# Predictability and diversity in immune repertoires

Aleksandra M Walczak

Laboratoire de Physique Théorique - ENS, CNRS

erc
European Research Council
Established by the European Commission

# Immune receptors

- T-cells important actors of immune system
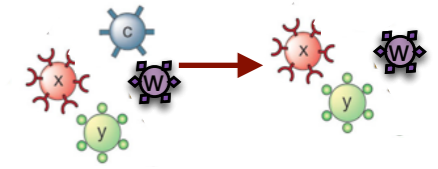


many unique receptors ⟷ many different pathogens (virus...)

- triggers immune response

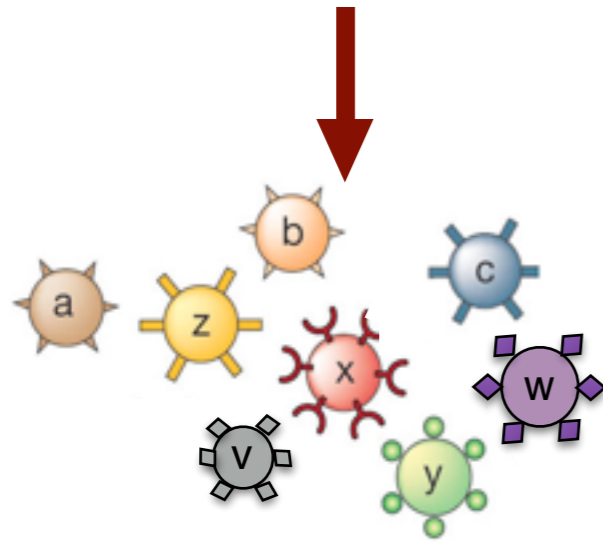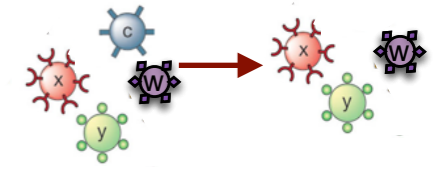natural, healthy (="normal") diversity of immune receptors?

optimal distribution ?

**RECEPTOR GENERATION**

combinatorics + *randomness* → diversity

# Repertoire evolution

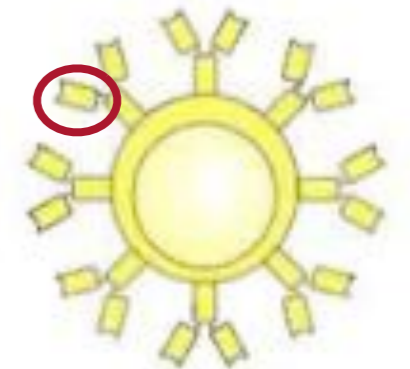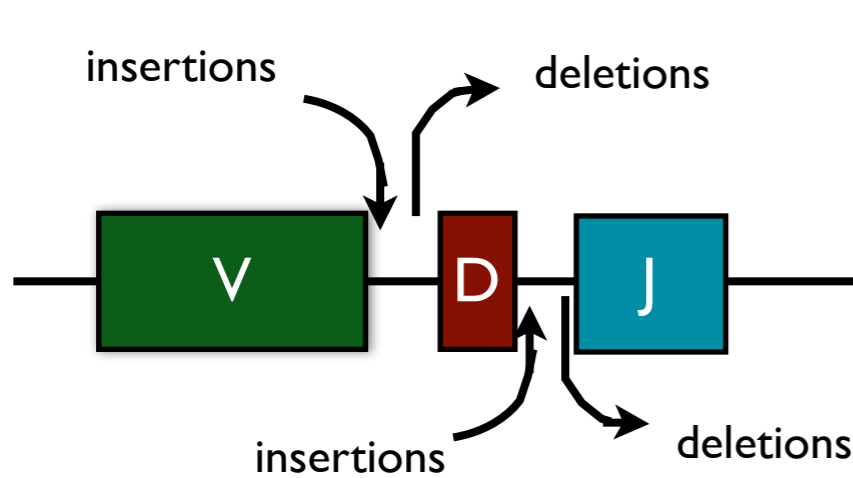**RECEPTOR GENERATION**

combinatorics + *randomness* → diversity

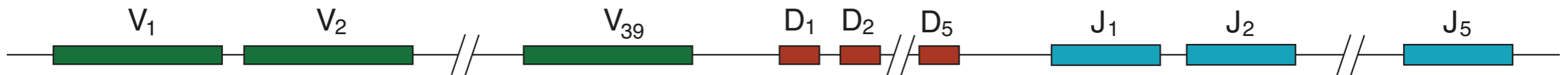$V_1$  $V_2$  $V_{39}$  $D_1$  $D_2$  $D_5$  $J_1$  $J_2$  $J_5$

insertions

deletions

V  D  J

insertions

deletions

$V_9$  $D_2$  $J_3$

• two chromosomes per cell = two attempts

RECEPTOR GENERATION

THYMIC SELECTION

bind to self?

bind to nothing?

bind too strongly to self?

SOMATIC SELECTION

constant somatic evolution

# Sequence data

**new data**

blood test

sequencing machine

**=**

*natural diversity* distribution

**+**

- human T-cell beta chain receptor sequences

- 9 people

- out of frame reads (~14% of each type of cells) = 35,000 unique reads → *generation*

- in frame reads  (~235,000 unique reads) → *selection*



CDR3

45 primers { → → → }   $V_9$   $D_2$   $J_3$   { ← ← ← } 18 reverse primers

DNA

PCR

480,000 unique reads per individual

64 -101 bp reads   { ← ← ← } 13 sequencing primers

Illumina reads

data from Robins lab and Chudakov lab

# Probabilistic VDJ recombination annotation
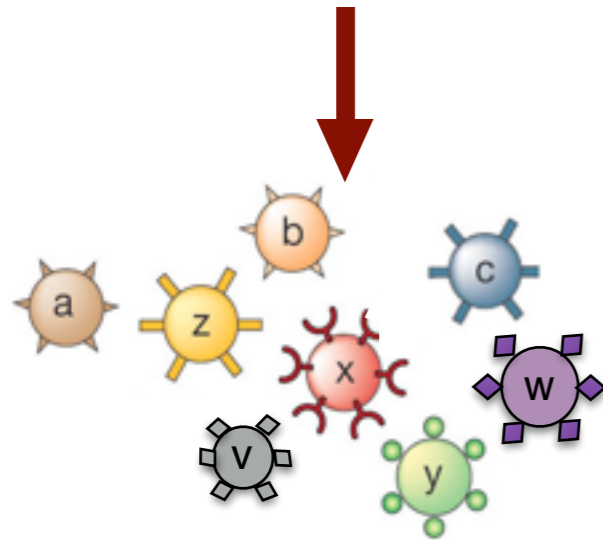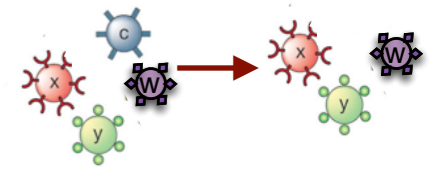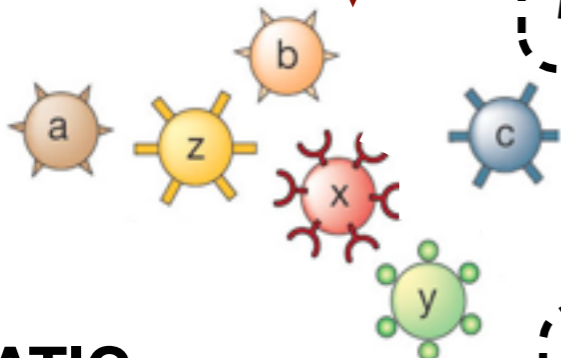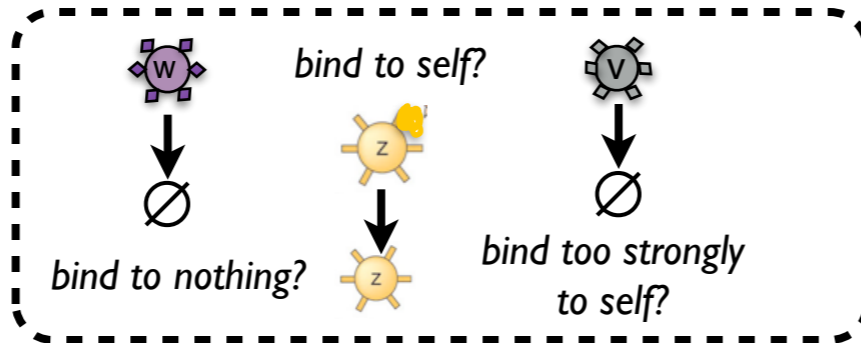
## sequence generation:

combinatorics + *randomness* → diversity

$\vec{\sigma}$ - receptor DNA sequence

V₁ V₂ V₃₉ D₁ D₁ₐ D₂ J₁ J₂ J₅

insertions → deletions

V₁₂ D₂ J₃

insertions → deletions

```
TAGGACCTCGGAAACCTCTTCTGCTGGCGCAGAGATACAGGGCTGAGTCTTCTT
TATCTGCCGATGTCGCGAAGGCCCTCCCGCTAAGATCACTGGTGGCACAGAAGT
TAGGGAGAGGTGCCTAACTGCTGGCACAGAAGTACAGAGAGGTCTGGTTGGGGT
TCCGCCGCTAGTCCCTGAAACTACTGGCACAGAGATAGAAAGCTGTCGGGTTCT
CGAAACTGCTGGCACAGAAGTACACAGATGTTTGGGAGGGAGCAGCCGACTCCA
TCTTGGCCGCTAGTCCGAGAAACTGCTGGCACAGAAGTACACAGATGTTTGGGA
TTGTAGGAGCCGGACCGGCCCCCTGTCCCCTTGGCTGCTGGCGCAGAGATACAG
TCATTTAACGTGCGGCCCGCCTGGCACAGAAGTAAAGAGCTGTCTGGTTGTGGT
TAGTAACTCCGCTTCACTGCTGGCACAGAAGTACACAGATGTCTGGGAGGGAGC
TCCCTCCGGTTTGAAGGGTCTGCTGGCACAGAAGTACACAGATGTTTGGGAGGG
CCGGTGTTCGCACAGCCCTGGGGACCCTGGCGCAAACCCCGCTTCCCTCGAGGA
```
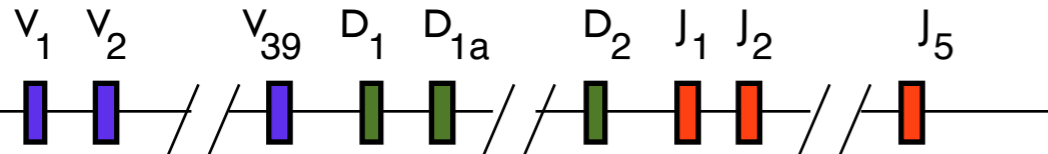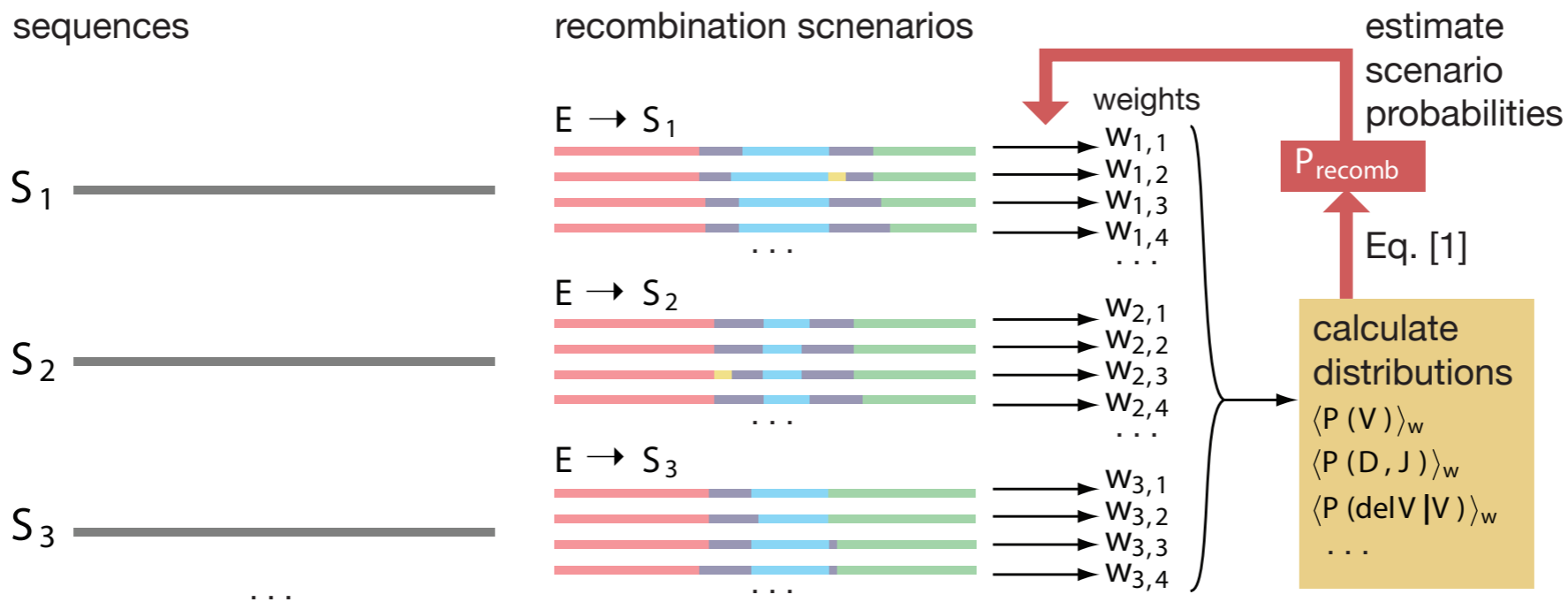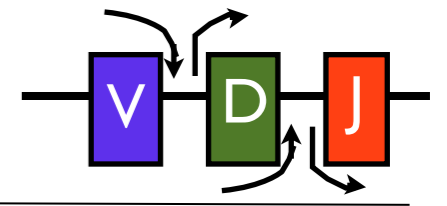
- ## probabilistic assignment of:

sequences

S₁ ─────────

S₂ ─────────

S₃ ─────────

. . .

recombination scnenarios

E → S₁

. . .

E → S₂

. . .

E → S₃

. . .

weights

w₁,₁
w₁,₂
w₁,₃
w₁,₄
. . .

w₂,₁
w₂,₂
w₂,₃
w₂,₄
. . .

w₃,₁
w₃,₂
w₃,₃
w₃,₄
. . .

estimate scenario probabilities

P_recomb

Eq. [1]

calculate distributions

$\langle P(V) \rangle_w$
$\langle P(D, J) \rangle_w$
$\langle P(\text{del}V \,|\, V) \rangle_w$
. . .

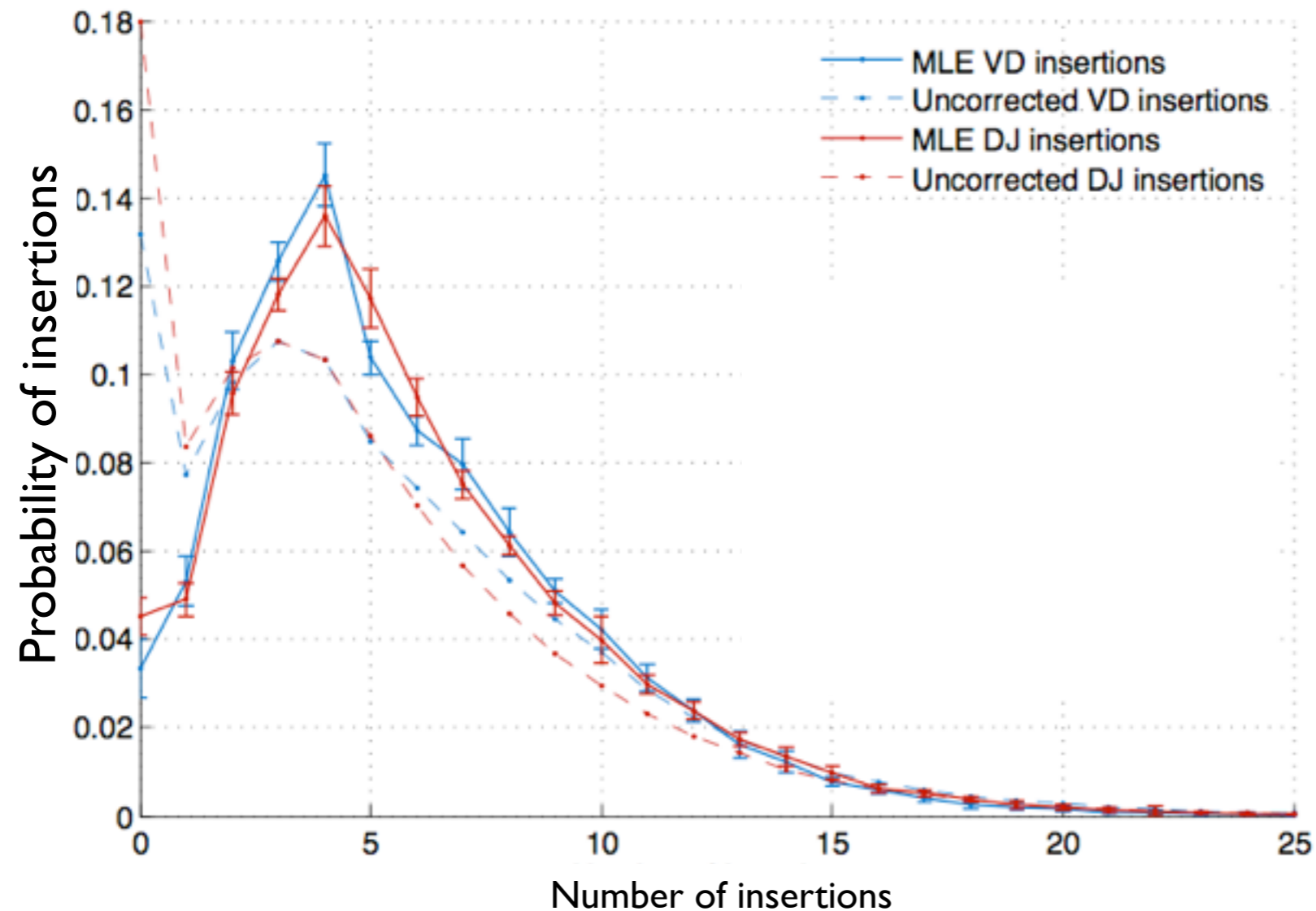$$P^{\text{recomb}}(\text{scenario}) = P(V)P(D, J)P(\text{deletions}V|V)P(\text{insertions}DJ)...$$

[1]

# Universal insertion profiles

universal mechanism for generation

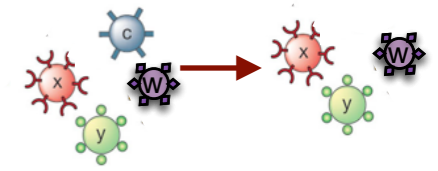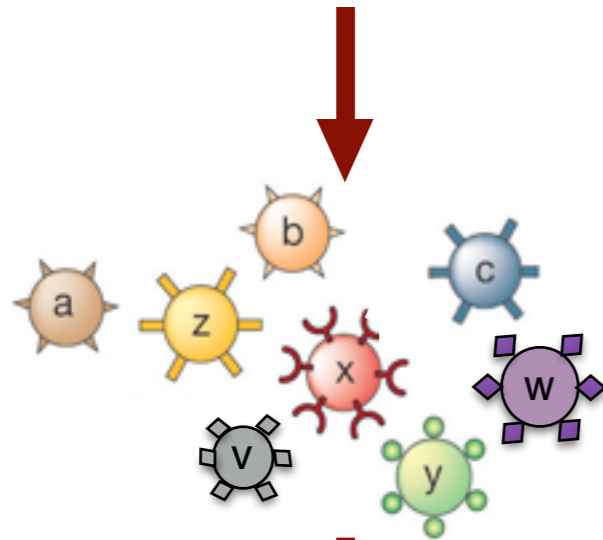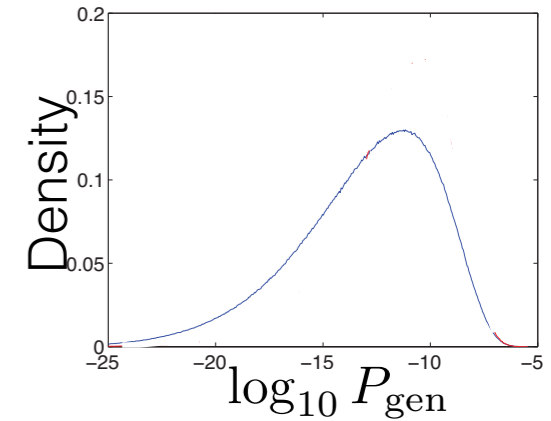- VD and DJ insertion profiles are identical
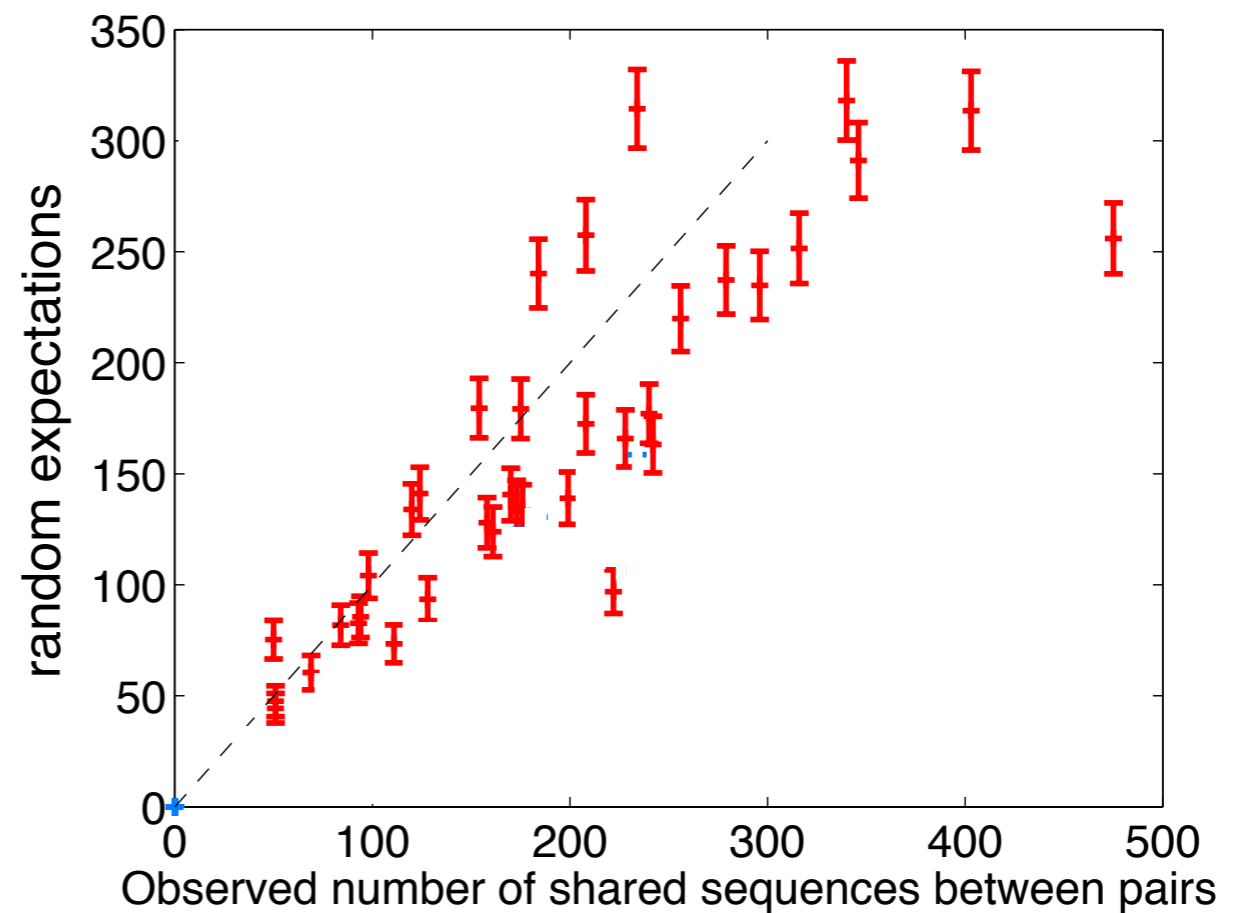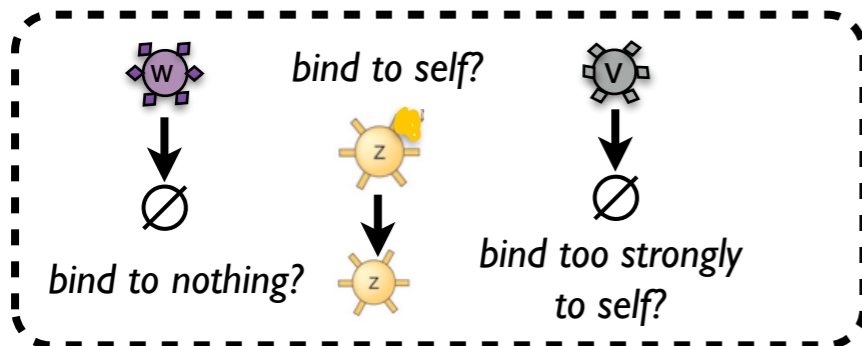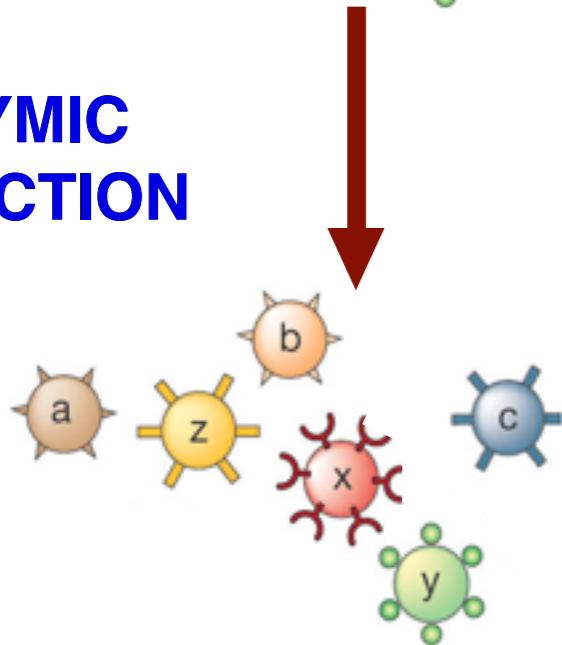
# Receptor sharing

**RECEPTOR GENERATION**

- quantify using selection factors

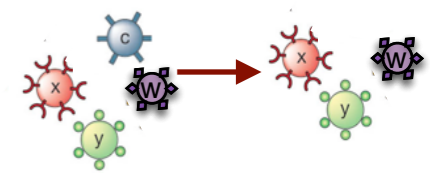$$Q(\{\sigma\}) = \frac{P_{\text{post-sel}}(\{\sigma\})}{P_{\text{gen}}(\{\sigma\})}$$



- how many shared receptors between 2 people?

**THYMIC SELECTION**



bind to self?

bind to nothing?     bind too strongly to self?

→ close to random expectations

- entropy of generated repertoire



$$\Longrightarrow \text{ repertoire size } 10^{13} \text{ sequences}$$

- entropy of post-thymic selection repertoire



$$\Longrightarrow \text{ repertoire size } 10^{11} \text{ sequences}$$

$\longrightarrow$ thymic selection gives 50-fold reduction in diversity

- thymic selection keeps ~15% of sequences but only 2% of diversity

$\longrightarrow$ thymic selection gets rid of rare clones

selection favours clones that are likely to be generated

limited number of encounters

**How should immune receptors be distributed
to minimize harm from infections?**

lymphocyte
repertoire

antigenic
environment

# Cross-reactivity



receptor distribution

$P_r$

antigen distribution

$Q_a$

receptors

r  a

✔ $f_{r,a}$

recognition probability

antigens

- cross-reactivity - recognition probability

- probability of immune response from encounter with a given antigen

$$\tilde{P}_a = \sum_r f_{r,a} P_r$$

recognition probability

receptor
antigen

shape space

# Receptors - antigens interactions

receptor distribution

antigen distribution

$P_r$

$Q_a$

r    a

✔ f$_{r,a}$

receptors

antigens

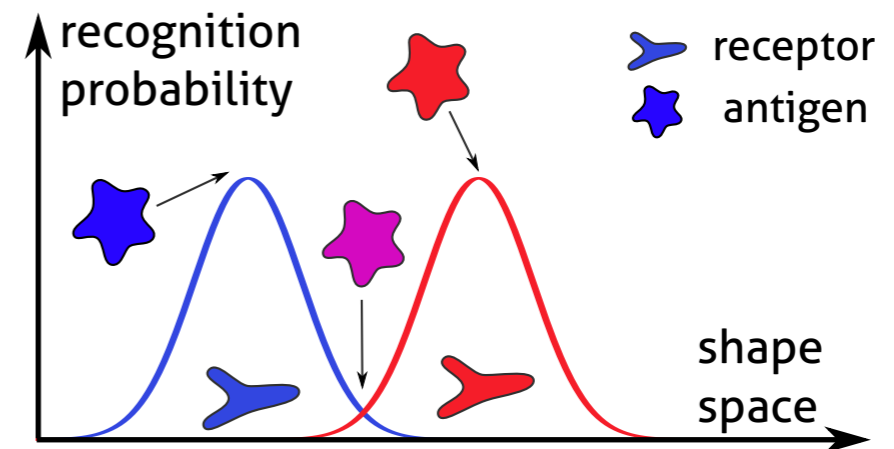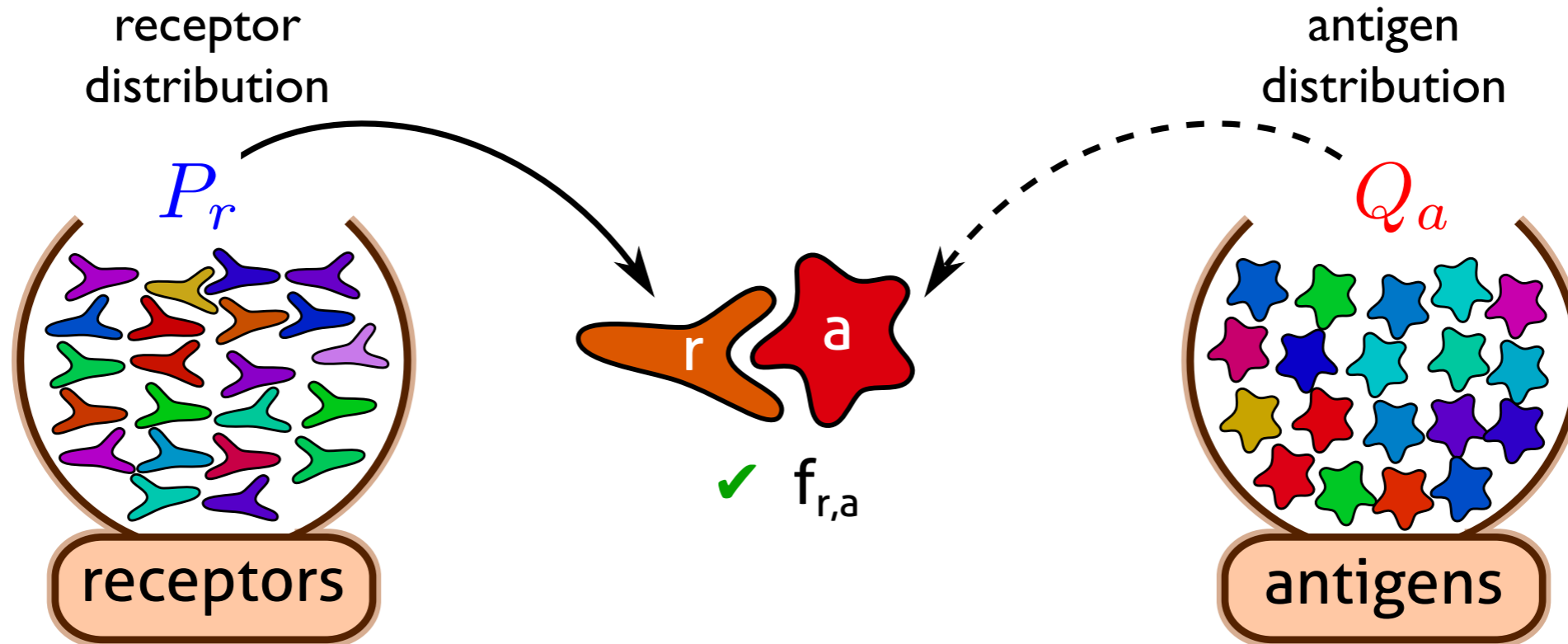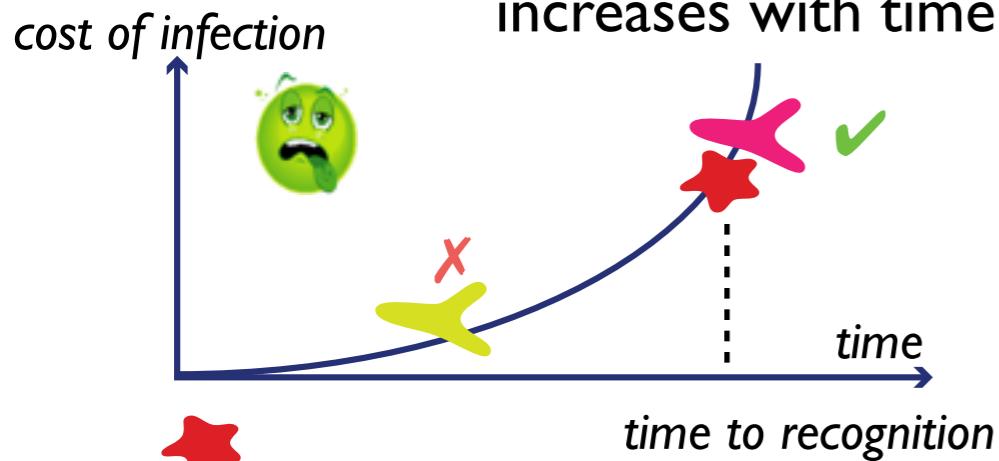- probability of immune response from encounter with a given antigen  $\tilde{P}_a = \sum_r f_{r,a} P_r$

- time measured in mean number of encounters  $m$

cost of infection

*harm* caused by a given antigen increases with time

✔

✗

time

time to recognition

virulence

effective cost of infection

$$\bar{F}_a(P_r) = \mu_a \int_0^{+\infty} dm\, F_a(m)\, \tilde{P}_a e^{-m\tilde{P}_a}$$

Poisson distributed recognition

$$\mathrm{Cost}(\{P_r\}) = \sum_a Q_a \bar{F}_a(P_r)$$

# Cost

receptor distribution

$P_r$

antigen distribution

$Q_a$

r a

✔ f<sub>r,a</sub>

receptors

antigens
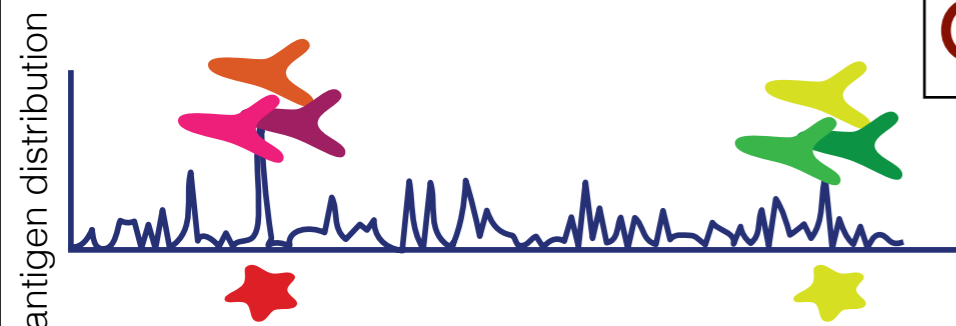
total harm caused by antigen increases with time

$$\mathrm{Cost}(\{P_r\}) = \sum_a Q_a \bar{F}_a(P_r)$$

trade-off: many antigens ⟷ limited resources

Optimal repertoire?

antigen distribution

antigen distribution

→ minimize cost given fixed antigen distribution
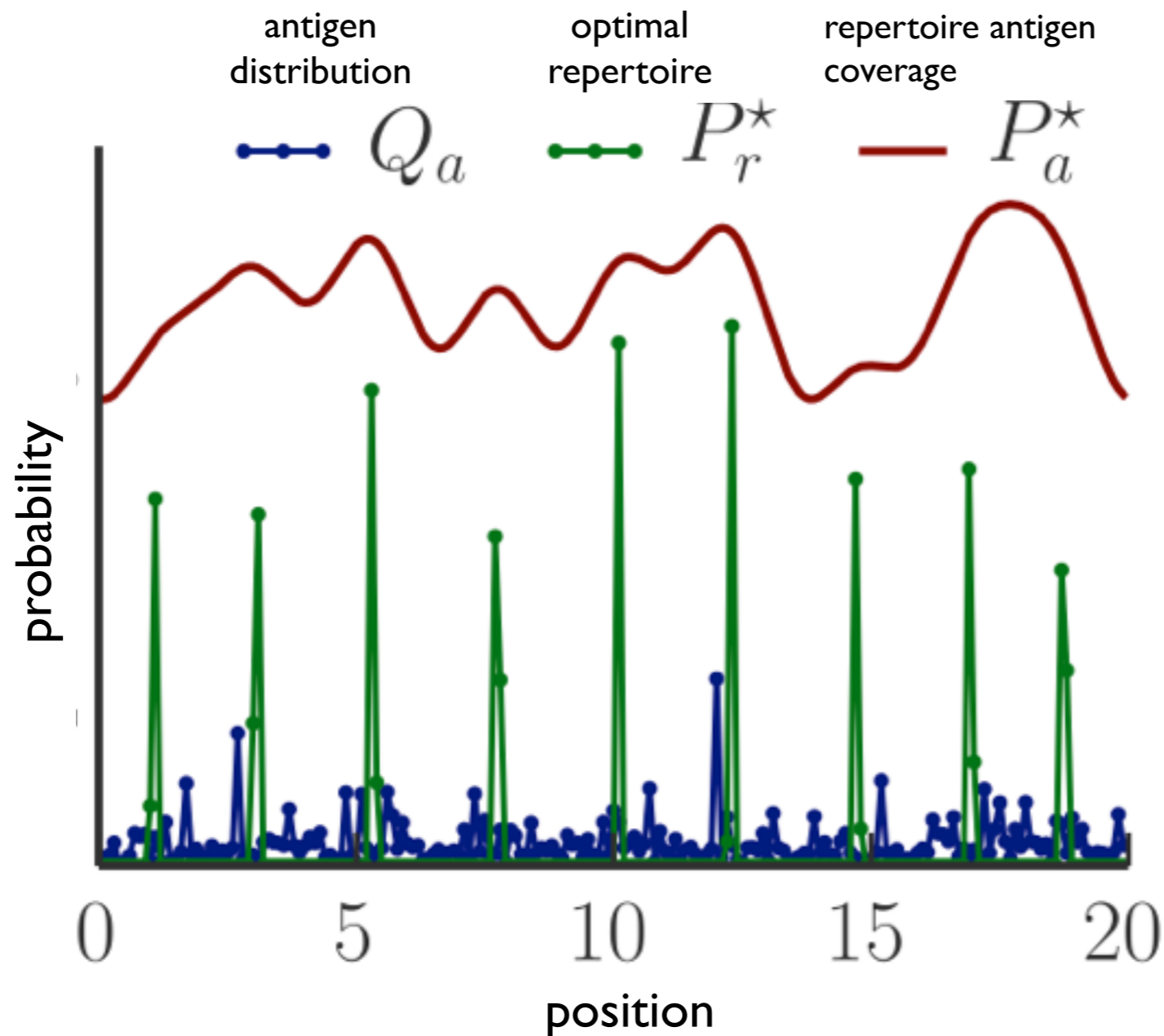
# Peaked optimal repertoires

- exponentially expanding antigen population
  + exponentially growing cost in time

$$F(m) = m$$

- peaked distributions
- tile space

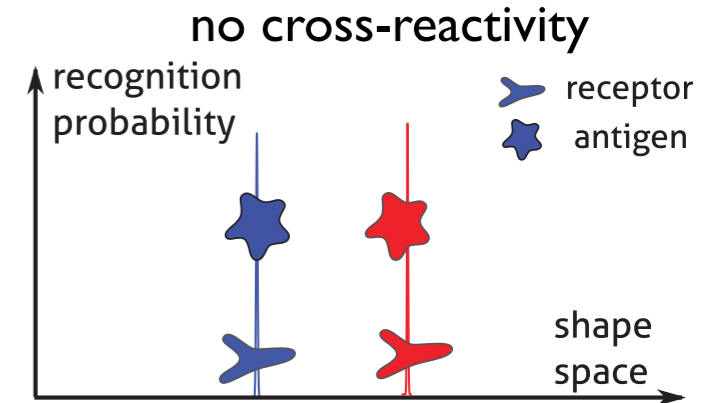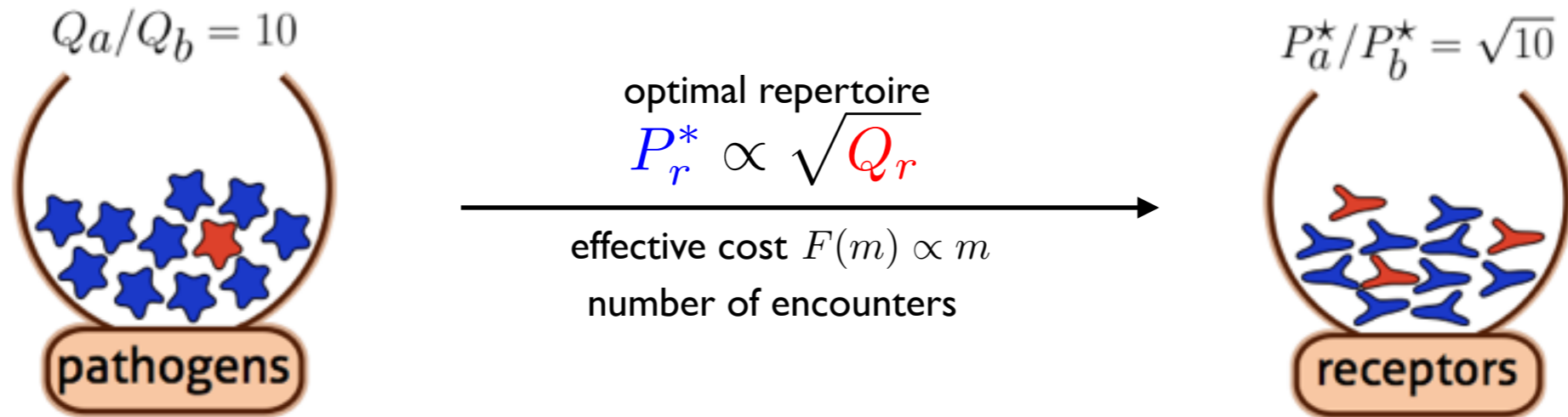- coverage follows antigen distribution
- but not exactly



antigen distribution $Q_a$    optimal repertoire $P_r^\star$    repertoire antigen coverage $P_a^\star$

# Covering rare pathogens

How many resources aimed at common/rare antigen?

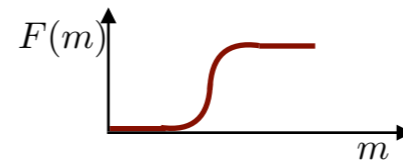depends on cost of late recognition $\longrightarrow$ effective cost function

- exponentially expanding antigen population

**no cross-reactivity**

recognition probability

$\longrightarrow$ receptor

$\bigstar$ antigen

shape space

$Q_a/Q_b = 10$

**pathogens**

optimal repertoire
$$P_r^* \propto \sqrt{Q_r}$$

effective cost $F(m) \propto m$

number of encounters
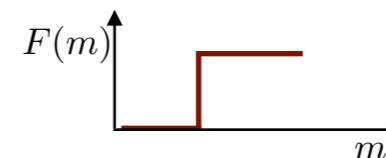
$P_a^\star/P_b^\star = \sqrt{10}$

**receptors**

- exponentially expanding antigen population + exponentially growing cost in time

$$F(m) = m^\alpha \qquad \longrightarrow \qquad P_r^* \propto Q_r^{\,1/(1+\alpha)}$$
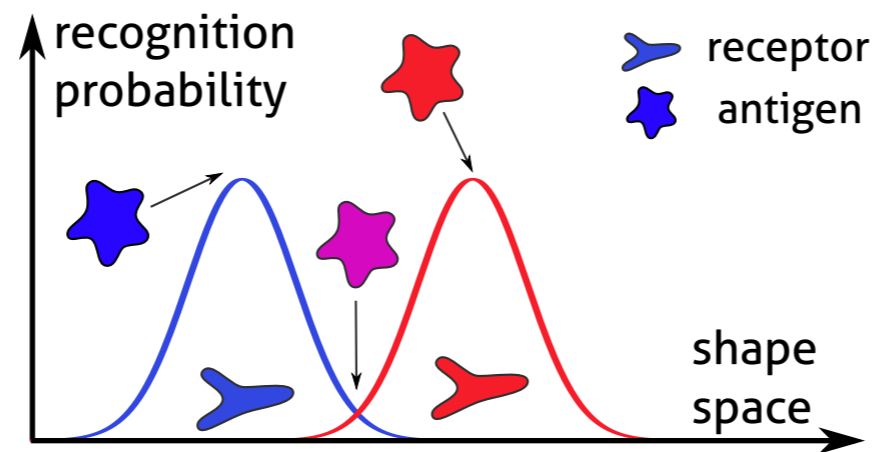
- saturated cost $\longrightarrow$ low frequency cut-off

$F(m)$ ... $m$

- harm past threshold $\longrightarrow$ flattened receptor distribution

$F(m)$ ... $m$

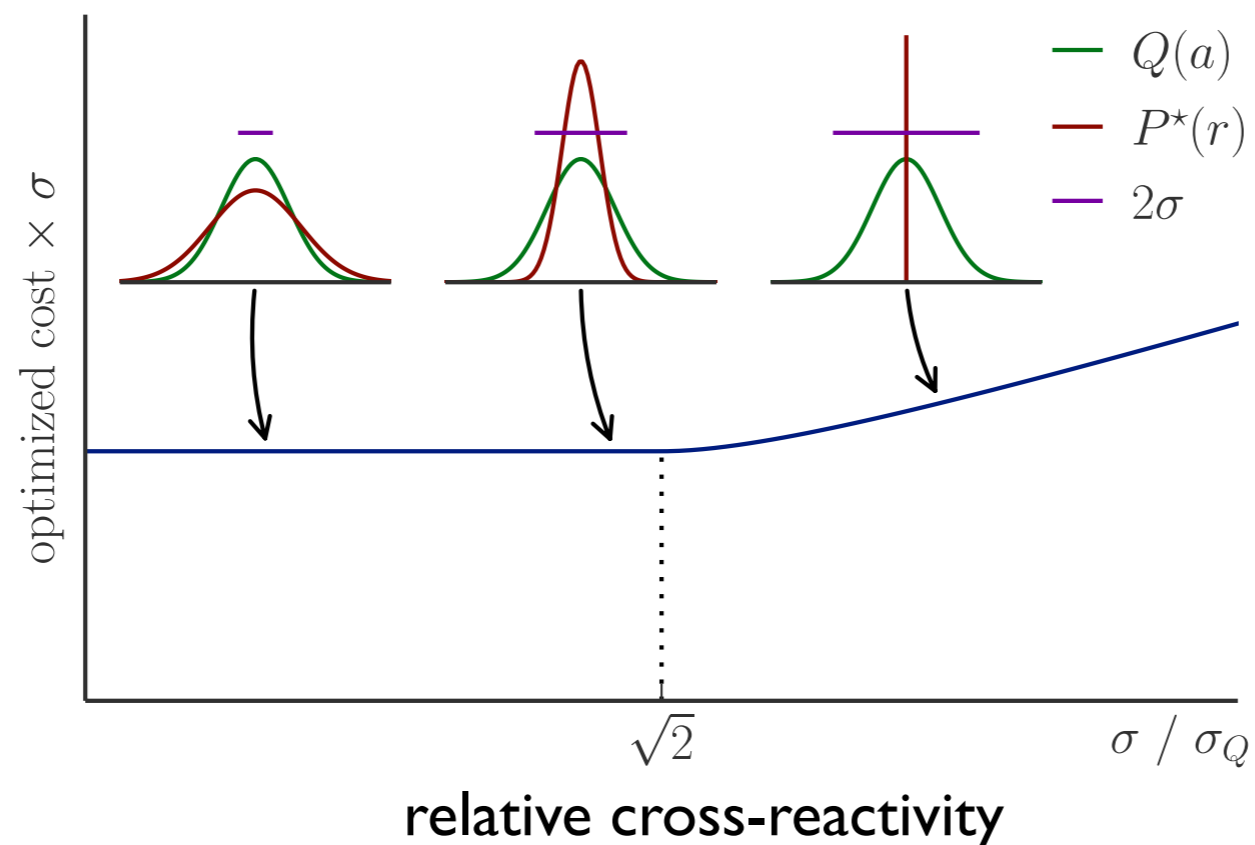- $P_r^* \propto Q_r \Longleftrightarrow$ very slowly increasing cost $\quad F(m) \propto \ln m$

- antigen distribution + cross-reactivity Gaussian


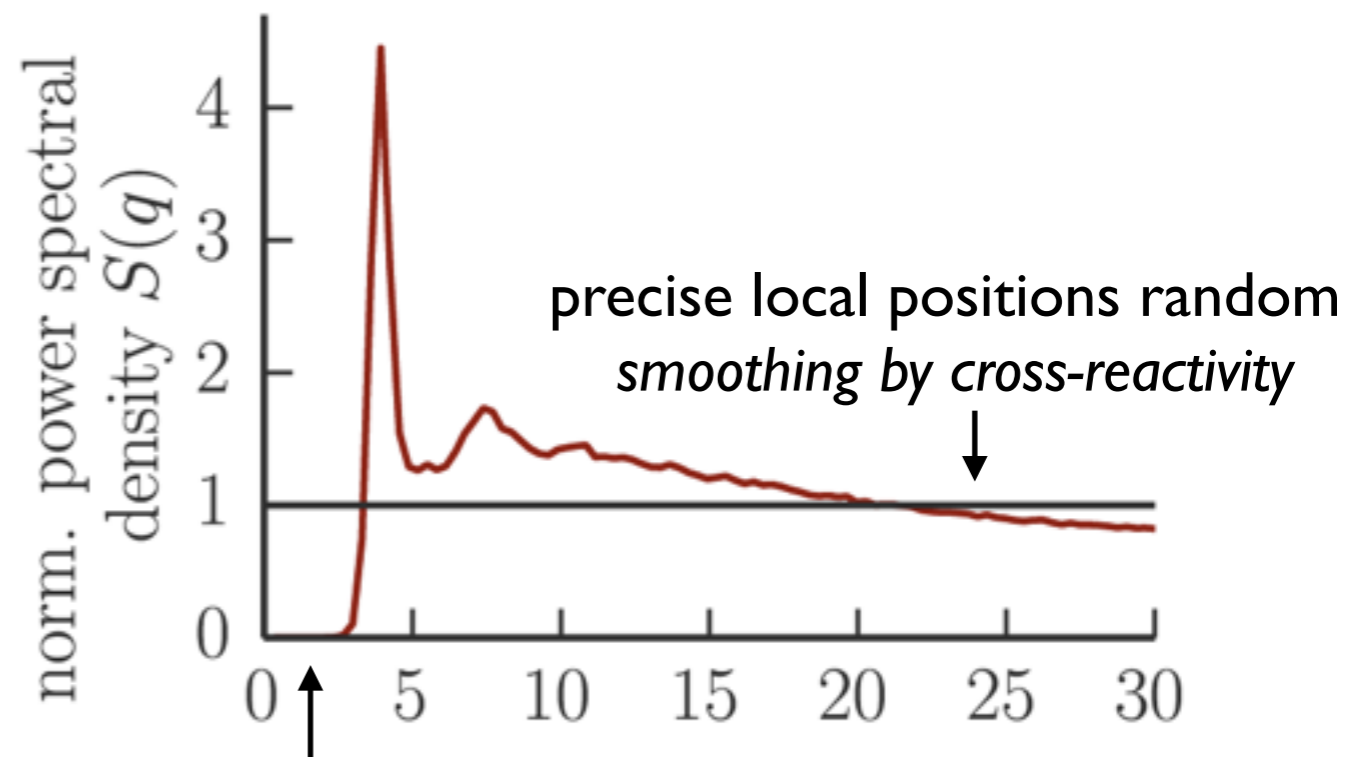
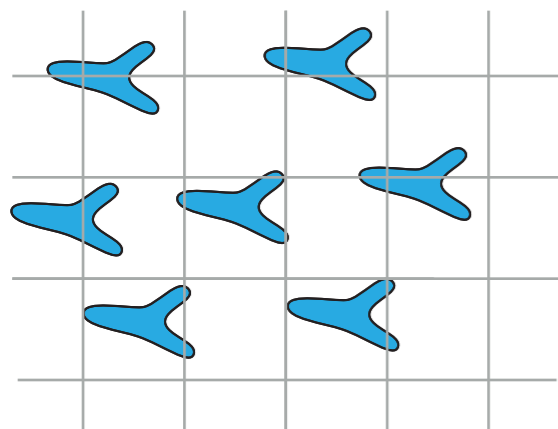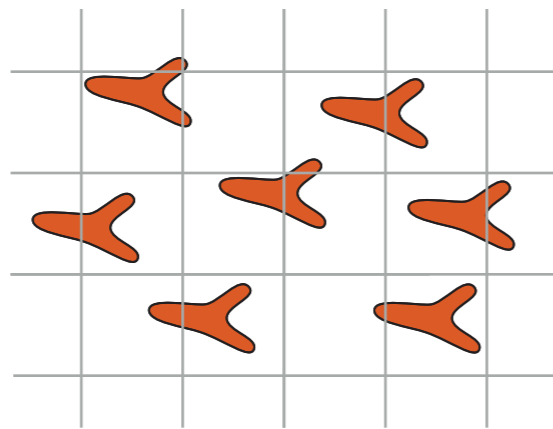$\longrightarrow$ large cross-reactivity concentrates distribution

# Peaked optimal repertoires

# Disordered hyperuniformity

receptors cannot be close to each other

disordered tiling -
as close as possible but excluding

radial distribution function $g(R)$

$R / \sigma$

precise local positions random
*smoothing by cross-reactivity*

norm. power spectral density $S(q)$
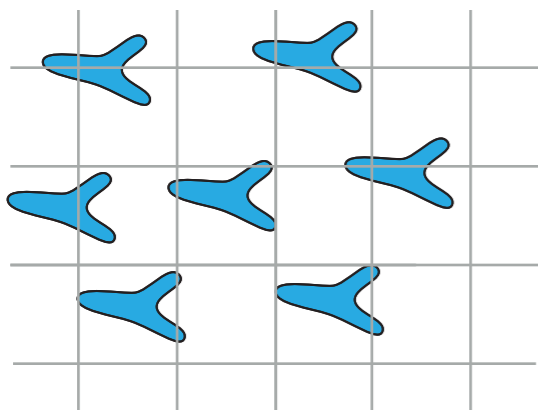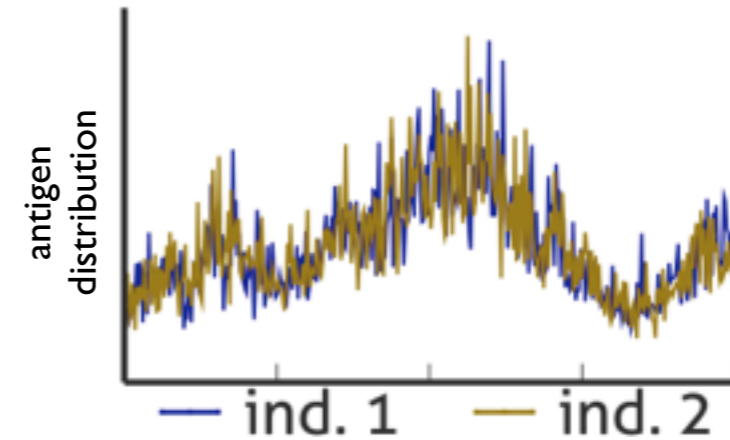
$q \sigma$

reproducible number of receptors in large space
*tracking of antigen*
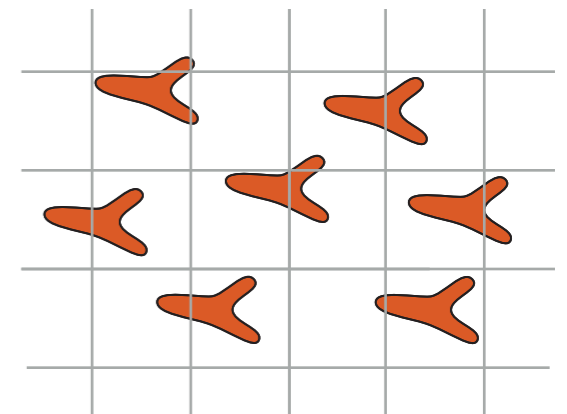
# Personalized responses

two individuals see the environment slightly differently
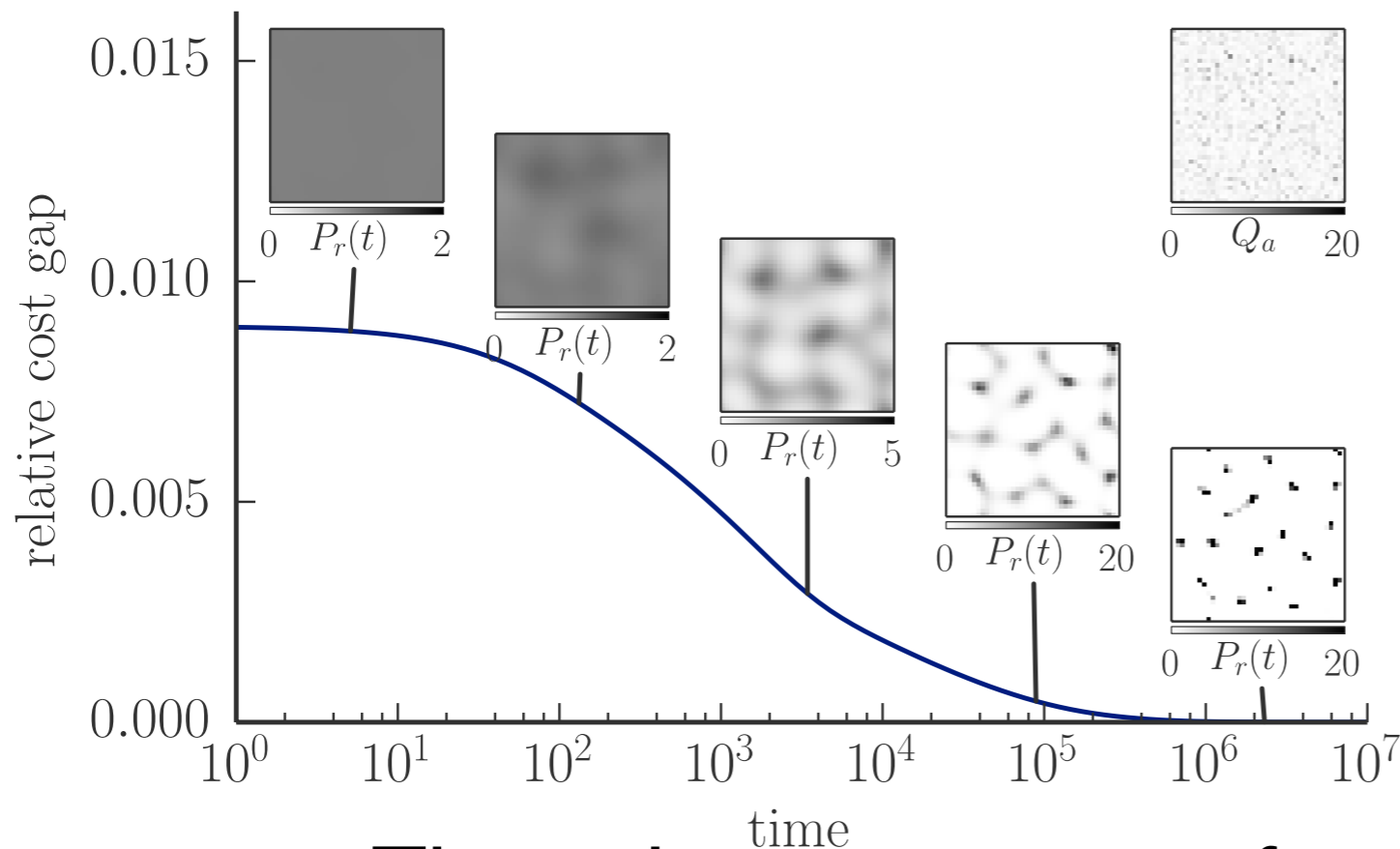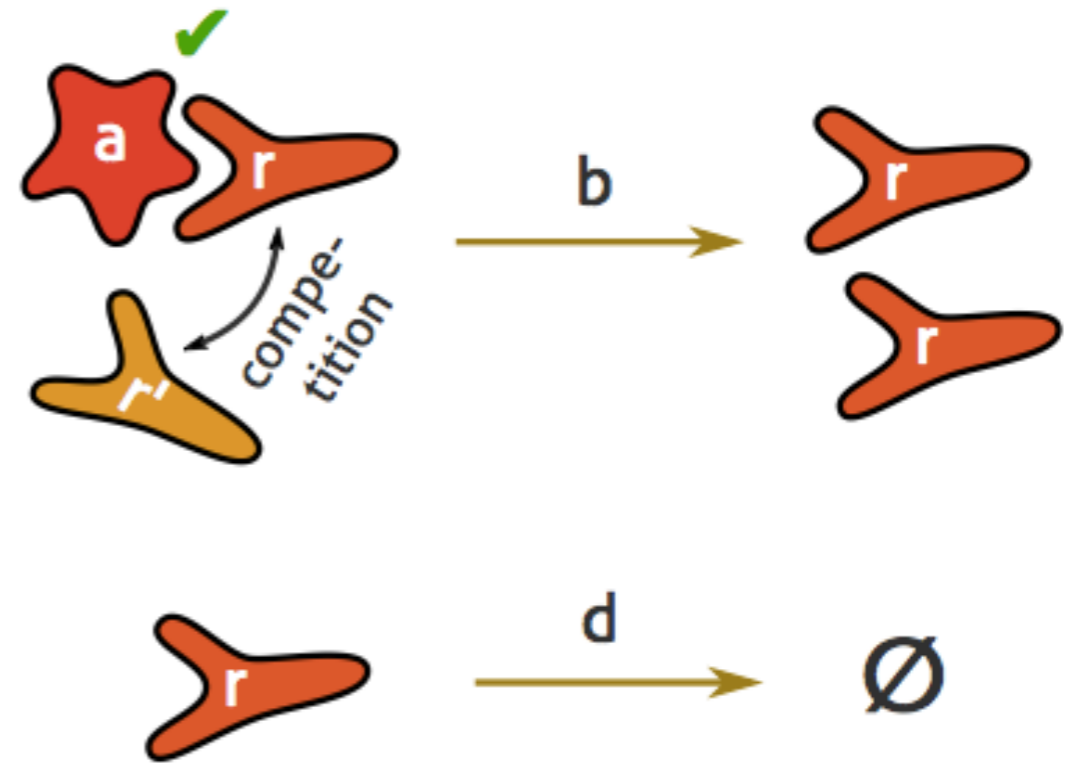


$\rightarrow$ very different repertoires

## Can optimal repertoires be reached via dynamics?
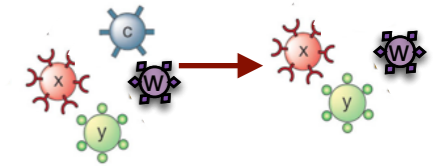
$$\dot{N}_r = N_r \left[ b \sum_p Q_p f_{r,a} \; A\left( \sum_{r'} N_{r'} f_{r',a} \right) - d \right]$$

population size

proliferation rate

detectable pathogen

availability of pathogen → reduced by competition

e.g. $A(\tilde{N}_a) = \frac{1}{(1+\tilde{N}_a)^2}$

death rate



**a** **r** **r'** competition **b** **r** **r**

**r** **d** Ø

relative cost gap

0.015

0.010

0.005

0.000

$0 \quad P_r(t) \quad 2$

$0 \quad P_r(t) \quad 2$

$0 \quad P_r(t) \quad 5$

$0 \quad P_r(t) \quad 20$

$0 \quad P_r(t) \quad 20$

$0 \quad Q_a \quad 20$

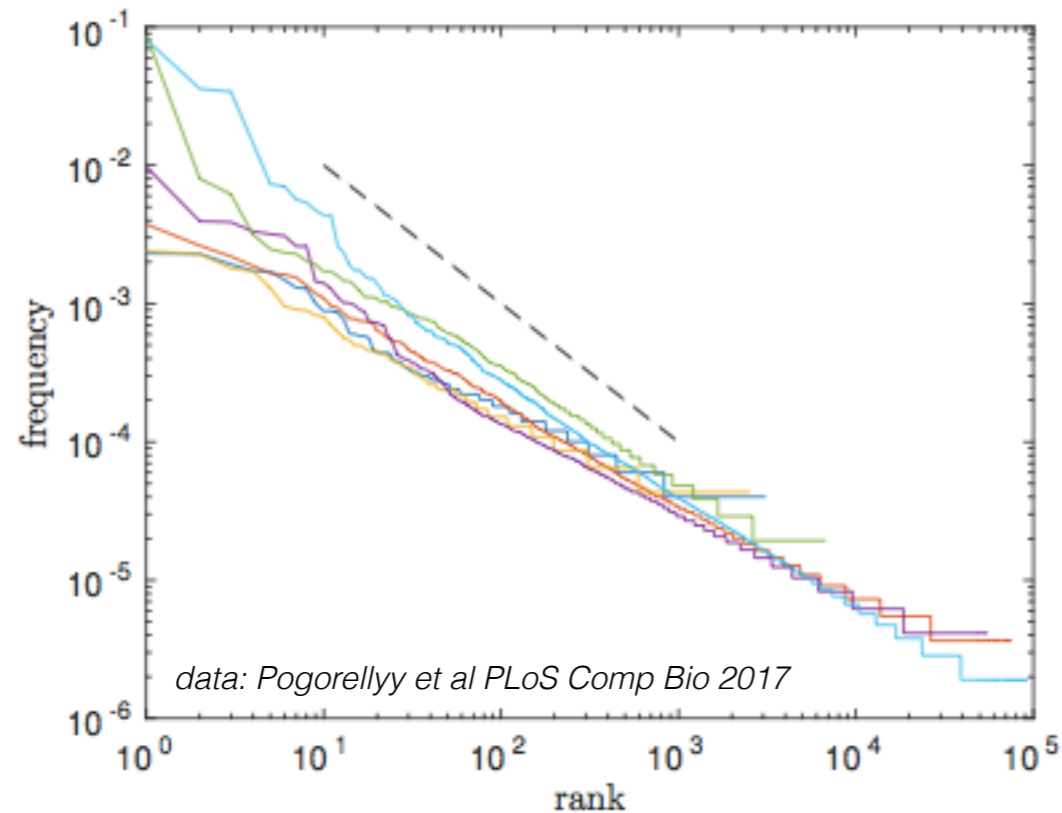$10^0 \quad 10^1 \quad 10^2 \quad 10^3 \quad 10^4 \quad 10^5 \quad 10^6 \quad 10^7$

time

→ Through competition of receptors for antigen

# Estimating frequencies

- trying to infer species frequencies

*human T-cells:*



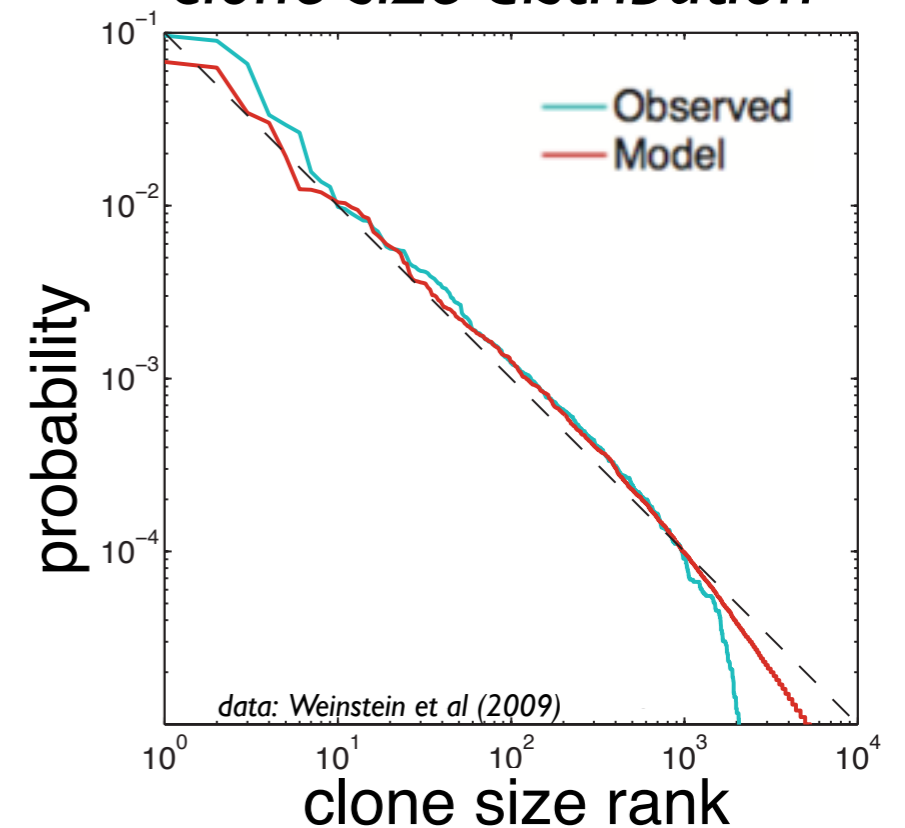data: Pogorellyy et al PLoS Comp Bio 2017
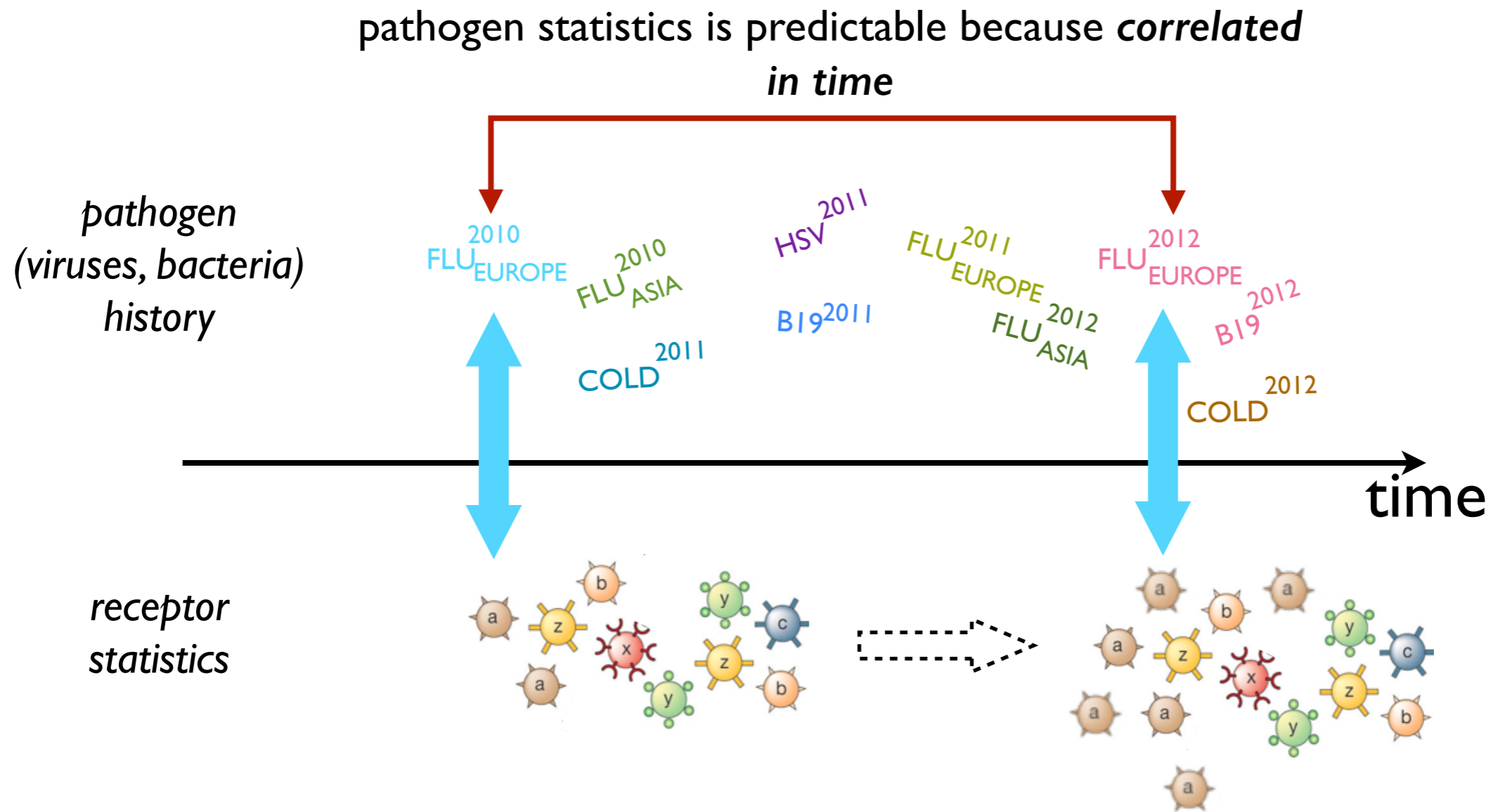
beta chain

Post-selection: 38 bits

*zebrafish B-cells:*
*clone size distribution*



data: Weinstein et al (2009)

clone size rank

- also in other environments (microbiome, ponds, forests)

# How to update receptor frequencies?

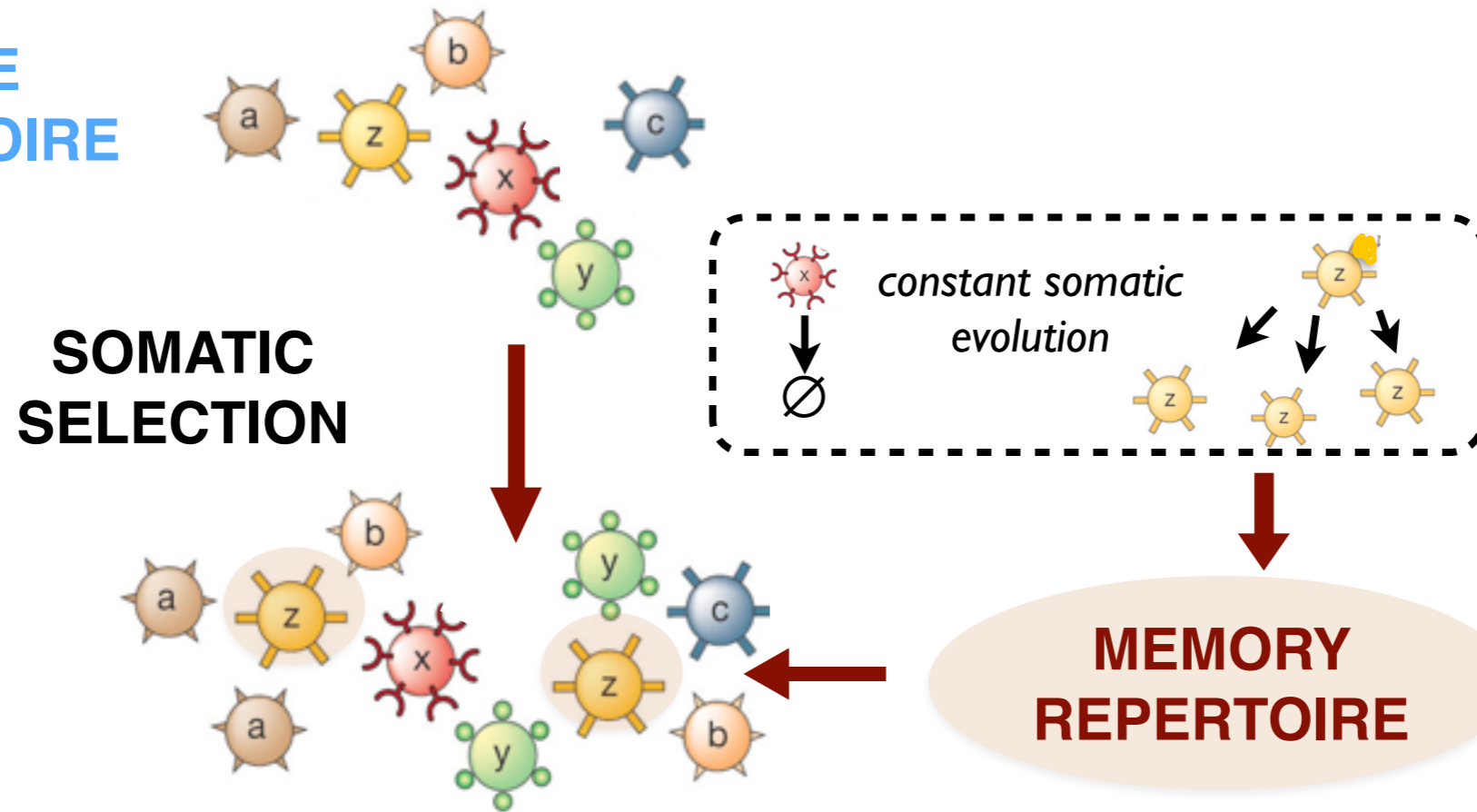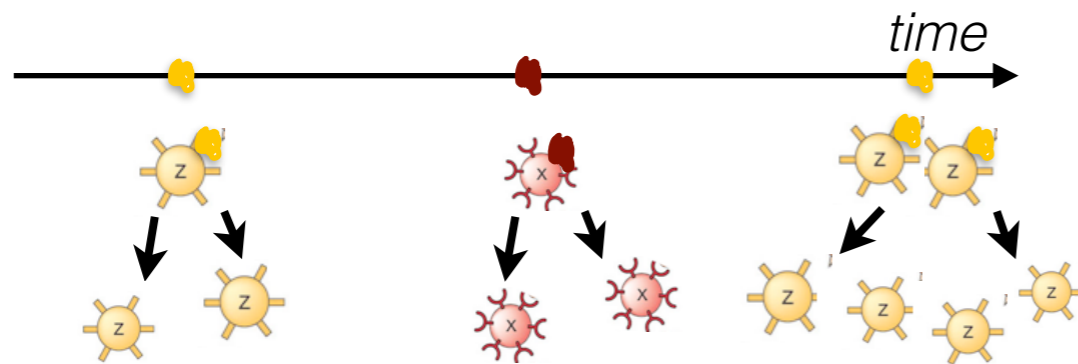- *adaptive* immune system - optimal predictor of future pathogens?



pathogen statistics is predictable because *correlated in time*

# How to remember?

**NAIVE REPERTOIRE**

**SOMATIC SELECTION**

*constant somatic evolution*
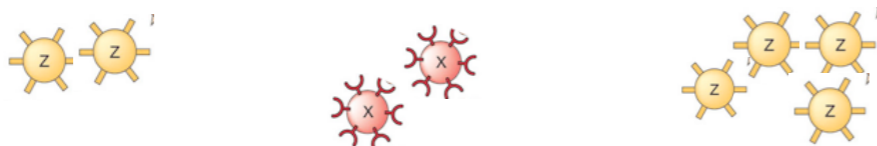
$\varnothing$

**MEMORY REPERTOIRE**

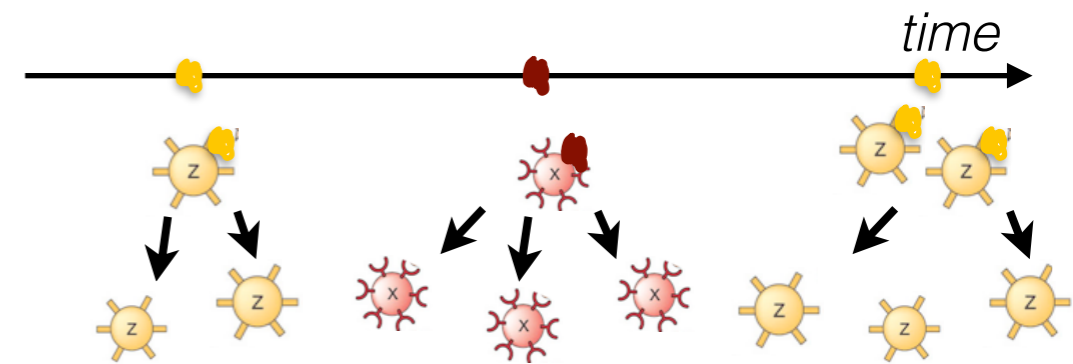is memory useful ?

*option 1:*

*time*

**proportional** *memory update:*

*option 2:*

*time*

**modulated** *memory update:*

# Estimate pathogen frequencies



$\lambda$ Poisson sampling rate of environment

$t_1$ $t_2$ $t_3$ $t_4$ $t_5$ $t_6$ time

pathogen frequencies $Q(t)$

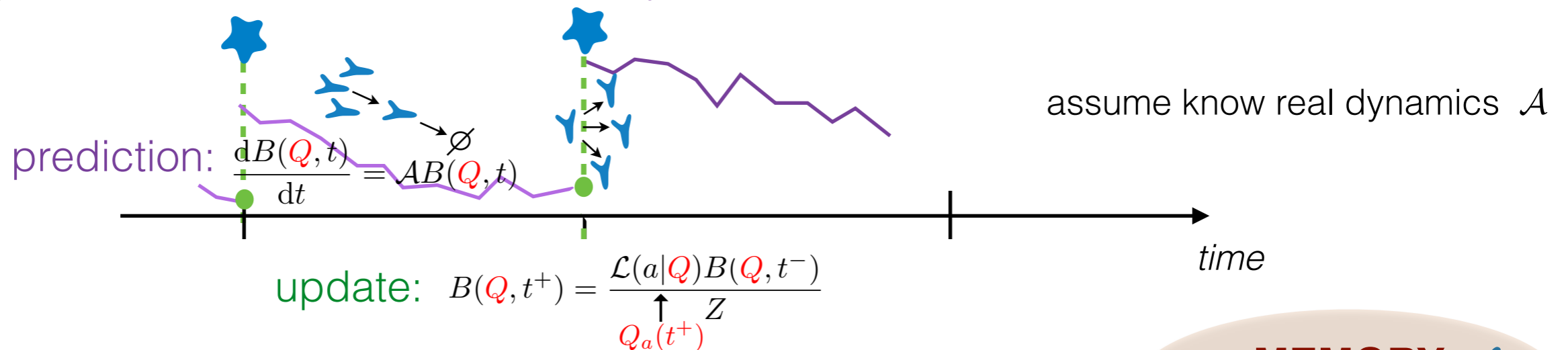receptor distribution $P(t)$

what receptor distribution P(t) minimizes expected harm?

belief of Q(t)

- Q(t) unknown → estimate it

$$\langle \text{Cost}\,(P(t), Q(t)) \rangle = \int dQ\, \text{Cost}\,(P(t), Q)\, B\,(Q, t) \xrightarrow{\text{min Cost}} P^\star(t) = f\,(\langle Q(t) \rangle)$$

expected cost of an infection

- propagate belief in time = encounters + prior

assume know real dynamics $\mathcal{A}$

prediction: $\dfrac{dB(Q,t)}{dt} = \mathcal{A}B(Q,t)$

time

update: $B(Q, t^+) = \dfrac{\mathcal{L}(a|Q)B(Q, t^-)}{Z}$

$Q_a(t^+)$

MEMORY REPERTOIRE

# Memory helps
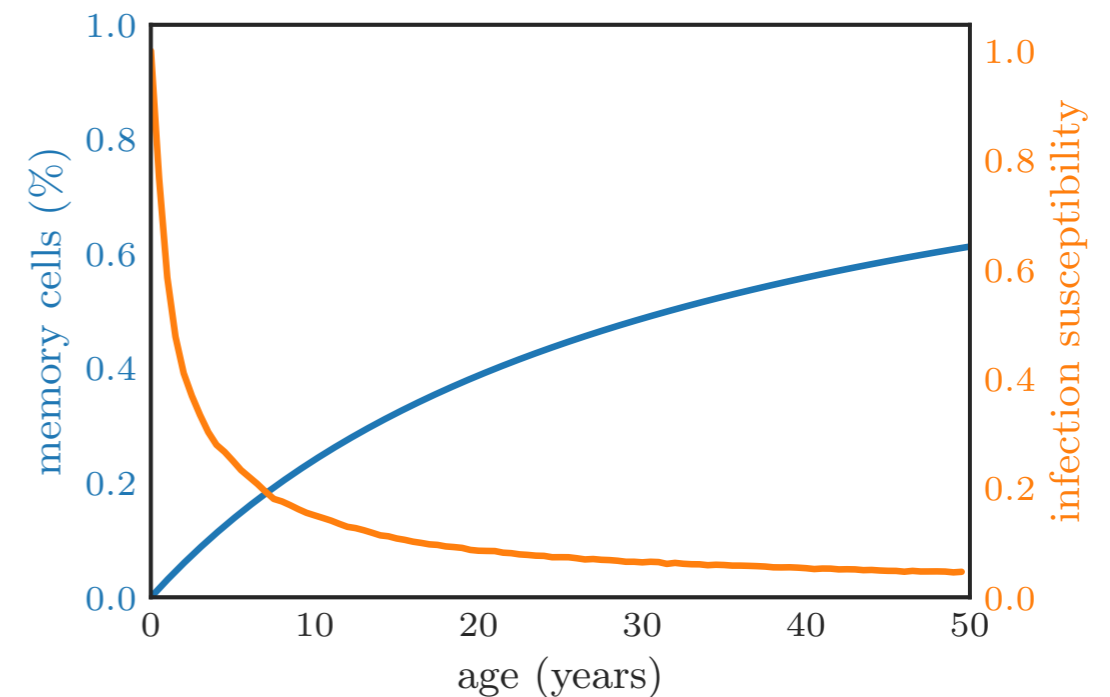
- memory helps in sparse environments

  → fast detection of few pathogens
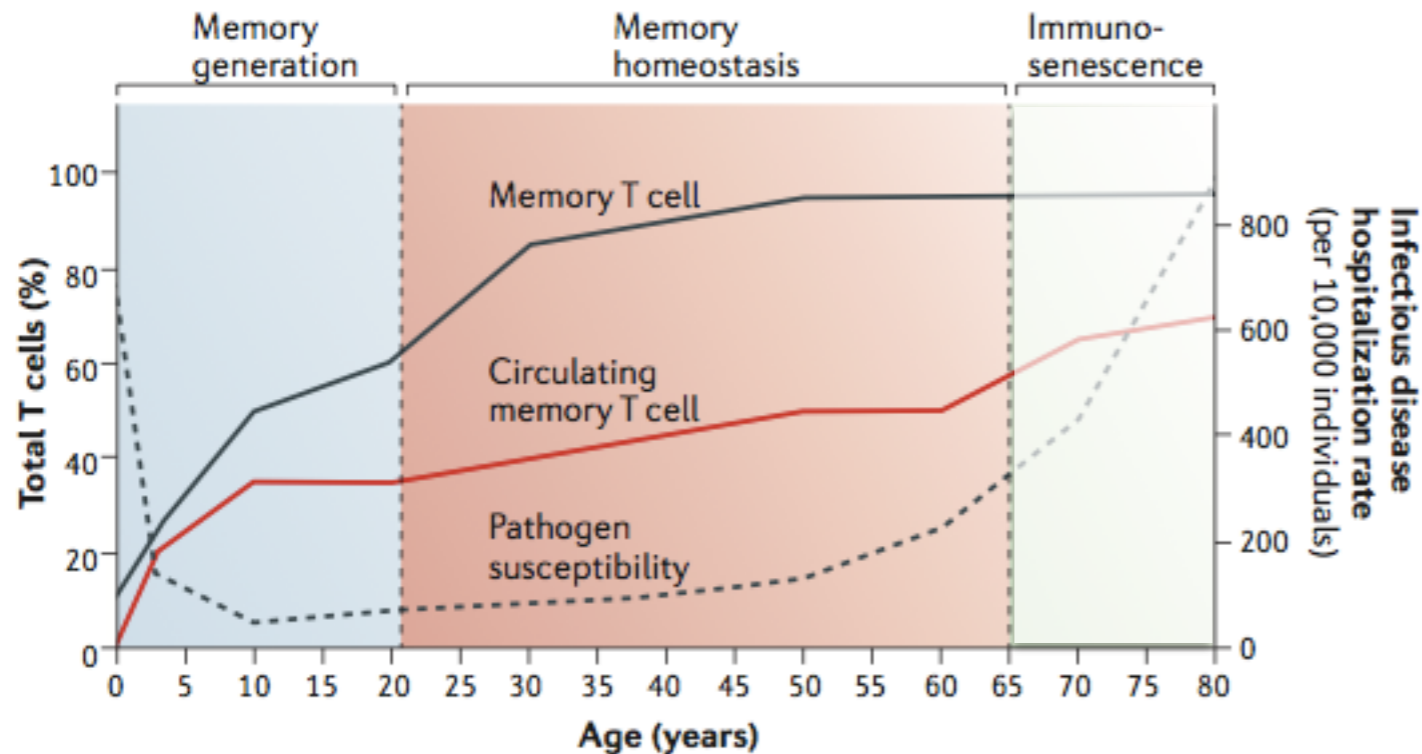


- advantage of memory - depends on sampling

→control theory rationalizes existence of immunological memory

- quickly learn global features of the distribution
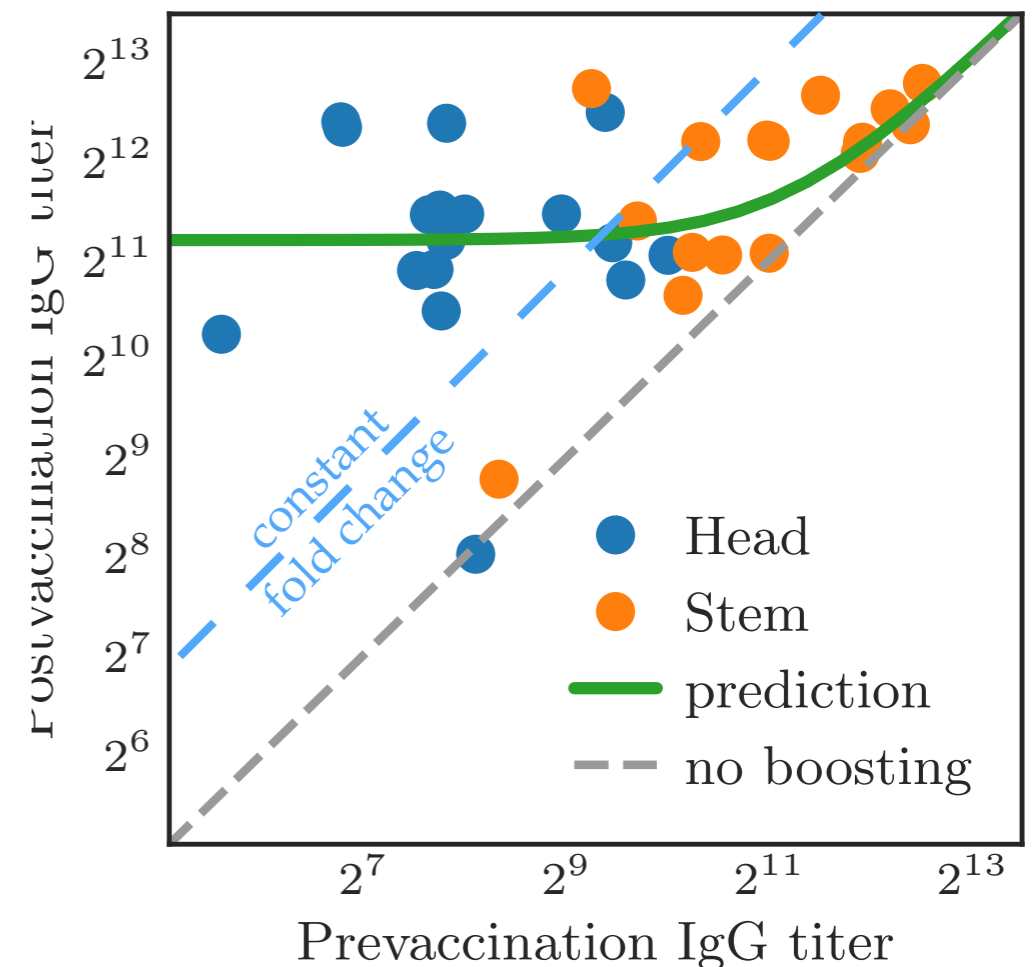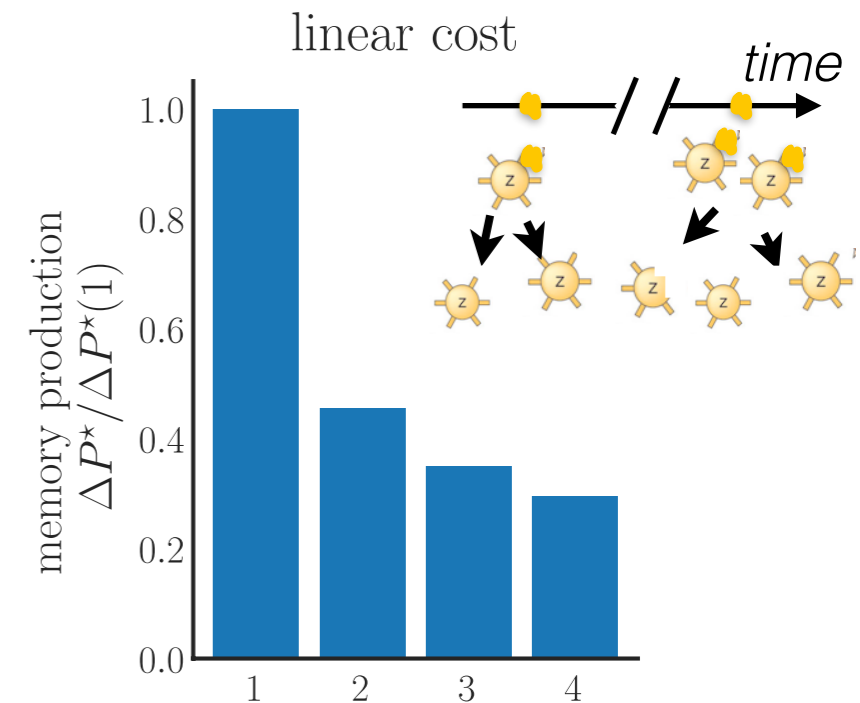


- predictive learning reproduces experimental features

- subsequent observations count as less evidence



→ *vaccination*

- booster vaccination titers for epitopes of hemaglutanin following vaccination with inactivated H5N1
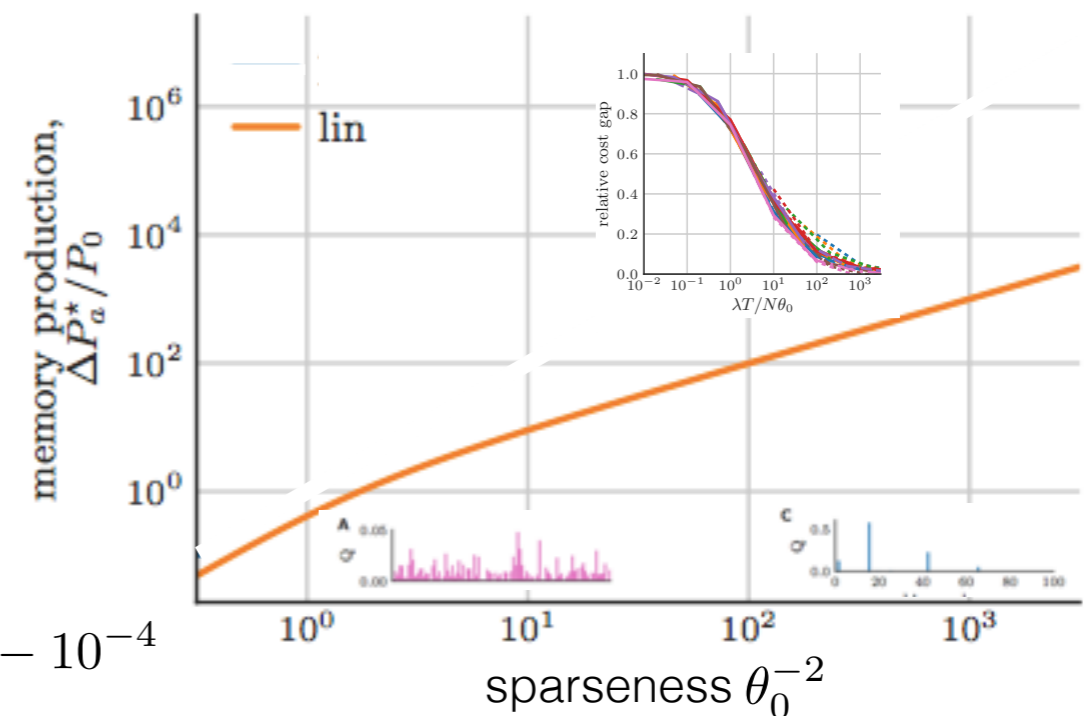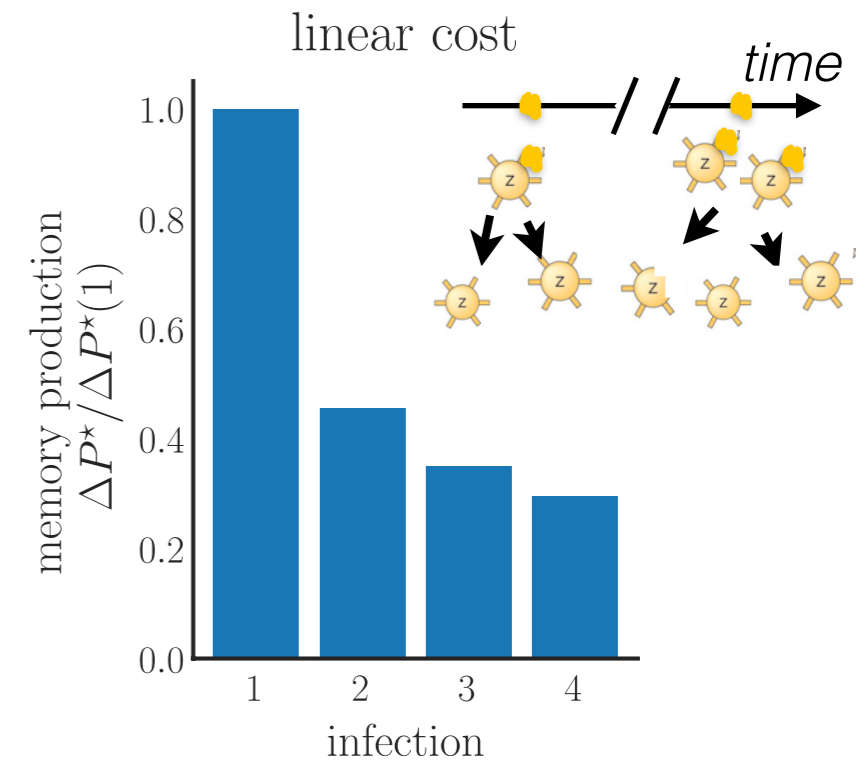


*data from Ellebedy et al. PNAS 2014*

# Learning

- subsequent observations count as less evidence



- stronger response in sparse environments

memory increase ~100-1000*



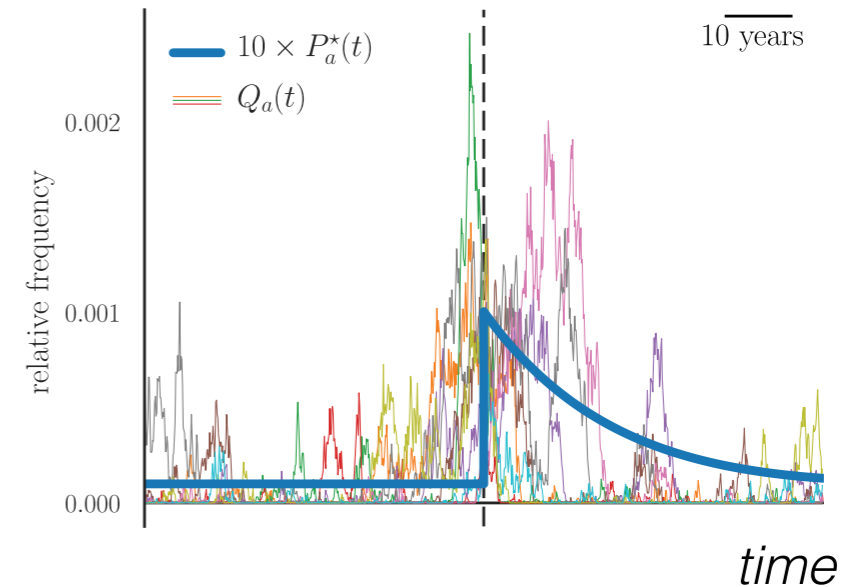→ very sparse environment $\theta_0 \sim 10^{-6} - 10^{-4}$

→ pathogen seen once - memory gives ~2 fold cost decrease
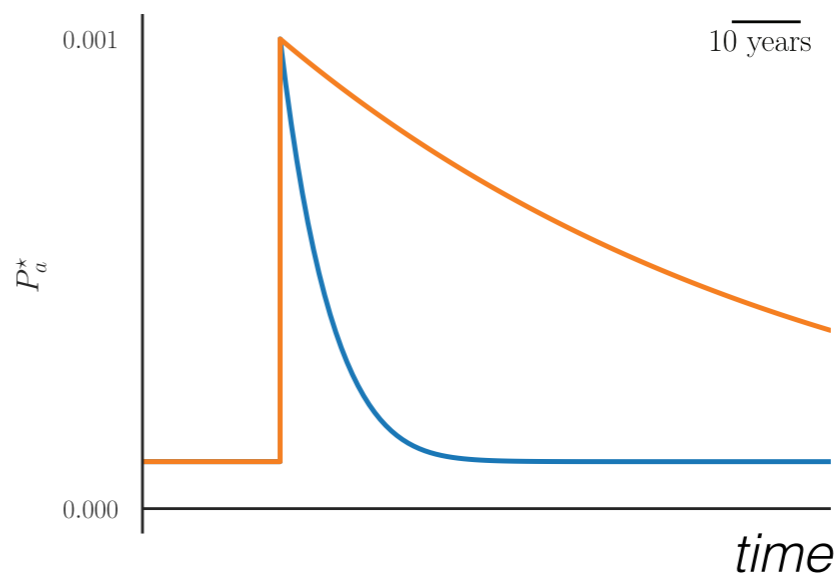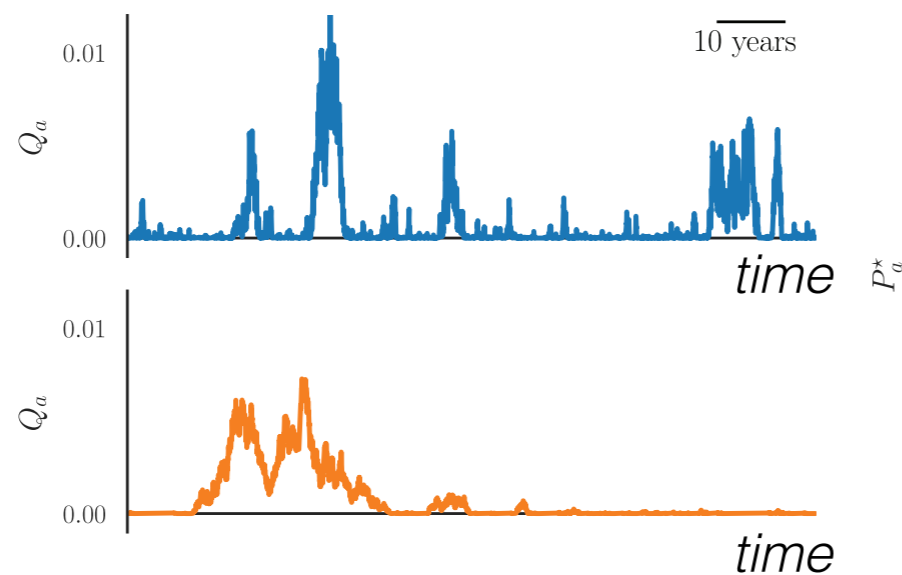
*Buchholz et al (2016) An. rev. imm.

# Forgetting

- changing environment → propagate belief between sampling

- Bayesian prediction of changing environment



- forget faster in rapidly changing environments
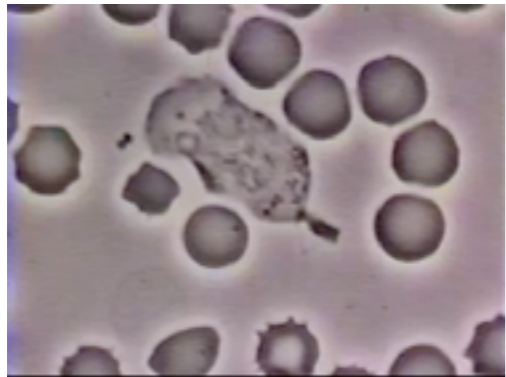
# Different immune strategies



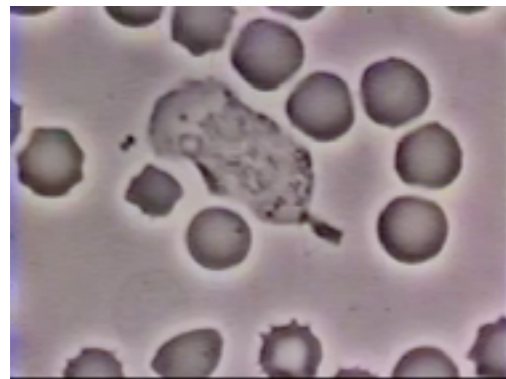adaptive immunity

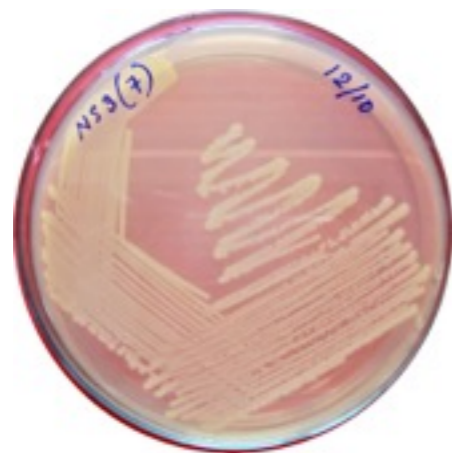# Other immune strategies



innate immunity



adaptive immunity



CRISPR immunity

# Common strategic choices



**randomly acquired**

*regulated*

innate immunity

adaptive immunity

**heritable**

**non-heritable**

*constitutive*

**actively acquired**

CRISPR immunity

Processing information about the environment on **evolutionary** timescales

Response during **organism** lifetime

# Common strategic choices



**randomly acquired**

*regulate*

in                    ty
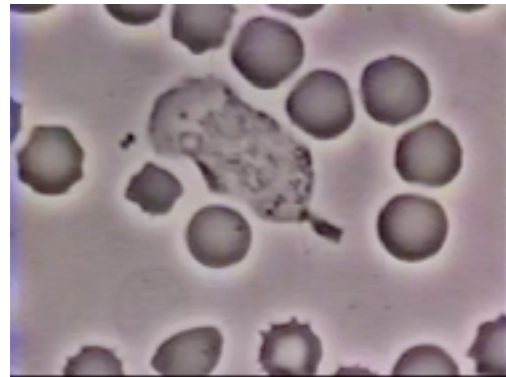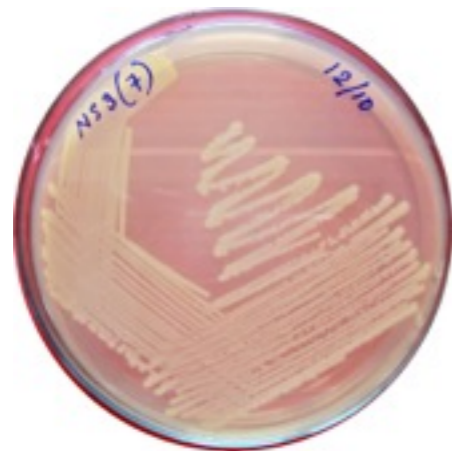
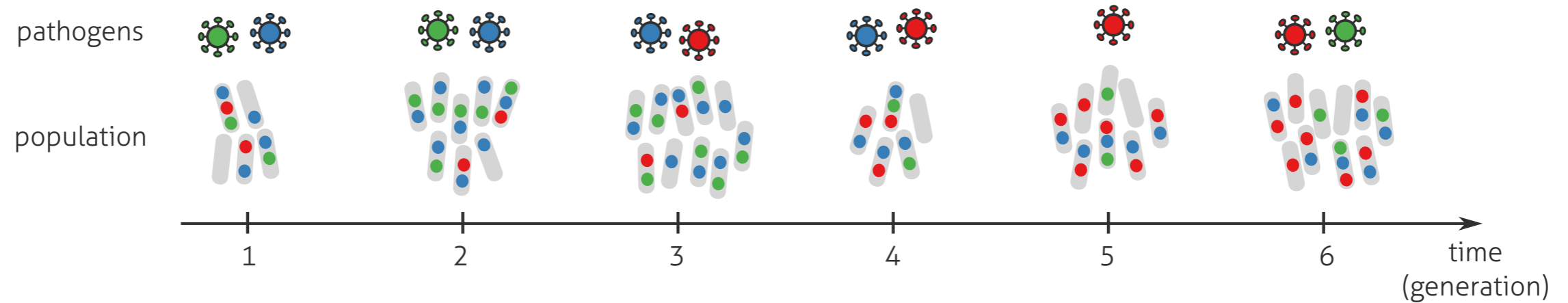## optimal immunity?

*constitutive*

**actively acquired**

Processing information about the environment on **evolutionary** timescales
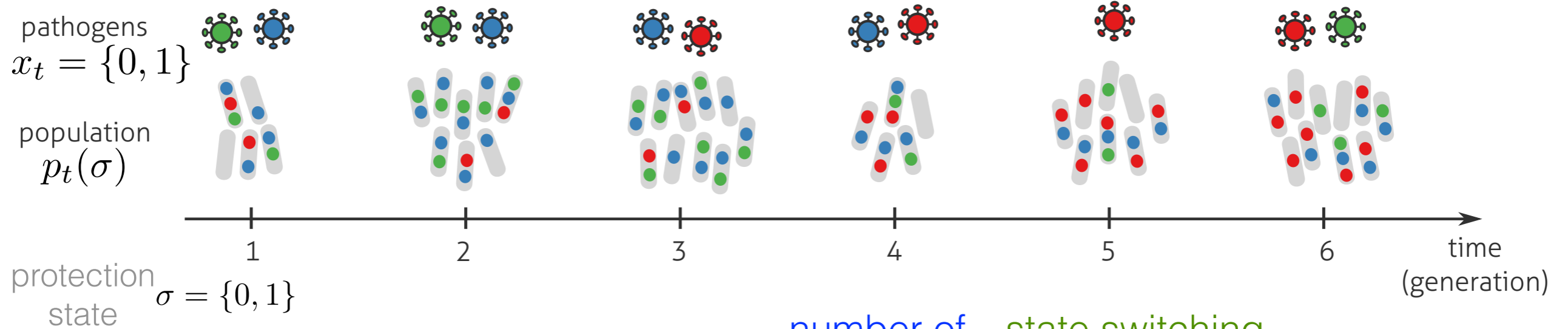
Response during **organism** lifetime

CRISPR immunity

# Optimal immunity



- match environment statistics

- ensure long term population growth

→ immunity as adaptation to pathogen statistics

- consider different strategies

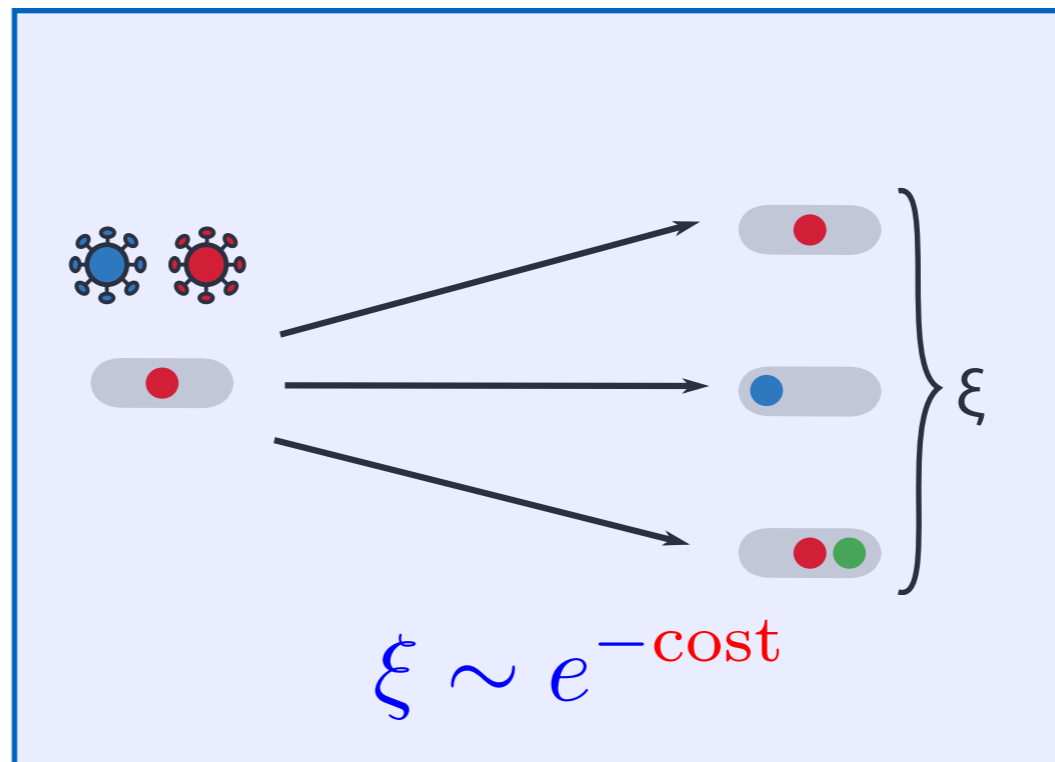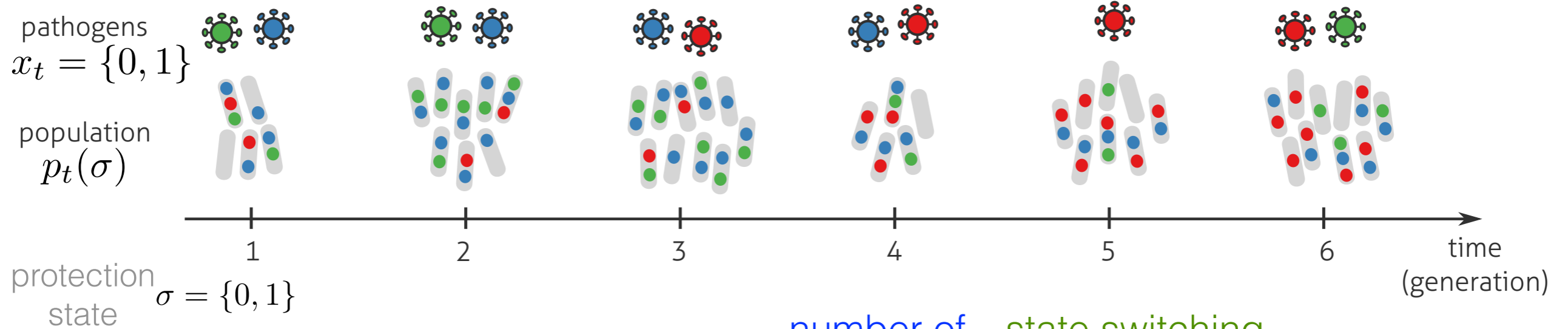- optimize long term population growth

# Population growth



pathogens
$x_t = \{0, 1\}$

population
$p_t(\sigma)$

protection
state $\quad \sigma = \{0, 1\}$

time
(generation)

number of
offspring

state switching
probability

$$p_{t+1}(\sigma) \;=\; \frac{1}{Z_t} \sum_{\sigma'} \xi(\sigma', x_t) \, \pi(\sigma | \sigma', x_t) \, p_t(\sigma')$$

$\xi \sim e^{-\text{cost}}$

# Population growth



$$p_{t+1}(\sigma) = \frac{1}{Z_t} \sum_{\sigma'} \underset{\substack{\text{number of} \\ \text{offspring}}}{\xi(\sigma', x_t)} \underset{\substack{\text{state switching} \\ \text{probability}}}{\pi(\sigma|\sigma', x_t)} p_t(\sigma')$$

Stochastic environment

$$e^{-\tau_{\text{env}}} = 1 - \alpha - \beta$$

$$\pi_{\text{env}} = \alpha/(\alpha + \beta)$$

Maximize long term growth rate

$$\Lambda(\text{strategy}, \text{environment}) = \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T} \log(Z_t)$$
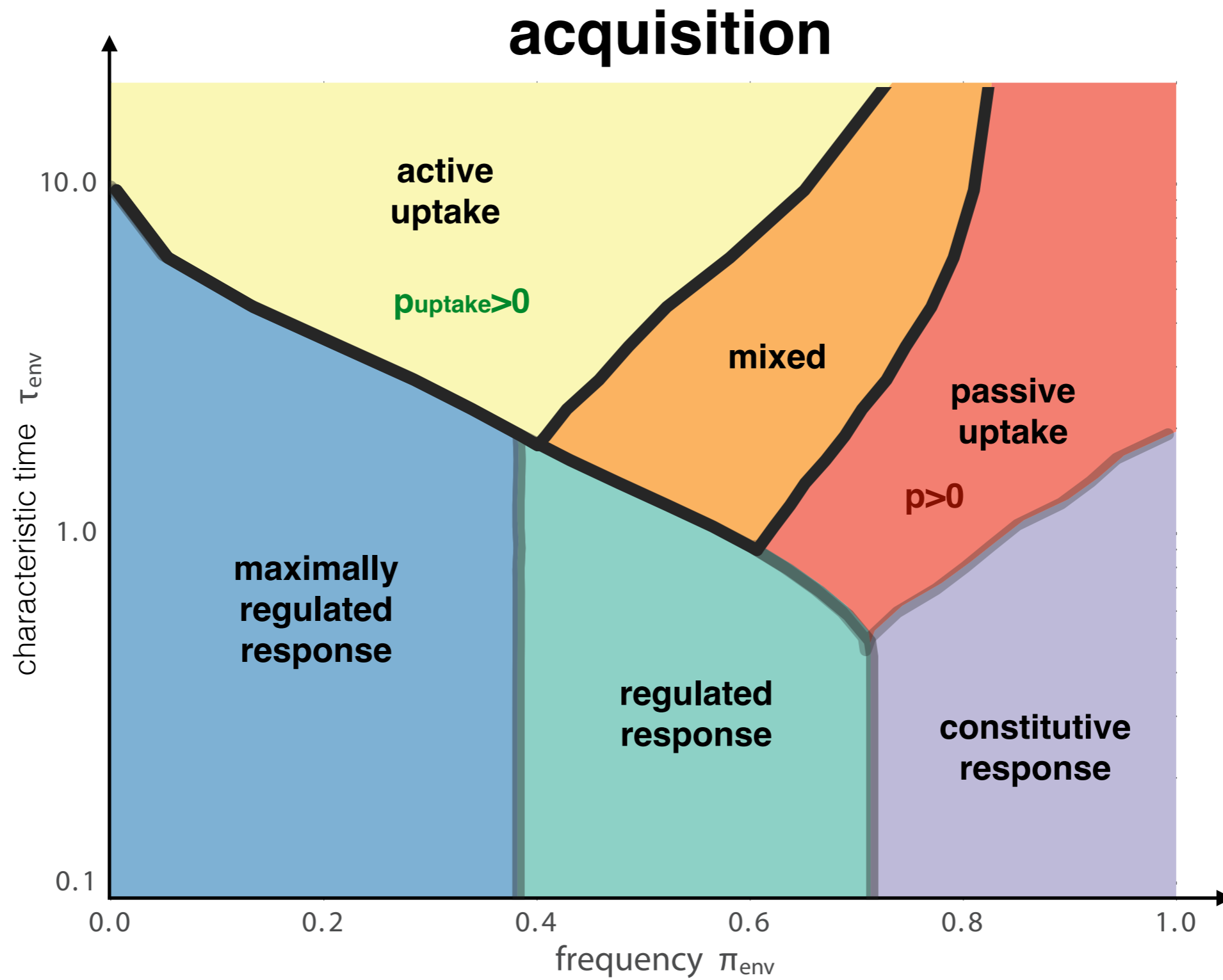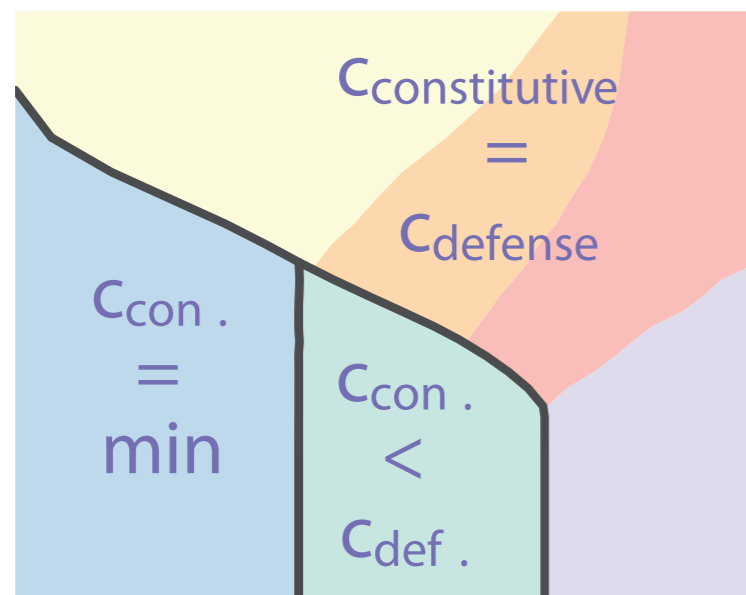
# Optimal strategies

# Optimal strategies

# Three strategy axes



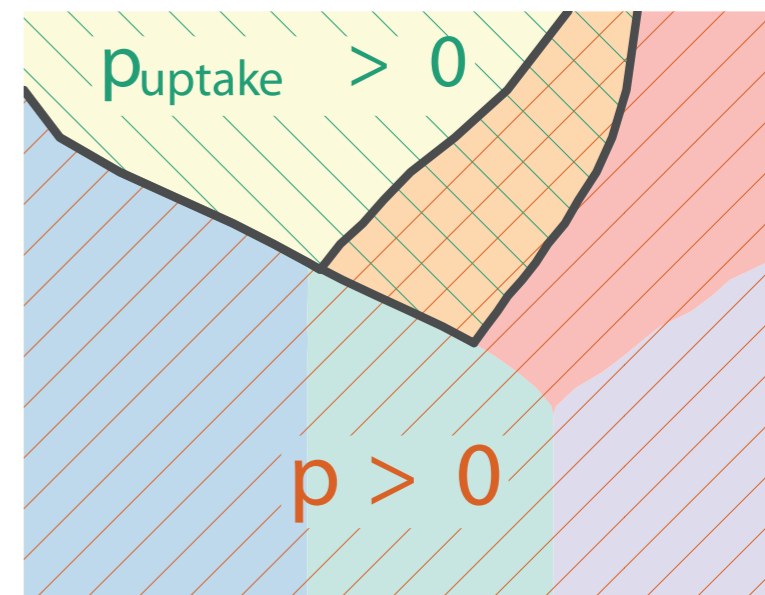adaptability
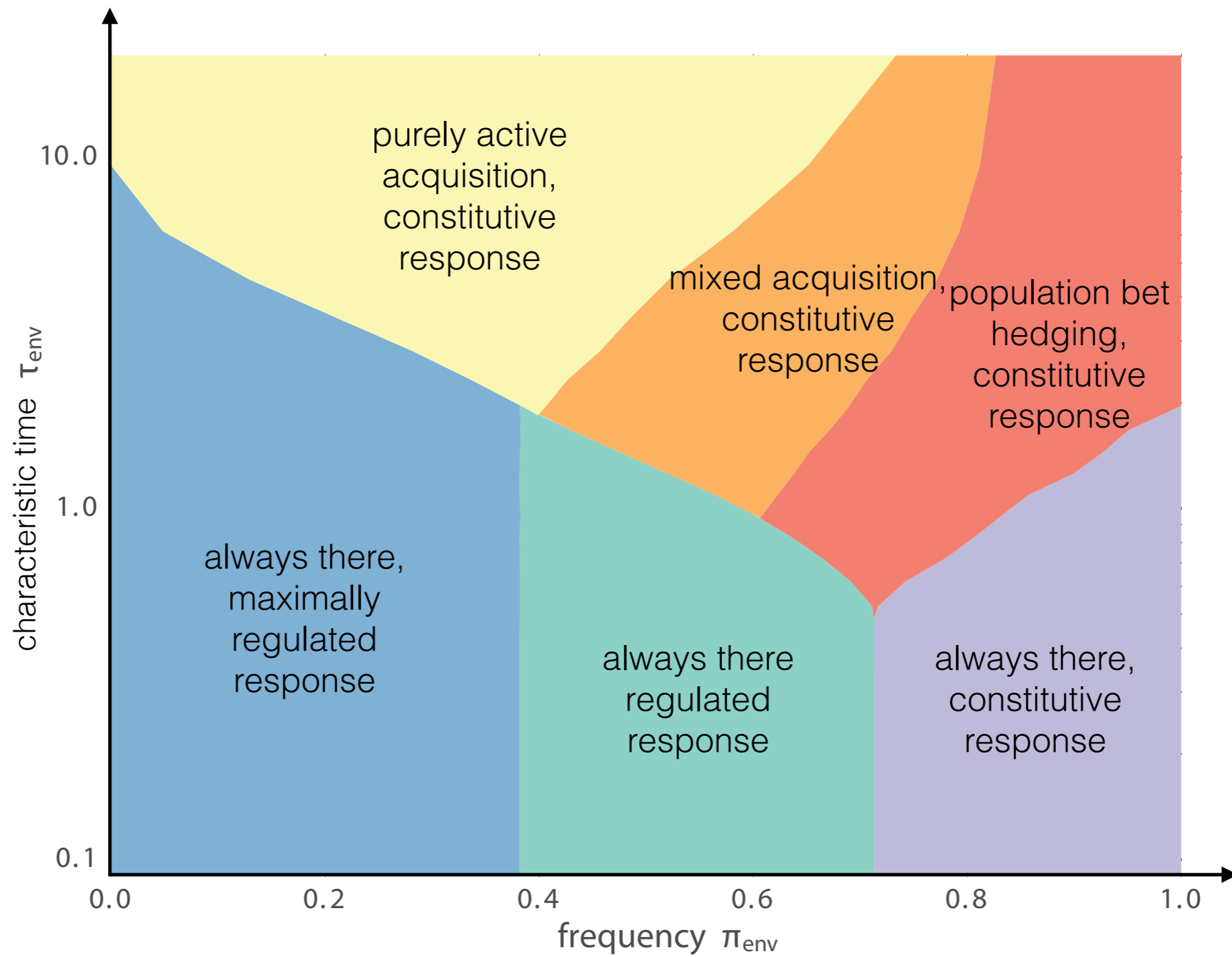
$c_{constitutive} = c_{defense}$

$c_{con.} = min$

$c_{con.} < c_{def.}$

heritability

$q > 0$

$q = 0$

acquisition mode

$p_{uptake} > 0$

$p > 0$
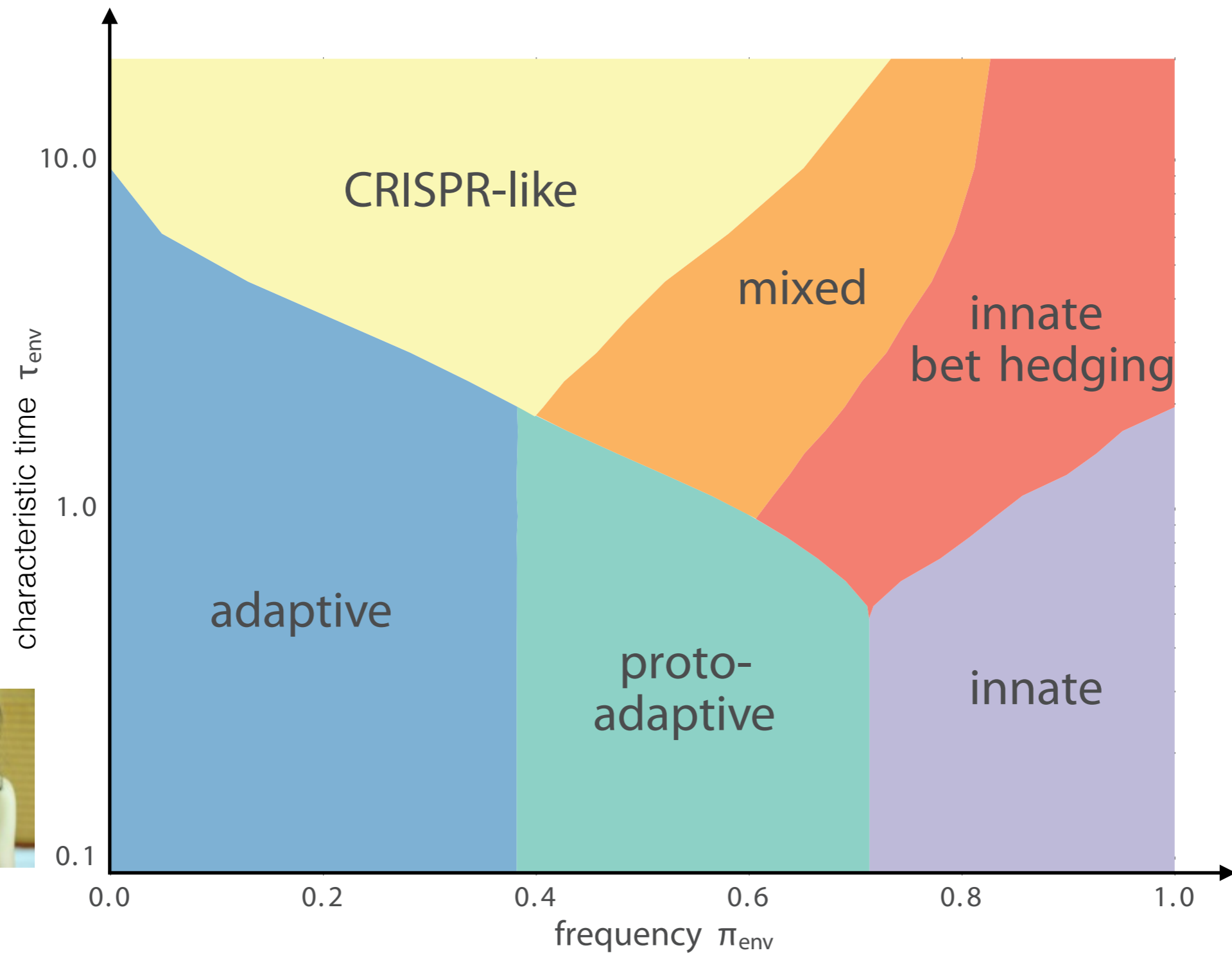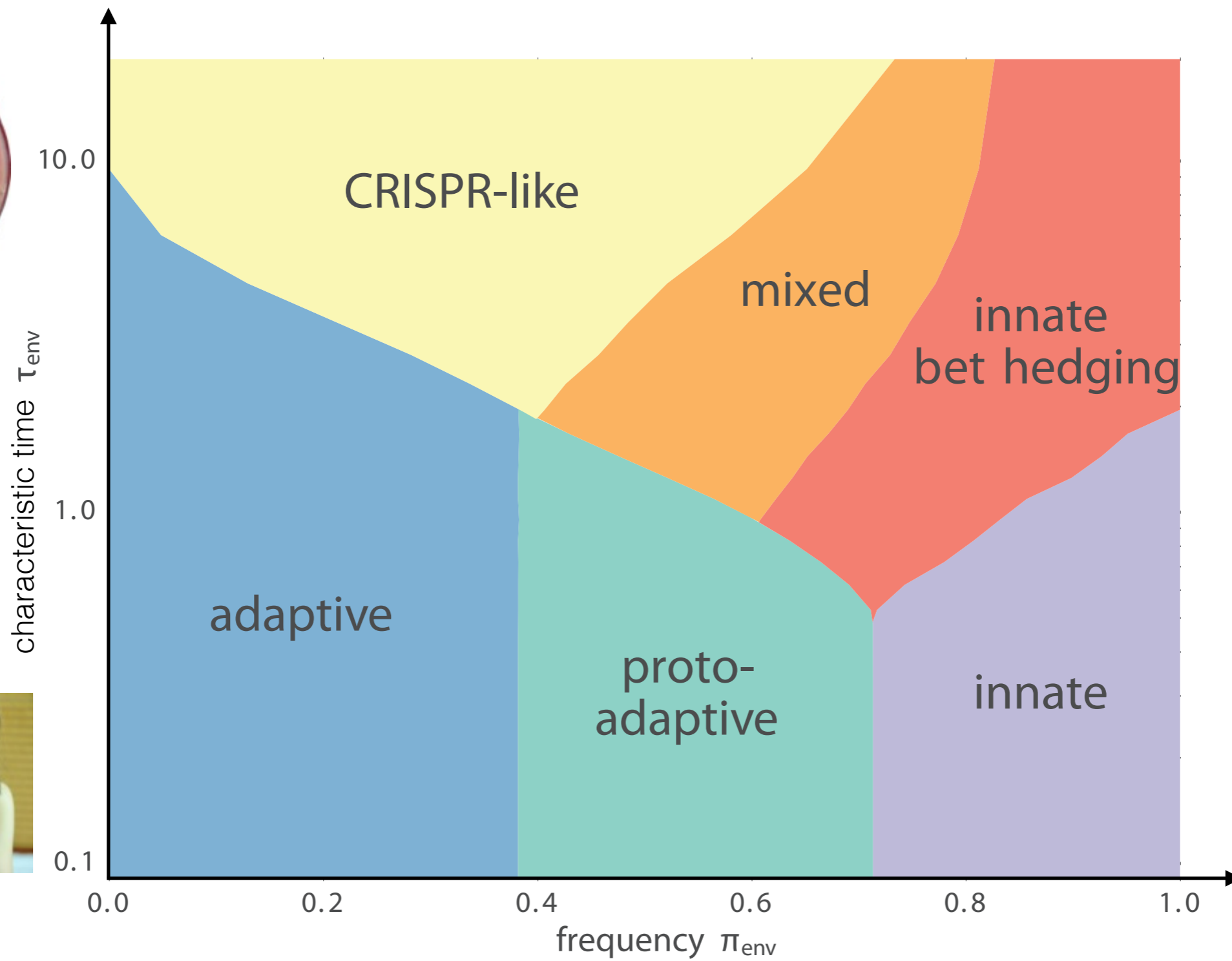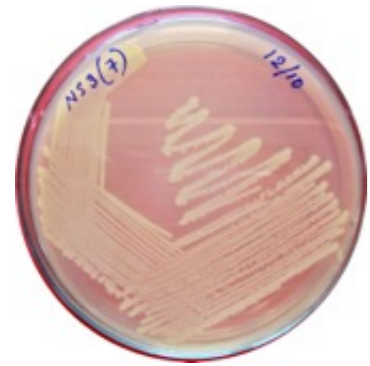
# Optimal strategies

# Optimal immune systems

# Optimal immune systems

# Optimal immune systems

## generating diversity and selection:

- random overlap between (most) individuals
- very long lived clones

## optimal repertoires:

- cover space but are random
- differ in two individuals

## predicting immune systems:

- use dynamics to anticipate frequencies
- memory useful in sparse environments

## optimal immunity:

- known immunity from evolutionary constraints
- depends on environment statistic

T Mora, AM Walczak, W Bialek, CG Callan, PNAS (2010)
A Murugan, T Mora, AM Walczak, CG Callan, PNAS (2012)
Y Elhanati, A Murugan, CG Callan, T Mora, AM Walczak, PNAS (2014)
A Mayer, V Balasubramanian, T Mora, AM Walczak, PNAS (2015)
Y Elhanati, Z Sethna, Q Marcou, CG Callan, T Mora, AM Walczak, Phil. Trans. B (2015)
J. Desponds, T. Mora, AM Walczak, PNAS (2015)
A Mayer, T Mora, O Rivoire, AM Walczak, PNAS (2016)
Y Elhanati, Q Marcou, T Mora, AM Walczak, Bioinformatics (2016)
RM Adams, JB Kinney, T Mora, AM Walczak, eLife (2017)
M. Pogorelyy et al, PloS CB (2017)
T Mora, AM Walczak, qbio/bioarxiv (2016)
Z Sethna, Y Elhanati, CG Callan, T Mora, AM Walczak, PNAS (2017)
M Laessig, V Mustonenen, AM Walczak Nature Ecology & Evolution (2017)
Q. Marcou. T. Mora, AM Walczak, qbio/bioarxiv (2017)
M. Pogorelyy et al, qbio/bioarxiv (2017)