

UC San Diego



UNIVERSITY OF
LOUISVILLE[®]

Computer Science meets Immunology: how computational analyses help to study diseases

Yana Safonova
postdoctoral researcher, PhD

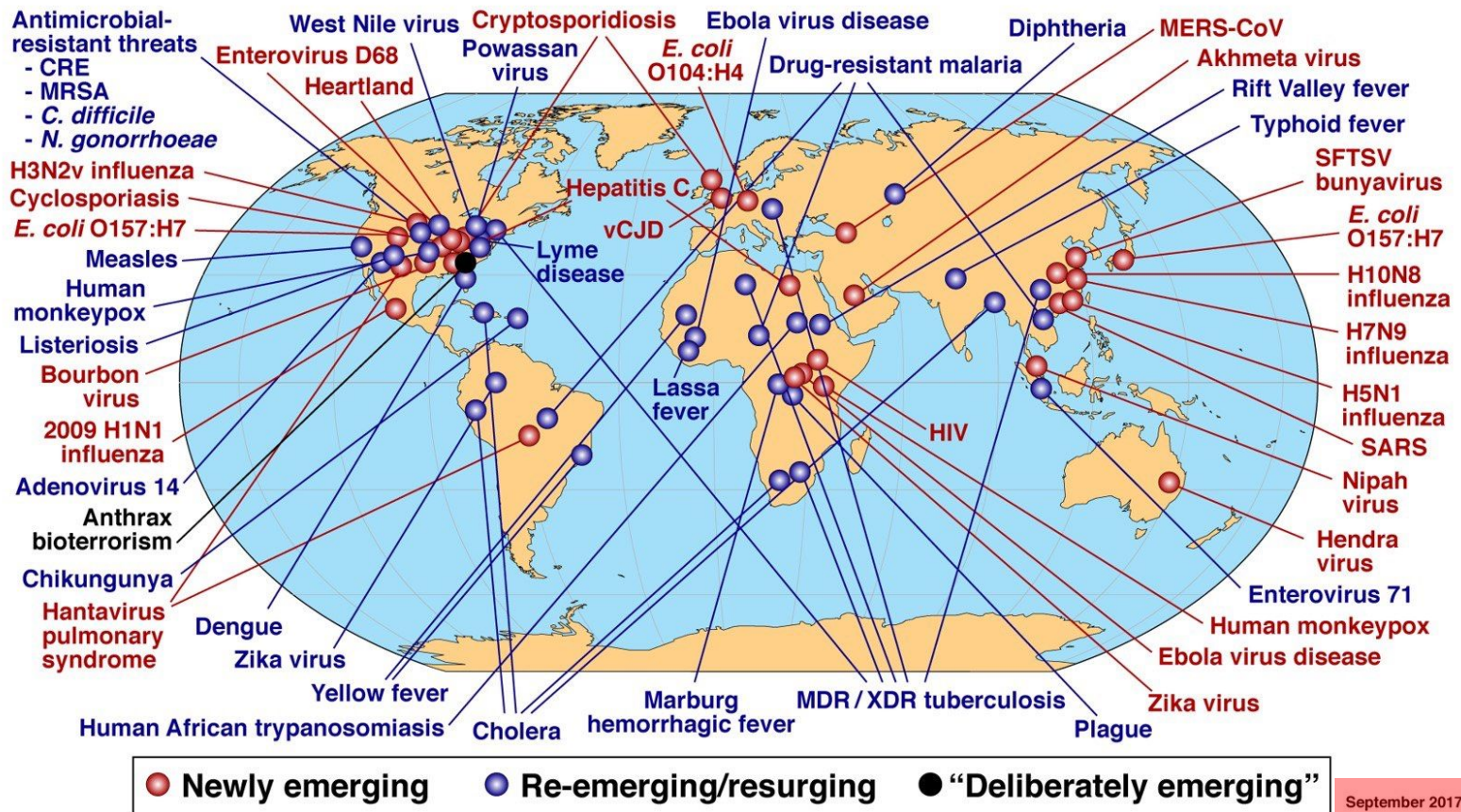
University of California San Diego
University of Louisville School of Medicine



@yana_safonova_

Newly emerging and re-emerging diseases

Global Examples of Emerging and Re-Emerging Infectious Diseases

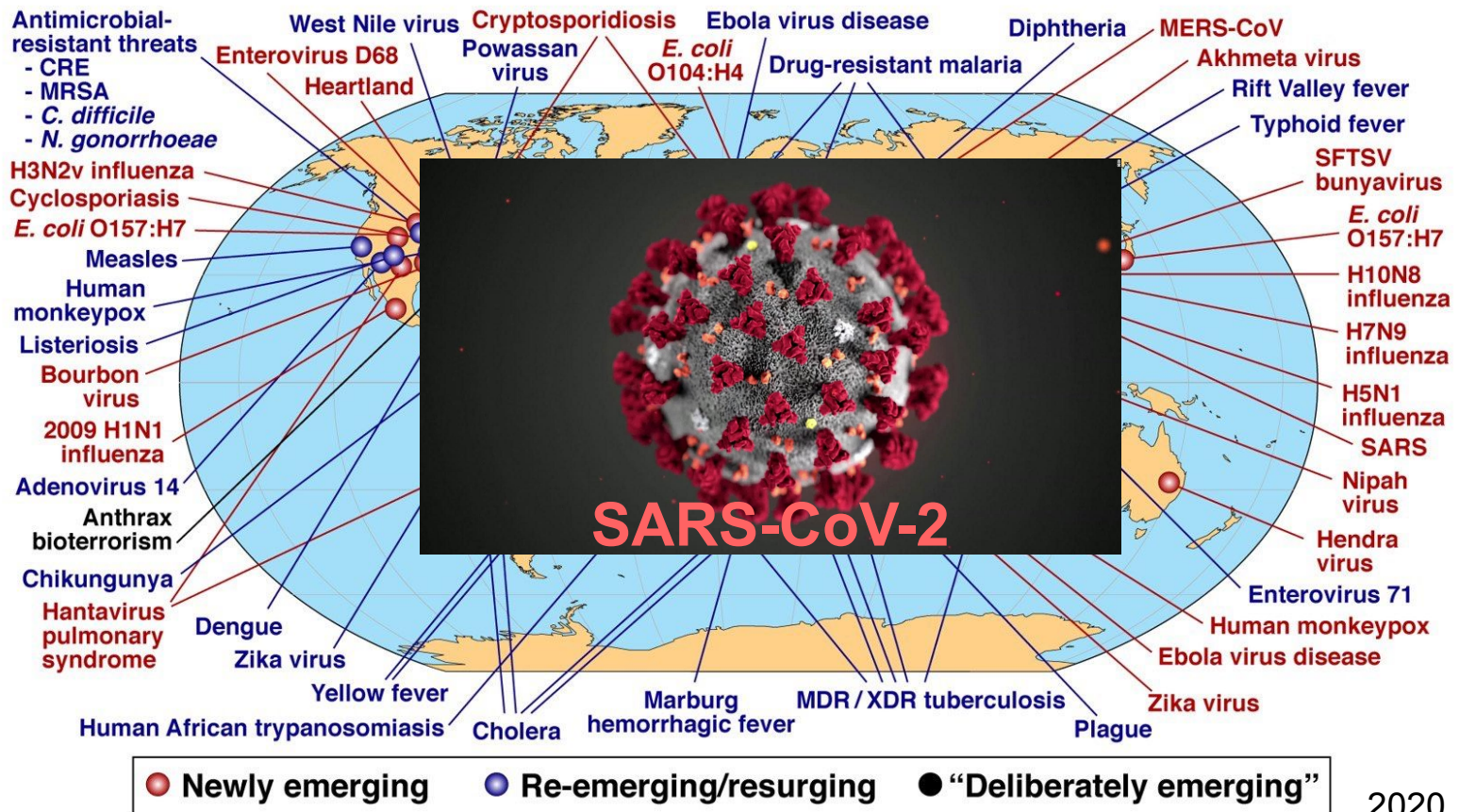


September 2017

https://en.wikipedia.org/wiki/Emerging_infectious_disease

Newly emerging and re-emerging diseases

Global Examples of Emerging and Re-Emerging Infectious Diseases



Response to (re)emerging diseases



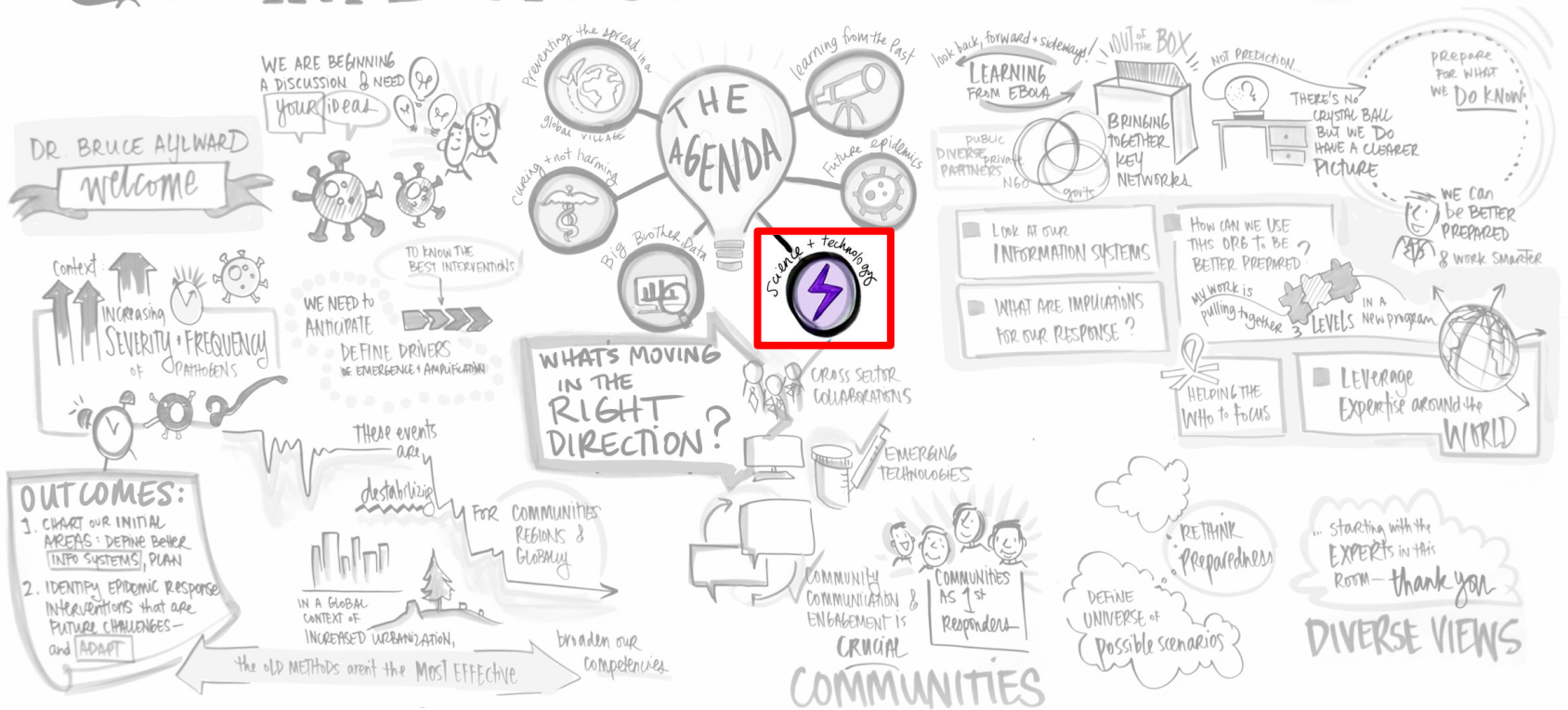
ANTICIPATING EMERGING INFECTIOUS DISEASE EPIDEMICS



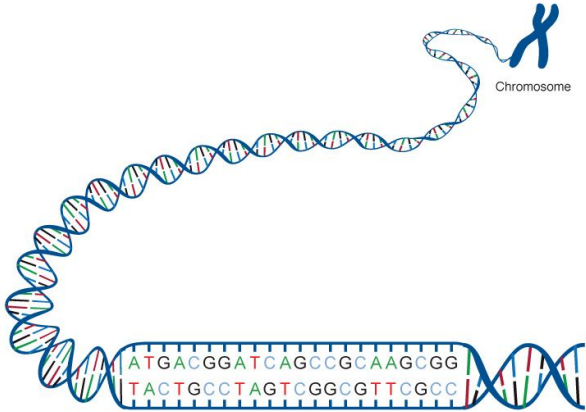
Response to (re)emerging diseases



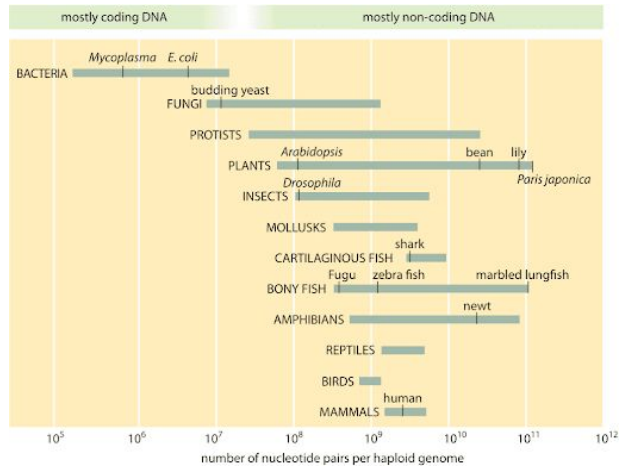
ANTICIPATING EMERGING INFECTIOUS DISEASE EPIDEMICS



Introduction to molecular biology

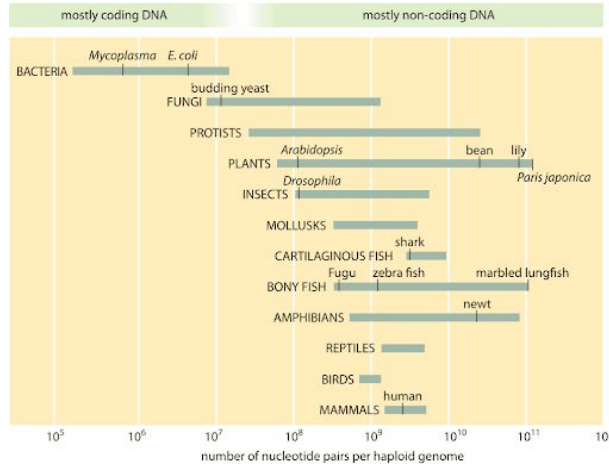
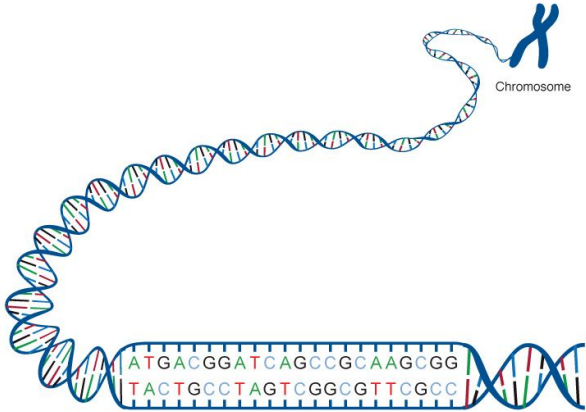


Introduction to molecular biology



Human genome ~ 3Gbp

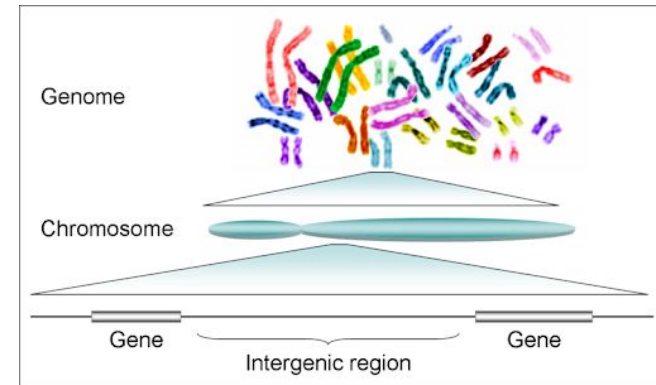
Introduction to molecular biology



Human genome ~ 3Gbp

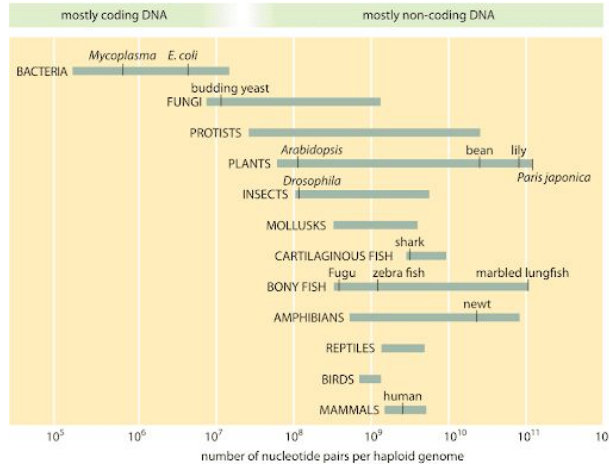


The Library of Babel (H. L. Borges) contains all possible combinations of symbols, mostly meaningless ones



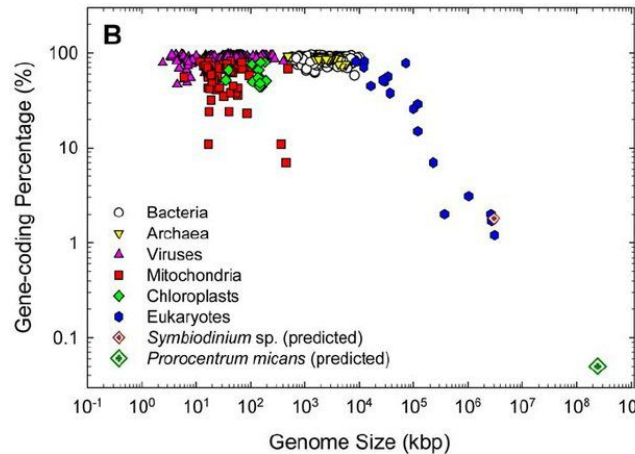
Genes in human genome take 1-2% of its length

Introduction to molecular biology

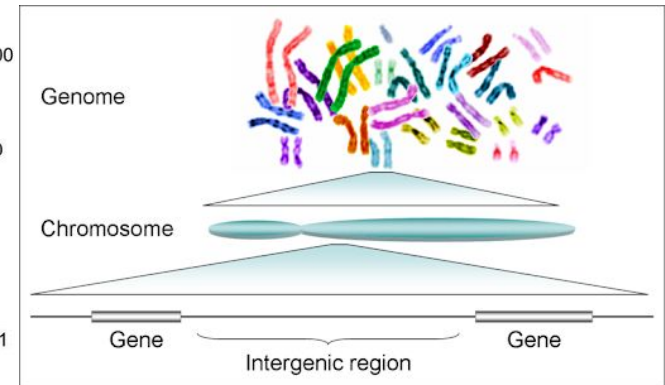


Human genome ~ 3Gbp

The Library of Babel (H. L. Borges) contains all possible combinations of symbols, mostly meaningless ones

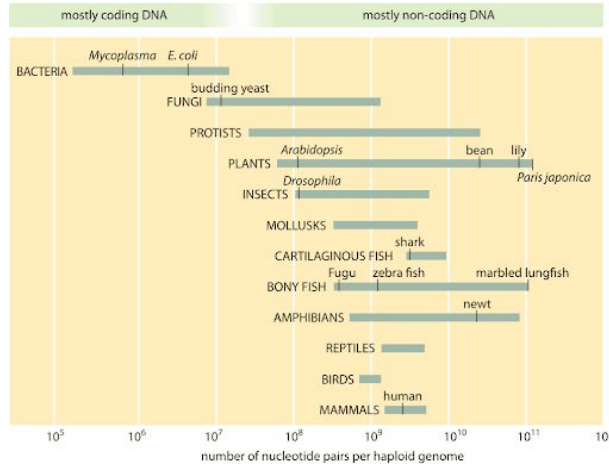
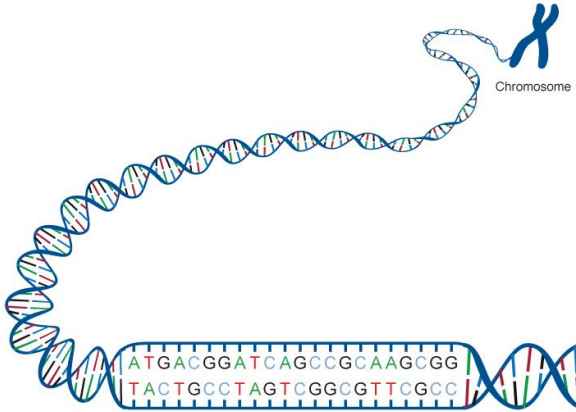


Hou and Lin, PLoS ONE, 2009



Genes in human genome take 1-2% of its length

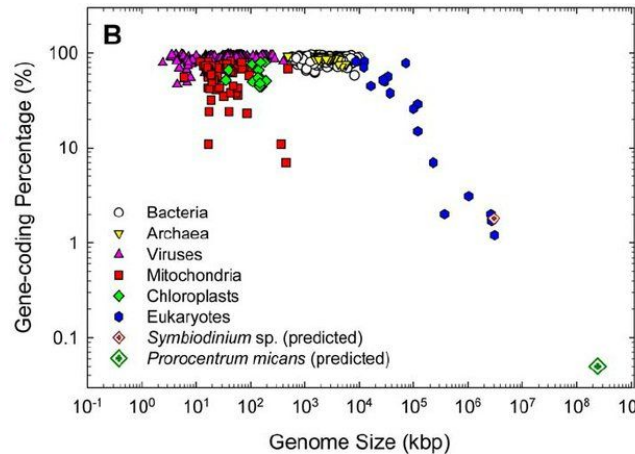
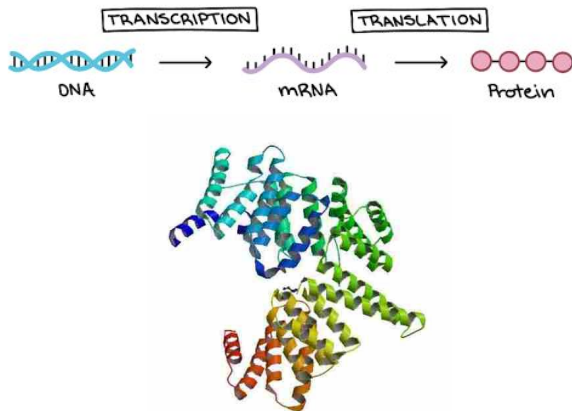
Introduction to molecular biology



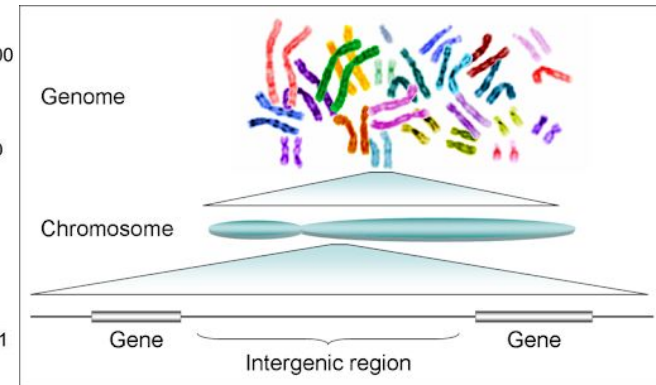
Human genome ~ 3Gbp

The Library of Babel (H. L. Borges) contains all possible combinations of symbols, mostly meaningless ones

The central dogma of biology

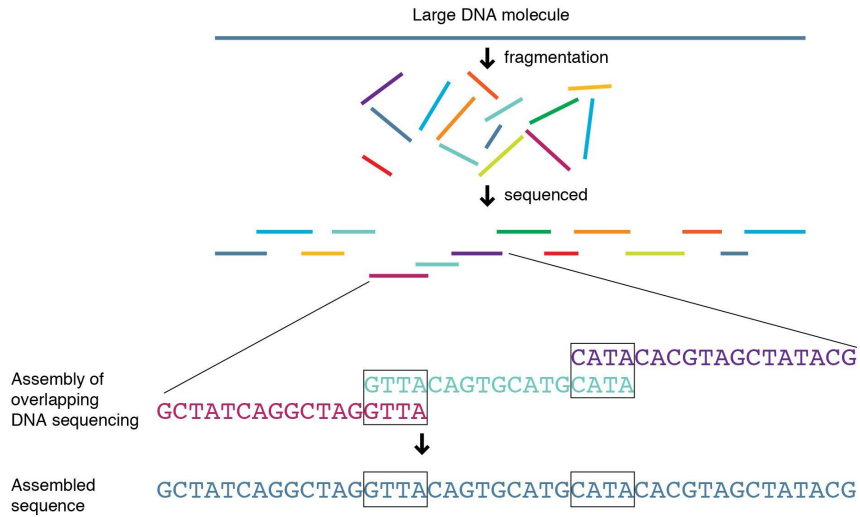


Hou and Lin, PLoS ONE, 2009

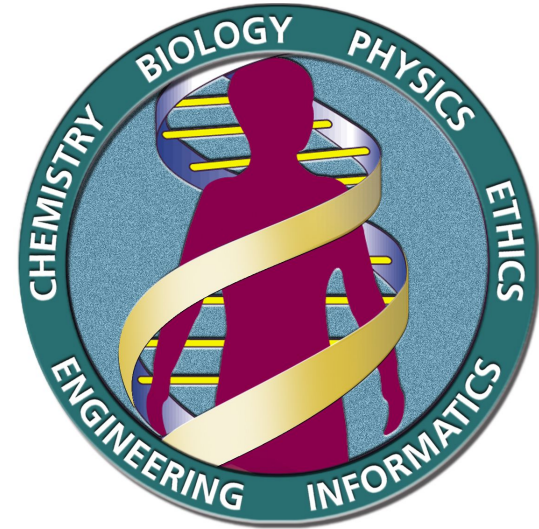


Genes in human genome take 1-2% of its length

The birth of bioinformatics

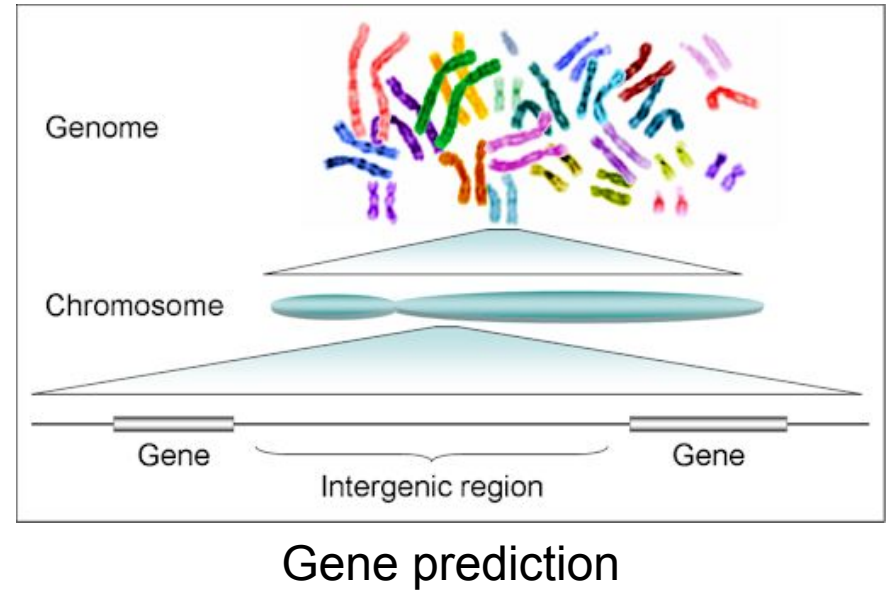
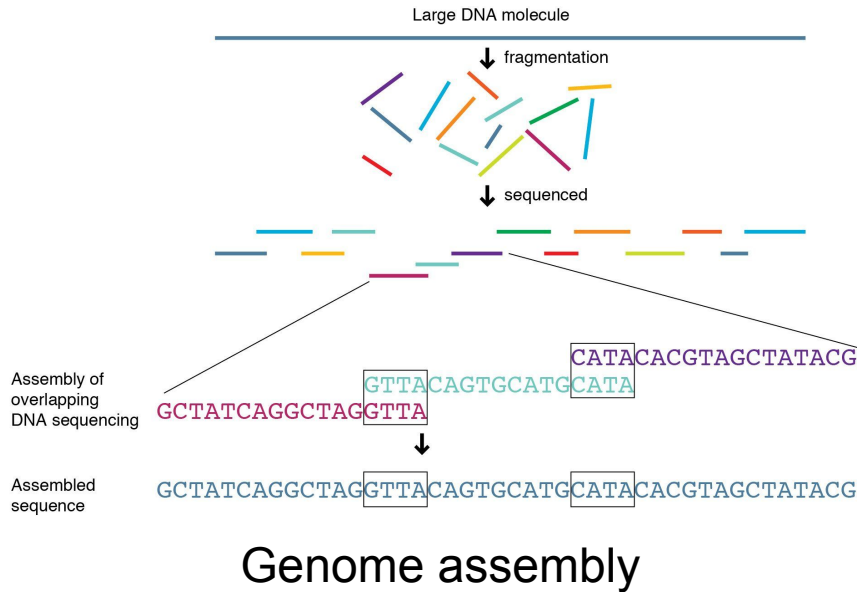


Genome assembly

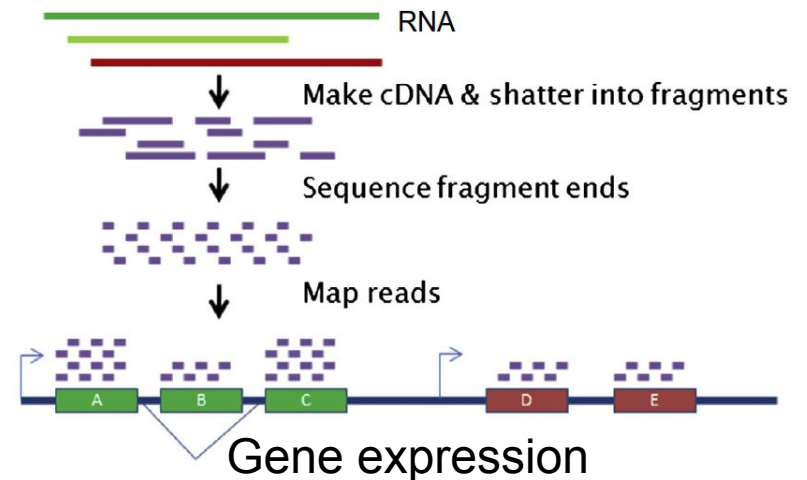
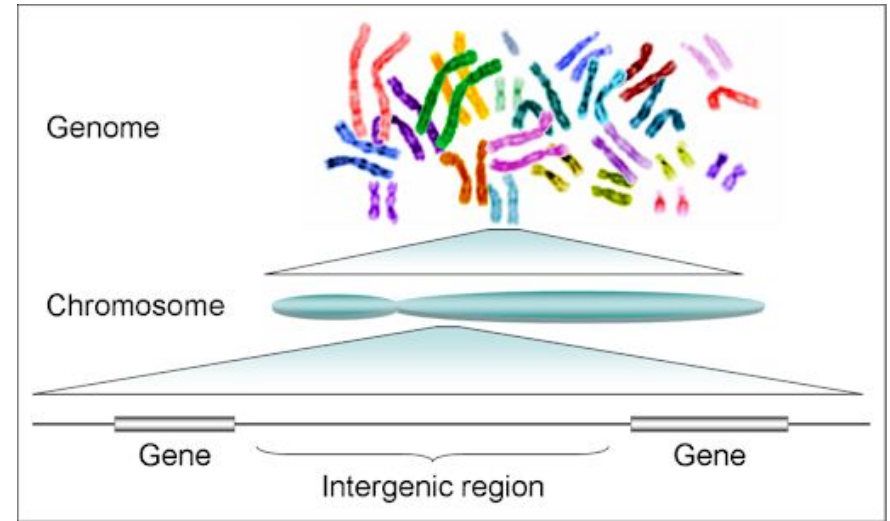
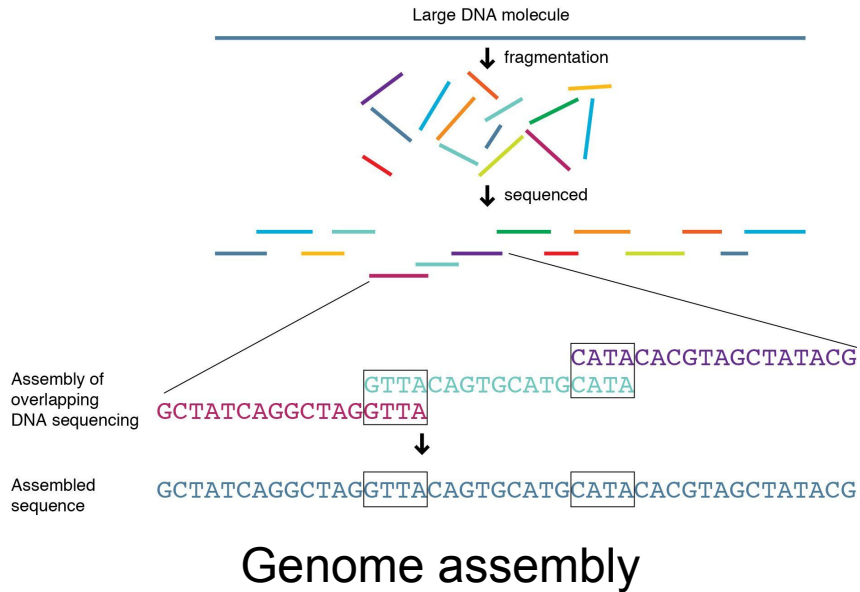


Human Genome Project
1990 – 2003

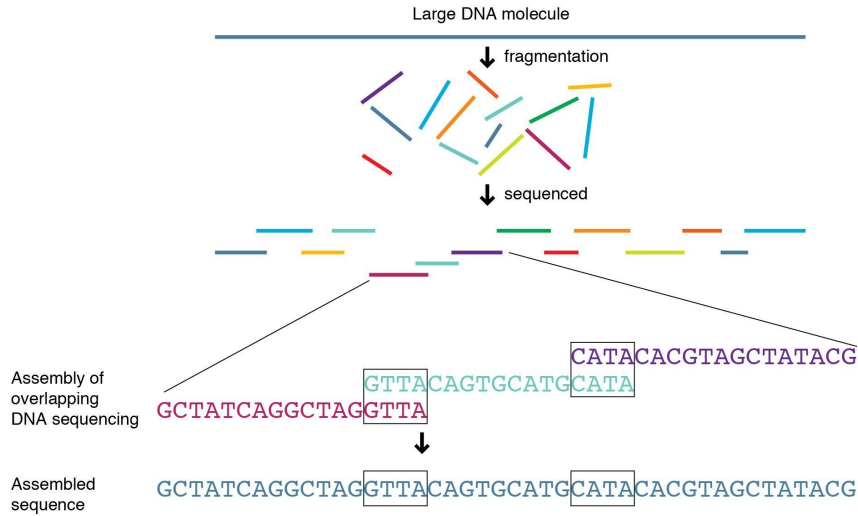
The birth of bioinformatics



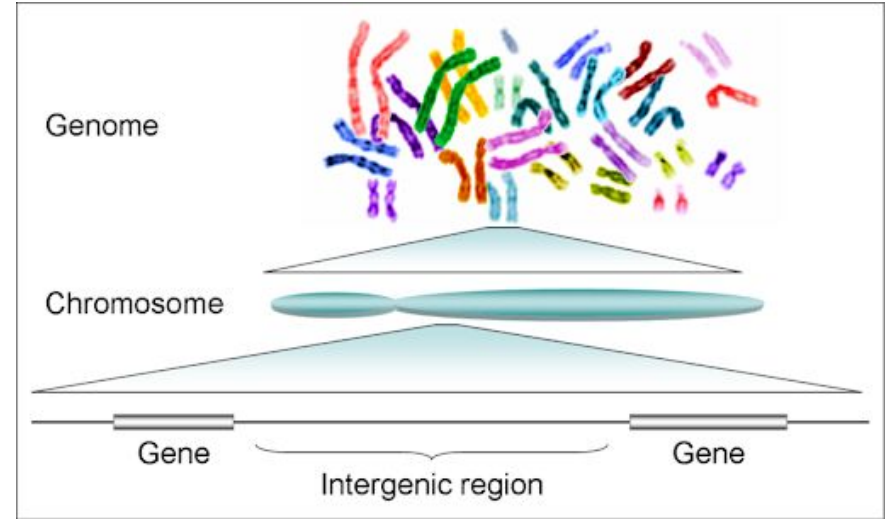
The birth of bioinformatics



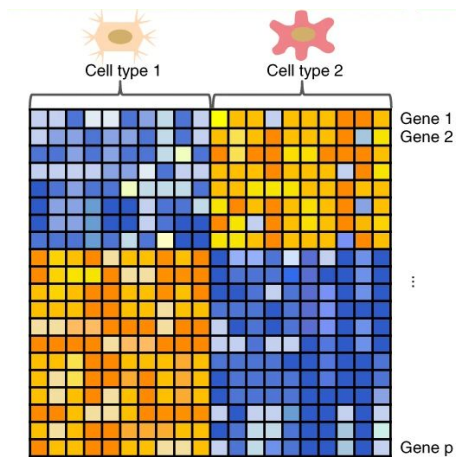
The birth of bioinformatics



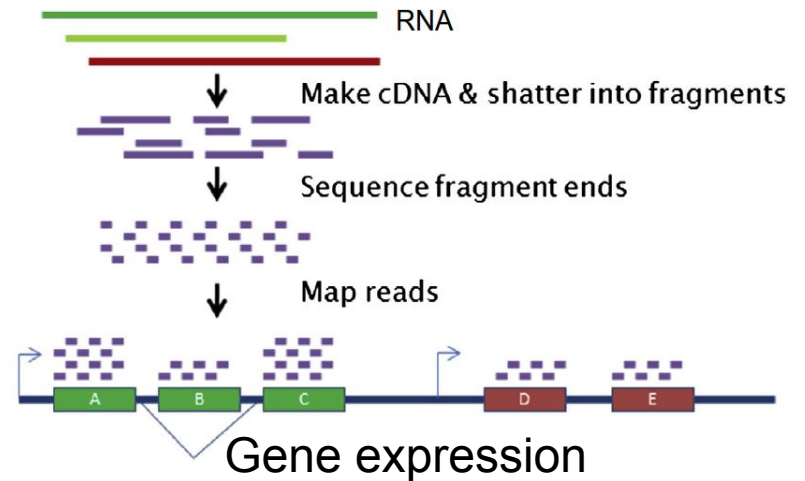
Genome assembly



Gene prediction



Differential gene expression

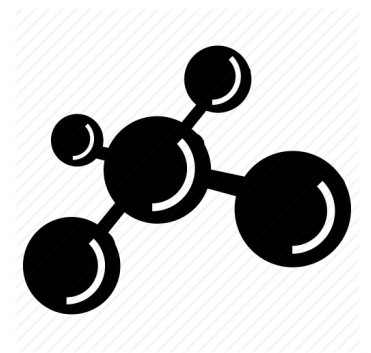
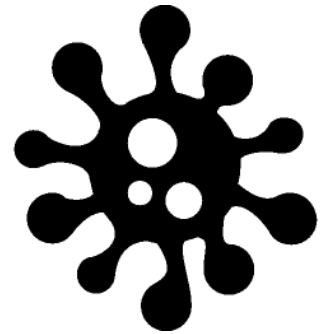
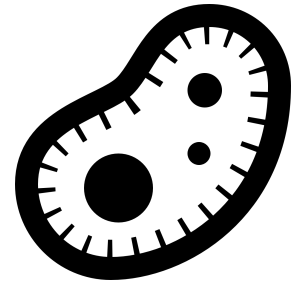


**Immune system = innate (or inherited) +
adaptive (or acquired) immune systems**

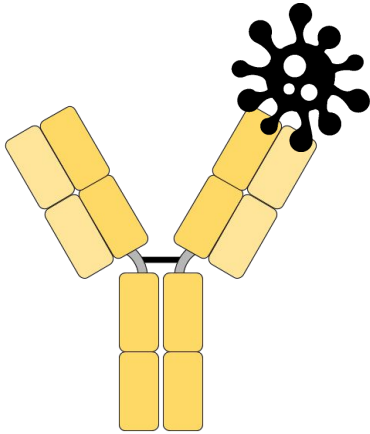


Adaptive immune system

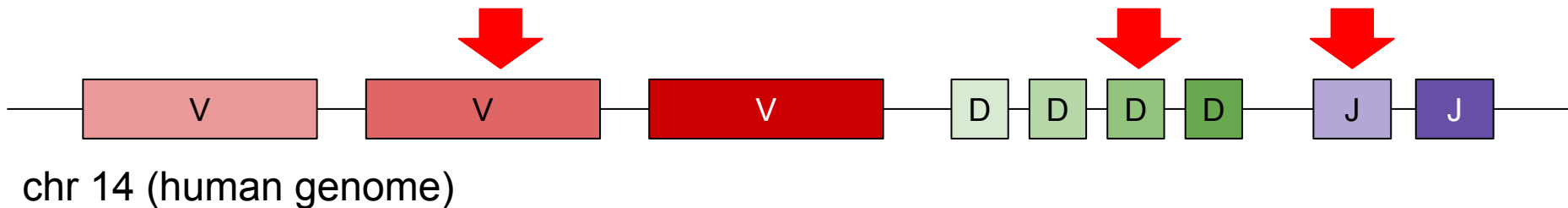
- Variety of threats to human body is huge and unpredictable
- Genome is too small to encode defences against all these threats
- Immune system has an ability to adapt to various threats using agents (e.g., antibodies) that **are not encoded** in the genome.



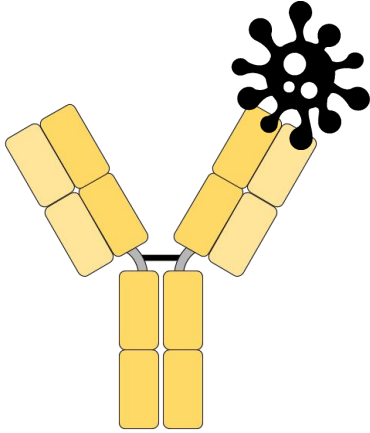
Antibodies are agents of the adaptive immune system



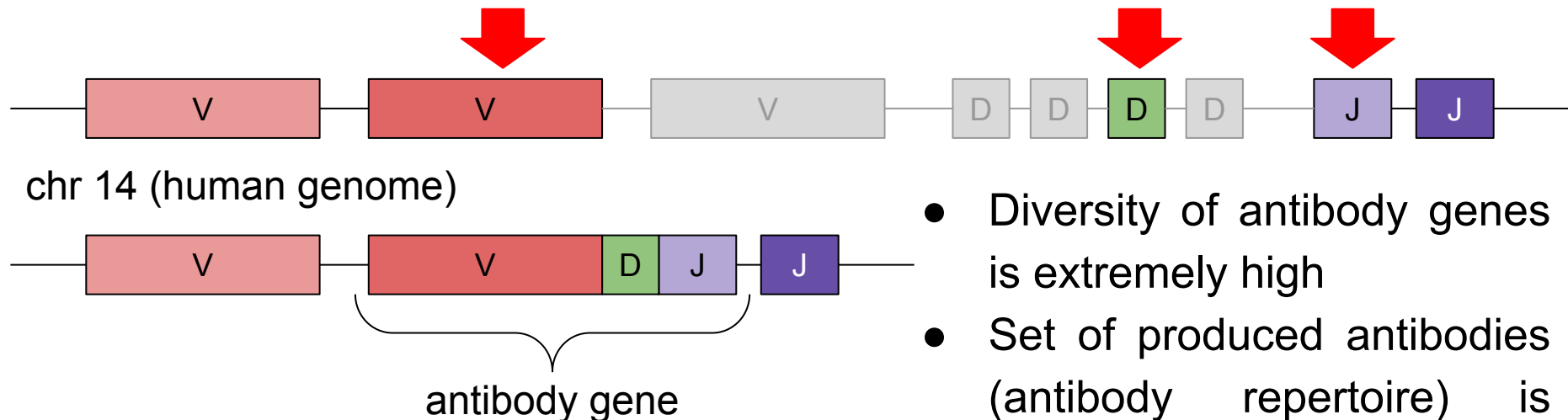
- Antibodies are proteins that bind to an antigen and cause its neutralization
- Antibodies are not encoded in the genome directly, but present a result of somatic genomic recombination



Antibodies are agents of the adaptive immune system

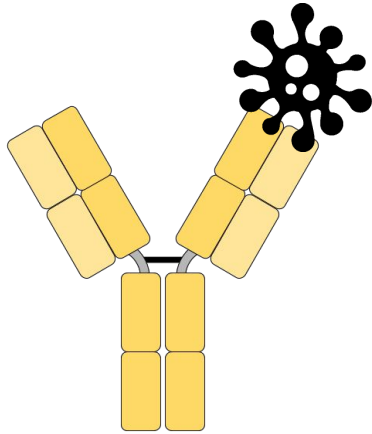


- Antibodies are proteins that bind to an antigen and cause its neutralization
- Antibodies are not encoded in the genome directly, but present a result of somatic genomic recombination

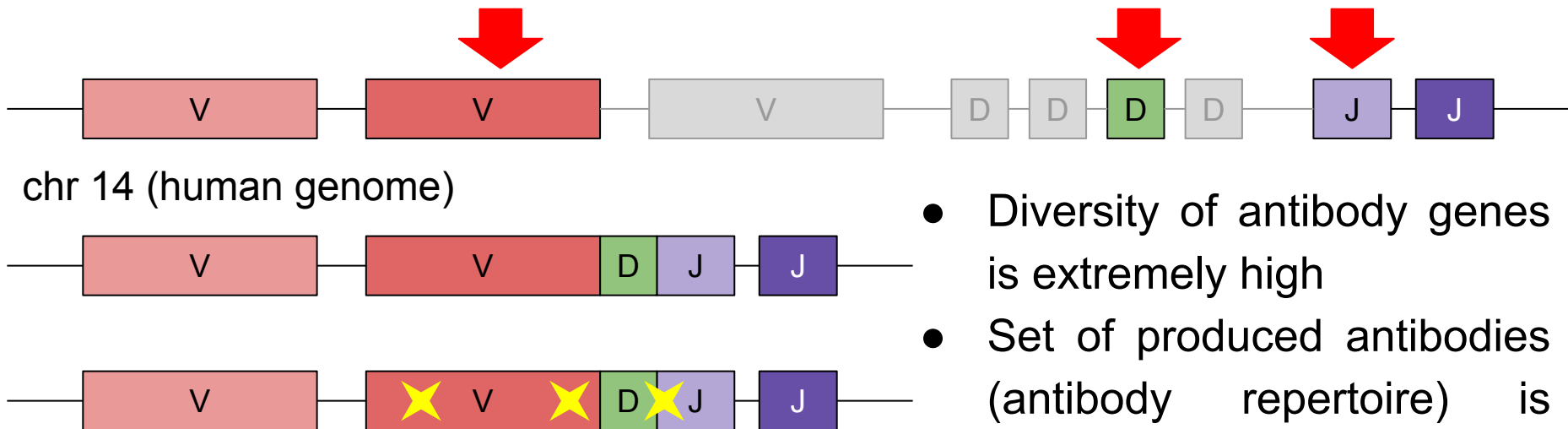


- Diversity of antibody genes is extremely high
- Set of produced antibodies (antibody repertoire) is unique for an individual

Antibodies are agents of the adaptive immune system

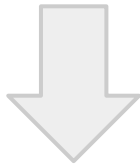
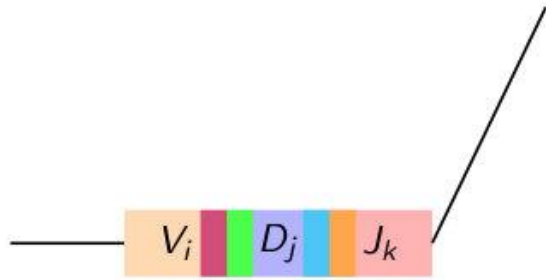


- Antibodies are proteins that bind to an antigen and cause its neutralization
- Antibodies are not encoded in the genome directly, but present a result of somatic genomic recombination

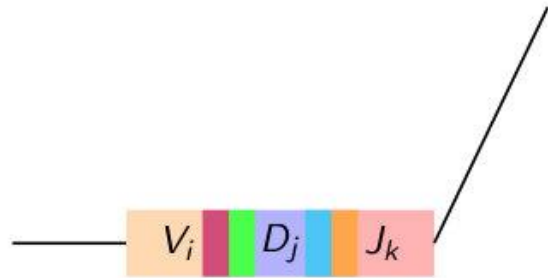


- Diversity of antibody genes is extremely high
- Set of produced antibodies (antibody repertoire) is unique for an individual

Why are antibodies so versatile if there are only $55 \times 23 \times 6$ VDJ recombinations?

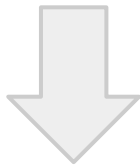


Why are antibodies so versatile if there are only $55 \times 23 \times 6$ VDJ recombinations?

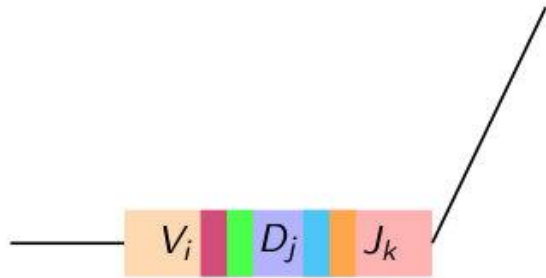


Recombination process is imperfect and includes many random processes:

- **Palindromic insertions**

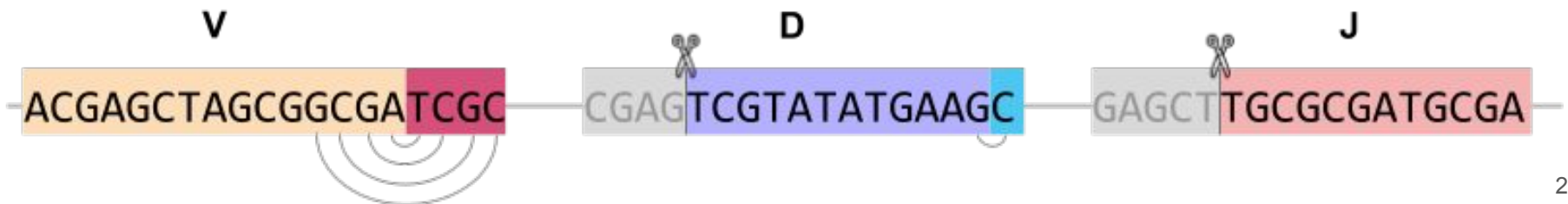


Why are antibodies so versatile if there are only $55 \times 23 \times 6$ VDJ recombinations?

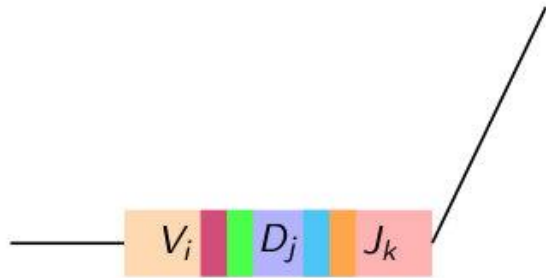


Recombination process is imperfect and includes many random processes:

- Palindromic insertions
- **Segment cleavage**

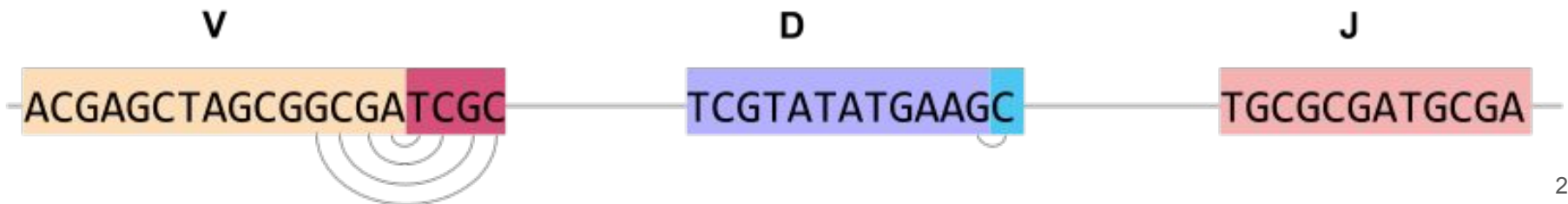
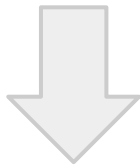


Why are antibodies so versatile if there are only $55 \times 23 \times 6$ VDJ recombinations?

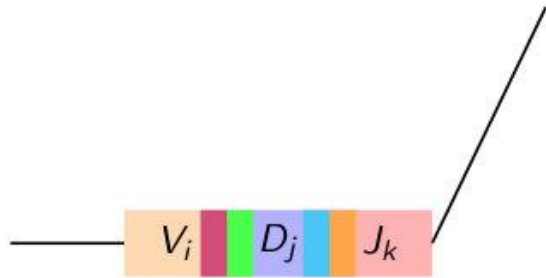


Recombination process is imperfect and includes many random processes:

- Palindromic insertions
- **Segment cleavage**

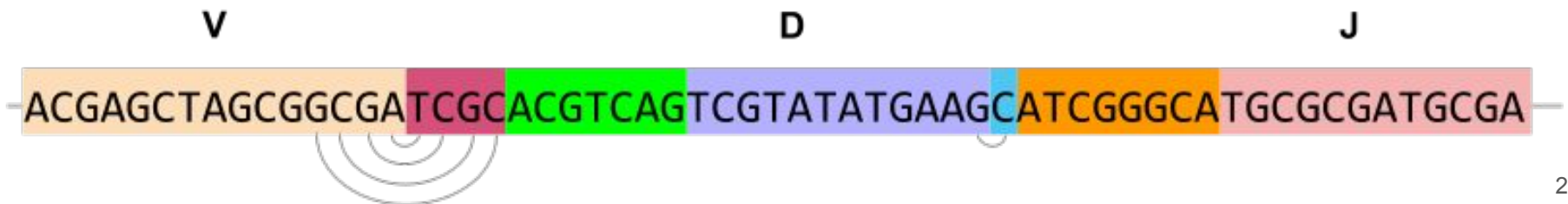
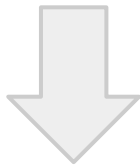


Why are antibodies so versatile if there are only $55 \times 23 \times 6$ VDJ recombinations?



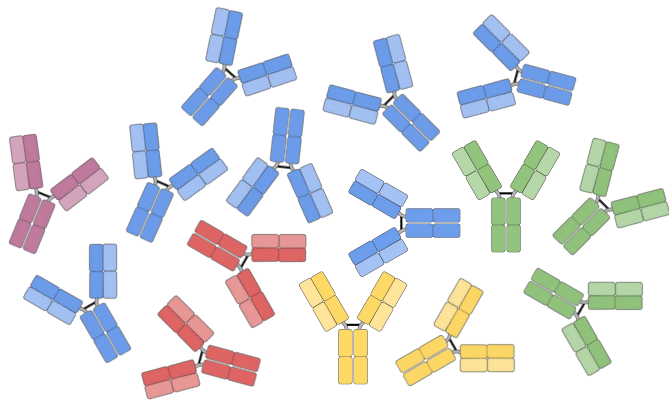
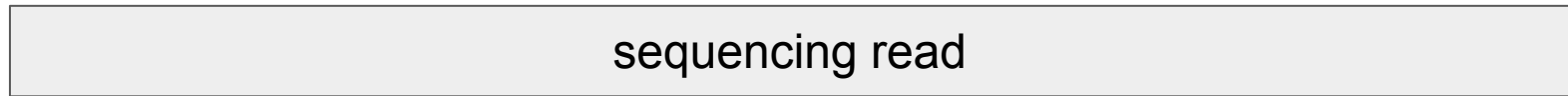
Recombination process is imperfect and includes many random processes:

- Palindromic insertions
- Segment cleavage
- **Non-genomic insertions**

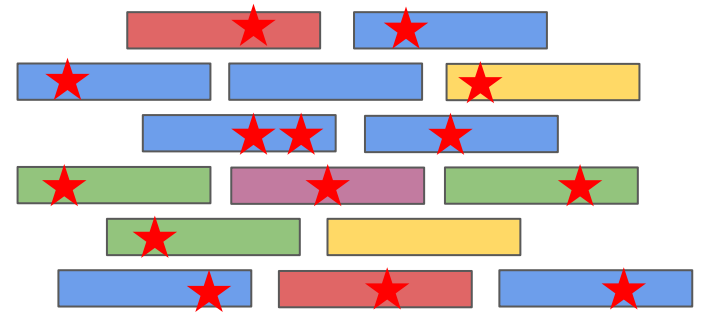


Antibody repertoire sequencing (Rep-seq)

Length: ~360 nt

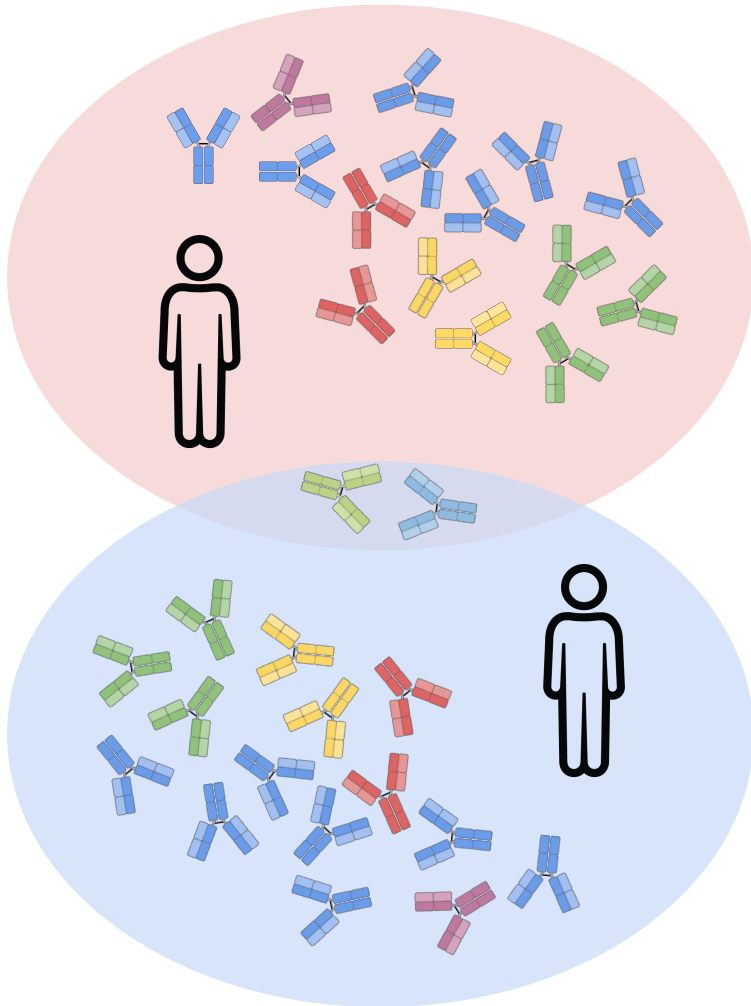


VDJ from DNA or RNA

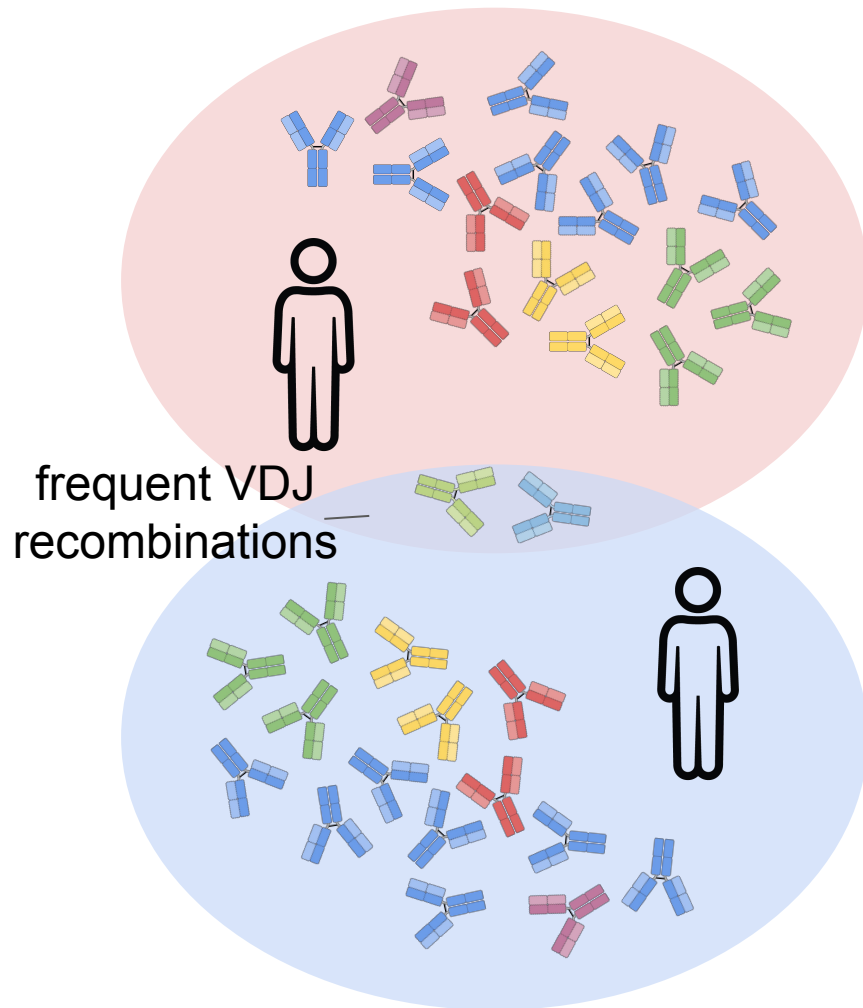


Error-prone immunosequencing reads

Antibody repertoire is unique for an individual



Antibody repertoire is unique for an individual

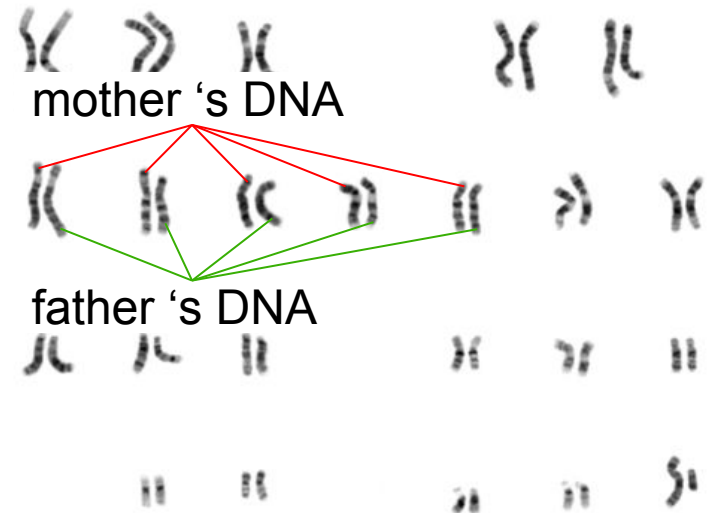
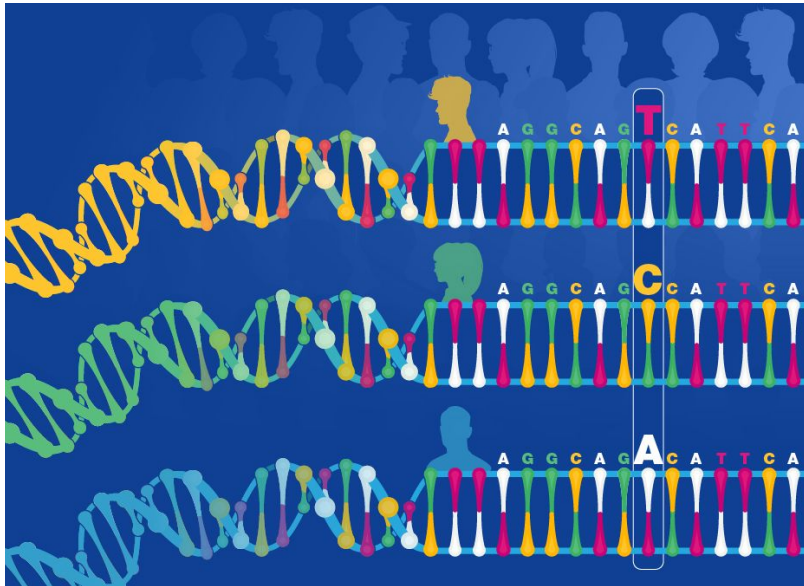


- VDJ sequences are extremely diverse
- If a VDJ sequence is shared between two individuals, it is a likely a frequent recombination rather a functionally important sequence
- **We cannot study antibody responses just comparing VDJ sequences**

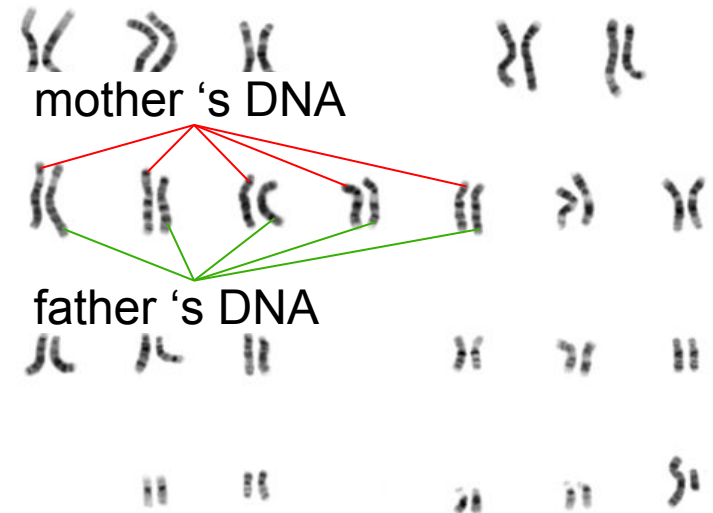
Genomic variations



Genomic variations

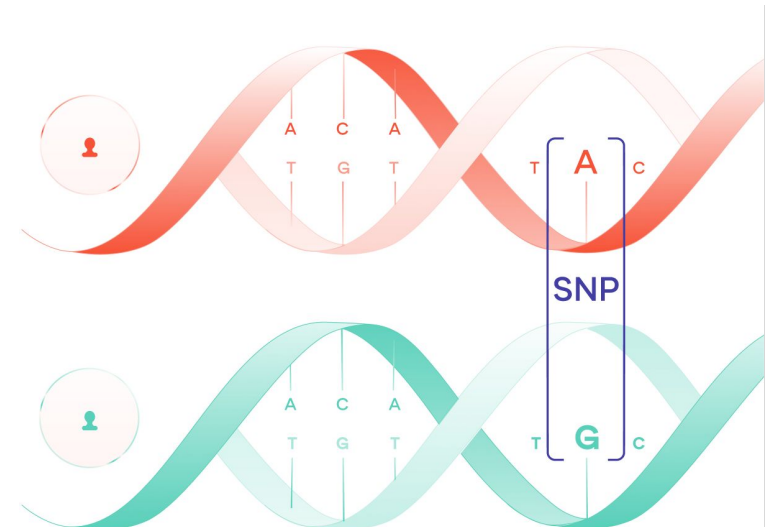


Genomic variations



Single Nucleotide Polymorphism (SNP) is associated with a position in the genome:

- A – both chromosomes have A
- A / G – one chromosome has A, another one has G
- G – both chromosomes have G

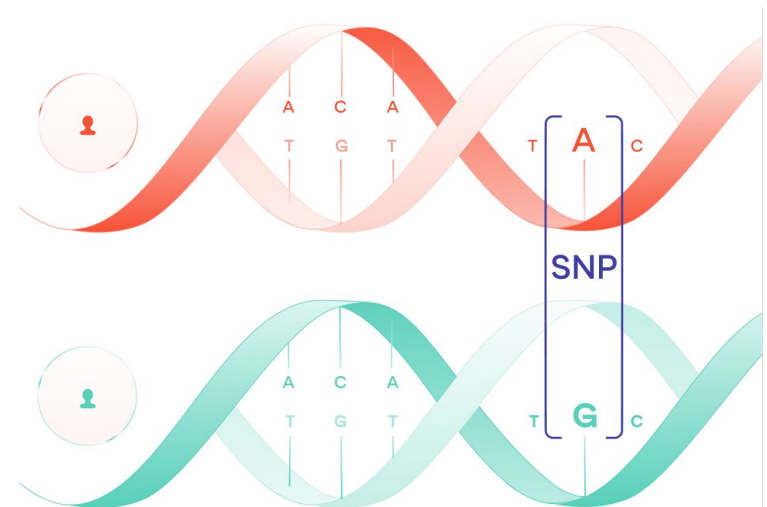


Altered genes produce altered proteins

		Second nucleotide				
		U	C	A	G	
U	U	UUU Phe	UCU	UAU Tyr	UGU Cys	U
	U	UUC	UCC Ser	UAC	UGC	C
	U	UUA Leu	UCA	UAA STOP	UGA STOP	A
	U	UUG	UCG	UAG STOP	UGG Trp	G
C	U	CUU Leu	CCU	CAU His	CGU	U
	C	CUC	CCC Pro	CAC	CGC Arg	C
	A	CUA	CCA	CAA Gln	CGA	A
	G	CUG	CCG	CAG	CGG	G
A	U	AUU Ile	ACU	AAU Asn	AGU Ser	U
	C	AUC	ACC Thr	AAC	AGC	C
	A	AUA	ACA	AAA Lys	AGA Arg	A
	G	AUG Met	ACG	AAG	AGG	G
G	U	GUU Val	GCU	GAU Asp	GGU	U
	C	GUC	GCC Ala	GAC	GGC	C
	A	GUA	GCA	GAA Glu	GGA	A
	G	GUG	GCG	GAG	GGG	G

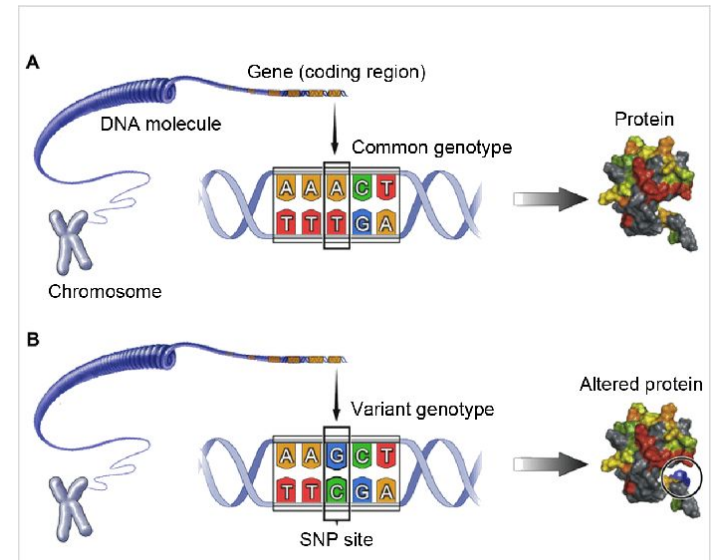
Single Nucleotide Polymorphism (SNP) is associated with a position in the genome:

- A – both chromosomes have A
- A / G – one chromosome has A, another one has G
- G – both chromosomes have G



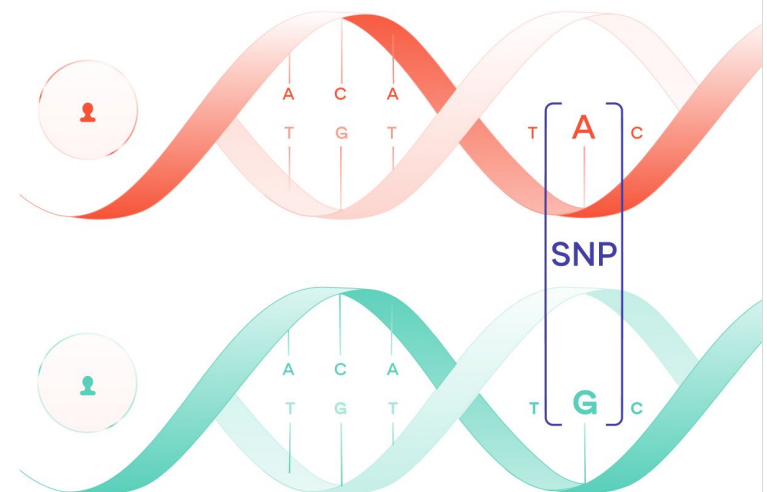
Altered genes produce altered proteins

		Second nucleotide				
		U	C	A	G	
U	UUU	Phe	UCU	Tyr	UGU	Cys
	UUC		UCC	Ser	UGC	
	UUA	Leu	UCA		UAA	STOP
	UUG		UCG		UAG	STOP
C	CUU		CCU	His	CGU	
	CUC	Leu	CCC	Pro	CGC	Arg
	CUA		CCA		CGA	
	CUG		CCG		CGG	
A	AUU	Ile	ACU	Asn	AGU	Ser
	AUC		ACC	Thr	AGC	
	AUA		ACA		AGA	Arg
	AUG	Met	ACG		AAG	
G	GUU		GCU	Asp	GGU	
	GUC	Val	GCC	Ala	GGC	Gly
	GUA		GCA		GGA	
	GUG		GCG		GAG	

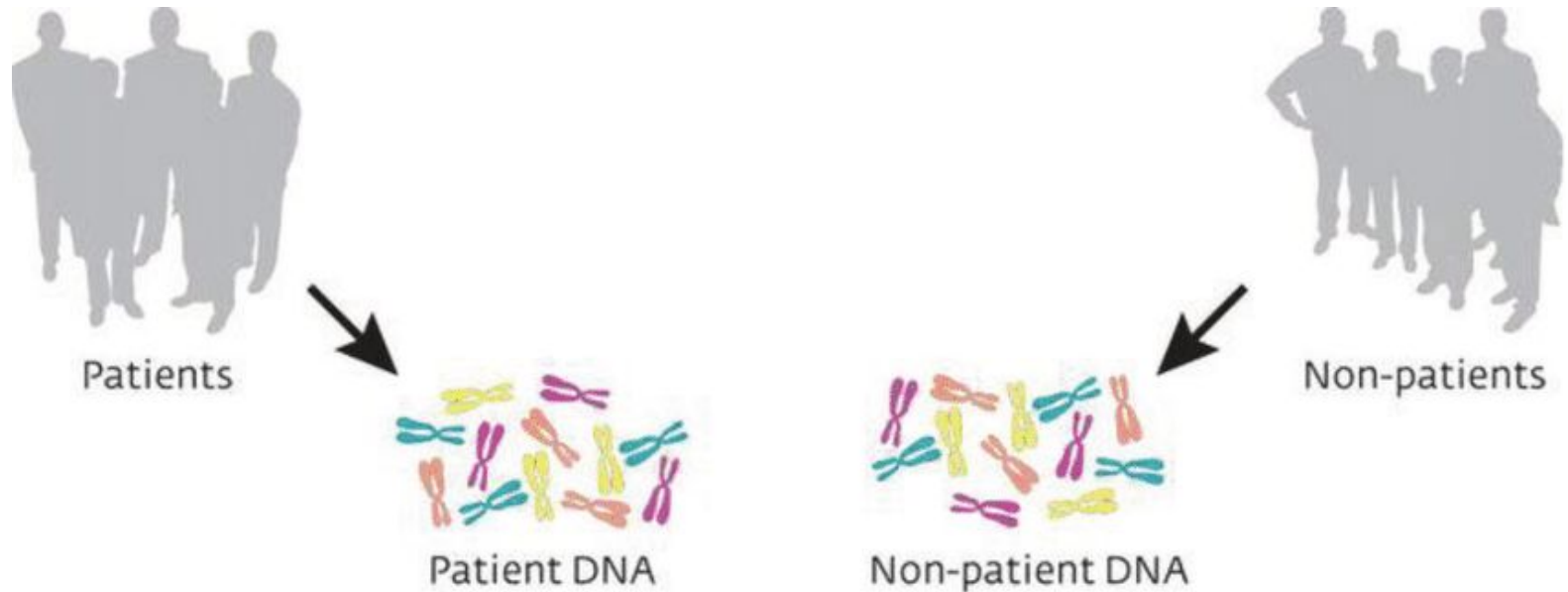


Single Nucleotide Polymorphism (SNP) is associated with a position in the genome:

- A – both chromosomes have A
- A / G – one chromosome has A, another one has G
- G – both chromosomes have G



Genome-wide association studies (GWAS)

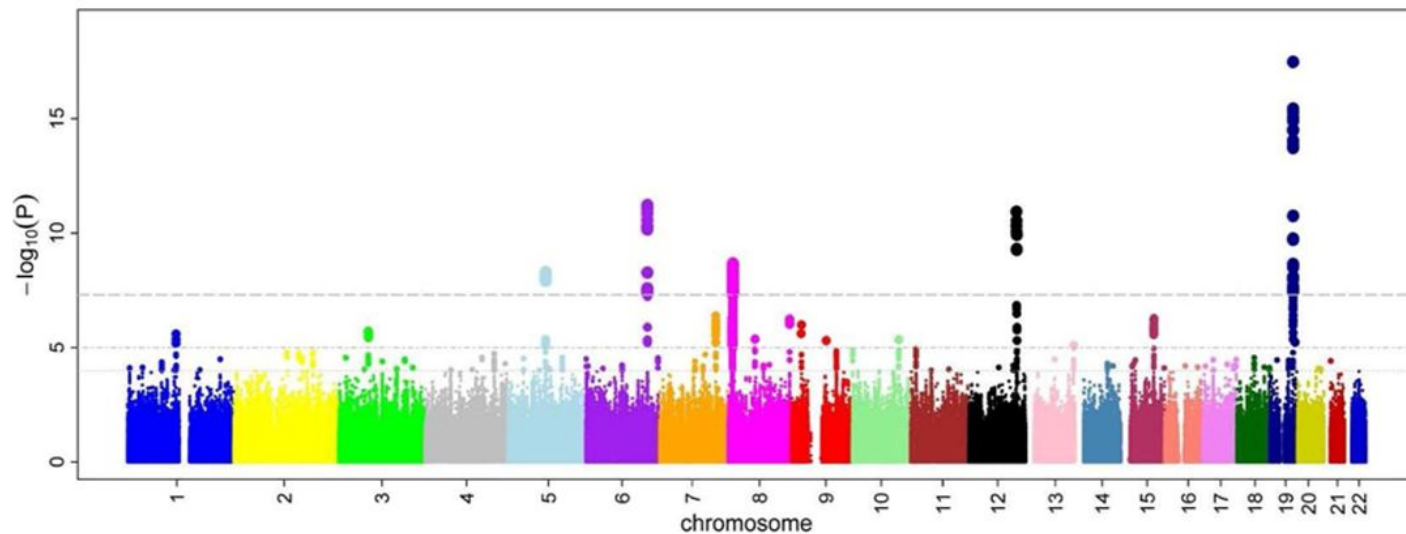
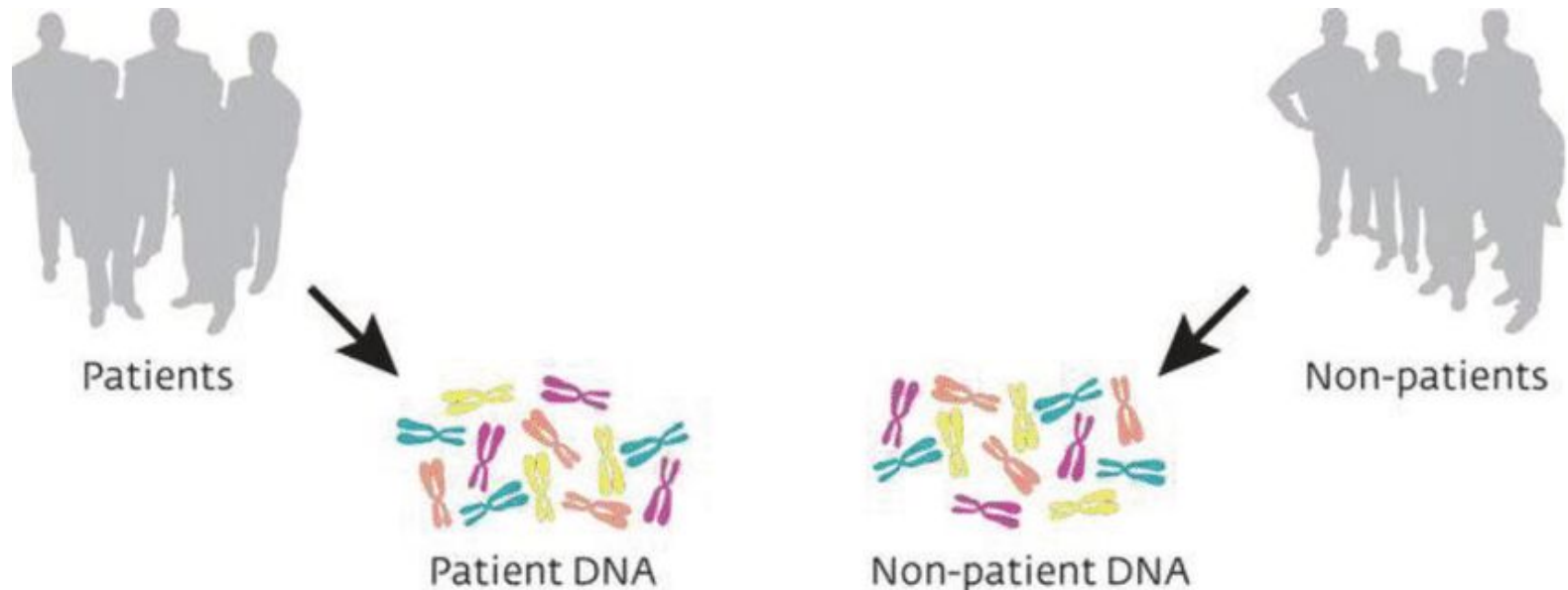


	A	A/G	G
Patients	15	10	8
Controls	2	6	25

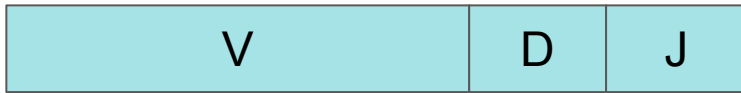
a single SNP

P-value (Fisher exact probability test) = 0.000026

Genome-wide association studies (GWAS)



Variants of IGHV1-69 shape Ab response to flu



IGHV1-69

IGHV1-69*01

54: F

IGHV1-69*02

54: L

IGHV1-69*03

54: F

IGHV1-69*04

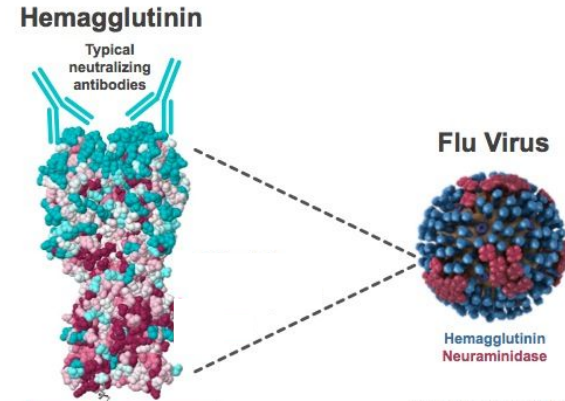
54: L

...

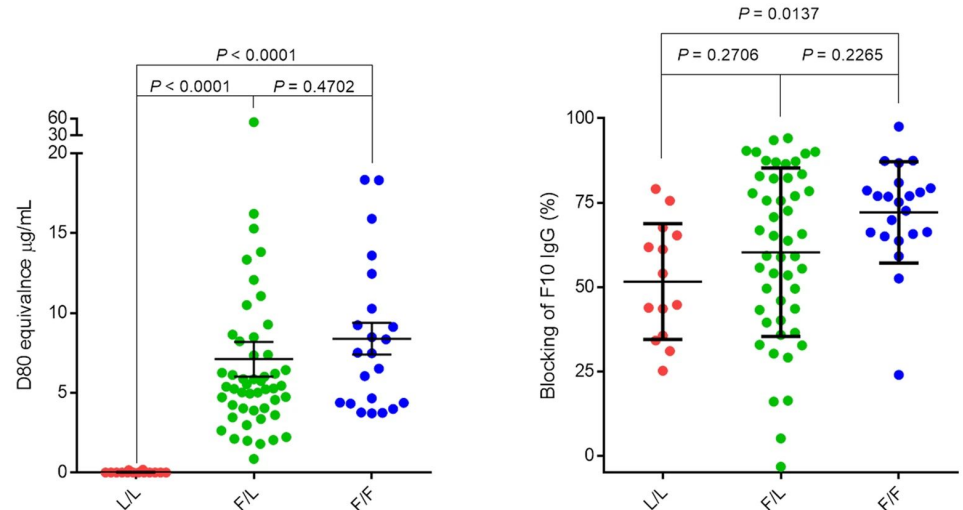
IGHV1-69*14

54: F

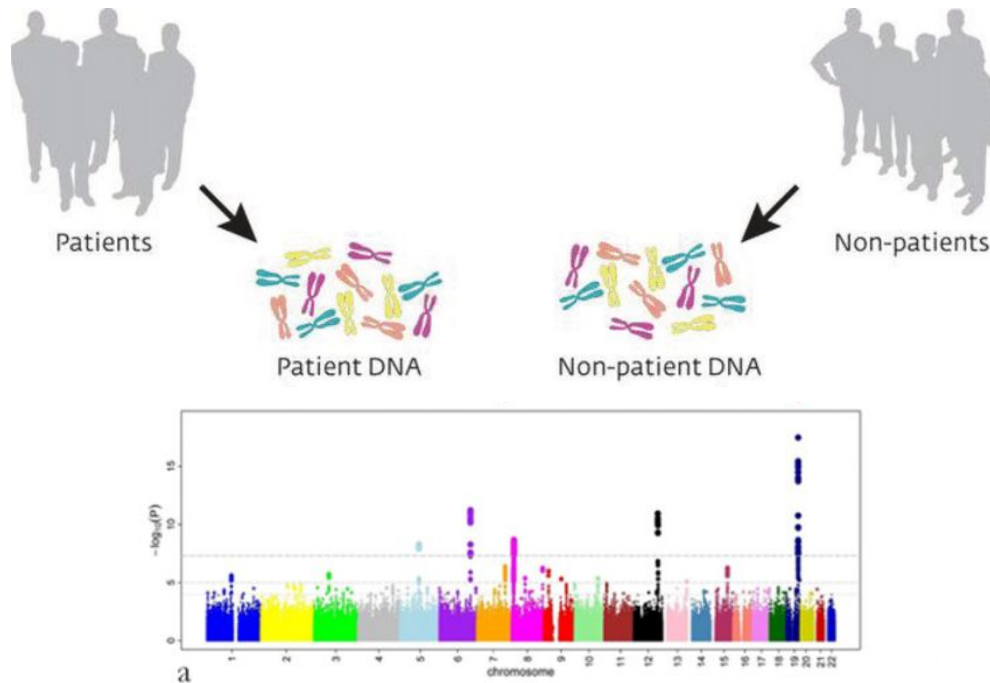
Successful binding to hemagglutinin
Loss of binding properties



Titers (= antibody counts) before and after immunization



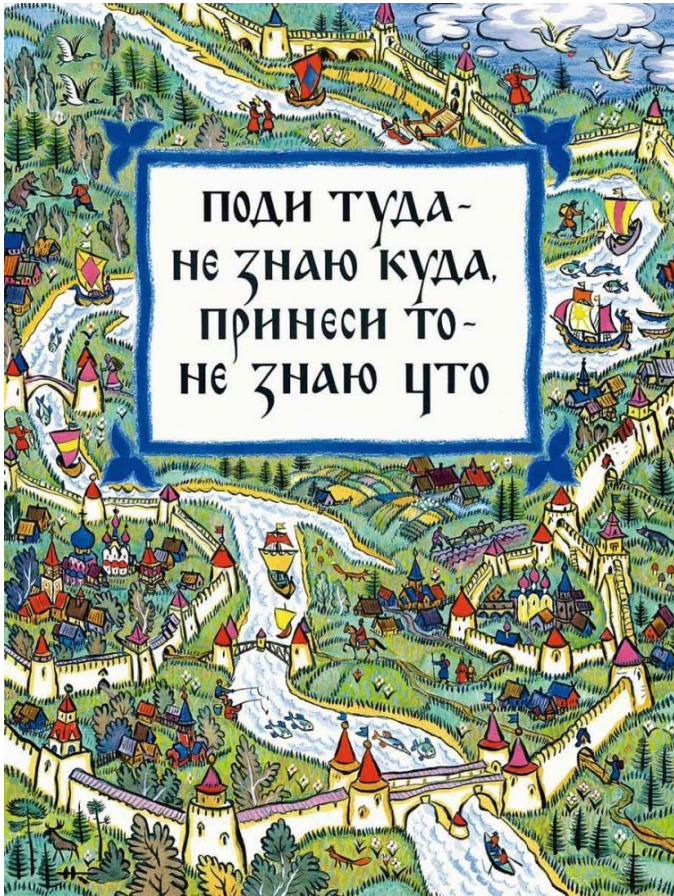
GWAS of antibody responses



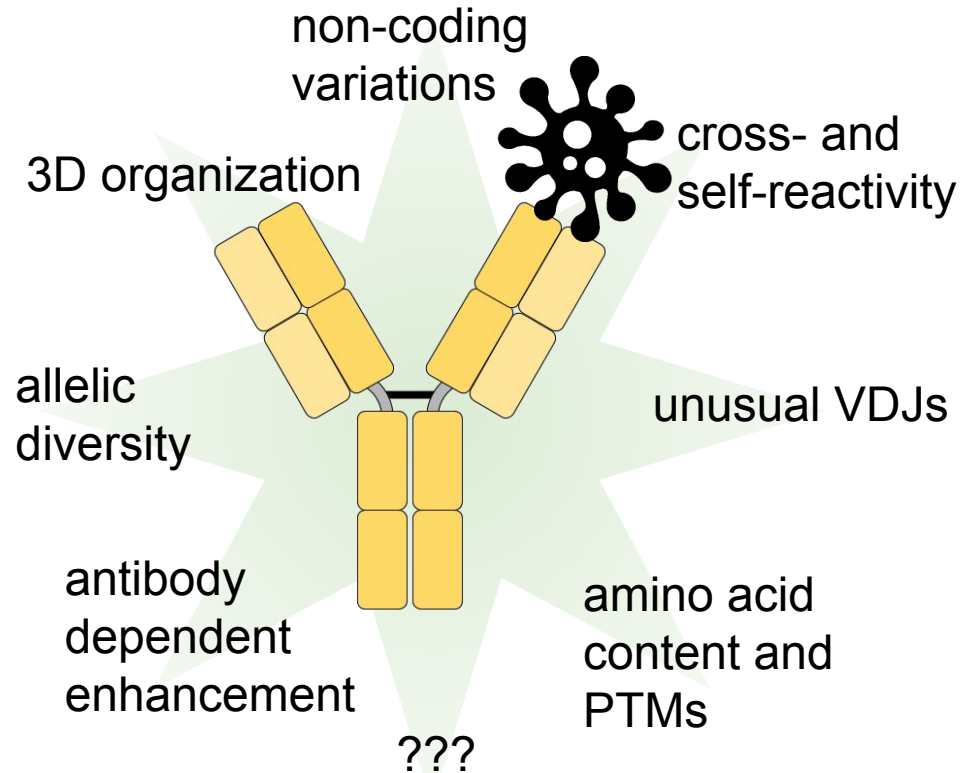
Challenges

- A single IG locus includes many V, D, and J genes. **If a single gene loses functionality, in most cases others can replace it**
- Recent studies report a lack of associations between genomic variations outside IG loci and adaptive immune responses
- Many other factors (age, diet, environment) influence adaptive immune responses
- **Features of antibody repertoires go far beyond genomic variations**

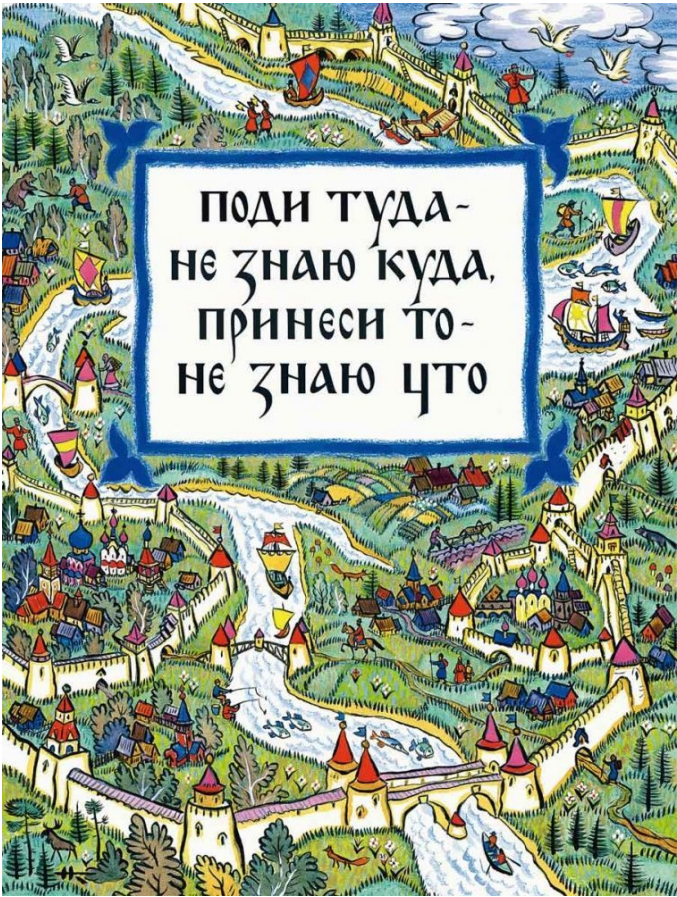
Bioinformatics + immunology = immunoinformatics



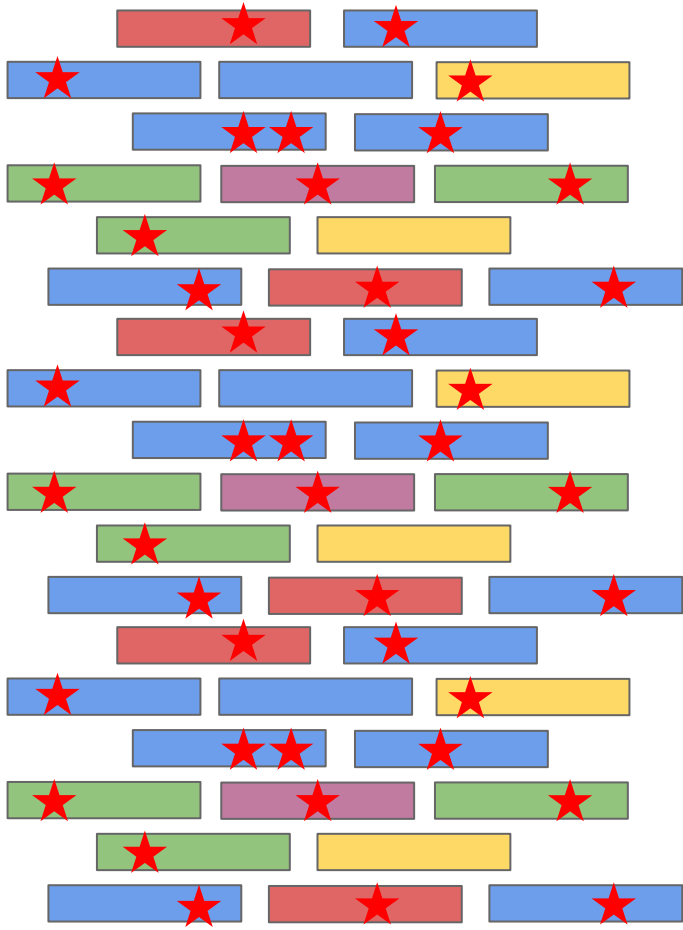
Go I Know Not Whither and
Fetch I Know Not What



Bioinformatics + immunology = immunoinformatics

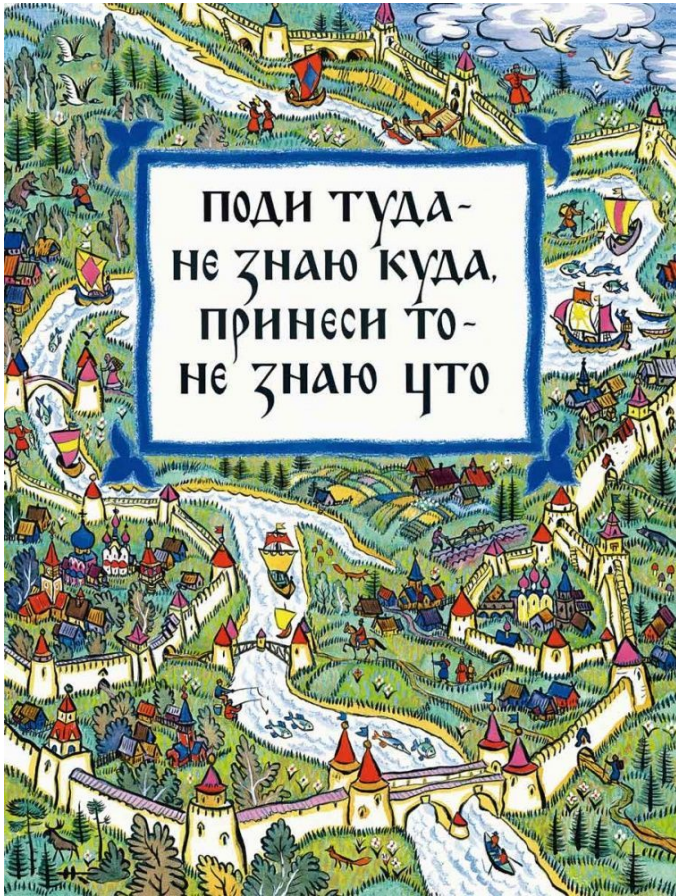


Go I Know Not Whither and Fetch I Know Not What

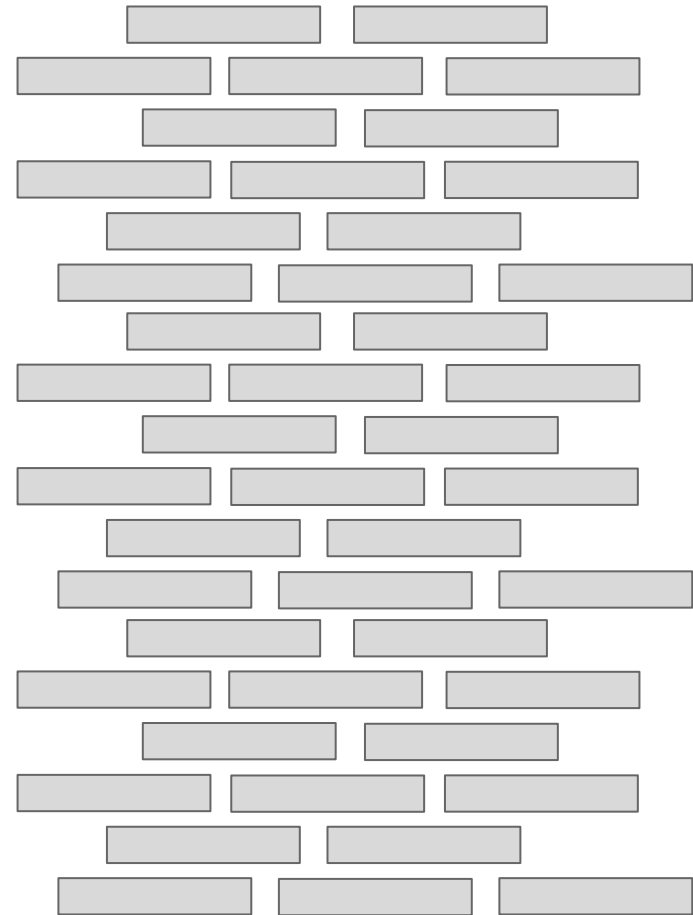


Error-prone immunosequencing reads

Bioinformatics + immunology = immunoinformatics



Go I Know Not Whither and Fetch I Know Not What

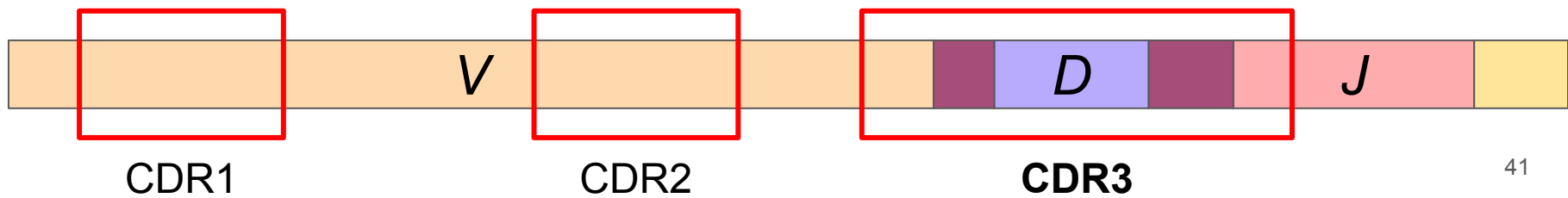
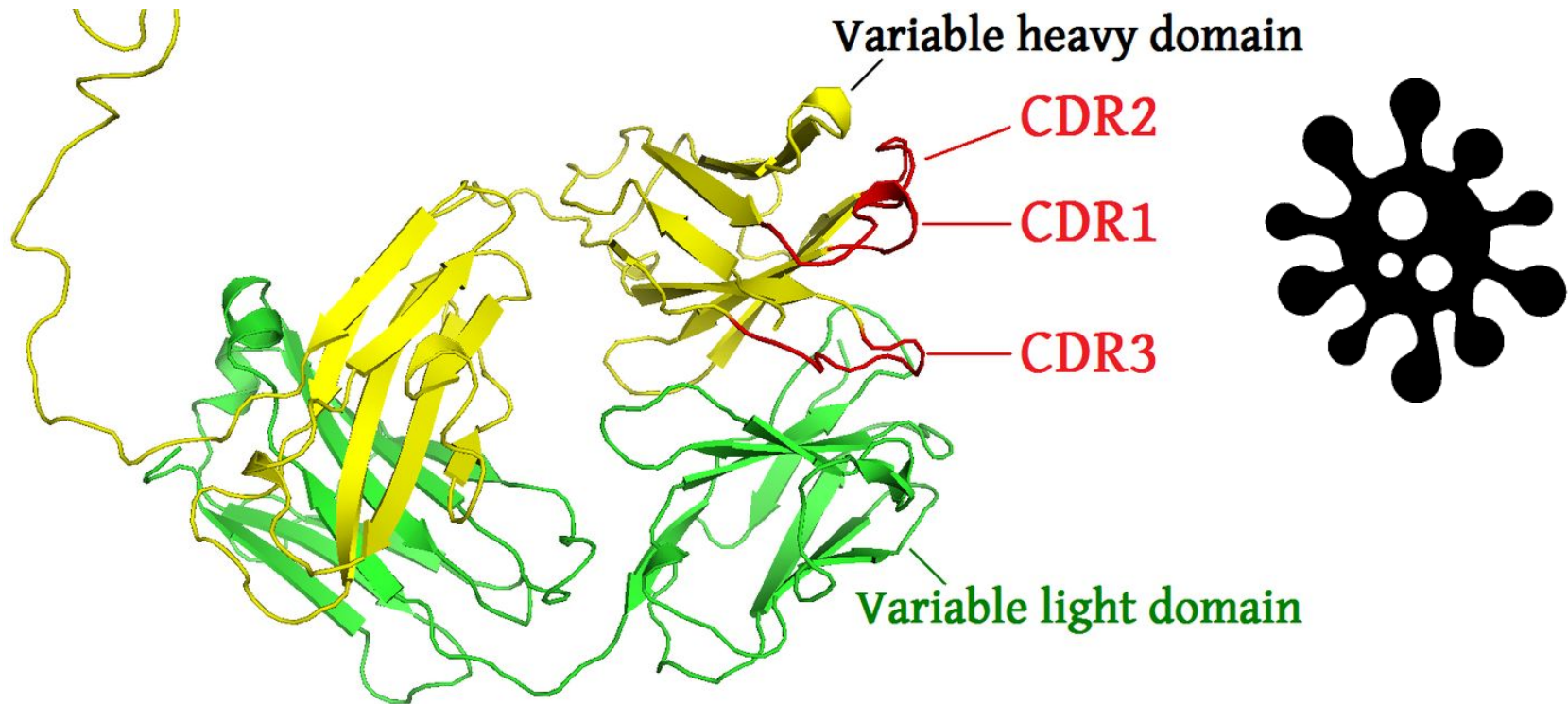


Error-prone immunosequencing reads



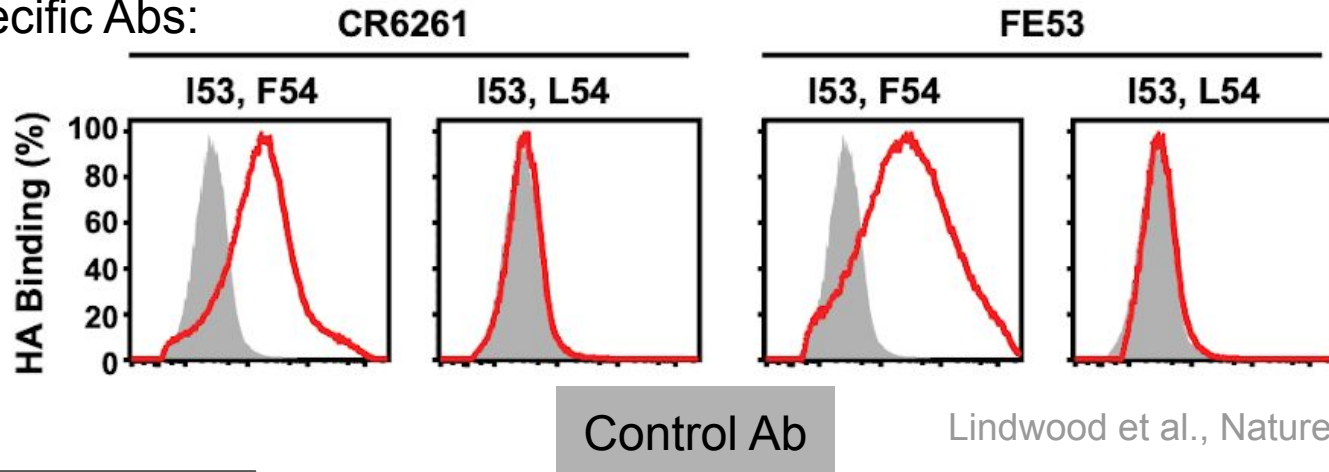
~ Flu ~

CDRs represent antigen-binding sites

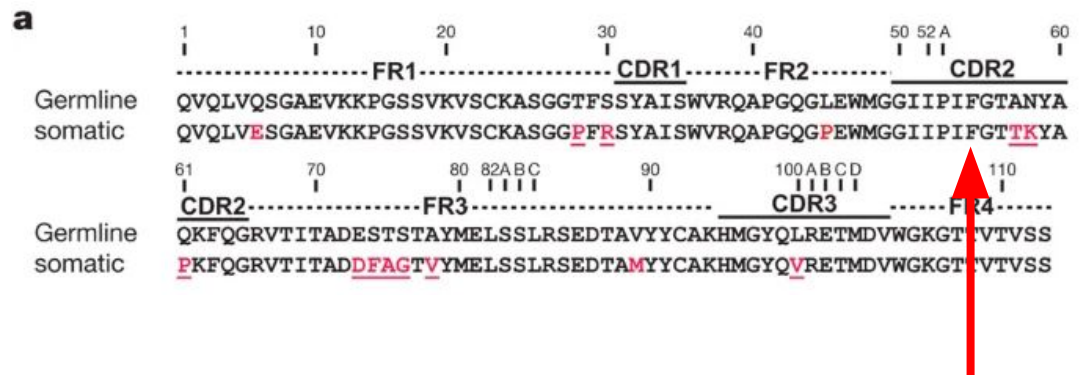


Anatomy of IGHV1-69-guided response to flu

Two flu-specific Abs:

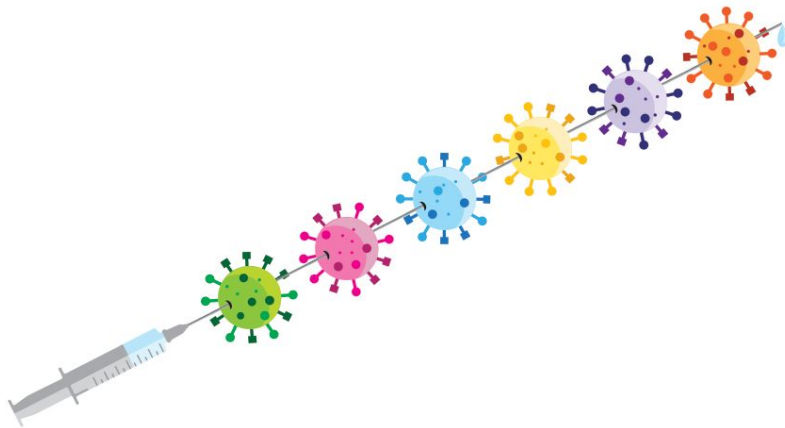


IGHV1-69*01	54: F
IGHV1-69*02	54: L
IGHV1-69*03	54: F
IGHV1-69*04	54: L
...	
IGHV1-69*14	54: F



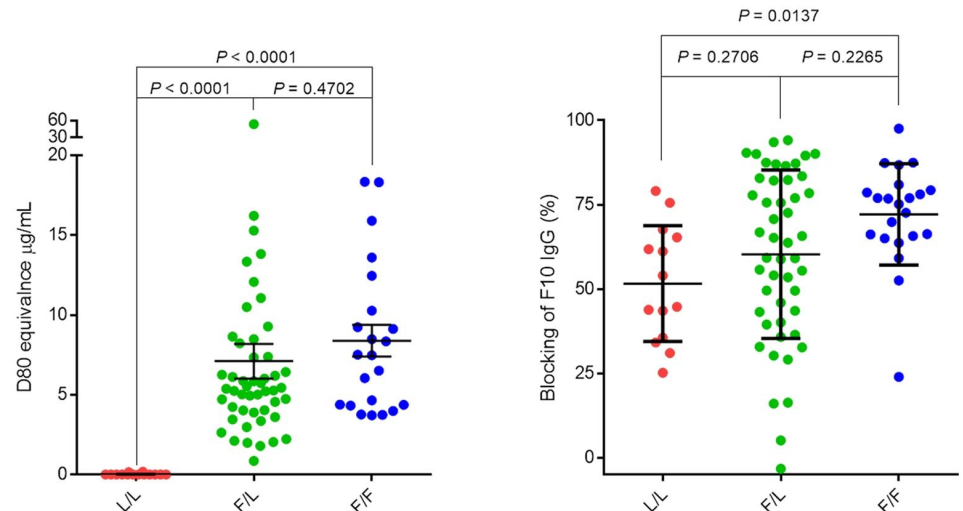
Antibody titers vs IGHV1-69 genotype

The titer data suggests that genotype of IGHV1-69 shapes the response to flu and other V genes do not fully replace the “bad” variant of IGHV1-69



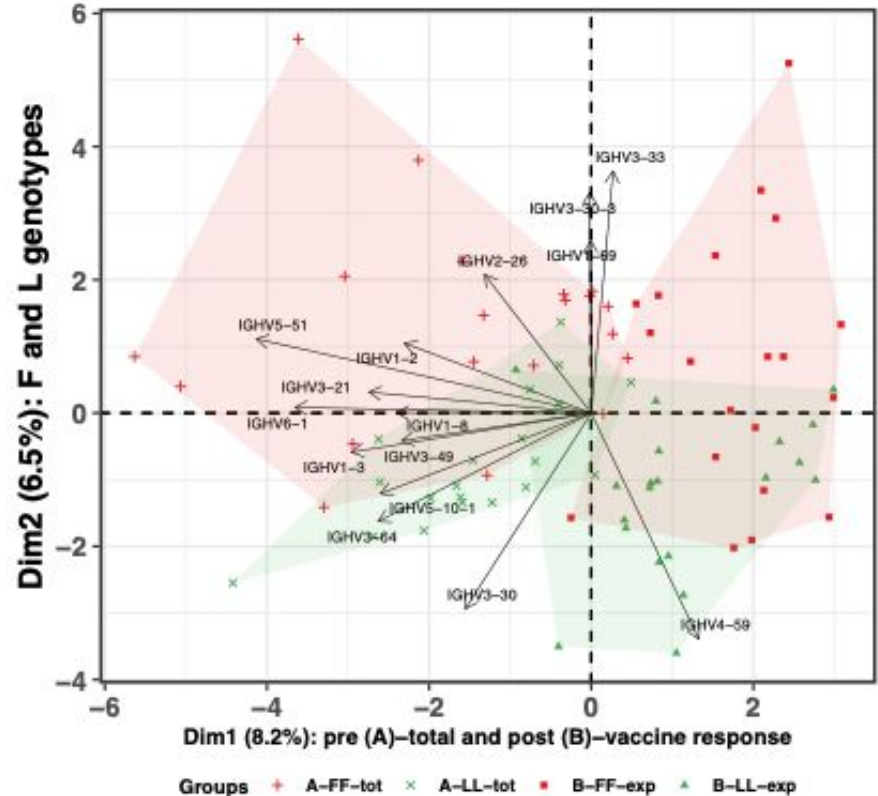
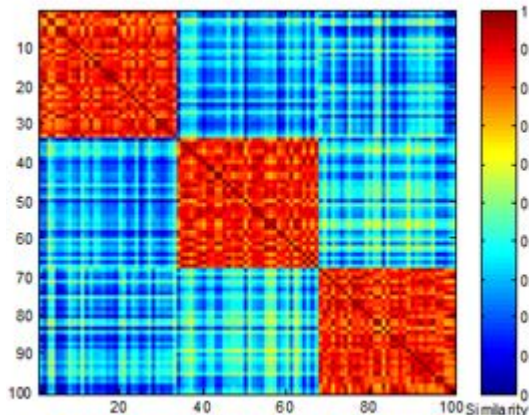
~20% of people have L/L variants of IGHV1-69

Titers (= antibody counts) before and after immunization



Bioinformatics analysis of flu response

- Usage of gene G = the fraction of VDJ sequences derived from G
- Individual antibody repertoire can be described as a *usage vector* for all existing genes
- We can compare usage vectors for F/F, F/L, and L/L individuals:



Ke, Nouri, Safonova, et al., in preparation

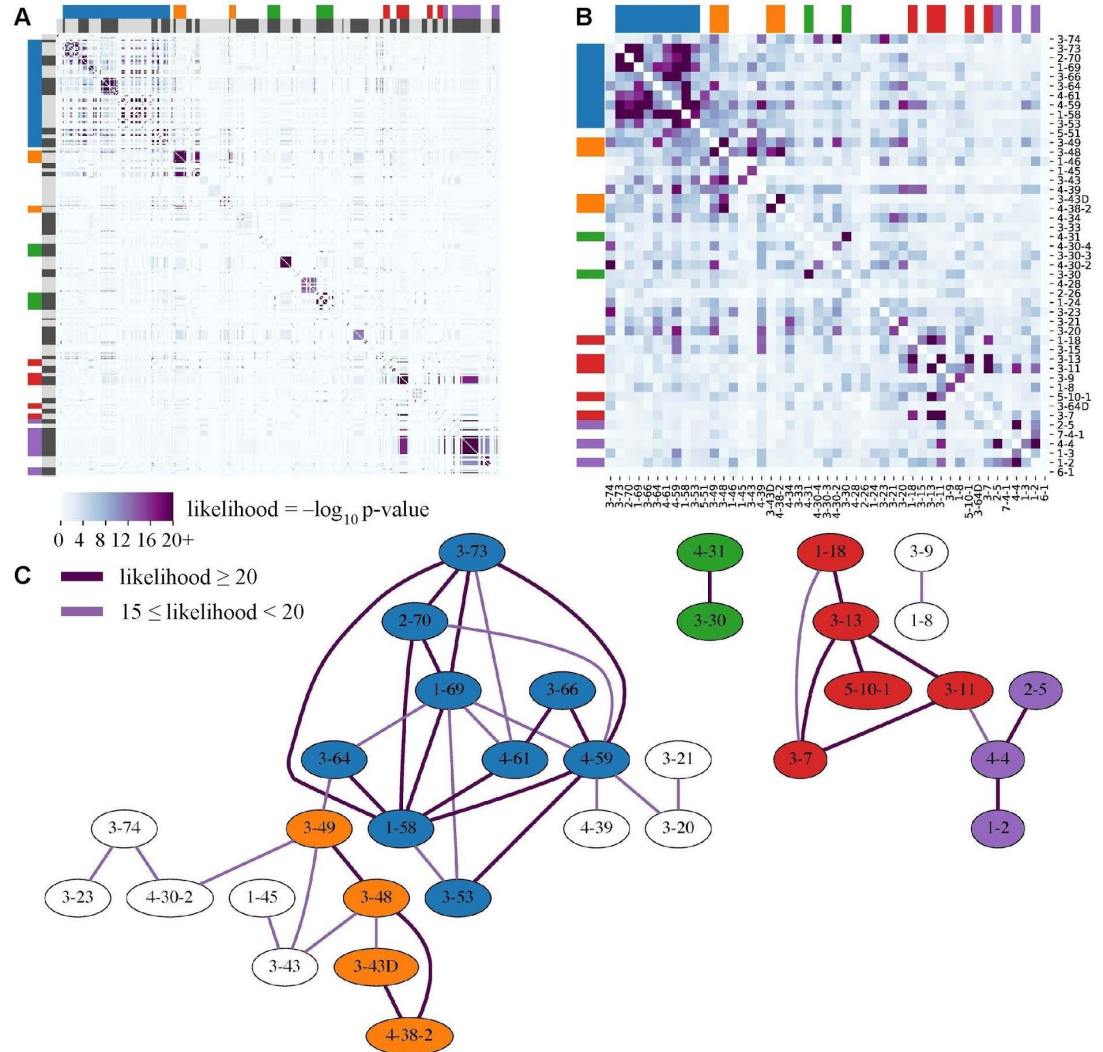
Linked variations of V genes

	A	A/G	G
C	15	10	8
A/C	6	8	3
A	2	5	17

P-value = 0.00967

Blue group is a set of V genes associated with flu response

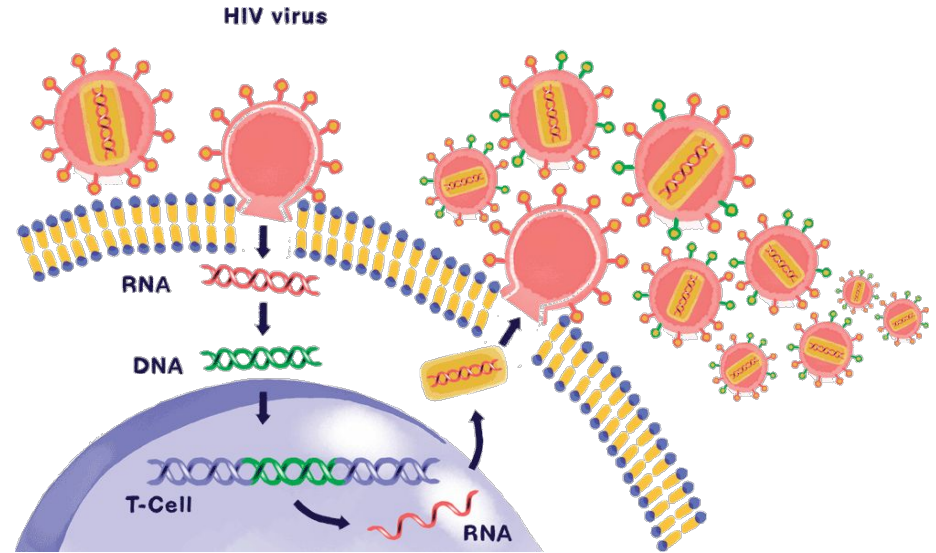
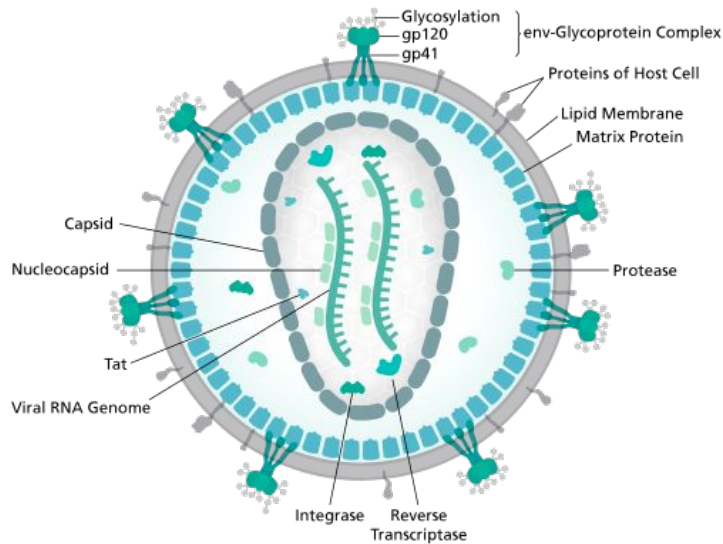
These genes have linked variations and perhaps close specificities



A microscopic view of blood cells and HIV virus particles. The background is a warm, reddish-orange color. Numerous red blood cells, appearing as biconcave discs, are scattered throughout. Several green, spherical HIV virus particles with a spiky surface are also visible. The text '~ HIV ~' is centered in the image.

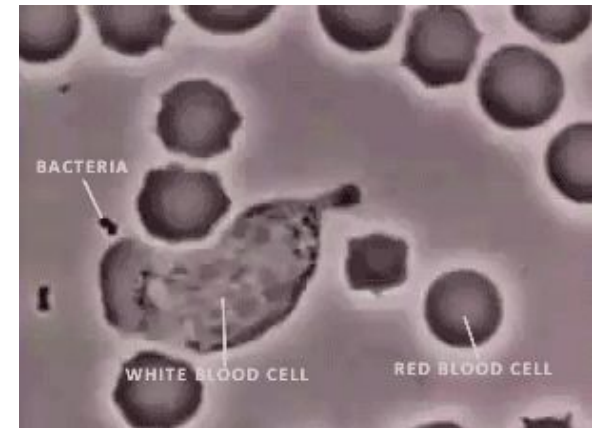
~ HIV ~

HIV and antibody response

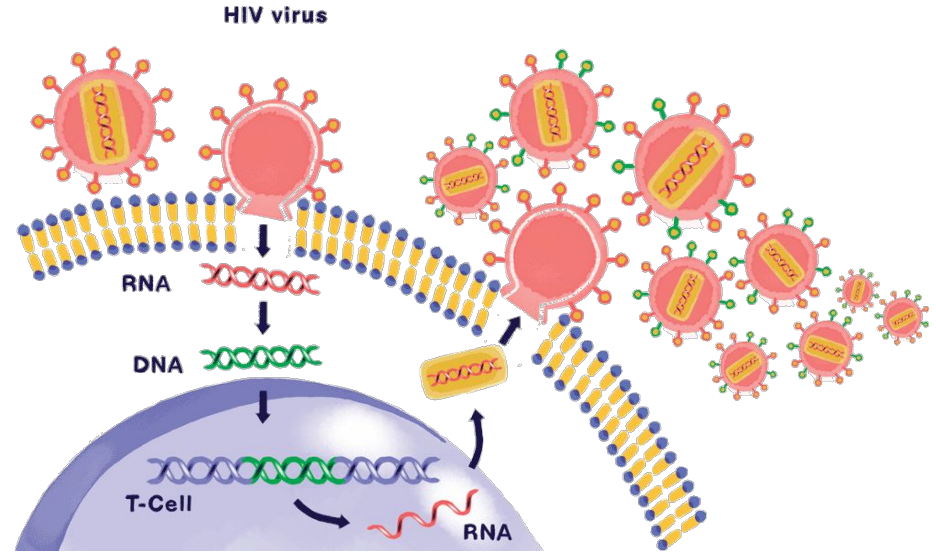
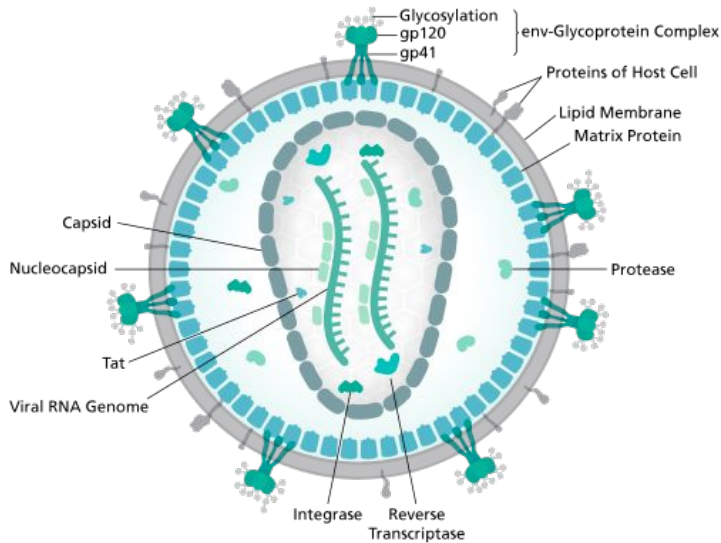


HIV infects immune cells:

- T cells
- Macrophages
- Dendritic cells



HIV and antibody response

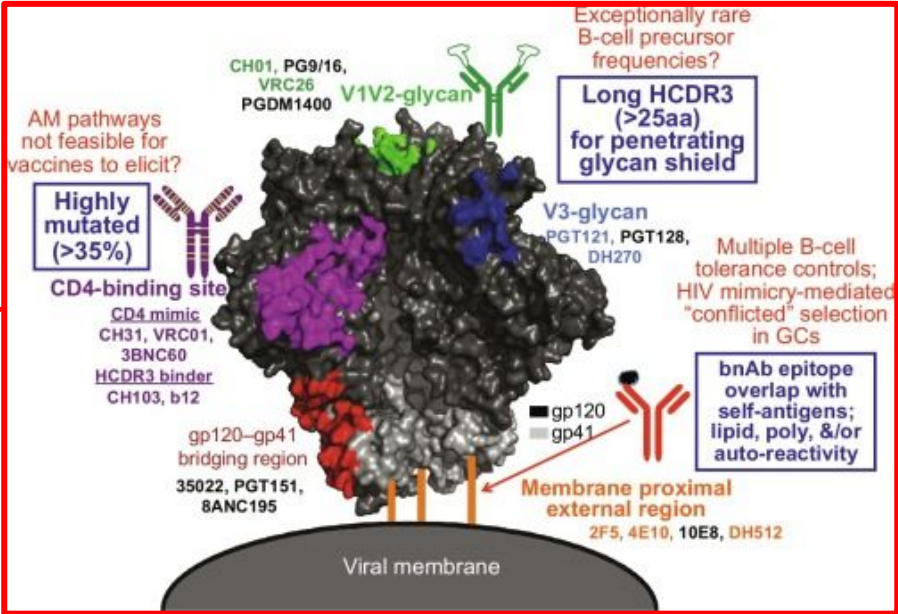
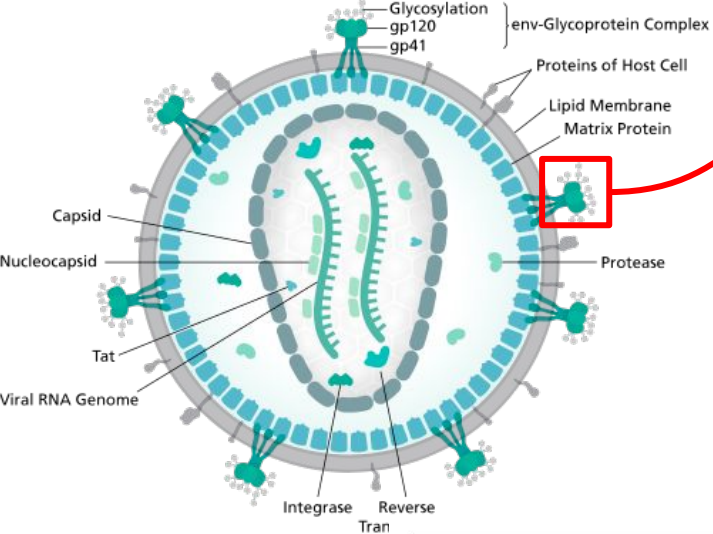


HIV infects immune cells:

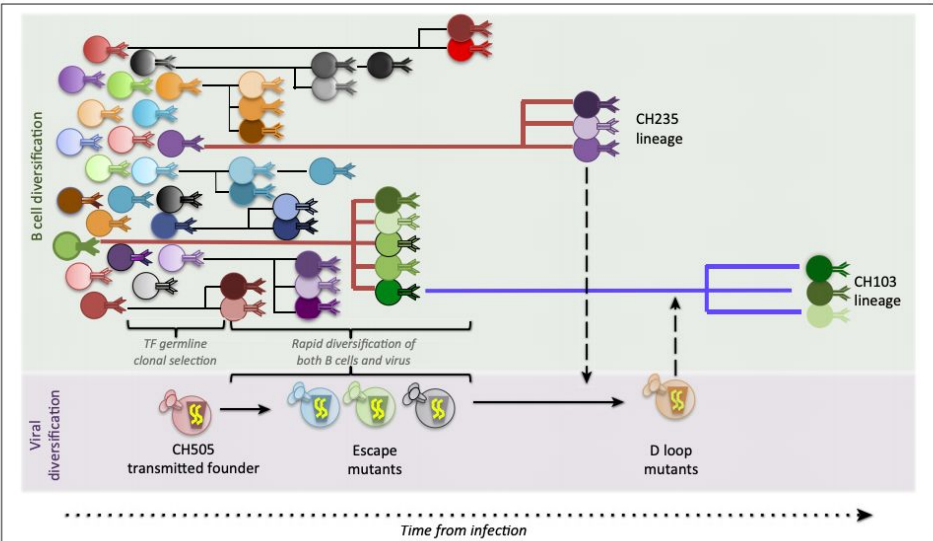
- T cells
- Macrophages
- Dendritic cells

B cells producing antibodies are not affected by HIV but cannot fight it because of a high mutation rate of the HIV genome

Recognition of HIV

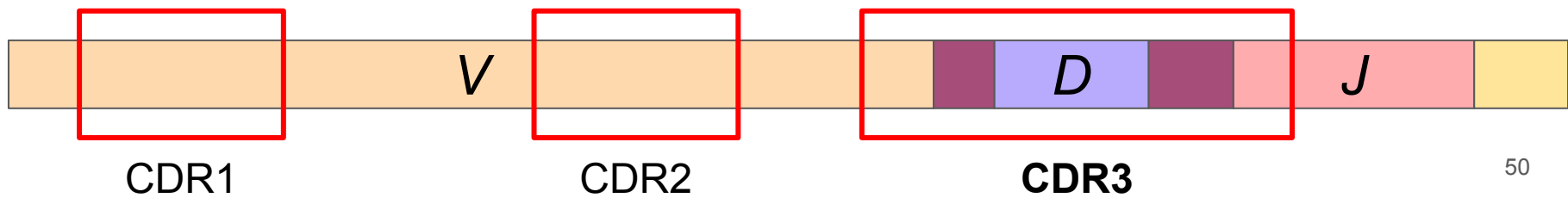
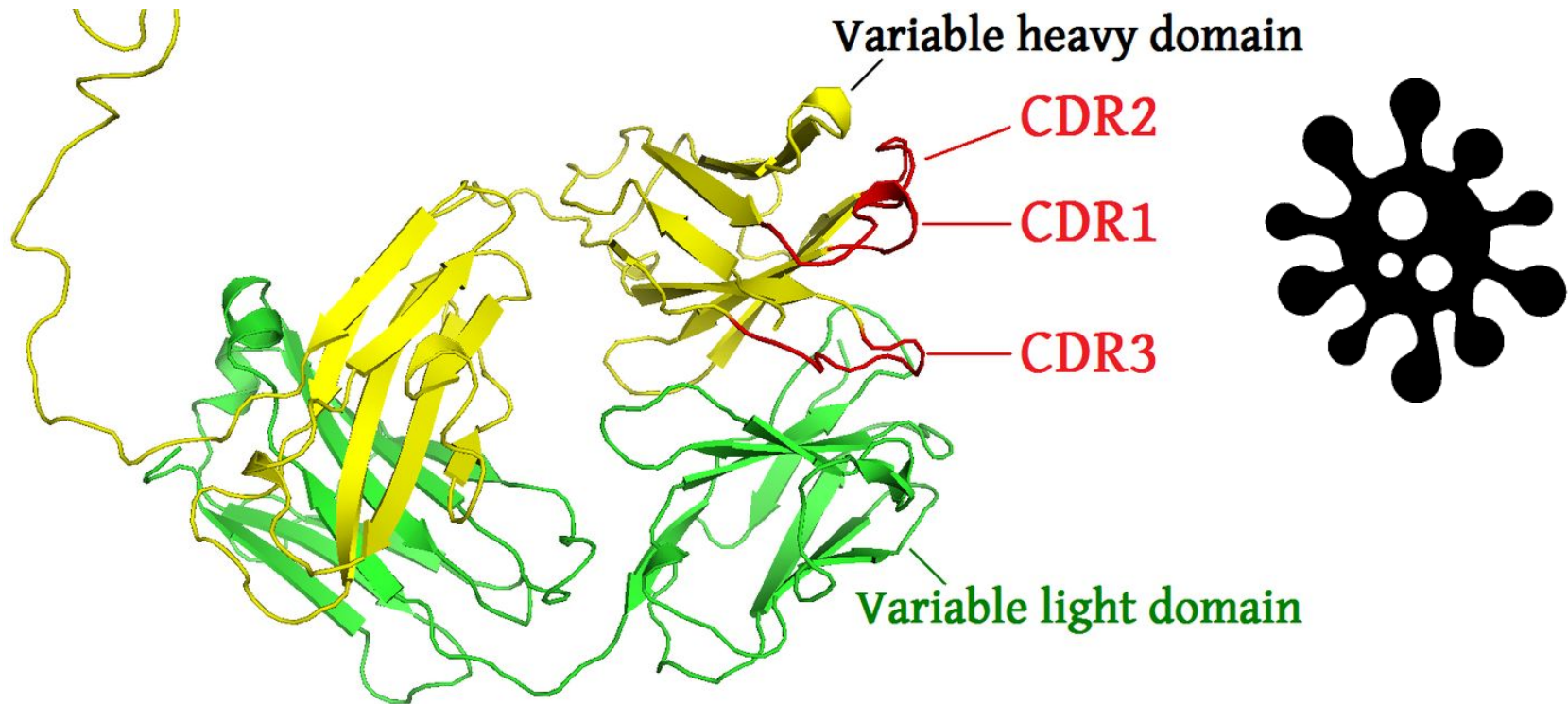


Verkoczy, Adv Immunol, 2017

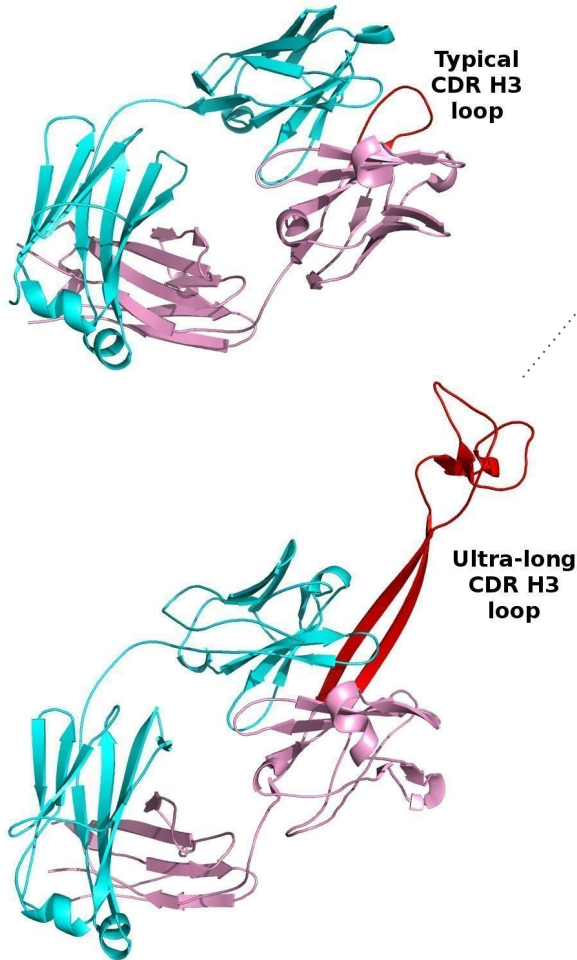


Alter & Ackerman, Cell, 2014

CDRs represent antigen-binding sites



Long CDR3s in response to HIV



Broadly neutralizing antibodies against HIV are characterized by **extremely long CDR3s**

Ultralong CDR3s in human antibodies present a result of **VDDJ** recombination:

D3: GTATTACGATTTTTGGAGTGGTTATtatacc
 D5: GTGGATACAGCTATGGttac
 CDR3: ACCACAGAACCGCTTCAGTTTAGTCCGTATTACGATTTTTGGAGTGGTTATCAGCCAGTGGATACAGCTATGGACCCGTTGACT
 span: GTATTACGATTTTTGGAGTGGTTATCAGCCAGTGGATACAGCTATGG
 inter-D insertion: CAGCCA

Finding VDDJ recombinations

D3: GTATTACGATTTTTGGAGTGGTTATtataacc
D5: GTGGATACAGCTATGGttac
CDR3: ACCACAGAACCGCTTCAGTTTAGTCCGTATTACGATTTTTGGAGTGGTTATCAGCCAGTGGATACAGCTATGGACCCGTTGACT
span: GTATTACGATTTTTGGAGTGGTTATCAGCCAGTGGATACAGCTATGG
inter-D insertion: CAGCCA

D10: GTATTACTATGGTTCggggagttattataac
D15: gtattatgattacgtttggGGGAGTTATGCttataacc
CDR3: GCGAGAGACACGTATTACTATGGTTCAGGGAGTTATGCGGCTAACAACTACTACTACTACGGTATGGACGTC
span: GTATTACTATGGTTCAGGGAGTTATGC
inter-D insertion: A

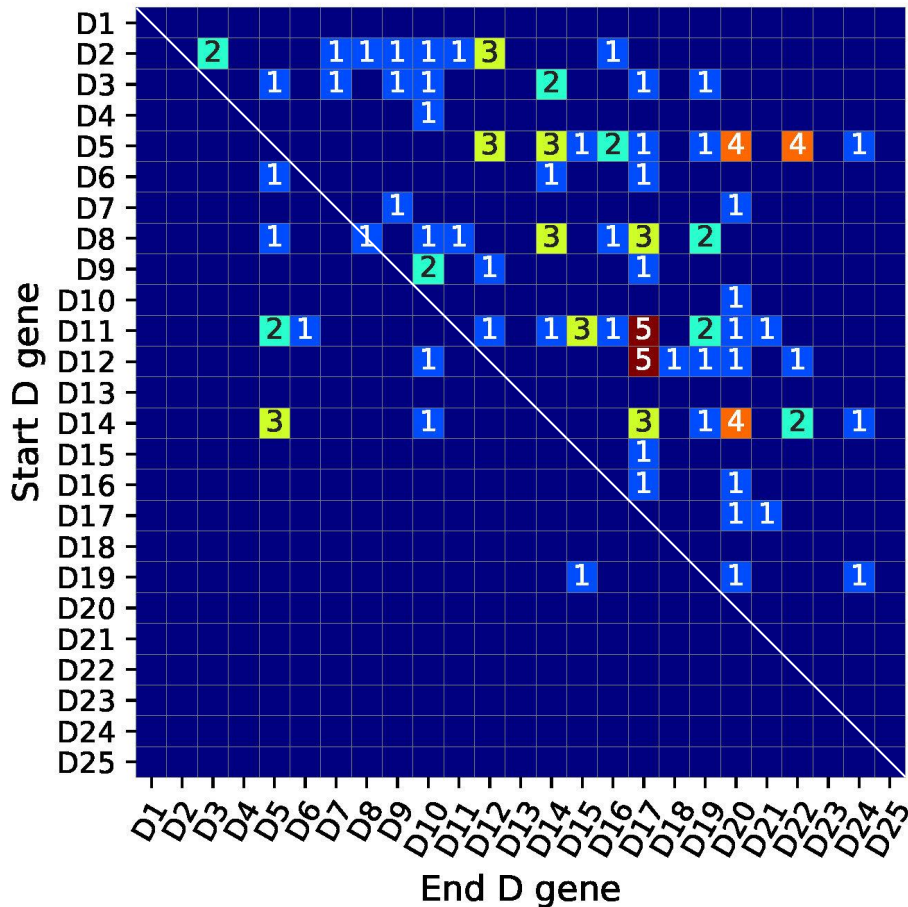
Finding VDDJ recombinations

D3: GTATTACGATTTTTGGAGTGGTTATtatacc
D5: GTGGATACAGCTATGGttac
CDR3: ACCACAGAACCCTTCAGTTTAGTCCGTATTACGATTTTTGGAGTGGTTATCAGCCAGTGGATACAGCTATGGACCCGTTGACT
span: GTATTACGATTTTTGGAGTGGTTATCAGCCAGTGGATACAGCTATGG
inter-D insertion: CAGCCA

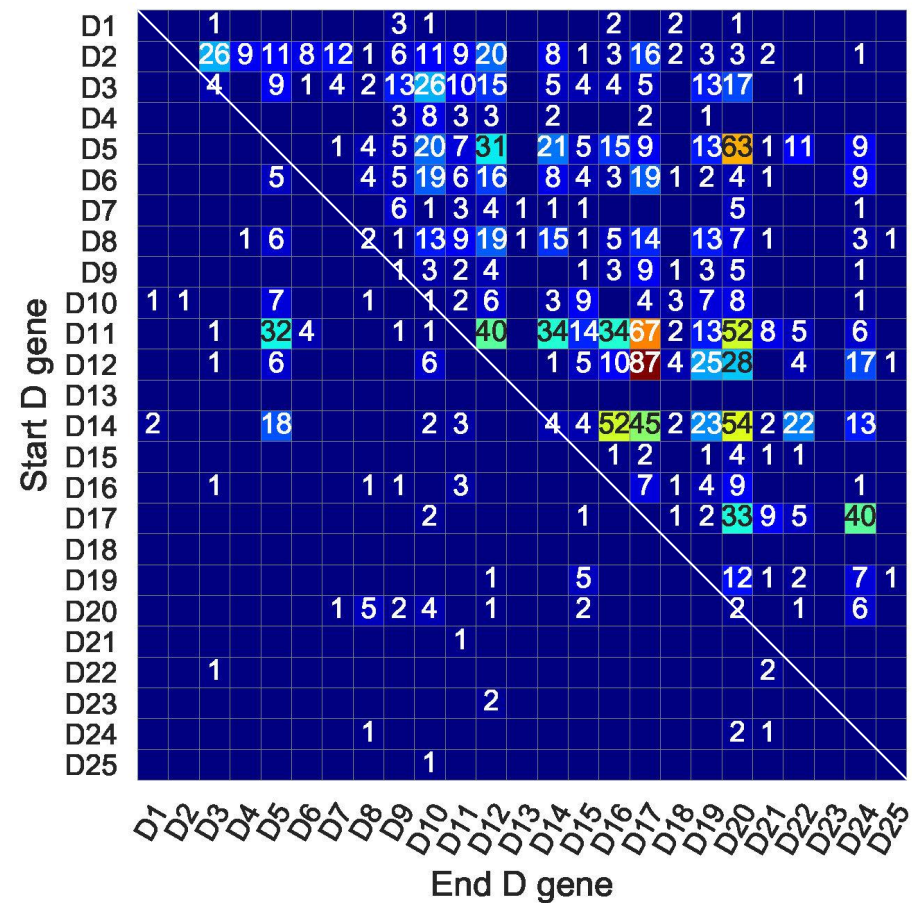
~~D10: GTATTACTATGGTTCggggagttattataac
D15: gtattatgatcaagtttggGGGAGTTATGCttataacc
CDR3: GCGAGAGACACGTATTACTATGGTTCAGGGGTTATGCTAACAACAACACTACTACTACTACGGTATGGACGTC
span: GTATTACTATGGTTCAGGGAGTTATGC
inter-D insertion: A~~

Most D-D pairs follow ordering in IGH locus

Set1

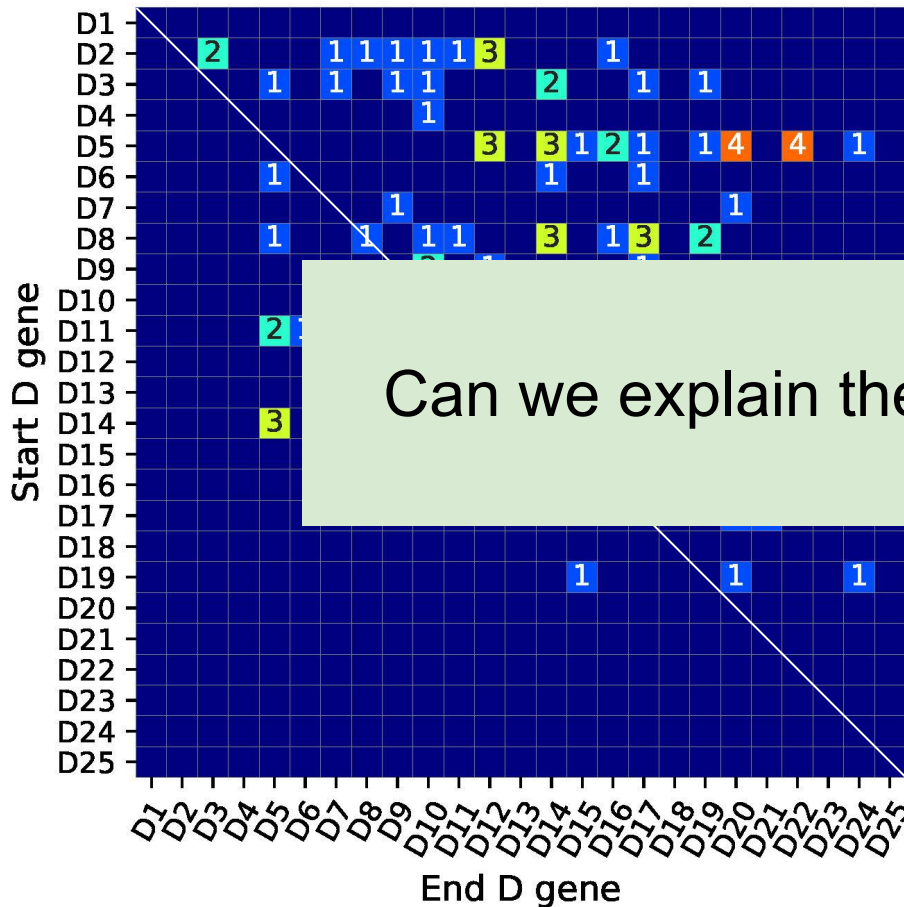


14 individuals

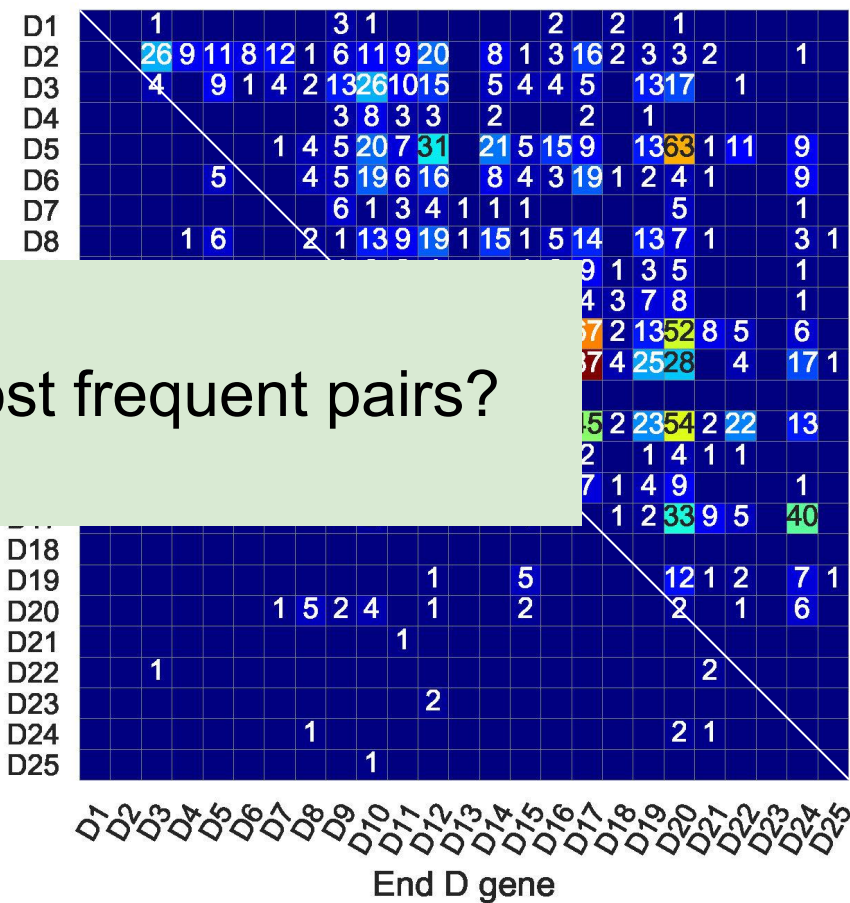


Most D-D pairs follow ordering in IGH locus

Set1

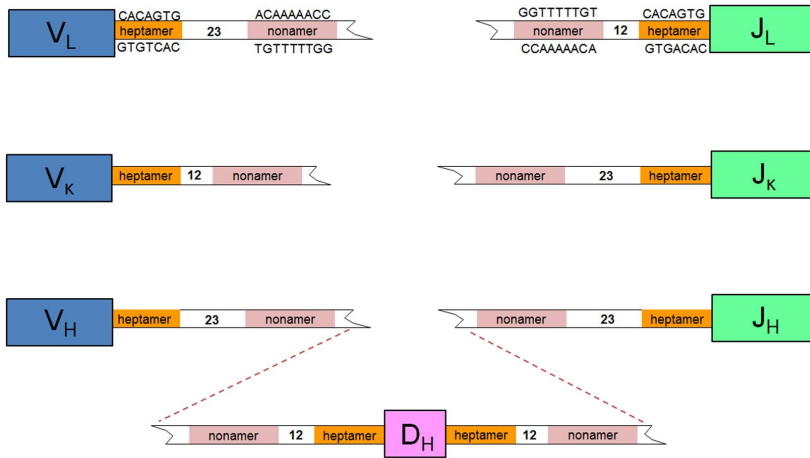


14 individuals



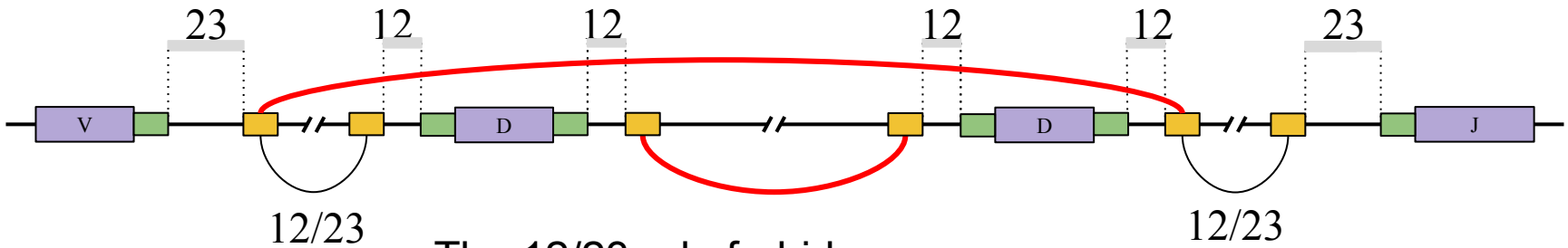
Can we explain the most frequent pairs?

Recombination signal sequences



V-RS		J-RS	
V-HEPTAMER	V-NONAMER	J-NONAMER	J-HEPTAMER
IGHV	TCAAAAACC 2	IGHJ	
	TCAGAAATC 2		
	ACATAAACT 1		
	ACAAAAATT 1		
CACAGCT 1	ACATAAACA 1		
CACAATG 4	ACAAATACT 1		
CACAGTG 29	ACATAAAACC 5		
	TCAGAAACC 5		
	CAGAAAACC 1		
	TCAGAAAC 5		
	ACACAAACC 1		
	ACAGAAACC 1		
		GGTTTTGT 2	CAATGTG 1
		GAGTTTTAG 1	GACTGTG 1
		ATTATTGT 1	TAGTGTG 1
			TATTGTG 1

5'D-RS		3'D-RS	
5'D-NONAMER	5'D-HEPTAMER	3'D-HEPTAMER	3'D-NONAMER
IGHD			
GATTTGAA 2	TACTGTG 13	CACGGTG 2	GCAAAAACC 2
GCTTTTGT 3	CACAGTG 5	CACTGTA 2	ACAAAAATC 1
GATTTTGT 10		CACAGTG 13	ACCAAAACT 2
GGATTTGA 1			ACAAAAACC 12
GGTTTTGAC 2			



The 12/23 rule forbids:

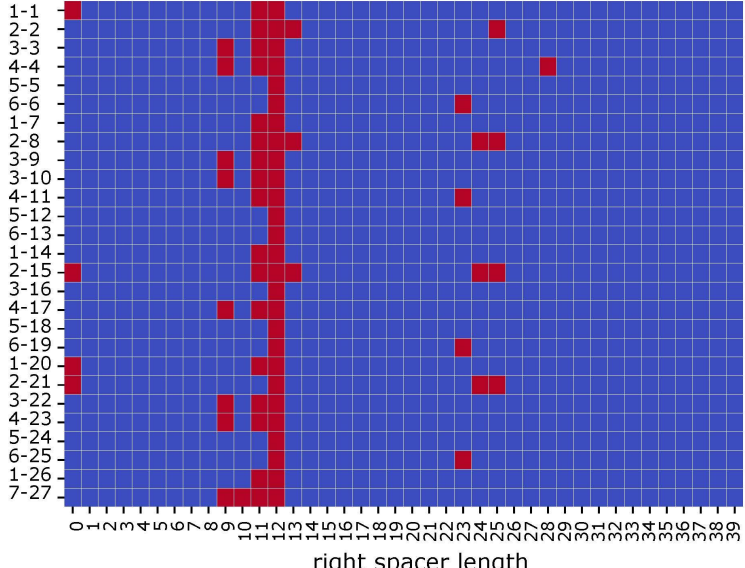
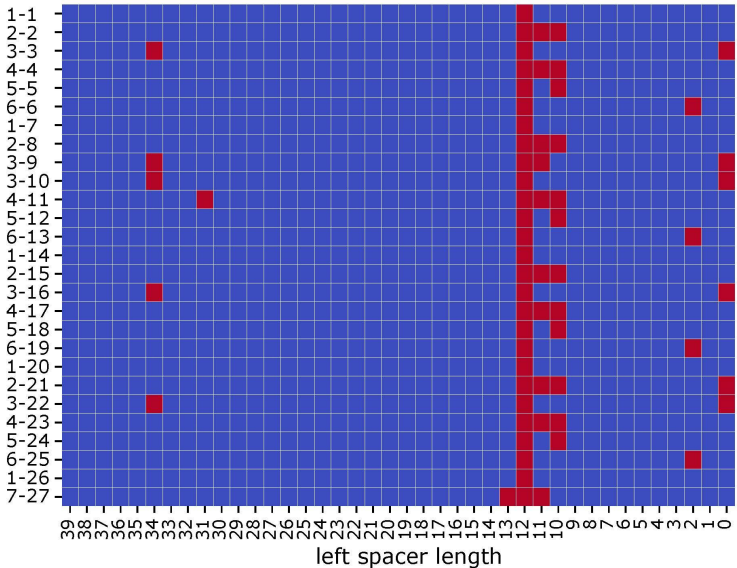
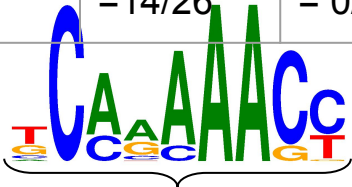
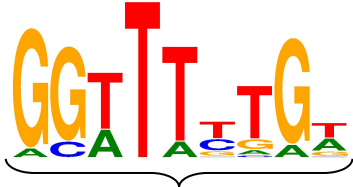
- V-J pairs
- D-D pairs
- V-V pairs
- J-J pairs

12 nt = 1 turn of DNA
23 nt = 2 turns of DNA

Non-canonical recombination signals

$$\text{Prob(ACGTACGTA)} = 11/26 * 25/26 * \dots$$

	Pos 1	Pos 2	...
A	= 11/26	= 1/26	...
C	= 1/26	= 25/26	...
G	= 0/26	= 0/26	...
T	= 14/26	= 0/26	...

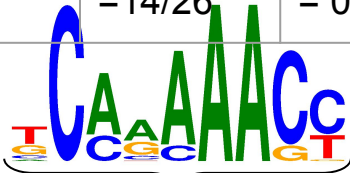
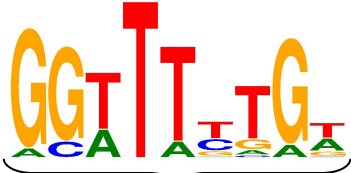


Non-canonical recombination signals

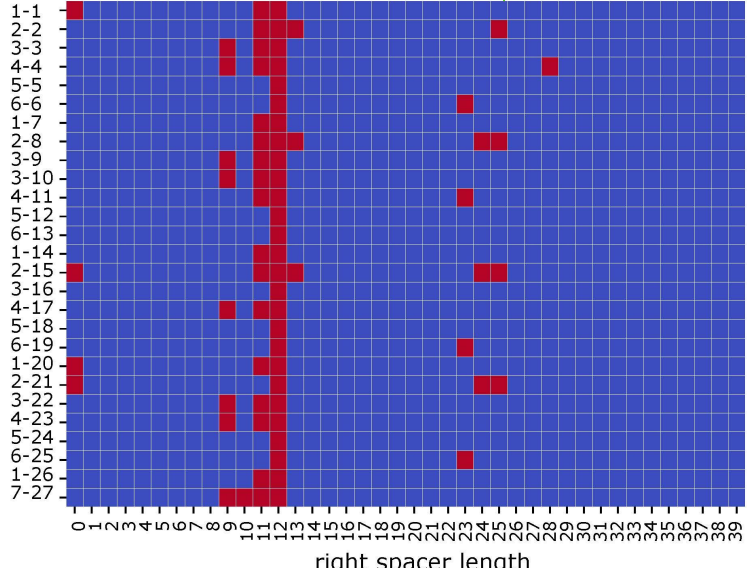
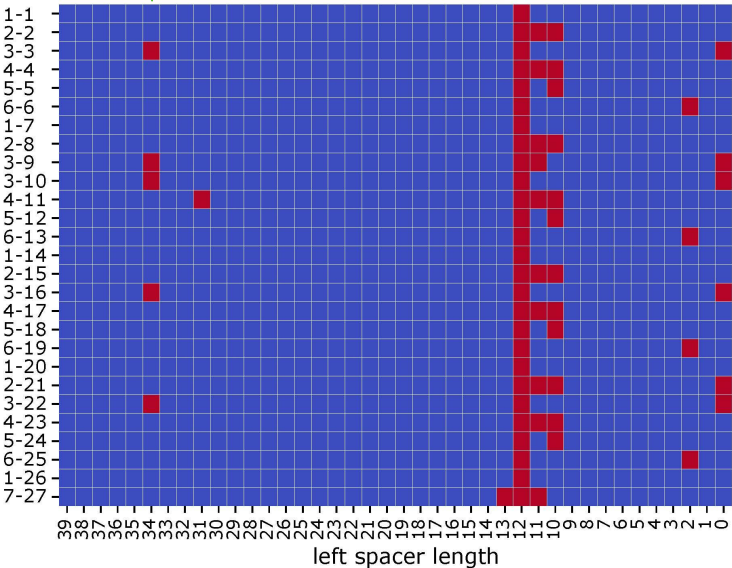
$$\text{Prob(ACGTACGTA)} = 11/26 * 25/26 * \dots$$

	Pos 1	Pos 2	...
A	= 11/26	= 1/26	...
C	= 1/26	= 25/26	...
G	= 0/26	= 0/26	...
T	= 14/26	= 0/26	...

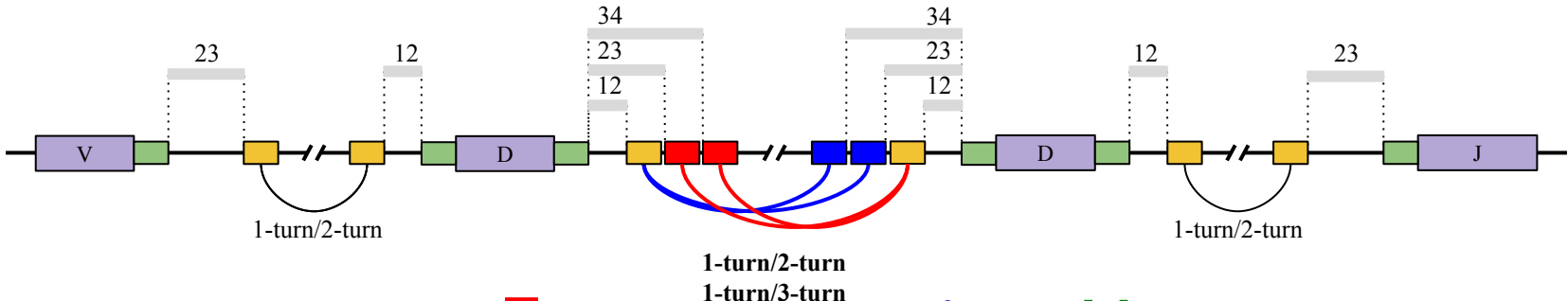
3 DNA turns



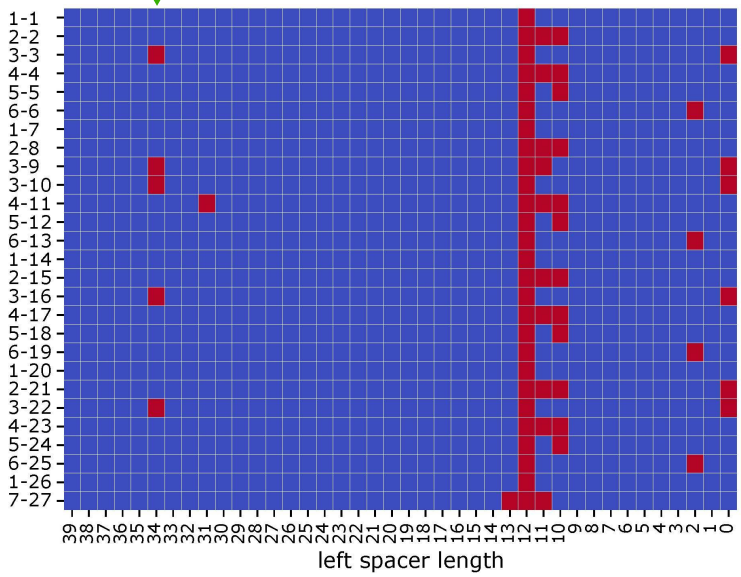
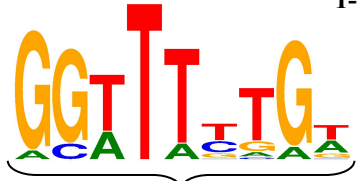
2 DNA turns



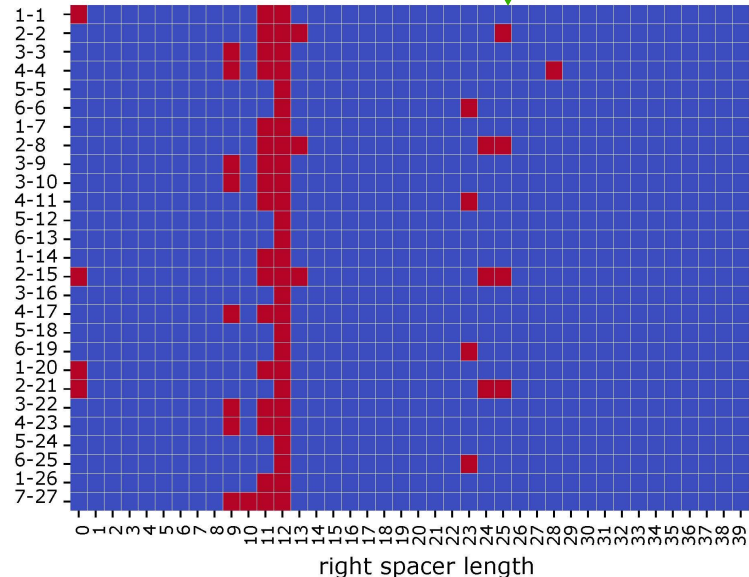
Non-canonical recombination signals



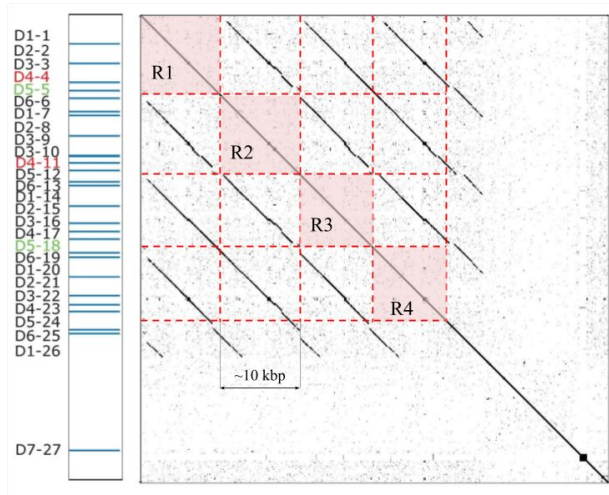
3 DNA turns



2 DNA turns



Tandem repeats correlate with D-D fusions

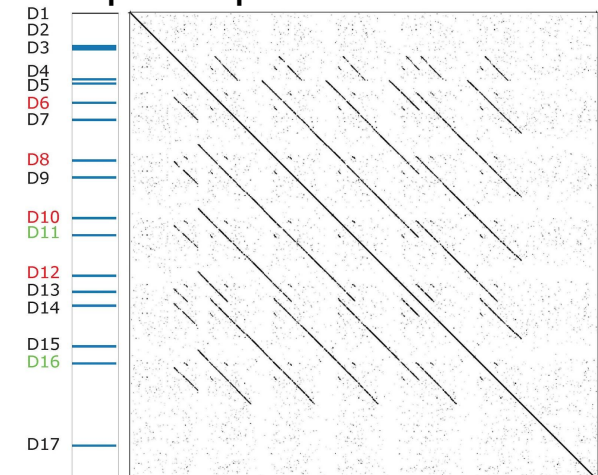
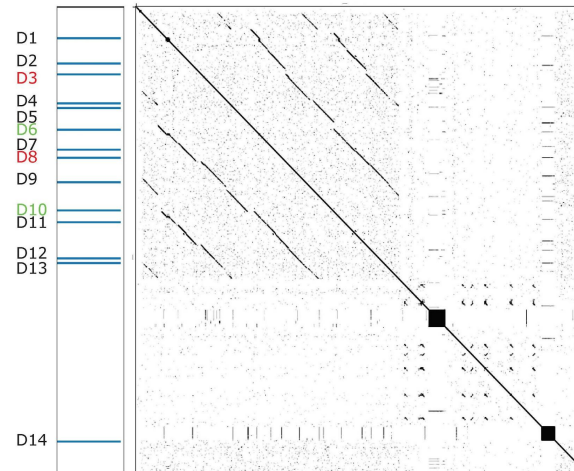
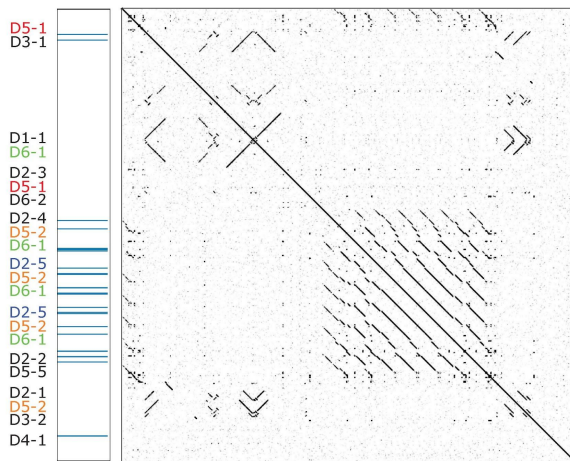


R1	R2	R3	R4
–	D6-6	D6-13	D6-19
D1-1	D1-7	D1-14	D1-20
D2-2	D2-8	D2-15	D2-21
D3-3	D3-9 D3-10	D3-16	D3-22
D4-4	D4-11	D4-17	D4-23
D5-5	D5-12	D5-18	D5-24

mouse

Common marmoset

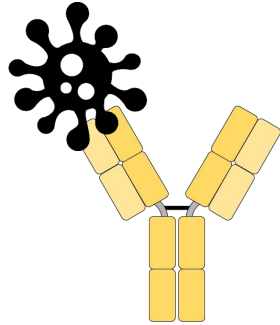
pale spear-nosed bat



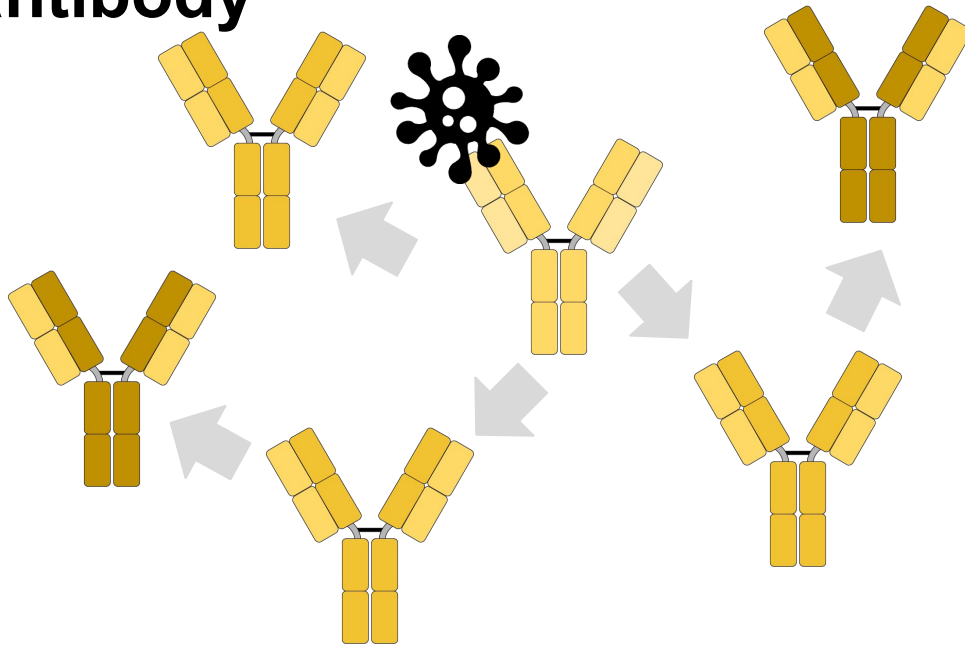
A transmission electron micrograph (TEM) showing several SARS-CoV-2 virus particles. The particles are roughly spherical with a distinct outer envelope and a darker, denser core. The background is dark, and the particles are scattered across the field of view. The text '~ SARS-CoV-2 ~' is overlaid in the center in white.

~ SARS-CoV-2 ~

Antibodies are subjects of fast evolution

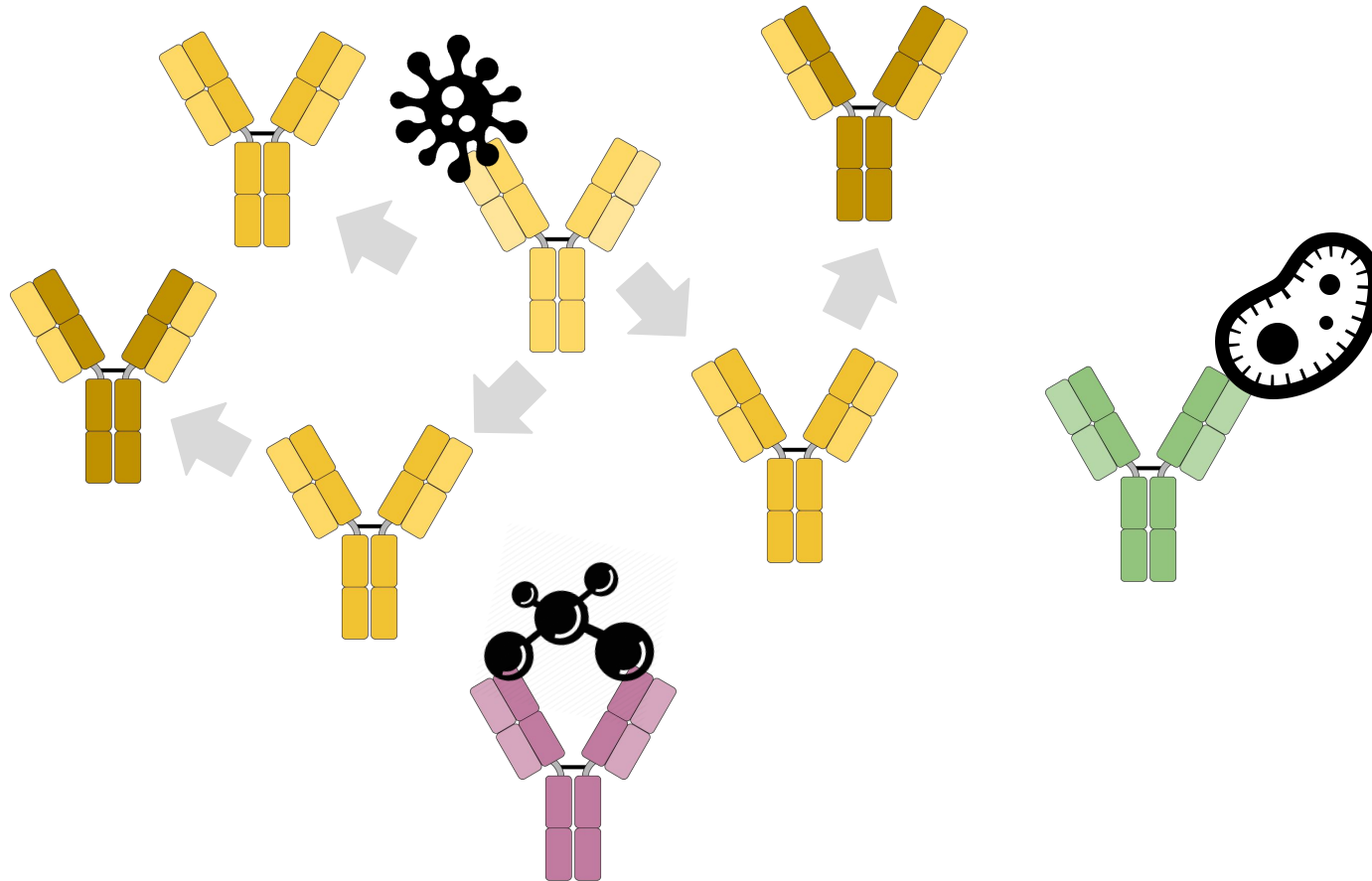


Immune system mutates and amplifies a binding antibody

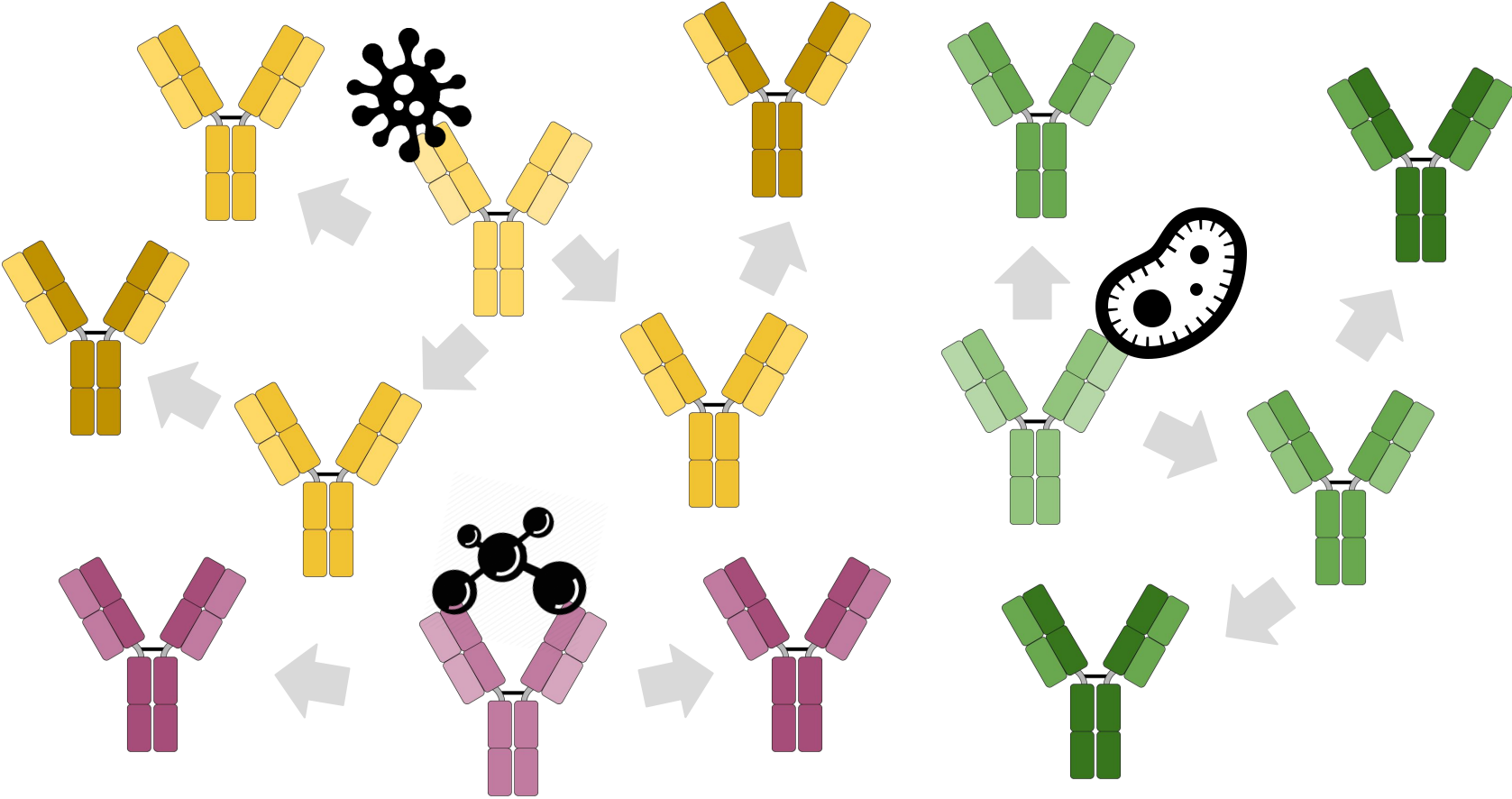


Mutation rate in antibody genes is 3-4 order of magnitude higher than in other genome

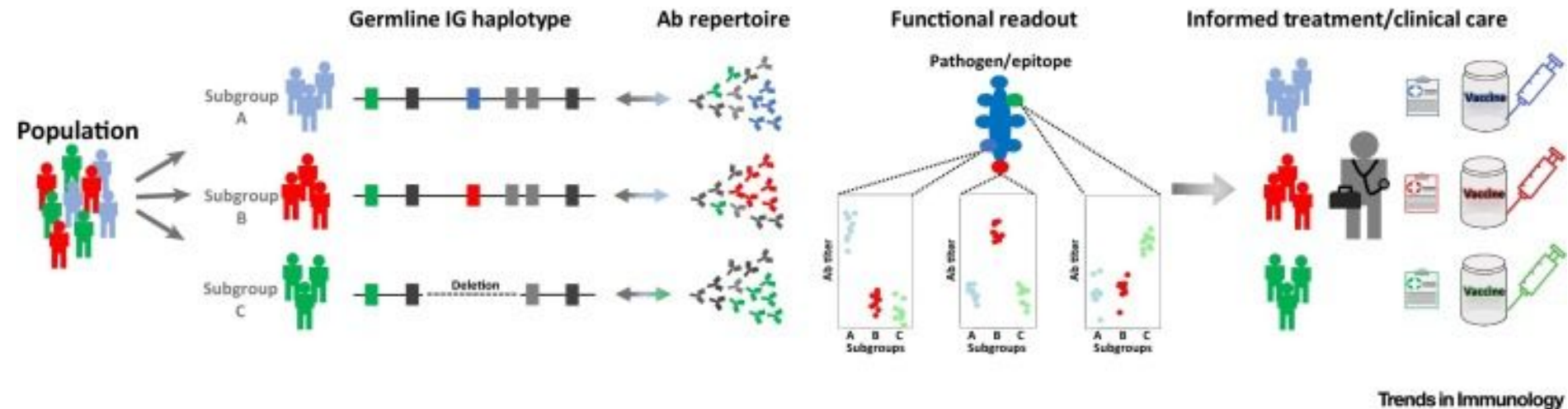
One antibody = one antigen



Antibody repertoire is a set of clonal lineages

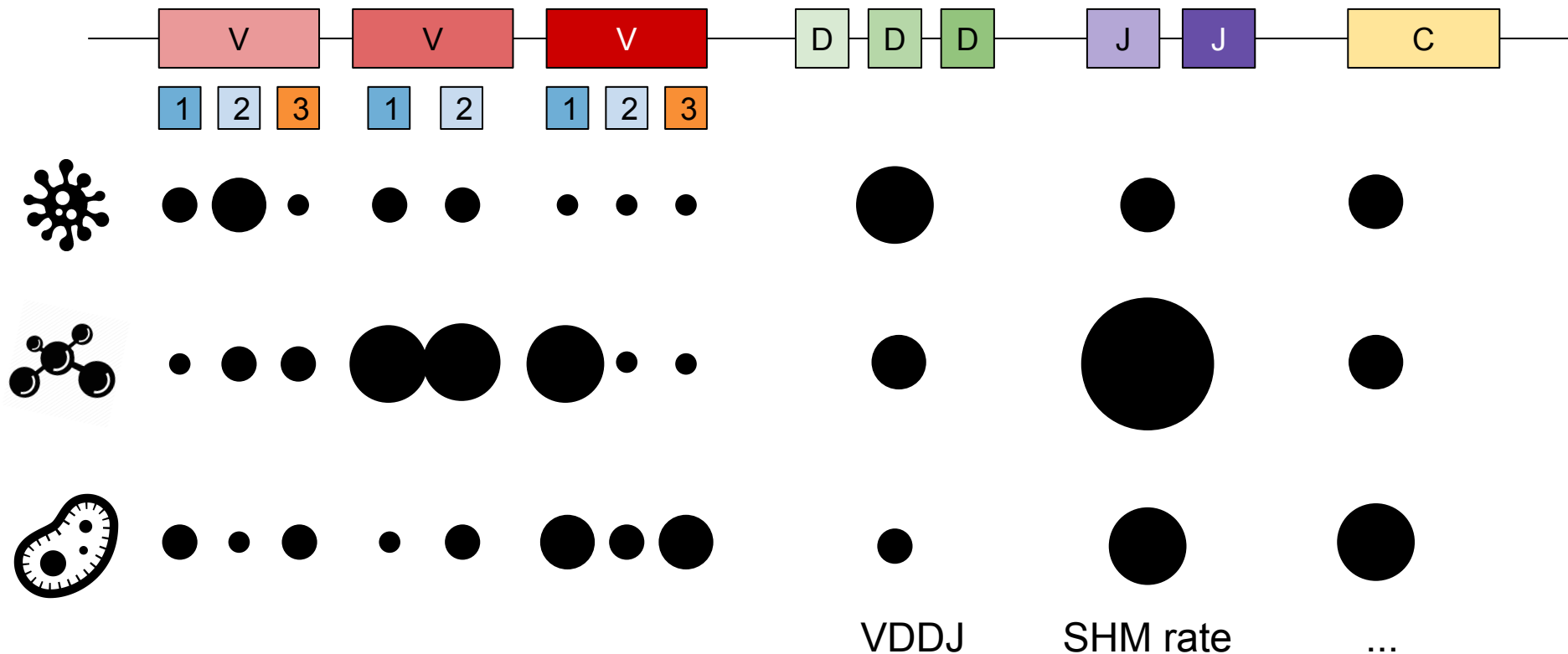


Immunogenomics approach to vaccine design



Watson, Glanville, Marasco, *Trends in Immunol*, 2017

Data science approach to predicting the efficiency of antibody response



Acknowledges



Pavel Pevzner
UCSD

Harvard Medical School
Wayne Marasco
Hanzhong Ke



Corey Watson
U of Louisville

Yale University
Steven Kleinstein
Nima Nouri

UC San Diego

Data Science postdoctoral
fellowship



Intersect fellowship for
computational immunologists

U of Louisville
William Gibson
Justin Kos
Oscar Rodriguez
David Tieri
Jun Yan

U of New South Wales
Andrew Collins
Katherine Jackson

UCSD
Massimo Franceschetti
Siavash Mirarab
Ramesh Rao
Andrey Bzikadze
Vinnu Bhardwaj
Chao Zhang

USDA
Tim Smith
Sung Bong Shin

Iowa State U
James Reecy
Luke Kramer

Scripps Institute
Raiees Andrabi
Vaughn Smider

Smithsonian Conservation Biology Institute
Klaus-Peter Koepfli

Institute for Bioorganic Chemistry
Ivan Zvyagin
Artem Mikelov
Mikhail Shugay