Assessing recombination detection methods using viral simulations

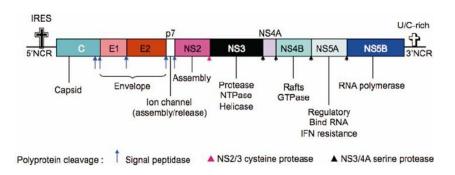
Fred Jaya - Stage 2

28/10/19





Hepatitis C virus (HCV)

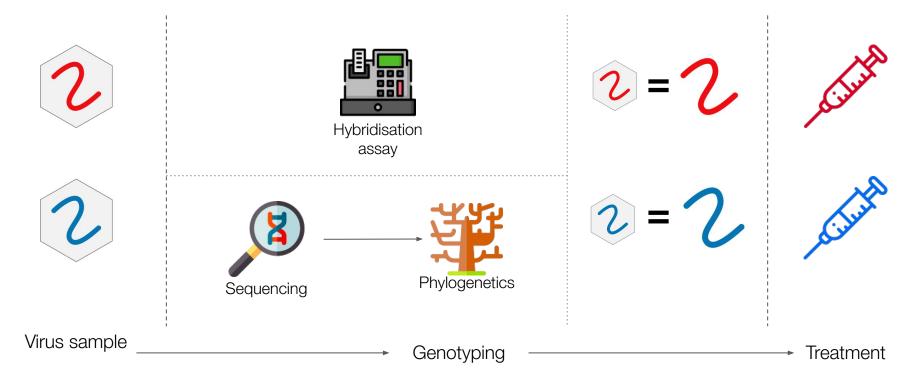


- ~71 million people with chronic infection
- Liver cirrhosis, carcinoma
- Single-stranded positive-sense RNA (+ssRNA)
- ~9000 bp ORF
- 7 genotypes, many subtypes

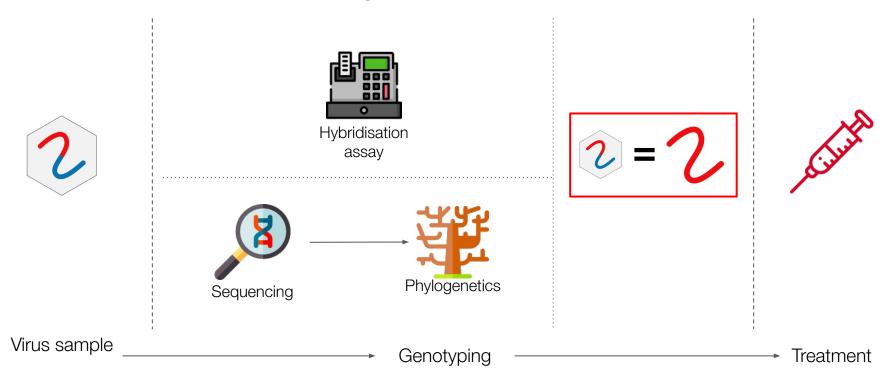
- Mutation + recombination → diversity
- Host expansion, drug and immune evasion
- High mutation rates → RdRp
- Infrequent recombination
- NGS has shown circulating recombinants



Genotyping HCV



Genotyping Recombinant HCV





Effect of recombination



Hybridisation assays

Commercial kits, VERSANT

Unable to identify recombinants



Phylogenetics

Obscures evolutionary relationships



Treatments

Effectiveness of treatments are type specific



Solution...



Next-Generation Sequencing (NGS)

Deep sequencing, WGS Higher-resolution for subtyping Surveillance of genotypes

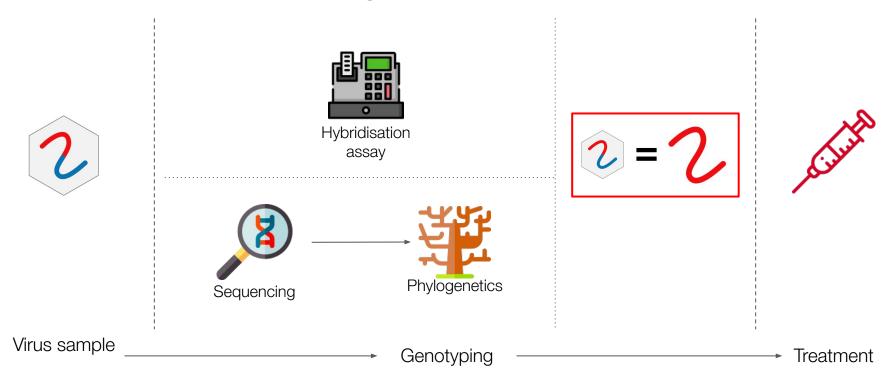


Bioinformatics

Recombination detection methods (RDMs)



Genotyping Recombinant HCV







Assessment of recombination detection methods

using viral simulations





Phylogenetics + bioinformatics



Recombination detection

Sequencing

Genotyping



RDM benchmarks

INVITED TECHNICAL REVIEW

Analysing recombination in nucleotide sequences

DARREN P. MARTIN,* PHILIPPE LEMEY+ and DAVID POSADA‡

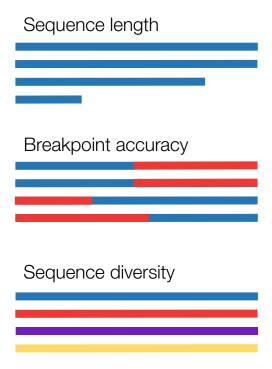
Detecting recombination in evolving nucleotide sequences Cheong Xin Chan, Robert G Beiko and Mark A Ragan*

Evaluation of methods for detecting recombination from DNA sequences: Computer simulations

David Posada* and Keith A. Crandall

A Simulation Study of the Reliability of Recombination Detection Methods

Carsten Wiuf,* Thomas Christensen,† and Jotun Hein‡





Limitations

- RDMs were not designed to process large (NGS) datasets
- No clear outline which RDMs are most effective for which datasets
 - Simulations are simplified



Research questions

- 1) Identify methods that are capable of being scaled to process NGS data
- 2) Explore the effect of evolutionary and sequence properties on RDM behaviour
 - a) Can we suggest effective RDMs based on viral properties i.e. mutation, recombination rate?



Next-generation sequencing (NGS)

Sanger ABI 3730xl	PacBio Sequel II HiFi
600 - 1000 bases	~ 15,000 bases (high accuracy)
96 reads	~ 4,000,000 reads (SMRT Cell 8M)



Methodology

- 1) Viral simulations and parametric sweep
- 2) Recombination detection method analyses
- 3) Pipeline design



Simulations

Mutation rate	Recombination rate	Dual infection probability	Sample size
10 ⁻⁸	10 ⁻⁸	0.00	100
10 ⁻⁷	10 ⁻⁷	0.05	1000
10 ⁻⁶	10 ⁻⁶	0.10	2500
10 ⁻⁵	10 ⁻⁵	0.25	5000
10-4	10 ⁻⁴	0.50	10000*
10 ⁻³	10 ⁻³	1.00	

Datasets were generated between all individual parameters

Simulation 1 - bold

Simulation 2 - all; * except n = 10000

SANTA-SIM - viral simulator



Pipeline design

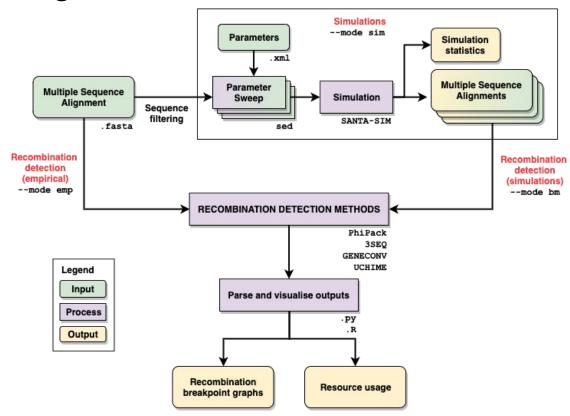
864 unique parameter combinations

- → 3 replicates
 - → 2592 total files to simulate
 - → 3 RDMs
 - \rightarrow 7776 analyses





Pipeline design





rec-bench



github.com/fredjaya/rec-bench





Results (initial)

Simulation stats



Breakpoint significance



Runtime (wall)



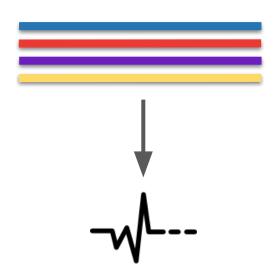
Recombination detection methods



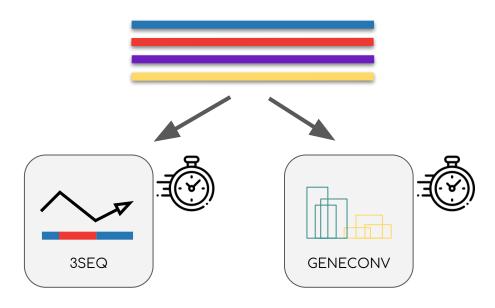




Key findings



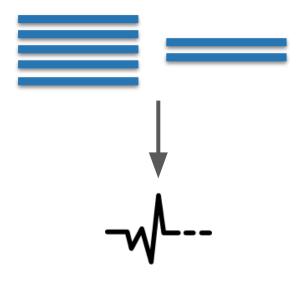
1. Diversity determines recombination detection



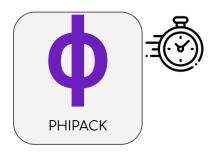
2. Diversity affects runtime in 3SEQ and GENECONV



Key findings



3. Sample size affects recombination detection



4. PhiPack is scalable

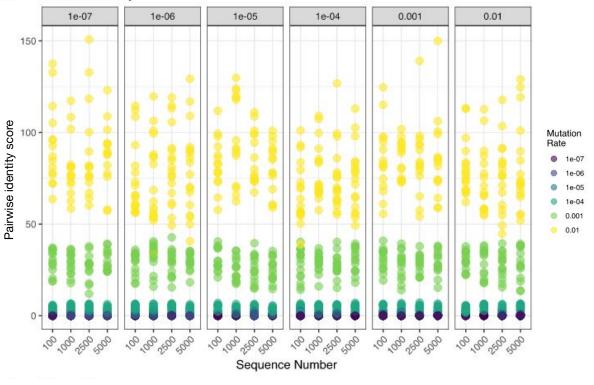


Results

Simulations

Mutation rate is correlated with output sequence diversity

A Mean diversity





√-- Recombination detection

Significance of detected signal is determined by algorithm/method



Pairwise Homoplasy Index

Observed number of states between two sites

No sites removed



Hypergeometric random walk

Triplets compared

No sites removed



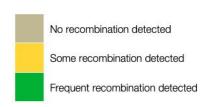
'BLAST-like'

Pairwise comparison

Monomorphic sites removed



-√--- Recombination detection





- Monomorphic sites = uninformative
- Sensitive to high divergence

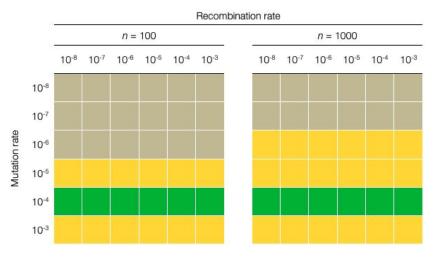
Recombination rate

Sensitive to recombination with n?

		Recomb							on rate					
		n = 100							n = 1000					
		10-8	10 ⁻⁷	10 ⁻⁶	10-5	10-4	10 ⁻³		10-8	10 ⁻⁷	10-6	10-5	10-4	10-3
	10-8													
	10 ⁻⁷													
on I are	10-6													
Mutation	10-5													
	10-4													
	10 ⁻³													

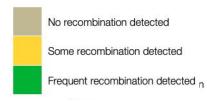


- Monomorphic sites = uninformative
- Very sensitive to a specific range
- n = affects diversity



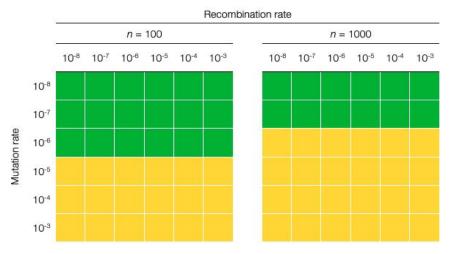


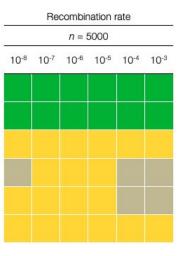






- Very sensitive
- Powerful at low divergence
- Increasing *n* reduces false positives?







25

√-- Recombination detection

Significance of detected signal is determined by algorithm/method



Pairwise Homoplasy Index

Observed number of states between two sites

No sites removed



Hypergeometric random walk

Triplets compared

No sites removed



'BLAST-like'

Pairwise comparison

Monomorphic sites removed





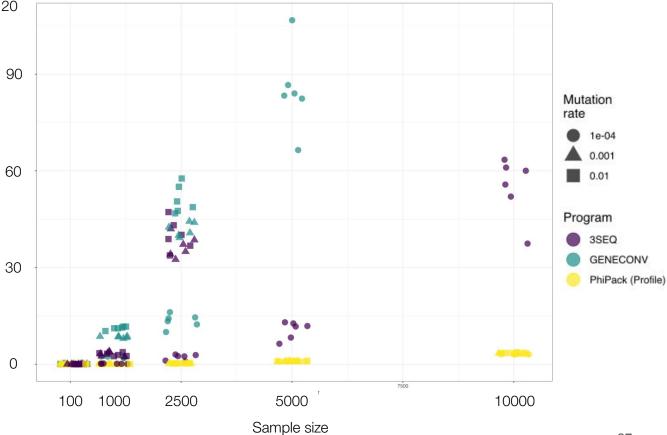




Hours



Mutation rate affects speed of 3SEQ and GENECONV







Wall time is determined by detection algorithm/method







Pairwise Homoplasy Index

Observed number of states between two sites

No sites removed

Hypergeometric random walk

Triplets compared

No sites removed

'BLAST-like'

Pairwise comparison

Monomorphic sites removed



Limitations

- Lack of detail
 - Presence/absence or positive/negative
 - Summarise data for comparison
- Unreliable reporting of simulated breakpoints
 - SANTA-SIM mod: cumulative breakpoints across generations
 - Doesn't account for false positives
 - No measure of accuracy
- Doesn't reflect empirical data
 - Generation of unfit sequences/populations
 - Mutation = diversity, what's the effect of recombination on sims?
 - Sequence diversity underrepresented
- Low replicates



Further work

Objective	NOV	DEC	JAN	FEB	MAR	APR
Data analysis (quantify, summarise, compare)						
Conferences (AusEvo, FoSTER, ABACBS)						
Empirical data (pipeline integration, analysis)						
UCHIME (RDM)						
T-RECs (RDM)						
HIV dataset						
Pre-print						
Tweak PhiPack						



Acknowledgements

Darling Lab

- Prof. Aaron Darling
- Dr. Barbara Brito
- Dr. Matt DeMaere
- Dr. Kay Anantanawat
- Dr. Mathieu Fourment
- Dr. Leigh Monahan
- Dr. Mike Imelfort
- Dr. Liza Kretzschmar
- Kevin Ying
- Daniela Gaio
- Nehleh Kargarfard
- Brent Bevear
- Sidaswar Krishnan
- Christian Cabato

the ithree institute

UTS eResearch

SHDRSA

Binfies

Level 7 HDRs



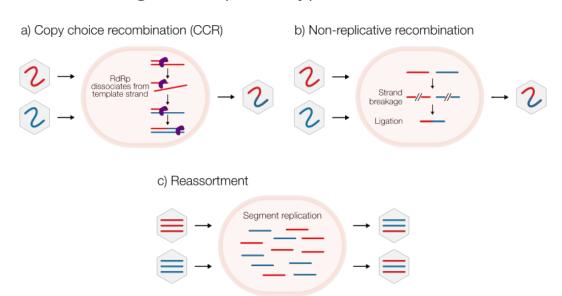
@fredjaya1



github.com/fredjaya

Recombination (viral)

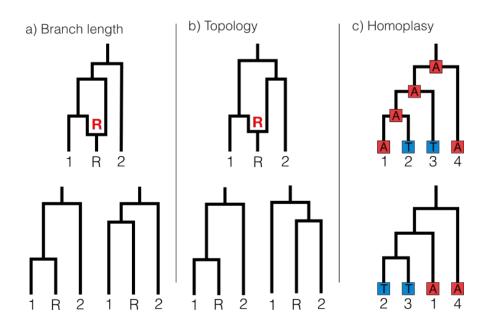
- Generation of a child genome with two parental genomes
- Co-infection of a host cell
- Affects evolution changes viral phenotype and fitness





Confounds phylogeny

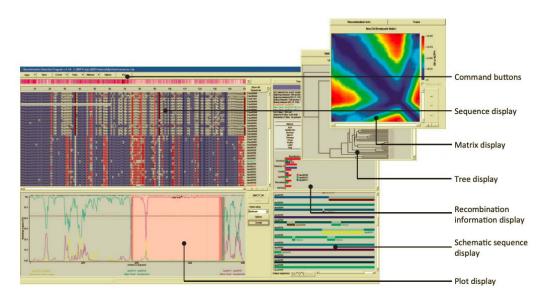
- Recombinants have two ancestral histories.
- Violates assumptions → inaccurate estimations





Recombination detection methods

- Test and account for recombination prior to analyses (sequences)
- Different statistical tests and algorithms
- Performance differs based on sequence properties
 - Recombination rate
 - Sequence diversity
 - Sample size





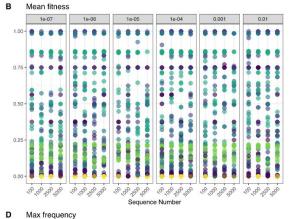
Simulations

	sim ₁							
n	PhiPack (Profile)	3SEQ	GENECONV	UCHIME				
100	Υ	Υ	Y	Ν				
1000	Y	Y	Y	Ν				
2500	Y	Y	Y	Ν				
5000	Y	Y	Y	Ν				
10000	Y	Y	Y	Ν				

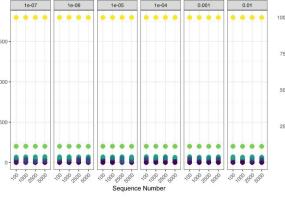
A	iviean diver	Sity			D IVIE	b wear inness			
	1e-07	1e-06	1e-05	1e-04	0.001	0.01		1e-07	1e-06
150					•		1.00 -		8
100					:		0.75 -	: : :	
							0.50		
50	!!!!		ilii				0.25 -		
0							0.00		
	'00'000 top top	'00' 000 too '00'	o o o o o o o o o o o o o o o o o o o	10,00,50,00	"00" 000 " 200 " OU. "	0,00,50,00	100	'000 Sea 2000	10,00,50,00
			Sequence	e Number					

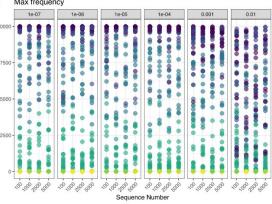
Mean diversity

Mean distance



	sim ₂							
n	PhiPack (Profile)	3SEQ	GENECONV	UCHIME				
100	Υ	Y	Y	Υ				
1000	Y	Y	Y	Υ				
2500	Ν	Ν	Ν	Ν				
5000	Y	Ν	Ν	Υ				
10000	N	N	N	N				





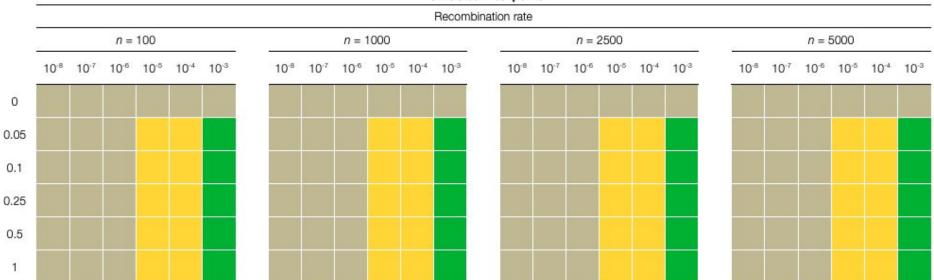
 $\mathbf{Y} = \text{Analysed}$

N = Not analysed



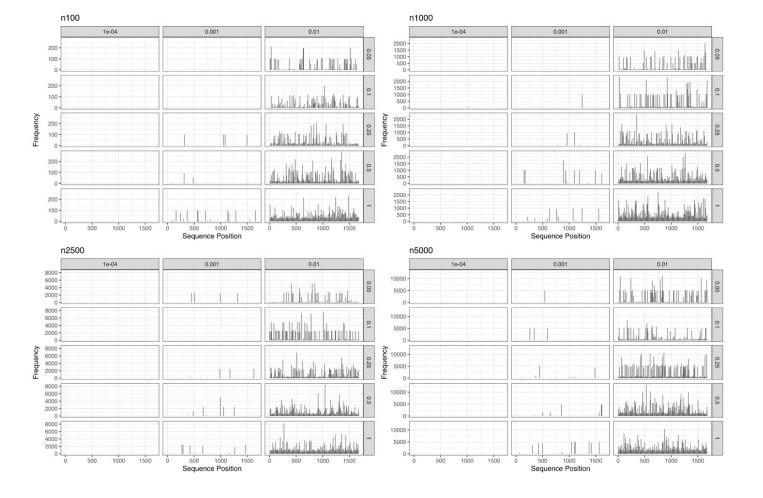
Mutation Rate

Simulated Breakpoints

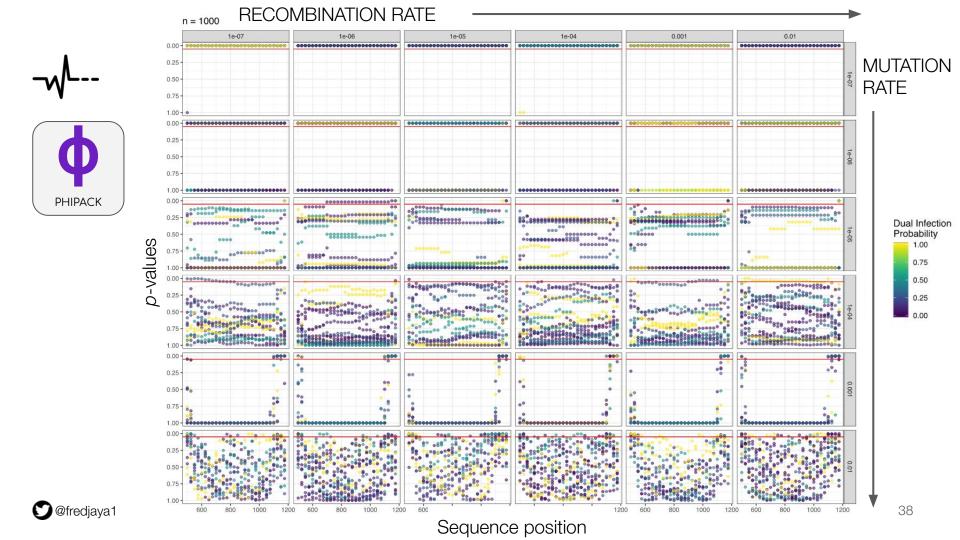




Dual infection probability







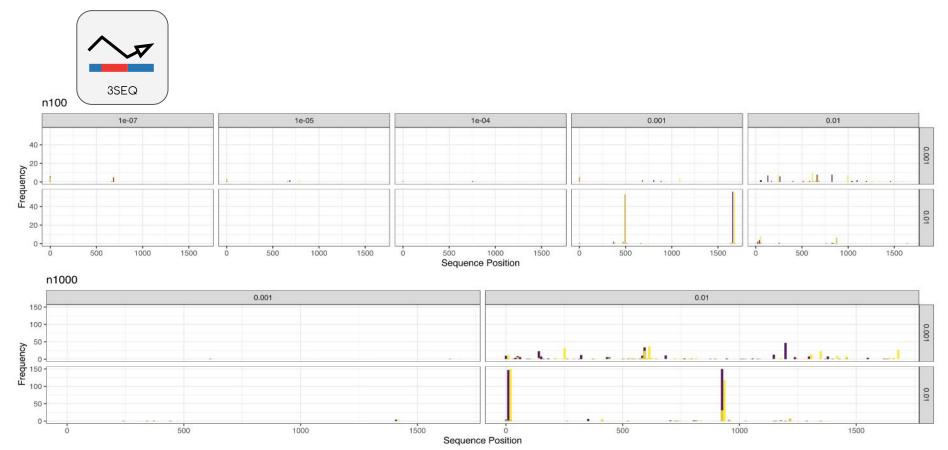


Figure 7. Frequency distribution of significant, paired breakpoint locations (start - purple, end - yellow) detected by 3SEQ from sim_2 . Horizontal facets denote recombination rate and vertical facets denote mutation rate.

fredjaya1

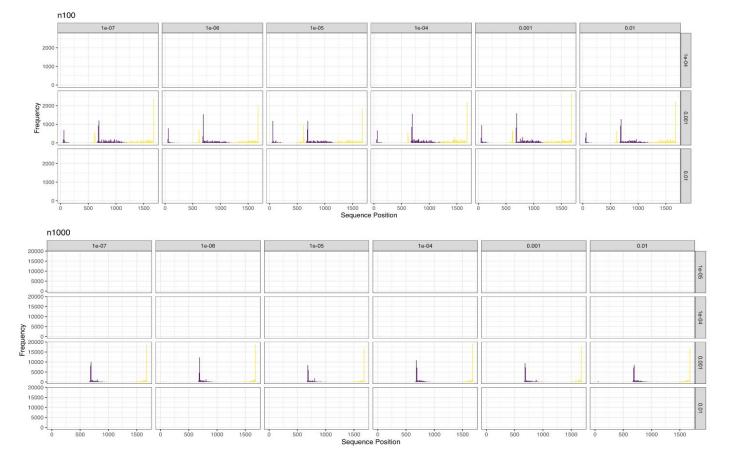


Figure 6. Frequency distribution of significant, paired breakpoint locations (start - purple, end - yellow) detected by GENECONV from sim_2 . Horizontal facets denote recombination rate and vertical facets denote mutation rate.

② @fredjaya1