



High-throughput sequencing in immune repertoire analysis

Nika Abdollahi

Supervised by :
Prof. Martin Weigt
Prof. Frédéric Davi
Dr. Juliana Bernardes

13 Mars 2020

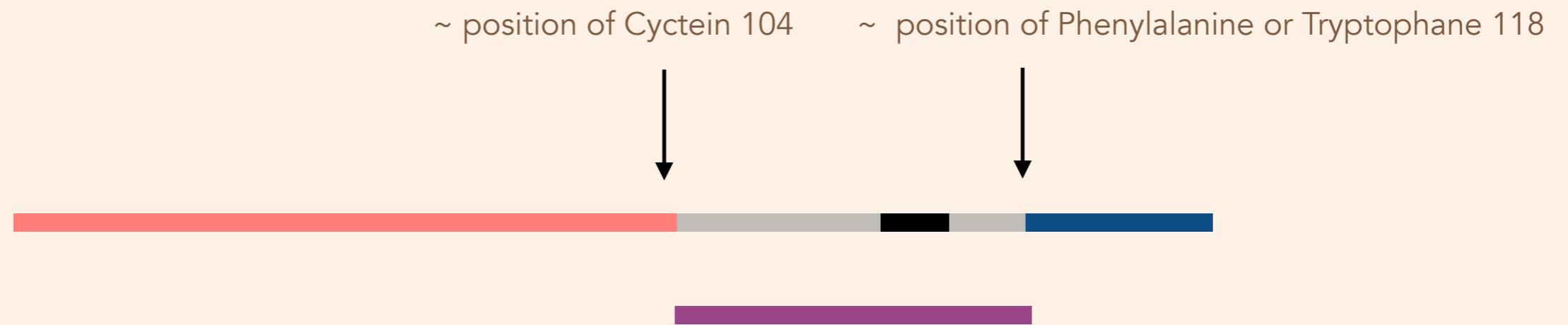
Plan

- Identifying clonal-related sequences in BCR repertoires without VDJ annotations
- Clustering T-cell beta-chain sequence receptor according to it's epitope preference
- The BCR intra-clonal diversity
- Benchmarking RepSeq analysis tools

Identifying clonal-related sequences in BCR repertoires without VDJ annotations

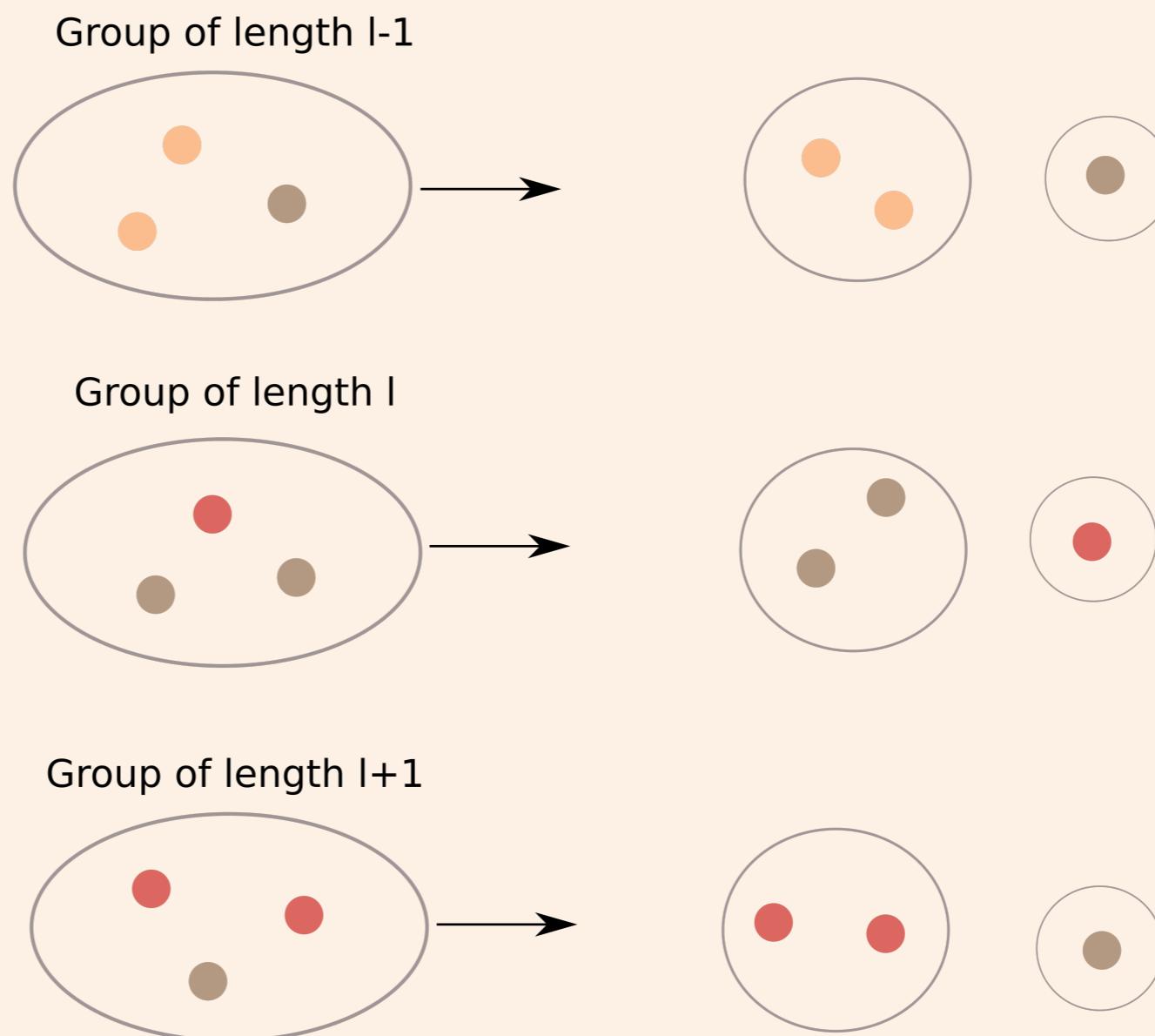
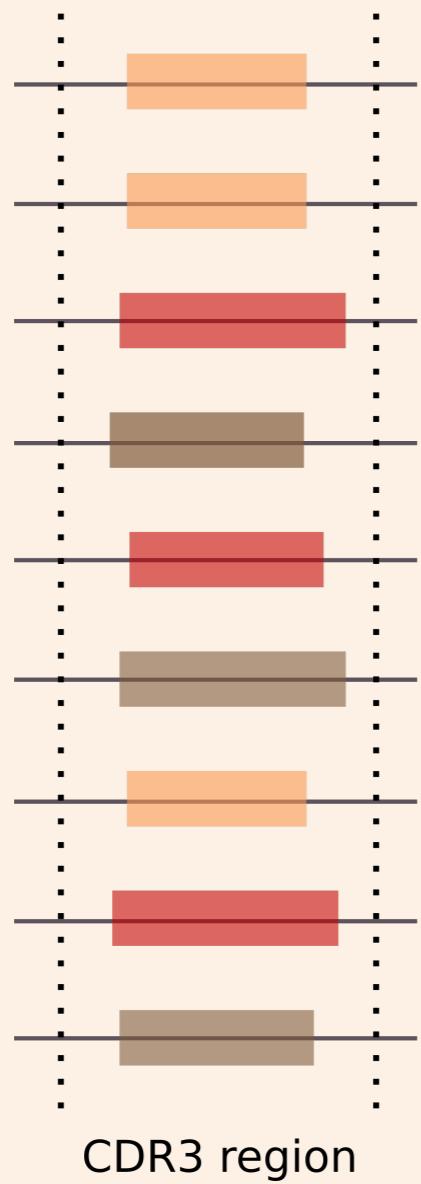
Steps of the algorithm

1- Junction extraction (optional)



Steps of the algorithm

2- Greedy clustering



Identifying clonal-related sequences in BCR repertoires without VDJ annotations

Improvements

Change the method's name from Mixclus to FaIR

Adapte the junction extraction part of the vidjil algorithm.

Design different methods for evaluating the performance of FaIR and its competitors.

Add few simple visualization tools for clinical use.

Extensively tested FaIR on simulated and experimental B-cell repertoires of both healthy individuals and patients with chronic lymphocytic leukemia.

Identifying clonal-related sequences in BCR repertoires without VDJ annotations

Improvements

Change the method's name from Mixclus to FaIR

Adapte the junction extraction part of the vidjil algorithm.

Design different methods for evaluating the performance of FaIR and its competitors.

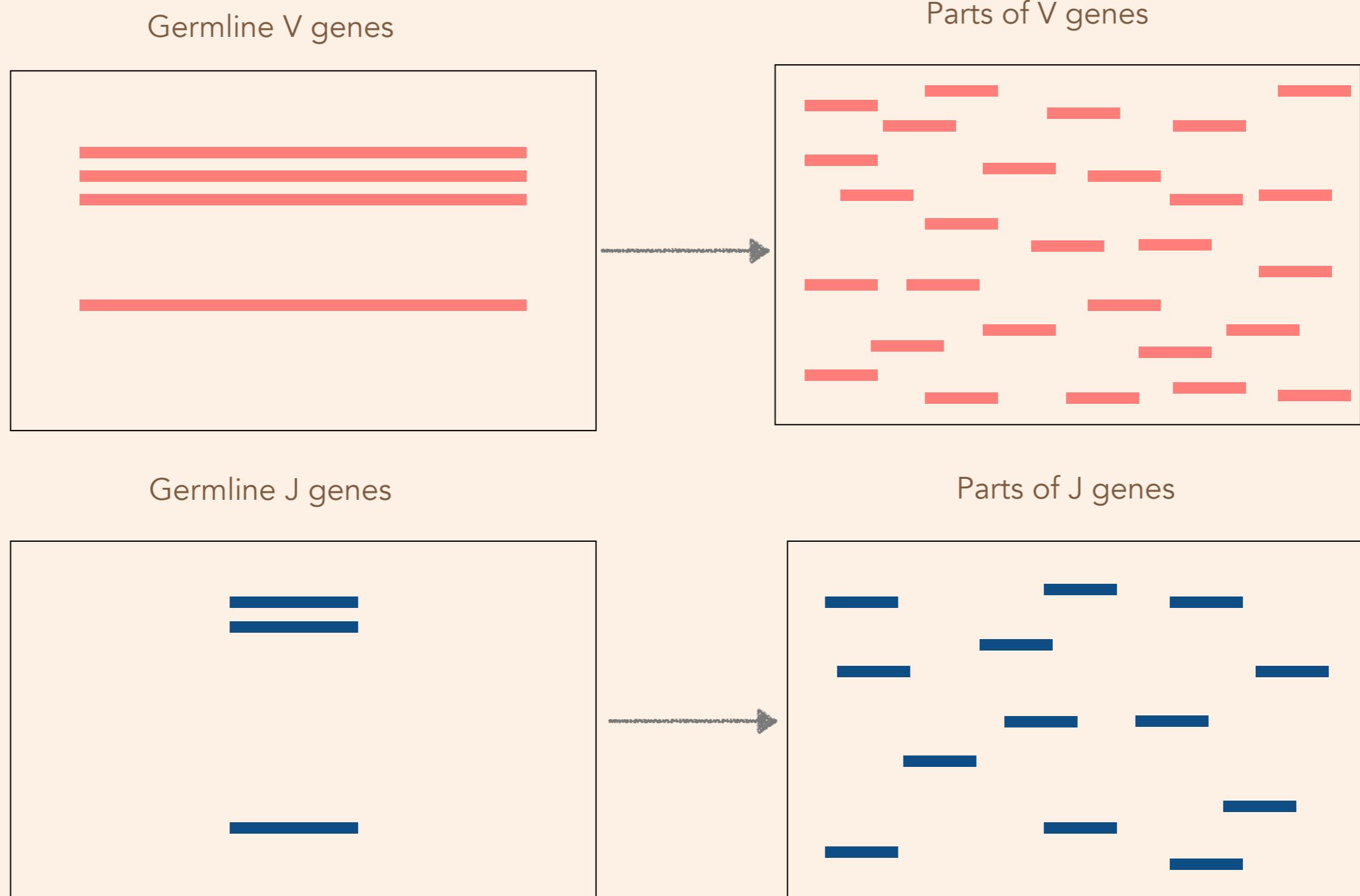
Add few simple visualization tools for clinical use.

Extensively tested FaIR on simulated and experimental B-cell repertoires of both healthy individuals and patients with chronic lymphocytic leukemia.

Adapte the junction extraction part of the vidjil algorithm

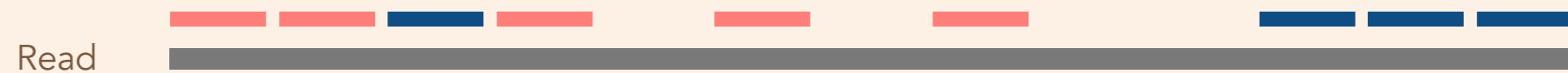
Vidjil_algo first phase

Indexing the germlines data bases

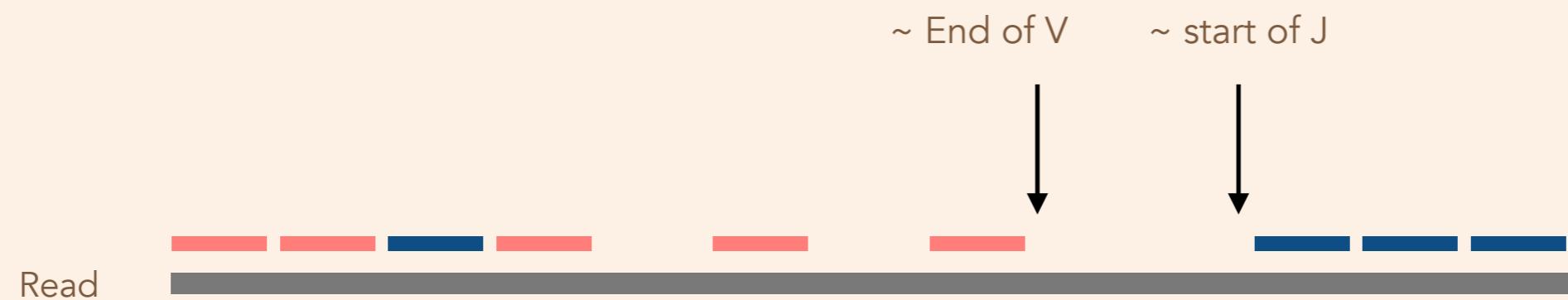


Vidjil_algo first phase

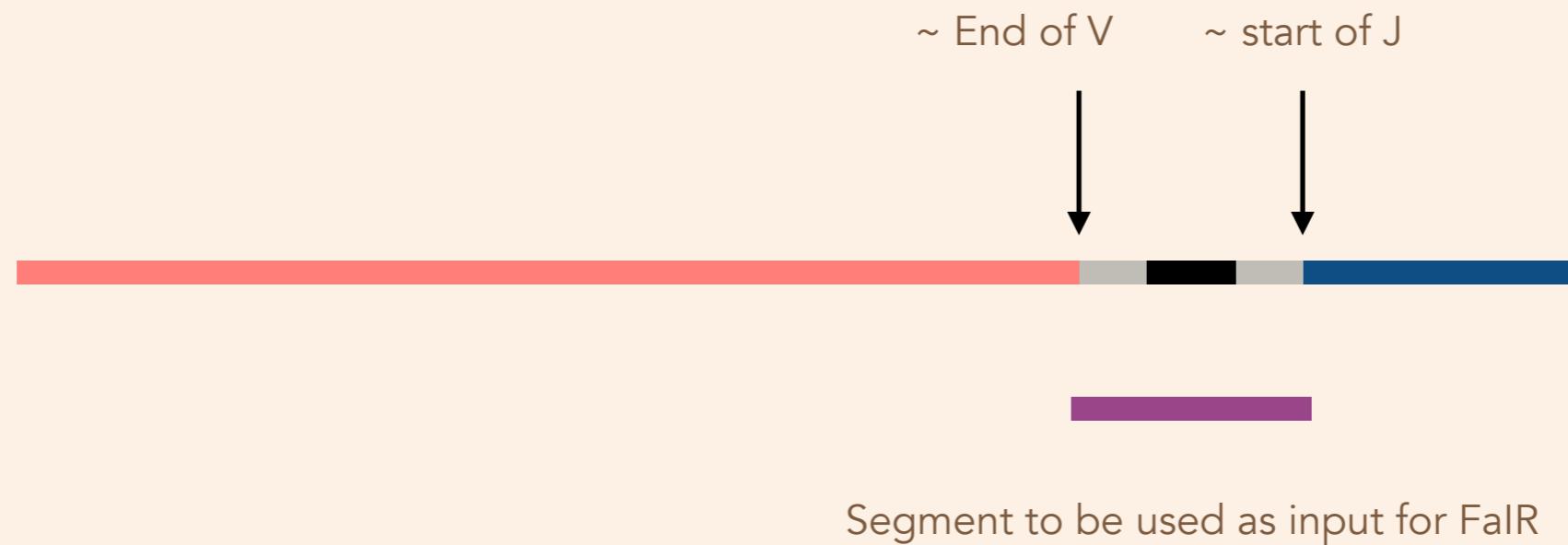
Matching the « k mers » on reads



Vidjil_algo first phase



Extracting the junction sequence



In general, It covers the two N regions and the D gene

It's shorter than CDR3

Not enough information for clustering

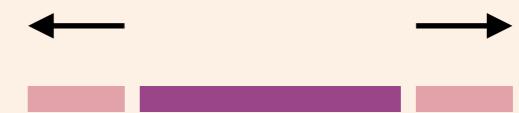
Extracting the junction sequence

Tradeoff between the sequence length ('effective information for clustering') and the computational time



Adjustment

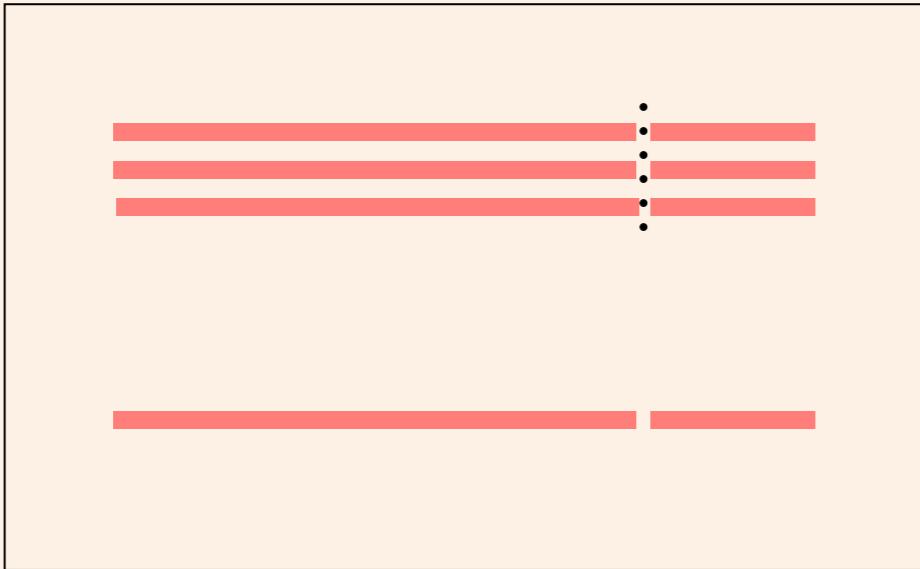
(to increase the chance of covering CDR3 region)



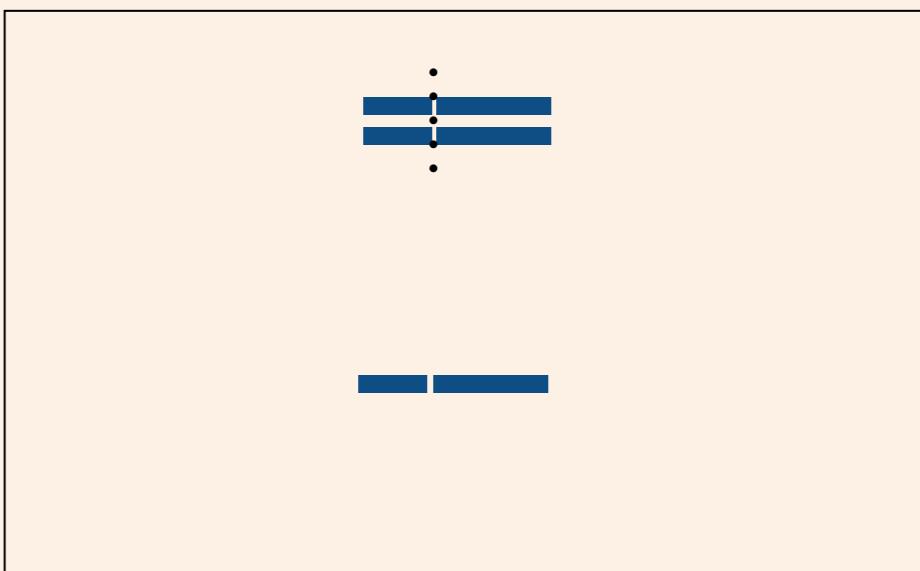
Vidjil_algo first phase

Trim the sequences based on the CDR3 anchors

Trim germline V genes after Cyctein 104

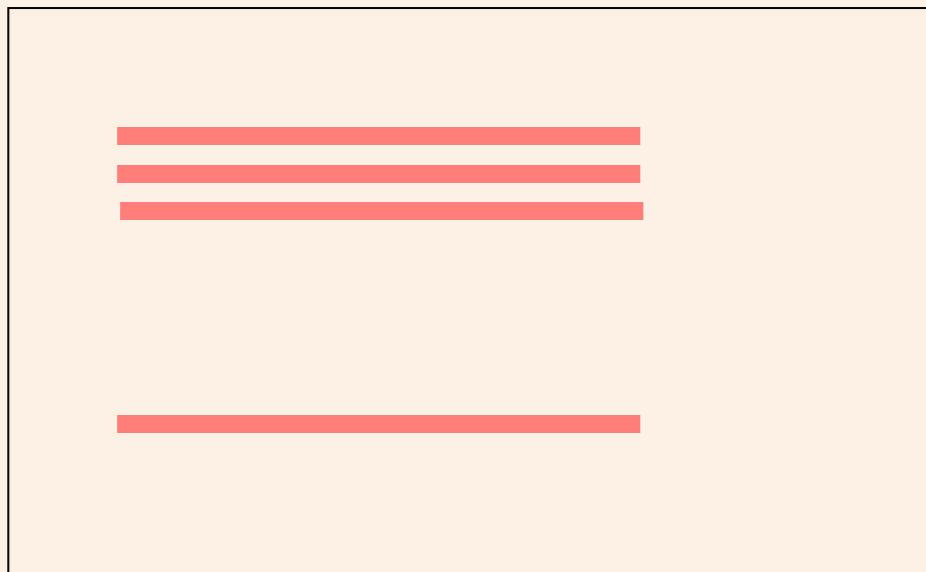


Trim germline J genes before Phenylalanine or Tryptophane 118

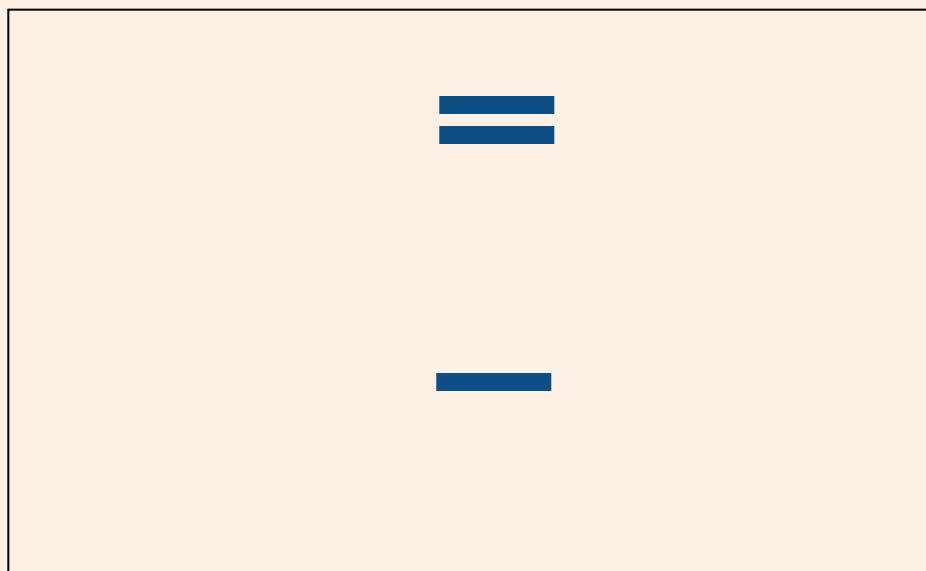


Vidjil_algo first phase

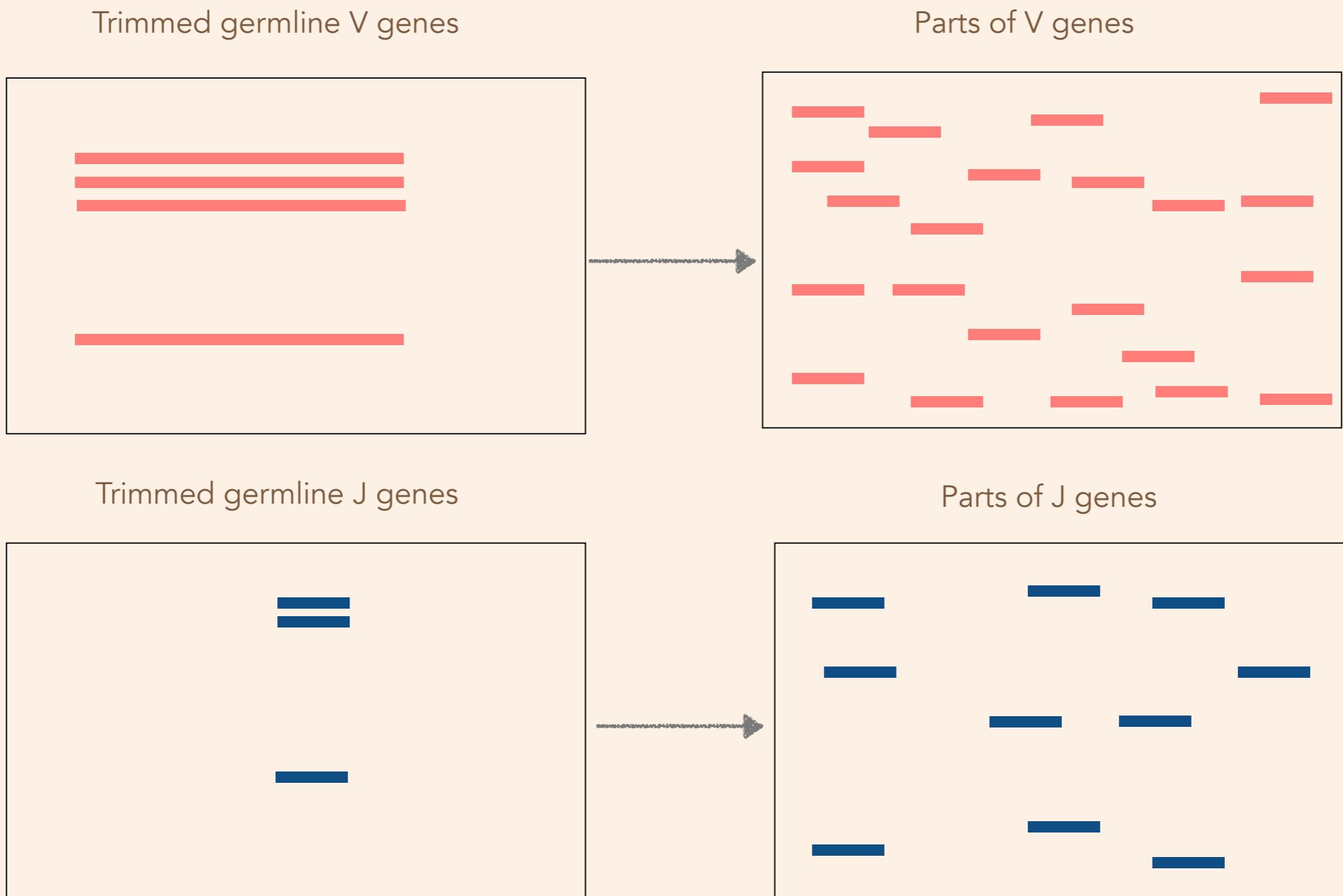
Trimmed germline V genes



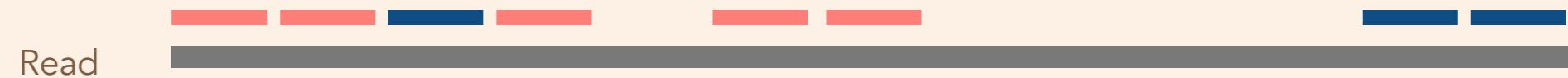
Trimmed germline J genes



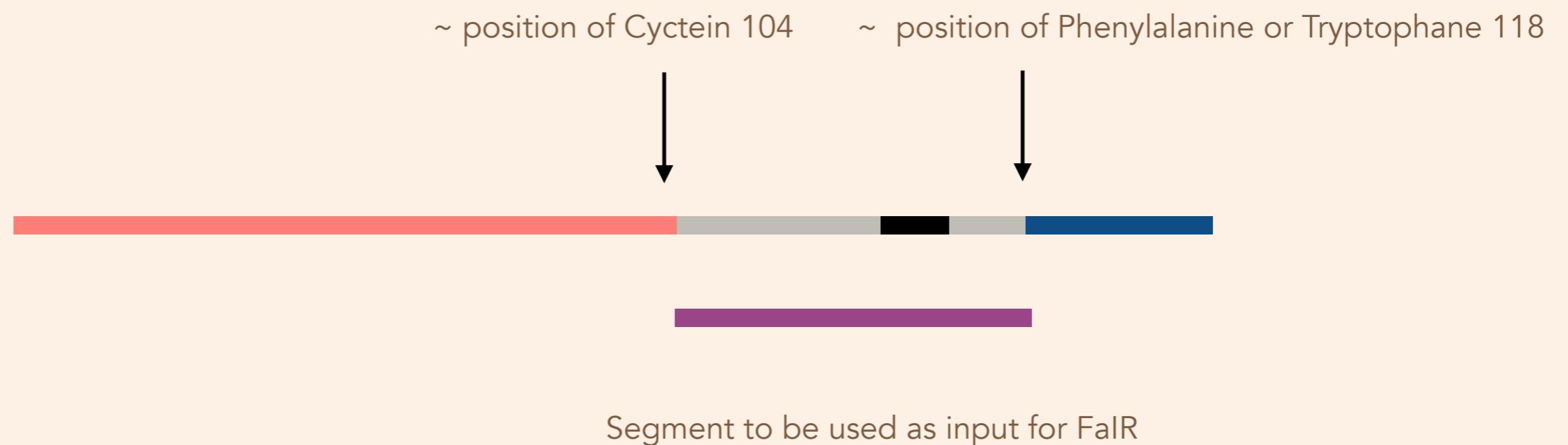
Vidjil_algo first phase



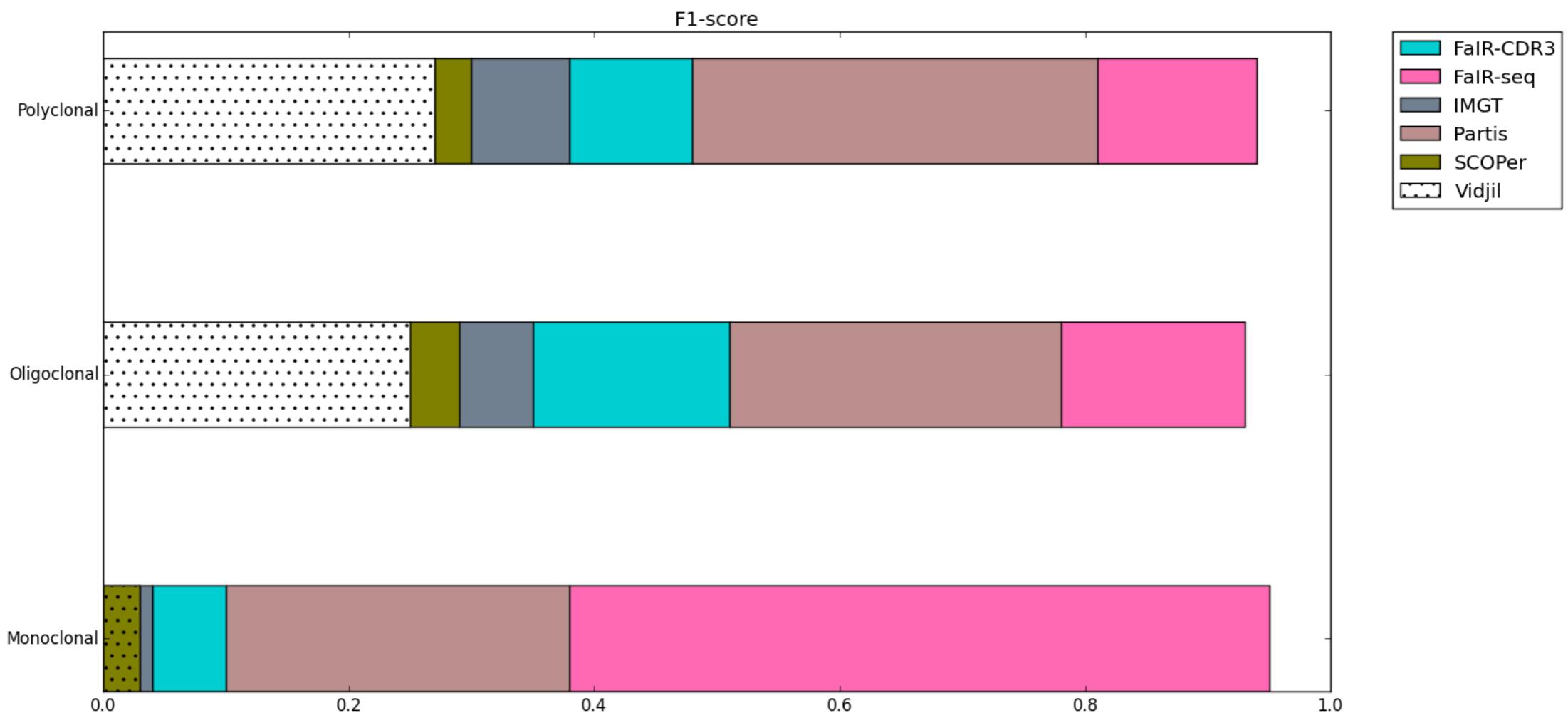
Vidjil_algo first phase

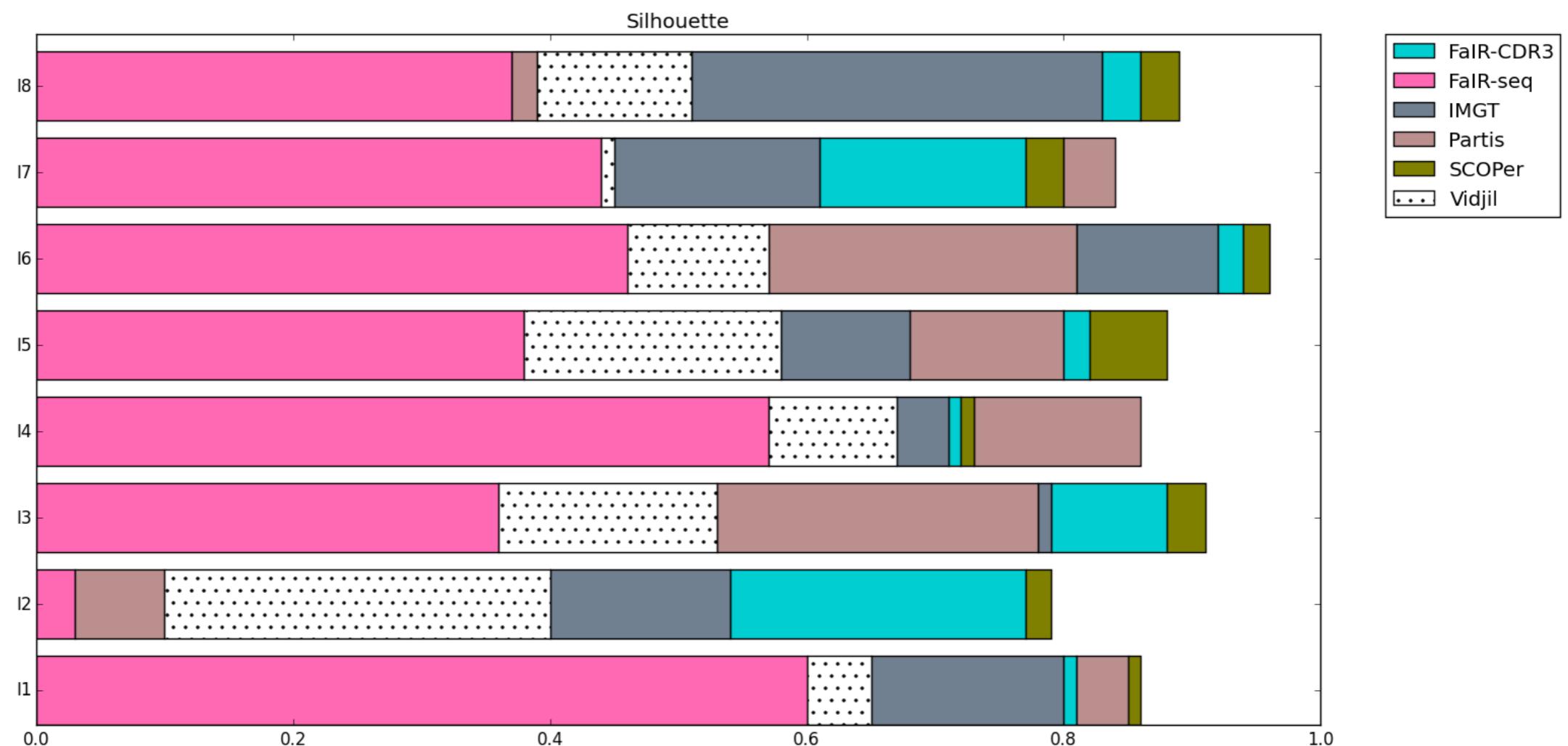


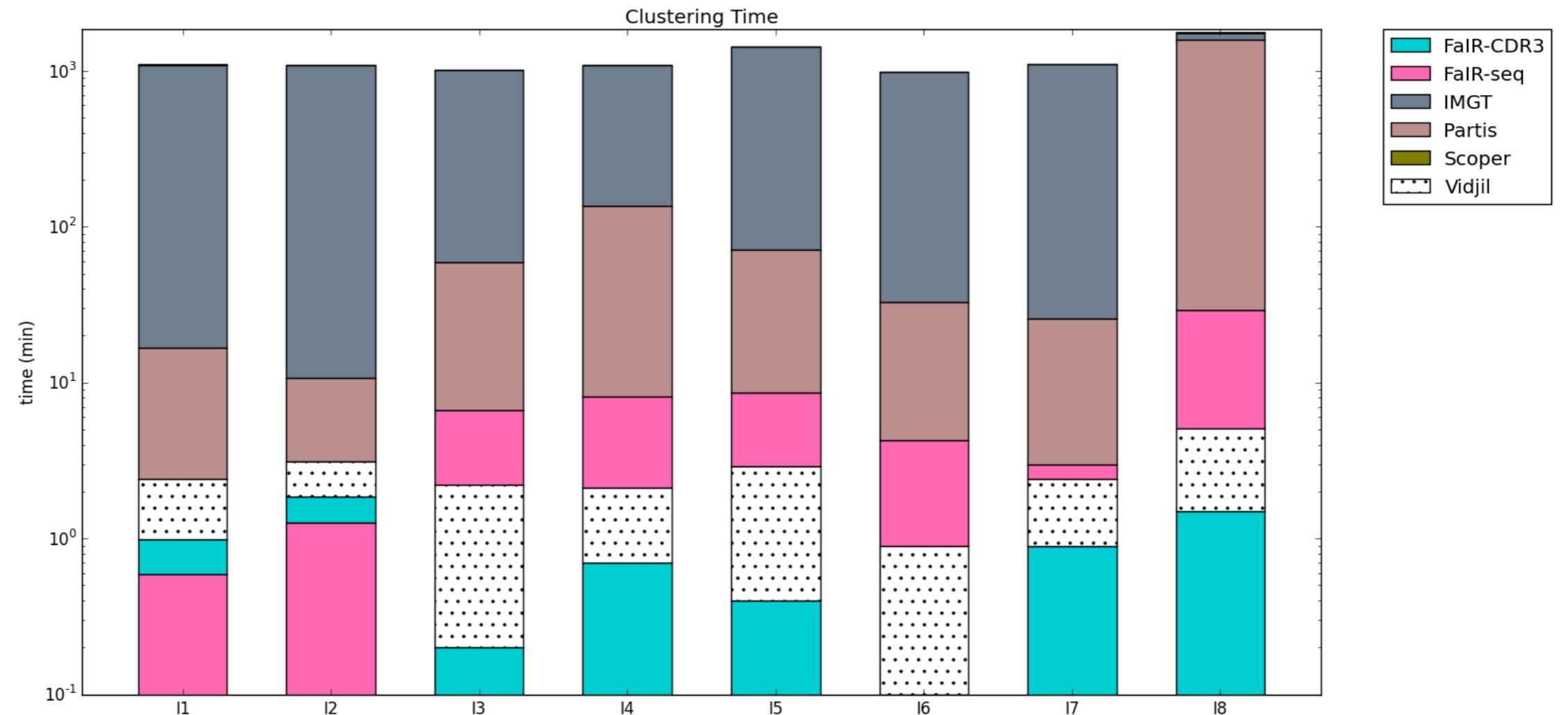
Vidjil_algo first phase



Design of different methods for evaluating the performance of FalR and its competitors on artificial and experimental repertoires







Silhouette

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

i : a datum

a(i) : the average distance between i and all other data within the same cluster

b(i) : the smallest average distance of to all points in any other cluster, of which i is not a member

Silhouette

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

i : a datum

a(i) : the average distance between i and all other data within the same cluster

b(i) : the smallest average distance of to all points in any other cluster, of which i is not a member

Distance 1 : Levenstein distance between 2 nucleotide sequences

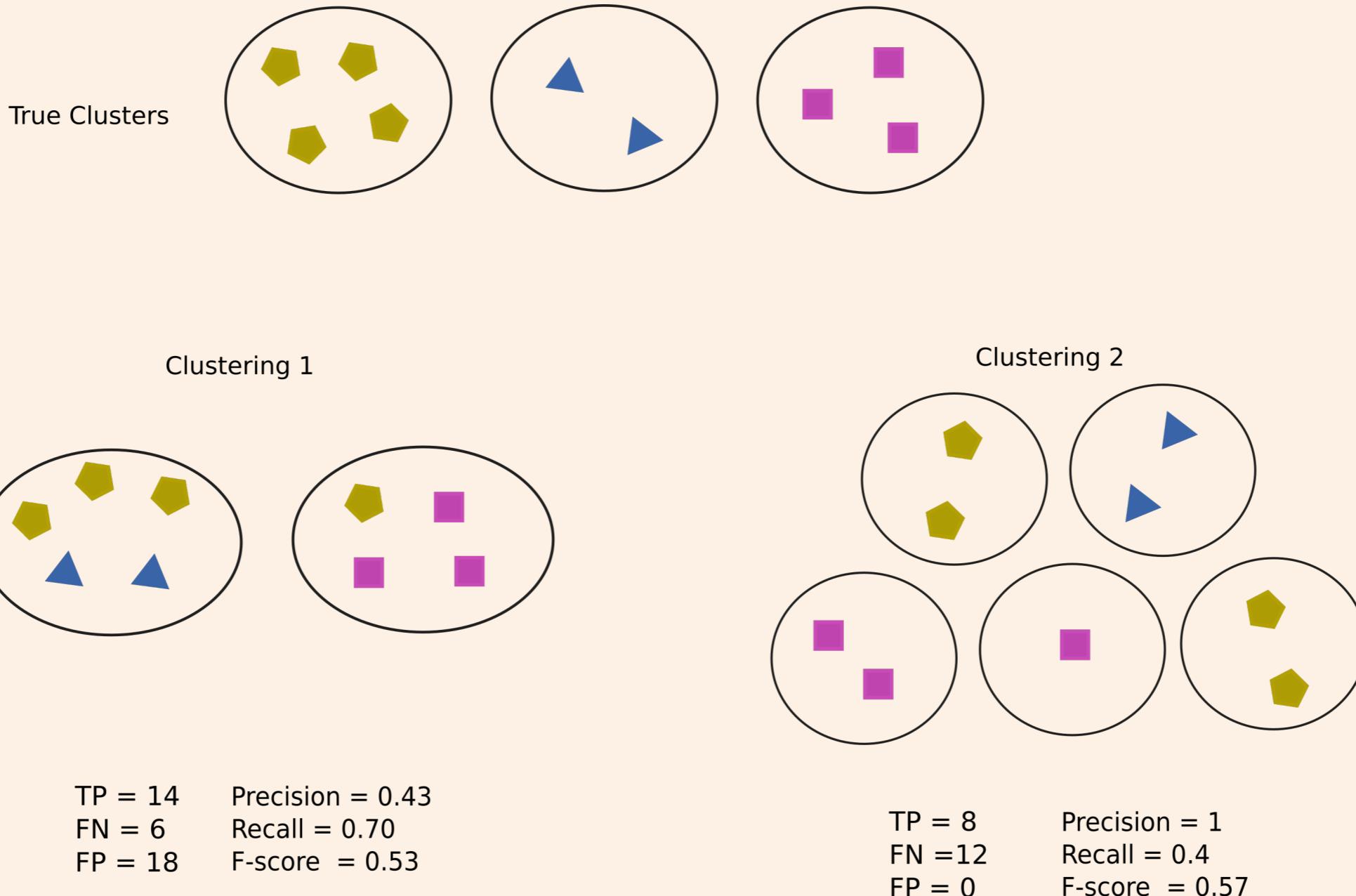
Distance 1 : Levenstein distance between 2 CDR3 sequences + V gene identity + J gene identity

Silhouette

Add The results

F-score

Sequence based

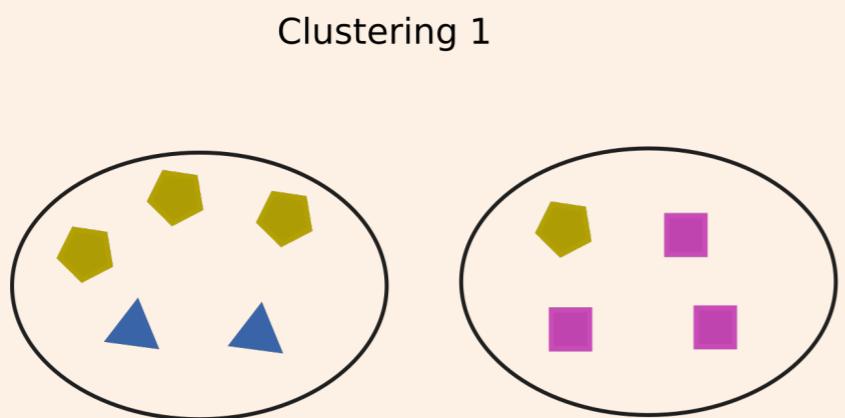
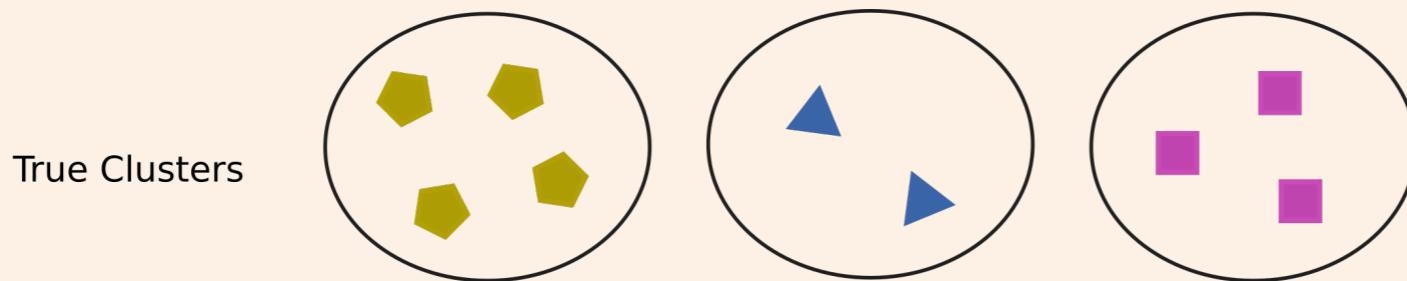


TP = 14 Precision = 0.43
FN = 6 Recall = 0.70
FP = 18 F-score = 0.53

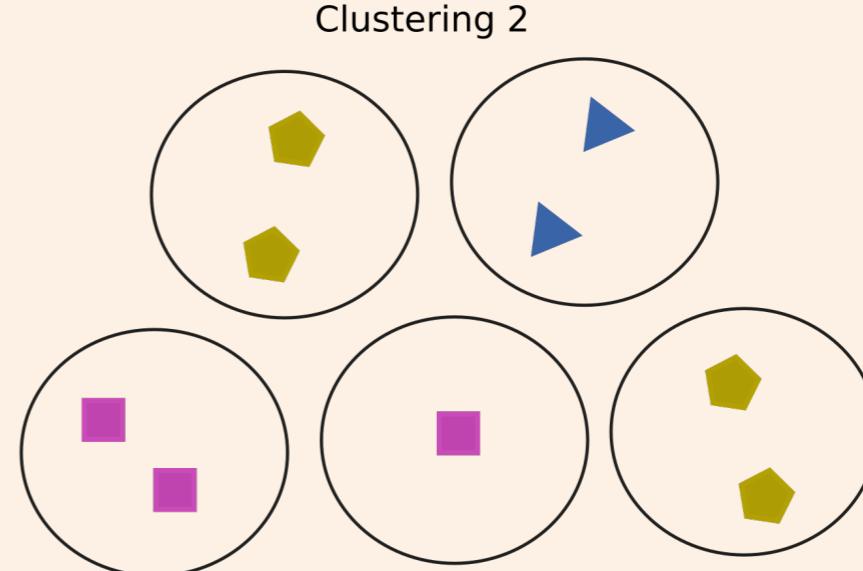
TP = 8 Precision = 1
FN = 12 Recall = 0.4
FP = 0 F-score = 0.57

F-score

Cluster based



TP = 6 Precision = 0.66
FN = 1 Recall = 0.85
FP = 3 F-score = 0.74



TP = 9 Precision = 1
FN = 7 Recall = 0.56
FP = 0 F-score = 0.71

F-score

Add The results

**Clustering T-cell beta-chain
sequence receptor according to
it's epitope preference**

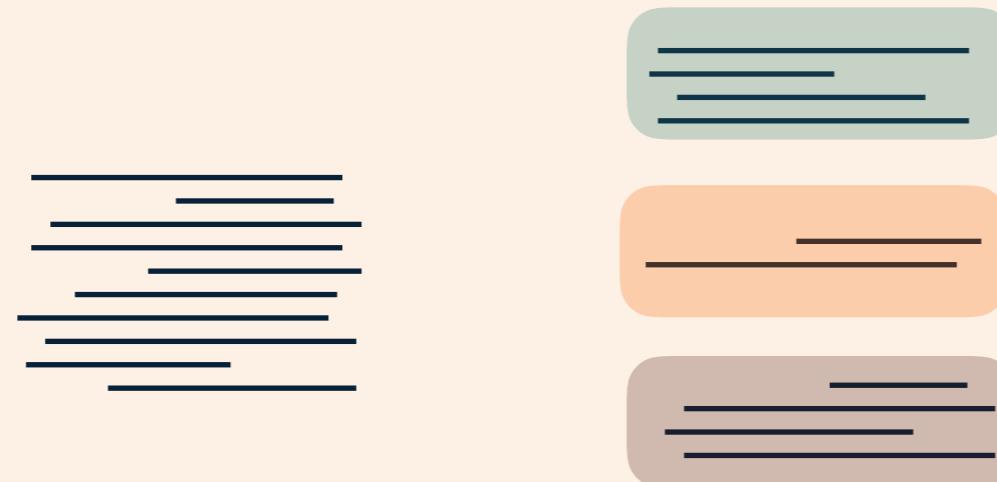
On the viability of unsupervised T-cell receptor sequence clustering for epitope preference

Pieter Meysman, Nicolas De Neuter, Sofie Gielis, Danh Bui Thi, Benson Ogunjimi, Kris Laukens

Structural distinct TCR groups ==> target identical epitopes

VDJdb

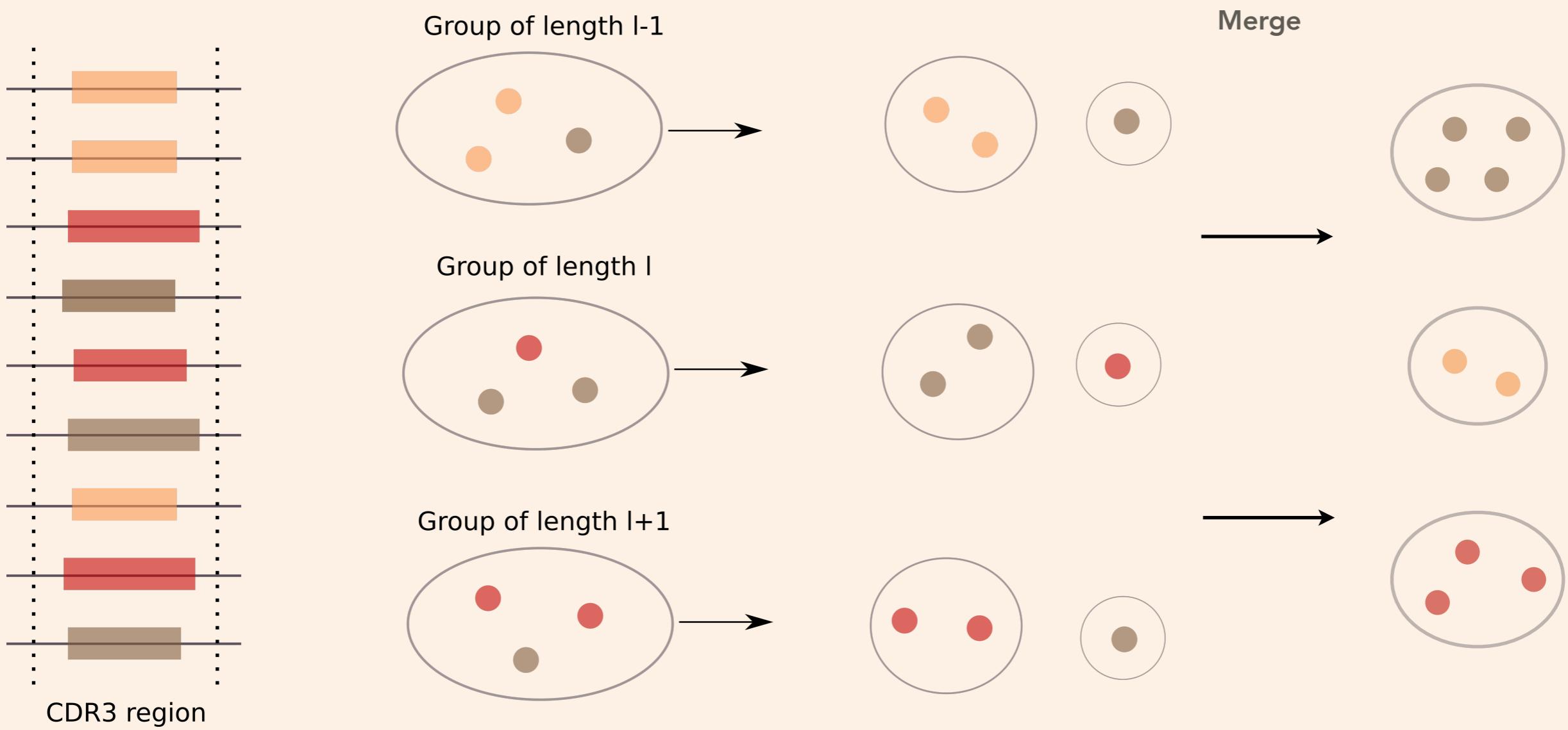
Cluster : CDR3 sequences associated to the same epitope-MHC
(#seq =1021, # true cluster= 78, #singleton = 11)



Challenge

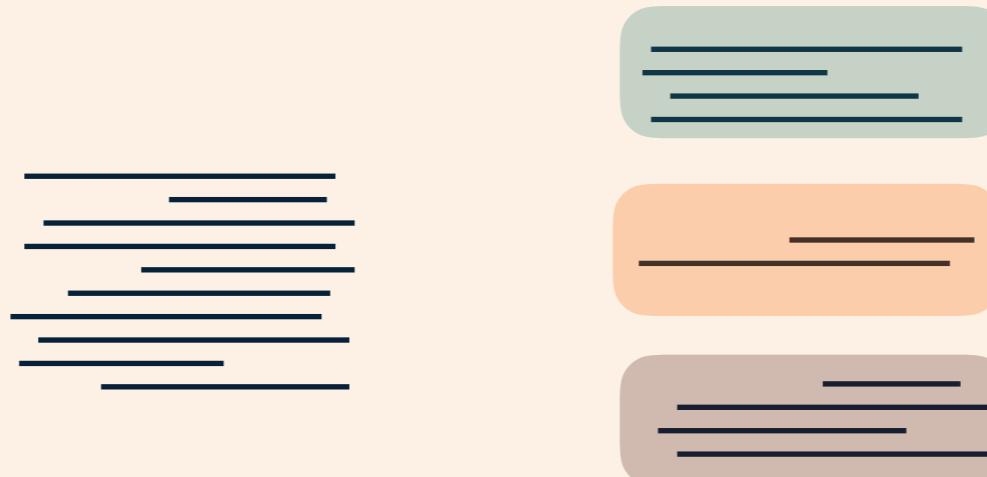


Adapted algorithm



Performance

Cluster : CDR3 sequences associated with the same epitope-MHC
 (#seq =1021, # true cluster= 78, #singleton = 11)



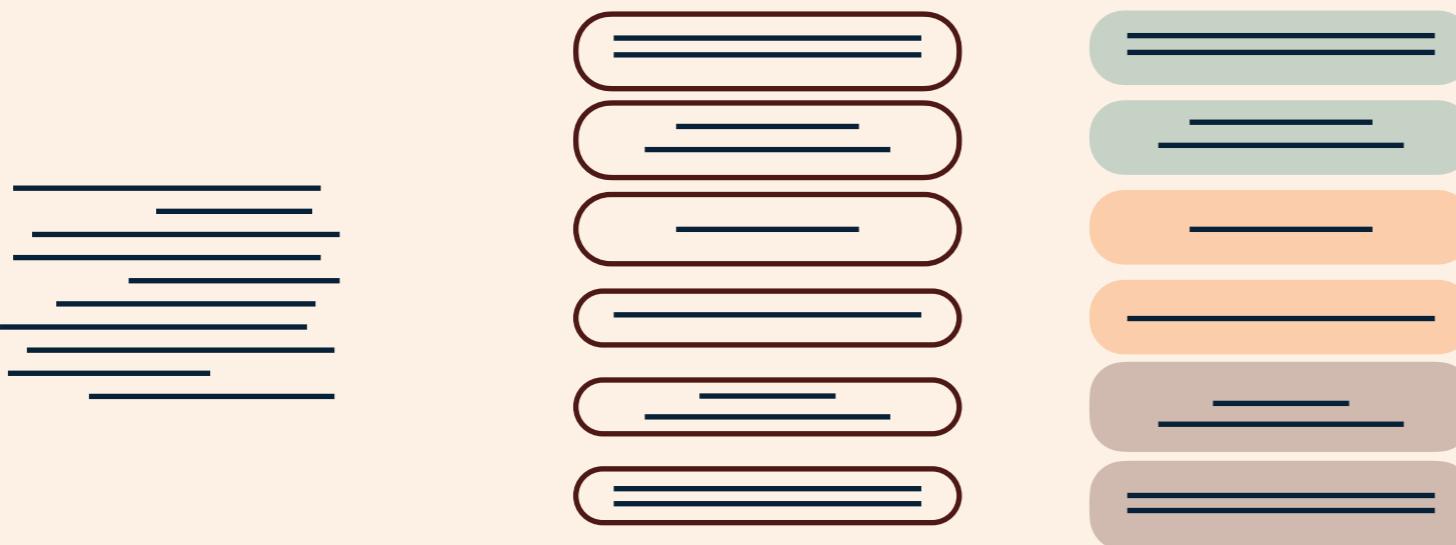
	#cluster	#singleton	Non assigned	Nearest cluster			# non singleton cluster	Pairwise sequences		
				Recall	precision	F-Score		Retention	Purity (Precision)	Consistency (Recall)
DBSCAN	594	441	441	0,025	0,88	0,049	153	0,568	0,868	0,204
Mixclus	93	9	0	0,06	0,78	0,11	84	0,91	0,76	0,48

VDJdb

Cluster : CDR3 sequences associated to the same epitope-MHC
(#seq =1021, # true cluster= 78, #singleton = 11)

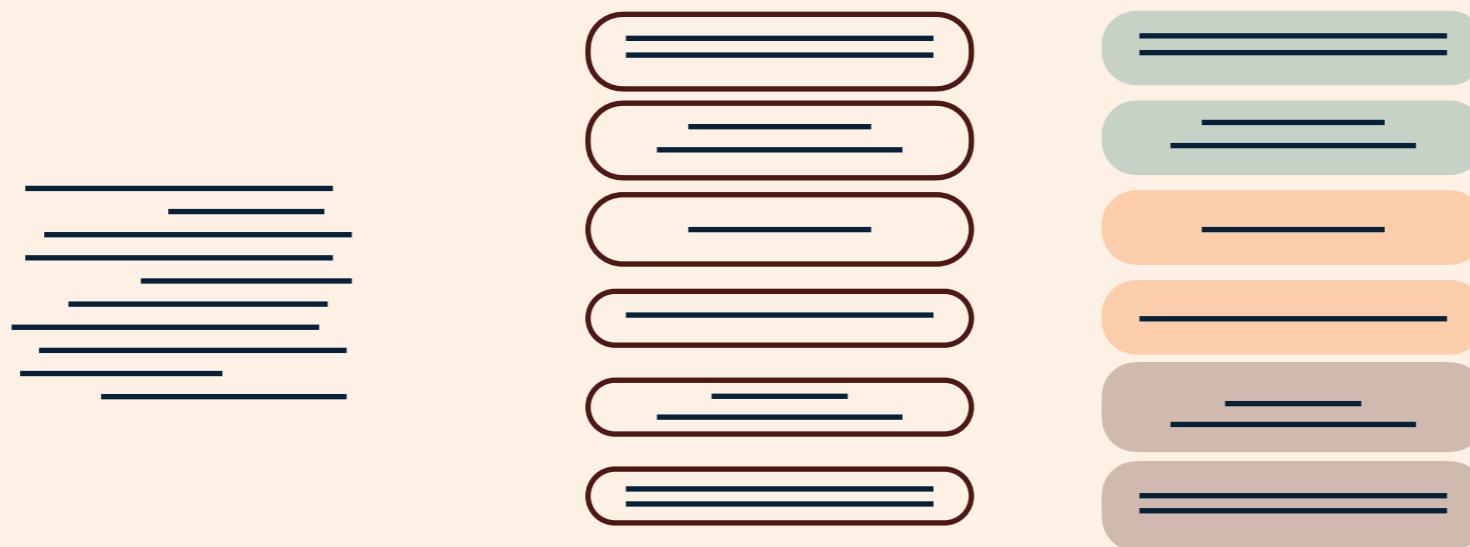


Cluster : CDR3 sequences with the same V and J genes associated to the same epitope-MHC
(#seq =1021, # true cluster= 544, #singleton =243)



Performance

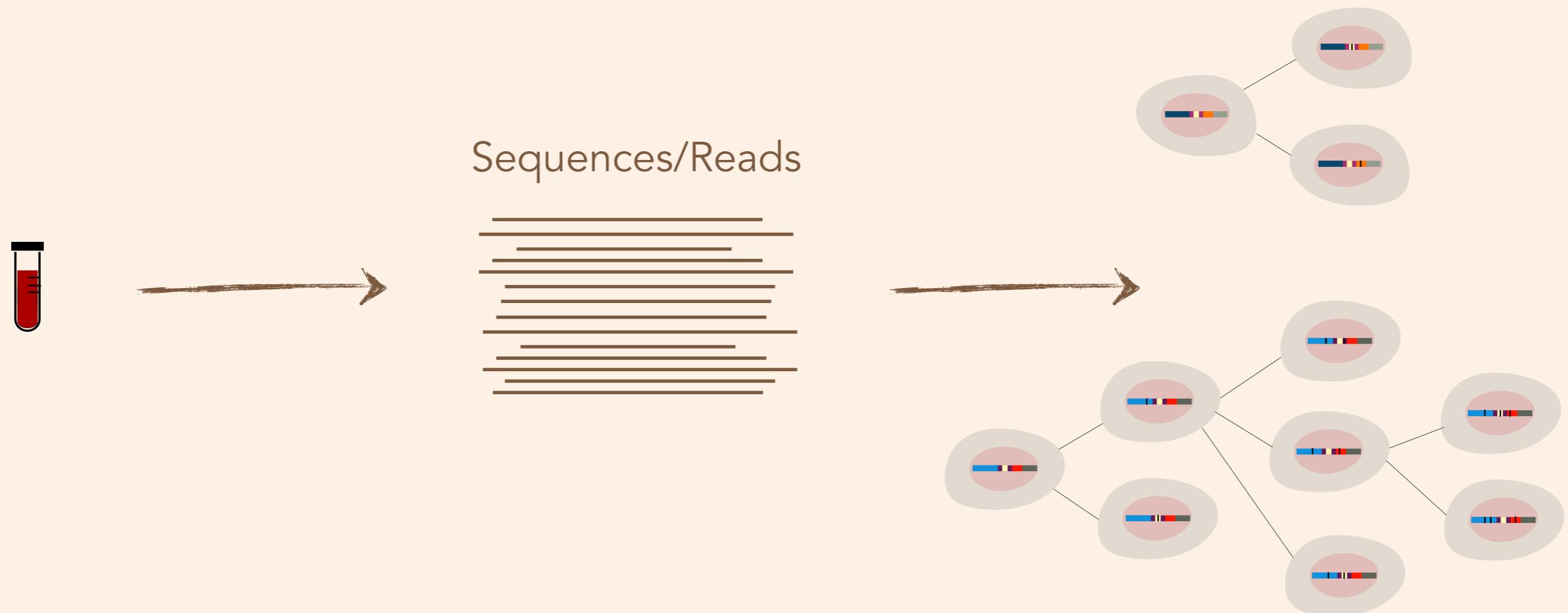
Cluster : CDR3 sequences with the same V and J genes associated to the same epitope-MHC
 (#seq = 1021, # true cluster= 544, #singleton = 367)



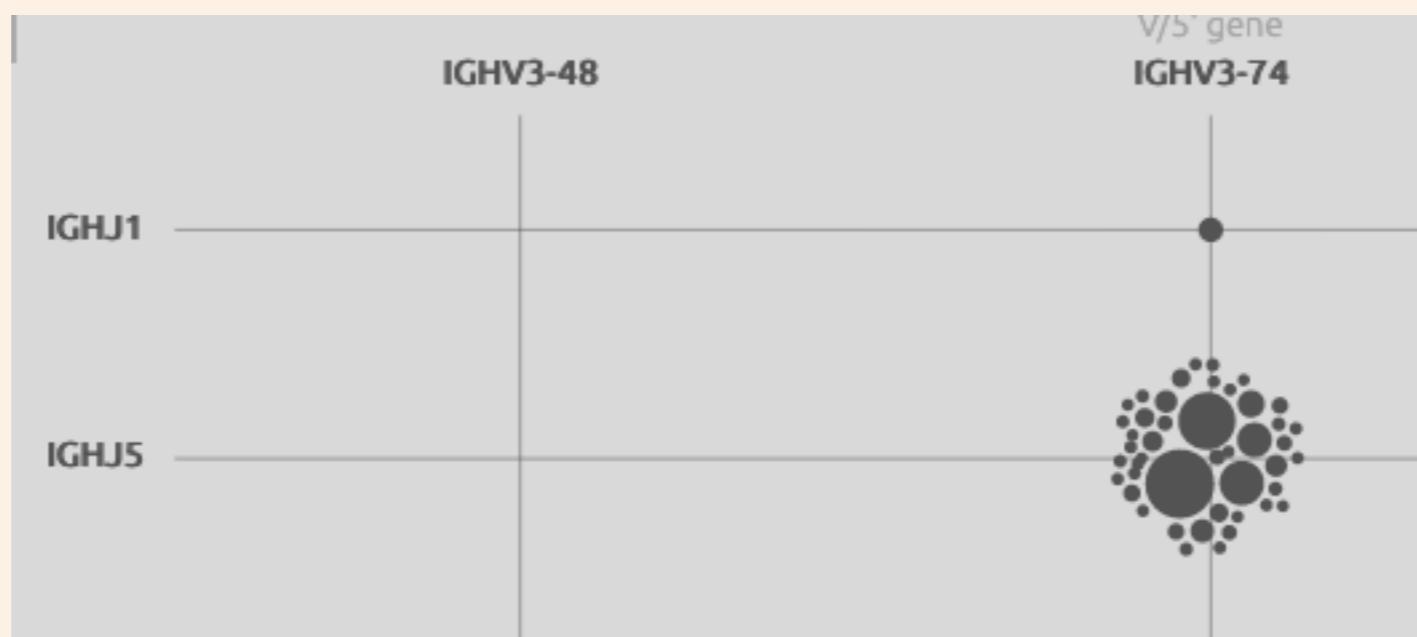
	#cluster	#singleton	Non assigned	Nearest cluster			# non singleton cluster	Pairwise sequences		
				Recall	precision	F-Score		Retention	Purity (Precision)	Consistency (Recall)
DBSCAN	632	470	470	0.43	0,90	0,55	162	0,539	0,94	0,55
Mixclus	112	33	0	0.78	0.66	0.71	112	0,95	0,60	0,7

The BCR intra-clonal diversity

Objectif



Intra-clonal diversity example

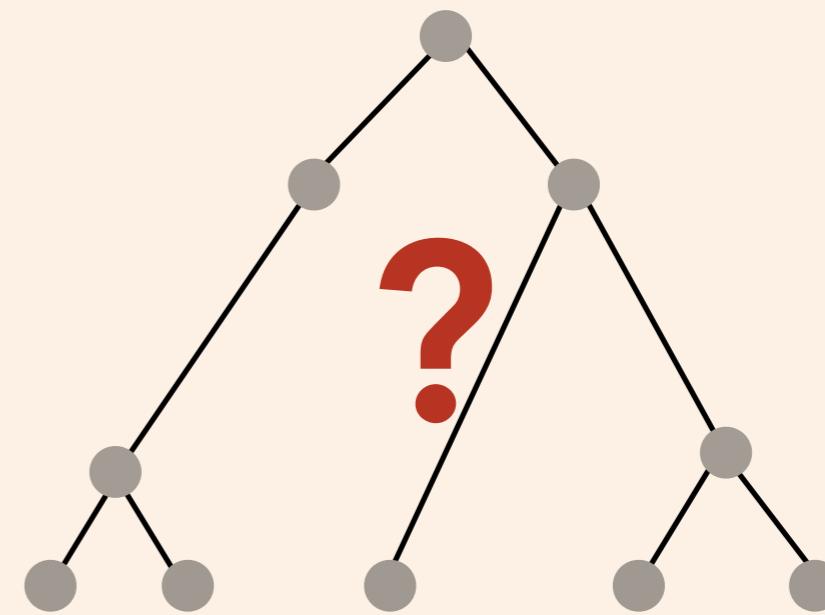
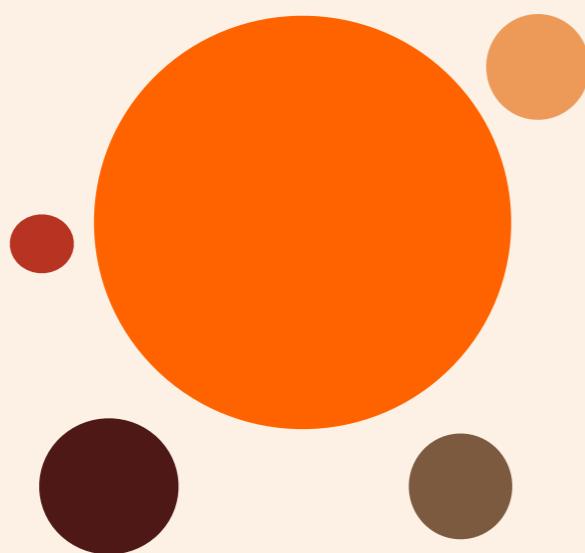


Intra-clonal diversity example

One clone :

2317 sequences

554 unique sequences

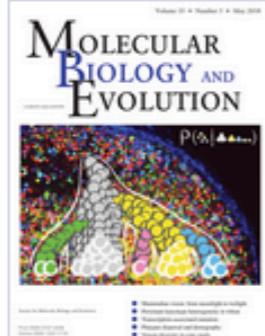


GC tree

Phylogenetic inference of genotype-collapsed tree

MOLECULAR BIOLOGY AND EVOLUTION

Issues Advance articles Submit ▾ Purchase Alerts About ▾ All Mo



Volume 35, Issue 5, May 2018

Volume 35, Issue 5
May 2018

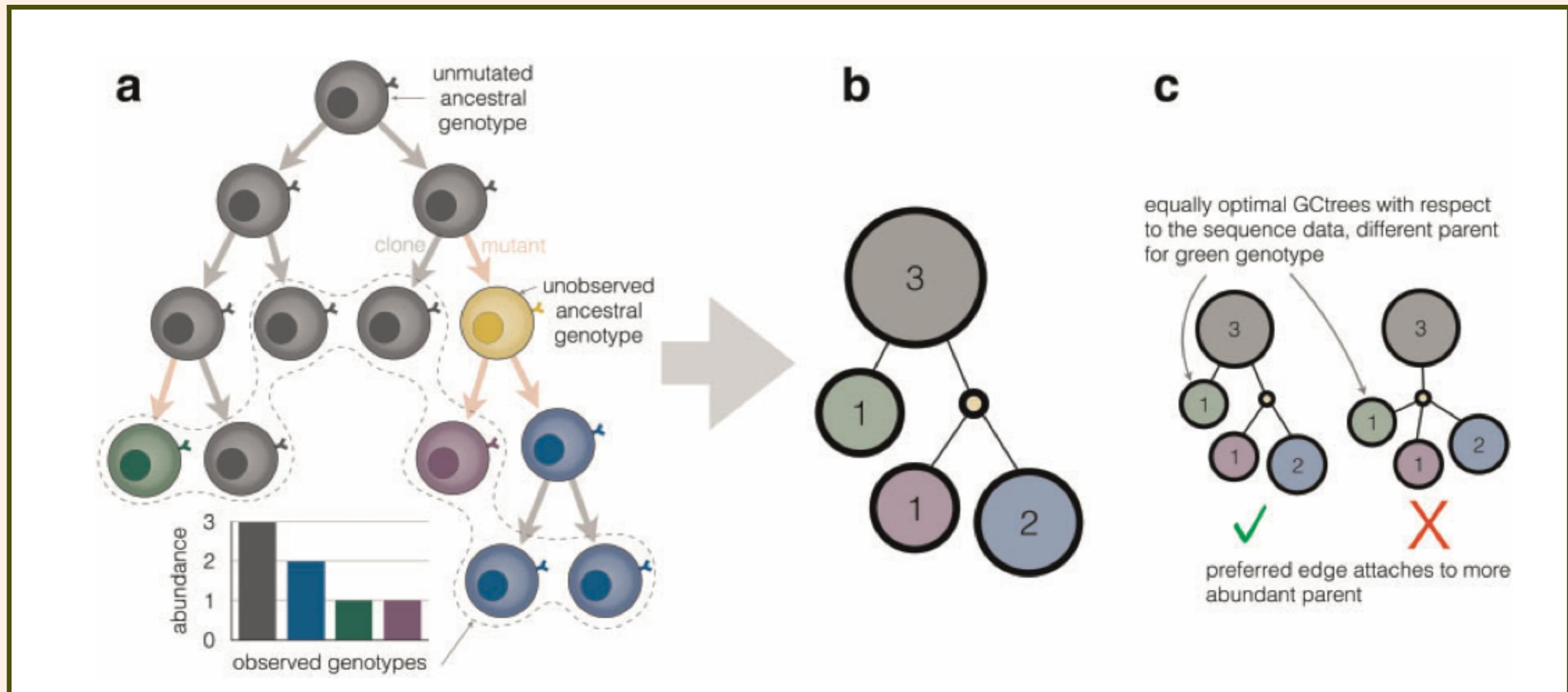
Using Genotype Abundance to Improve Phylogenetic Inference 

William S DeWitt, III, Luka Mesin, Gabriel D Victora, Vladimir N Minin ,
Frederick A Matsen, IV 

Molecular Biology and Evolution, Volume 35, Issue 5, May 2018, Pages 1253–1265,
<https://doi.org/10.1093/molbev/msy020>

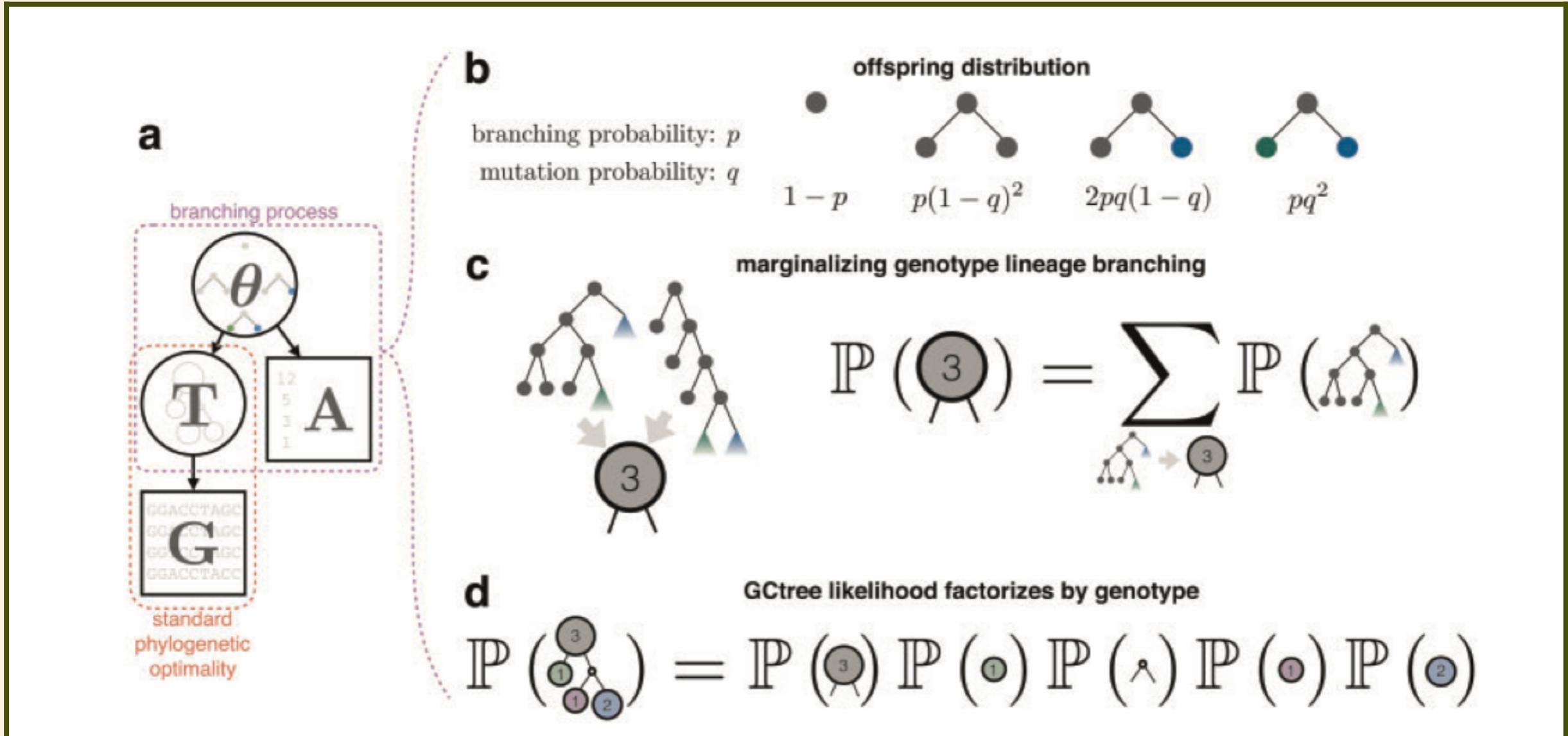
Published: 20 February 2018

GC tree



The more abundant parent is more likely to have generated mutant descendants

GC tree

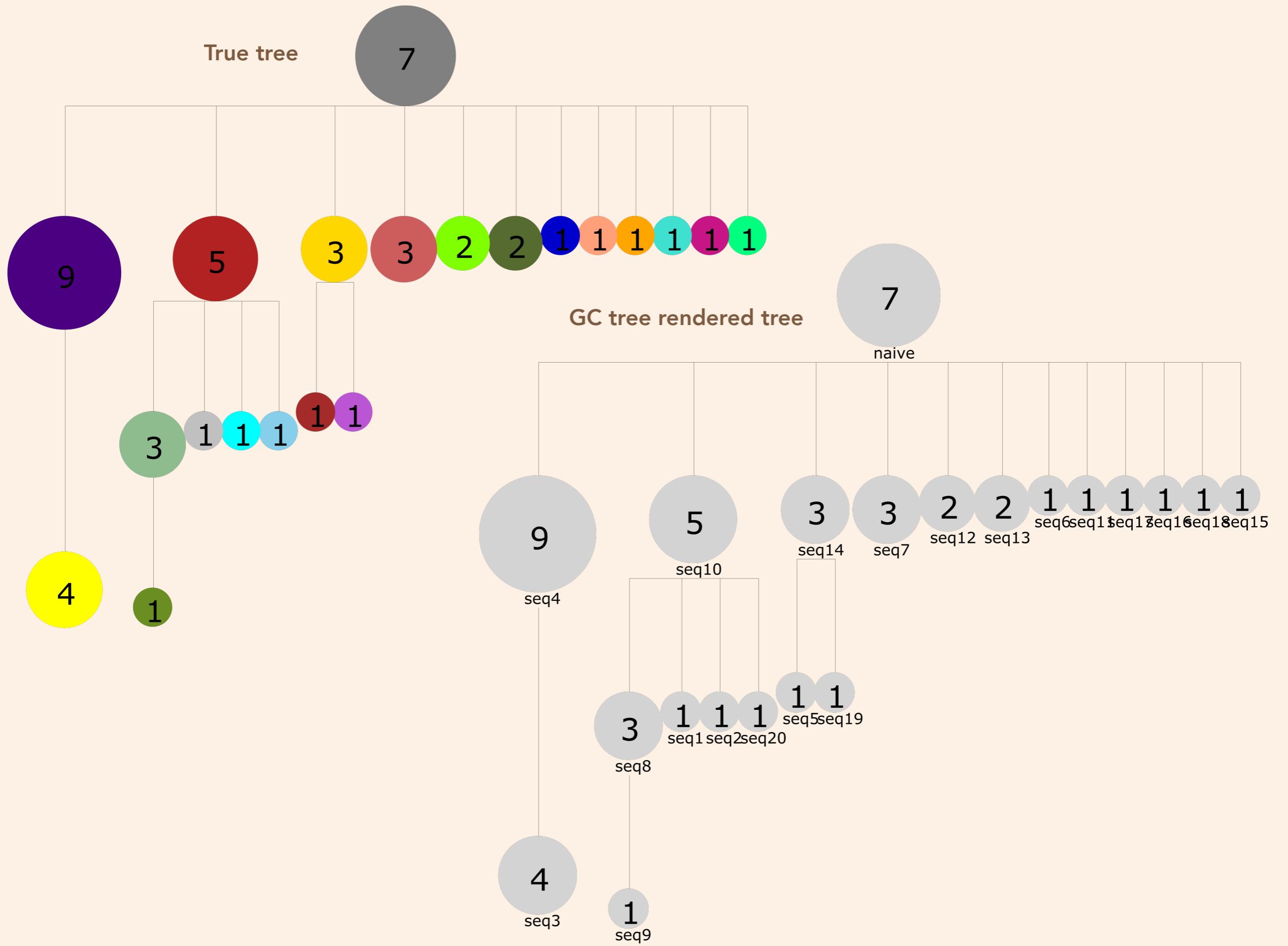


GC tree

simulated dataset

GCtree-simulated dataset

- ▶ containing 50 heavy chain V gene sequences
- ▶ germline : the germline sequence of the V gene used in the V(D)J rearrangement that defines this clonal family



GC tree

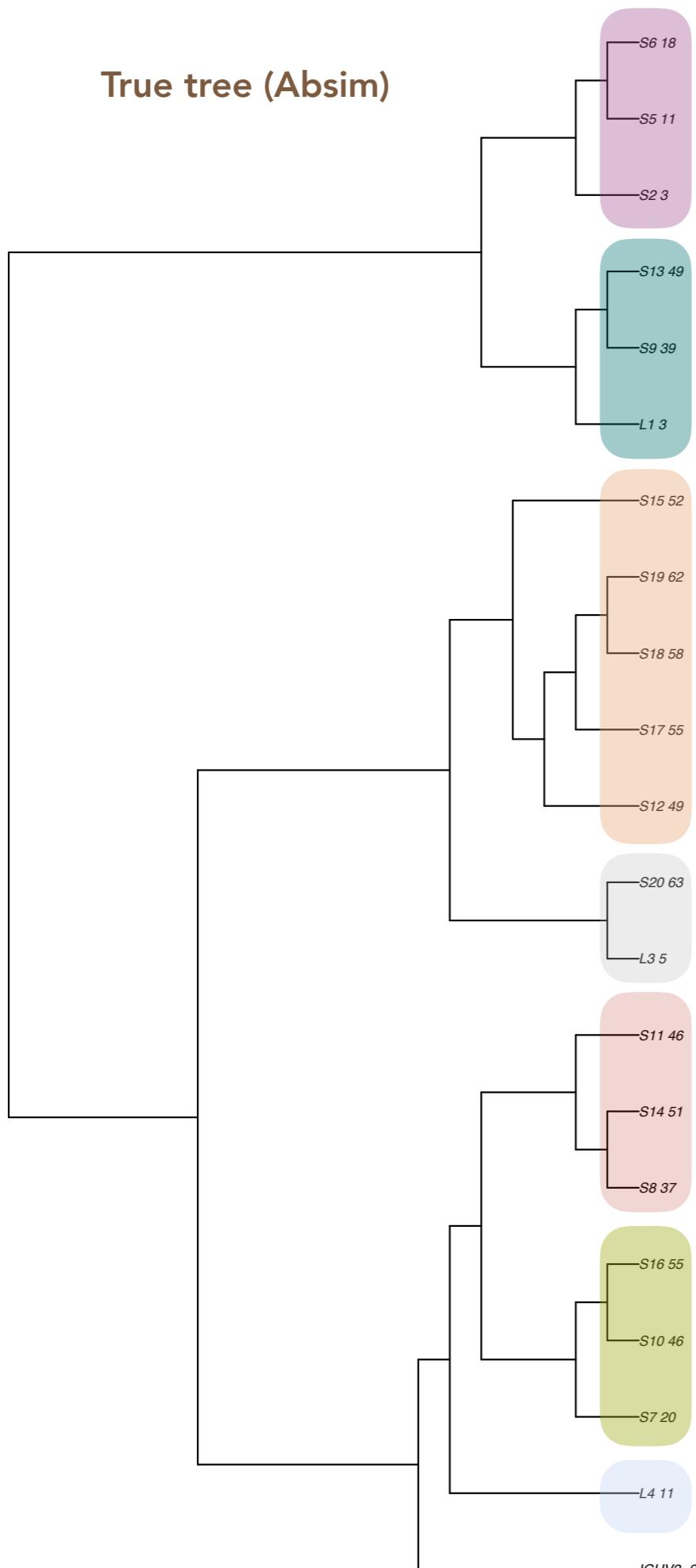
simulated dataset

Absim-simulated dataset

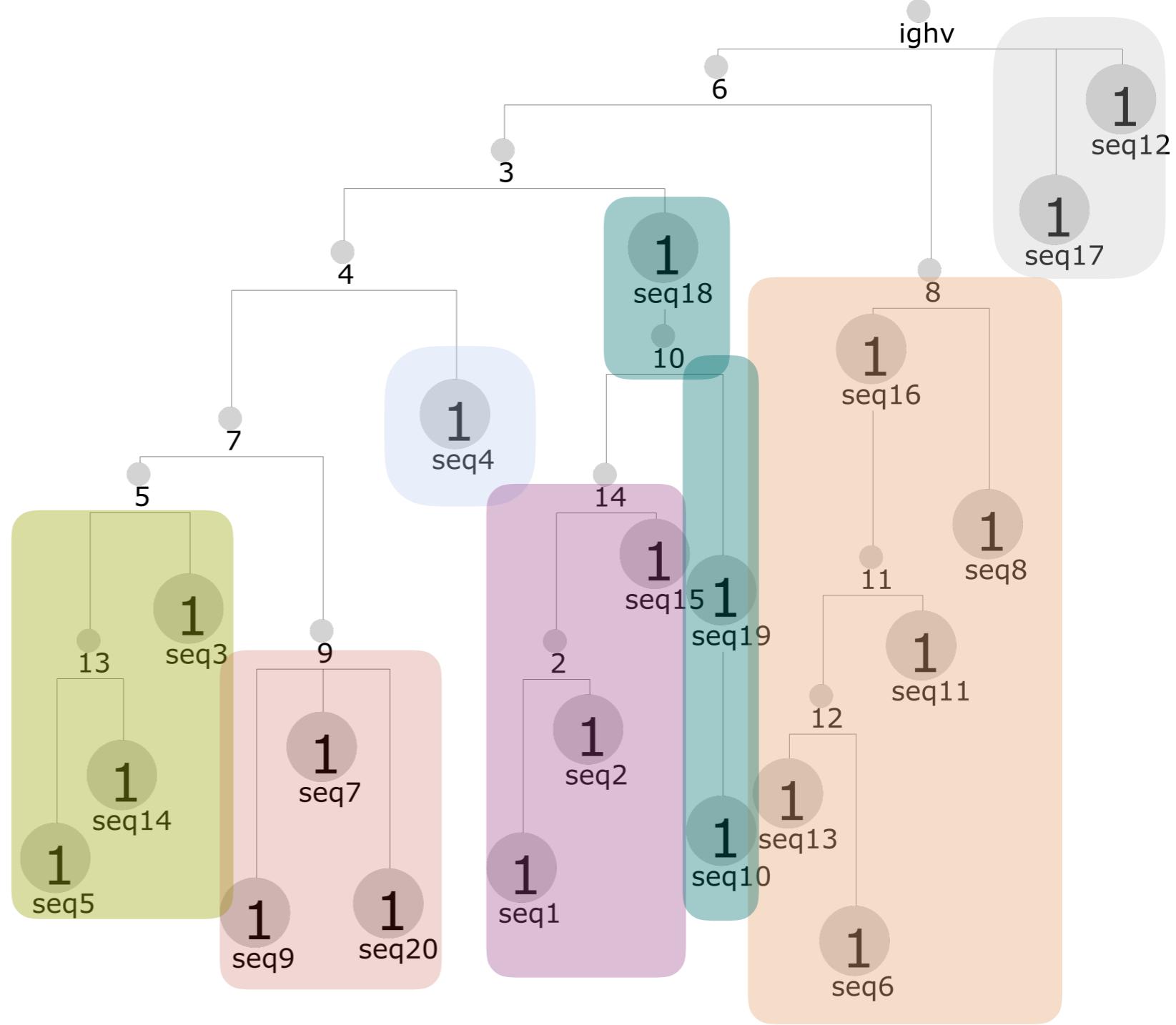
- ▶ containing 20 IGH chain sequences
- ▶ germline : the germline sequence used in the V(D)J rearrangement that defines this clonal family

Gap as un unknown nucleotide

True tree (Absim)

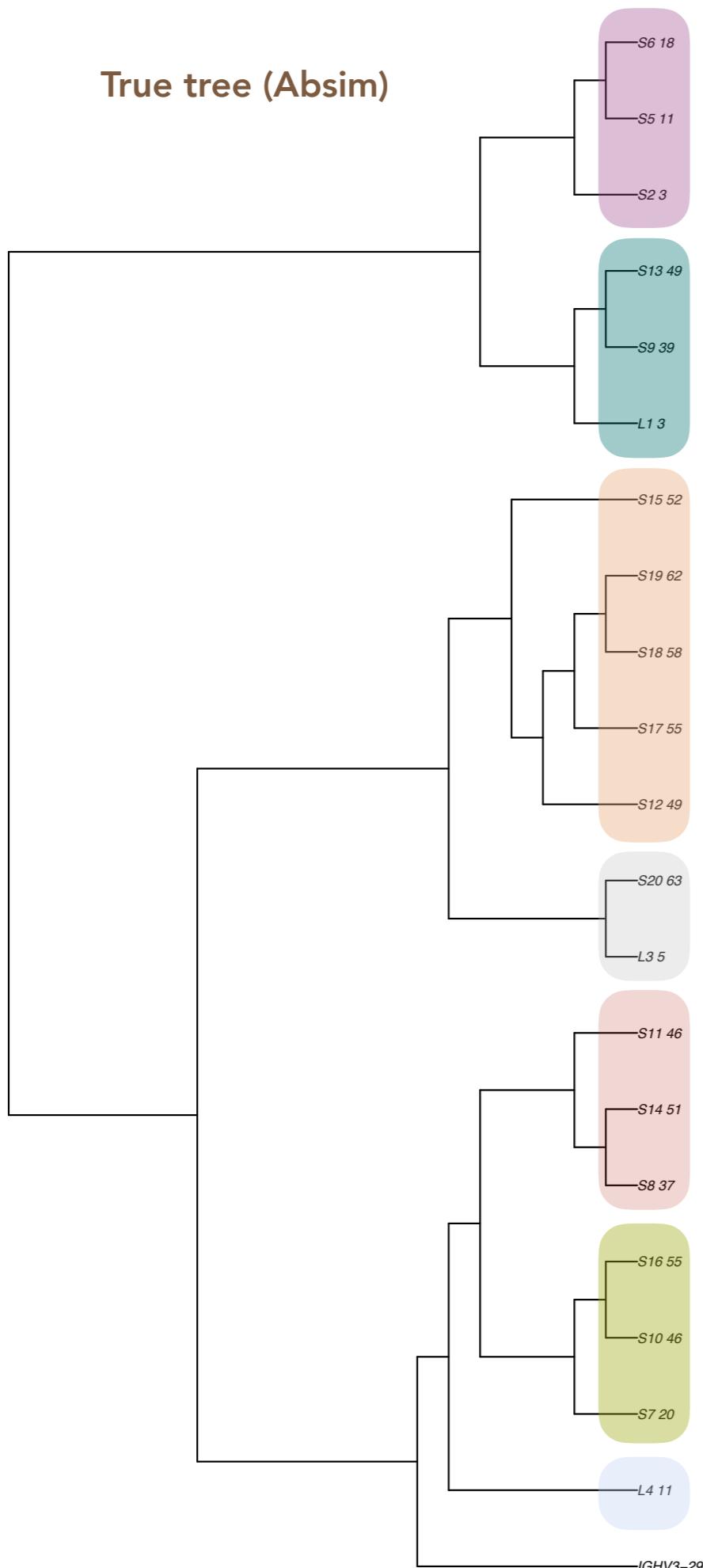


GC tree rendered tree

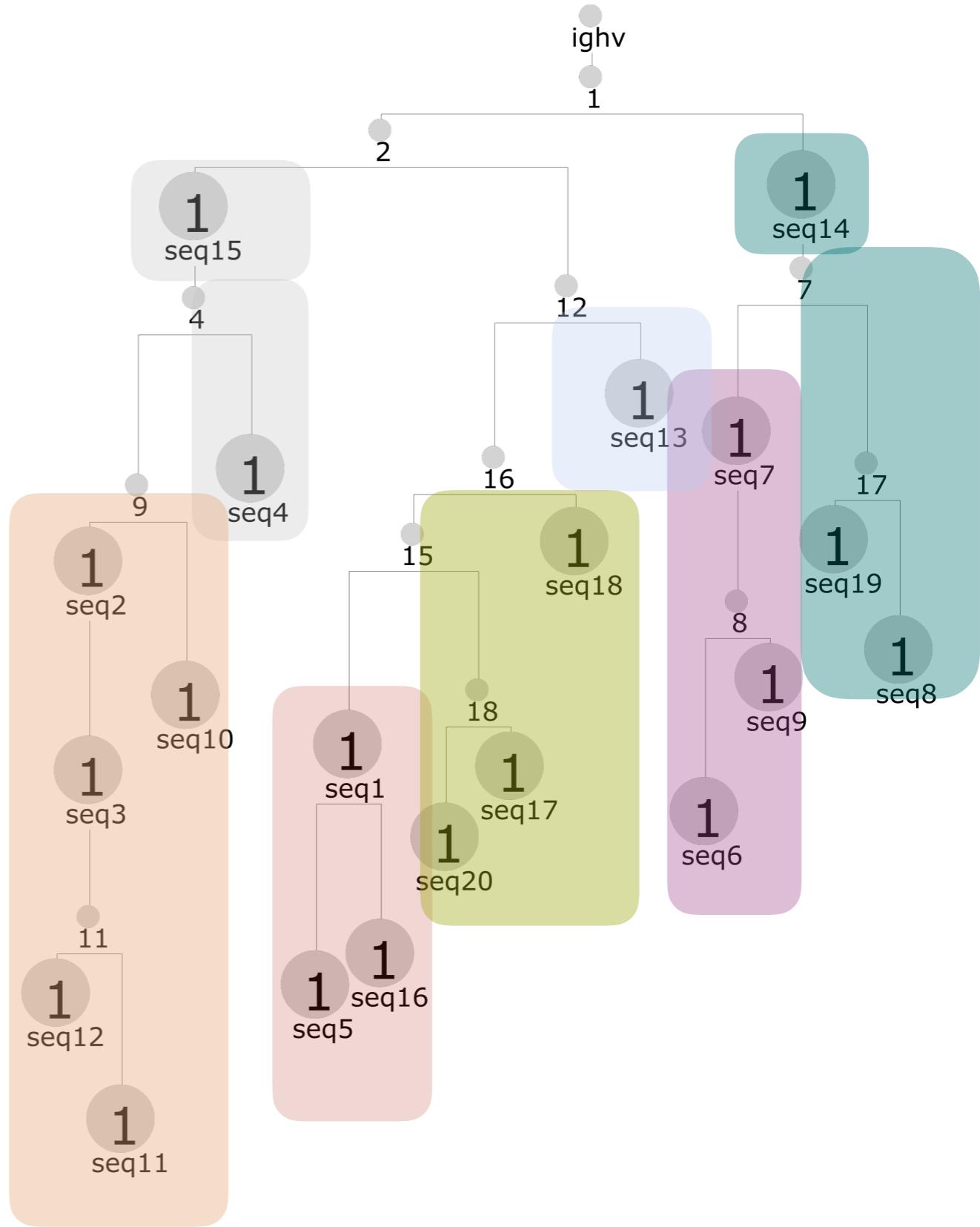


Gap as a fifth nucleotide state

True tree (Absim)



GC tree rendered tree



Study of intraclonal-diversity in a chronic lymphocytic leukaemia case

Using Vidjil and GCtree

Pipeline

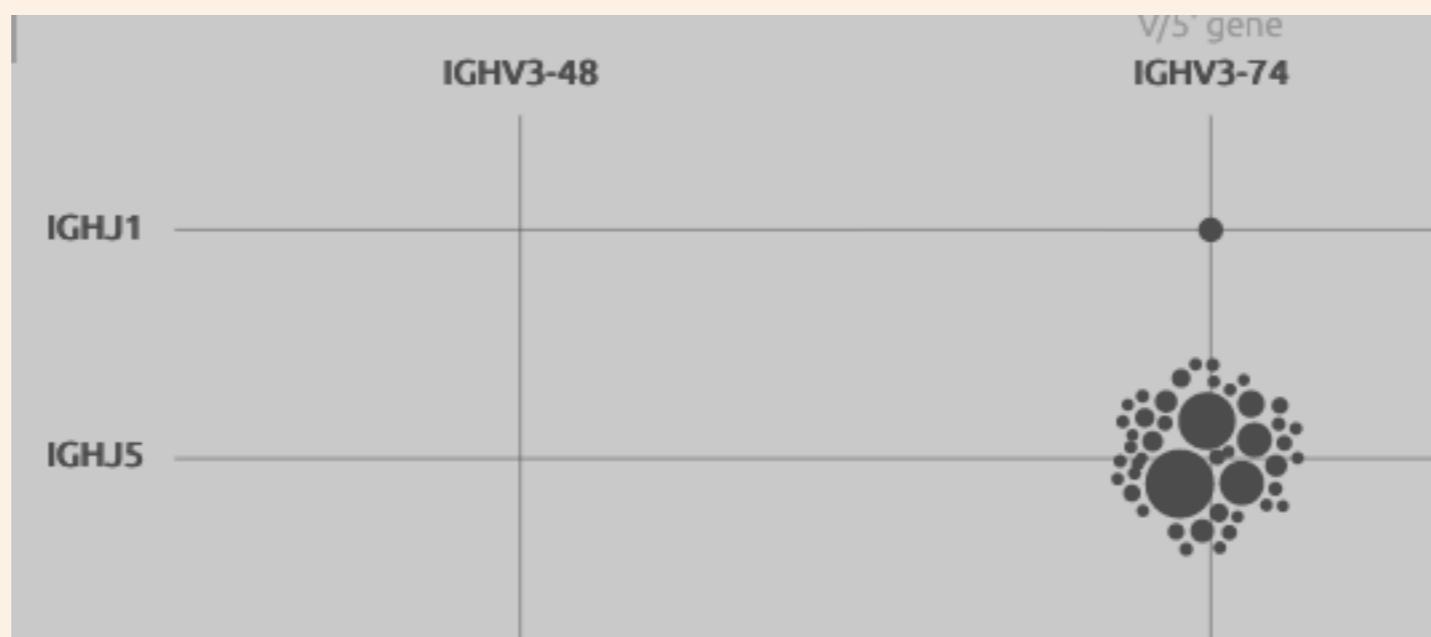
P2 dataset with 33599 seq

Run Vidjil

Select 10 most abundant clones

Run GCtree on the selection (consensus sequence)

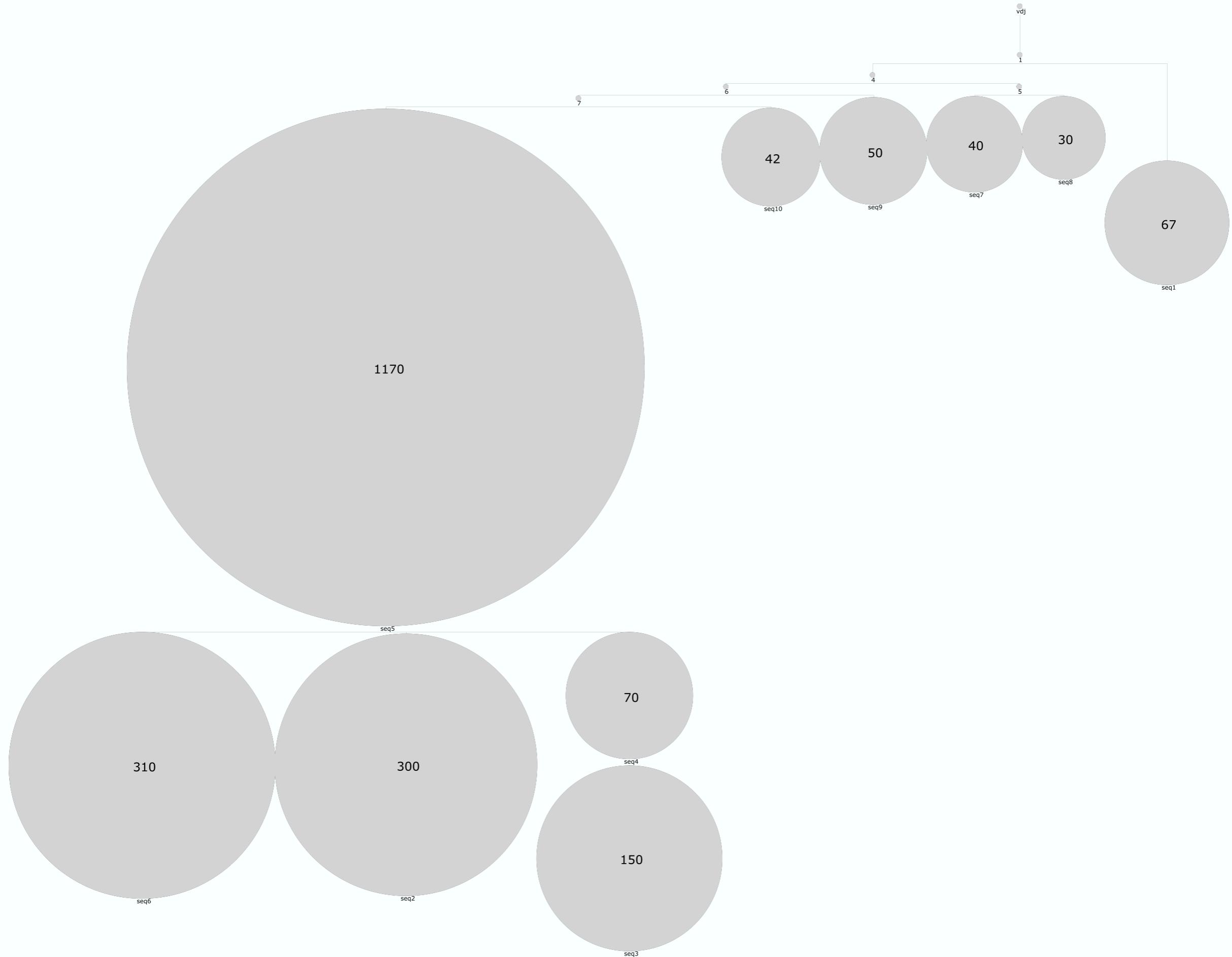
P2 major clones and satellites



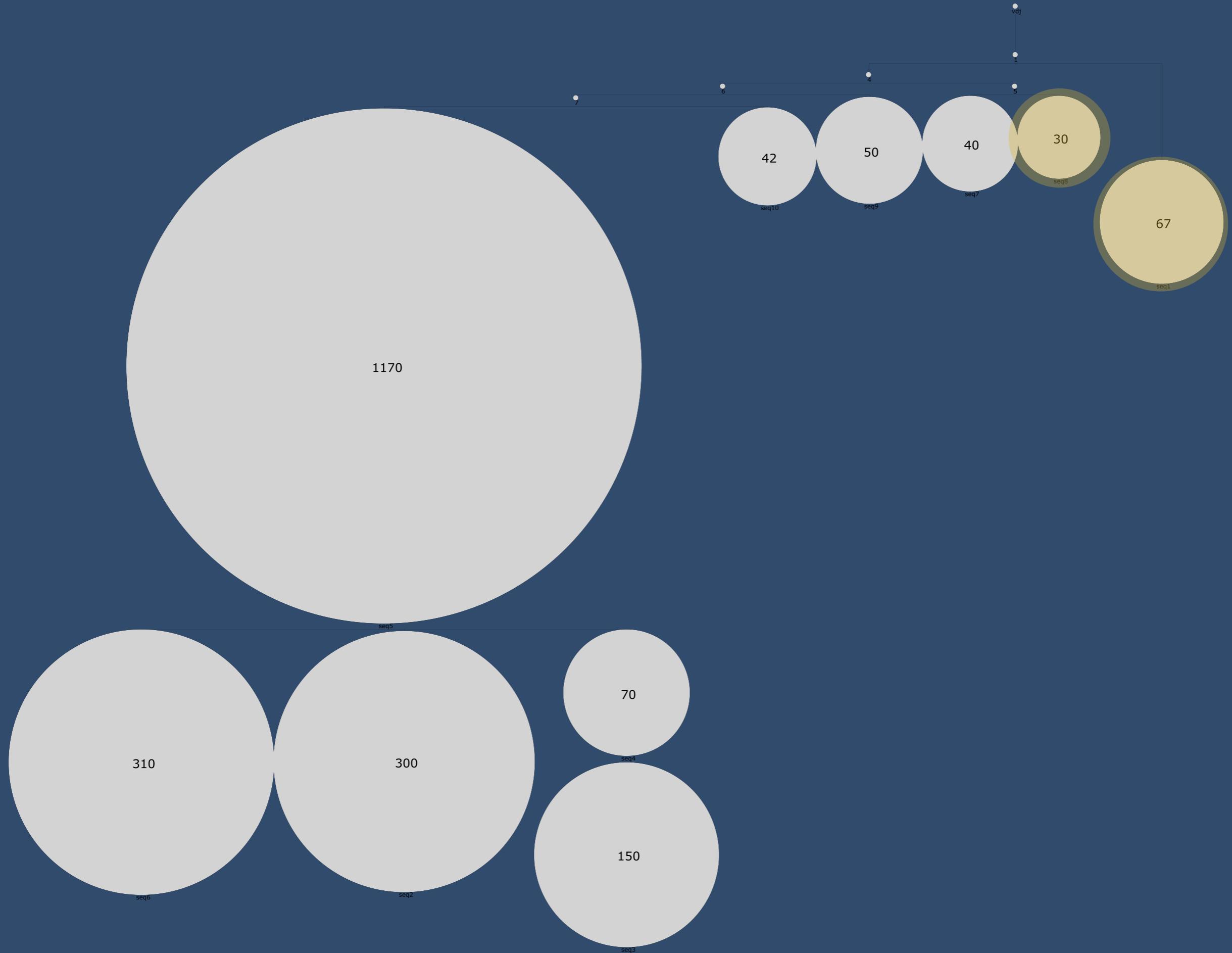
P2 major clones and satellites

Clone rank	# read	% read	IGHV	IGHJ
clone 1	11772	35 %	V3-74*01	J5*02
clone 2	6778	20 %	V3-74*02	J5*02
clone 3	3191	10 %	V3-74*01	J5*02
clone 4	1494	4 %	V3-74*01	J5*02
clone 5	699	2 %	V3-74*01	J5*02
clone 6	513	2 %	V3-74*03	J1*01
clone 7	464	1 %	V3-74*01	J5*02
clone 8	395	1 %	V3-74*01	J5*02
clone 9	369	1 %	V3-74*01	J5*02
clone 10	301	0,08 %	V3-74*01	J5*02

Intraclonal diversity via GCTree



Intraclonal diversity



Intraclonal diversity

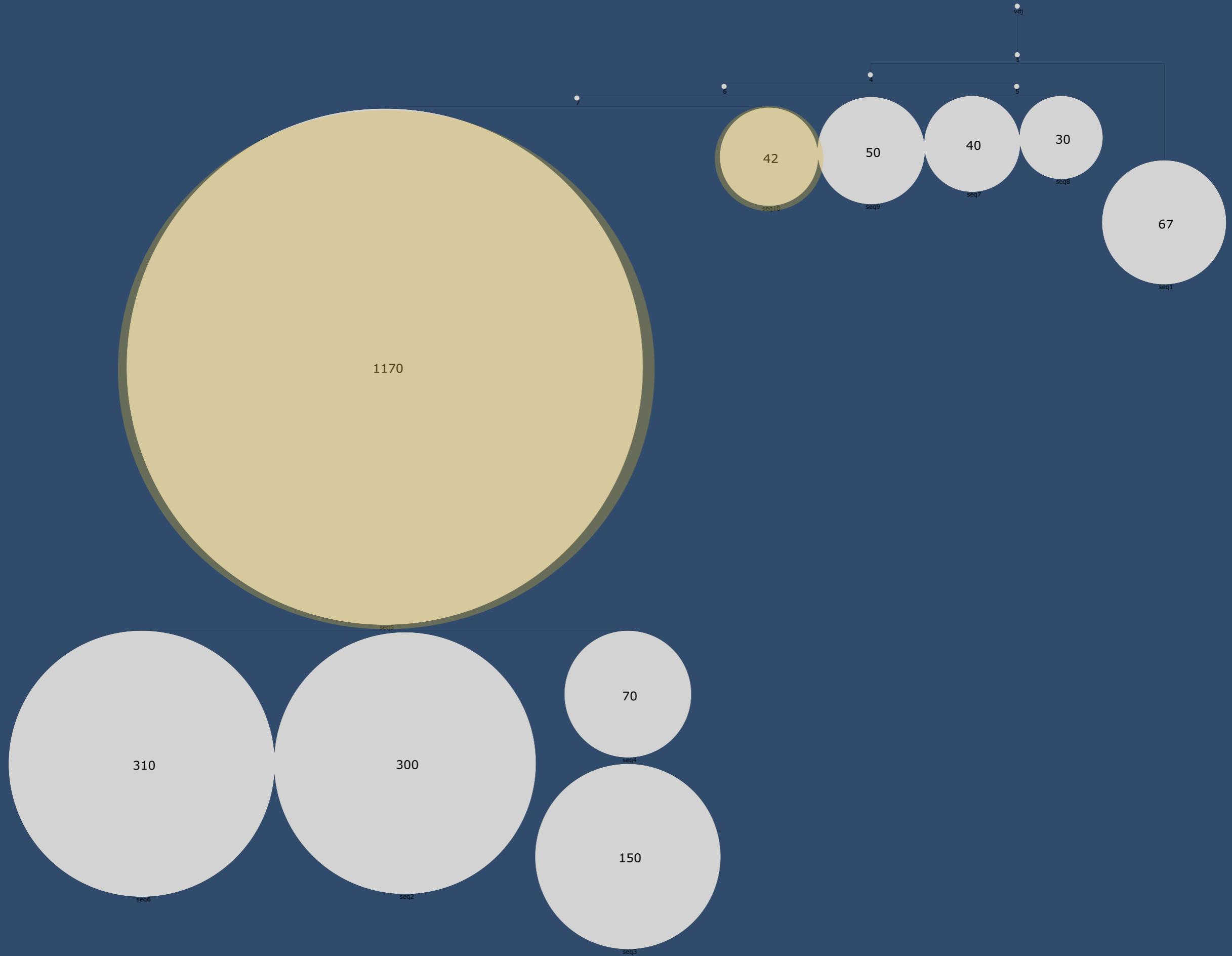
30

104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)	
C	A	R	D	Q	G	S	A	D	T	G	V	G	E	A	V	P	F	D	S	W						
30	tgt	gca	aga	gat	cag	ggg	tca	gcg	gat	aca	ggt	gtt	gga	aca	gcc	gta	ccc	ttt	gac	tcc	tgg	+	19	2,140.28	4.15	CARDQGSADTGVTAVPFDSW

67

104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)	
C	A	R	D	Q	G	S	A	D	T	G	V	G	T	A	V	P	F	D	S	W						
67	tgt	gca	aga	gat	cag	ggg	tca	gcg	gat	aca	ggt	gtt	ggt	aca	gcc	gta	ccc	ttt	gac	tcc	tgg	+	19	2,140.28	4.15	CARDQGSADTGVTAVPFDSW

Intraclonal diversity



Intraclonal diversity

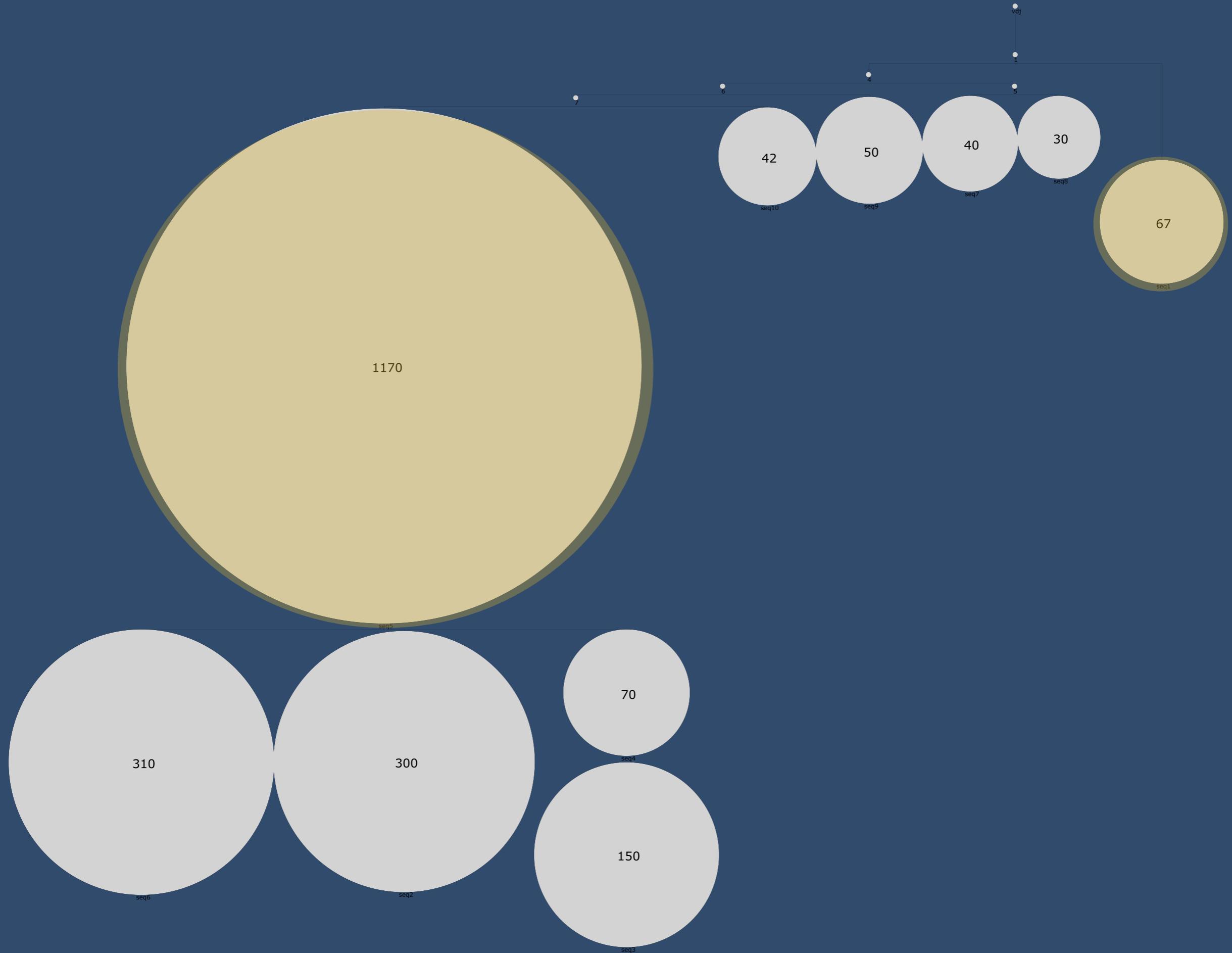
1170

104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)	
C	V	R	D	Q	G	S	A	D	T	G	V	G	T	A	V	P	F	D	S	W						
1170	tgt	gta	aga	gat	cag	ggg	tca	gcg	gac	aca	ggt	gtg	ggc	aca	gcc	gta	ccc	ttt	gac	tca	tgg	+	19	2,168.33	4.15	CVRDQGSADTGVTAVPFDSW

42

104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)	
C	A	R	D	Q	G	S	A	D	T	G	V	G	T	A	V	P	F	D	S	W						
42	tgt	gca	aga	gat	cag	ggg	tca	gcg	gac	aca	ggt	gtg	ggc	aca	gcc	gtc	ccc	ttt	gac	tca	tgg	+	19	2,140.28	4.15	CARDQGSADTGVTAVPFDSW

Intraclonal diversity



Intraclonal diversity

1170

104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)	
C	V	R	D	Q	G	S	A	D	T	G	V	G	T	A	V	P	F	D	S	W						
1170	tgt	gta	aga	gat	cag	ggg	tca	gcg	gac	aca	ggt	gtg	ggc	aca	gcc	gta	ccc	ttt	gac	tca	tgg	+	19	2,168.33	4.15	CVRDQGSADTGVGTAVPFDSW

67

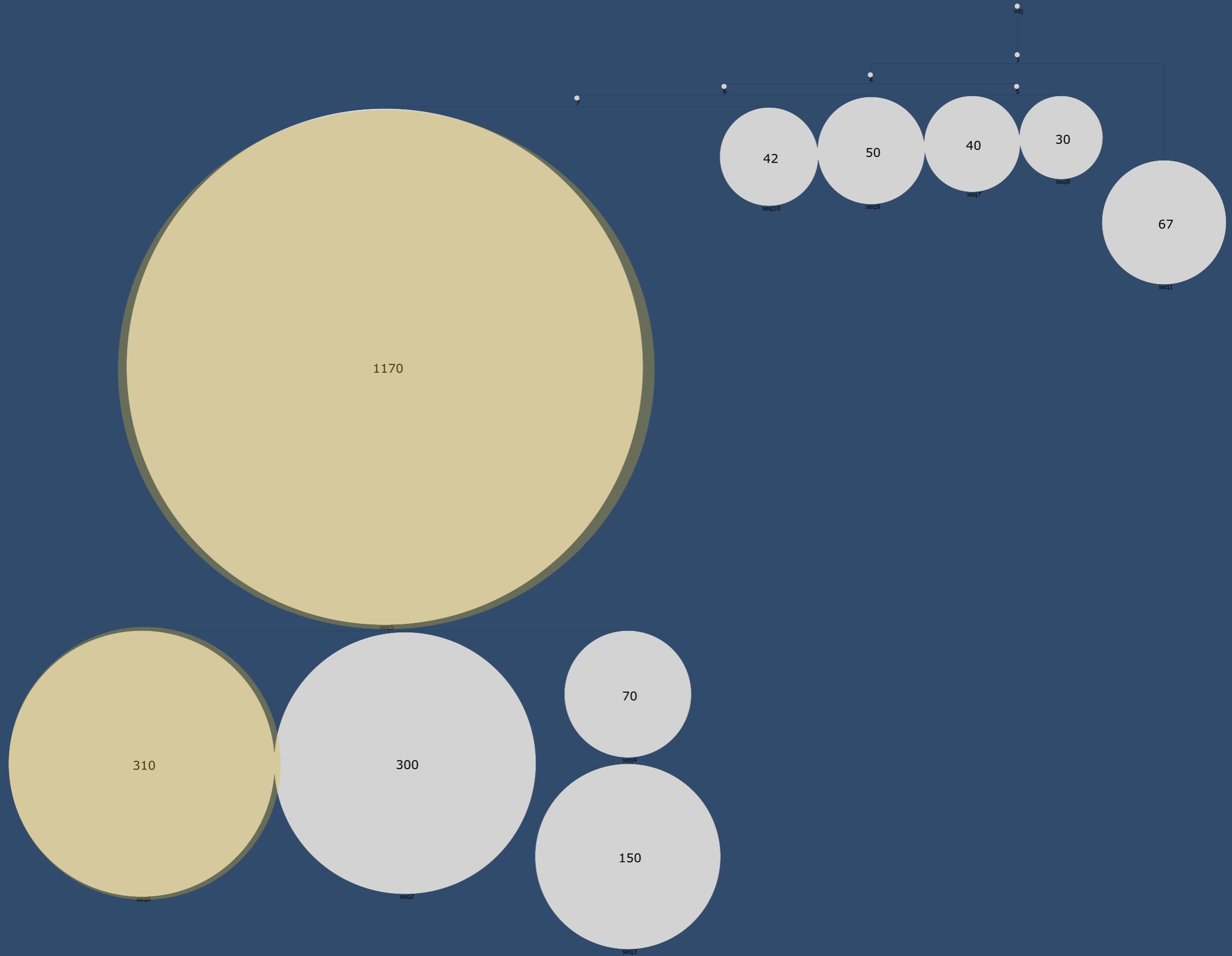
104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)	
C	A	R	D	Q	G	S	A	D	T	G	V	G	T	A	V	P	F	D	S	W						
67	tgt	gca	aga	gat	cag	ggg	tca	gcg	gat	aca	ggt	gtg	ggt	aca	gcc	gta	ccc	ttt	gac	tcc	tgg	+	19	2,140.28	4.15	CARDQGSADTGVGTAVPFDSW

Intraclonal diversity

7. V-REGION translation

7. V-REGION translation

Intraclonal diversity



Intraclonal diversity

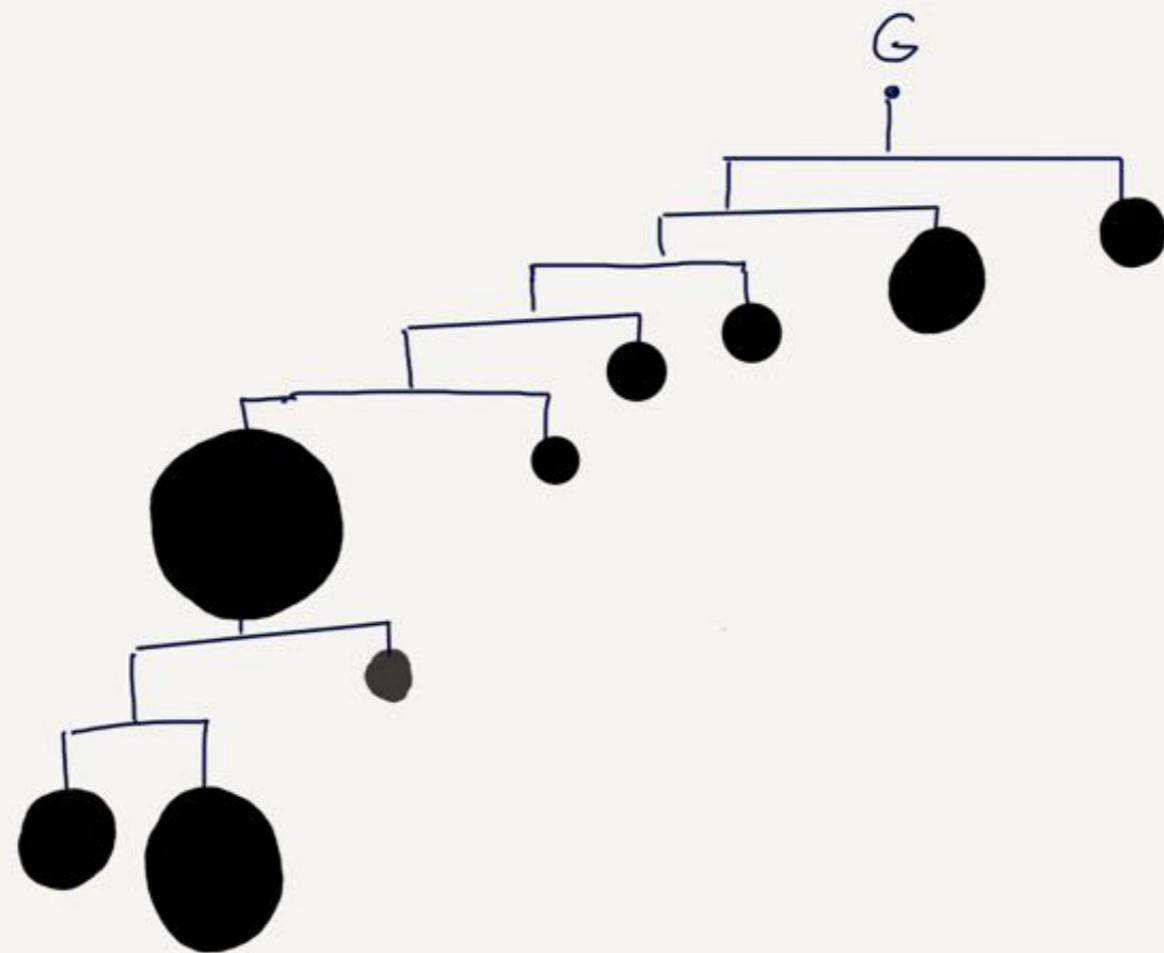
1170

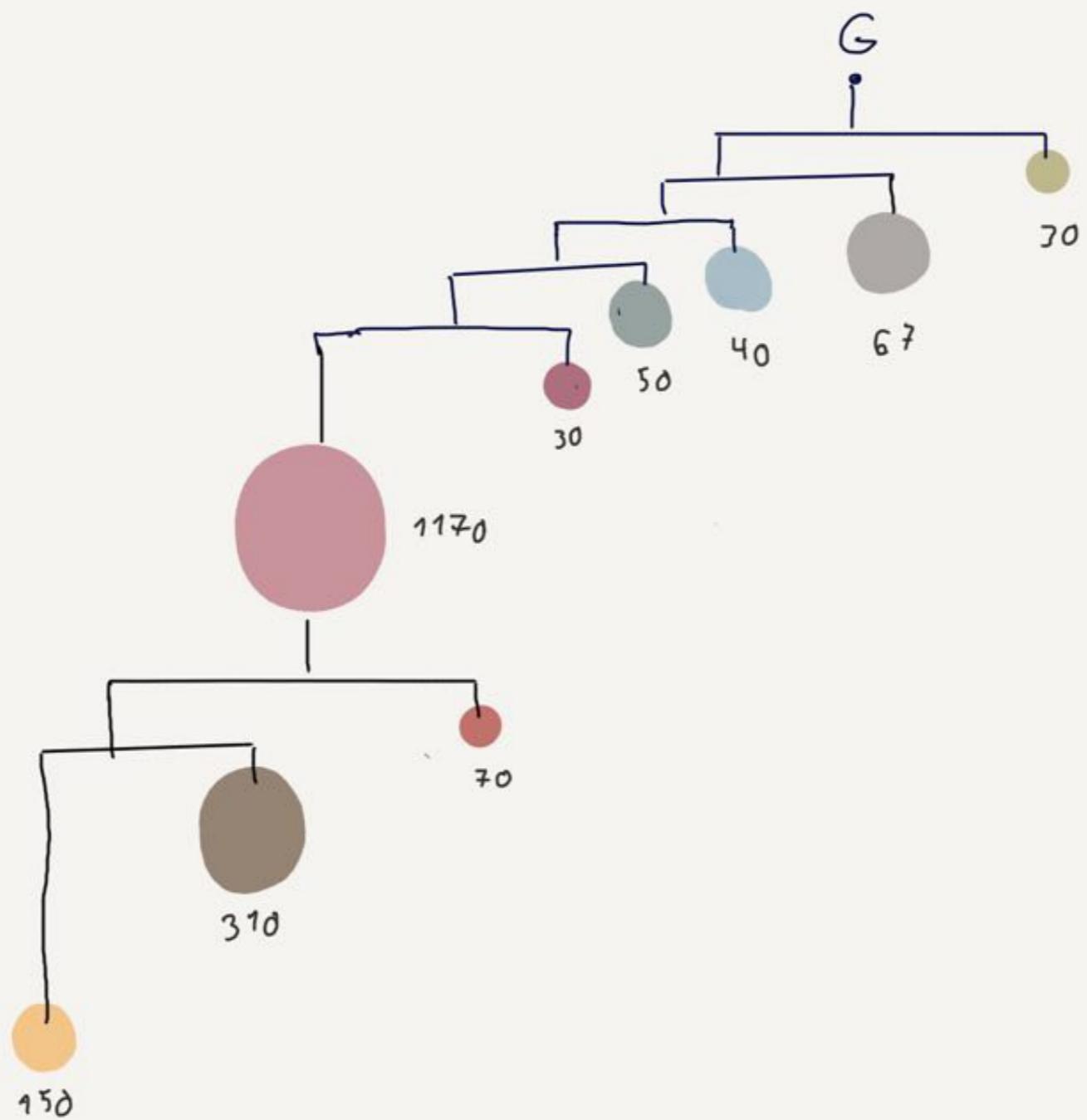
	104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)
	C	V	R	D	Q	G	S	A	D	T	G	V	G	T	A	V	P	F	D	S	W					
1170	tgt	gta	aga	gat	cag	ggg	tca	gcg	gac	aca	ggt	gtg	ggc	aca	gcc	gta	ccc	ttt	gac	tca	tgg	+	19	2,168.33	4.15	CVRDQGSADTGVTAVPFDSW

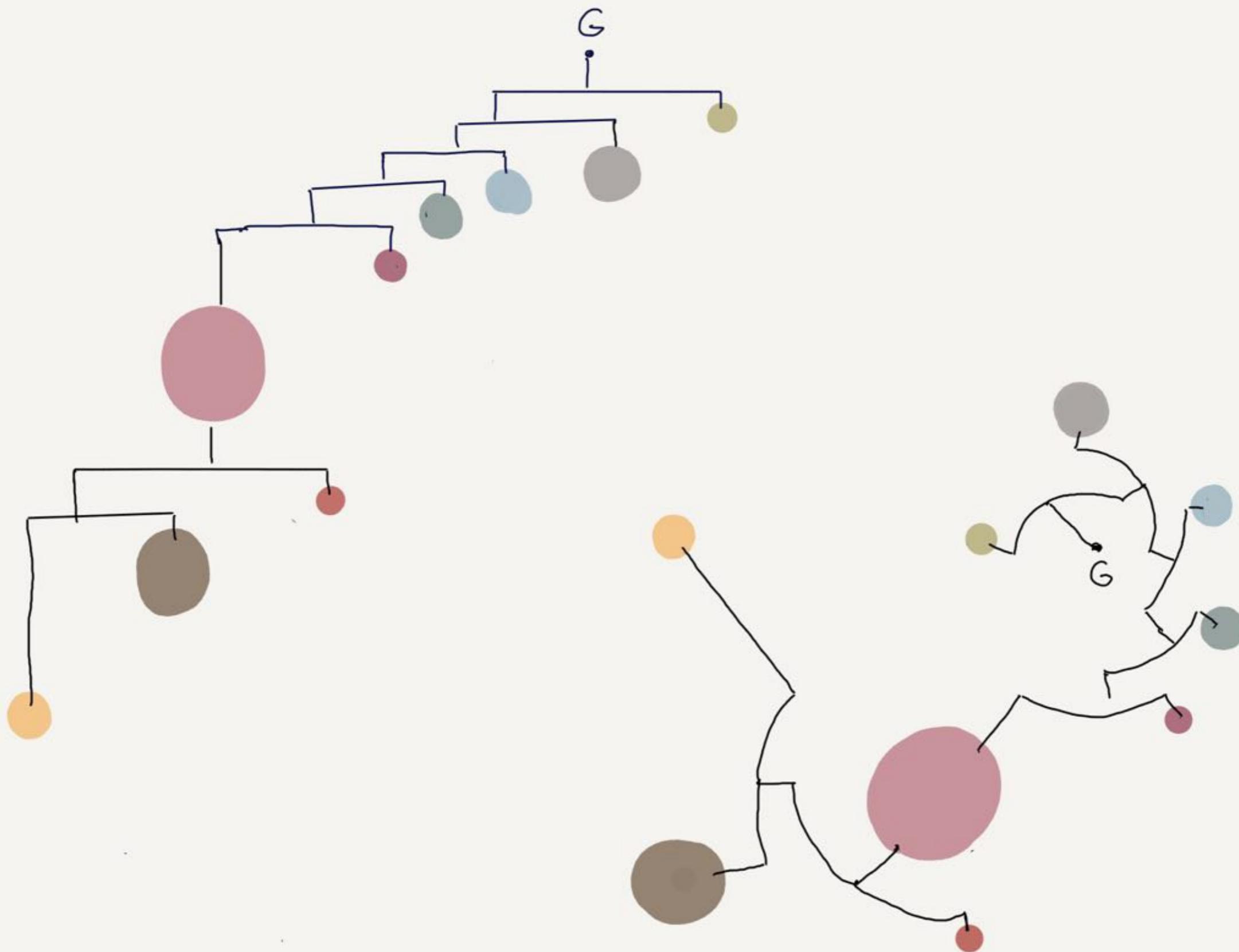
310

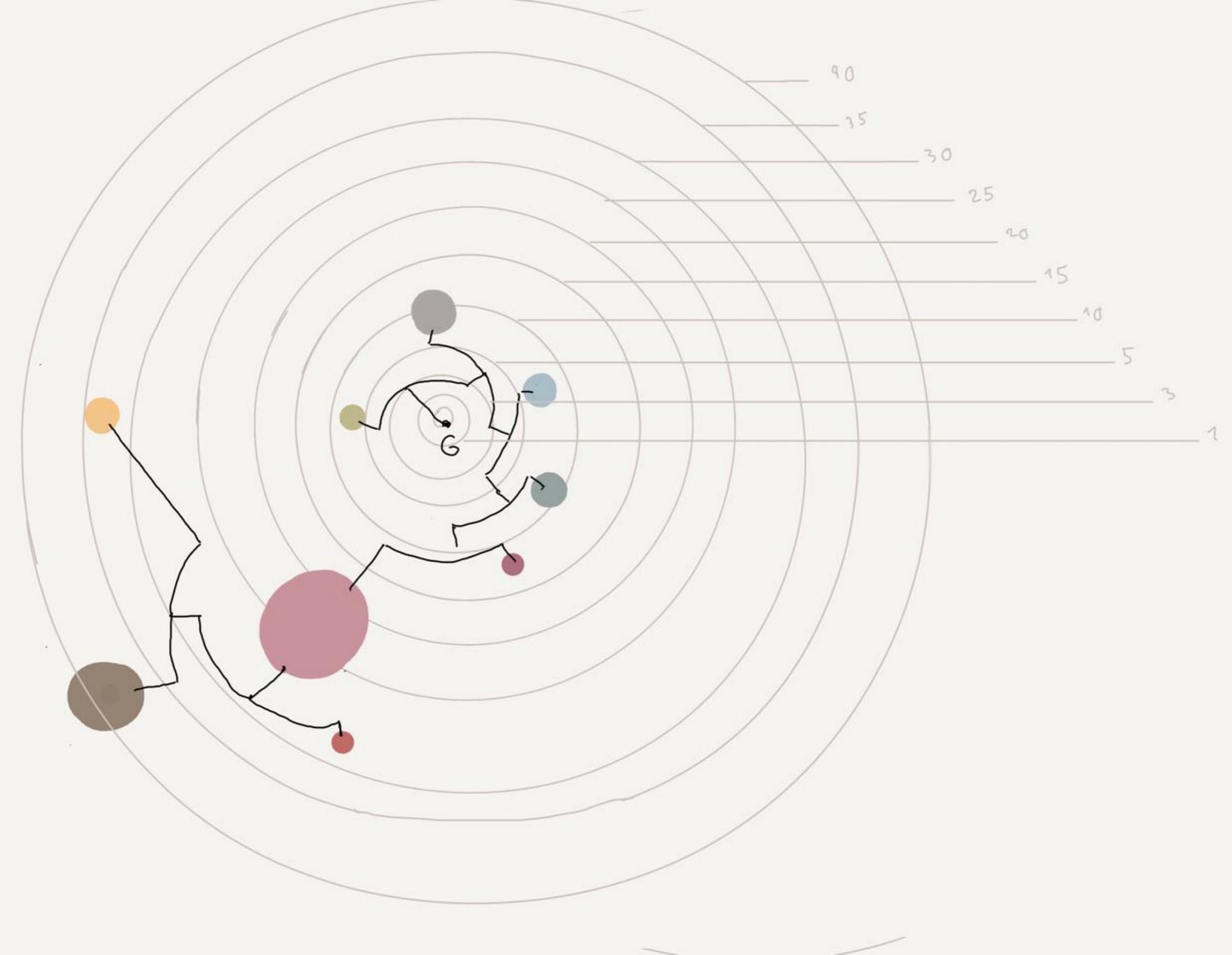
	104	105	106	107	108	109	110	111	111.1	111.2	111.3	112.3	112.2	112.1	112	113	114	115	116	117	118	Frame	CDR3- IMGT length	Molecular mass	pI	PhysicoChemical Descriptor (by BRFAA)
	C	V	R	D	Q	G	S	A	D	T	G	V	G	T	A	V	P	F	D	S	W					
310	tgt	gta	aga	gat	cag	ggg	tca	gcg	gac	aca	ggt	gtg	ggc	aca	gcc	gtg	ccc	ttt	gac	tca	tgg	+	19	2,168.33	4.15	CVRDQGSADTGVTAVPFDSW

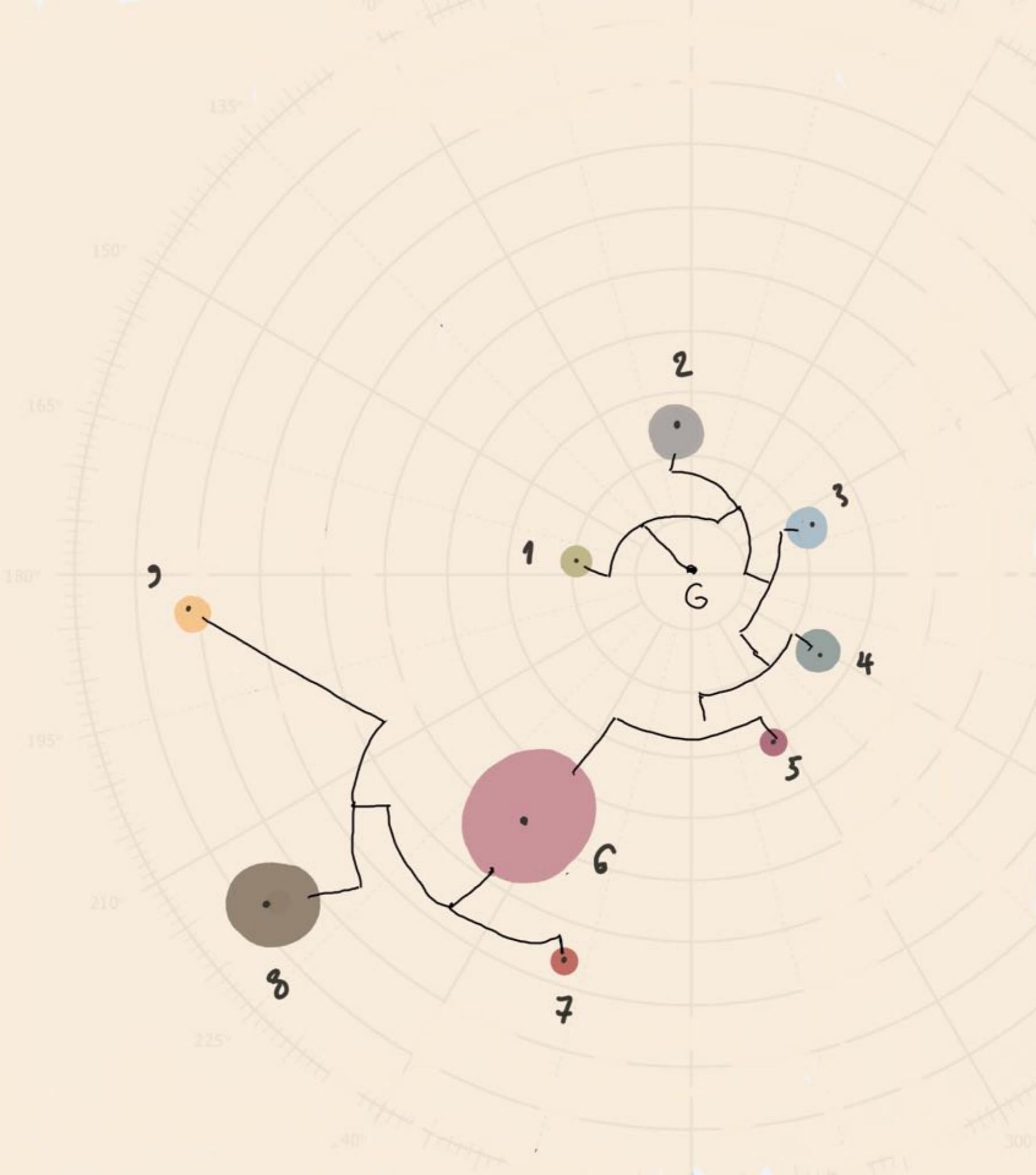
Visualisation prototype











Benchmarking RepSeq analysis tools

Benchmarking RepSeq analysis tools

Improvements

Studied **eight** tools

Analyzed different repertoire simulators

Studied different performance evaluation measures and found out that for the same input data, using different tools results in having

- Variable number of clusters
- Variable cluster sizes (especially for the major clones of each repertoire)
- Various standards to accept a sequence as an input.

Prepared an *in vitro* benchmark to evaluate the clone detection caliber of each tool. Data set is constructed by a serial dilution of a known clone (1\%, 0.1\%, and 0.01\%) in a polyclonal background

Thank you