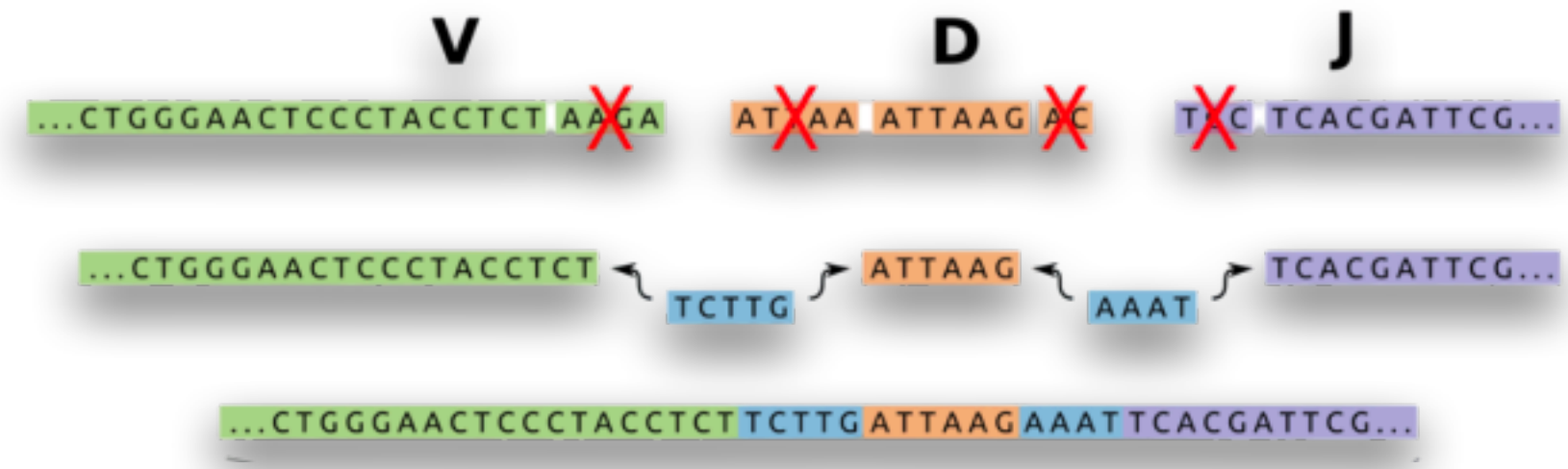# Modelling and predicting the overlap of B- and T-cell receptor repertoires in healthy and SARS-CoV-2 infected individuals

**María Ruiz Ortega,**

**Aleksandra Walczak and Thierry Mora**
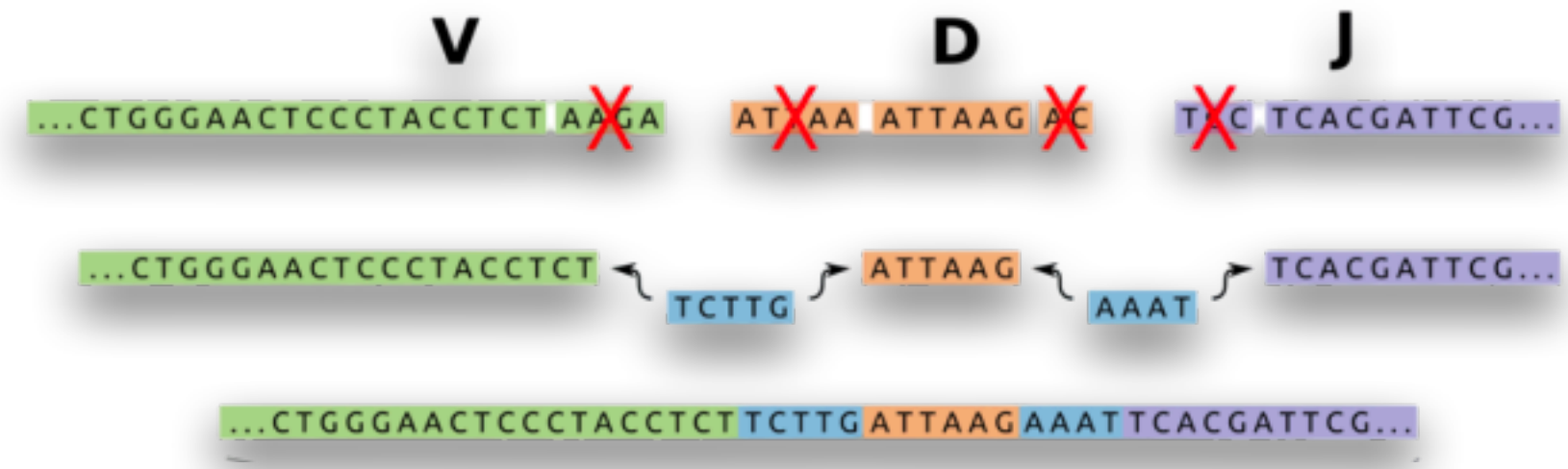
*Rencontre des Jeunes Physicien.ne.s 2022*

ENS | PSL ★

PLOS Computational Biology | DOI:10.1371/journal.pcbi.1004409 January 11, 2016

PLOS Computational Biology | DOI:10.1371/journal.pcbi.1004409 January 11, 2016

$10^{61}$ different sequences

**IgM**

$10^{12}$ unique naive sequences

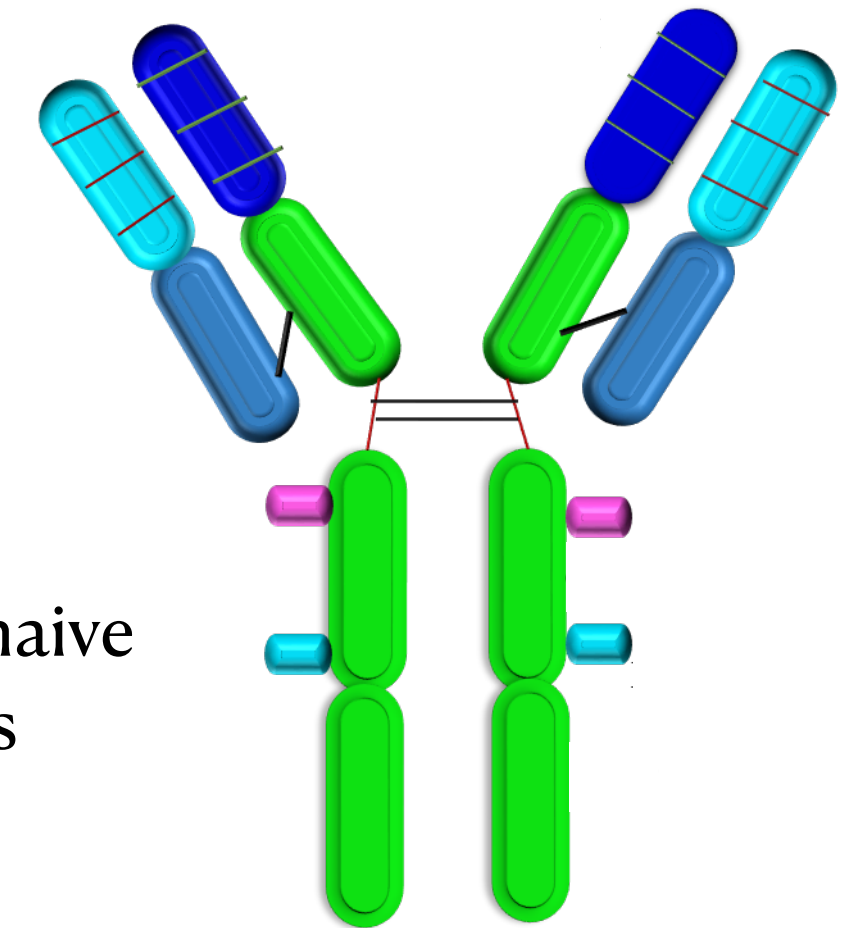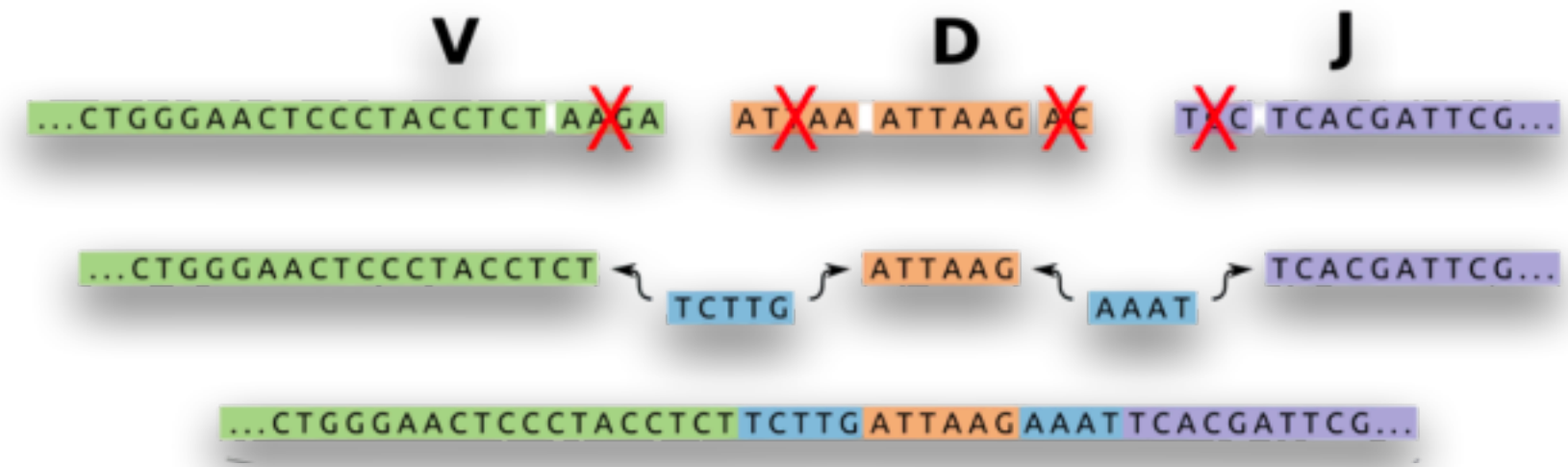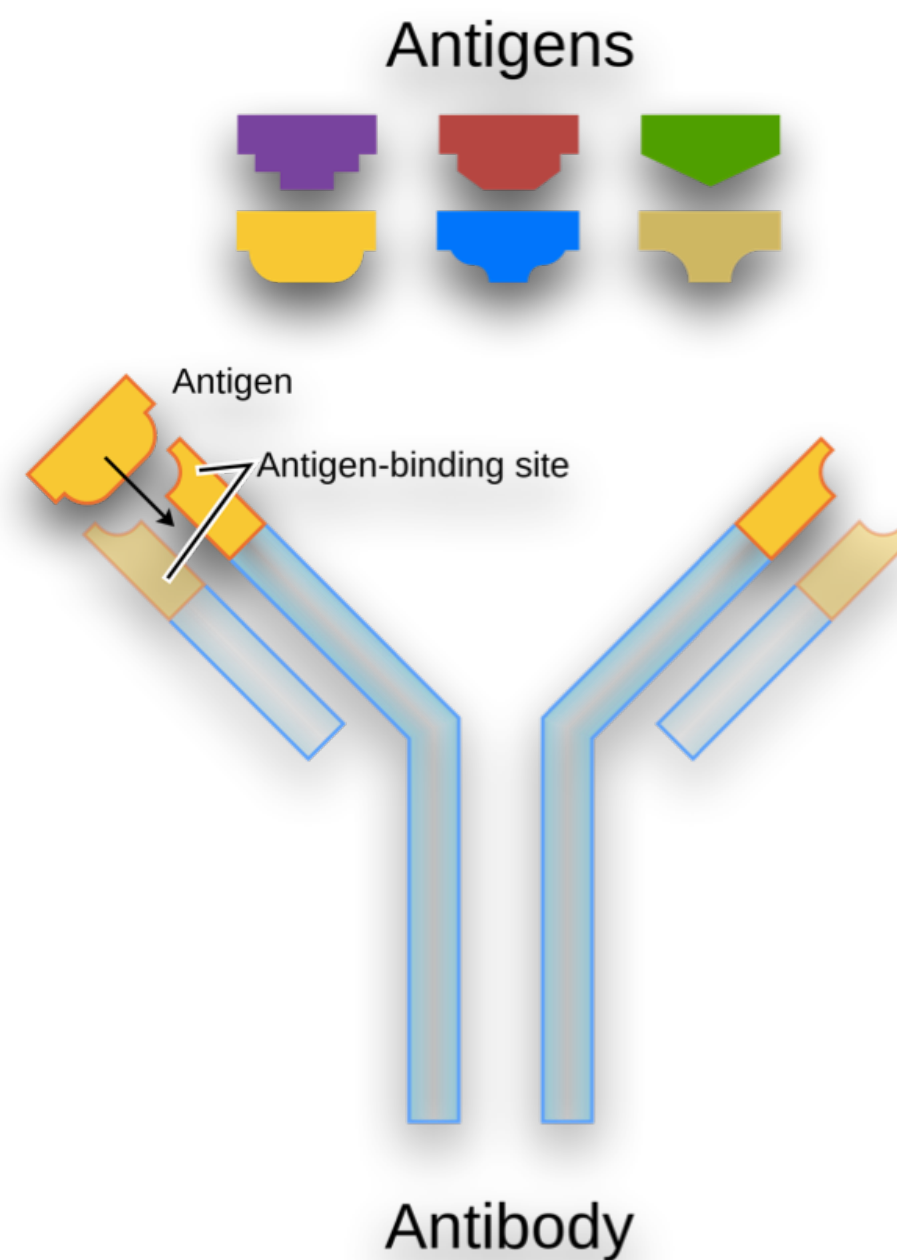PLOS Computational Biology | DOI:10.1371/journal.pcbi.1004409 January 11, 2016

$10^{61}$ different sequences

IgM

$10^{12}$ unique naive sequences

Antigens

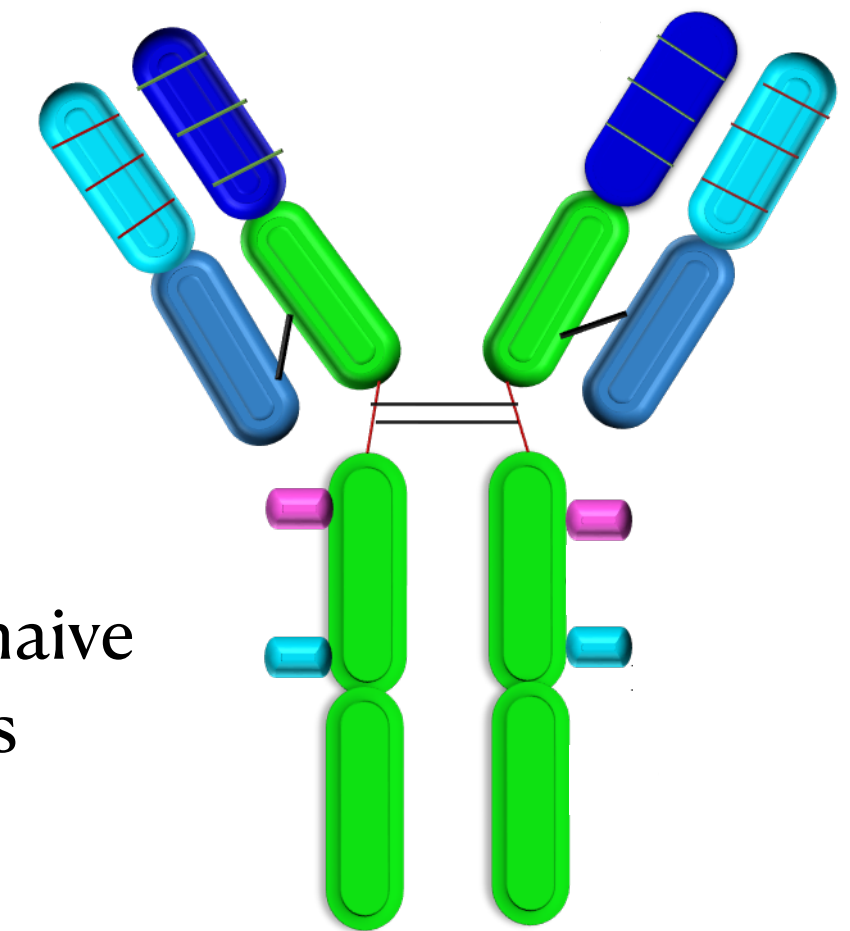Antigen

Antigen-binding site

Antibody

PLOS Computational Biology | DOI:10.1371/journal.pcbi.1004409 January 11, 2016

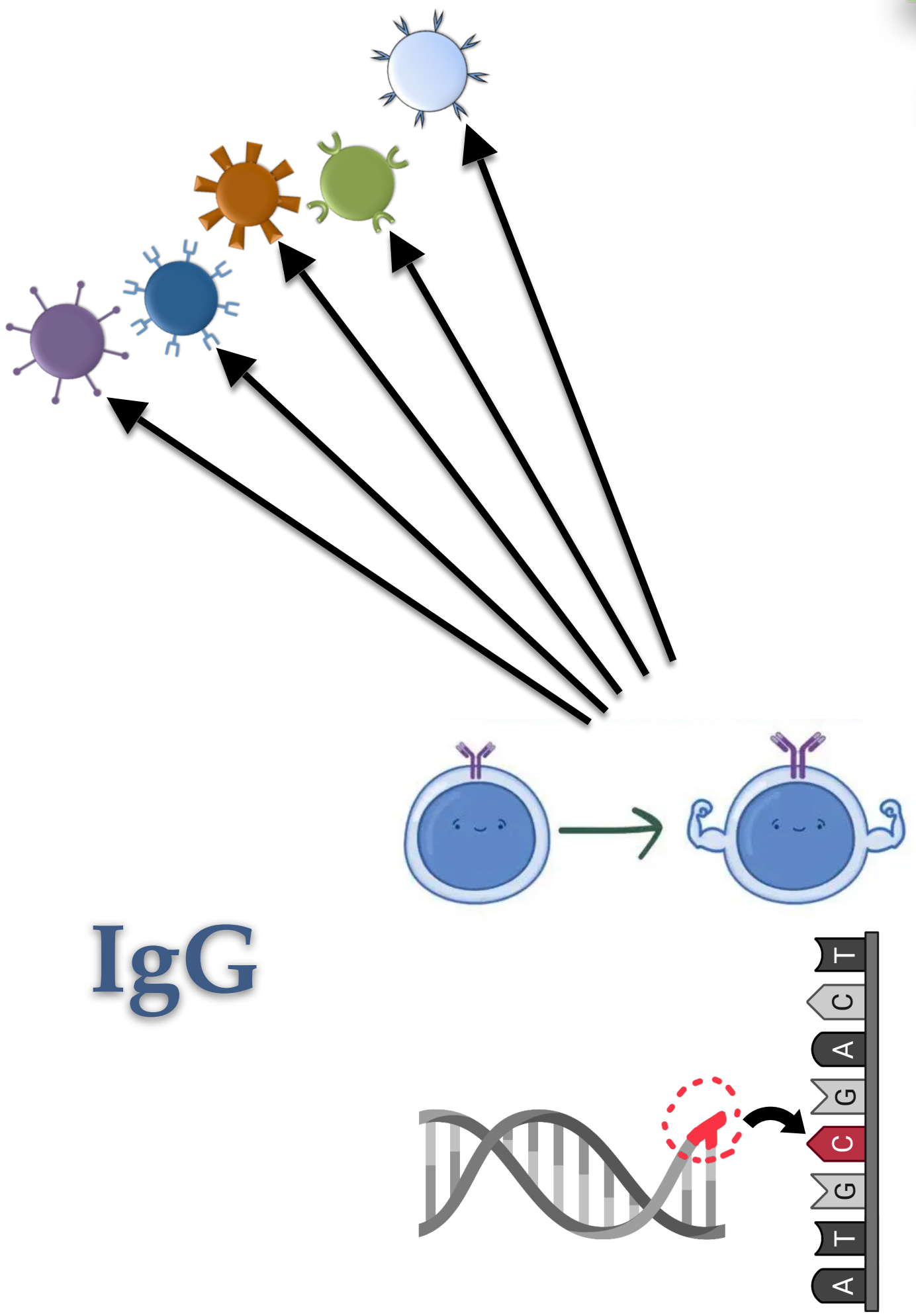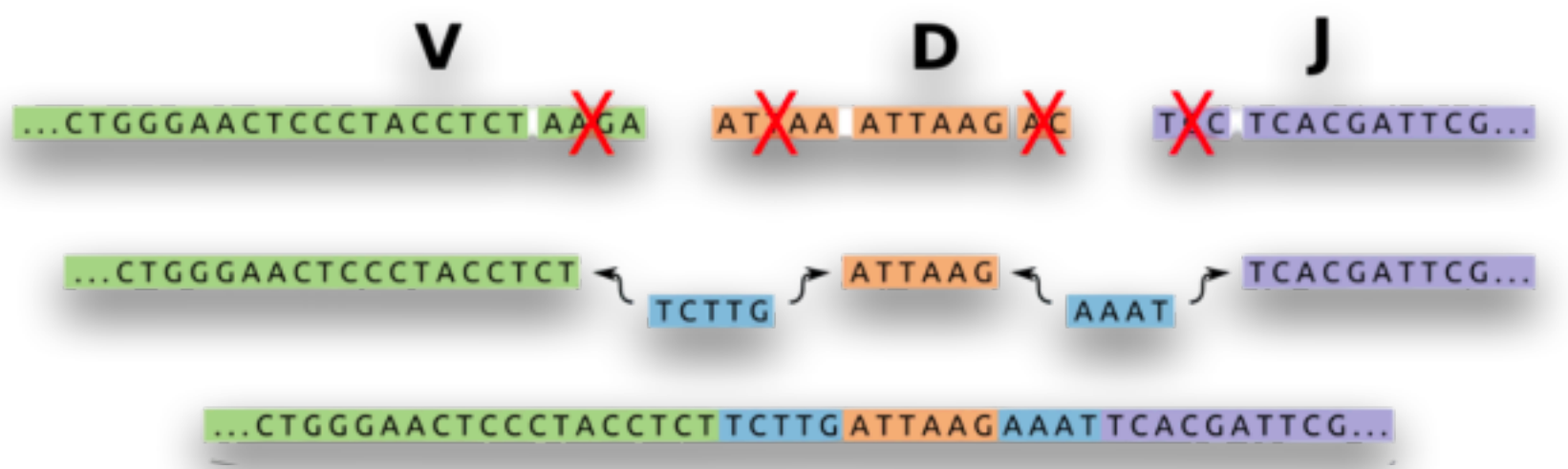$10^{61}$ different sequences
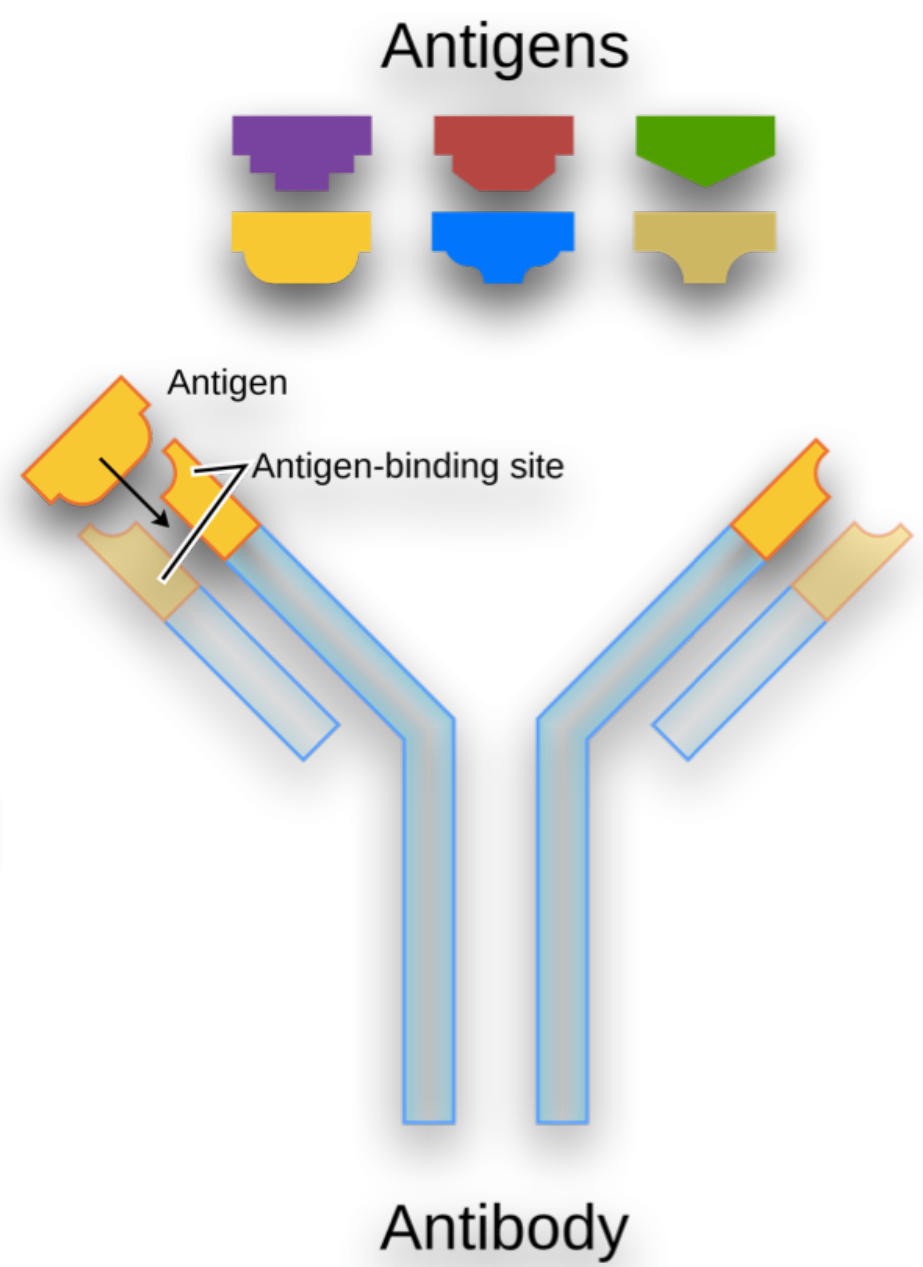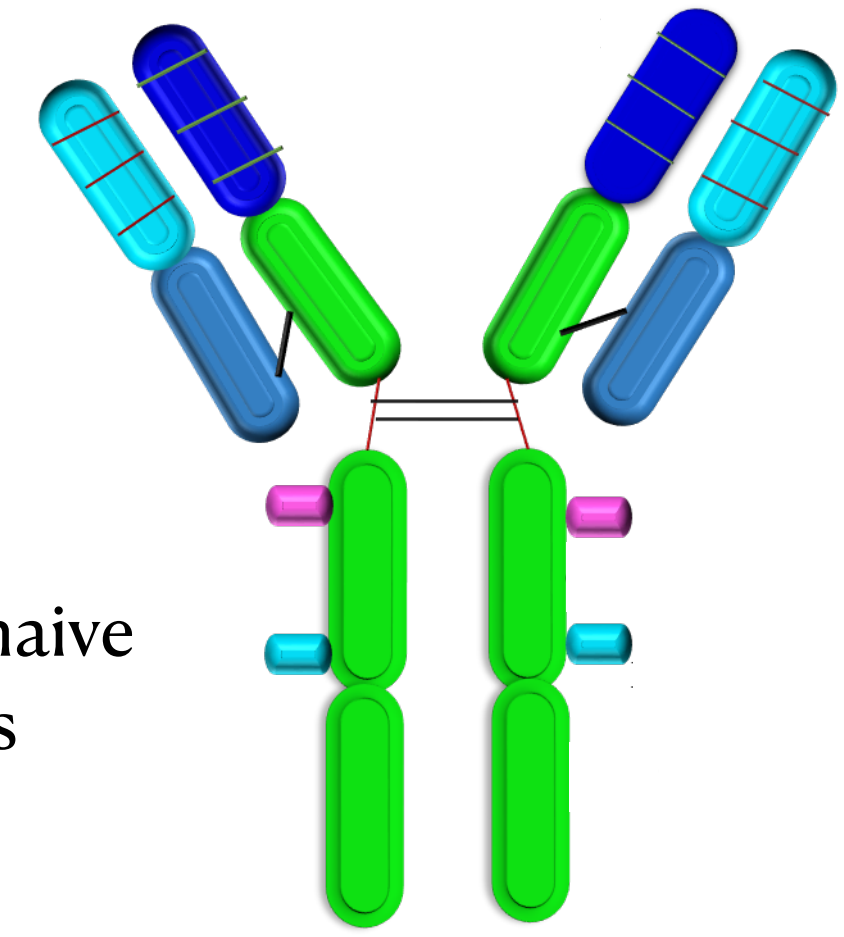
**IgM**

$10^{12}$ unique naive sequences

Antigens

Antigen

Antigen-binding site

Antibody

**IgG**

# MOTIVATION

❖ B cell receptor (BCR) repertoires are highly diverse thanks to V(D)J recombination process + somatic hypermutations.
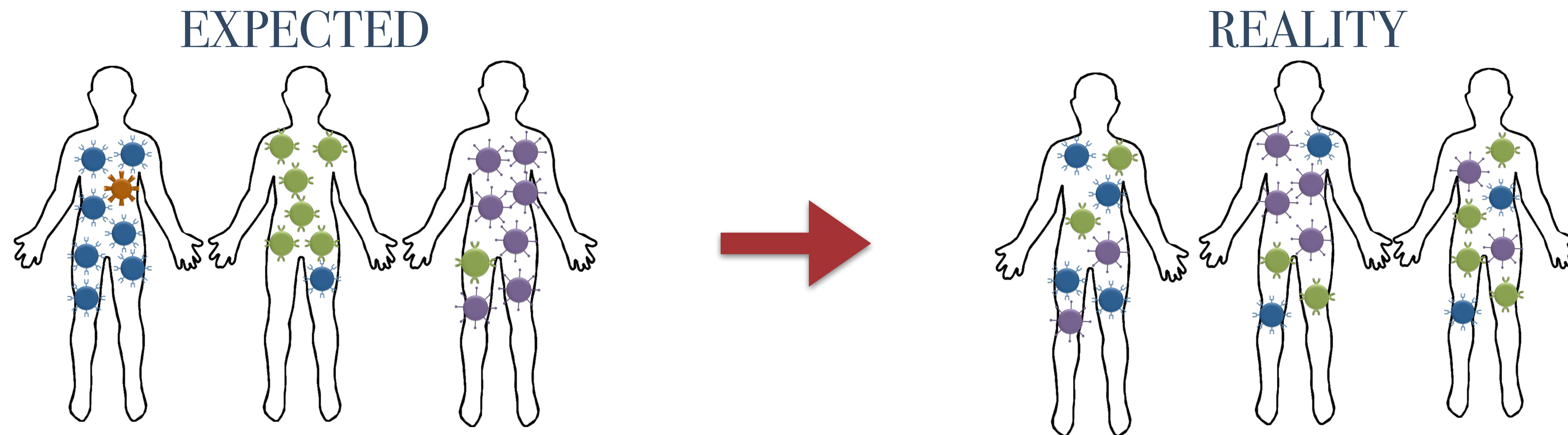
# MOTIVATION

❖ B cell receptor (BCR) repertoires are highly diverse thanks to V(D)J recombination process + somatic hypermutations.

❖ A higher than expected by chance overlap of receptors is observed when repertoires from different individuals are compared.



EXPECTED

REALITY

❖ **OBJECTIVE**: design a statistical model that is able to predict the number of sequences that will be shared among individuals.

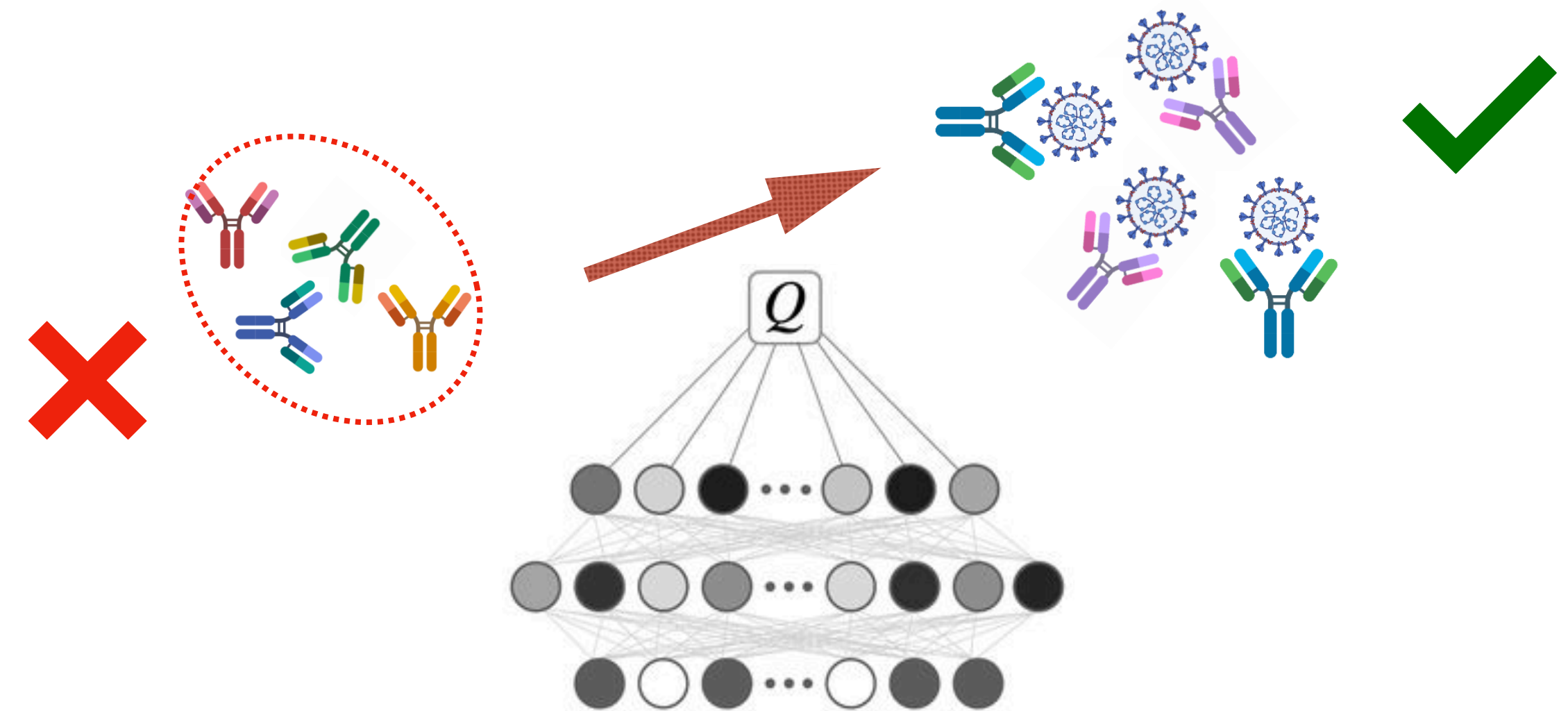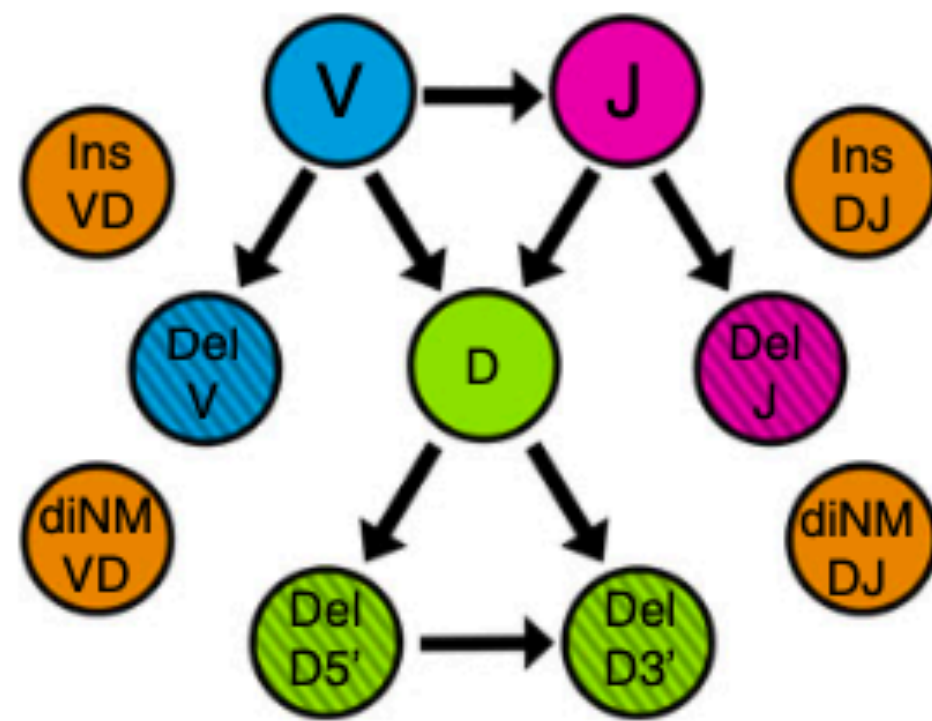# MODELLING THE SHARING MECHANISMS

❖ **Recombination biases** → Not all scenarios all equally likely.

The statistics of V(D)J recombination process are well captured by $P_{gen}$ model.

❖ **Selection processes** → Like central tolerance or clonal selection.

$P_{post}$ is a selection model that identifies sequence features characteristic of functional lymphocyte repertoires



$$P(\text{scenario}) = P(V)P(J|V)P(D|V,J)P(\text{del}V|V)$$
$$\times P(\text{del}J|J)\ P(\text{del}D5'|D)\ P(\text{del}D3'|\text{de}$$
$$\times P(\text{insVD}) \prod_i^{\text{InsVD}} P(n_i|n_{i-1})$$
$$\times P(\text{insDJ}) \prod_i^{\text{InsDJ}} P(n_i|n_{i-1})$$

ENS | PSL★

Human IgM BCRs

The sharing number indicates the number of different individuals in which a sequence is found

| Individual 1 | Individual 2 | Individual 3 |
|---|---|---|
| CASSENIQYF | CASSLTEAGEYF | CASSEDNNEQFF |
| CASSEDNNEQFF | CAWTWGGTGGEKLFF | CASNVQGSTEAFF |
| CASSLVLNTEAFF | CASSPPAGGVREQFF | CASLLTDTQYF |
| CASSELDTQYF | CSASVAVSGNQPQHF | CASAAEGLNTEAFF |
| CASSPPGELFF | CARCFTGFSLREQYF | CSAKGFGTEAFF |
| CASSLGTGARQPQHF | CASLLTDTQYF | CASSQGDRHQPQHF |
| CASSLGQGGSPLHF | CASSEDNNEQFF | CASSPPGELFF |
| CASTVGVDGYYEQYF | CASSELDTQYF | |
| CASSLTEAGEYF | CASSLTGNNSPLHF | |
| | CASSLAAREGSSQYF | |

SHARING NUMBER = 3

## Human IgM BCRs



Theoretical prediction using the generation function:

$$G(x, \{N_i\}) = \sum_{m=0}^{n} M_m(N_i)x^m = \sum_{s \in S} \prod_{i=1}^{n} \left[ e^{-p(s)N_i} + (1 - e^{-p(s)N_i})x \right]$$

Using $p(s) = P_{gen}$ underestimates how many sequences are shared among individuals. But the model of convergent recombination + selection, $p(s) = P_{post}$, accurately predicts this quantity.
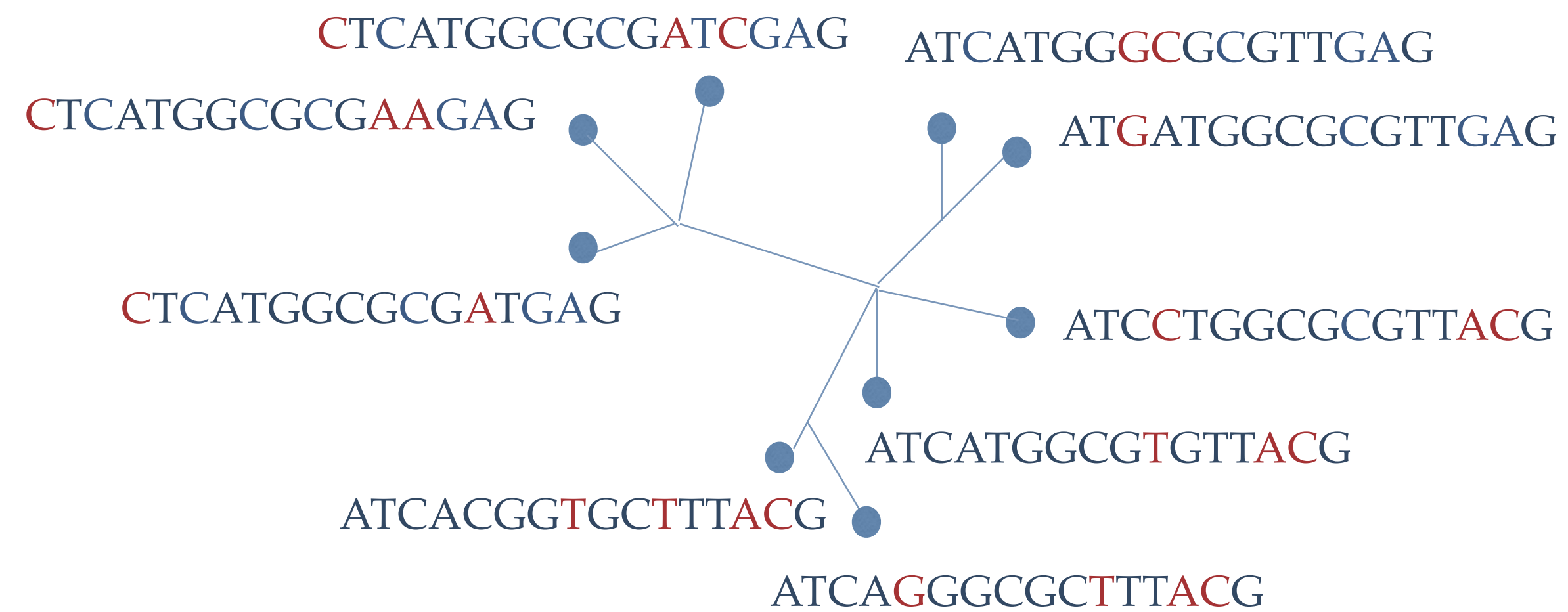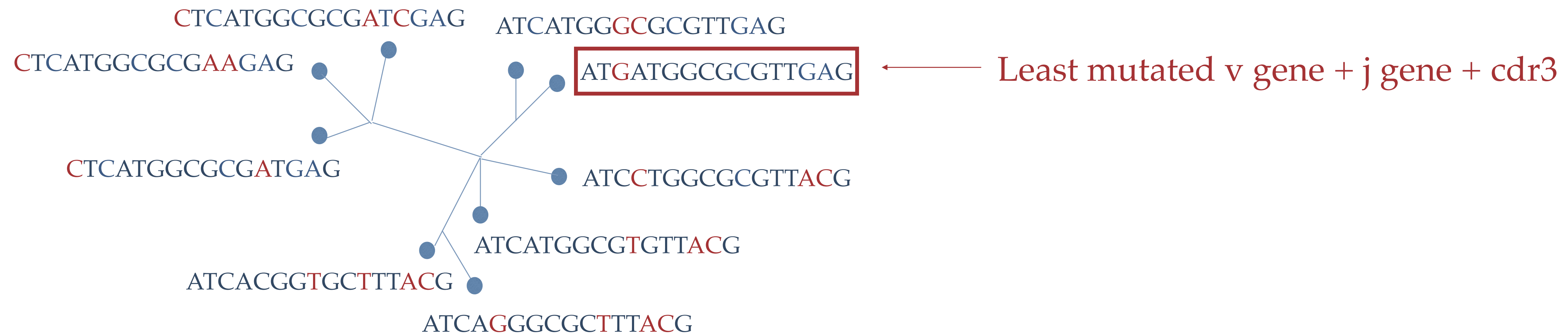
# PREDICTING SHARING IN IGG REPERTOIRES

❖ The same pipeline can be applied to non mutated IgG repertoires. We need to recover the ancestors of clonal families:

# PREDICTING SHARING IN IGG REPERTOIRES

❖ The same pipeline can be applied to non mutated IgG repertoires. We need to recover the ancestors of clonal families:

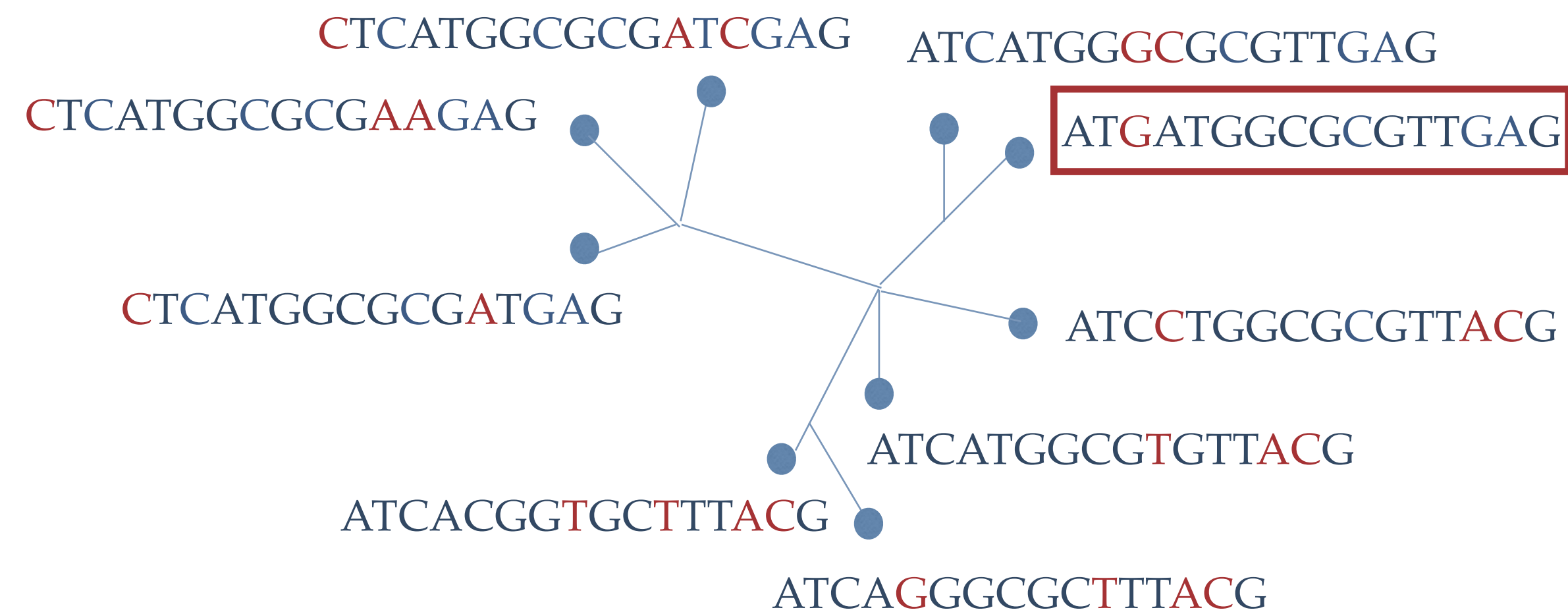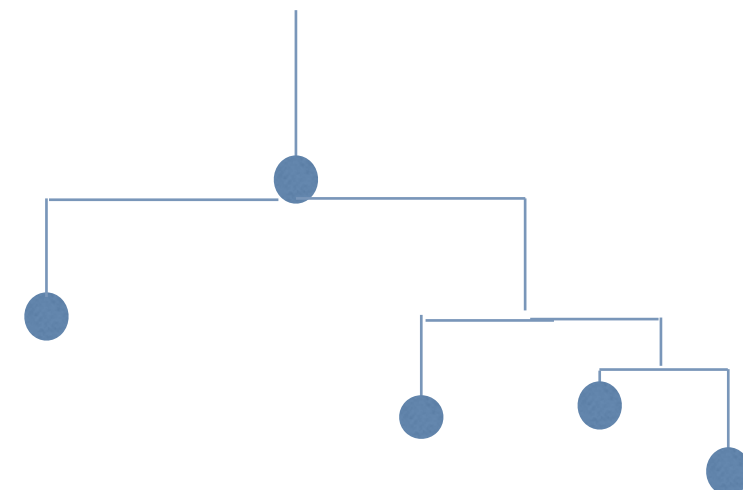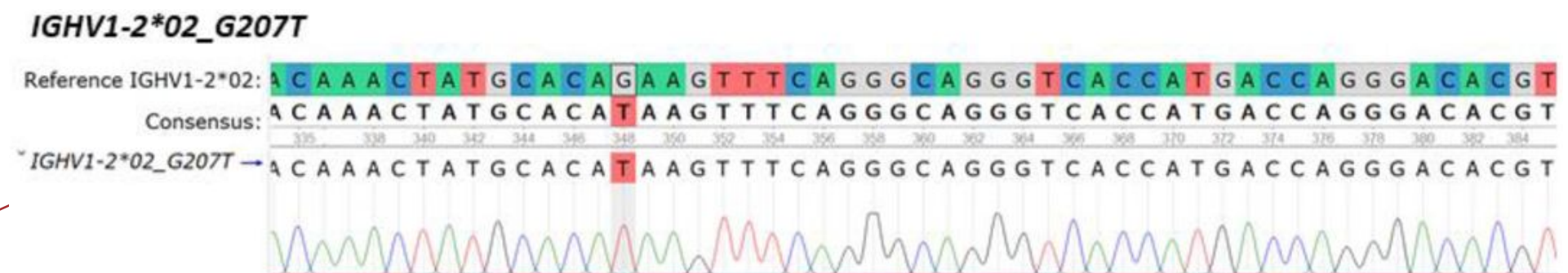❖ The same pipeline can be applied to non mutated IgG repertoires. We need to recover the ancestors of clonal families:
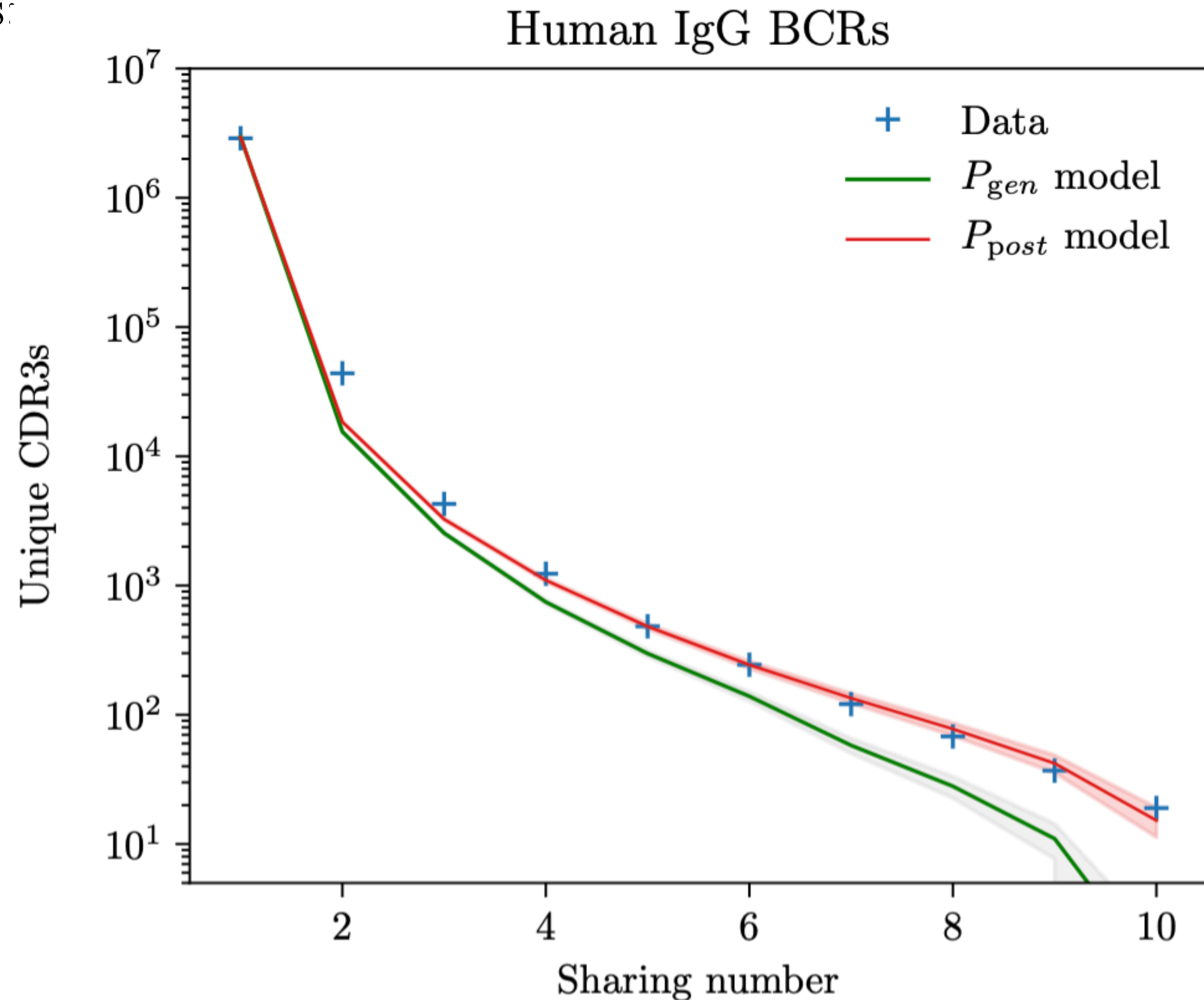


Human IgG BCRs

❖ Comparison of sharing number distribution in a cohort of 43 individuals **COVID-19 positive** and predictions from $P_{post}$ :

❖ Comparison of sharing number distribution in a cohort of 43 individuals **COVID-19 positive** and predictions from $P_{post}$ :



**Are these overshared sequences indicative of a specific antibody response to an antigen?**

Expected frequency from
the sharing pattern

$$P^*_{data} = \text{argmax}_{P_{data}} \, \mathbb{P}(x_1, \ldots, x_n \mid P_{data})$$

Level of certitude on $P_{data}$ given the observations $\{x_1, x_2, \ldots, x_n\}$:

$$\mathbb{P}(P_{post} > P_{data}) = \int_0^{P_{post}} \frac{\mathbb{P}(x_1, \ldots, x_n \mid P_{data}) \rho_{prior}(P_{data})}{\int_0^1 \mathbb{P}(x_1, \ldots, x_n \mid P_{data}) \rho_{prior}(P_{data}) dP_{data}} dP_{data}$$

Expected frequency from the sharing pattern

$$P^*_{data} = \text{argmax}_{P_{data}} \mathbb{P}(x_1, \ldots, x_n \mid P_{data})$$

Expected frequency from
the sharing pattern

$$P^*_{data} = \text{argmax}_{P_{data}} \, \mathbb{P}(x_1, \ldots, x_n \mid P_{data})$$

Level of certitude on $P_{data}$ given the observations $\{x_1, x_2, \ldots, x_n\}$:

$$\mathbb{P}(P_{post} > P_{data}) = \int_0^{P_{post}} \frac{\mathbb{P}(x_1, \ldots, x_n \mid P_{data})\rho_{prior}(P_{data})}{\int_0^1 \mathbb{P}(x_1, \ldots, x_n \mid P_{data})\rho_{prior}(P_{data})dP_{data}} \, dP_{data}$$
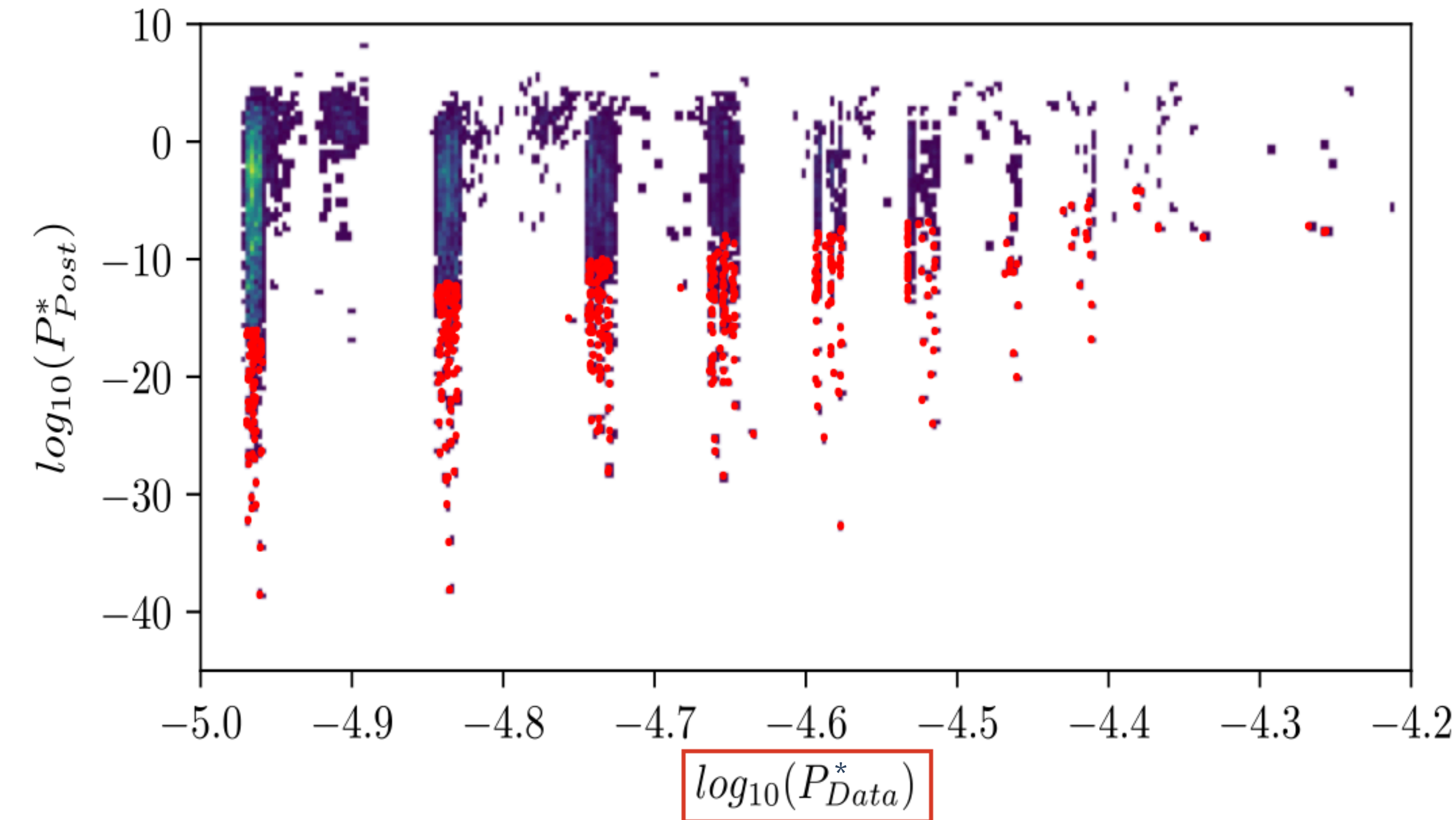
Bayesian analogous to p-value
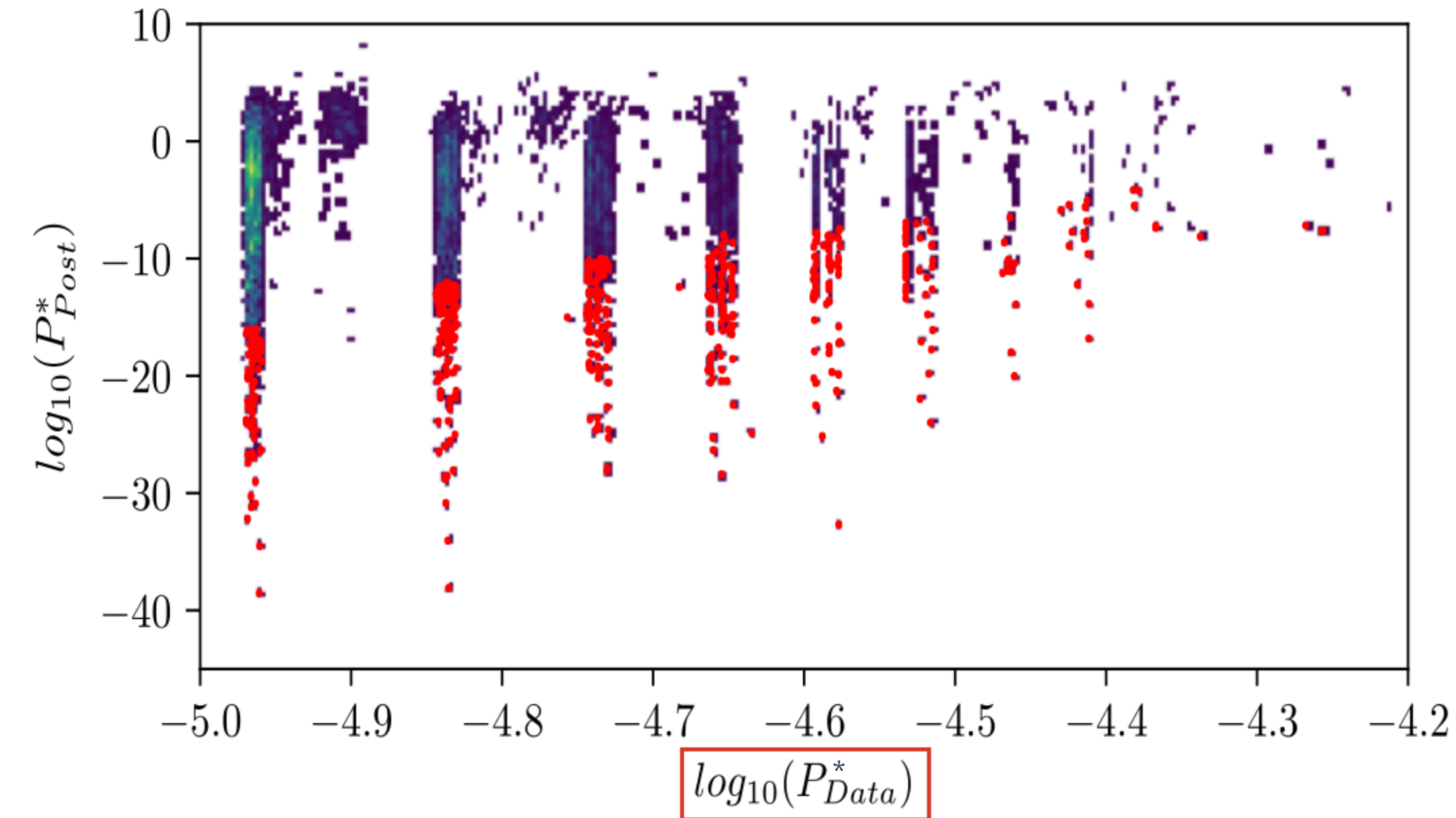
Expected frequency from
the sharing pattern

$$P^*_{data} = \mathrm{argmax}_{P_{data}} \, \mathbb{P}(x_1, \ldots, x_n \mid P_{data})$$

Level of certitude on $P_{data}$ given the observations
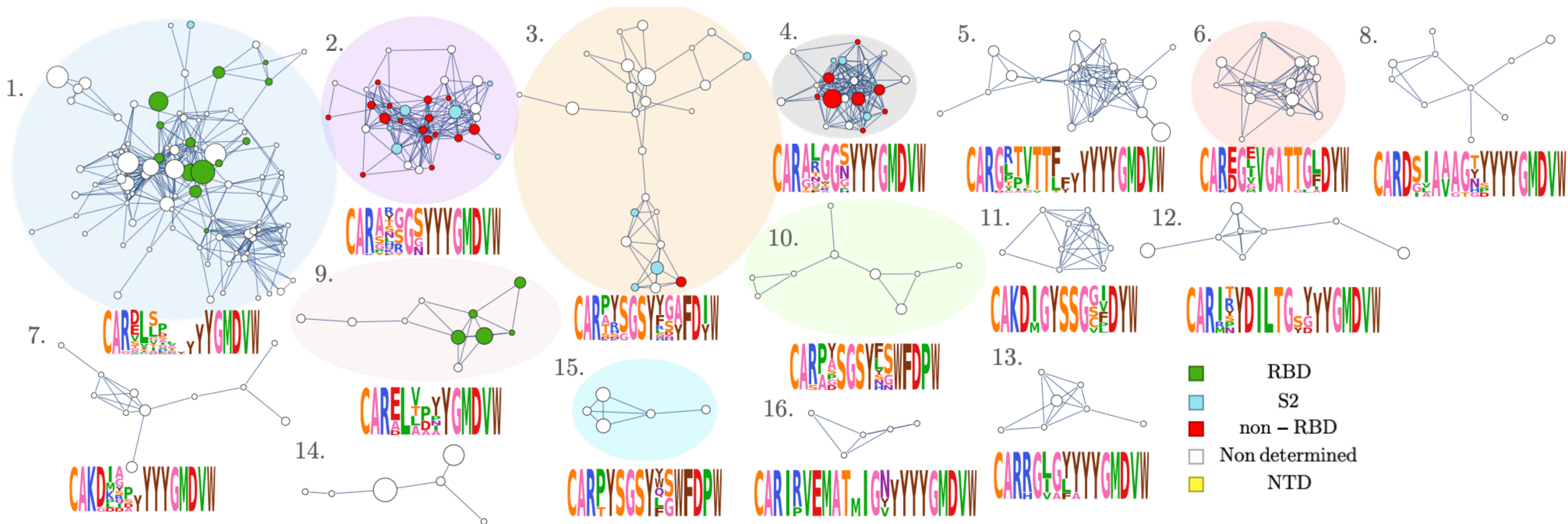$\{x_1, x_2, \ldots, x_n\}$:

$$\mathbb{P}(P_{post} > P_{data}) = \int_0^{P_{post}} \frac{\mathbb{P}(x_1, \ldots, x_n \mid P_{data})\rho_{prior}(P_{data})}{\int_0^1 \mathbb{P}(x_1, \ldots, x_n \mid P_{data})\rho_{prior}(P_{data})dP_{data}} dP_{data}$$

**Their frequency can't be explained by a high recombination likelihood.**

**→ Potential candidates of SARS-CoV-2 antibodies.**
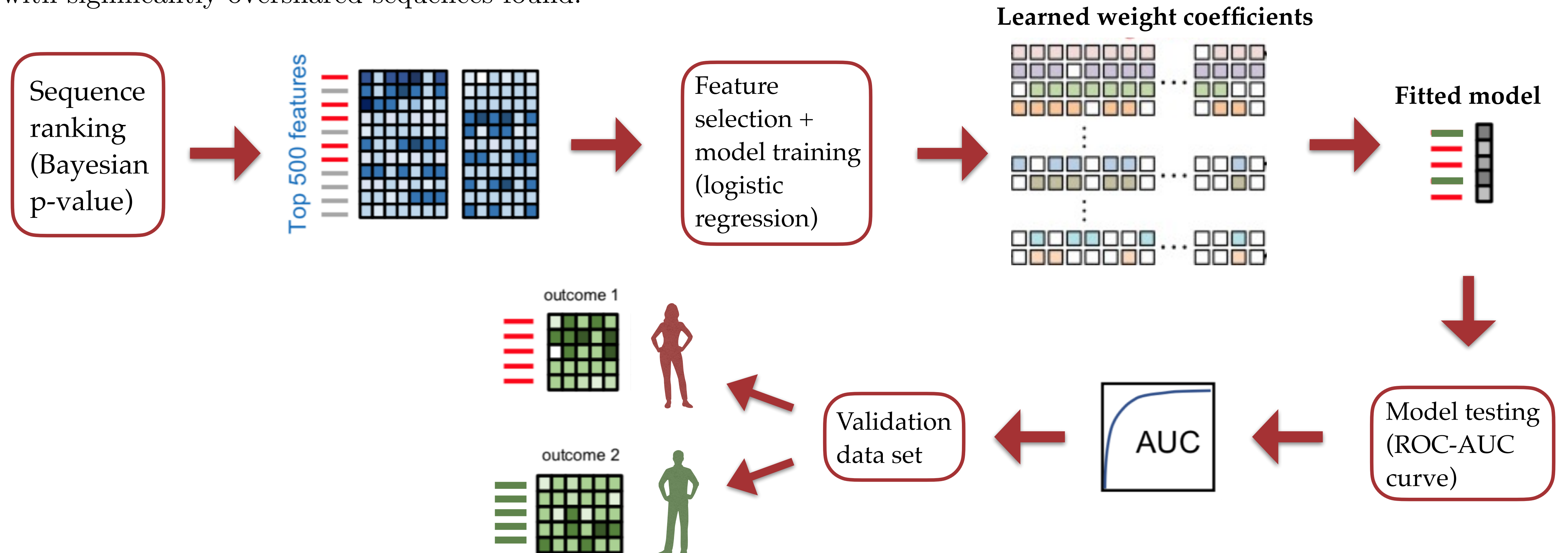
(i) **Clustering + annotation** → Overshared IGH sequences are clustered in networks with Levenshtein distance < 2 and are matched against reported SARS-CoV-2 antibodies.

(ii) **Learning of regression models for COVID-19 diagnosis** → Prediction of COVID-19 status based on the overlap with significantly overshared sequences found.

# CONCLUSIONS

❖ The statistical model here presented accurately estimates the probability of observing a productive B cell receptor in a repertoire.

❖ The model accurately predicts how many sequences will be shared among n healthy individuals but it fails at capturing the selection pressure existing after antigen encounter. The significance of this effect is measured by defining a Bayesian analogous to p-value.

❖ The sharing analysis here presented might be particularly useful to help designing a vaccine that ellicitates a more transverse immune response since the antibodies that have been isolated have been already produced by a large number of individuals.