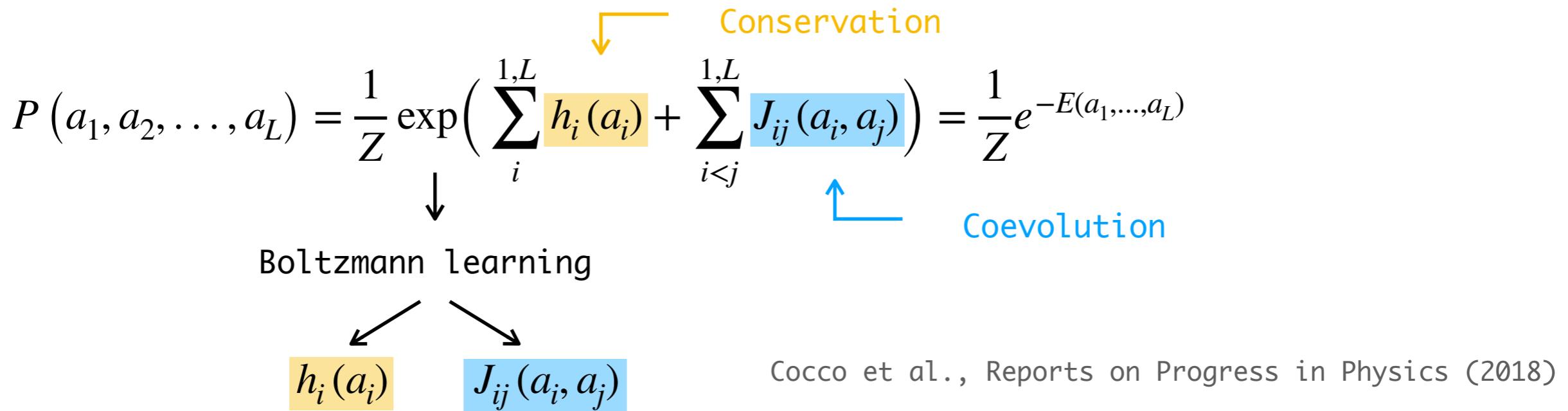
Hawkins-Hooker et al. PLoS Comp. Bio. (2021)

Typical protein with $L = 100$

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

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

Generative modeling



Statistical physics approach: maximum likelihood, disordered Potts models

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

Coevolution is at the basis of all structure prediction methods (DCA, AlphaFold, etc.)

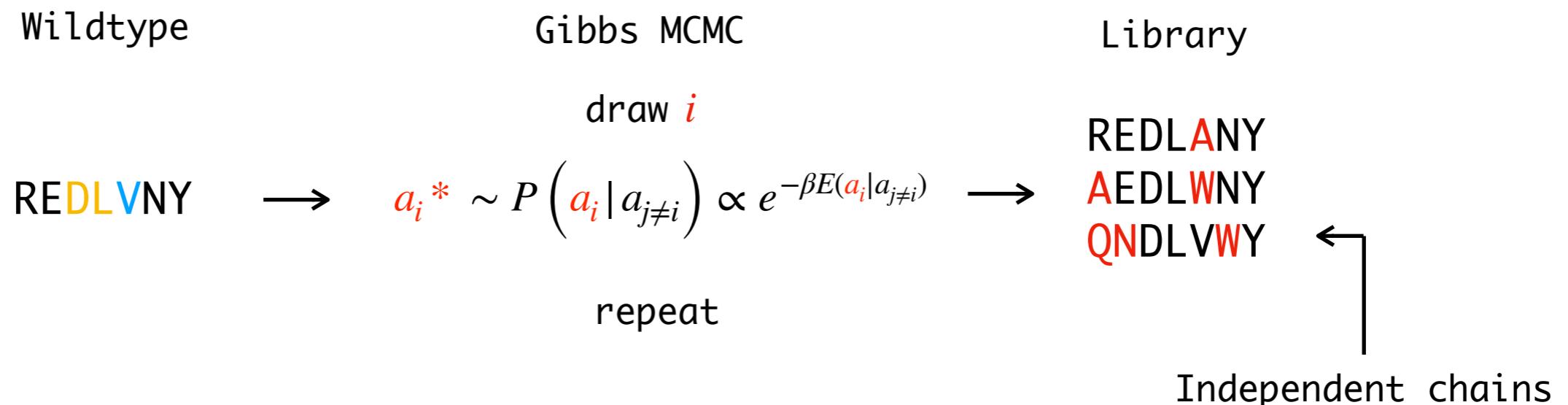
Some contributions from our group:

- Increase alignment accuracy Muntoni et al. PRE (2020)
- Simple and computationally efficient architectures Trinquier et al. Nat. Comm. (2021)
- Information-based procedure for parameter reduction Barrat-Charlaix et al. PRE (2021)
- Two publicly available packages: arDCA and bmDCA Muntoni et al. BMC Bioinformatics (2021)

Local sampling ~ evolution?

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

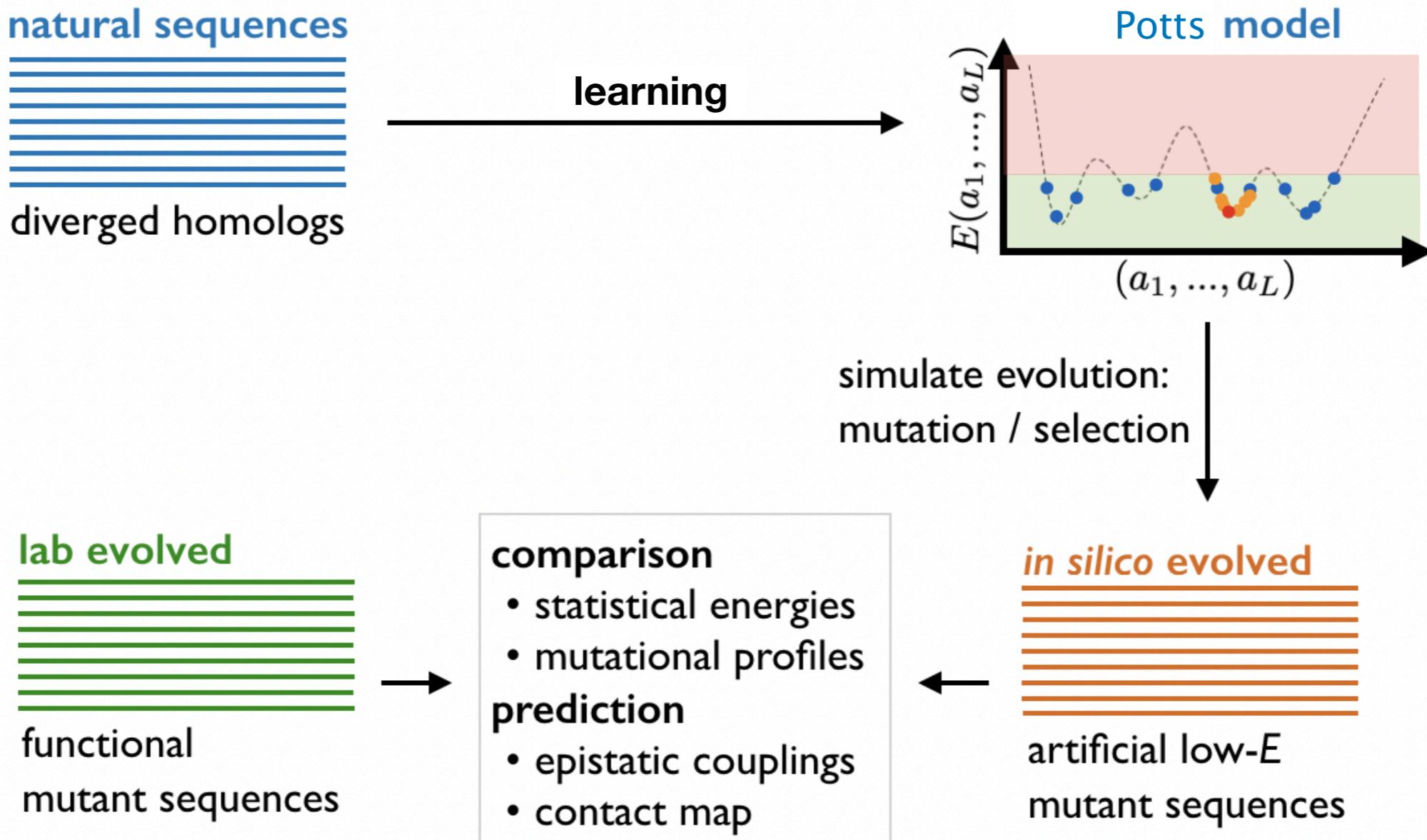
Boltzmann learning



* a_i is drawn from the set of amino acids at distance 1
in terms of nucleotides (DNA sequence)

De la Paz et al., PNAS (2020)

Local sampling ~ evolution?



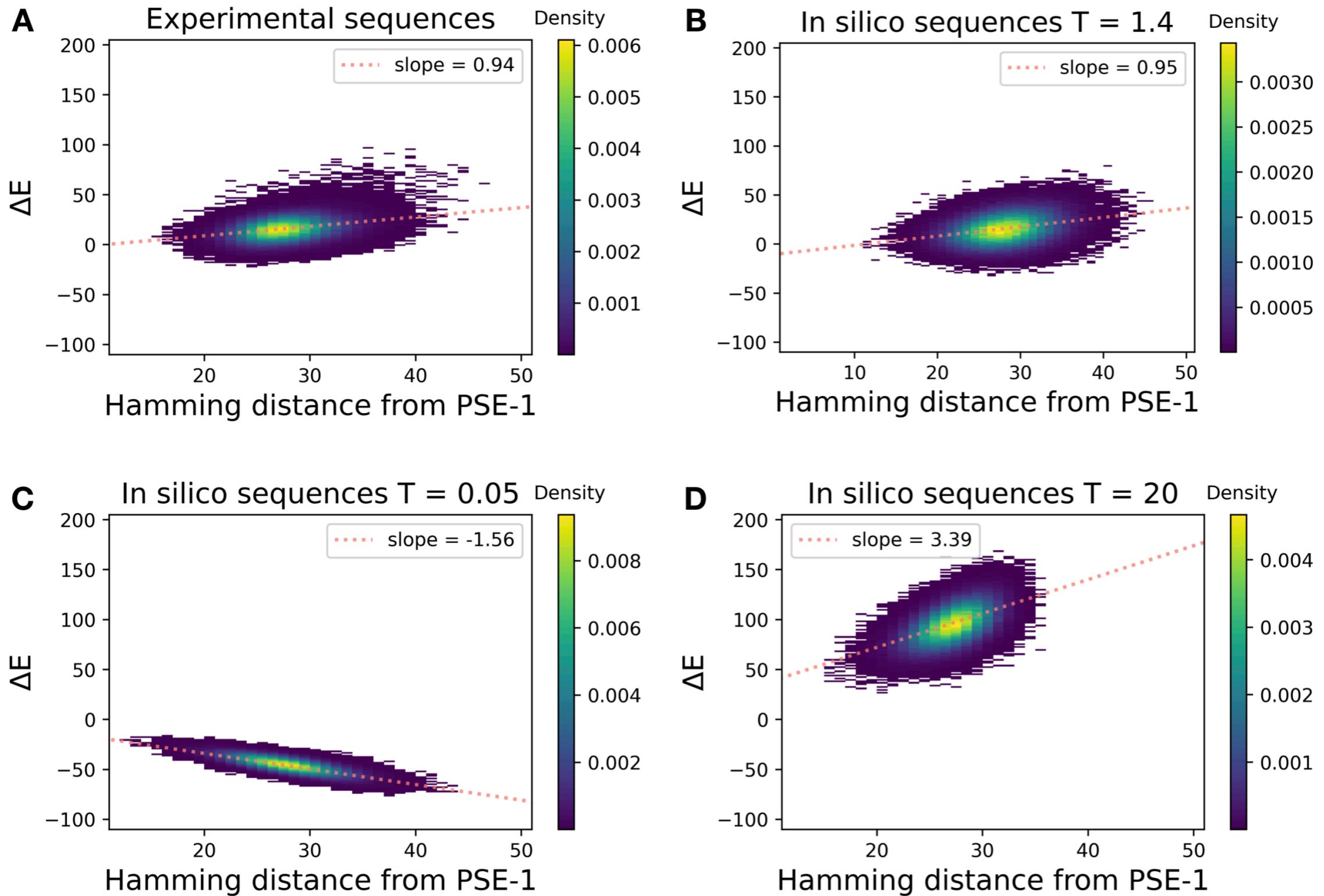
Global learning \neq Local sampling

Natural evolution \neq In-vitro evolution

Phylogenetic effects \neq Independent chains (star phylogeny)



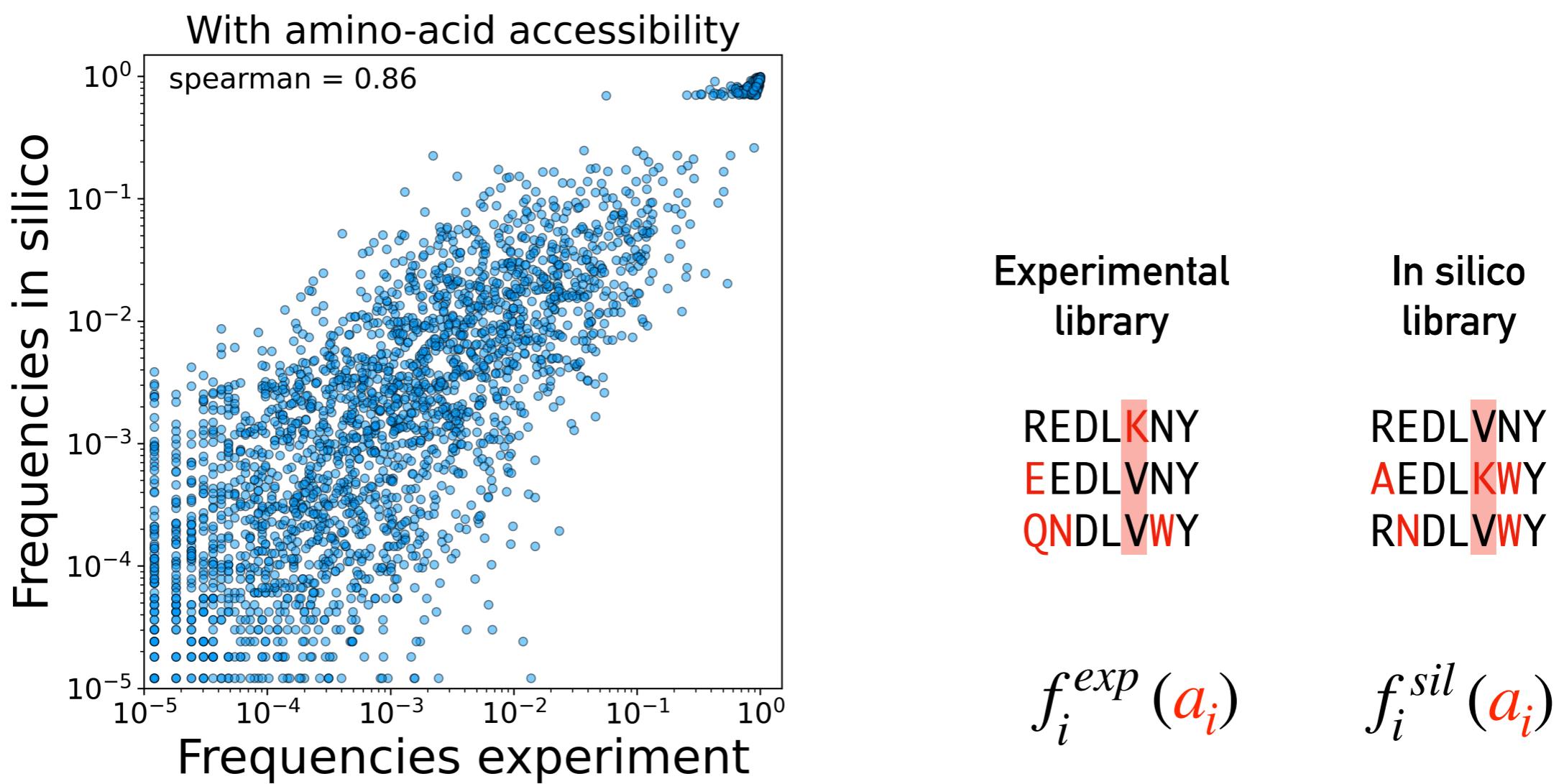
Local sampling ~ evolution?



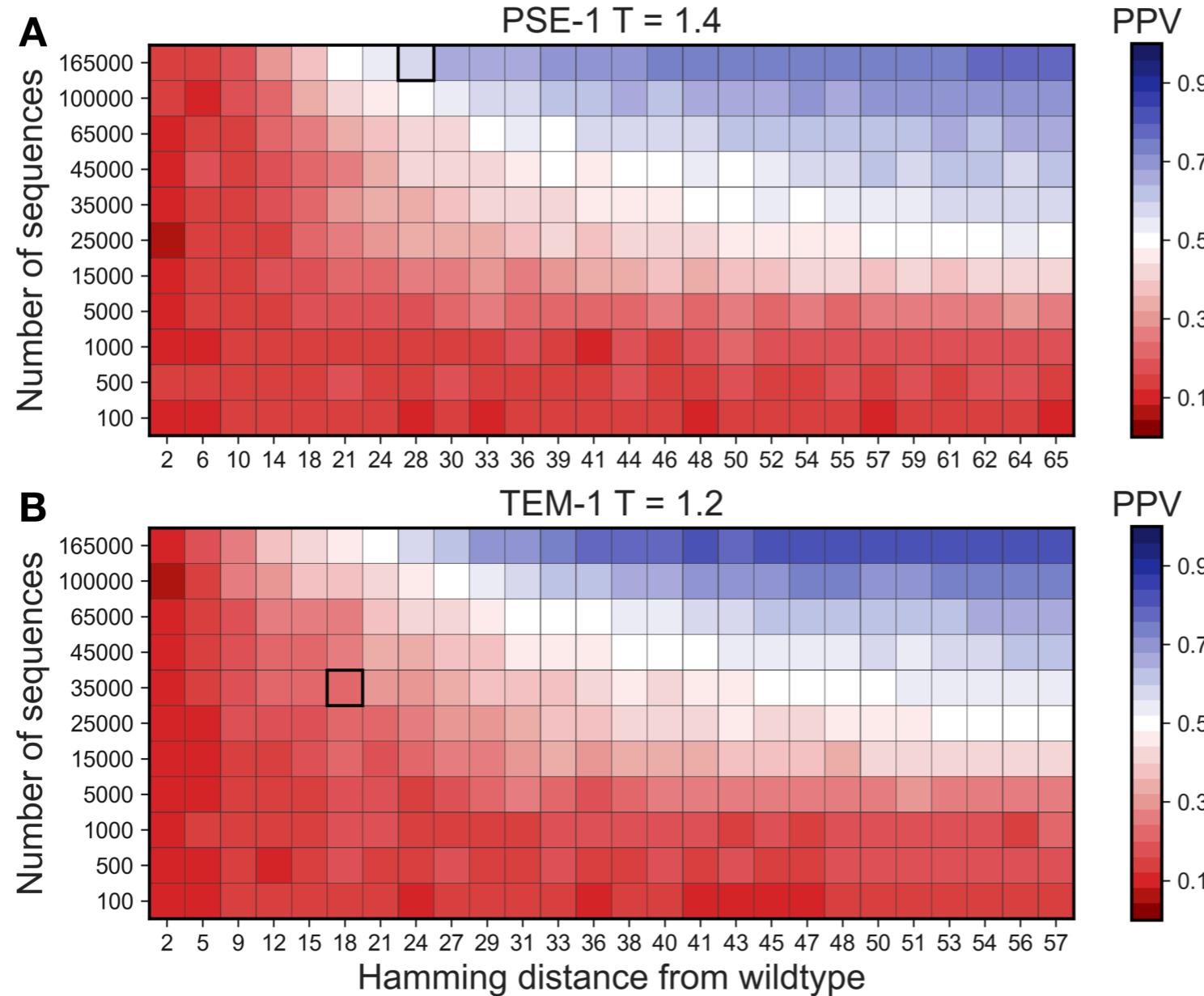
Artificial temperature T can model selection strength

Local sampling ~ evolution?

The model can quantitatively reproduce statistical features of experiments:
e.g. high correlation (86%) of amino acid frequencies



In silico evolution: design new experiments



Stiffler et al., Cell Systems (2020)

Fantini et al., Mol.Bio.Evo. (2020)



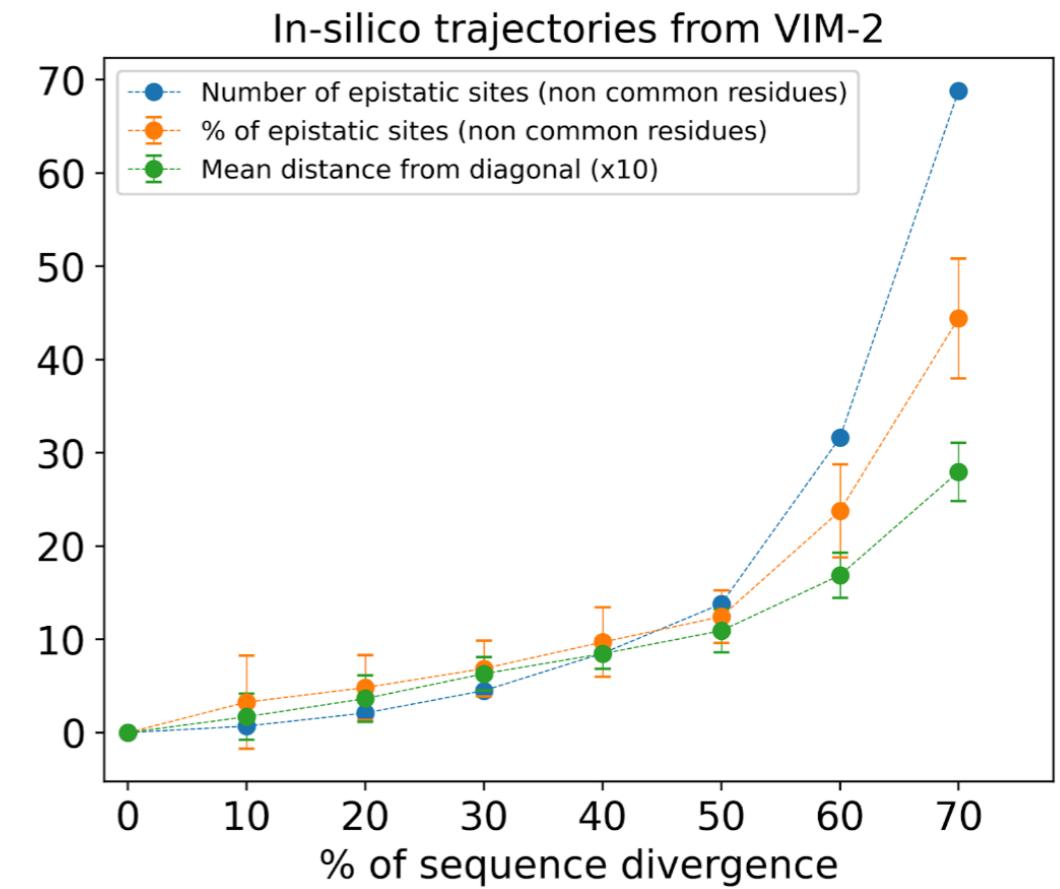
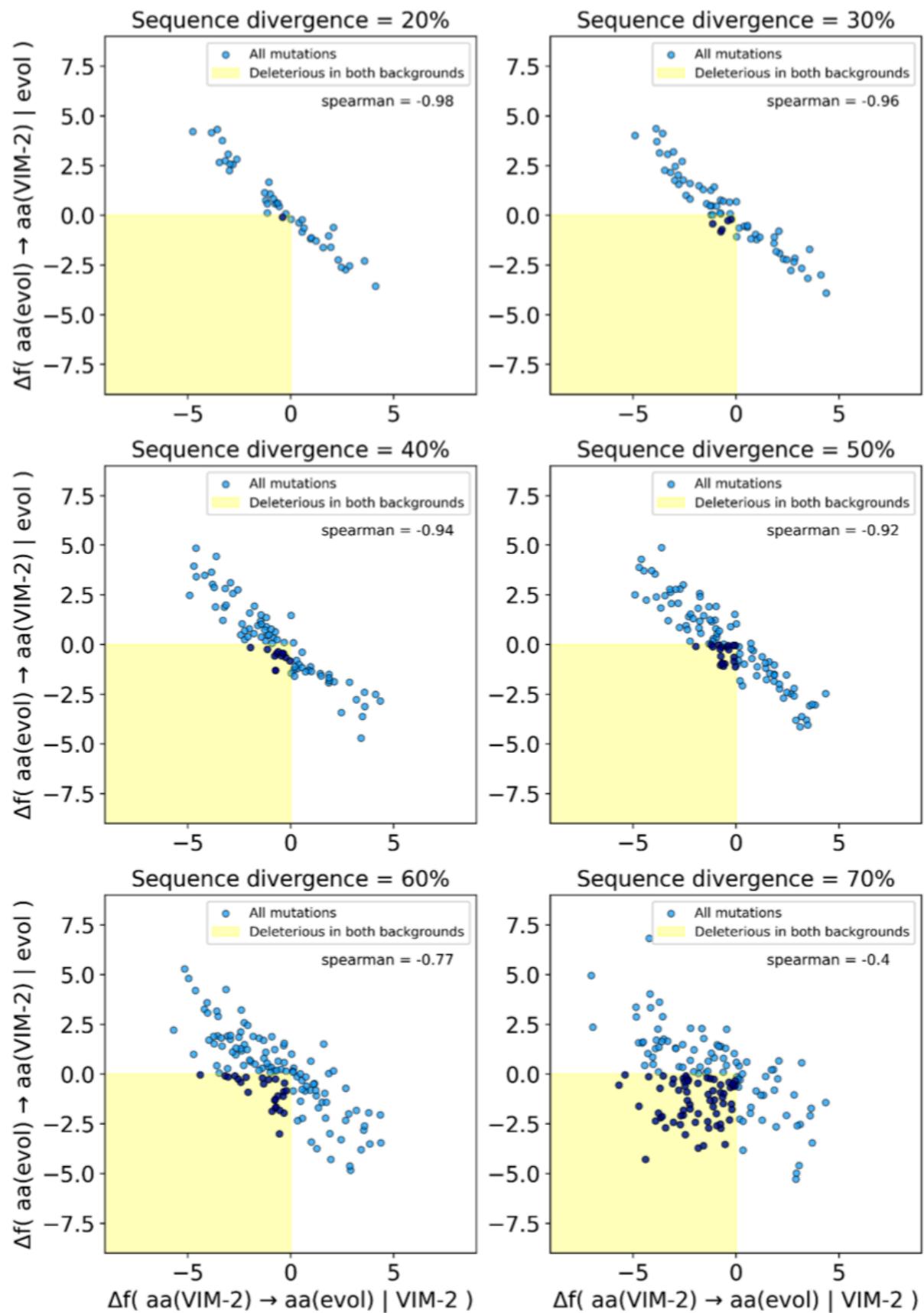
Artificial temperature T can model selection strength

Emergence of coevolution as a function of distance and number of sequences

Optimize and design new experiments

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

In silico evolution: emergence of coevolution



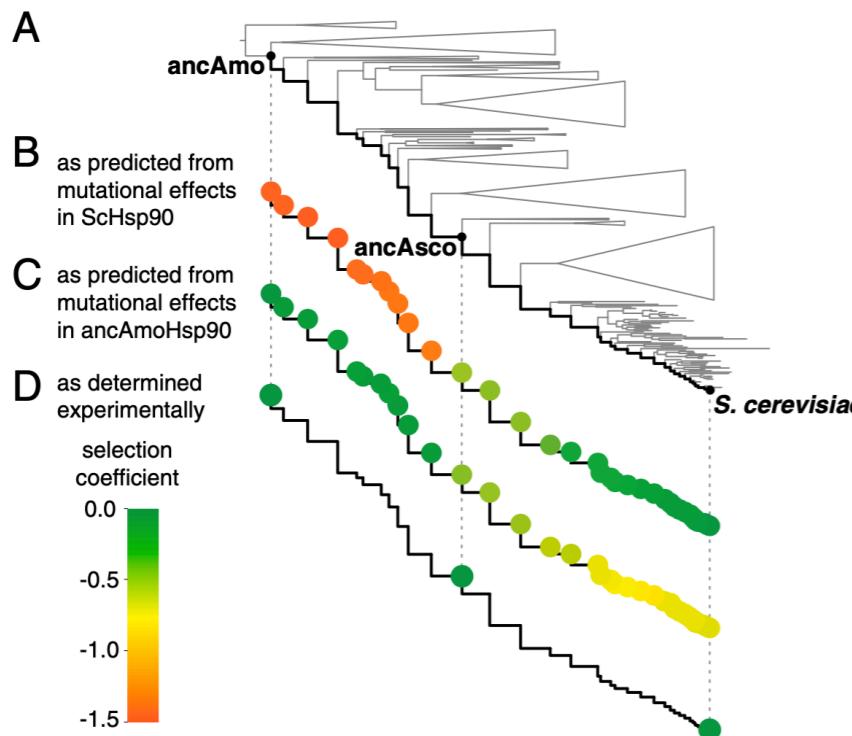
**Emergence of coevolution (aka epistasis)
along in silico evolutionary trajectories**



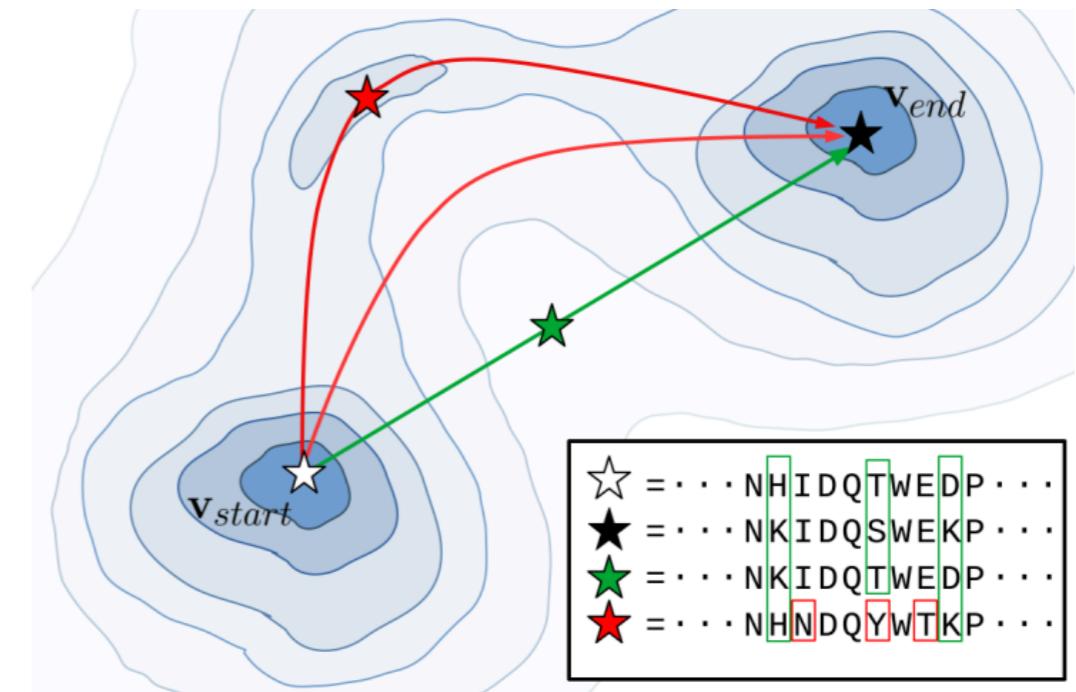
Bisardi, Cotogno, Weigt, FZ
+ Tokuriki lab (in preparation)



Path sampling



Starr et al. PNAS (2017)



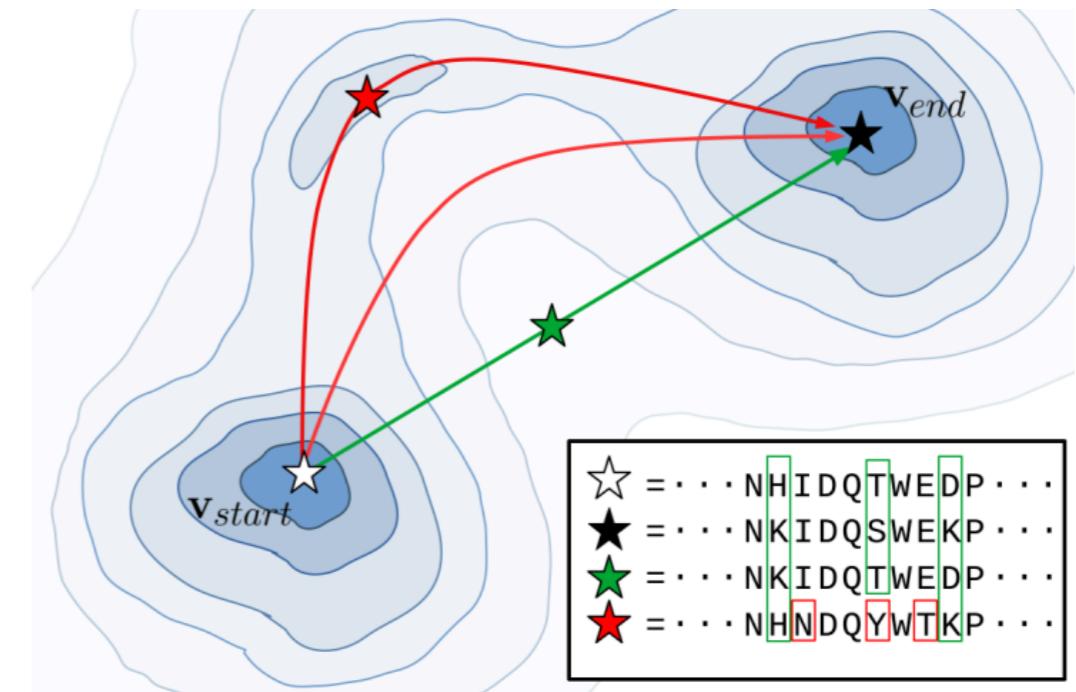
Mauri, Cocco, Monasson, preprint (2022)

- Navigability of the sequence landscape
- Evolutionary paths connecting two wild types
- Intermediate sequences between two wild types with desired properties
- Direct versus wandering paths
- Well defined computational framework: transition path sampling

Dellago et al. JCP (1998), Mora-Walczak-FZ PRE (2012)

$$P(\mathbf{a}_1, \dots, \mathbf{a}_T) = P(\mathbf{a}_1)P(\mathbf{a}_1 \rightarrow \mathbf{a}_2)P(\mathbf{a}_2 \rightarrow \mathbf{a}_3)\cdots P(\mathbf{a}_{T-1} \rightarrow \mathbf{a}_T)$$

Path sampling



Metallo β -lactamases enzymes confer broad antibiotic resistance in bacteria

VIM-2 and NDM-1: well studied “superbugs” in this class
~64% sequence divergence

Tokuriki’s lab characterized all single-mutants of both enzymes (Deep Mutational Scanning, DMS)

We constructed a model based on an alignment of natural sequences and the integration of the experimental DMS

Generate intermediate artificial sequences along a mutational path VIM-2 → NDM-1
Express them in bacteria and measure their fitness

Path sampling

1 m f k l i s k l i l v y l t a s i m a i a s p l a f s v d s g e y p t v s e i p v - - 41
 1 m e l p n i m h p v a k l s t a l a a l m l s g c m p g e i r p t i g q q m e t g d 43

42 - - - G E V R L Y Q I A D G V W S H I A T Q S F D G - A V Y P S N G L I V R D G D E L 80
 44 Q R F G D L V F R Q L A P N V W Q H T S Y L D M P G F G A V A S N G L I V R D G G R V 86

81 L L I D T A W G A K N T A A L L A E I E K Q I G L P V T R A V S T H F H D D R V G G V 123
 87 L V V D T A W T D D Q T A Q I L N W I K O E I N L P V A L A V V T H A H Q D K M G G M 129

124 D V L R A A G V A T Y A S P S T R R L A E V E G N E I P T H S L E G L S S S G D A V R 166
 130 D A L H A A G I A T Y A N A L S N Q L A P Q E G M V A A Q H S L T F A A N G W V E P A 172

167 - - - F G P V E L F Y P G A A H S T D N L I V Y V P S A S V L Y G G C A I Y E L S R 205
 173 t a p n F G P L K V F Y P G P G H T S D N I T V G I D G T D I A F G G C L I K D S K A 215

206 T S A G N V A D A D L A E W P T S I E R I Q Q H Y P E A Q F V I P G H G L P G G L D L 248
 216 K S L G N L G D A D T E H Y A A S A R A F G A A F P K A S M I V M S H S A P D S R A A 258

249 L K H T T N V V K A H T N - r s v v e 266
 259 I T H T A R M A D K - - - I r - - - 270



Bisardi, Cotogno, Chen, Lee
Work in progress

Preliminary result: a functional path

Sequence	Position model	Position VIM	EC50 µg/mL		ΔFitness log2(ΔEC50)	ΔFitness DMS VIM	-ΔE integrated model	-ΔE non-integ. model	Residue distance
			Old AA	New AA			VIM		
VIM2-LE									
					406.7				
MeanE1 1	144	181	T	S	37.86	-3.43	0.05	-1.39	0.59
MeanE1 2	30	67	Y	V	101.18	1.42	1.27	1.14	-3.23
MeanE1 3	27	ins	I	F	36.01	-1.49	missing	0.49	1.22
MeanE1 4	205	242	L	A	22.4	-0.68	0.83	1.13	0.79
MeanE1 5	20	58	A	S	123.52	2.46	1.14	2.57	3.73
MeanE1 6	220	257	K	D	1.56	-6.31	-1.08	-1.82	-1.01
MeanE1 7	108	145	V	Q	53.61	5.1	missing	-0.06	1.55
MeanE1 8	19	57	I	T	22.79	-1.23	1.25	1.4	2.25
MeanE1 9	65	102	Q	E	43.59	0.94	0.55	1.15	1.11
MeanE1 10	101	138	S	L	157.57	1.85	0.24	-0.29	3
MeanE1 11	17	55	S	Q	163.89	0.06	0.95	0.59	-0.3
MeanE1 12	192	229	H	A	315.52	0.95	0.57	1.03	1.61
MeanE1 13	53	90	K	D	125.61	-1.33	0.87	0.62	1.22
MeanE1 14	209	246	L	R	113.19	-0.15	0.62	1.27	1.37
MeanE1 15	143	180	S	T	190.9	0.75	1.09	1.17	1.05
MeanE1 16	141	178	A	G	153.31	-0.32	-0.46	1.04	2.58
MeanE1 17	135	172	L	V	16.04	-3.26	-0.14	0.46	4.87
MeanE1 18	162	199	A	L	299.84	4.22	-0.6	-0.86	-2.89
MeanE1 19	148	185	I	T	123.13	-1.28	-0.32	-1.8	0.19
MeanE1 20	164	201	Y	K	11.46	-3.43	-3.05	-0.96	6.57
MeanE1 21	158	195	Y	F	77.75	2.76	0.38	-0.47	1.74
MeanE1 22	23	61	S	D	1.74	-5.48	-0.54	-1.67	1.24
MeanE1 23	61	98	E	W	1.74	0	missing	-0.3	0.96
MeanE1 24	153	190	S	G	2.24	0.36	0.15	0.03	0.54
MeanE1 25	213	250	K	T	1.86	-0.27	0.19	0.39	1.9
MeanE1 26	124	161	S	G	1.83	-0.02	0.38	0.2	2.69
MeanE1 27	169	206	T	K	2.04	0.16	0.44	0.56	-0.53
MeanE1 28	168	205	R	A	17.17	3.07	0.68	1.41	-0.15

Generative modeling for protein families leads to disordered models with “rough” fitness landscapes
Basins of attraction, transition paths, barriers, topology

Study of the dynamics (evolution) in these landscapes:

- Validate the model against *in vitro* evolution
- Study the emergence of coevolution (epistasis)
- Construct evolutionary paths *in silico* and *in vitro*
- Ultimate goal: design artificial proteins with desired properties

Thank you for your attention!