

IgDiscover analysis of adaptive immune receptor germline gene diversity

Martin Corcoran

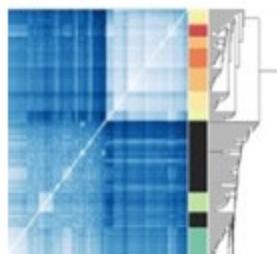
Gunilla Karlsson Hedestam Group
Karolinska Institute

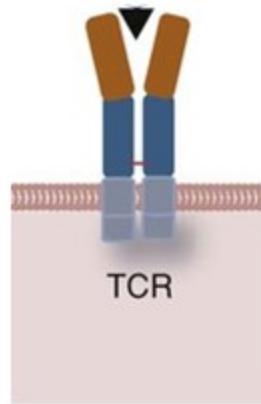
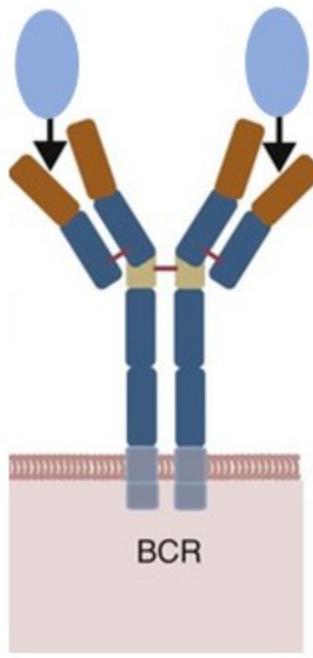
Germline Inference

Computational Identification of germline Ig or TCR sequences from expressed Rep-Seq data

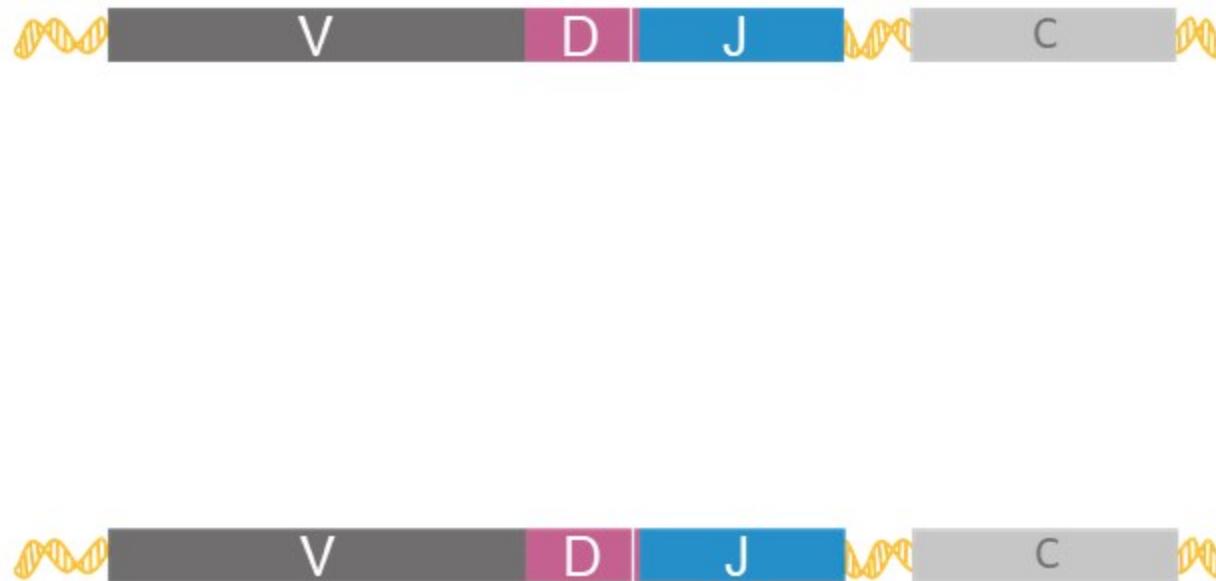
Isolation of naïve germline sequences that are used in multiple independent antibody or TCR rearrangements

IgDiscover

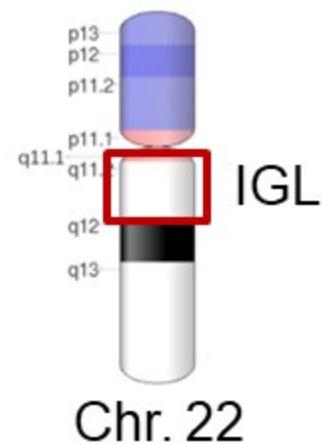
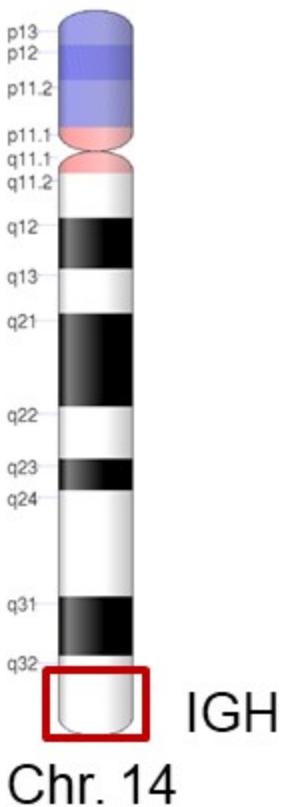
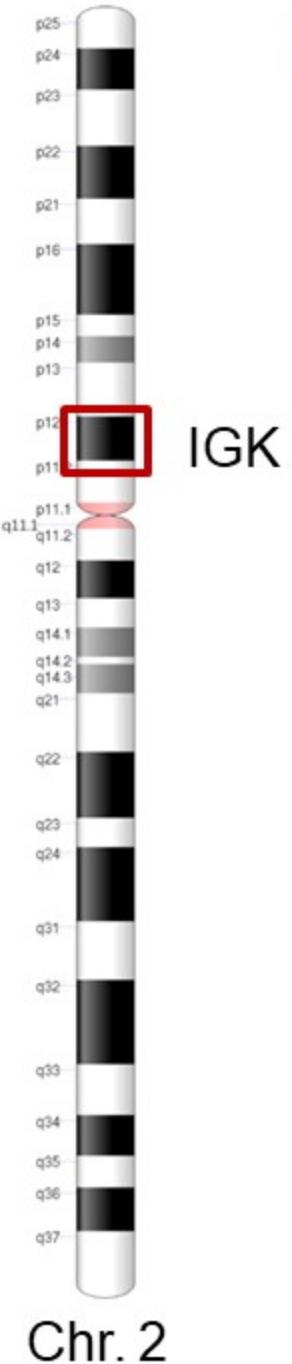




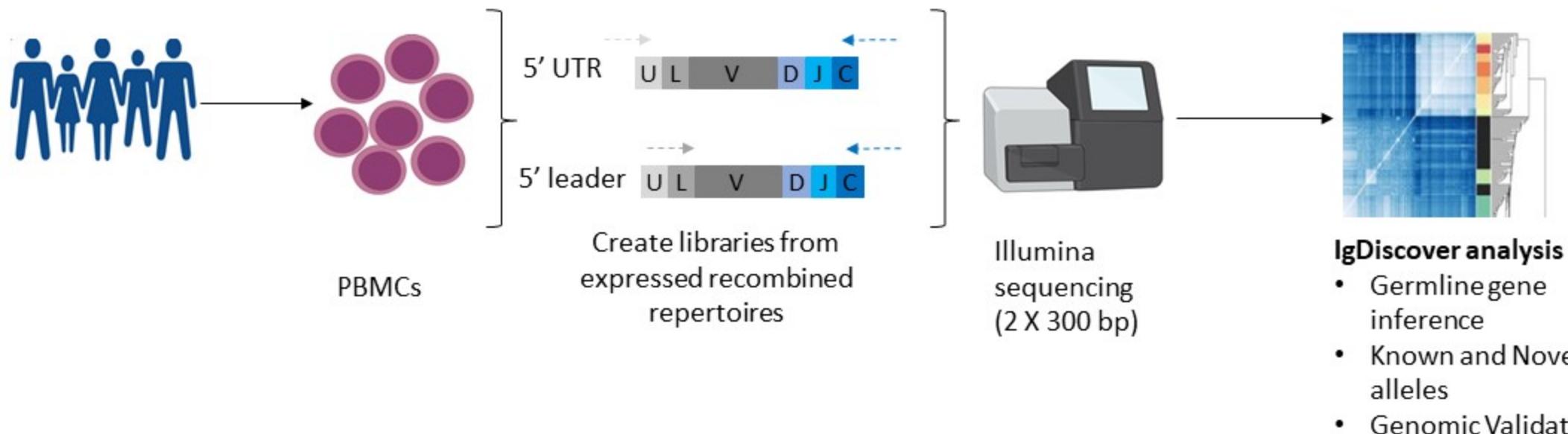
Adaptive immune receptor loci are among the most polymorphic in the human genome



Genomic localization



IgDiscover library production

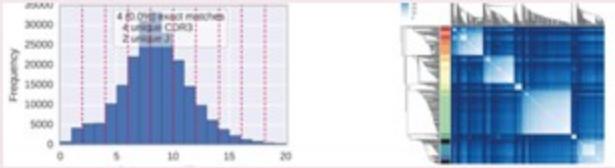


IgDiscover analysis overview

High frequency alleles

Rep-Seq Library > MiSeq 2 x 300 bp

- Assignment of sequences to closest reference database sequence
- Windowed and linkage clustering



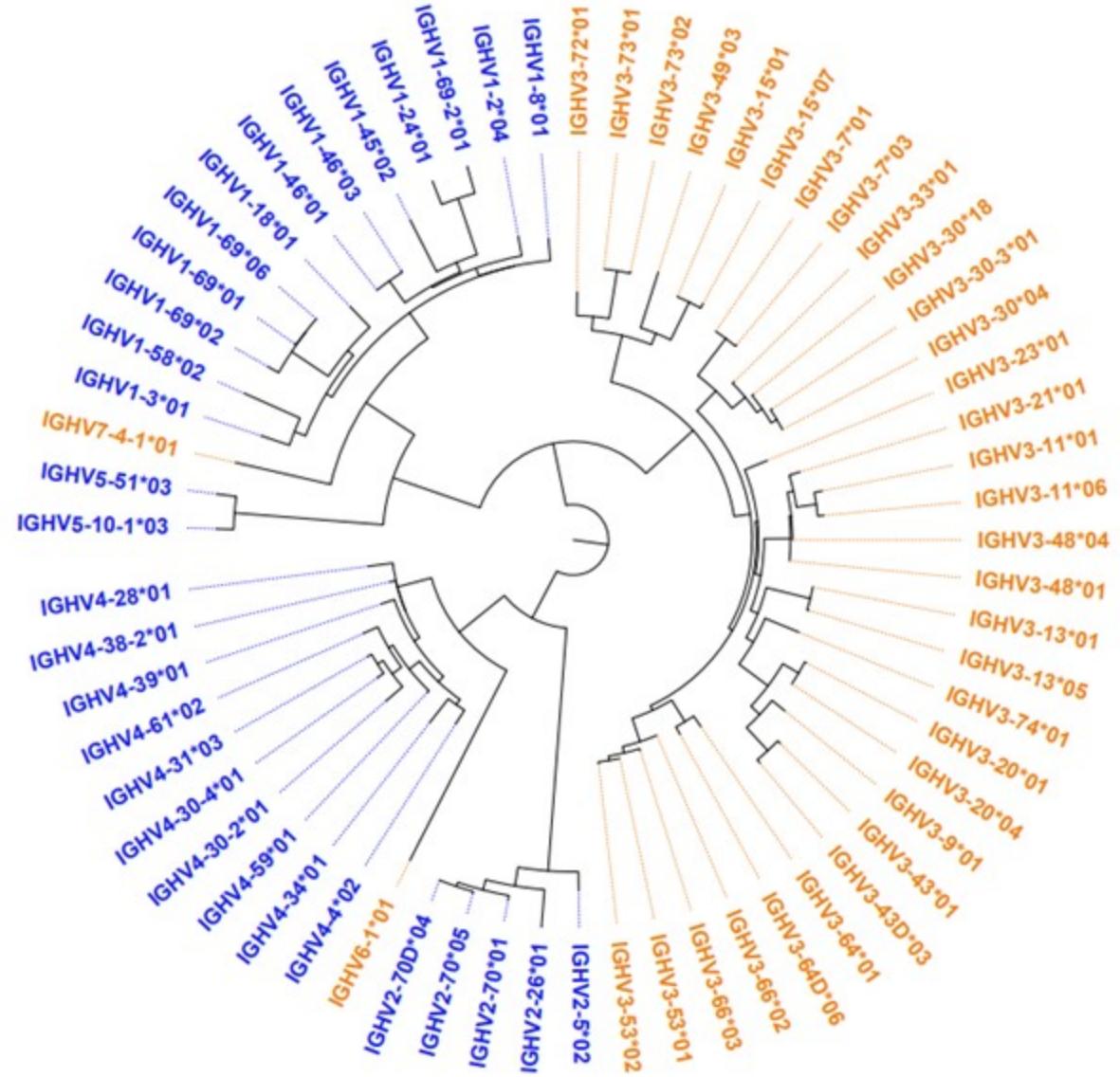
Low frequency alleles

- Further iterations, replacement of previous database, whitelisting

V-gene "Starting" Database

Replacement Database Iteration *n*
Individualized database

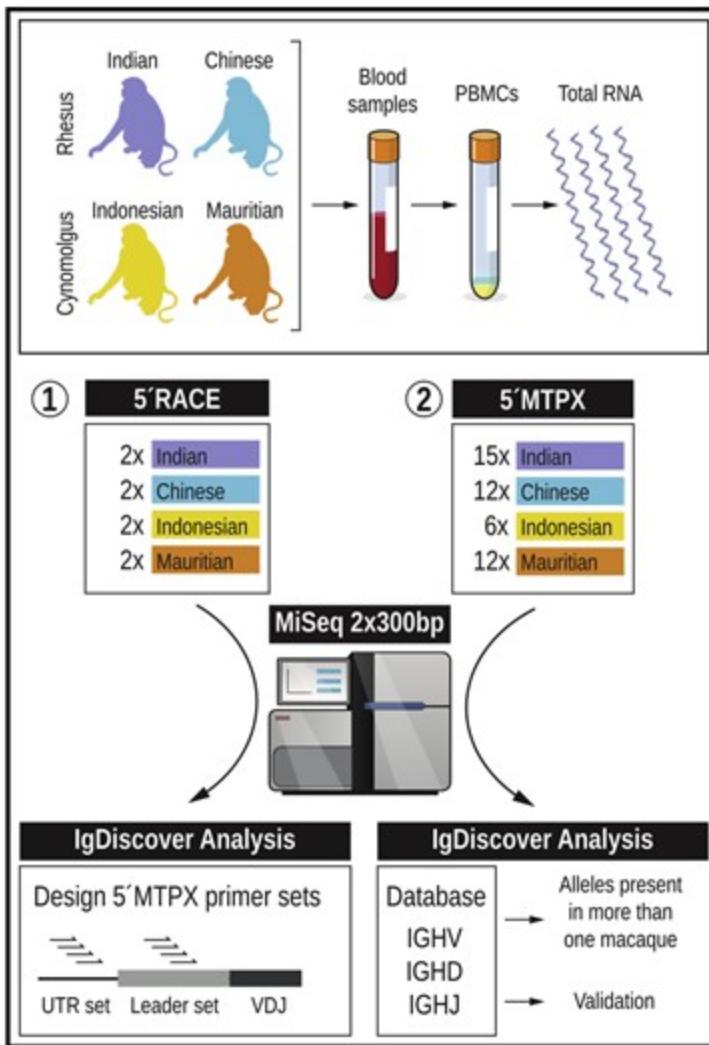
Individualized IGHV Genotype



Individualized database usage

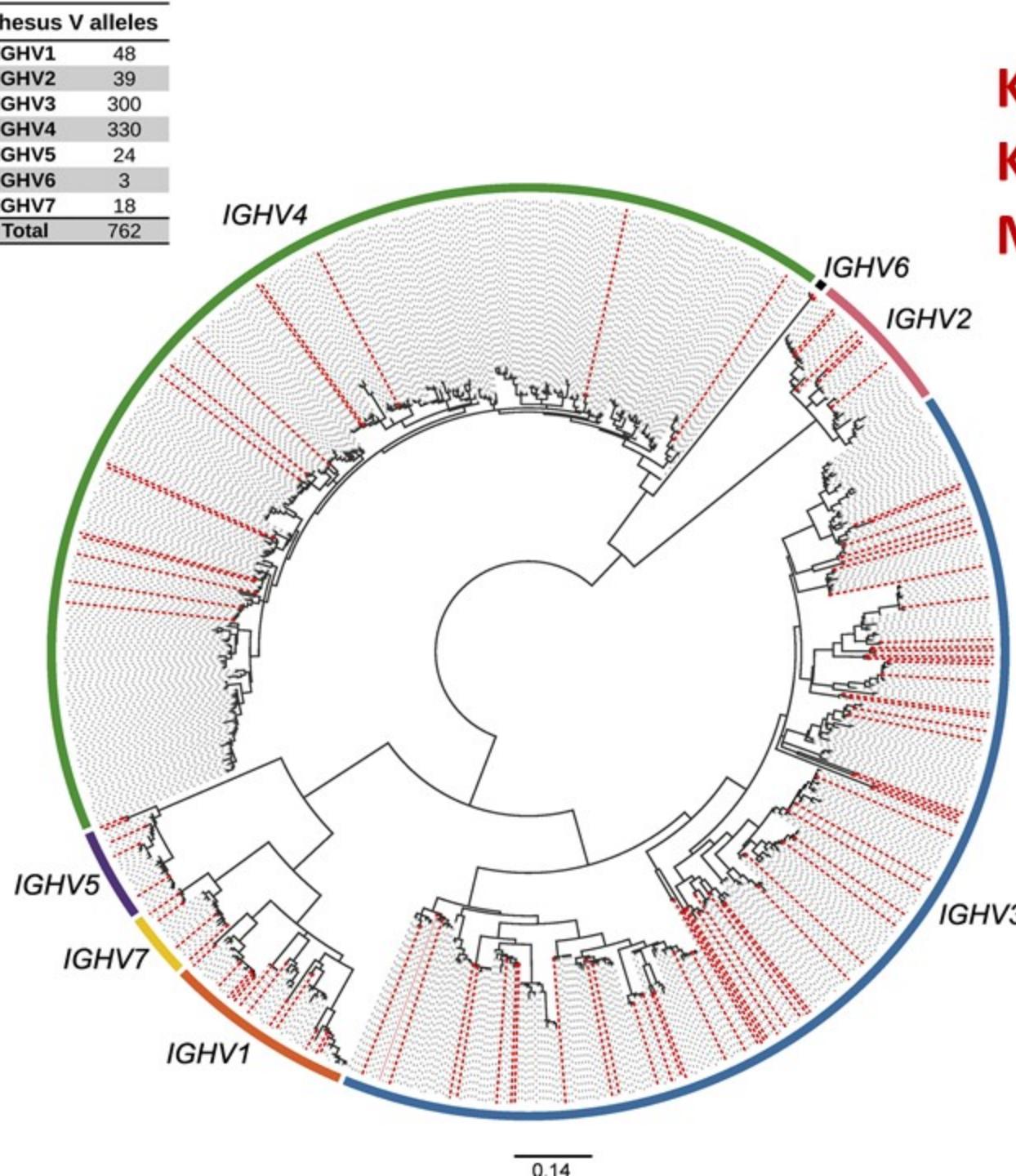
- Accuracy of SHM estimation
- Lineage tracing of diverse antibodies
- Identification of gene/allele content to enable association studies

Using Germline Inferral to assemble a Rep-Seq suitable database



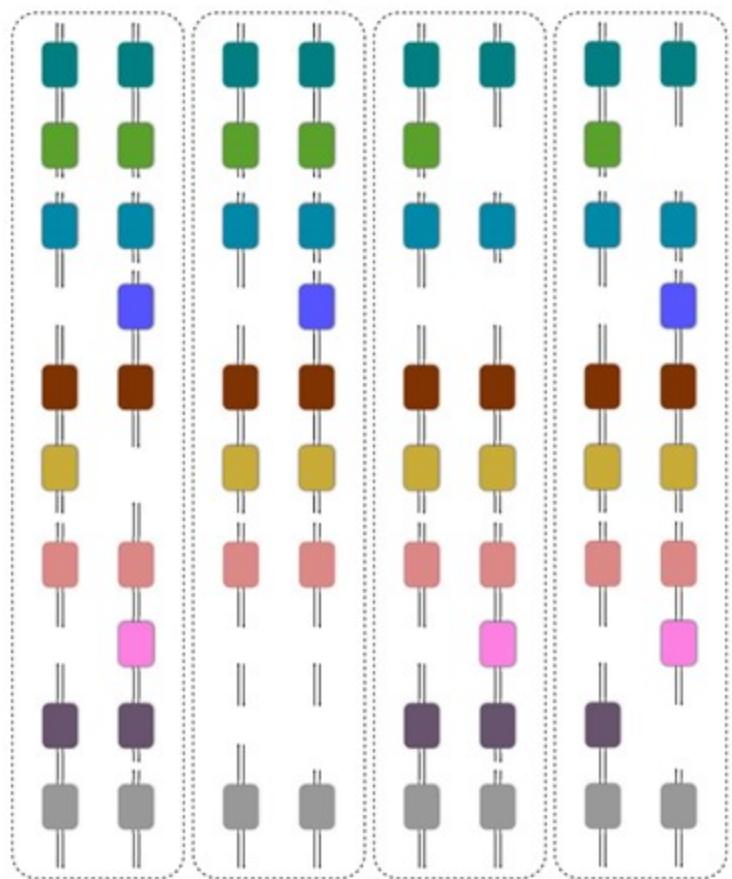
Vázquez Bernat N, et al. Rhesus and cynomolgus macaque immunoglobulin heavy-chain genotyping yields comprehensive databases of germline VDJ alleles.
Immunity. 2021

KIMDB
Karolinska Institute
Macaque database

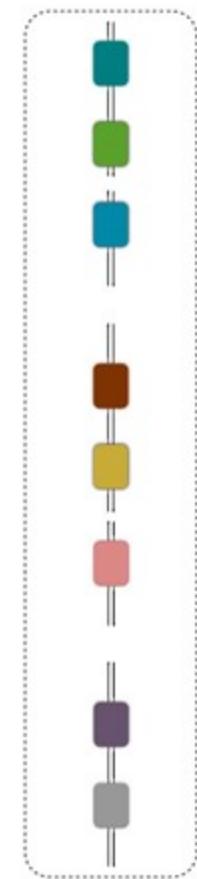


0.14

Multiple animals



Haploid assembly



Database

Database

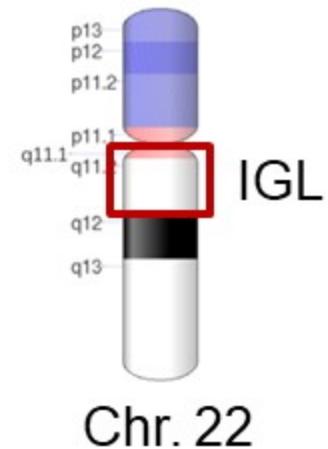
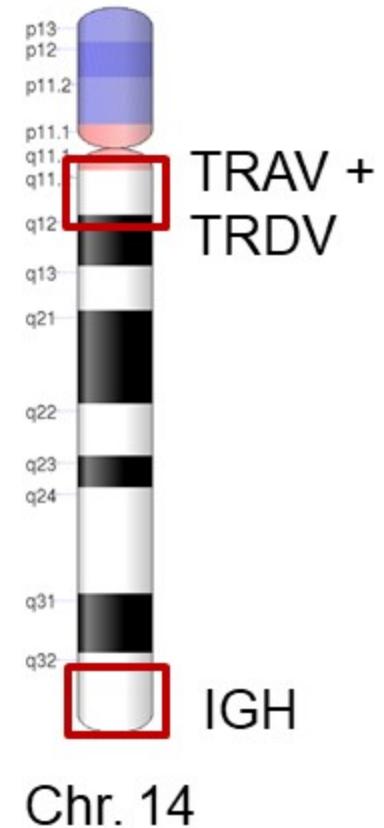
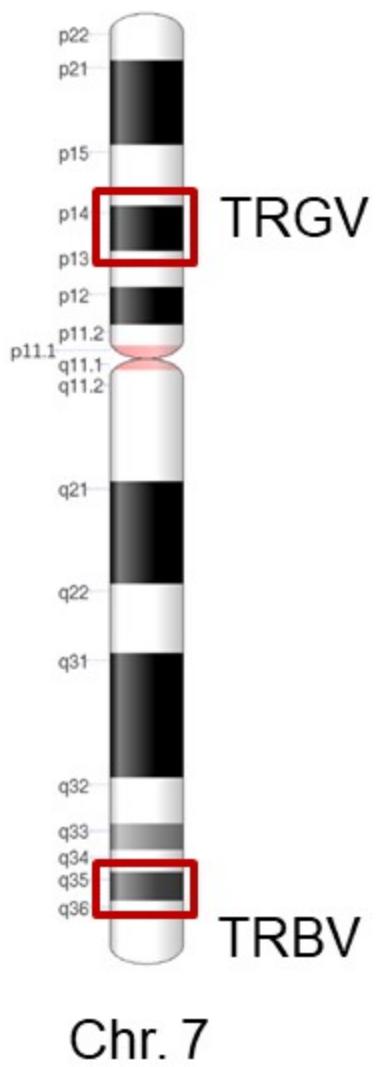
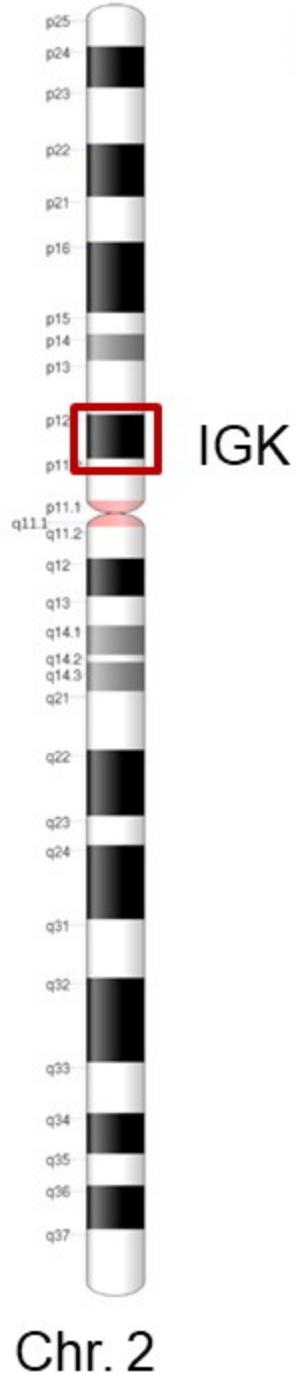
Inferred germline database

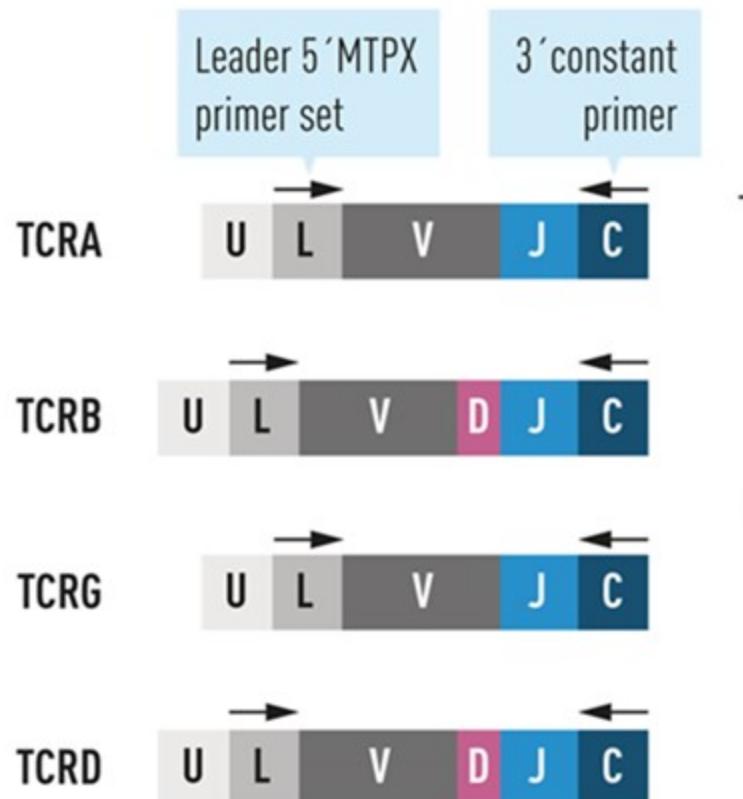
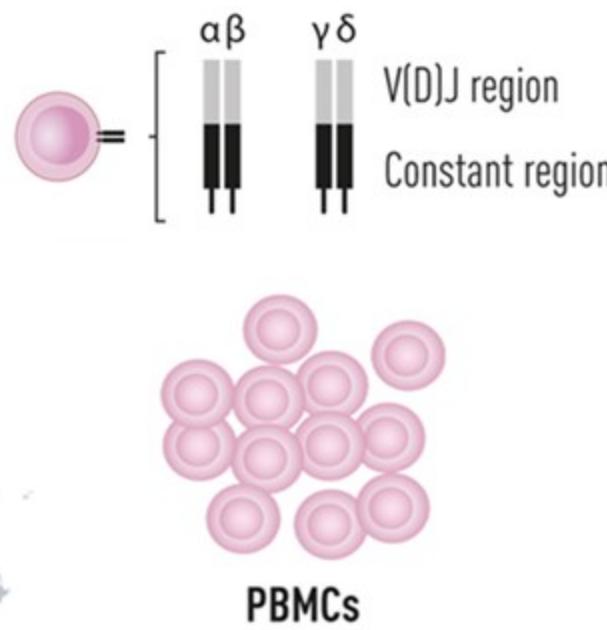
- Highly comprehensive in terms of germline sequences
- Problems with gene nomenclature remain (non-human species)

IgDiscover updated version

- MIARR compliant
- TCR discovery
- Improved haplotyping facility
- Genotype module – core count
- J and D discovery

Genomic localization



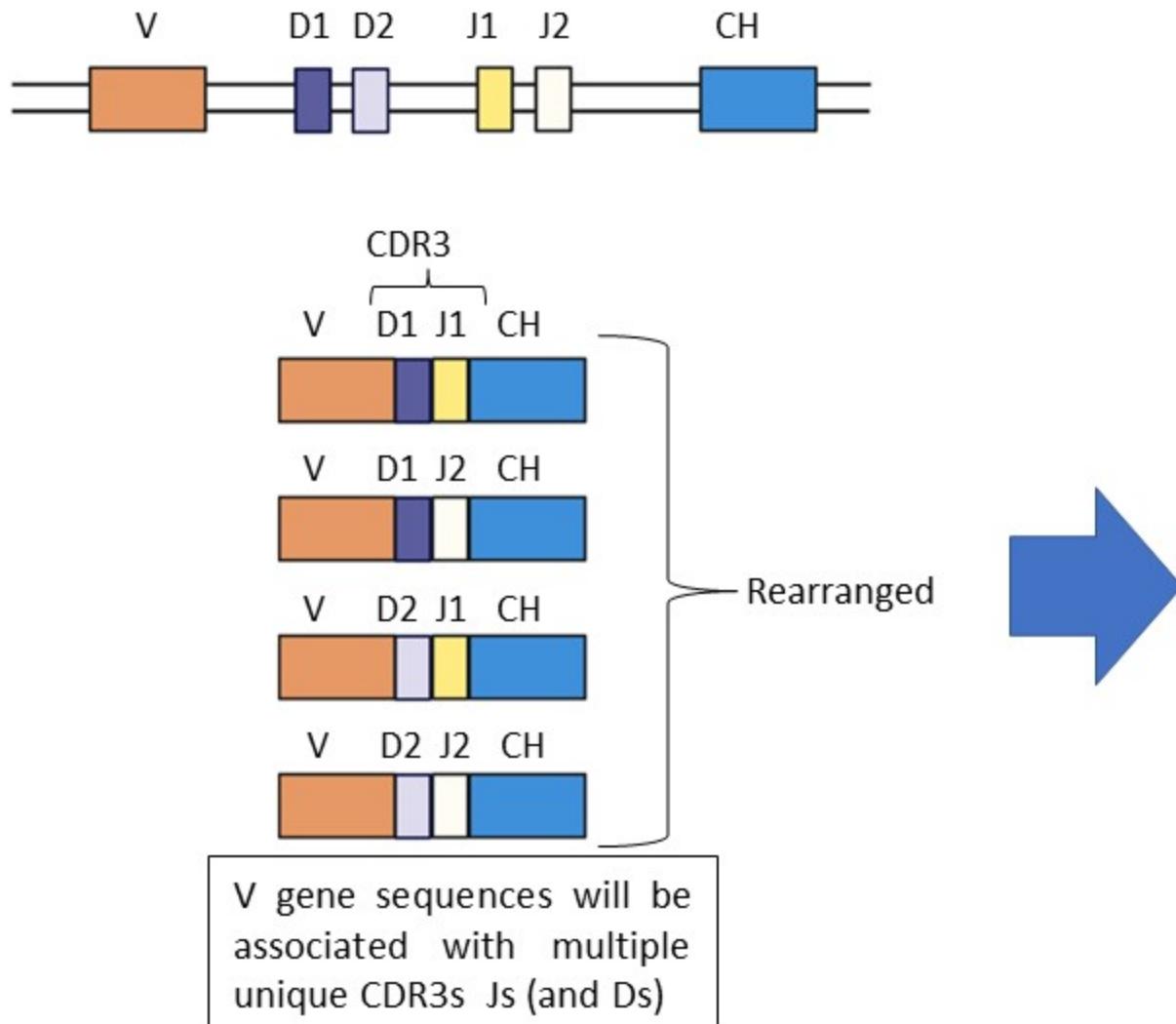


Computational analysis

Germline gene inference, known and novel alleles

Core count, genotype validation

TCR Germline Discovery



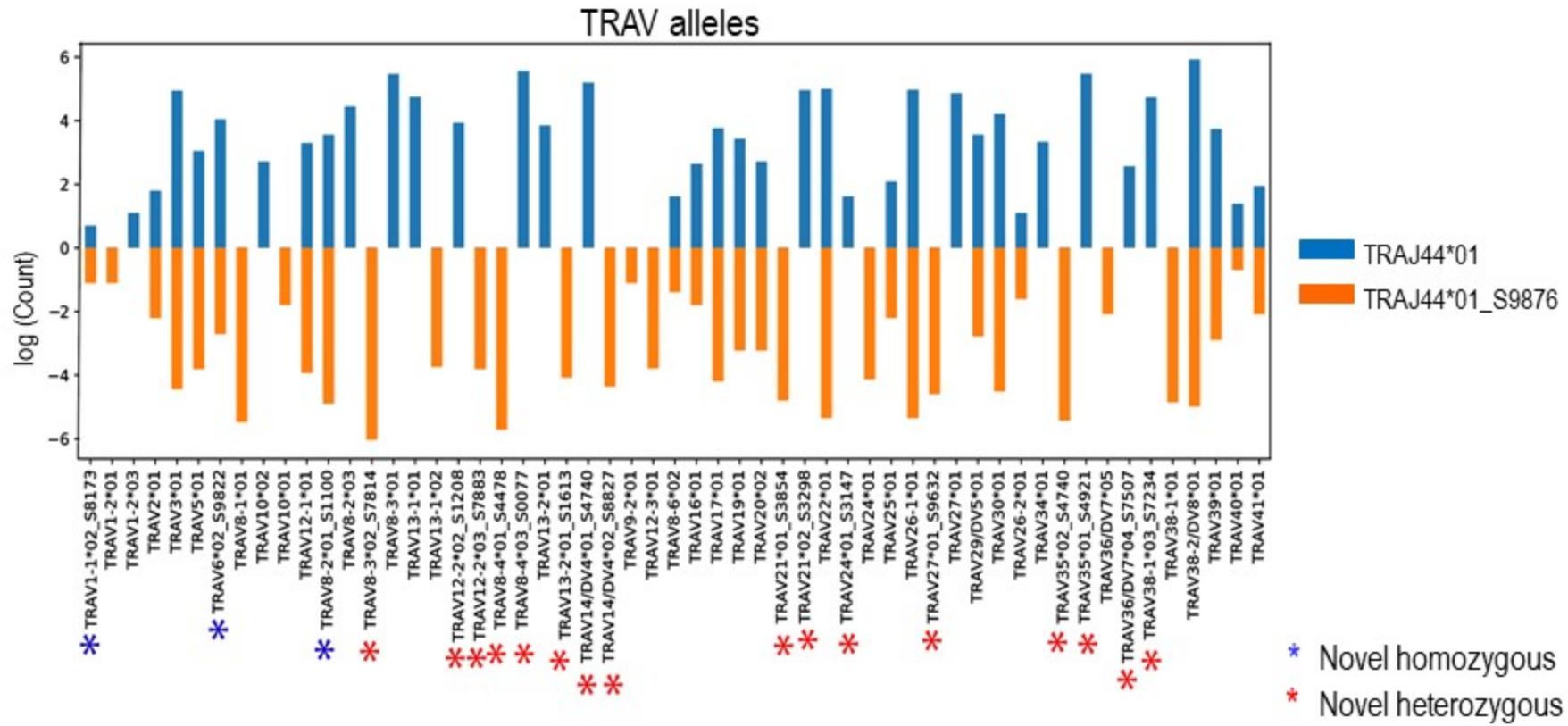
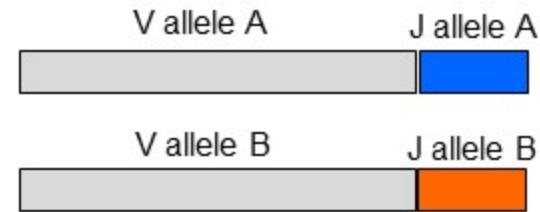
TCR specific

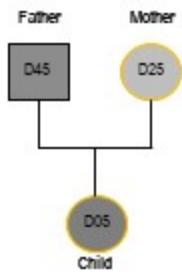
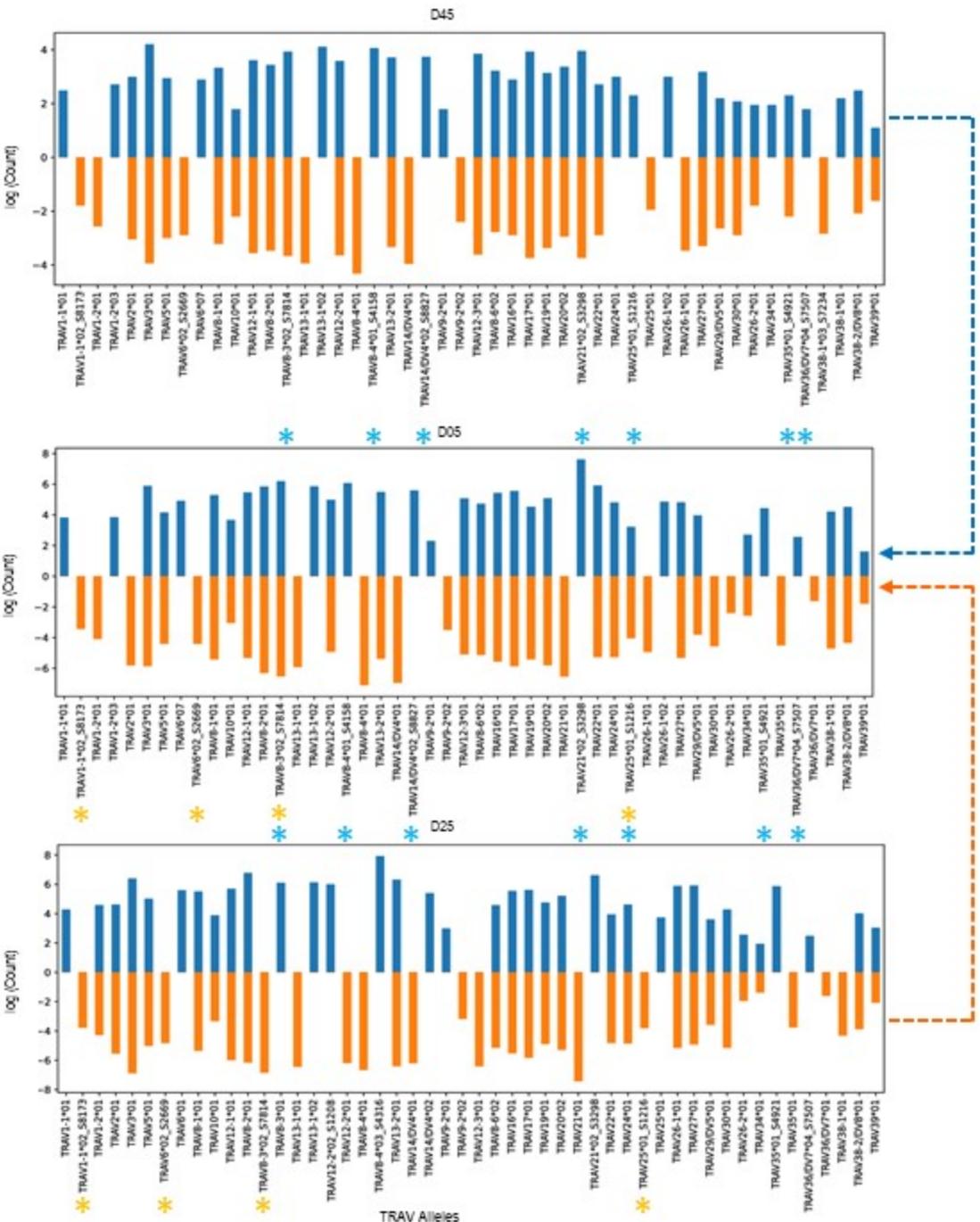
Validation

- PCR and Sanger sequencing
- Haplotype analysis
- corecount genotype analysis

Validation by inferred haplotype analysis

V allele status can be determined through linkage to heterozygous J alleles

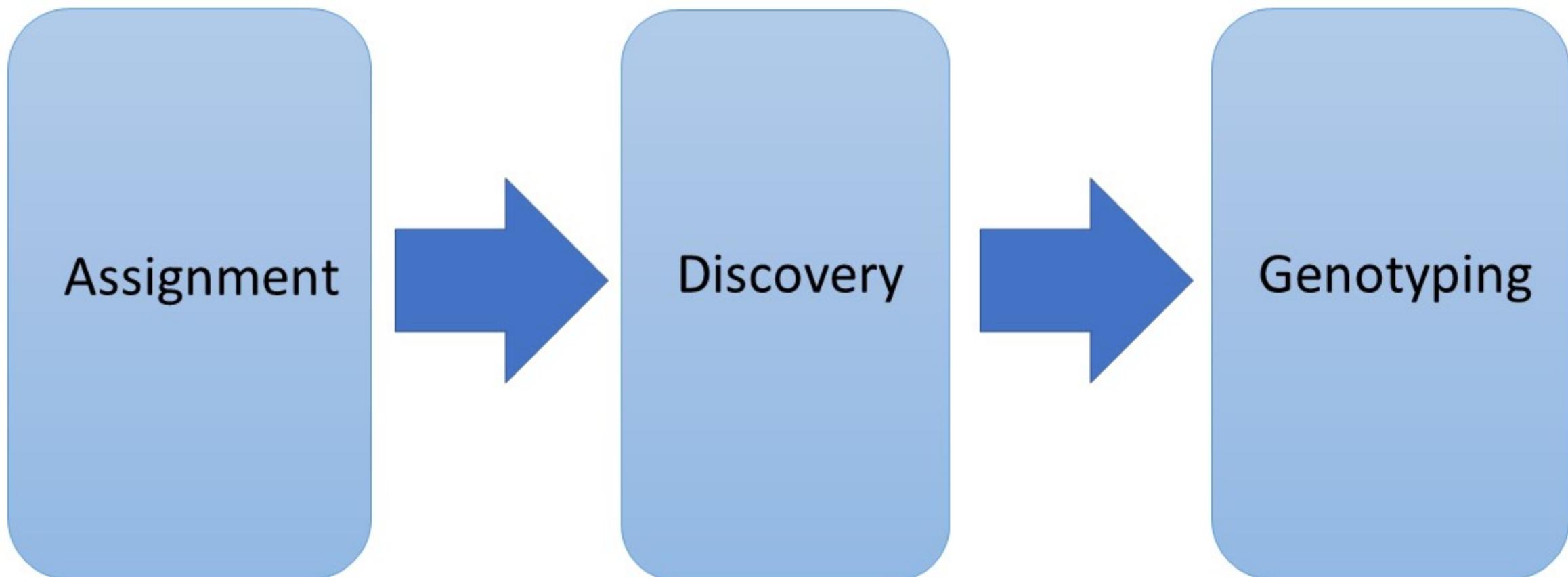




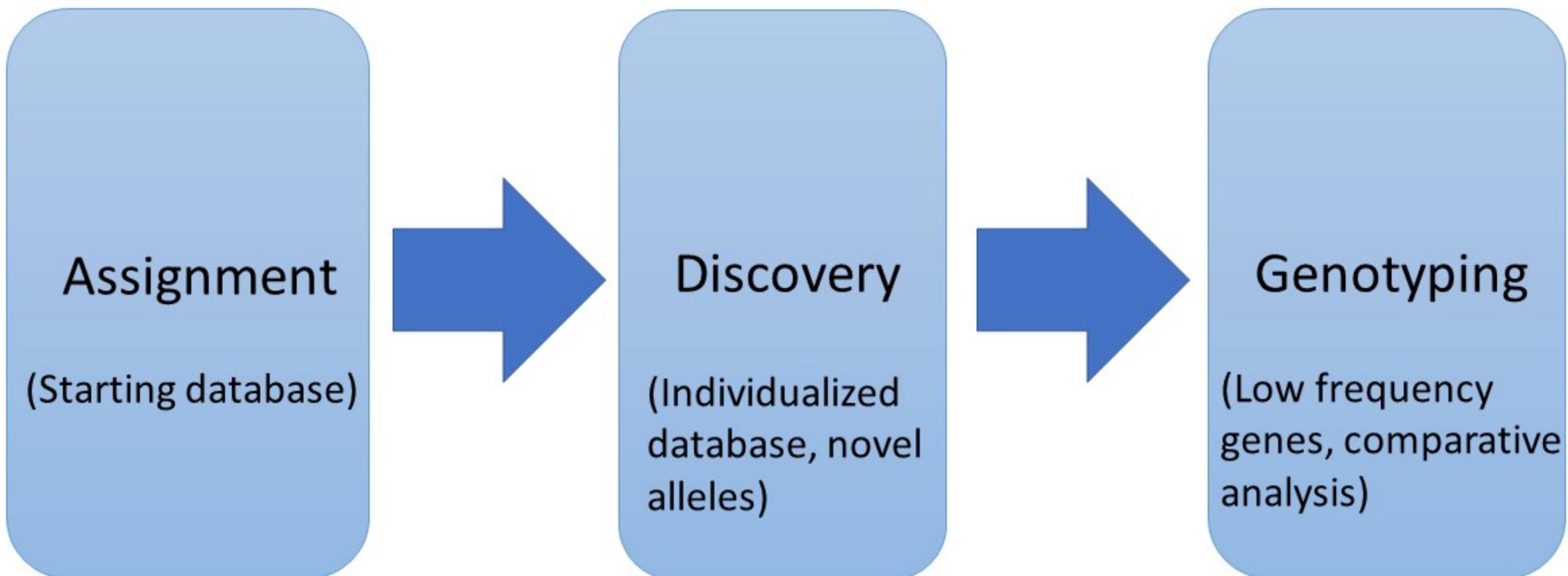
Familial analysis

- 9 novel alleles inherited
- 7 novel alleles from father
- 4 novel alleles from mother
- 2 novel alleles from both

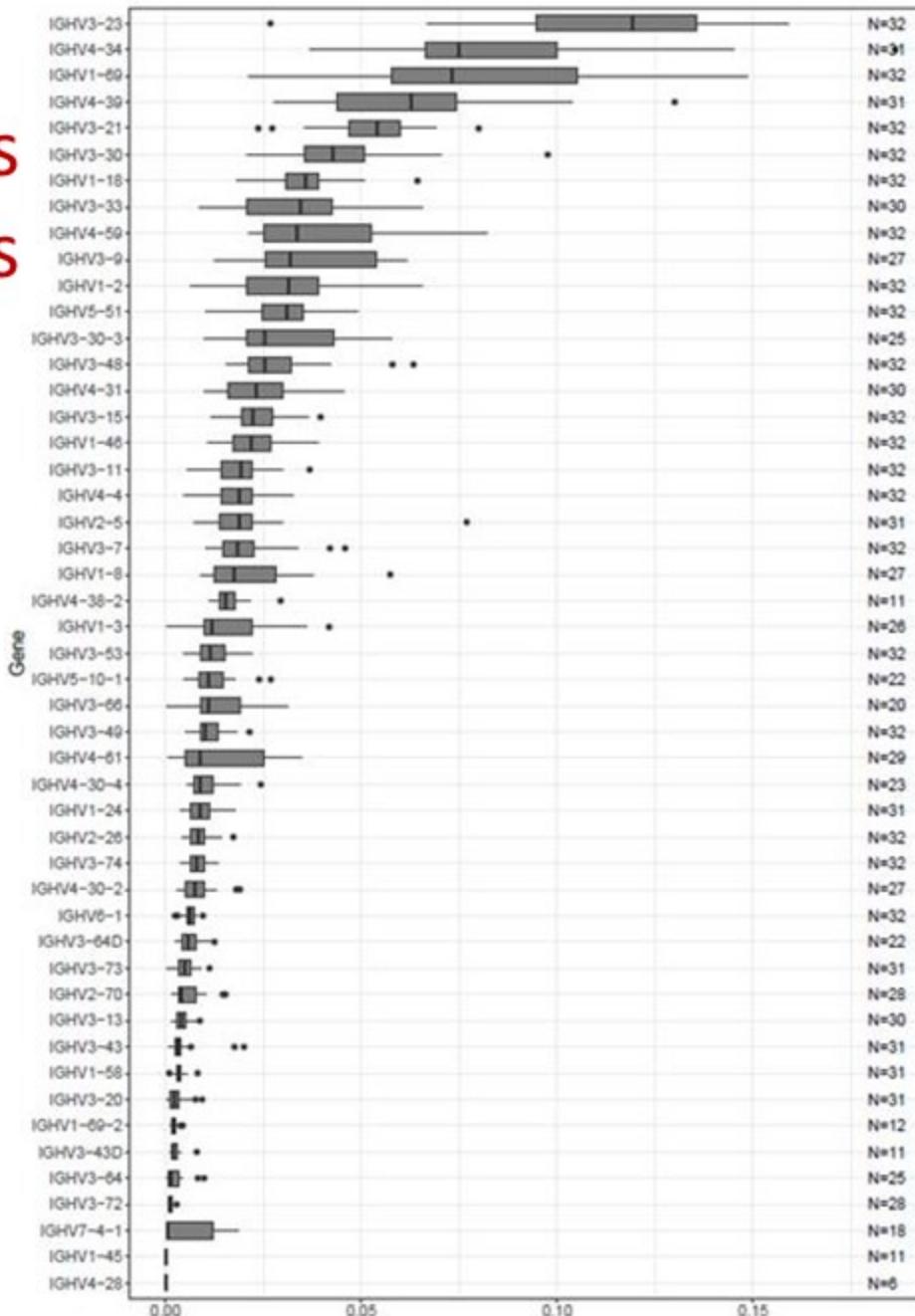
Genotype comparison corecount analysis module



Genotype comparison corecount analysis module



Unmutated sequences in human IgM libraries



Genotype comparison corecount analysis module



- Most recombined VDJ sequences will not contain full length Vs, Ds or Js
- They will contain a 'core' of the appropriate V, D and J allele
- Using a modified 'core' database allows highly accurate identification of even low frequency genes and alleles from the filtered.tab output of IgDiscover.
- core sequence = full length allele minus nucleotides frequency lost during recombination

Reference database sequence

```
>IGHD2-21*01  
AGCATATTGTGGTGGTGATTGCTATTCC  
>IGHD2-21*02  
AGCATATTGTGGTGGTGACTGCTATTCC
```

Corecount database sequence

```
>IGHD2-21*01  
CATATTGTGGTGGTGATTGCTATT  
>IGHD2-21*02  
CATATTGTGGTGGTGACTGCTATT
```

Genotype comparison corecount analysis module

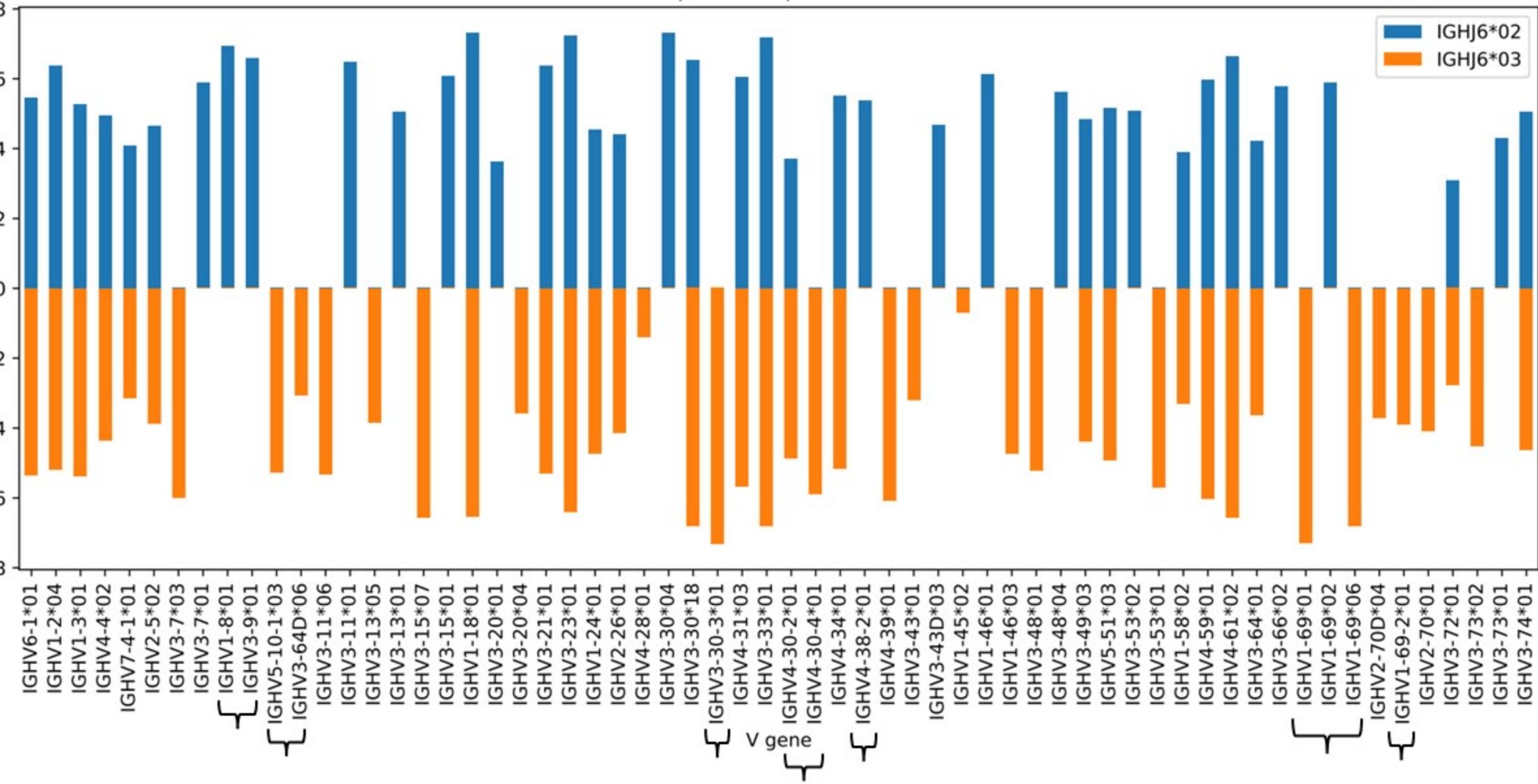
Requirements

- Comprehensive database
- Rep-Seq library with sufficient unmutated sequences (IgM, IgD, IgK, IgL, TRA, TRB, TRD, TRG)
- High quality sequence library (MiSeq, 2 x 300 V3, avoid high clustering density)

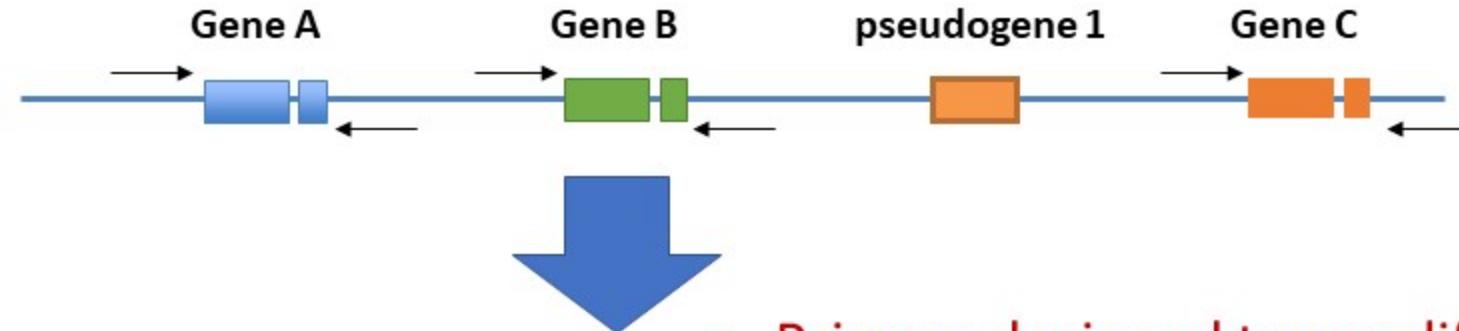
Testing the accuracy of the corecount genotyping

- Analysis of subject of "known" genotype
- High gene content individual

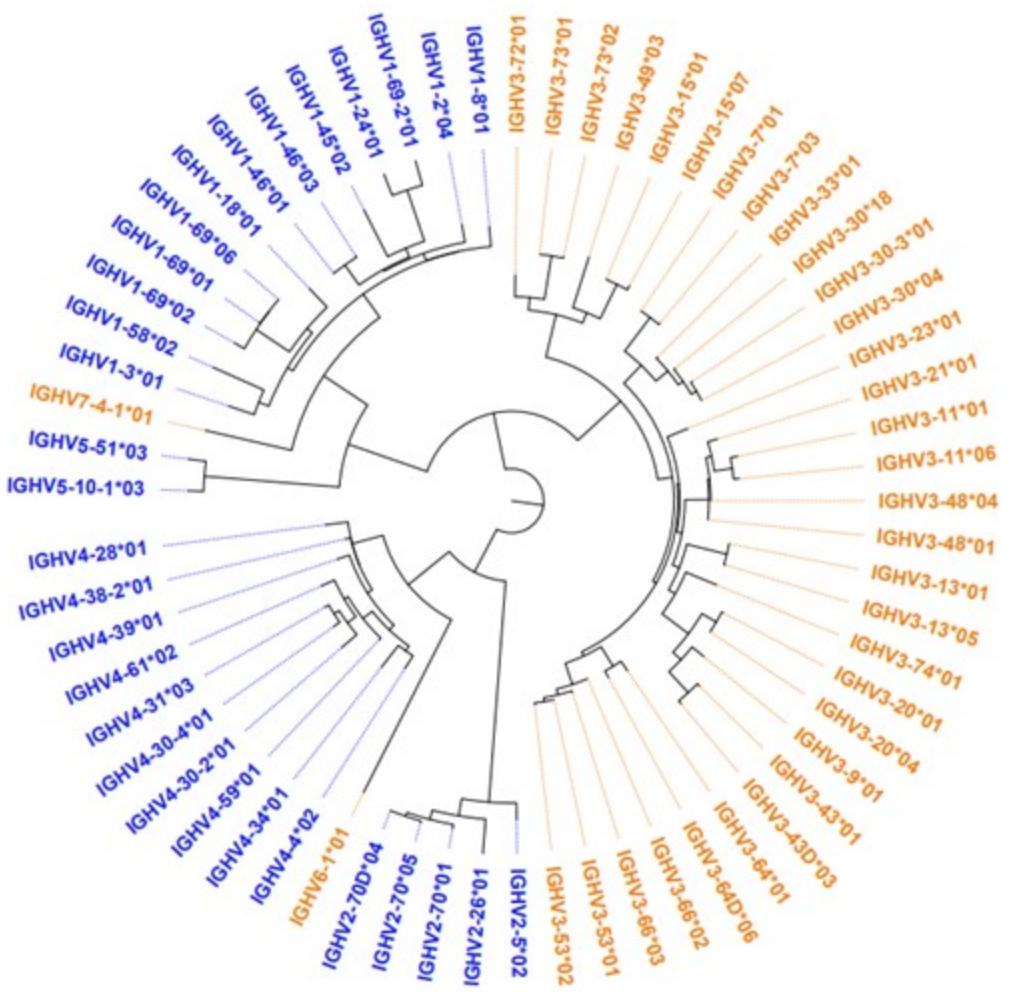
Allele-specific expression counts



Genomic validation of all expressed sequences

- 
- Primers designed to amplify all expressed genes
 - Products cloned into plasmids and Sanger sequenced
 - All expressed V genes, D genes and J genes targeted
- 50 IGHV genes, 64 alleles
- 27 IGHD genes
- 6 IGHJ genes, 7 alleles

V (50 genes, 64 alleles)



- Sanger sequenced genotype showed 100% identity to the corecount genotype output

Corecount genotyping summary

- Performed on filtered output of IgDiscover
- Batch analysis is **rapid**
- **Filtering is simplified** (allelic ratio, expected frequency, CDR3 diversity)
- **Highly accurate in genotyping low frequency genes**
- Can **discriminate between alleles with 3' end variation**

Conclusion

- Diversity is extensive in outbred species
- Involves both allele and gene content
- Combination of Rep-Seq and genomic methods help define diversity
- Databases will become more comprehensive and accurate
- High throughput, rapid and consistent methods will enable future association studies

Acknowledgements

Karlsson Hedestam lab

Sanjana Narang

Mateusz Kaduk

Mark Chernyshev

Marco Mandolesi

Christopher Sundling
Anna Färnert



**Karolinska
Institutet**



Mateusz Kaduk

Sanjana Narang



Mark
Chernyshev

Marco
Mandolesi



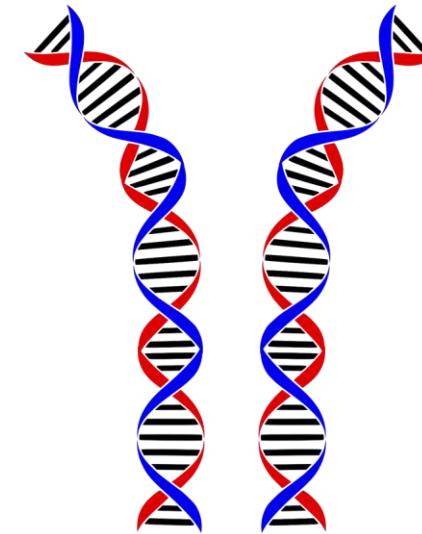
THANK YOU!

Revealing Ig/TR germline variations by AIRR-seq analysis

VDJbase



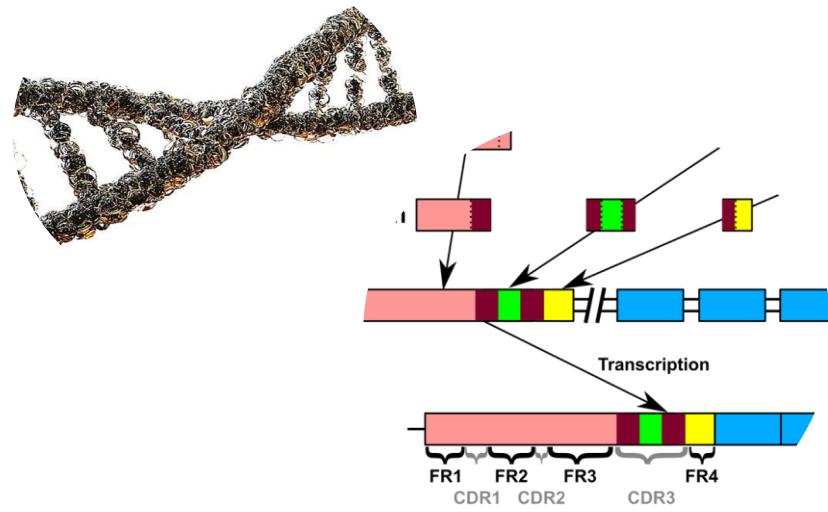
RAbHIT



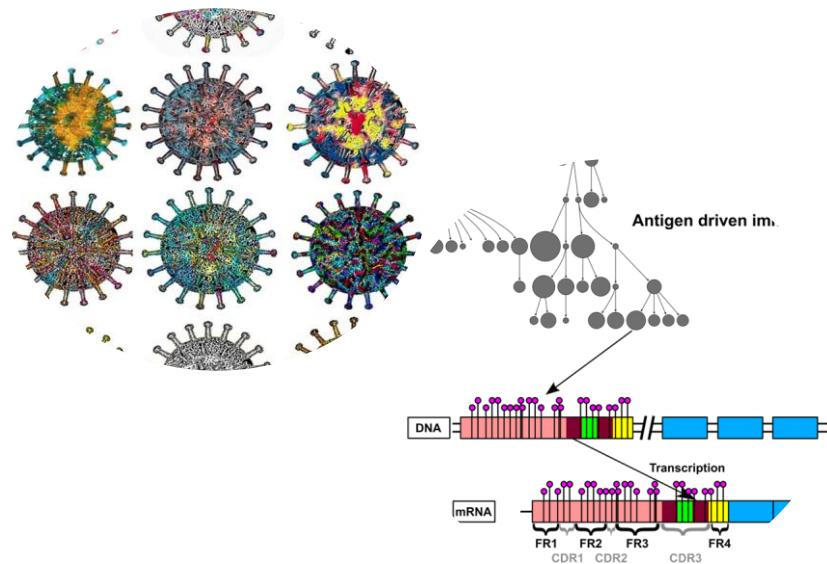
Ayelet Peres
Prof. Gur Yaari
Bar Ilan university

Adaptive immune system

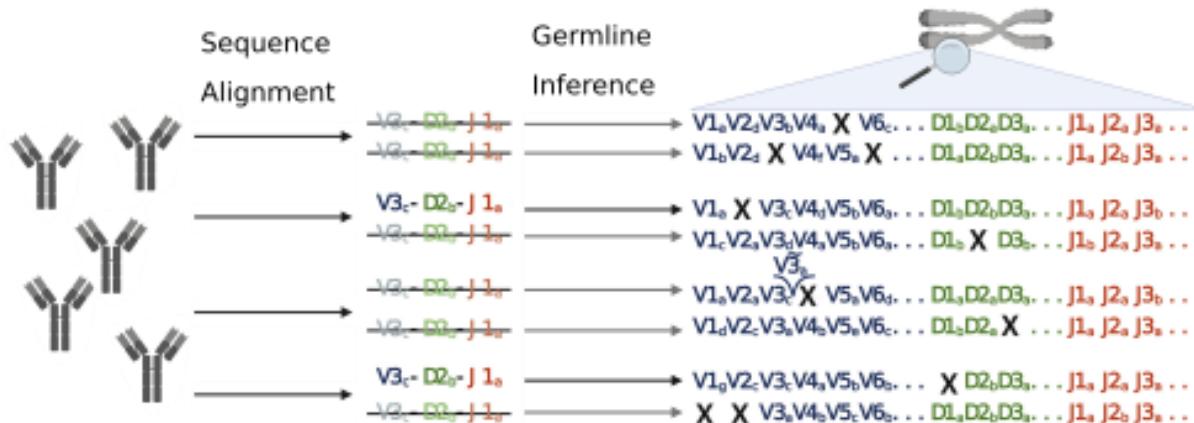
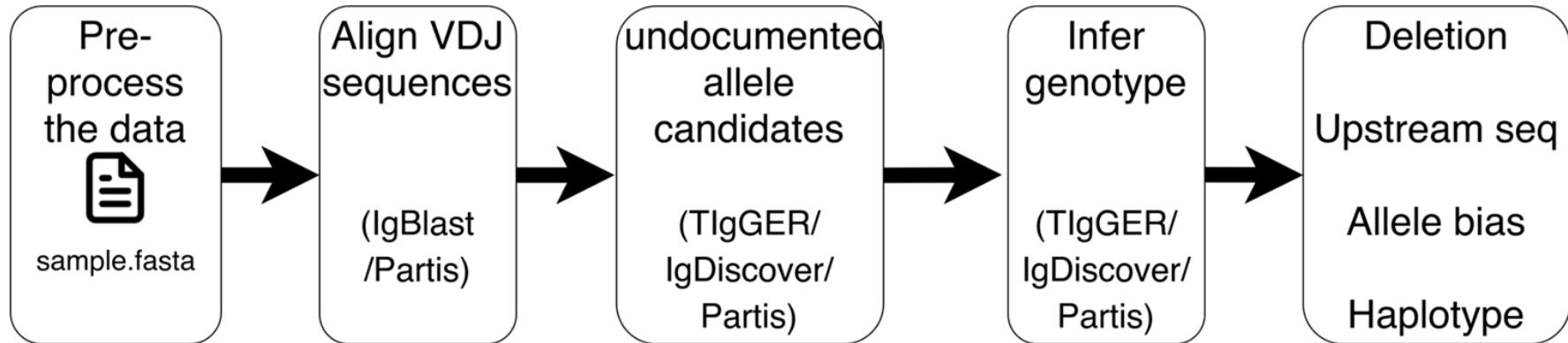
Genetics



Environment



Repertoire inference steps



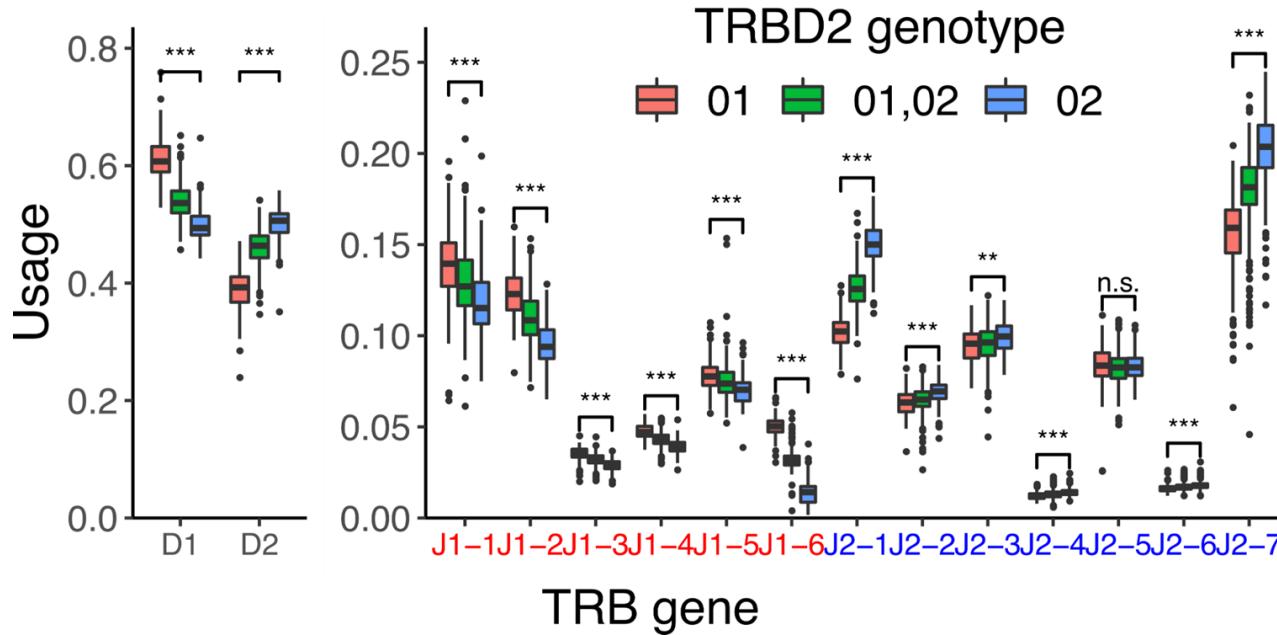
Undocumented allele inference

Chain	Known Alleles	Undocumented discovered	
IGH	287	25	Mikocziova at el. 2020, NAR
IGK	64	25	Mikocziova and Peres at el. 2021, Iscience
IGL	75	23	
TRB	64	38	Omer and Peres at el. 2021, genome medicine

Genotype influence [Omer* and Peres* et al., 2021 *Genome Medicine*]

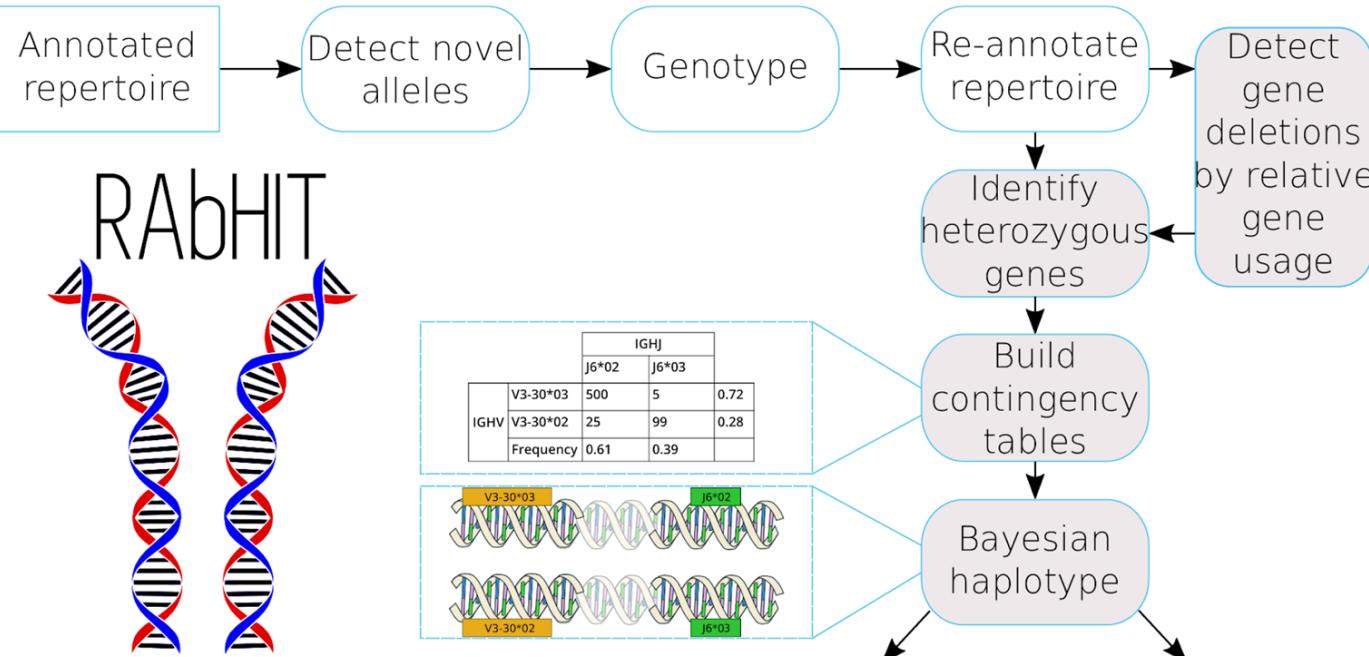
TRBD2*01
TRBD2*02

GGGACTAGCGGGGGGG
GGGACTAGCGGGAGGG
***** · ***

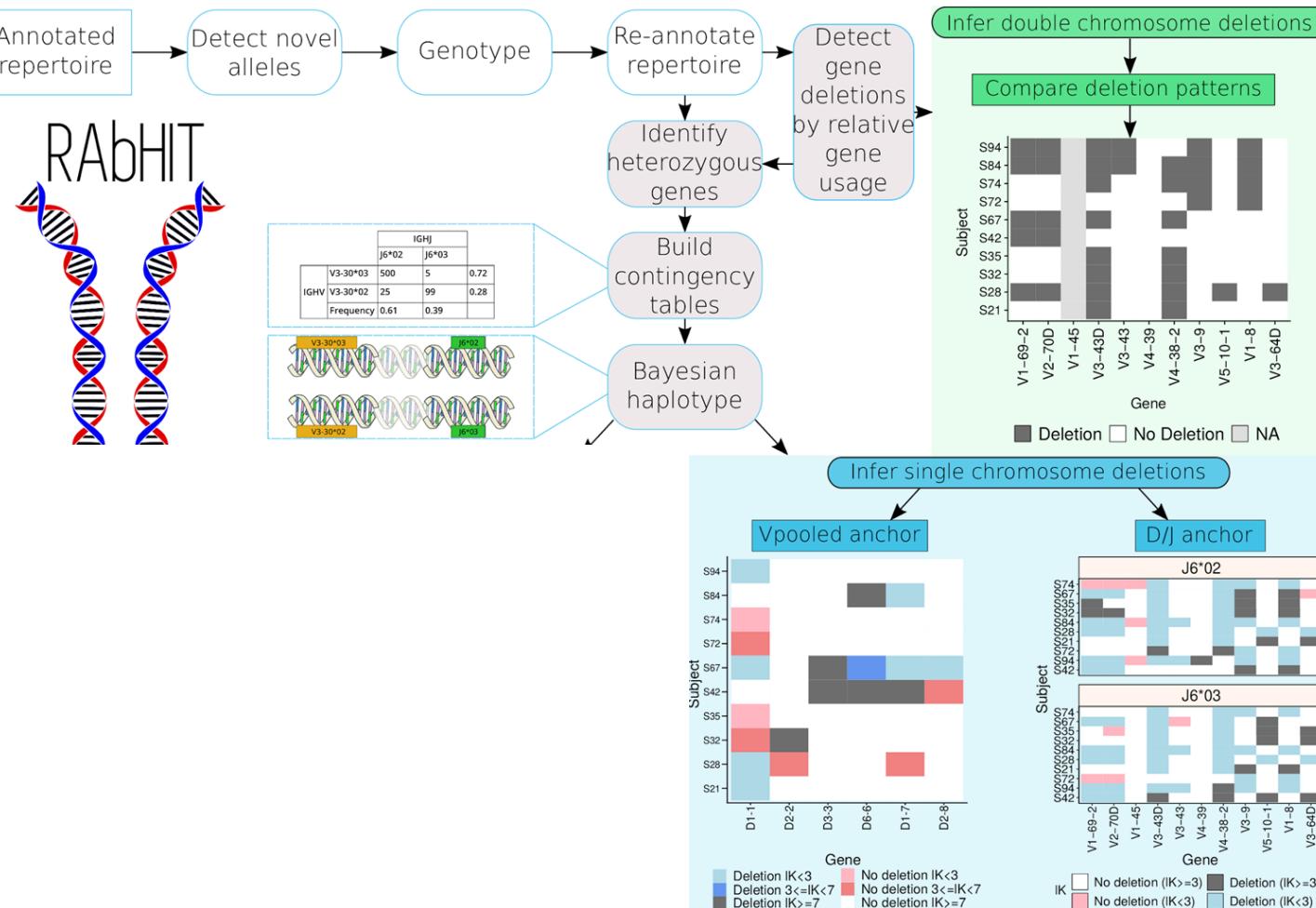


Strong bias between TRBD2 genotype and TRBJ1-6 usage

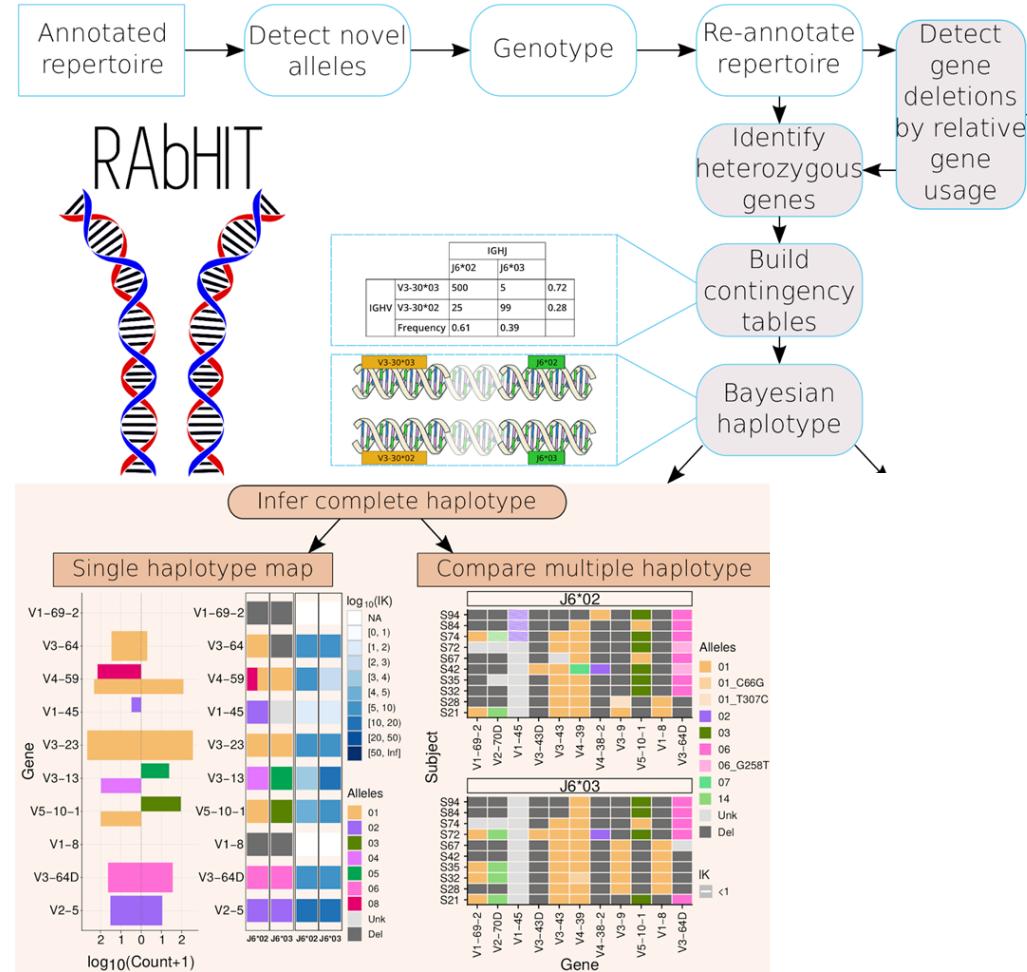
RAbHIT: R Antibody Haplotype Inference Tool [Peres et al., 2019 *Bioinformatics*]



RAbHIT: R Antibody Haplotype Inference Tool [Peres et al., 2019 *Bioinformatics*]



RAbHIT: R Antibody Haplotype Inference Tool [Peres et al., 2019 *Bioinformatics*]



Haplotype inference

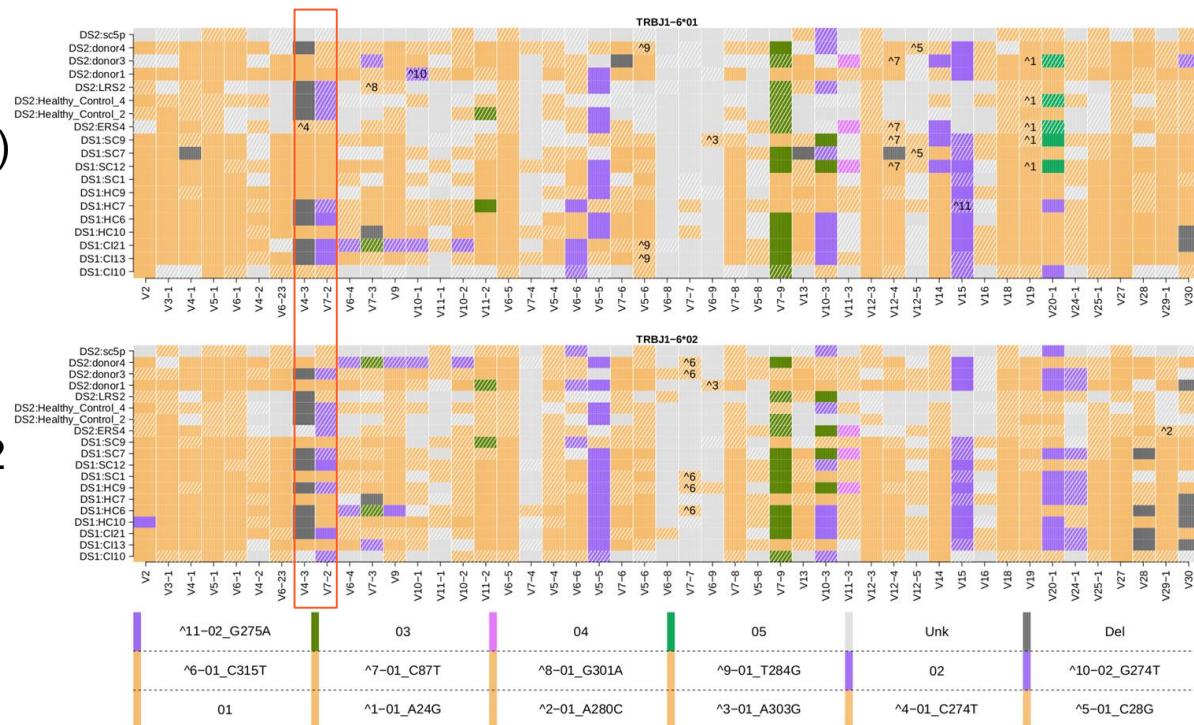
- J genes as anchors (TRBJ1-6)

- Allele patterns and deletion

associations

- Linkage between TRBV7-2*02

and TRBV4-3 deletion



VDJbase.org [Omer*, Shemesh*, Peres* et al., 2020 NAR]

- Publicly available Ig and TRB database search engine
- Describing complete sets of genotype and haplotype data
- Gene usage data

- Available for local installation
- Allows you to privately add your own datasets
- VDJbase processing pipeline is available in dockerhub [peresay/suite](#)

Species Human Genomic Sets Select Datasets AIRR-seq Sets IGH x

VDJbase Reports

Currently available reports run against AIRR-Seq data sets. Reports that run against genomic datasets, or that combine data of both types, will be added in due course.

Please select above the species and datasets on which you would like to report. Multiple datasets from the same species can be selected.

You can refine the samples against which the reports will run by using filters in the Samples windows to display just those samples that you wish to be included.

To run the report, click on an icon in the Run column corresponding to the output format that you would like (e.g. pdf, on-screen, Excel)

[Click here](#) for explanation of reports.

Reports using just AIRR-seq Sample Data

Thumbnail	Name	Description	Run
	Genotype List	Heatmap-style report of the inferred genotypes for up to 20 samples	
	Genotype Heatmap	Heatmap showing the inferred genotype derived from each selected sample	
	Haplotype Heatmap	Heatmap showing the inferred haplotype derived from each selected sample. Only samples that can be haplotyped with the selected gene will be shown.	
	Allele Appearances	Charts the frequency at which each allele appears (on either or both chromosomes) in the selected samples	
	Heterozygosity	Charts the extent to which alleles of each gene are heterozygous in the selected samples	
	Allele Usage	Charts the number of alleles of each gene that appear in the selected samples	
	Gene Frequencies	Plots the expression frequency of genes in selected samples	
Overview comparison of two datasets.			

Future plans

- Adapting TigGER, RAbHIT and VDJbase to other species.
- Adding more projects to VDJbase.
- Adding and improving comparative visualizations.
- Creating interface between different servers and VDJbase
(e.g OGRDB)



Acknowledgements

Yaari Lab

- ★ Gur Yaari
- ★ Pazit Polak
- ★ Aviv Omer
- ★ Moriah Gidoni

Oslo Group

- ★ Ludvig Sollid
- ★ Victor Greiff
- ★ Ivana Mikocziova
- ★ Omri snir
- ★ Ida Lindeman

Watson Lab

- ★ Corey T Watson
- ★ Oscar L Rodriguez

IARC Group

- ★ William Lees
- ★ Andrew M Collins
- ★ Martin Corcoran
- ★ Mats Ohlin

Detecting immunoglobulin heavy chain locus genetic variation by targeted long-read sequencing

Oscar L. Rodriguez, PhD

Watson Lab

Postdoctoral fellow | Department of Biochemistry and Molecular Genetics
University of Louisville School of Medicine

AIRR Community webinars
February 10, 2022



Outline of presentation

- 1) Introduction
 - a) Genetic variation in the immunoglobulin heavy chain locus (IGH)
 - b) Advantages of long-read sequencing over short-read sequencing
- 2) Framework for resolving the immunoglobulin heavy locus in a high throughput fashion using long read sequencing
- 3) Resolving the immunoglobulin heavy chain locus in a cohort with adaptive immune repertoire sequencing
- 4) Resolving the T-cell receptor loci using long read sequencing
- 5) Application of framework to the immunoglobulin loci in rhesus macaque
- 6) Conclusion

Genetic variation defined by genetic differences between individuals

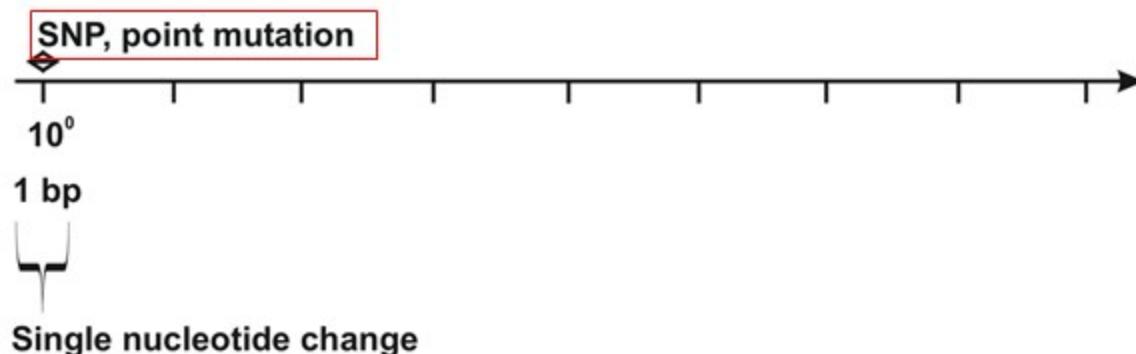
Person 1: ATGCTAGCTAGTCAG

Person 2: ATGCTA**C**CTAGTCAG

Genetic variation ranges in size from base differences to large structural variants

Person 1: ATGCTAGCTAGTCAG

Person 2: ATGCTA**C**CTAGTCAG



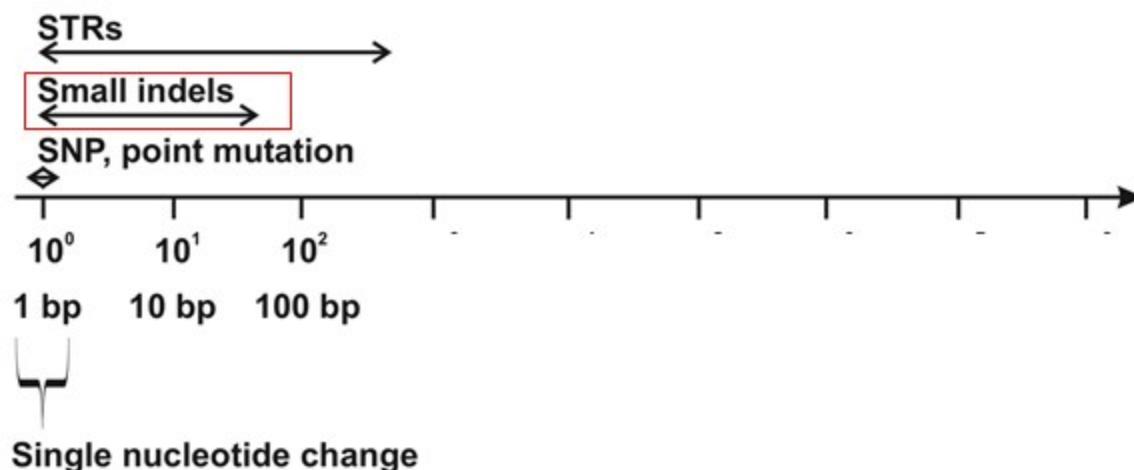
Genetic variation ranges in size from base differences to large structural variants

Person 1: ATGCTAGCTAGTCAG

Person 2: ATGCTAC**C**CTAGTCAG

Person 1: ATGCTAGCTAGTCAG

Person 2: ATGCTA**T**TCAG



Genetic variation ranges in size from base differences to large structural variants

Person 1: ATGCTAGCTAGTCAG

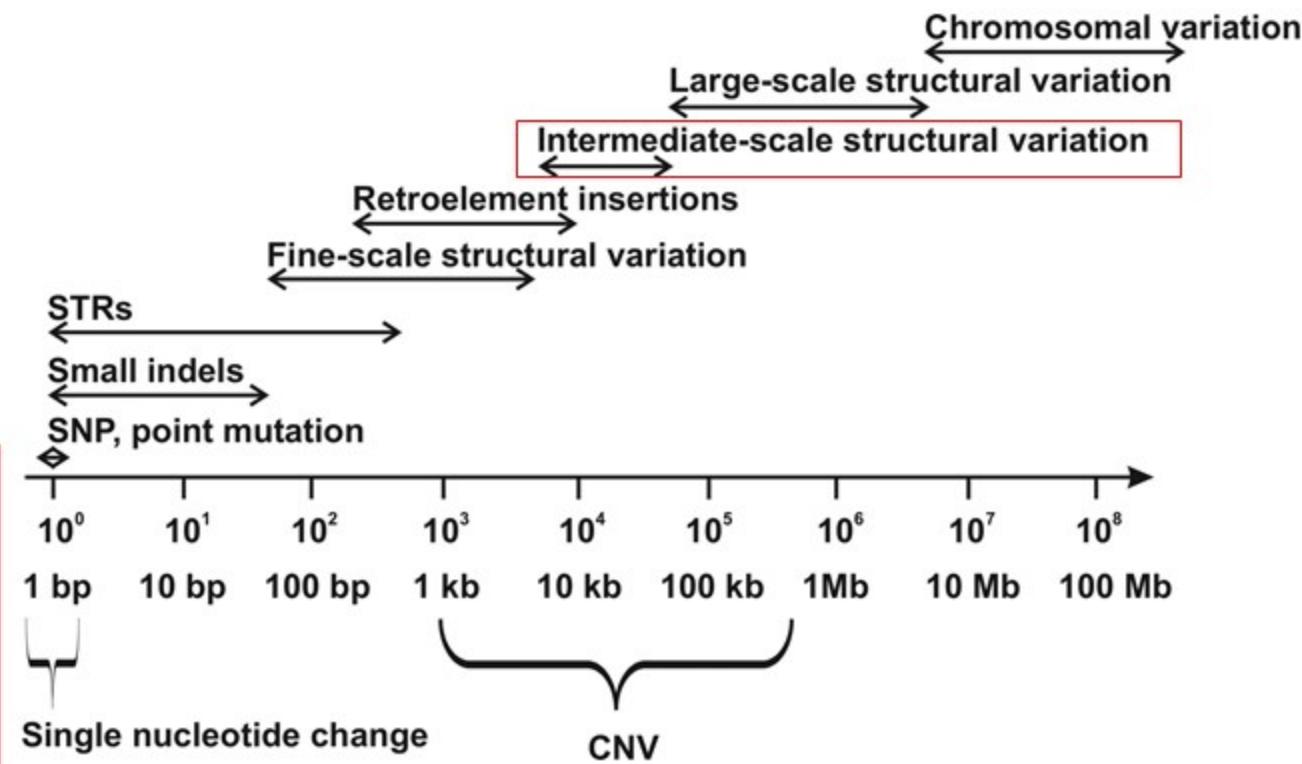
Person 2: ATGCTA~~C~~CTAGTCAG

Person 1: ATGCTAGCTAGTCAG

Person 2: ATGCTA~~T~~TCAG

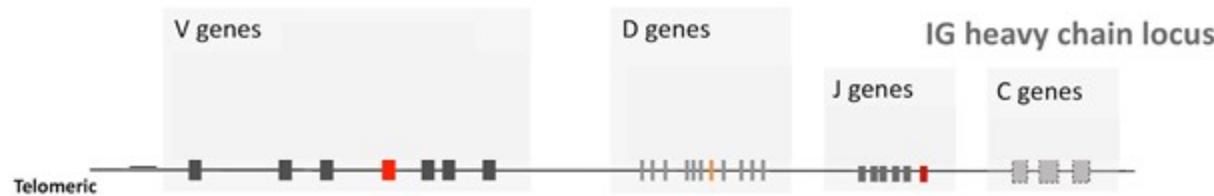
Person 1: ATGCTAGCTAGTCAG

Person 2: ATGCTACCTAGTCAG

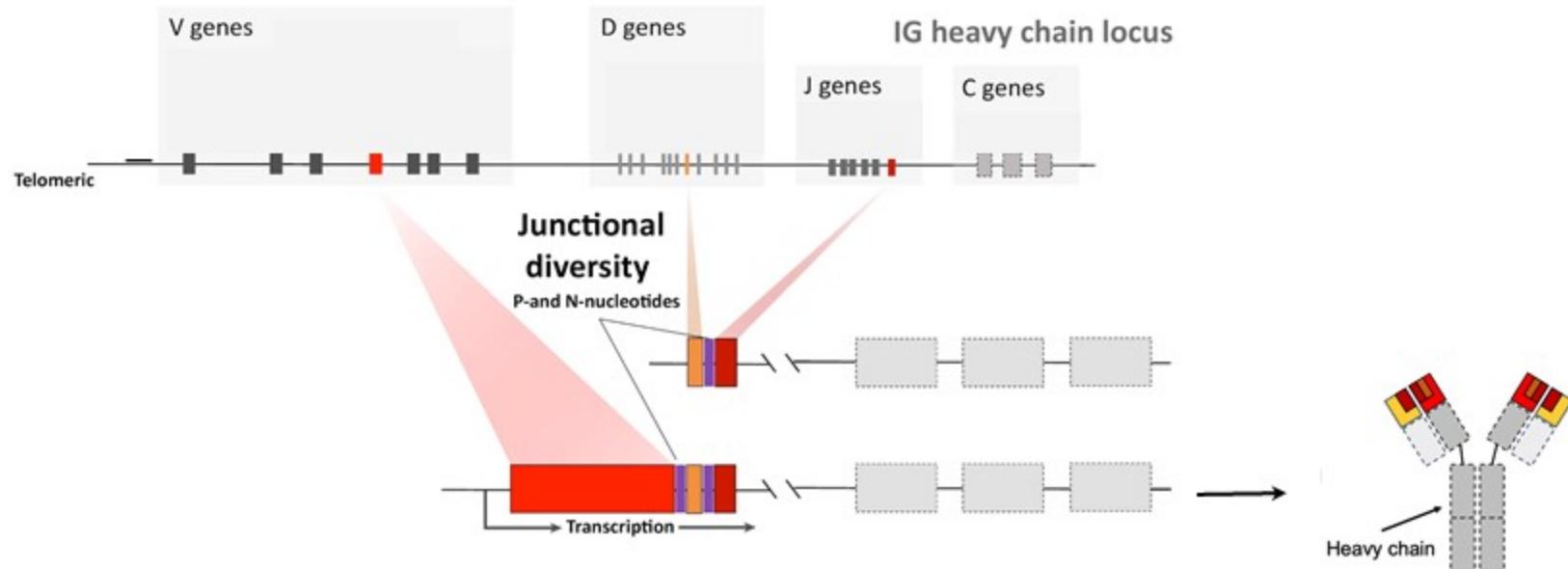


ATGCTACCTAGTCAGATGCTAC
CTAGTCAGATGCTACCTAGTCA
GATGCTACCTAGTCAGATGCTA
CCTAGTCAGATGCTACCTAGTC
AGAGTCAGATGCTACCTAGTCA
GAGTCAGATGCTACCTAGTCAG
AGTCAGATGCTACCTAGTCAGA
CAGTCAGATGCTACCTAGTCAG
AGTCAGAT

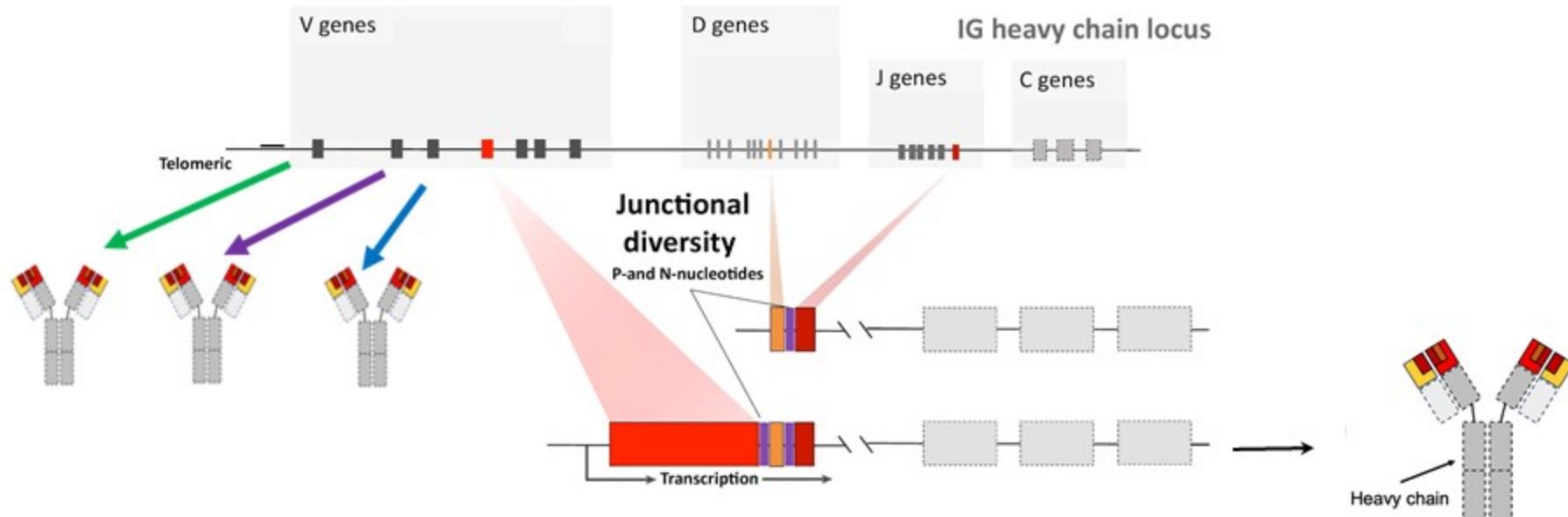
IGH locus is located at the telomeric end of chromosome 14



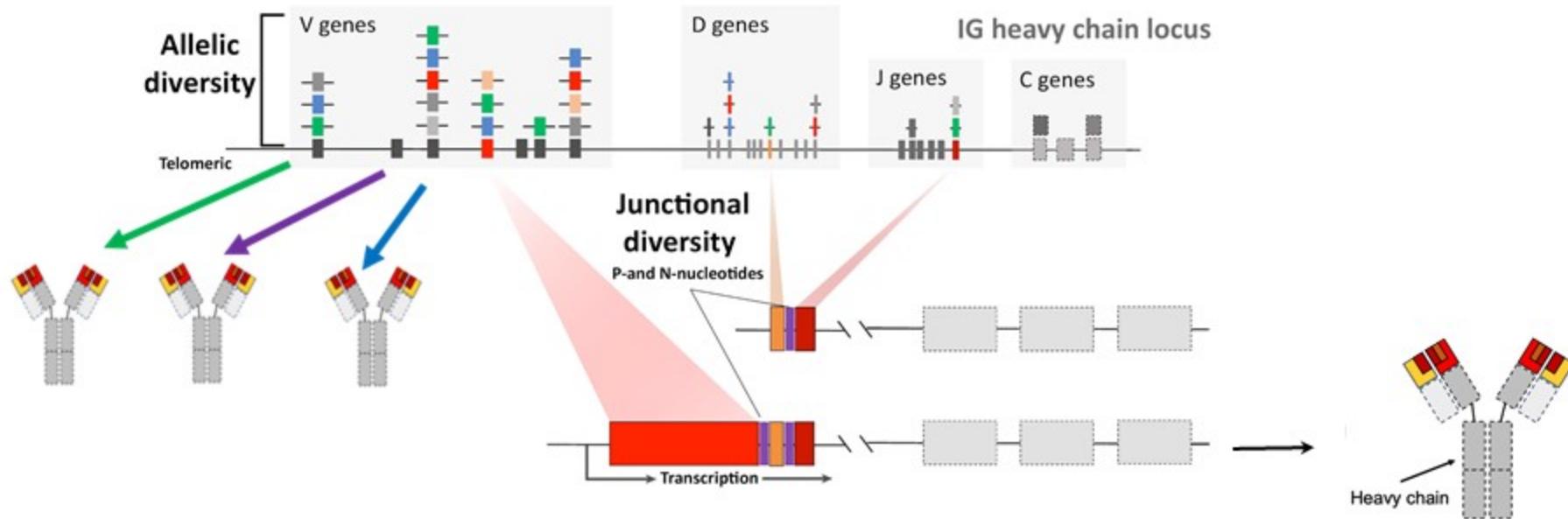
Rearranged DNA gets transcribed and translated to produce the heavy chain of antibodies



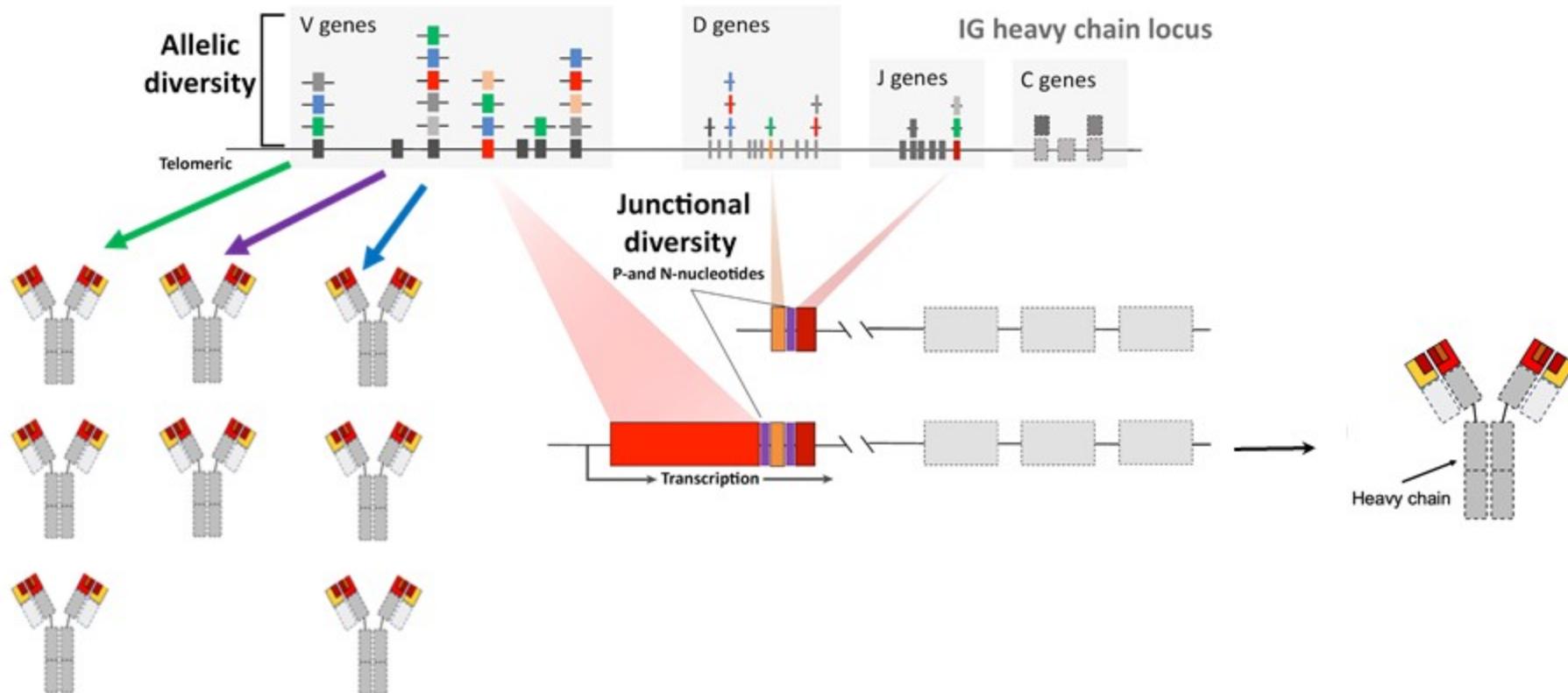
By selecting different genes, B cells can create a diverse antibody repertoire



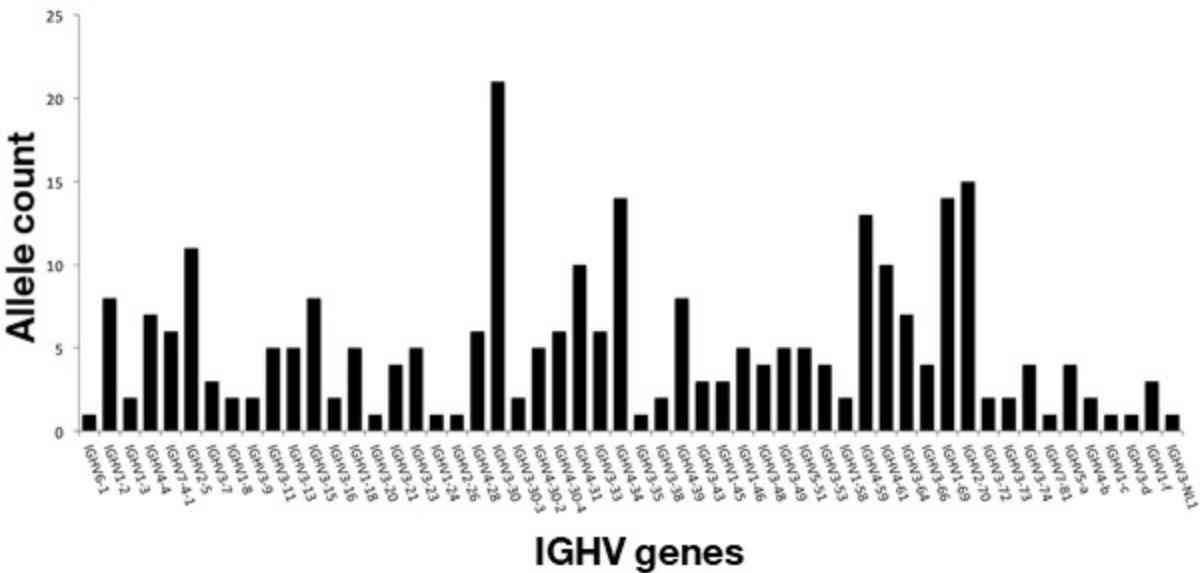
At the population level, there also exist haplotype diversity



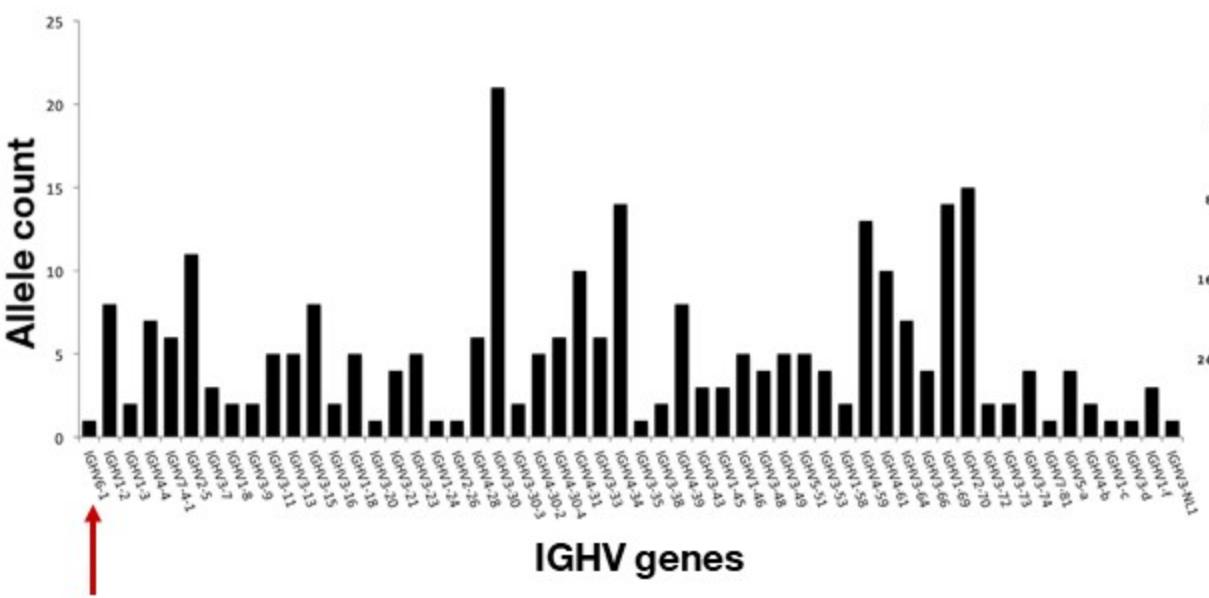
At the population level, there also exist haplotype diversity



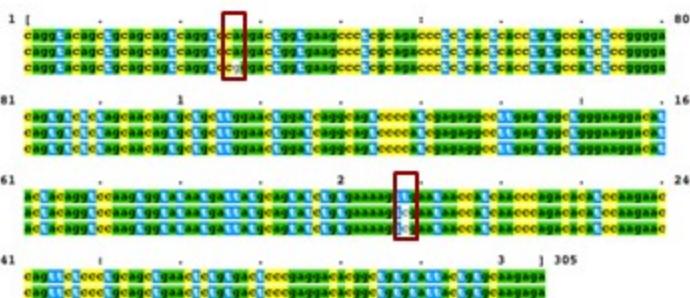
Current data shows extensive allelic diversity



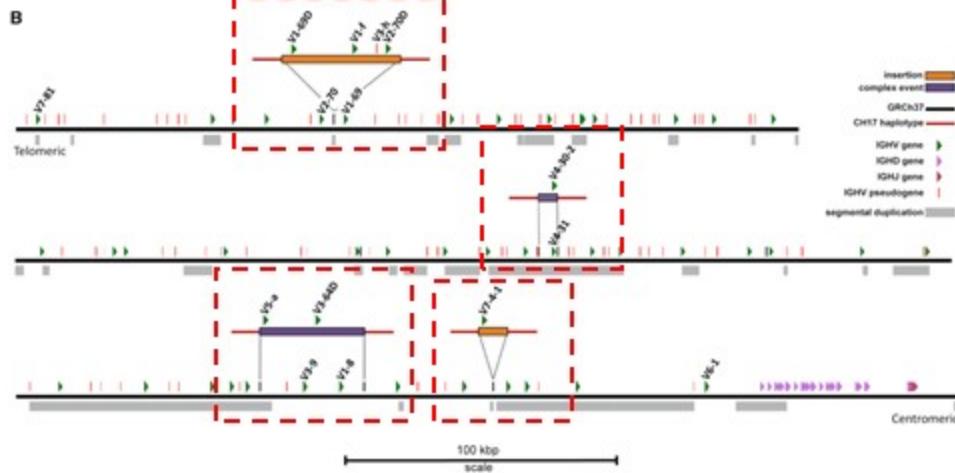
Current data shows extensive allelic diversity



IGHV6-1*01
IGHV6-1*02
IGHV6-1*03



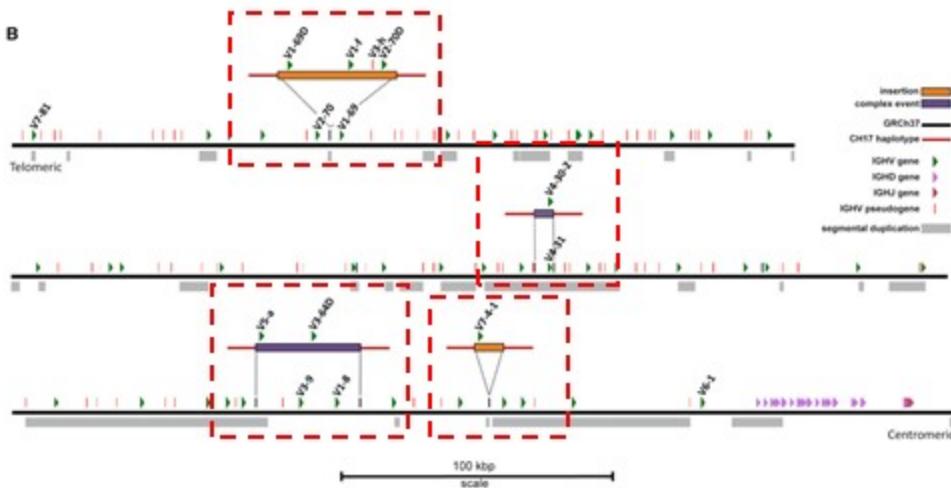
Large number of structural variants in IGH



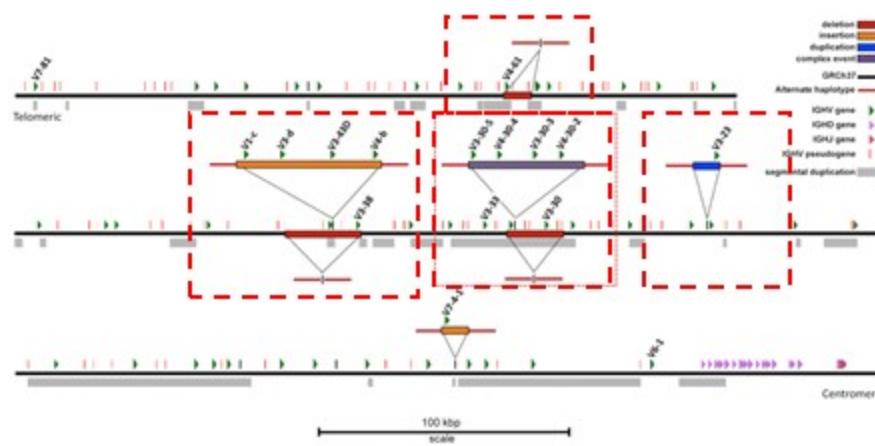
Structural variants identified in a
single individual

Large number of structural variants in IGH

B



Structural variants identified in a single individual



Structural variants identified in a many individuals across different populations

Large number of structural variants in IGH

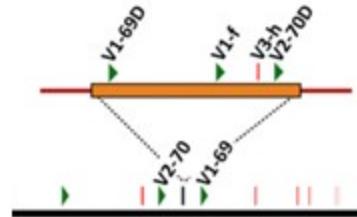
Table 1. CNVs Identified from BAC and Fosmid Clones

Individual	Population	CNV Type	IGHV Genes Included in CNV ^a	GRCh37 Outer-Start (Breakpoint) ^b	GRCh37 Outer-End (Breakpoint) ^b	Event Size (-kbp)
CH17	nd	Insertion	V1-69D, V1-f, V3-h, V2-70D (gain)	107174927	107174941	46.6
CH17	nd	Complex event	V4-30-2 (gain) V4-31 (loss)	106804332	106810878	6.5 ^c /48.8 ^d
CH17	nd	Complex event	V5-a, V3-64D (gain) V3-9, V1-8 (loss)	106531320	106569343	38 ^c /37.7 ^e
CH17	nd	Insertion	V7-4-1 (gain)	106483362	106484225	9.5
NA12156	CEPH	Deletion	V4-39, V3-38 (loss)	106866357	106899042	32.7
NA15510 and NA19240	nd and Yoruba	Insertion	V1-c, V3-d, V3-43D, V4-b (gain)	106877146	106877535	61.1
NA18555 (haplotype A)	Han Chinese	Complex event	V3-30-5, V4-30-4, V3-30-3, V4-30-2 (gain)	106804332	106804333	49.2
NA18555 (haplotype B)	Han Chinese	Deletion	V4-31, V3-30 (loss)	106786254	106811213	24.9 ^c /73.9 ^d
NA18507	Yoruban	Complex event ^g	V4-30-4, V3-30-3 (gain) V3-30 (loss)	106784242	nd	25.2
NA18502	Yoruban	Complex event ^g	V3-30-5 (gain) V3-33, V4-31 (loss)	nd	106820685	24.7
NA18956 and NA12156	Japanese and CEPH	Duplication	V3-23D (gain)	106716650	106727861	10.8
NA19240 and NA12878	Yoruban and CEPH	Insertion	V7-4-1 (gain)	106483362	106484225	9.5

Large number of structural variants in IGH

Table 1. CNVs Identified from BAC and Fosmid Clones

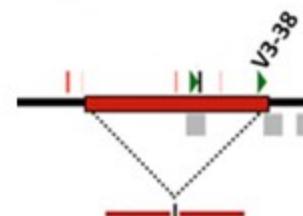
Individual	Population	CNV Type	IGHV Genes Included in CNV ^a	GRCh37 Outer-Start (Breakpoint) ^b	GRCh37 Outer-End (Breakpoint) ^b	Event Size (-kbp)
CH17	nd	Insertion	V1-69D, V1-f, V3-h, V2-70D (gain)	107174927	107174941	46.6
CH17	nd	Complex event	V4-30-2 (gain) V4-31 (loss)	106804332	106810878	6.5 ^c /48.8 ^d
CH17	nd	Complex event	V5-a, V3-64D (gain) V3-9, V1-8 (loss)	106531320	106569343	38 ^c /37.7 ^e
CH17	nd	Insertion	V7-4-1 (gain)	106483362	106484225	9.5
NA12156	CEPH	Deletion	V4-39, V3-38 (loss)	106866357	106899042	32.7
NA15510 and NA19240	nd and Yoruba	Insertion	V1-c, V3-d, V3-43D, V4-b (gain)	106877146	106877535	61.1
NA18555 (haplotype A)	Han Chinese	Complex event	V3-30-5, V4-30-4, V3-30-3, V4-30-2 (gain)	106804332	106804333	49.2
NA18555 (haplotype B)	Han Chinese	Deletion	V4-31, V3-30 (loss)	106786254	106811213	24.9 ^c /73.9 ^d
NA18507	Yoruban	Complex event ^g	V4-30-4, V3-30-3 (gain) V3-30 (loss)	106784242	nd	25.2
NA18502	Yoruban	Complex event ^g	V3-30-5 (gain) V3-33, V4-31 (loss)	nd	106820685	24.7
NA18956 and NA12156	Japanese and CEPH	Duplication	V3-23D (gain)	106716650	106727861	10.8
NA19240 and NA12878	Yoruban and CEPH	Insertion	V7-4-1 (gain)	106483362	106484225	9.5



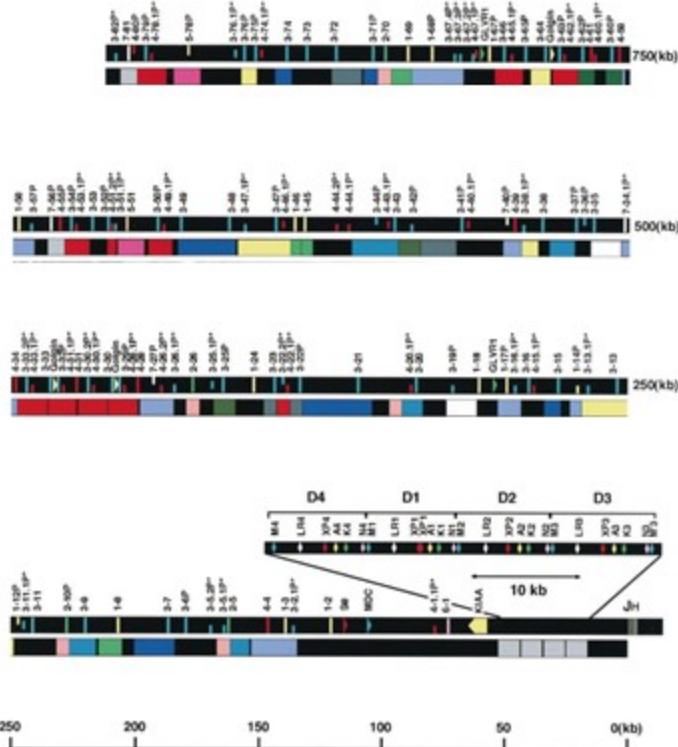
Large number of structural variants in IGH

Table 1. CNVs Identified from BAC and Fosmid Clones

Individual	Population	CNV Type	IGHV Genes Included in CNV ^a	GRCh37 Outer-Start (Breakpoint) ^b	GRCh37 Outer-End (Breakpoint) ^b	Event Size (-kbp)
CH17	nd	Insertion	V1-69D, V1-f, V3-h, V2-70D (gain)	107174927	107174941	46.6
CH17	nd	Complex event	V4-30-2 (gain) V4-31 (loss)	106804332	106810878	6.5 ^c /48.8 ^d
CH17	nd	Complex event	V5-a, V3-64D (gain) V3-9, V1-8 (loss)	106531320	106569343	38 ^c /37.7 ^e
CH17	nd	Insertion	V7-4-1 (gain)	106483362	106484225	9.5
NA12156	CEPH	Deletion	V4-39, V3-38 (loss)	106866357	106899042	32.7
NA185510 and NA19240	nd and Yoruba	Insertion	V1-c, V3-d, V3-43D, V4-b (gain)	106877146	106877535	61.1
NA18555 (haplotype A)	Han Chinese	Complex event	V3-30-5, V4-30-4, V3-30-3, V4-30-2 (gain)	106804332	106804333	49.2
NA18555 (haplotype B)	Han Chinese	Deletion	V4-31, V3-30 (loss)	106786254	106811213	24.9 ^c /73.9 ^d
NA18507	Yoruban	Complex event ^g	V4-30-4, V3-30-3 (gain) V3-30 (loss)	106784242	nd	25.2
NA18502	Yoruban	Complex event ^g	V3-30-5 (gain) V3-33, V4-31 (loss)	nd	106820685	24.7
NA18956 and NA12156	Japanese and CEPH	Duplication	V3-23D (gain)	106716650	106727861	10.8
NA19240 and NA12878	Yoruban and CEPH	Insertion	V7-4-1 (gain)	106483362	106484225	9.5

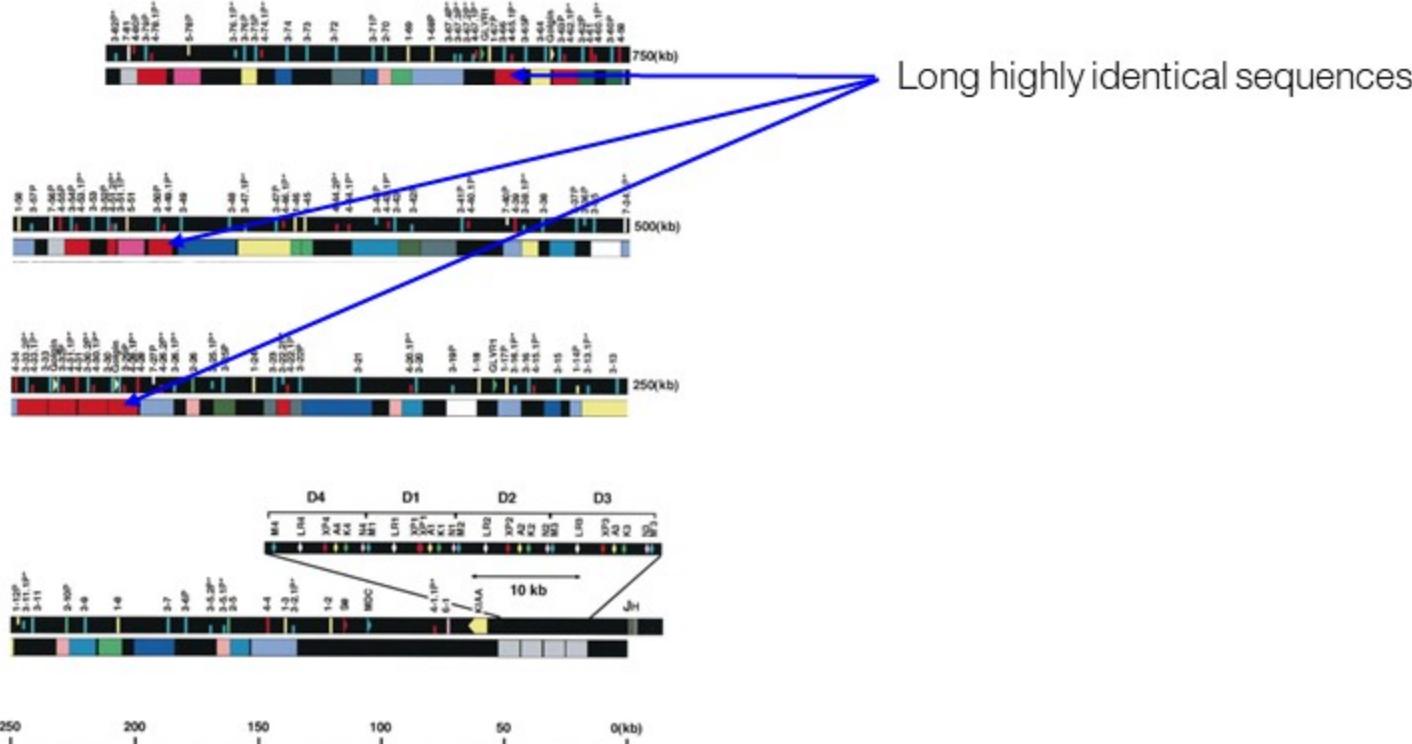


Complexity of locus has hindered previous sequencing methods to genotype IGH

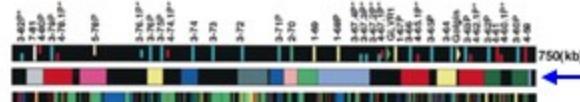


Matsuda et al. J. Exp. Med. 1998

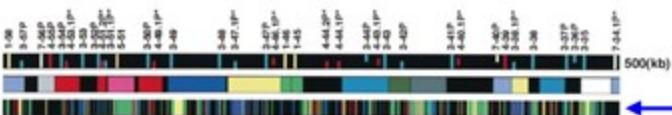
Complexity of locus has hindered previous sequencing methods to genotype IGH



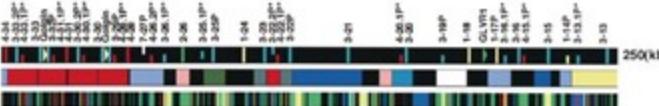
Complexity of locus has hindered previous sequencing methods to genotype IGH



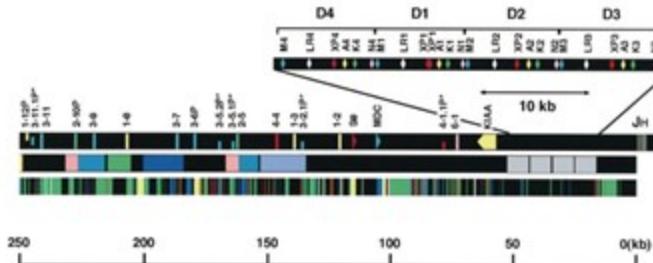
Long highly identical sequences



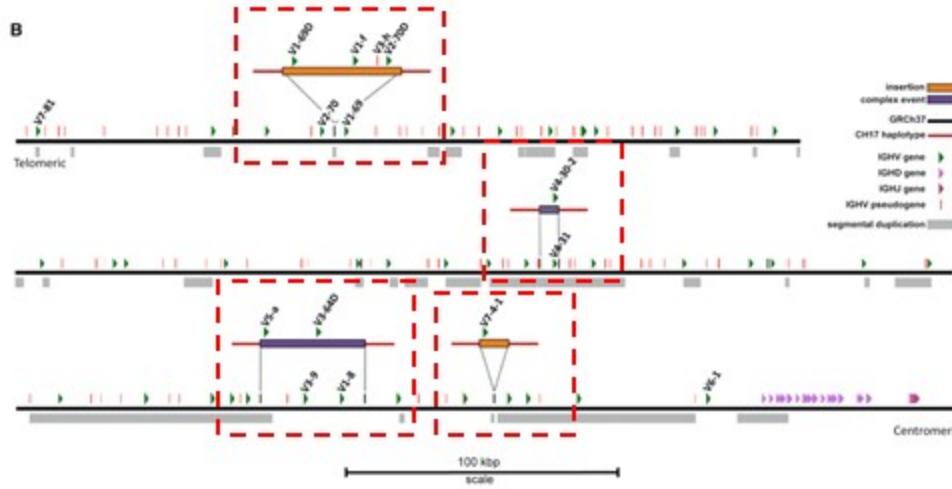
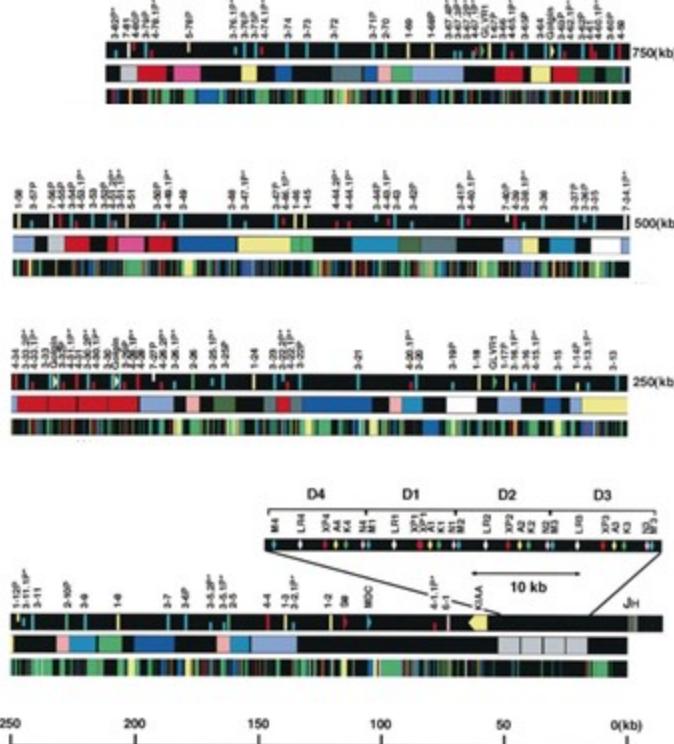
Different repeat classes



1. Alu
 2. MIR
 3. LINE
 4. Simple repeats



Complexity of locus has hindered previous sequencing methods to genotype IGH



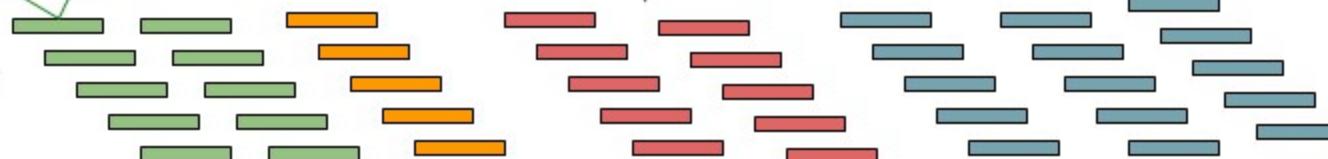
Next-generation sequencing technology generates reads from 50 to 250 bases

Input DNA:



gtgggcagtagttacatcagaattcaaggtaatt
atgcgggtggccctaaaccctaattccgagtatt
agatgttaatcatggtatgttataactatttggcaaat
ggtcacatacatacgccgccttcgttgca
ct

150bp reads:



Reads are aligned to a reference genome

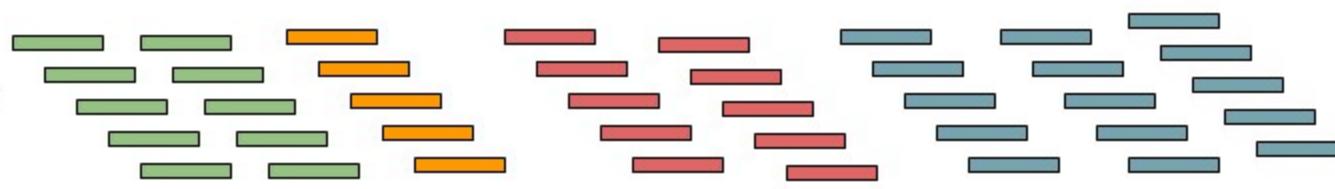
Input DNA:



Reference genome:



150bp reads:



Difference between reads and reference genome correspond to genetic variants

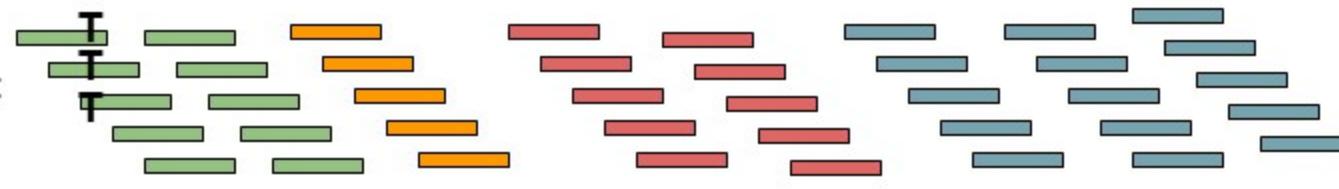
Input DNA:



Reference genome:
A



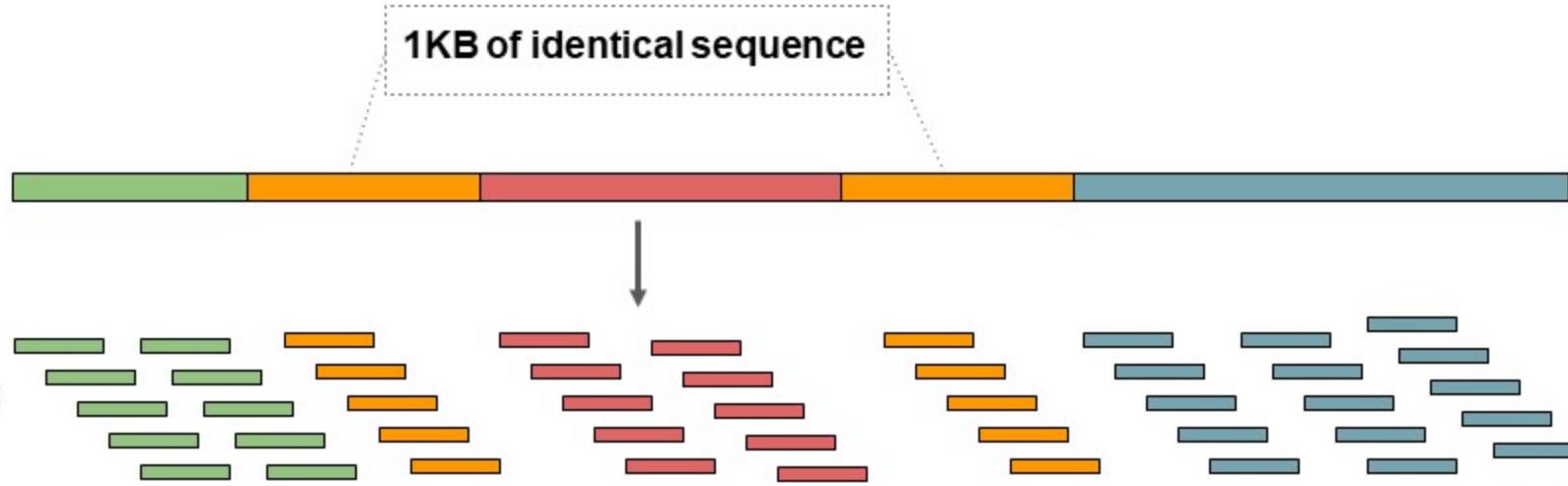
150bp reads:



SNV

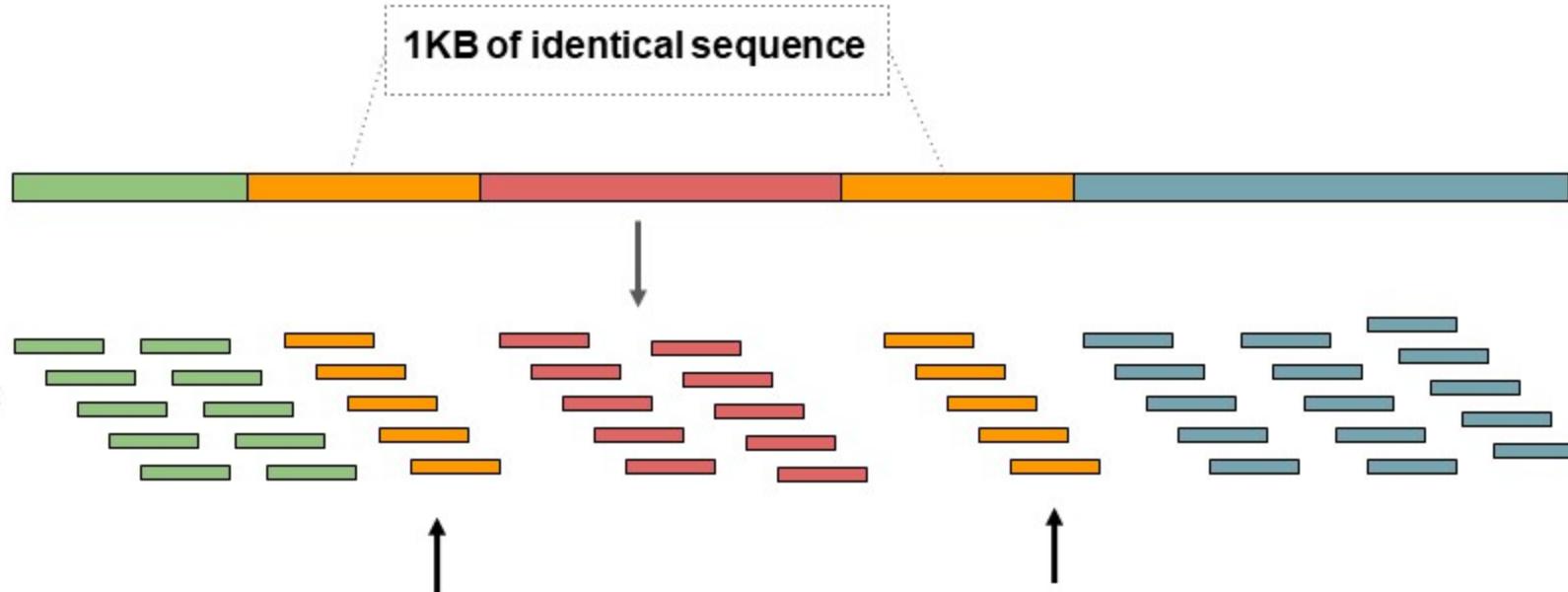
Limitation of next-generation sequencing technology

Input DNA:

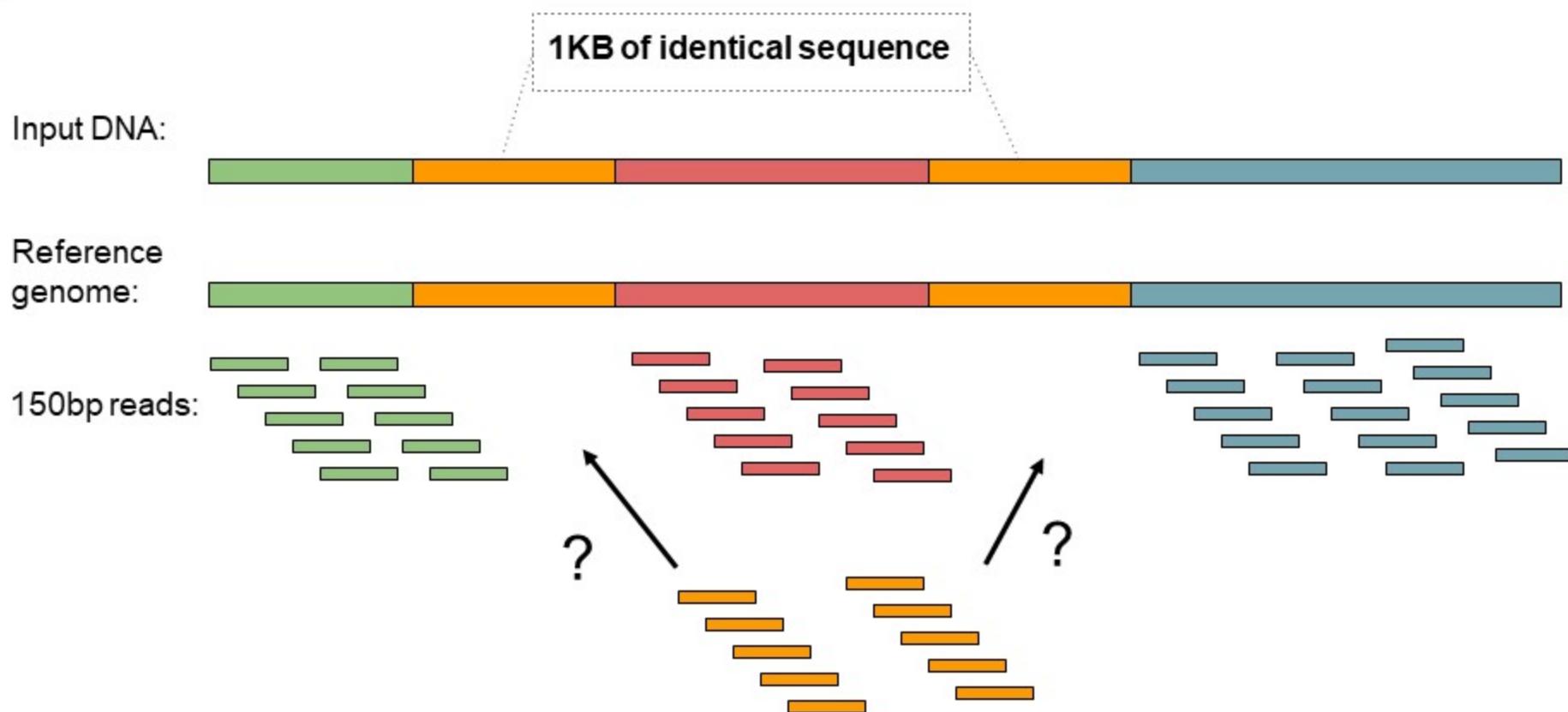


Input DNA with repetitive or identical sequence generates reads with identical sequence

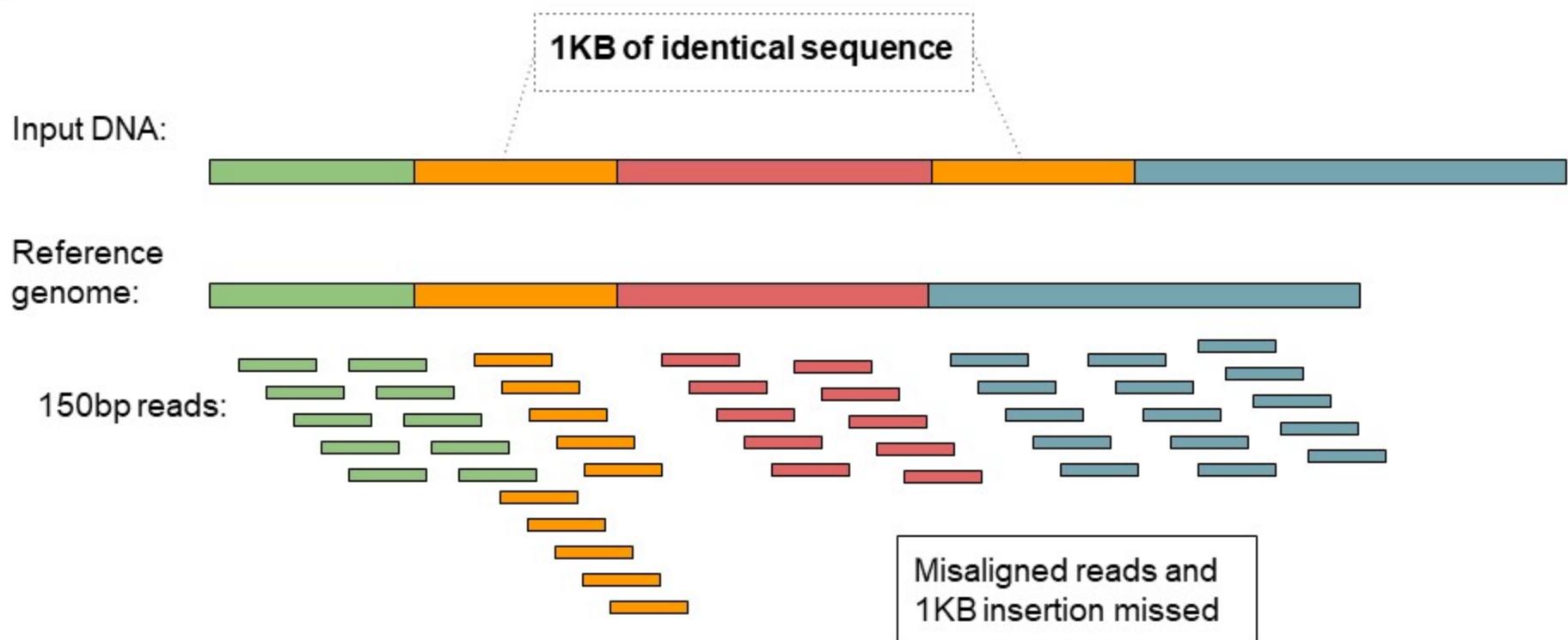
Input DNA:



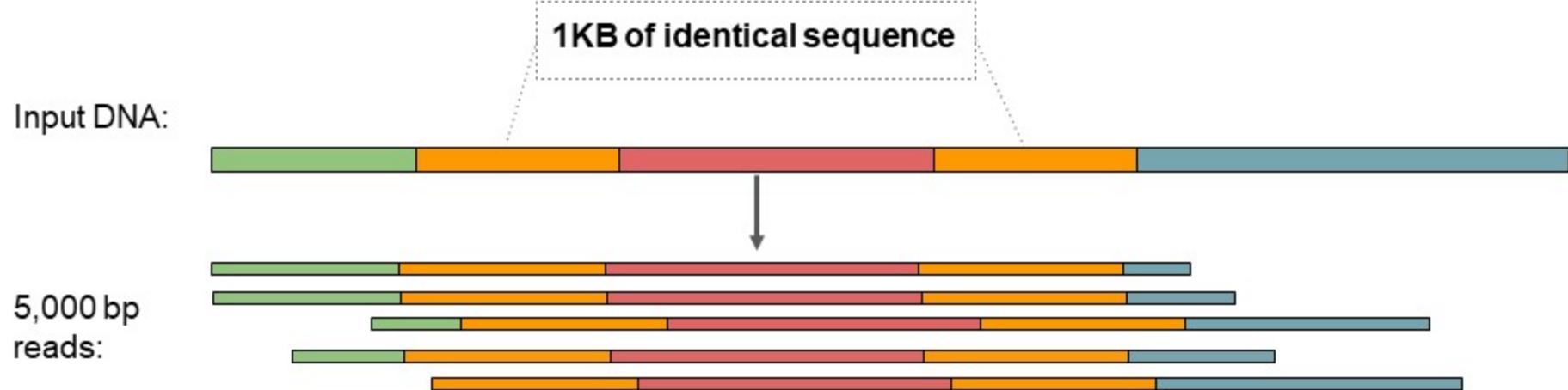
Aligning reads from repetitive sequence misses genetic variation



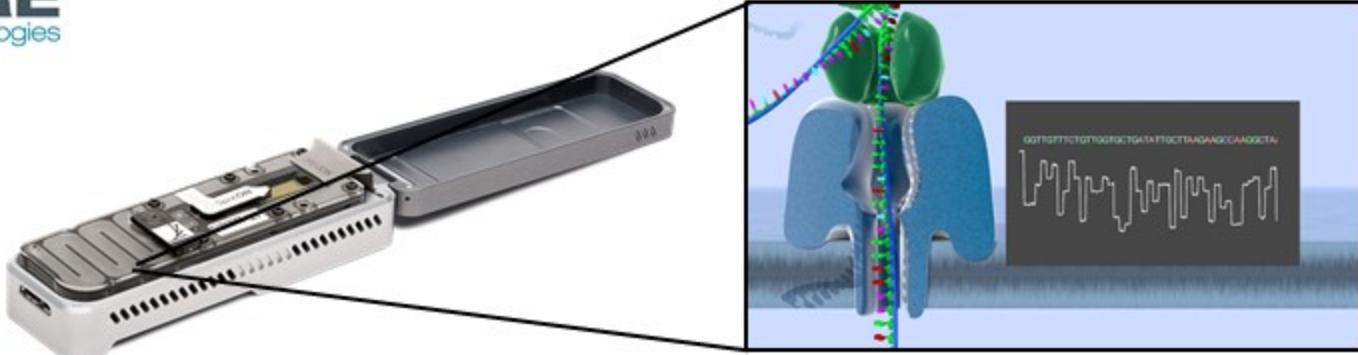
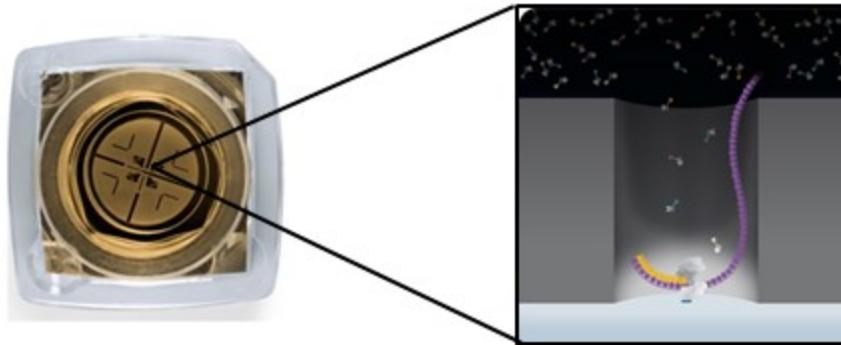
Aligning reads from repetitive sequence misses genetic variation



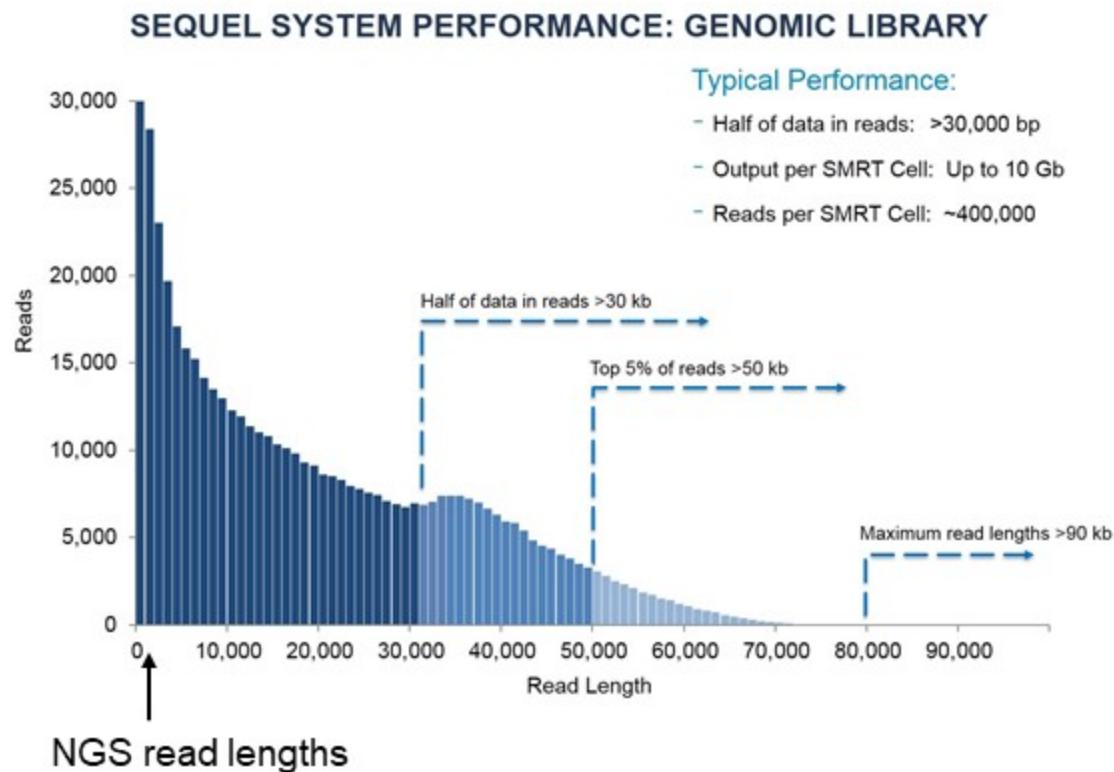
Long-read sequencing can detect more genetic variants



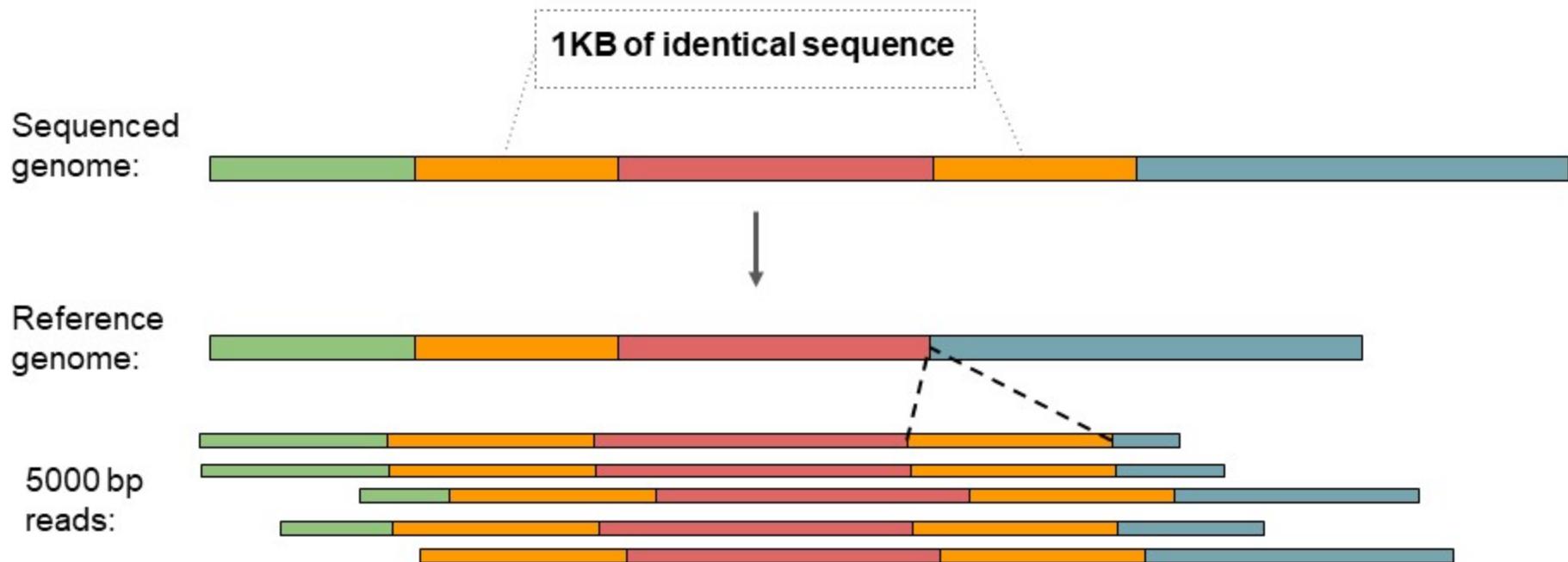
Third generation sequencing technologies are capable of generating long reads



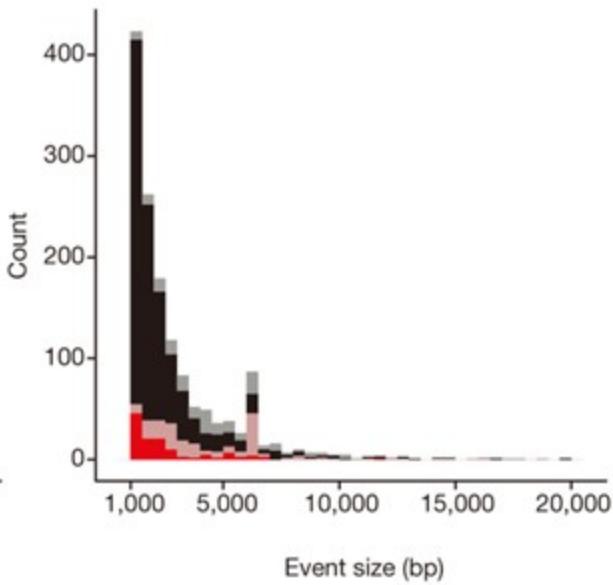
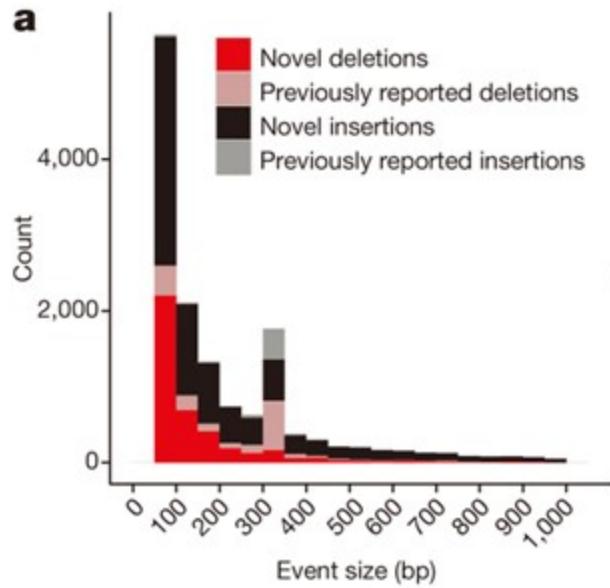
Third generation sequencing technologies are capable of generating long reads



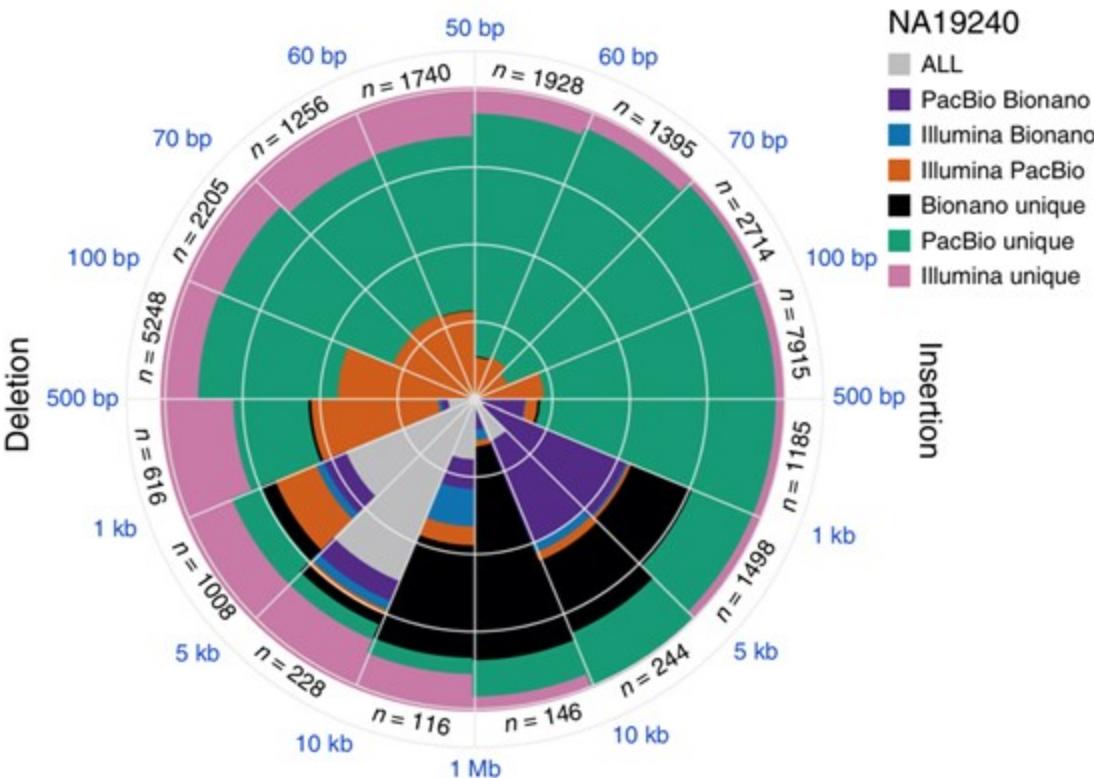
Aligning long reads to the reference identifies previously missed insertion



Large number of novel deletions and insertions identified in complete hydatidiform mole (CHM)



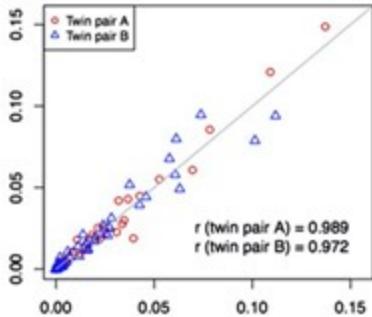
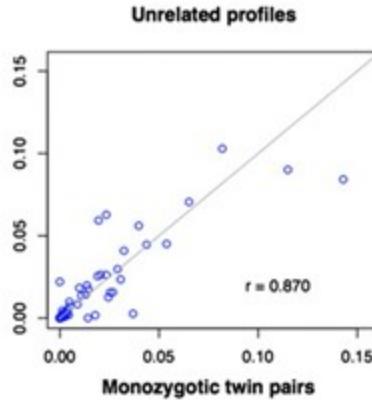
Large number of SVs only detected by long read sequencing methods



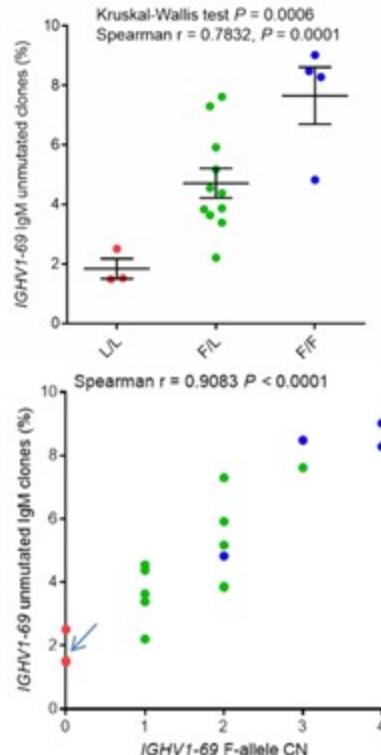
1. ~83% of insertions are being missed by short-read-calling algorithms.
2. 3 to 7 fold more variation detected than with just Illumina-based methods

Evidence that genetic diversity affects the antibody repertoire

Twin studies



Functional studies



GWAS

Seven diseases/traits

- Kawasaki disease
- HDL cholesterol
- Rheumatic heart disease
- Blood protein levels
- Alzheimer's disease (late onset)
- Reaction time
- Myopia (age of diagnosis)

Identification of novel susceptibility loci for kawasaki disease in a Han Chinese population by a genome-wide association study. Tsai et al, 2011 PLoS One

Exploration of haplotype research consortium imputation for genome-wide association studies in 20,032

Generation Scotland participants. Nagy et al, 2017 Genome Medicine

Association between a common immunoglobulin heavy chain allele and rheumatic heart disease risk in Oceania. Parks et al, 2017 Nature Communications

Genomic atlas of the human plasma proteome. Sun et al, 2018 Nature

Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. Lambert JC et al, 2013 Nature Genetics

Study of 300,486 individuals identifies 148 independent genetic loci influencing general cognitive function. Davies G et al, 2018 Nature Communications

Genome-wide association meta-analysis highlights light-induced signaling as a driver for refractive error. Tedja MS et al, 2018 Nature Genetics

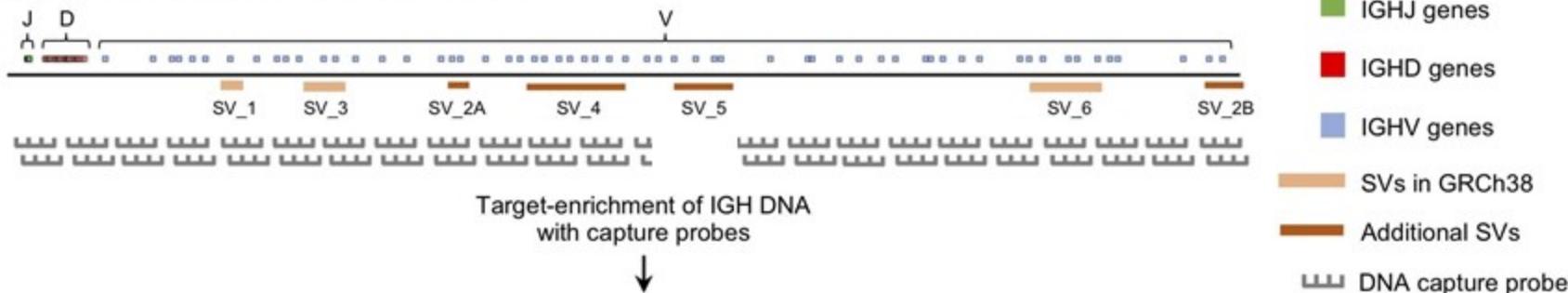
Resolving the IGH locus in a higher throughput and scalable fashion using long read sequencing

O. L. Rodriguez, W. S. Gibson, T. Parks, M. Emery, J. Powell, M. Strahl, G. Deikus, K. Auckland, E. E. Eichler, W. A. Marasco, R. Sebra, A. J. Sharp, M. L. Smith, A. Bashir, C. T. Watson, A novel framework for characterizing genomic haplotype diversity in the human immunoglobulin heavy chain locus, *Frontiers in Immunology*. (2020)

Specifically targeting the IGH locus

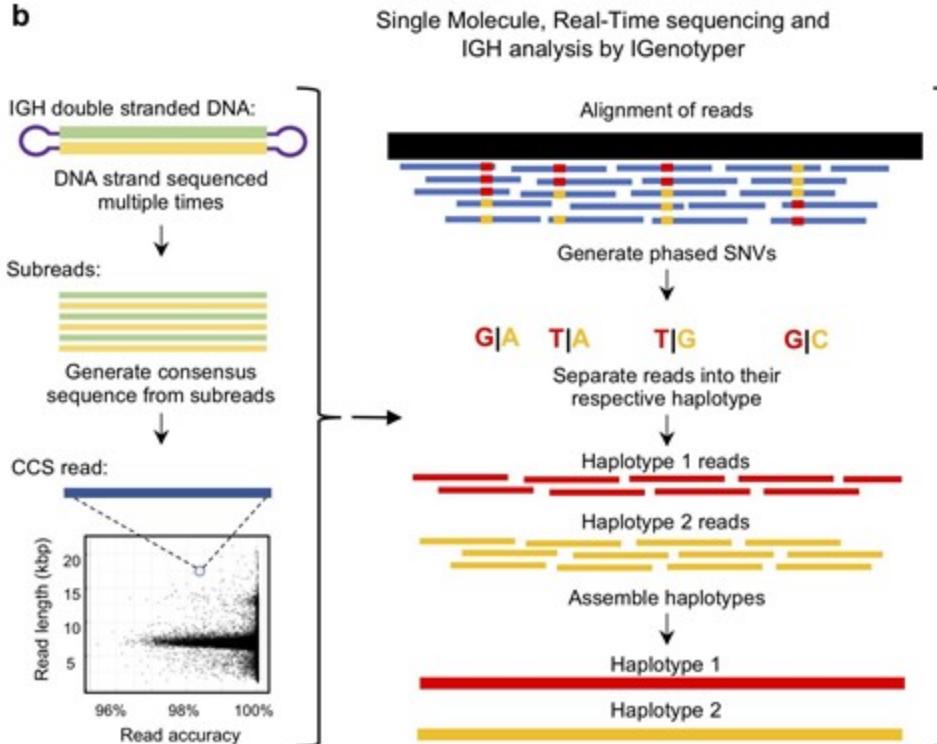
a Custom IGH-reference

(GRCh38 IGH locus with additional SV sequence)



Long read sequencing of targeted DNA

b



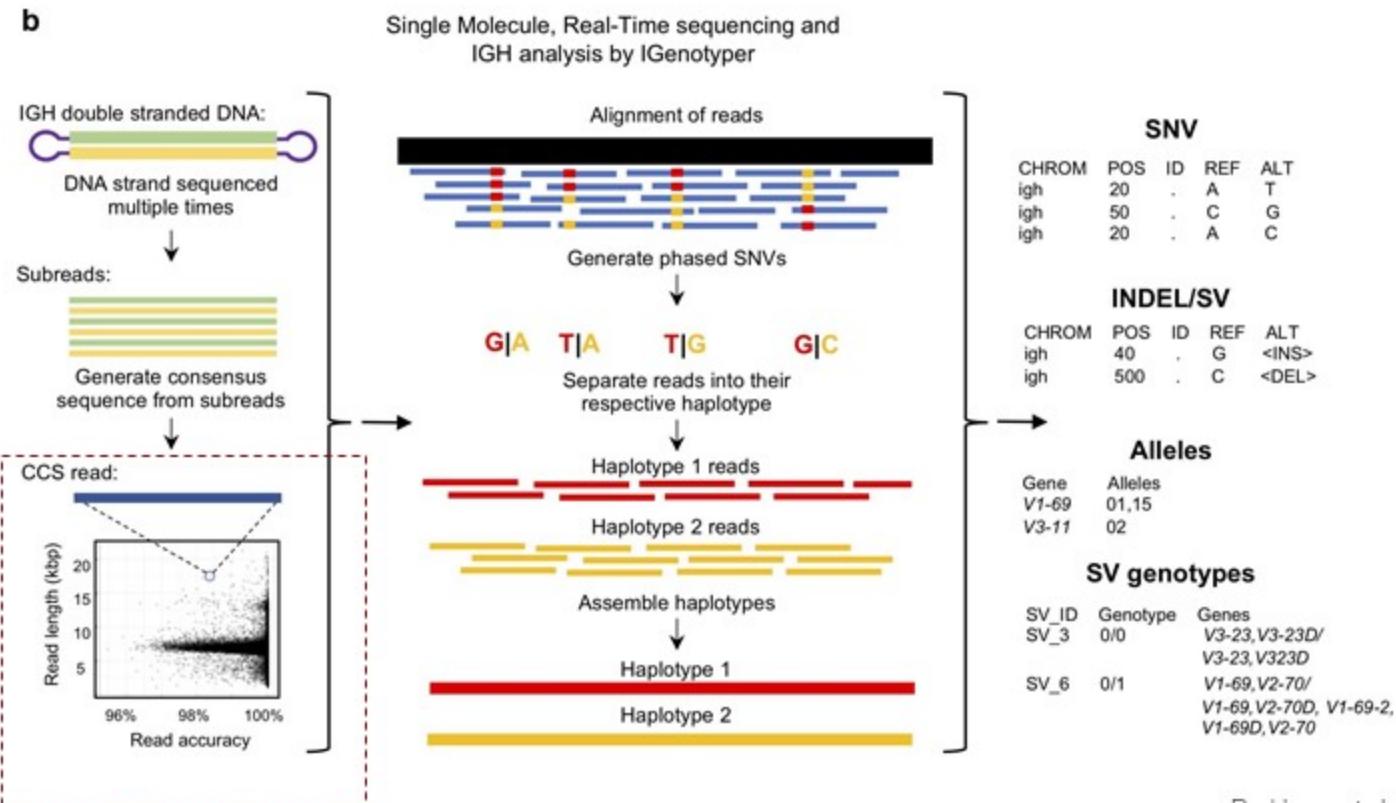
SNV				
CHROM	POS	ID	REF	ALT
igh	20	.	A	T
igh	50	.	C	G
igh	20	.	A	C

INDEL/SV				
CHROM	POS	ID	REF	ALT
igh	40	.	G	<INS>
igh	500	.	C	

Alleles	
Gene	Alleles
V1-69	01,15
V3-11	02

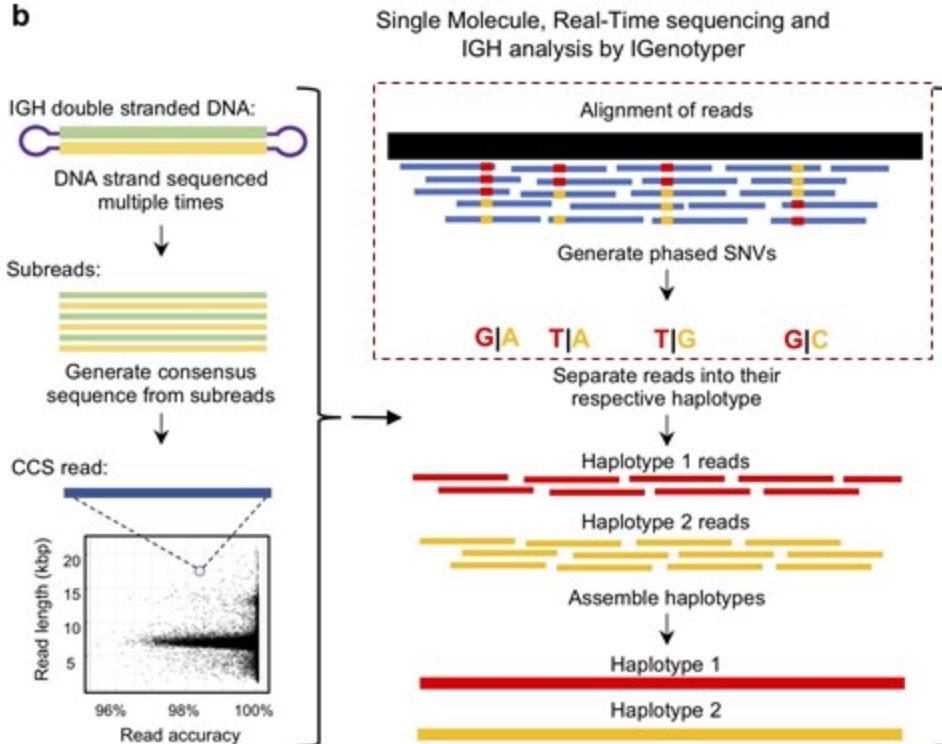
SV genotypes		
SV_ID	Genotype	Genes
SV_3	0/0	V3-23,V3-23D/ V3-23,V323D
SV_6	0/1	V1-69,V2-70/ V1-69,V2-70D, V1-69-2, V1-69D,V2-70

Input to IGenotyper are the highly accurate CCS/HiFi reads



First step in IGenotyper is to detect phased SNVs

b



SNV

CHROM	POS	ID	REF	ALT
igh	20	.	A	T
igh	50	.	C	G
igh	20	.	A	C

INDEL/SV

CHROM	POS	ID	REF	ALT
igh	40	.	G	<INS>
igh	500	.	C	

Alleles

Gene	Alleles
V1-69	01,15
V3-11	02

SV genotypes

SV_ID	Genotype	Genes
SV_3	0/0	V3-23,V3-23D/ V3-23,V323D
SV_6	0/1	V1-69,V2-70/ V1-69,V2-70D, V1-69-2, V1-69D,V2-70

Input reads are aligned to the reference genome



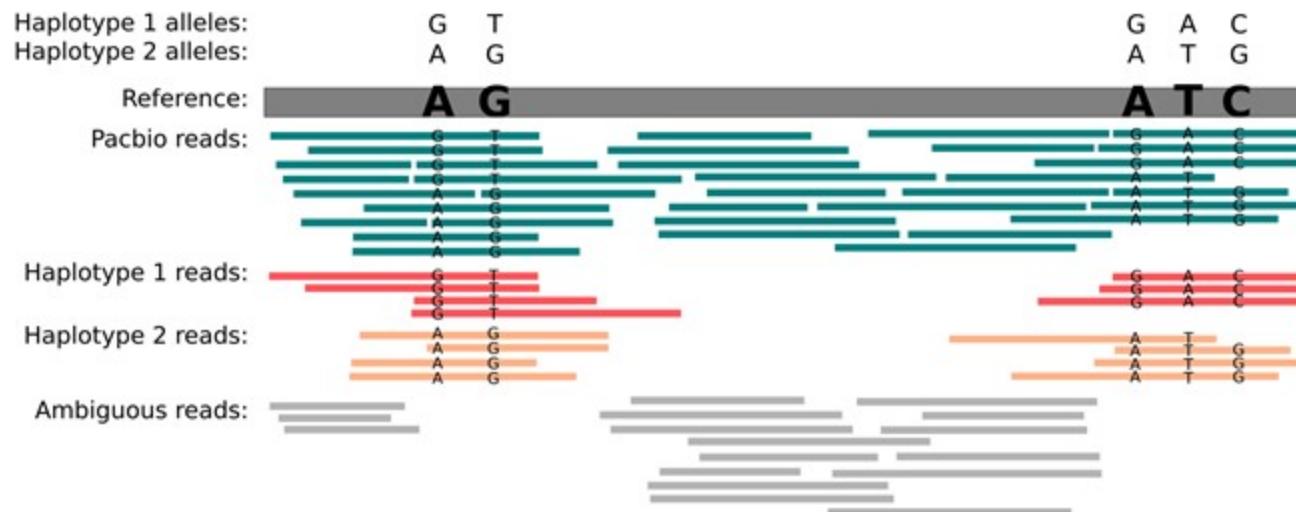
Phased heterozygous SNVs detected using WhatsHap



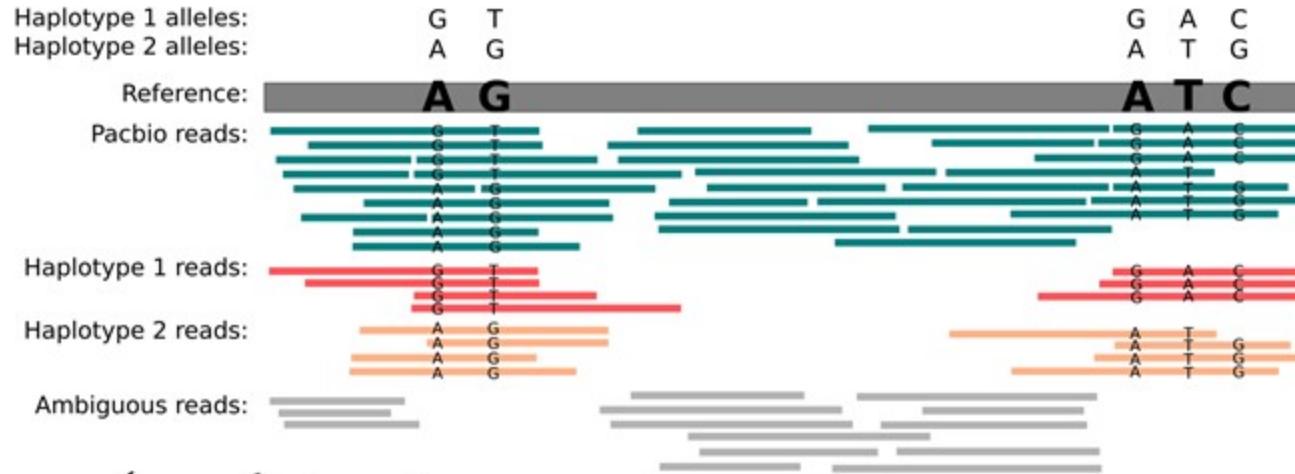
Phased heterozygous SNVs are used to partition reads into their respective haplotype



PacBio reads with haplotype-specific bases are partitioned into their respective haplotype



PacBio reads with haplotype-specific bases are partitioned into their respective haplotype



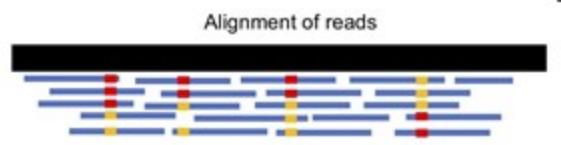
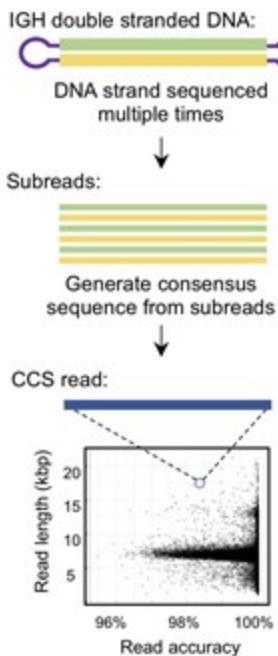
$$P(r | h_r = b) = \prod_{(s_i^1, s_i^2) \in S_r} \begin{cases} 1 - 10^{-\frac{q_{r_i}}{10}} & \text{if } r_i = s_i^b \\ 10^{-\frac{q_{r_i}}{10}} & \text{otherwise} \end{cases} \quad (1)$$

$$b_r = \begin{cases} 1 & \text{if } \frac{P(r|h_r = 1)}{P(r|h_r = 1) + P(r|h_r = 2)} > \tau \\ 2 & \text{if } \frac{P(r|h_r = 2)}{P(r|h_r = 1) + P(r|h_r = 2)} > \tau \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

IGenotyper produces haplotype-specific assemblies

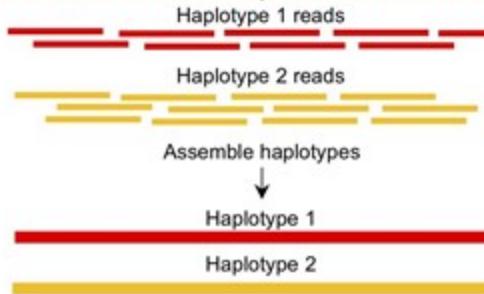
b

Single Molecule, Real-Time sequencing and
IGH analysis by IGenotyper



G|A T|A T|G G|C

Separate reads into their respective haplotype



SNV

CHROM	POS	ID	REF	ALT
igh	20	.	A	T
igh	50	.	C	G
igh	20	.	A	C

INDEL/SV

CHROM	POS	ID	REF	ALT
igh	40	.	G	<INS>
igh	500	.	C	

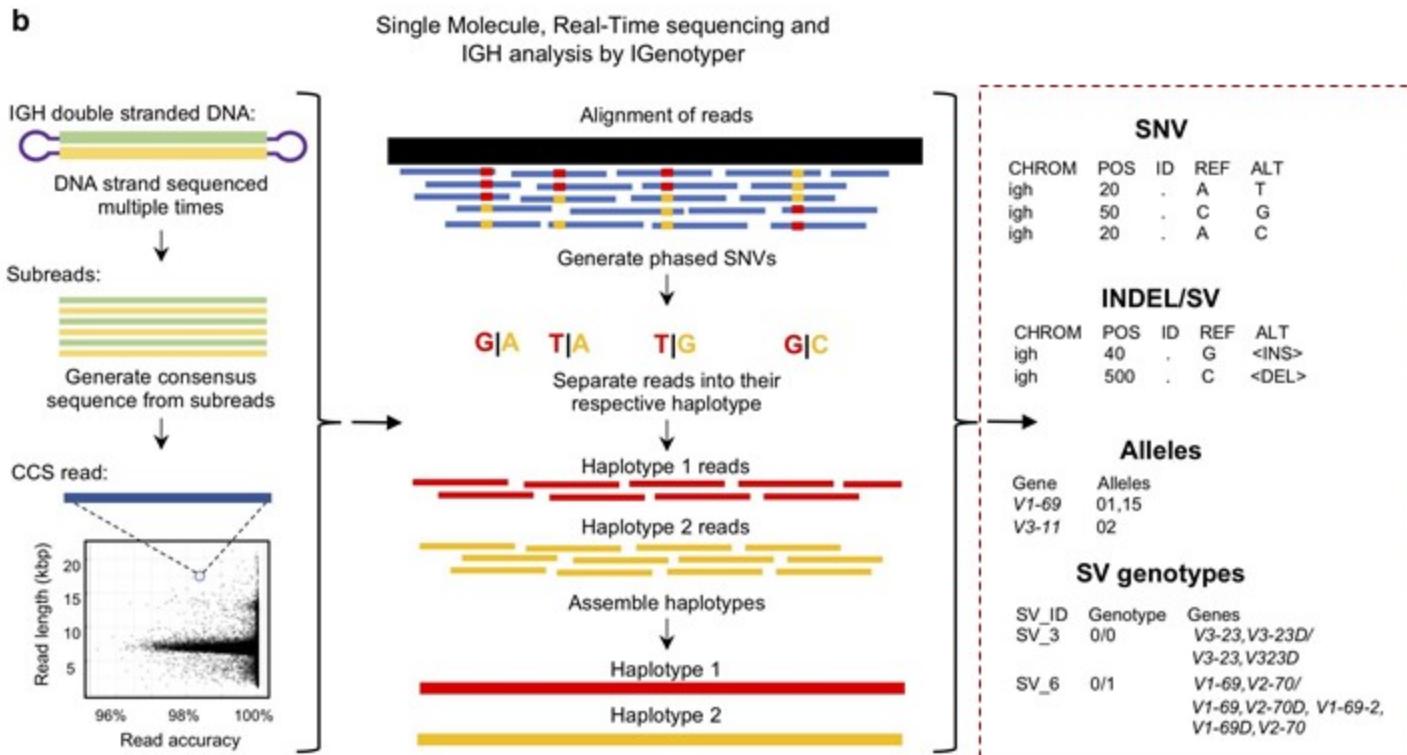
Alleles

Gene	Alleles
V1-69	01,15
V3-11	02

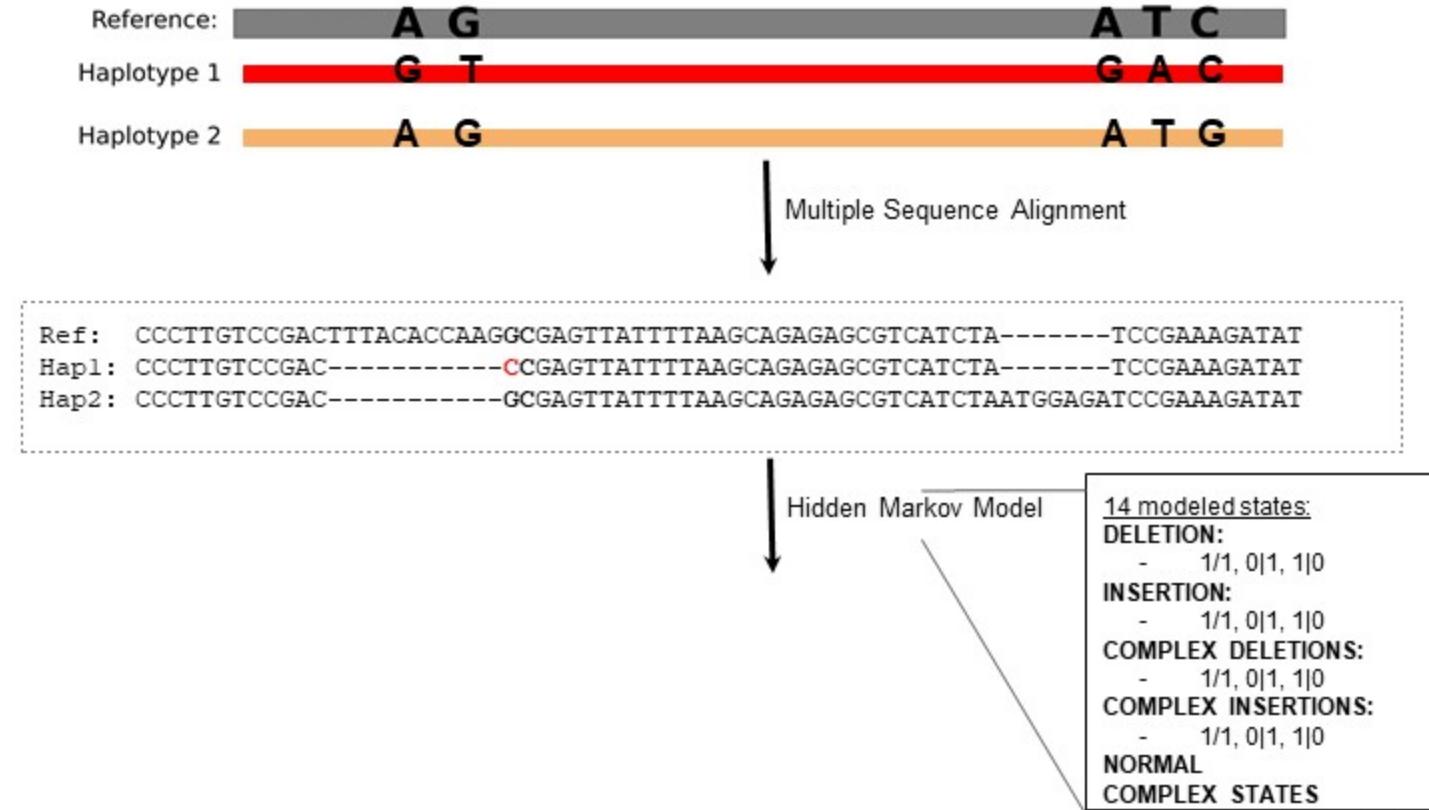
SV genotypes

SV_ID	Genotype	Genes
SV_3	0/0	V3-23,V3-23D/ V3-23,V323D
SV_6	0/1	V1-69,V2-70/ V1-69,V2-70D, V1-69-2, V1-69D,V2-70

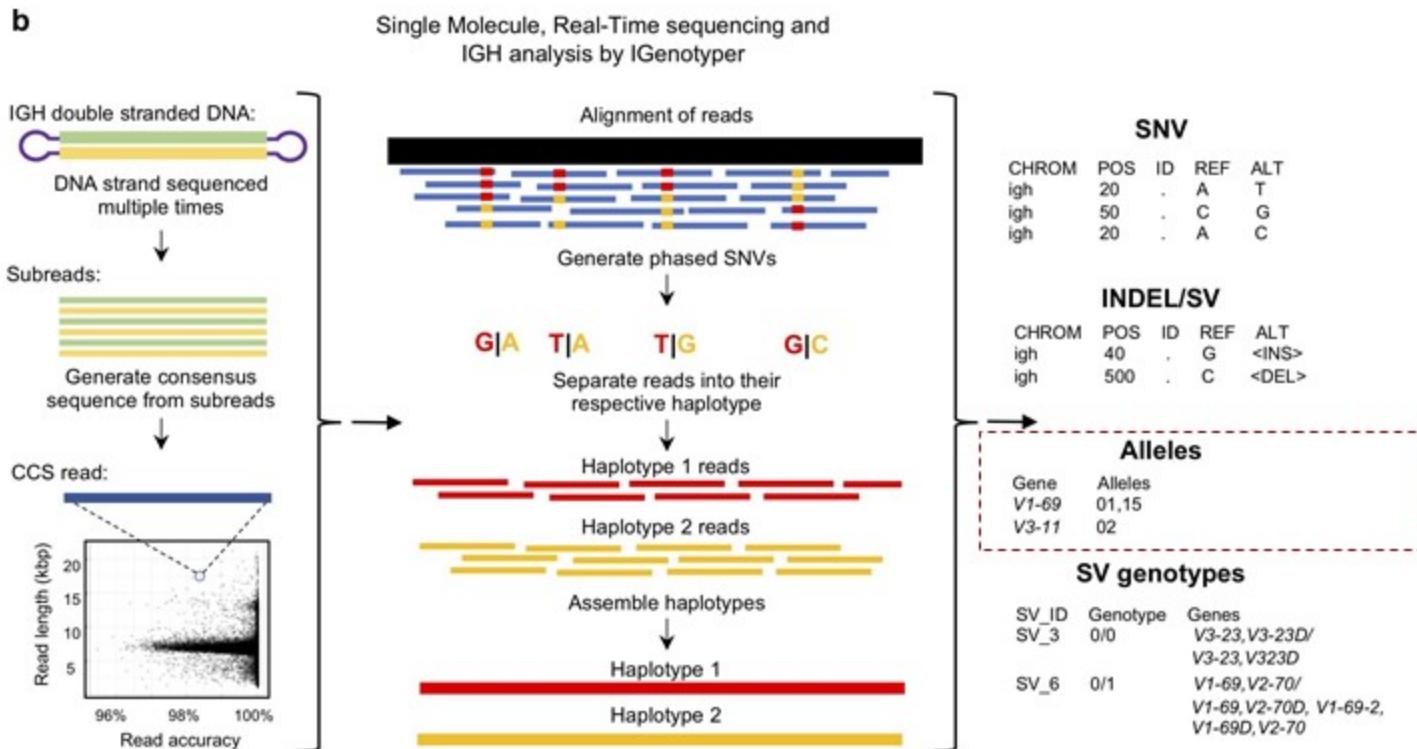
Haplotype-resolved assemblies are used to detect genetic variation



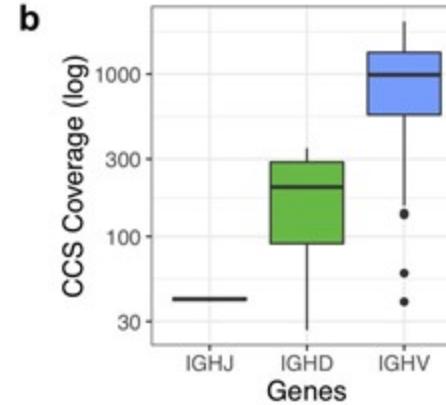
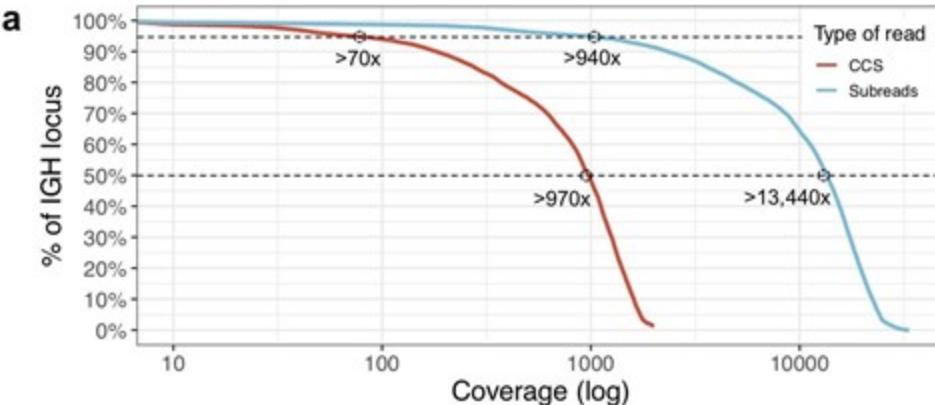
After assembling each haplotype, next step is to detect indels and SVs



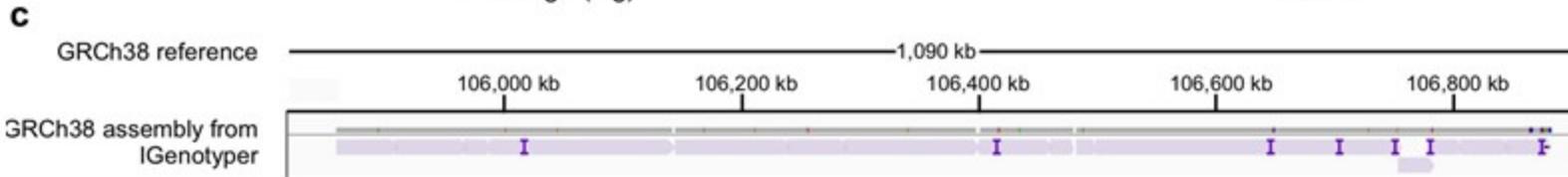
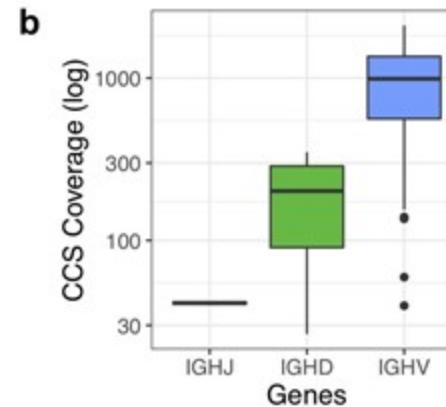
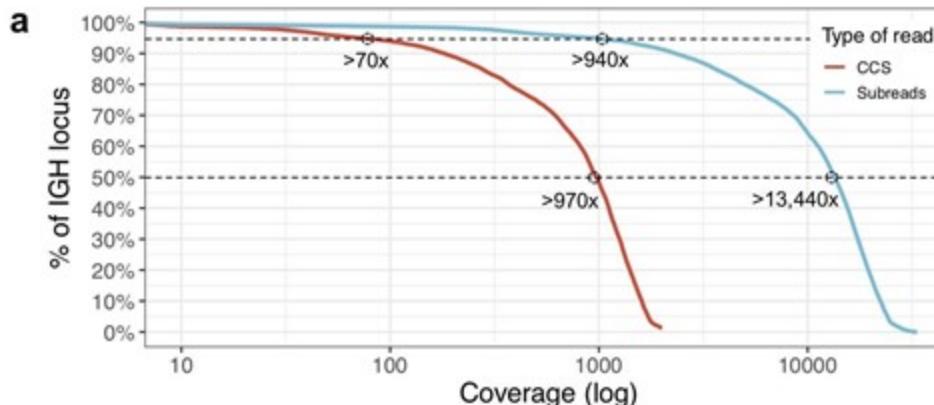
Genes and alleles are directly annotated from assemblies



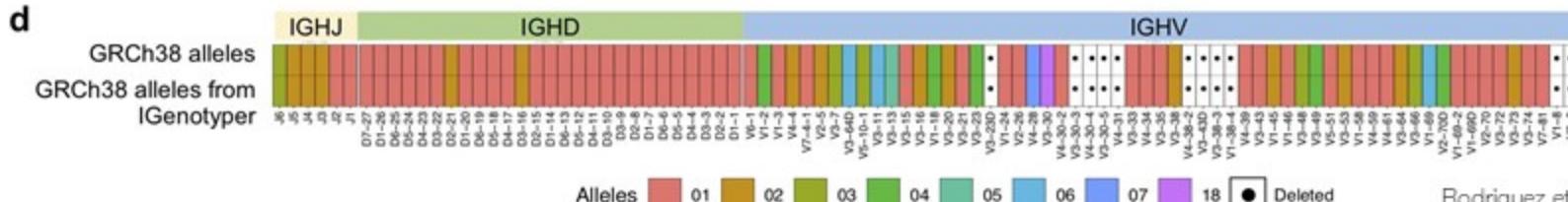
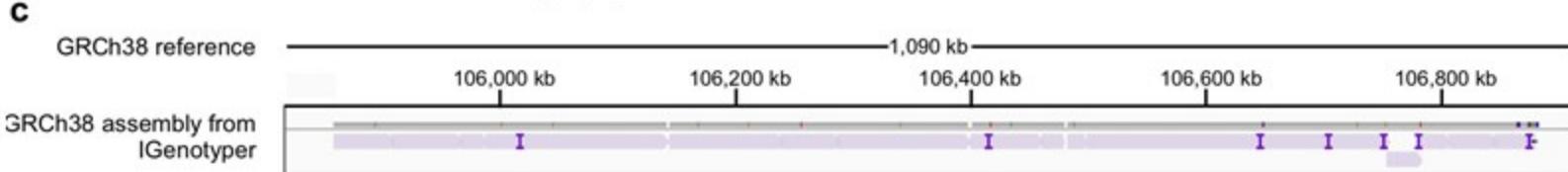
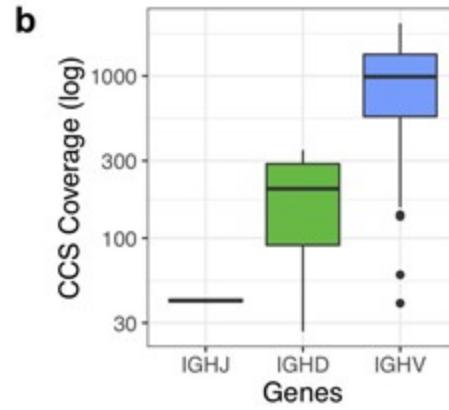
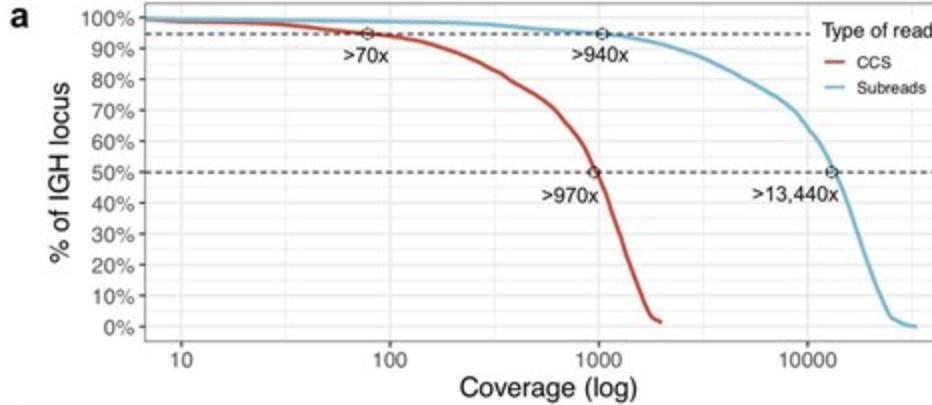
IGH locus is captured at an extremely high coverage



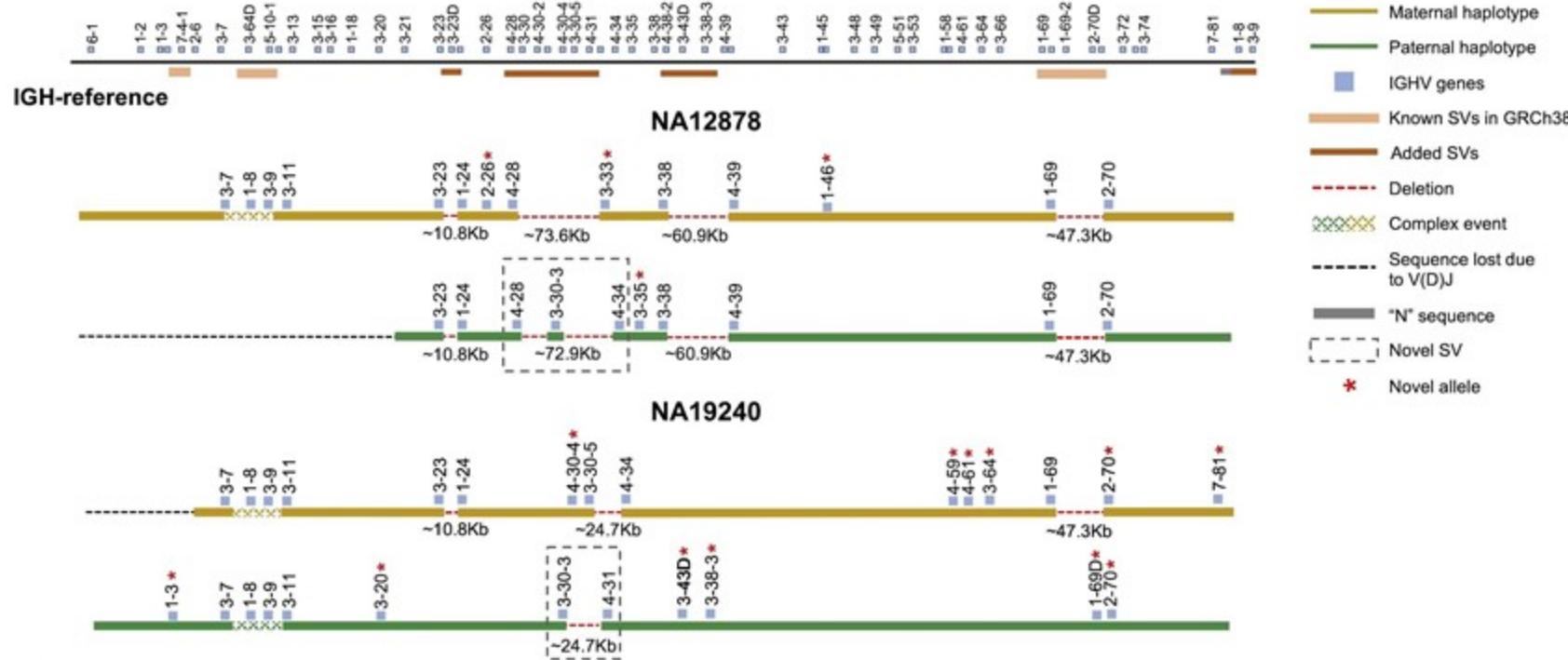
IGenotyper assembled 98.7% of IGH with only 37 errors producing an accuracy > 99.99%



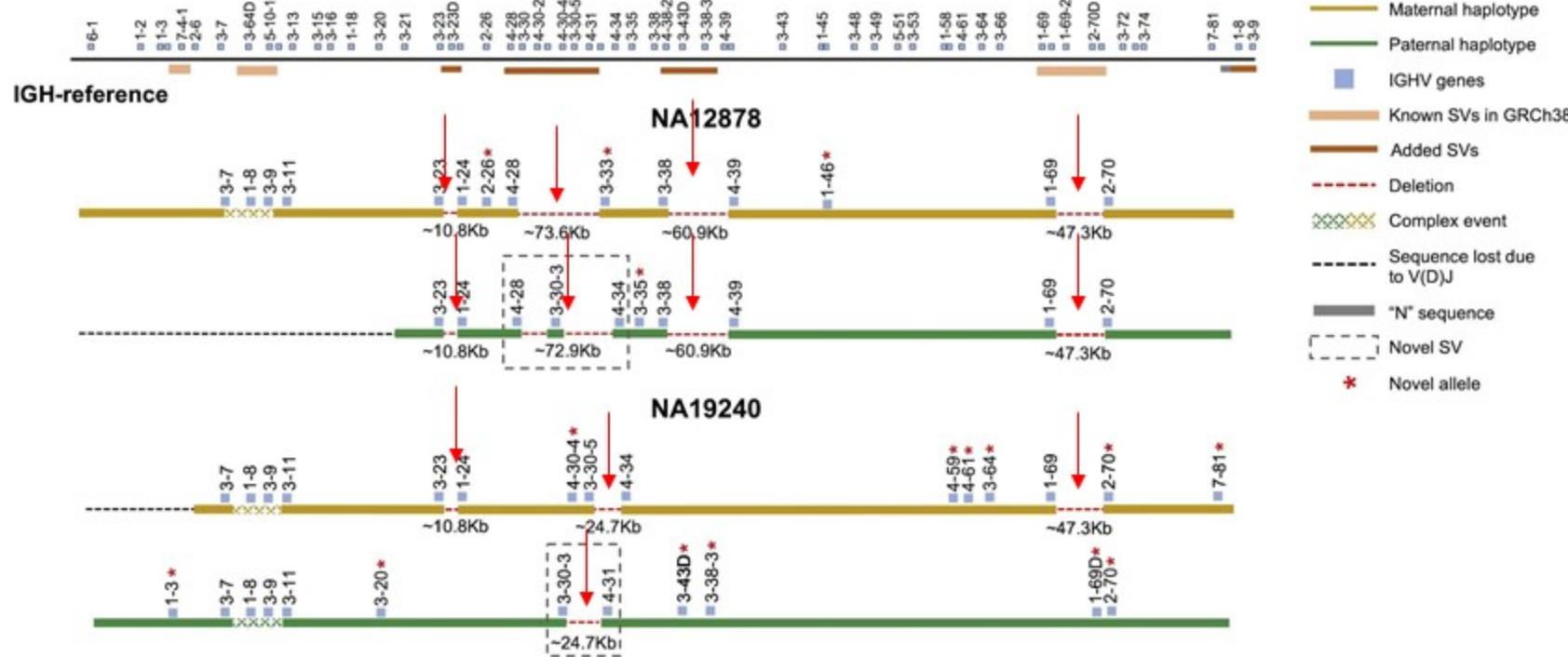
All alleles are correctly assigned



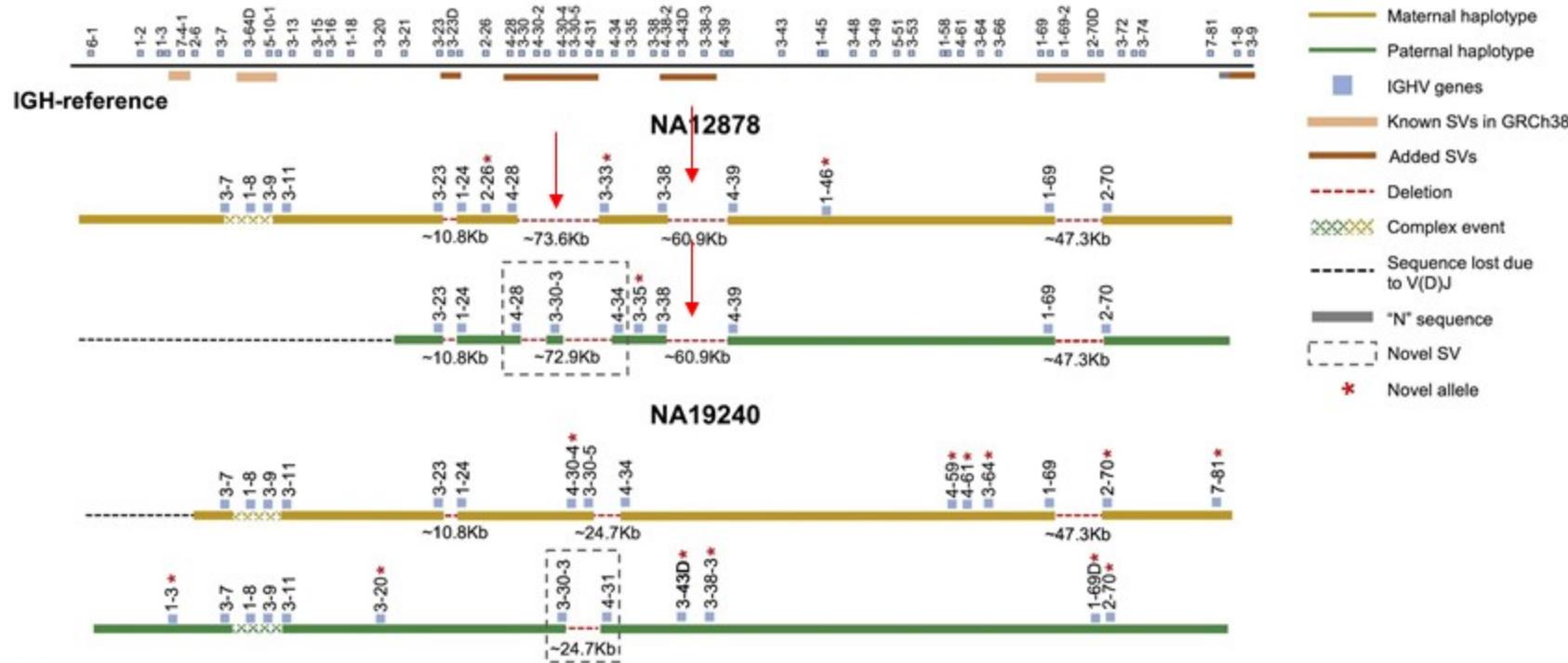
Resolving the IGH locus in two diploid samples



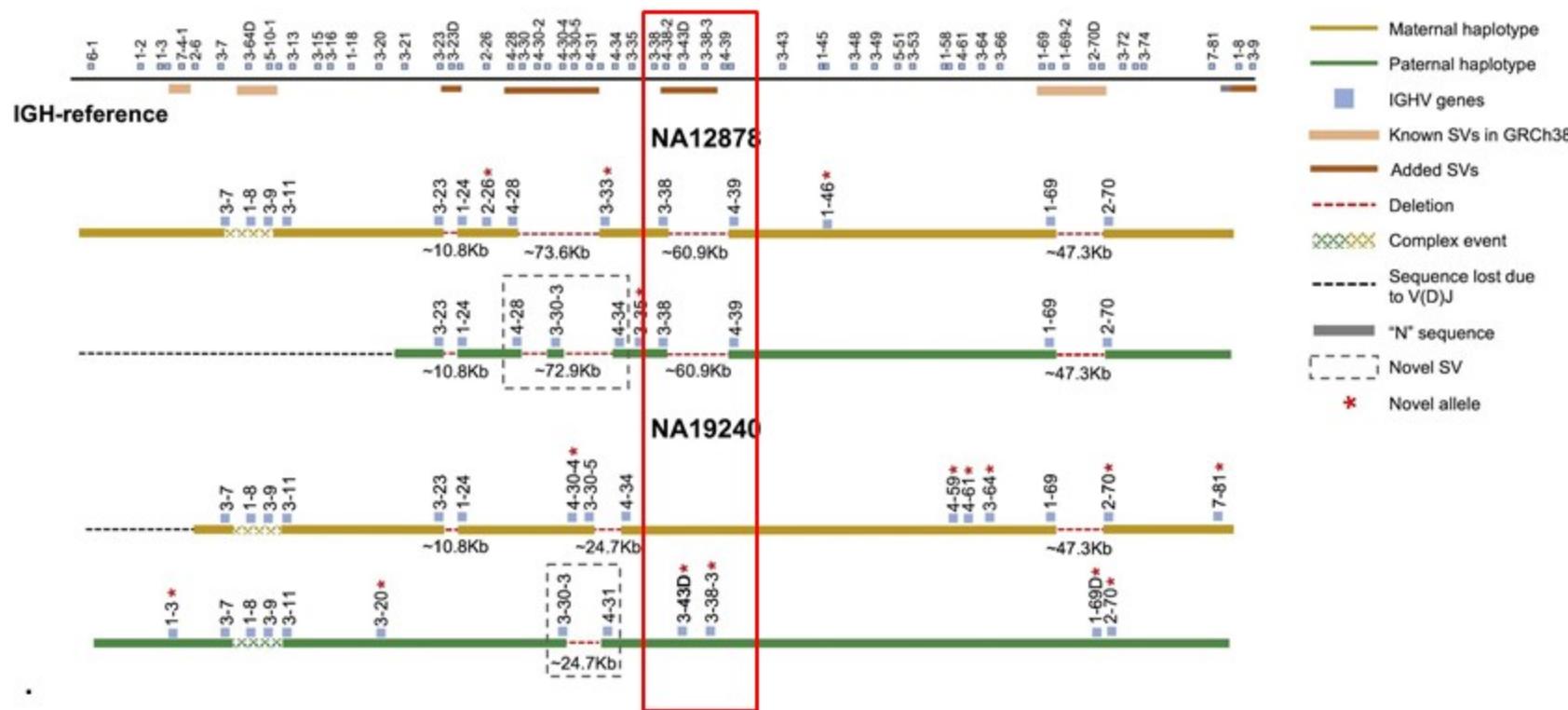
Resolving the IGH locus in two diploid samples identifies several SVs



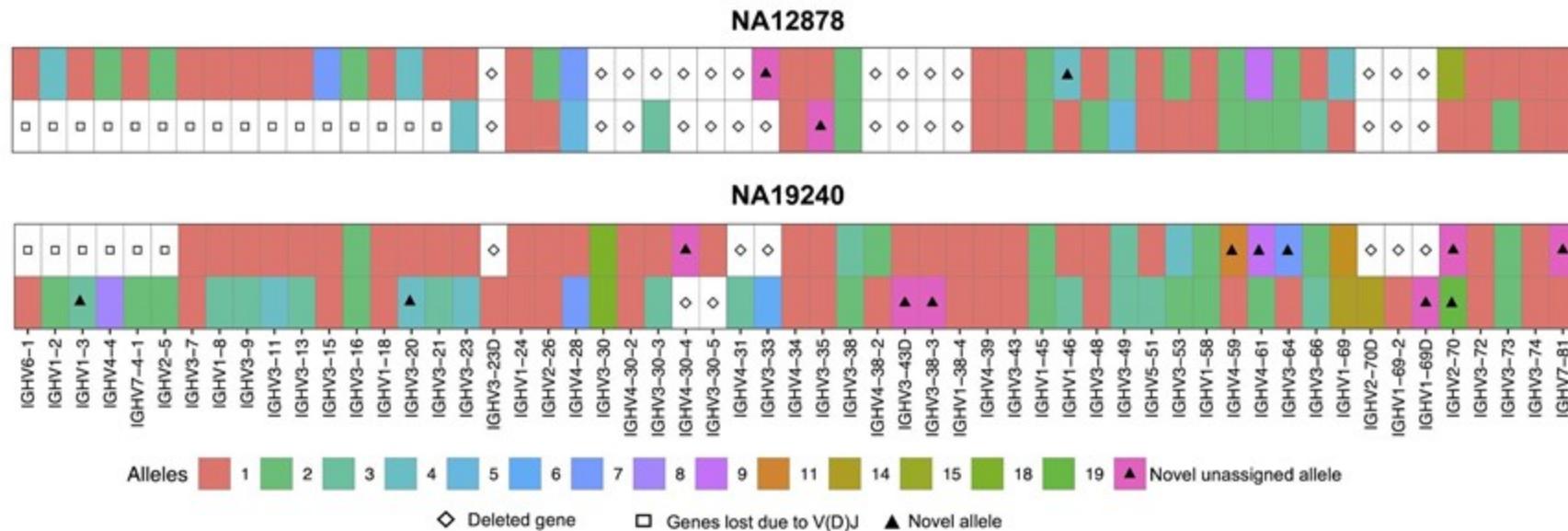
Resolving the IGH locus in two diploid samples identifies large SVs



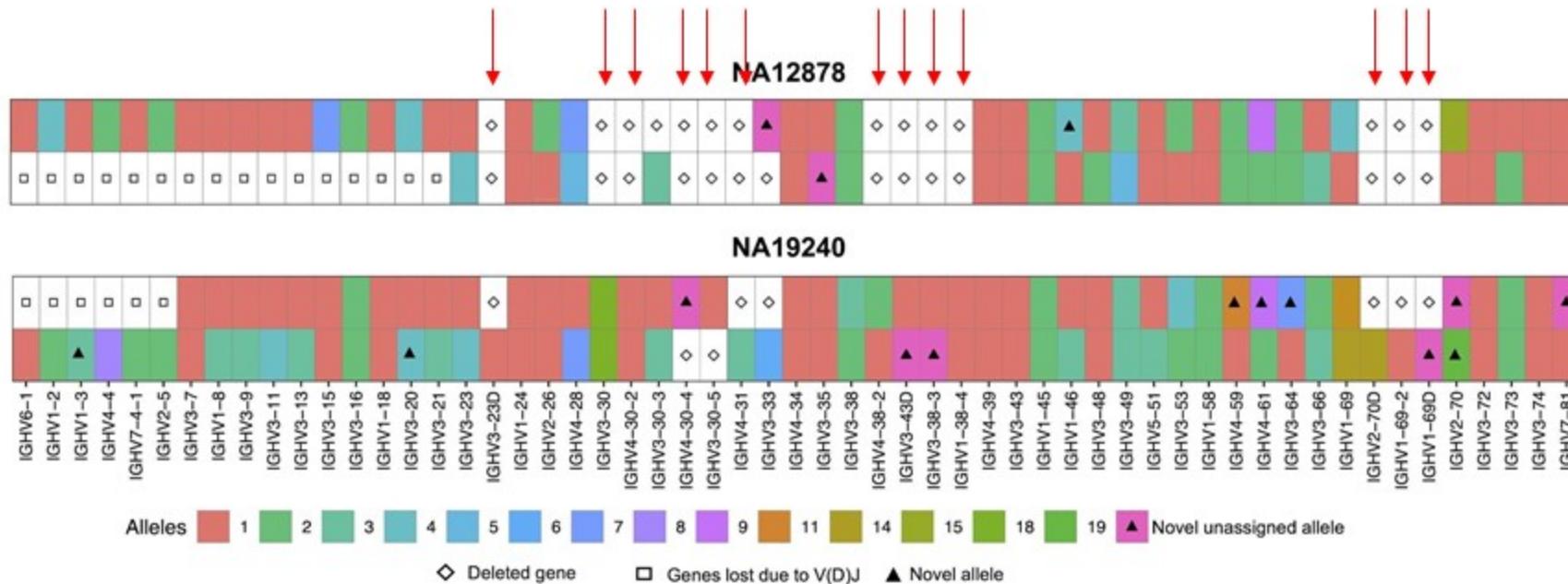
Resolving the IGH locus in two diploid samples identifies inter-individual haplotype variation



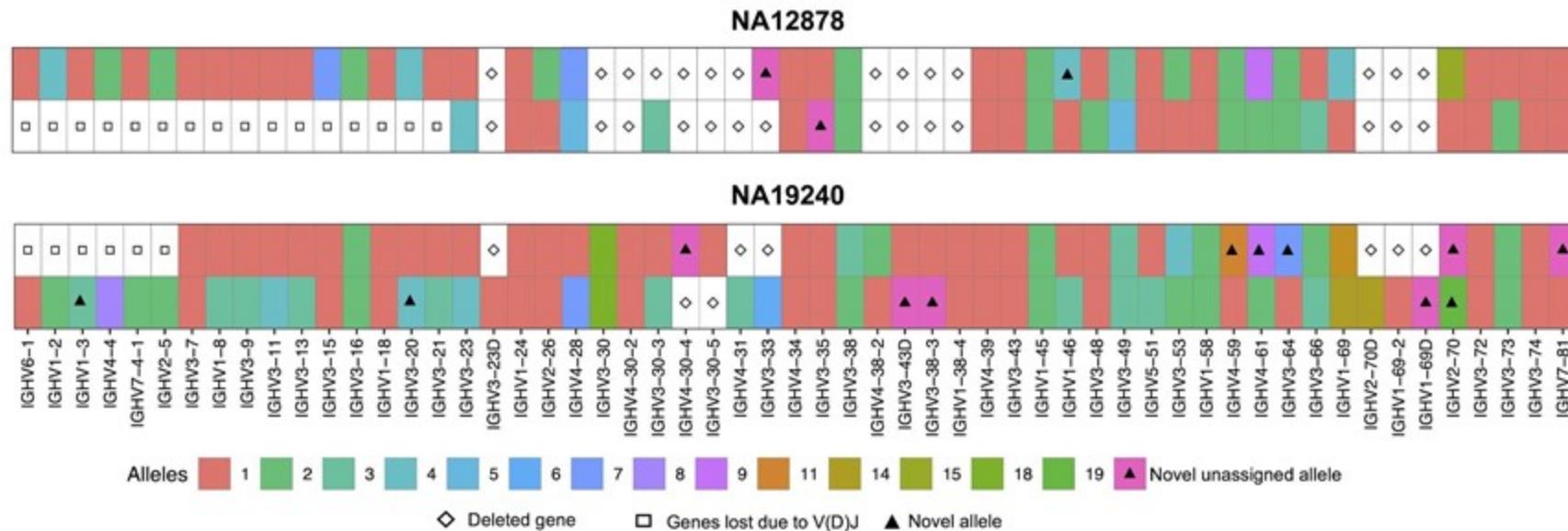
Alleles from each haplotype are resolved using the IGH assemblies



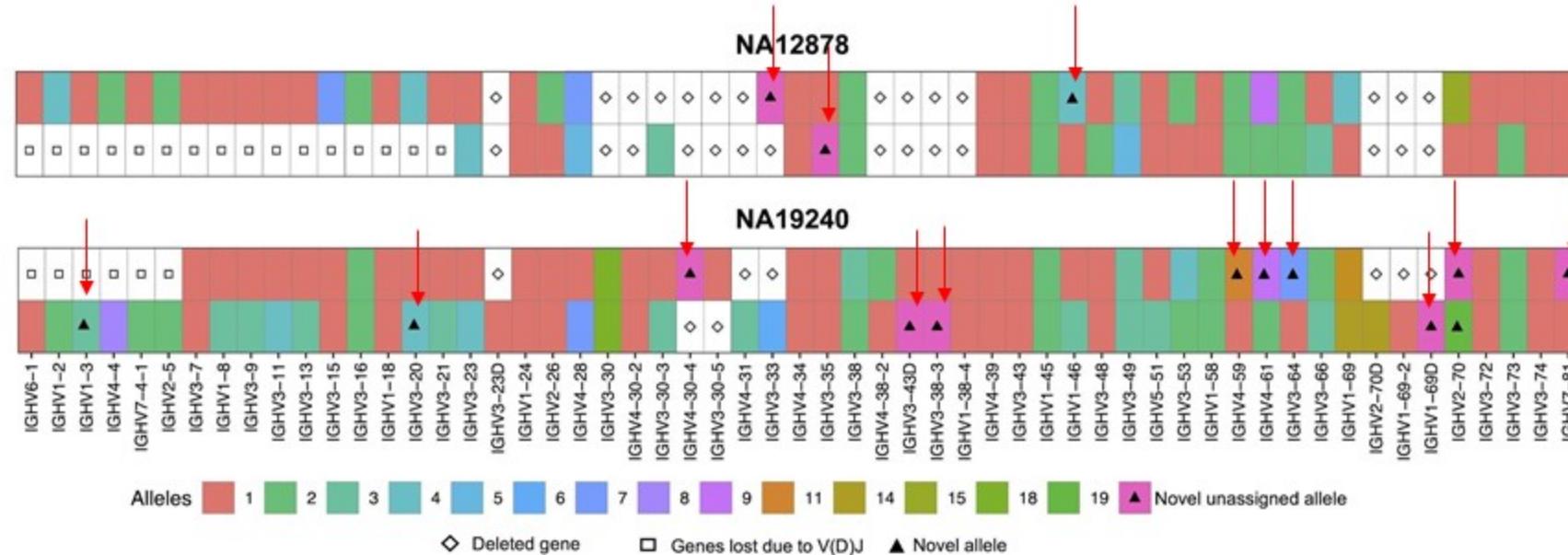
Alleles from each haplotype are resolved using the IGH assemblies



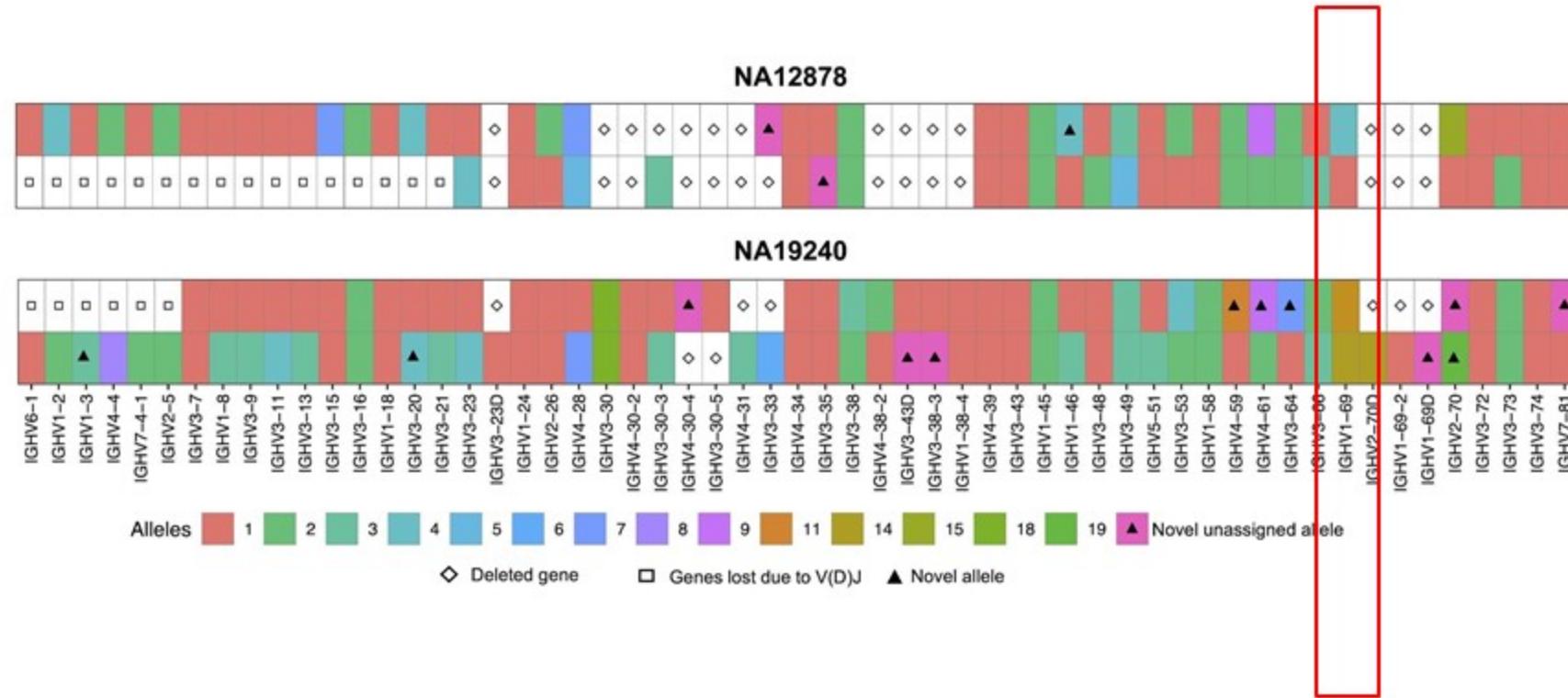
Alleles from each haplotype are resolved using the IGH assemblies



Alleles from each haplotype are resolved using the IGH assemblies



Alleles from each haplotype are resolved using the IGH assemblies



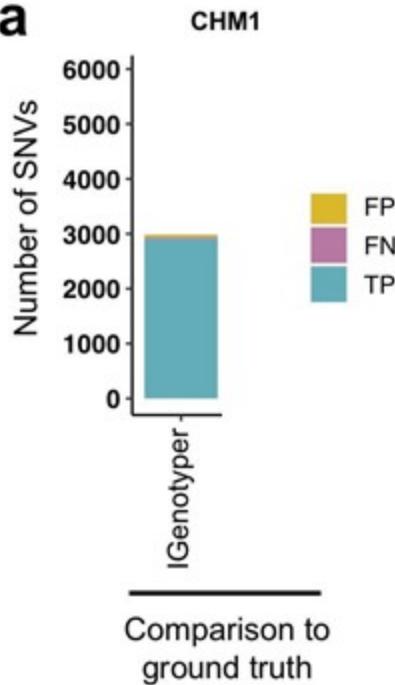
IGH locus is accurately resolved

TABLE 1 | Assembly statistics and evaluation of the accuracy of the haplotype-specific assemblies.

Sample	Contigs (n)	Assembly size (bp)	Assembly validation		
			Concordance with fosmids (SMRT sequencing)	Concordance with BACs or fosmids (Sanger sequencing)	Concordance with Pilon/Illumina
CHM1	16	1,026,385	NA	99.996%	NA
NA19240	38	1,829,616	99.996%	99.99%	99.99%
NA12878	45	1,442,310	99.995%	100.0%	99.99%

IGenotyper SNVs are more accurate than short-read data SNVs

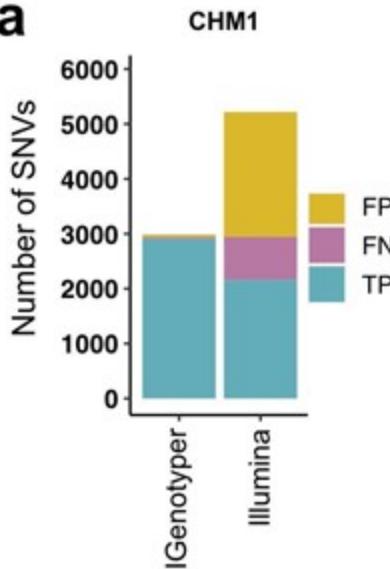
a



1. 2,958 SNVs using IGenotyper
 - a. 2,912 (99.0%) true SNVs
 - b. 46 false-positive SNVs

IGenotyper SNVs are more accurate than short-read data SNVs

a

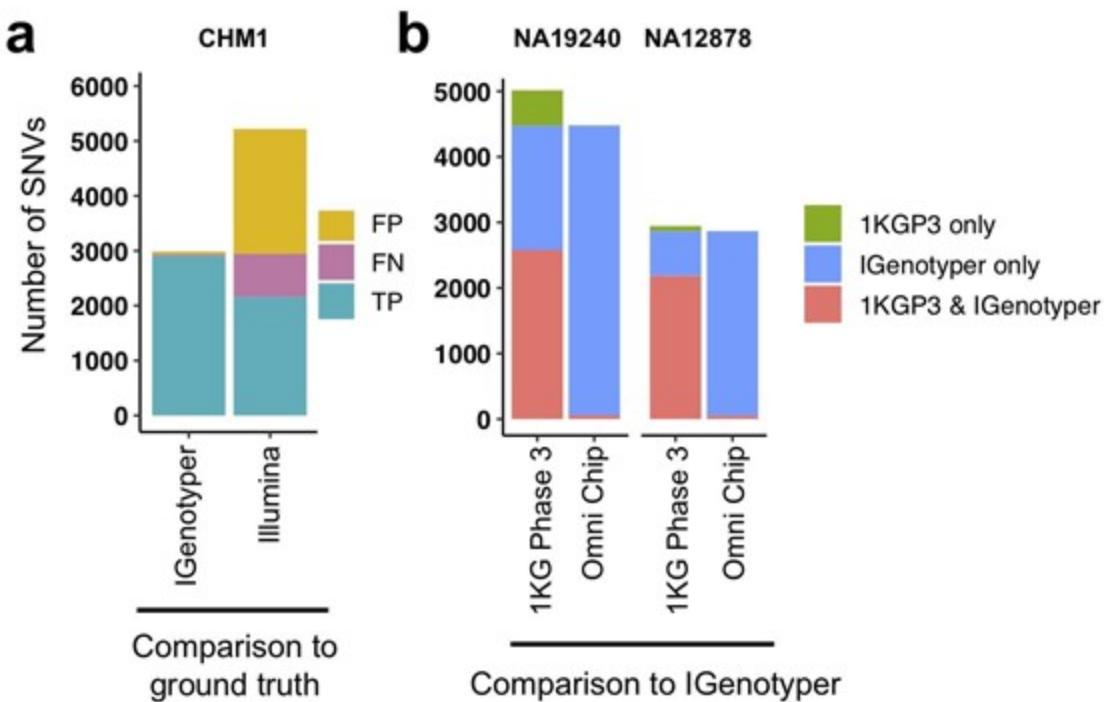


1. 2,958 SNVs using IGenotyper
 - a. 2,912 (99.0%) true SNVs
 - b. 46 false-positive SNVs

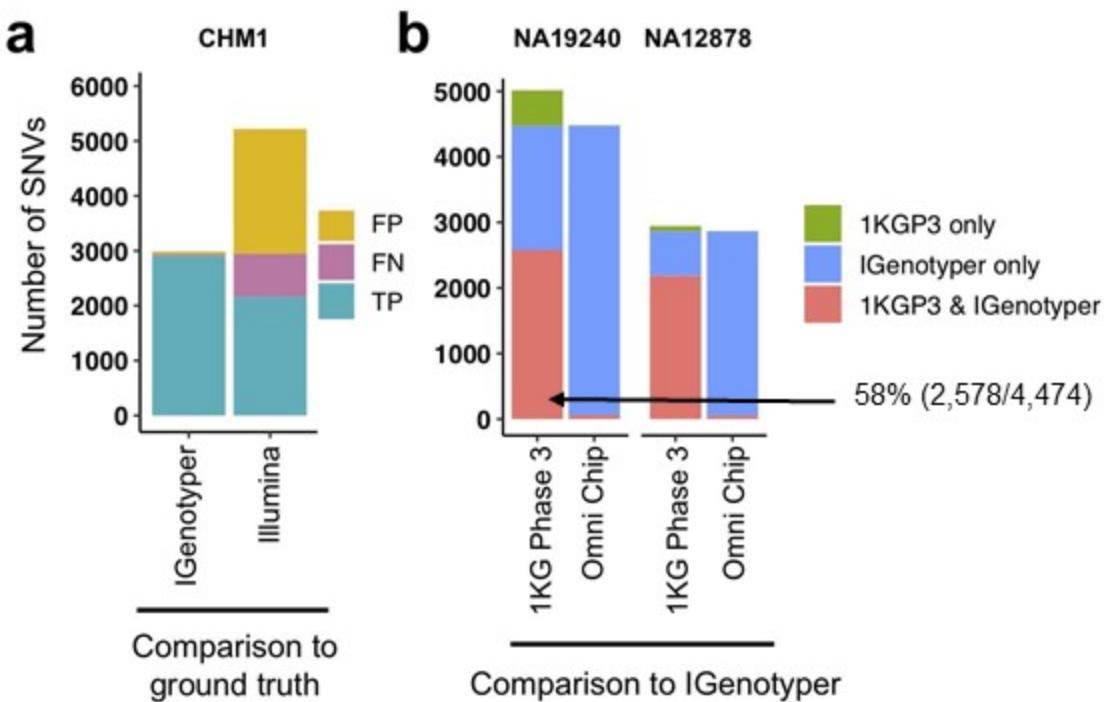
1. 4,433 SNVs from short read data
 - a. 2,159 (49%) true SNVs
 - b. 2,274 false-positive SNVs

Comparison to
ground truth

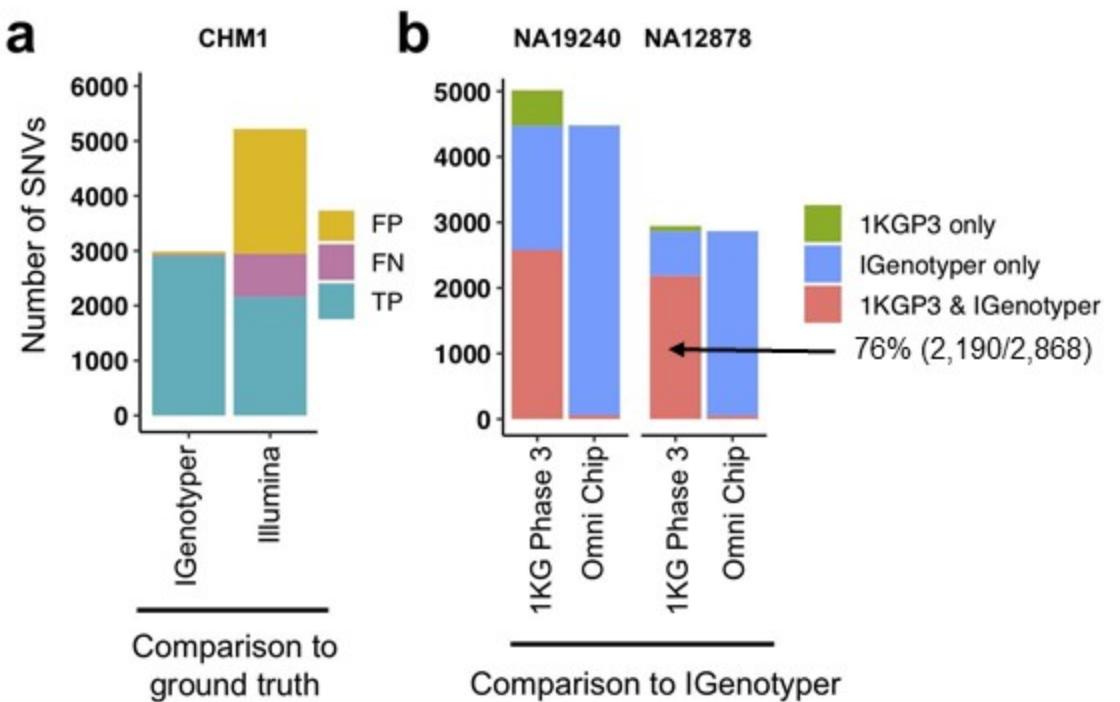
IGenotyper SNVs are more accurate than short-read data SNVs



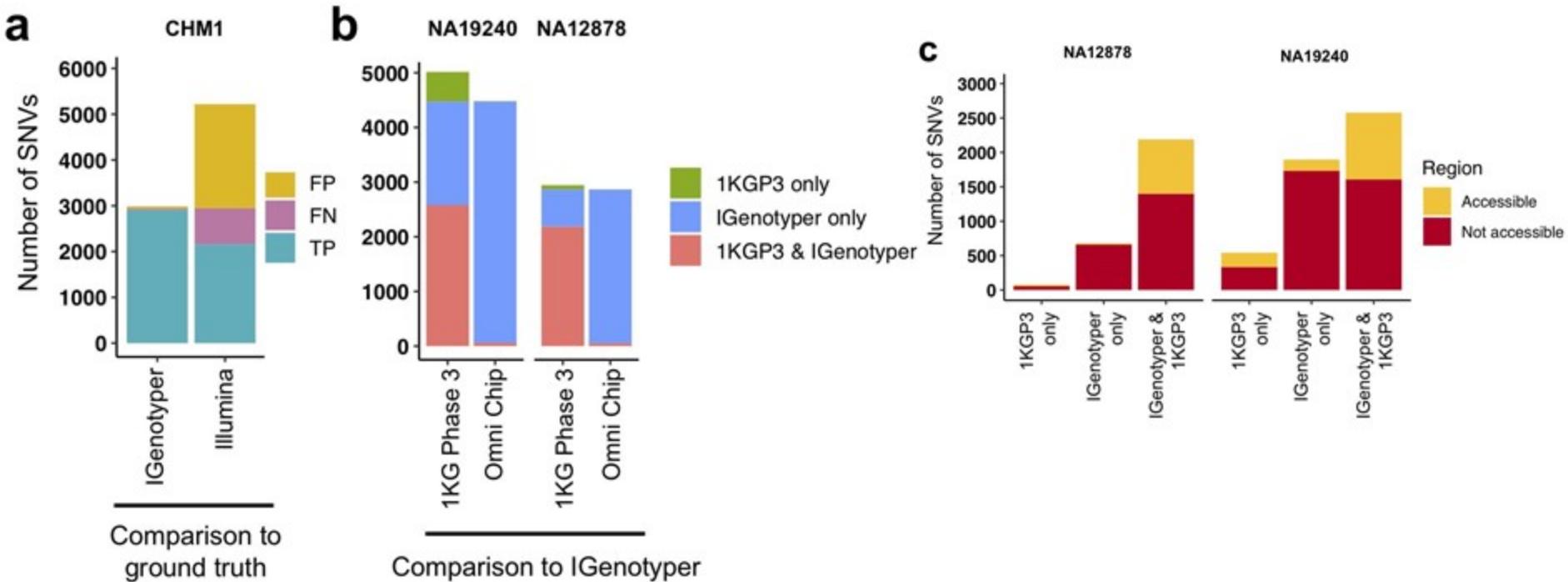
IGenotyper SNVs are more accurate than short-read data SNVs



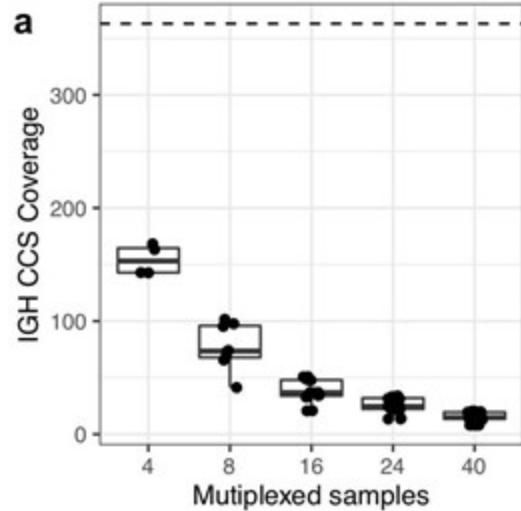
IGenotyper SNVs are more accurate than short-read data SNVs



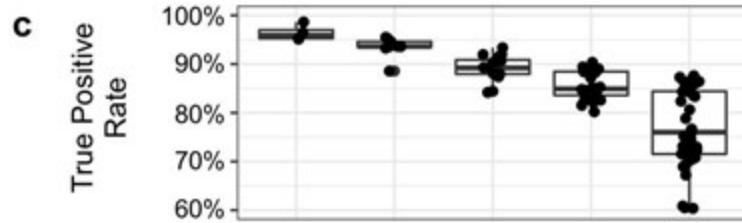
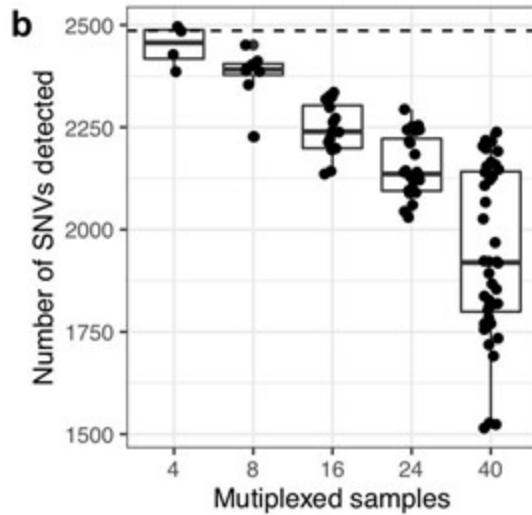
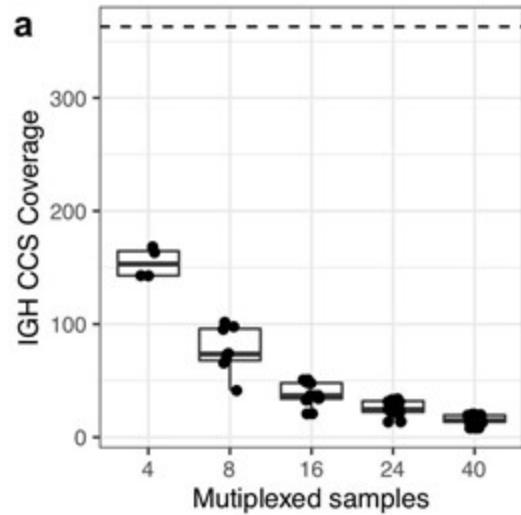
SNVs found just by IGenotyper were within inaccessible regions



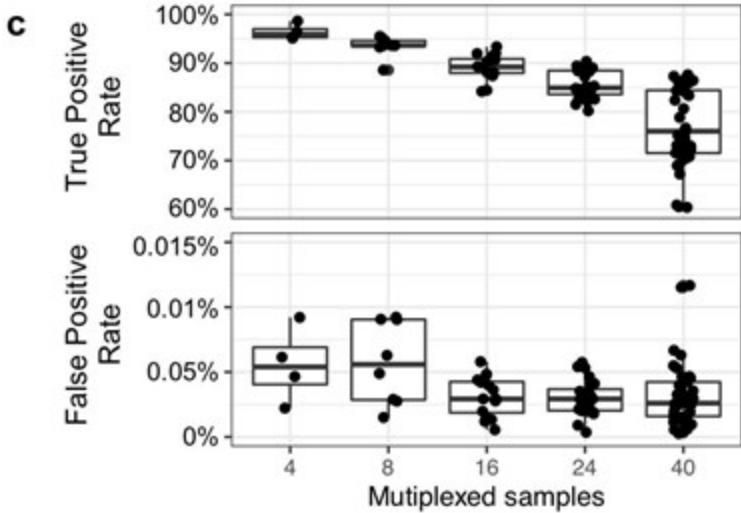
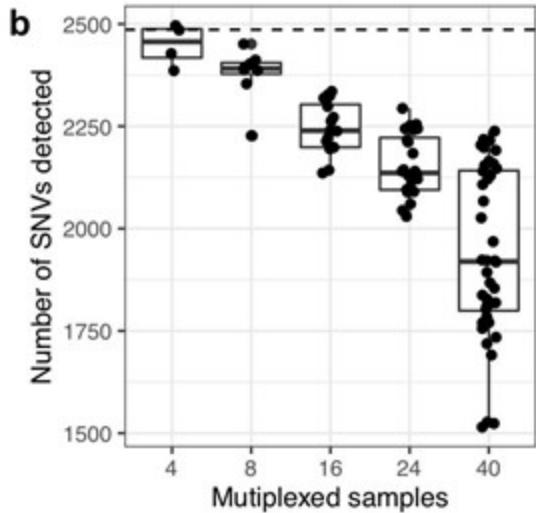
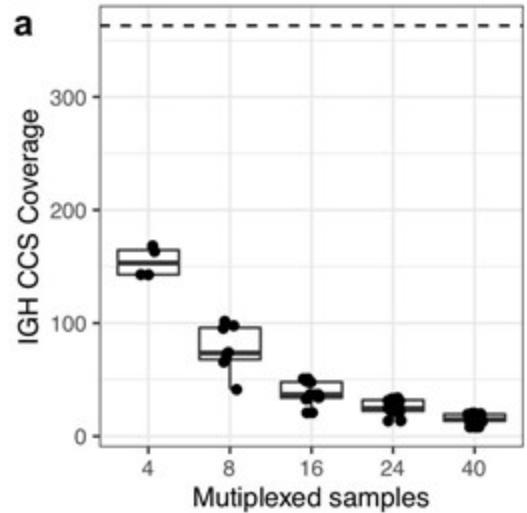
Coverage decreases when multiplexing samples



Less SNVs are detected when multiplexing samples



However, the false positive rate stays consistent



Output from IGenotyper

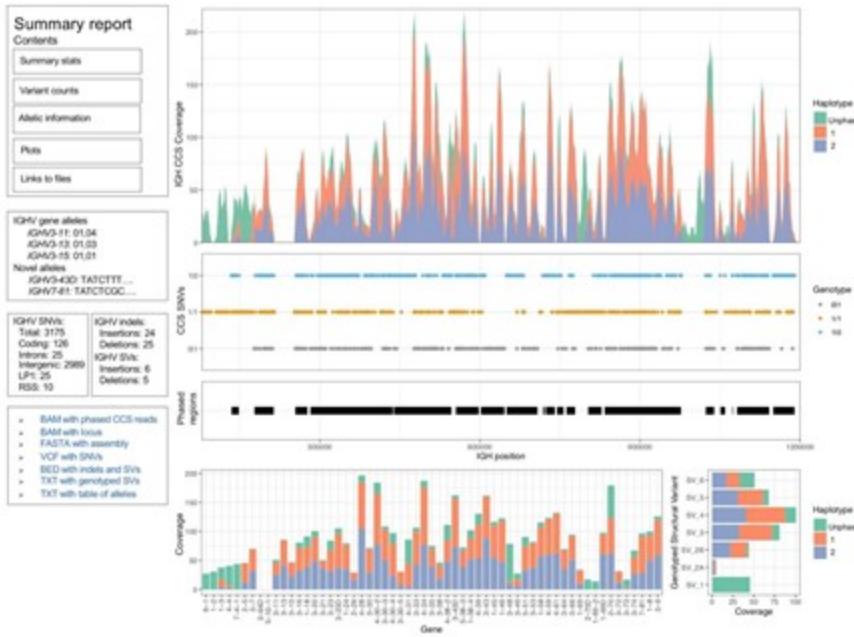
Output directories

alignments alleles assembly logs plots preprocessed report.html tmp variants

Directories	Description
<output>/alignments	Alignments of CCS, subreads and contigs (phased and unphased)
<output>/assembly	Assembly of IGH locus
<output>/variants	SNVs, indels and SVs
<output>/alleles	Alleles in sample
<output>/logs	Log files with input parameters
<output>/tmp	Temporary files. Could be deleted.

Output files

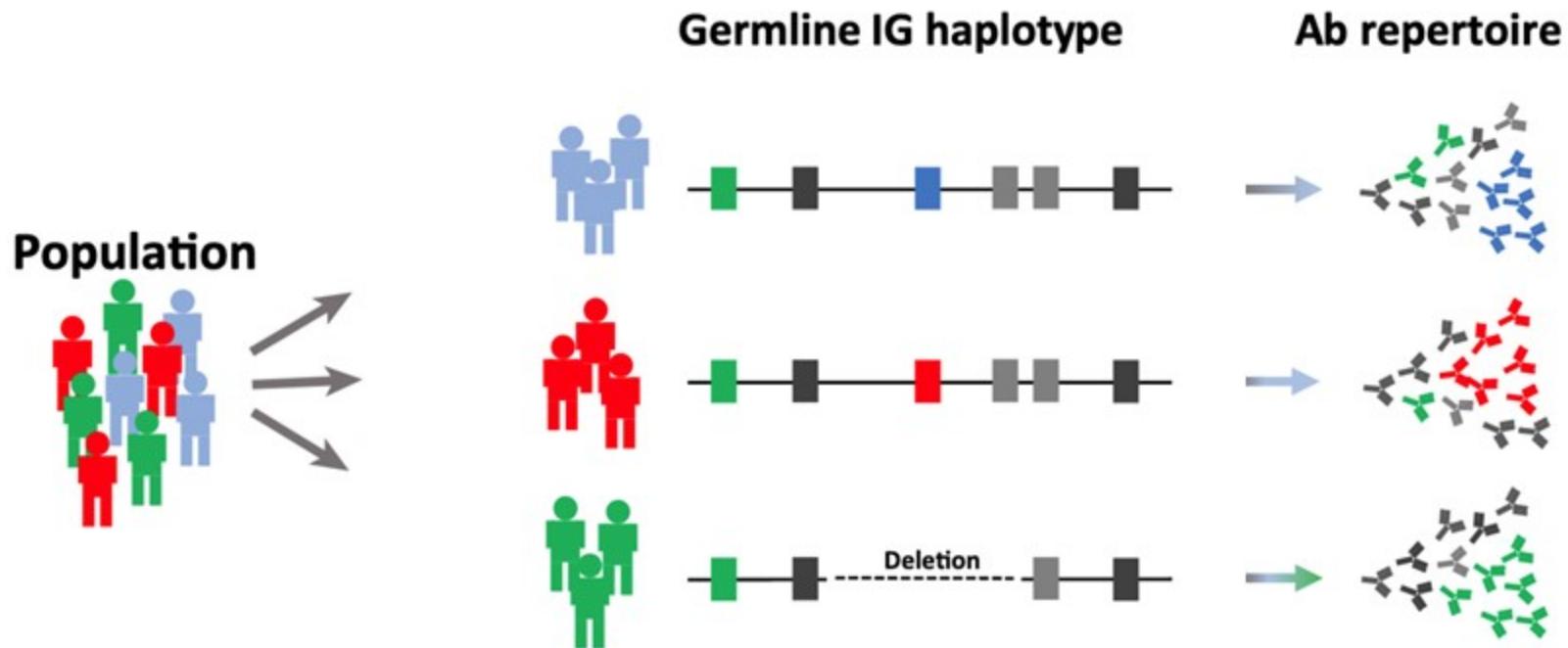
1. alignments/
 - i. ccs_to_ref* : CCS reads aligned to reference
 - ii. contigs_to_ref* : All assembled contigs aligned to reference
 - iii. igh_contigs_to_ref* : IGH assembled contigs aligned to igh reference
2. assembly/
 - i. contigs.fasta : All assembled contigs
 - ii. igh_contigs.fasta : IGH assembled contigs
3. alleles/
 - i. assembly_alleles.bed : Alleles extracted from the assembly for each gene
 - ii. assembly_genes.fasta : Fasta sequence from the assembly for each gene/allele
 - iii. ccs_alleles.bed : Alleles extracted from the CCS reads for each gene
 - iv. ccs_genes.fasta : Fasta sequence from the CCS reads for each gene/allele
4. logs/
 - i. gene_cov.txt : Haplotype coverage for each gene
5. variants/
 - i. snvs_phased_from_ccs.vcf : Phased SNVs detected from the CCS reads
 - ii. snvs_assembly.vcf : SNVs detected from the assembly
 - iii. indel_assembly.bed : Indels detected from the assembly
 - iv. sv_assembly.bed : SVs detected from the assembly
 - v. phased_blocks.txt : Phased haplotype blocks



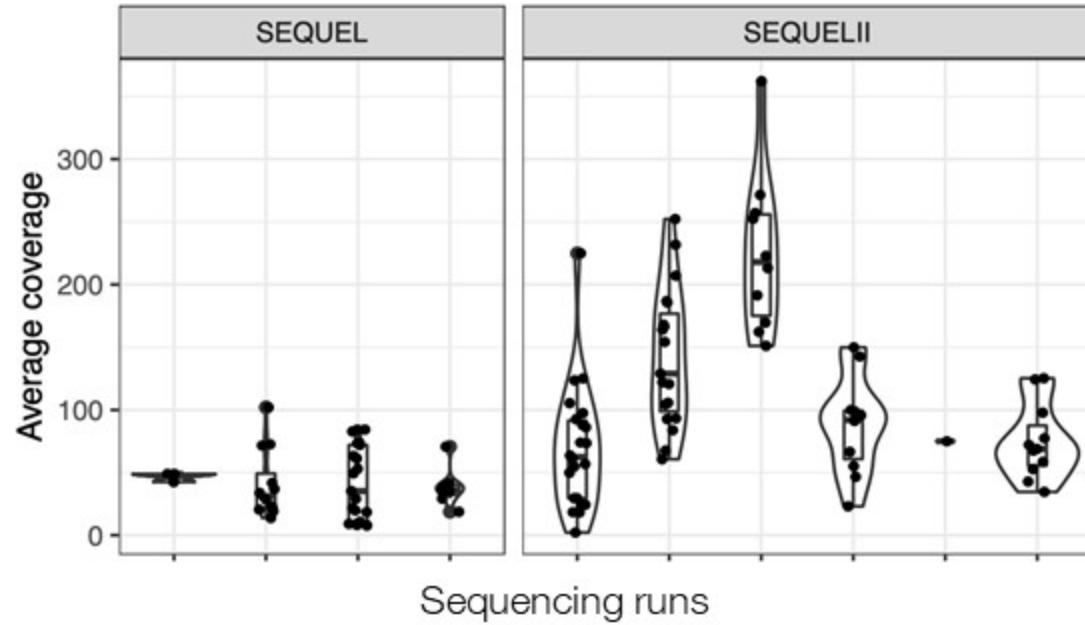
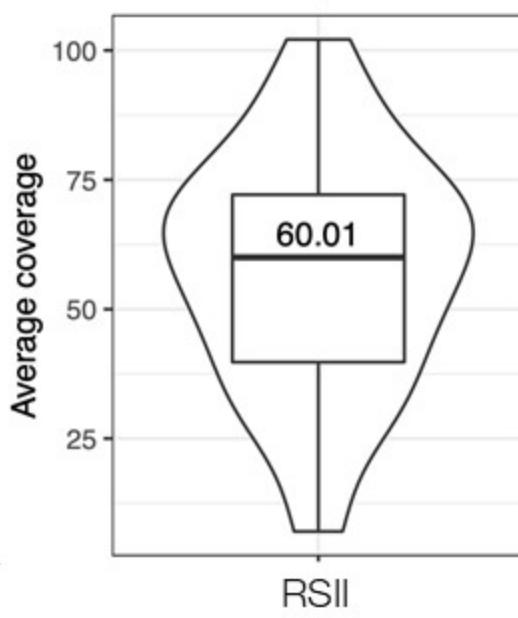
Outline of presentation

- 1) Introduction
 - a) Genetic variation in the immunoglobulin heavy chain locus (IGH)
 - b) Advantages of long-read sequencing over short-read sequencing
- 2) Framework for resolving the immunoglobulin heavy locus in a high throughput fashion using long read sequencing
- 3) Resolving the immunoglobulin heavy chain locus in a cohort with adaptive immune repertoire sequencing data
- 4) Resolving the T-cell receptor locus using long read sequencing
- 5) Application of framework to the immunoglobulin heavy chain locus in rhesus macaque
- 6) Conclusion and perspective

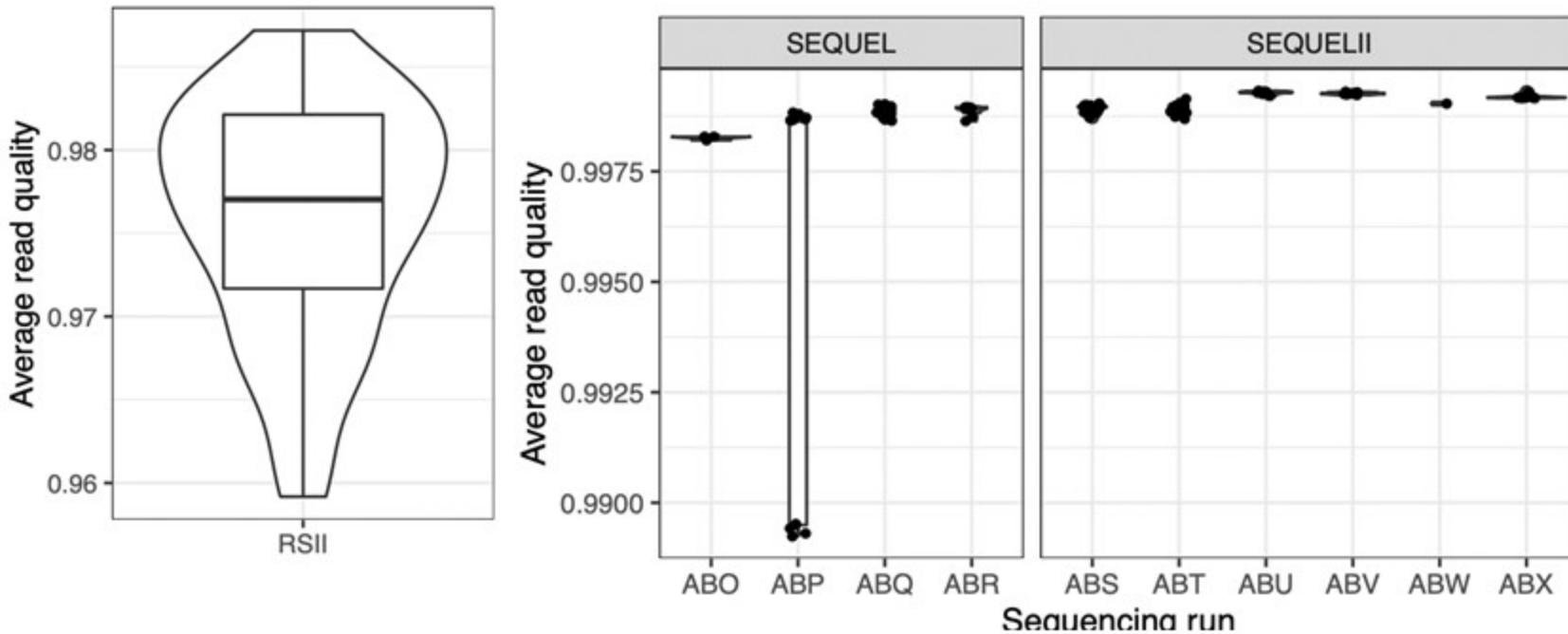
Using framework on a study with a cohort of 154 individuals with Ab repertoire data



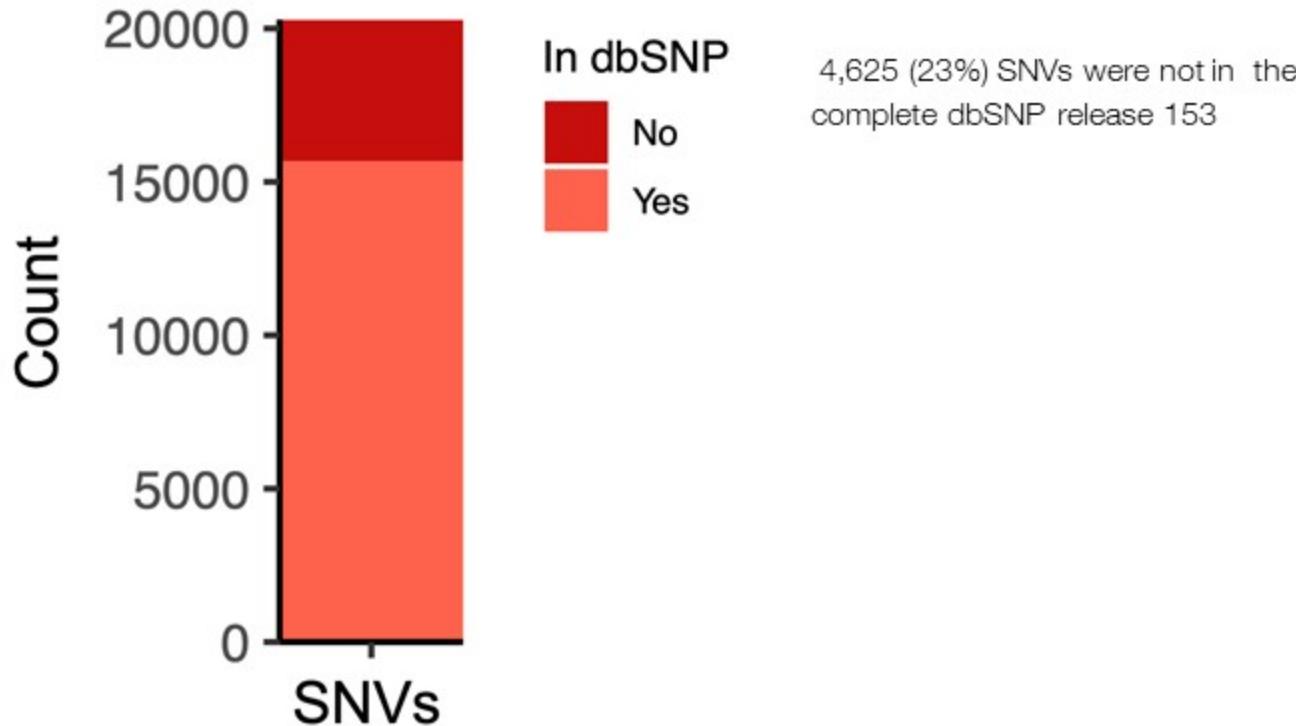
Average long-read sequencing coverage across cohort was 76X



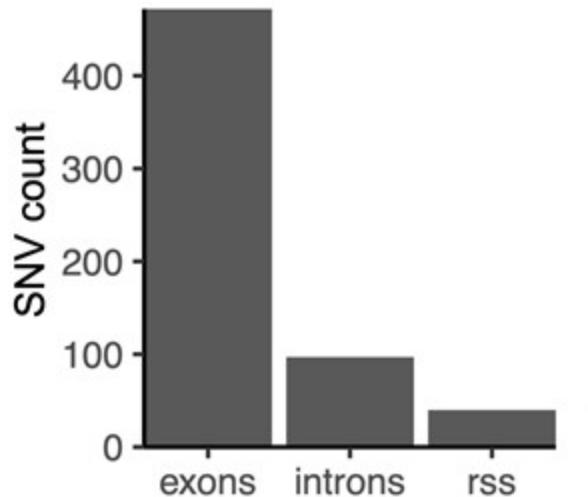
Largest improvement is observed in read quality



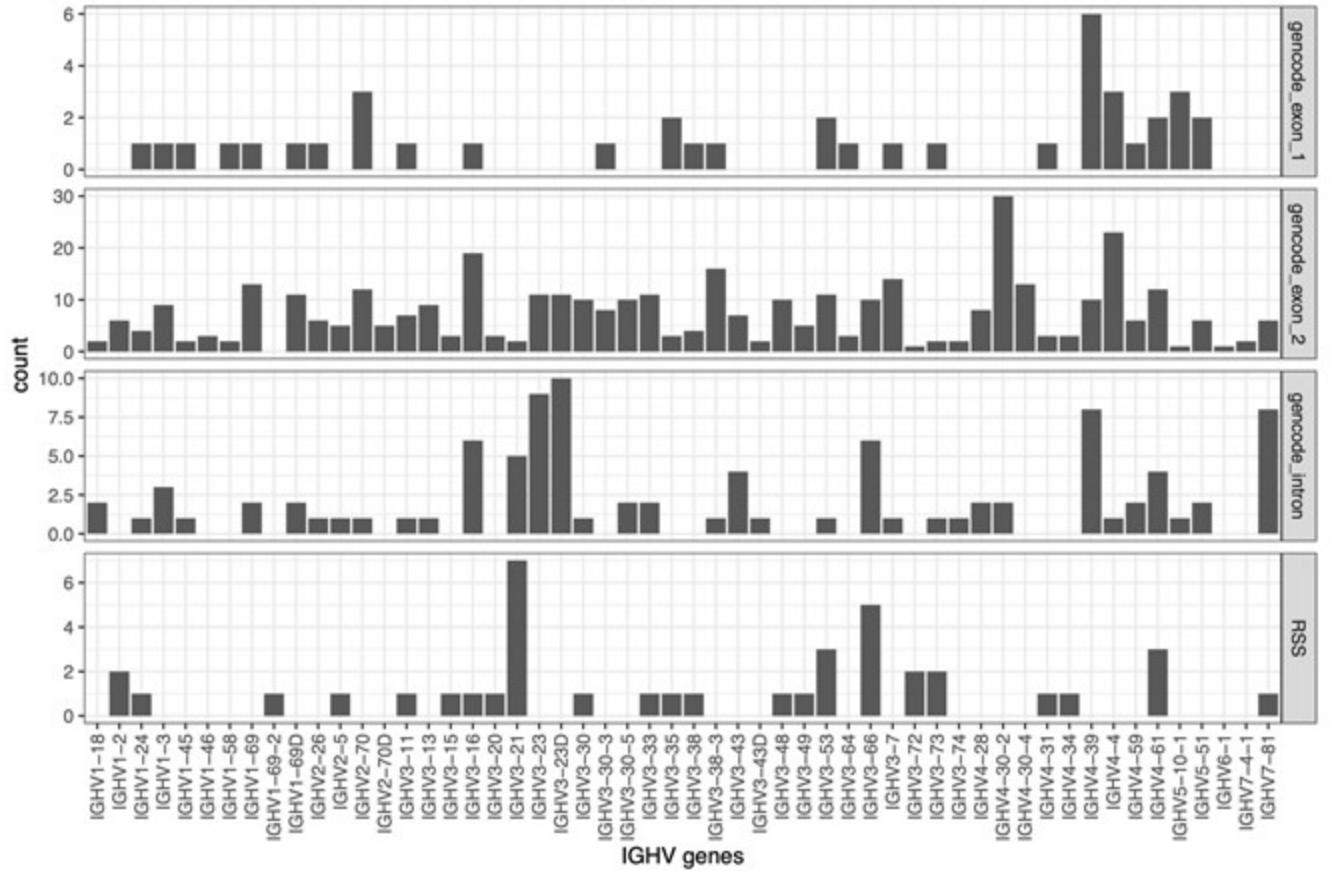
A total of 20,510 SNVs were found in one or more individuals



SNVs are found in different gene components



SNVs are found in different gene components



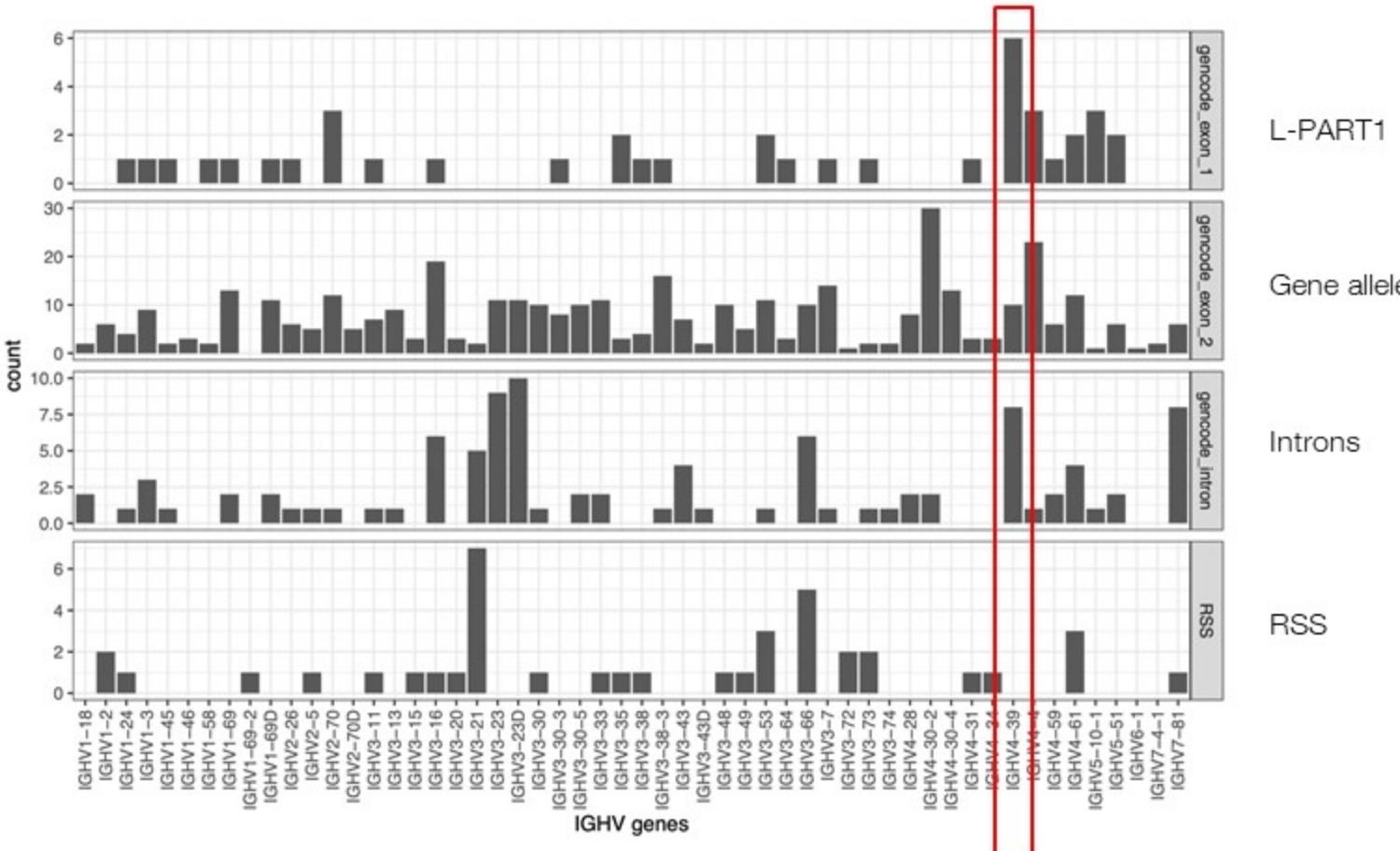
L-PART1

Gene alleles

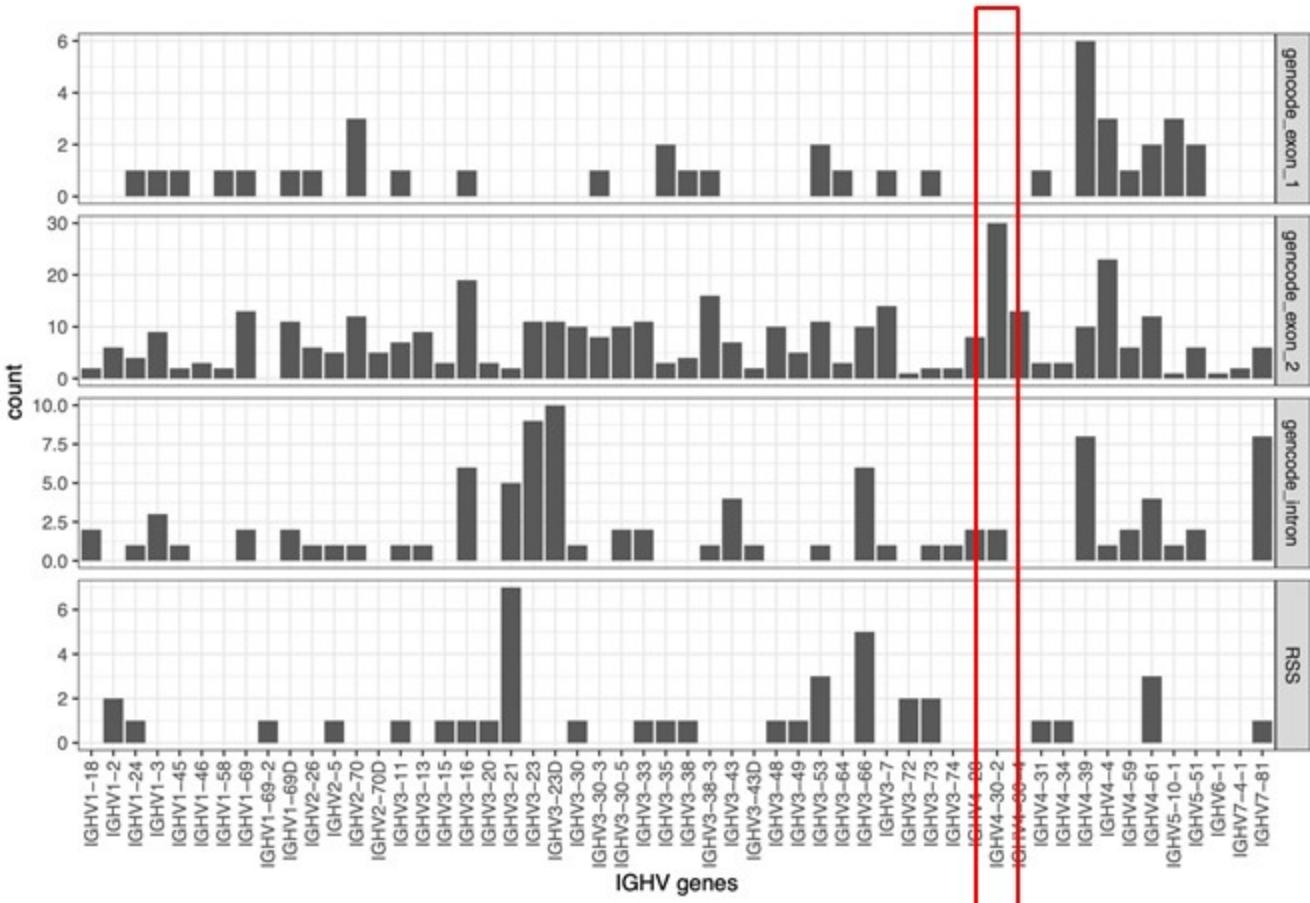
Introns

RSS

SNVs are found in different gene components



SNVs are found in different gene components



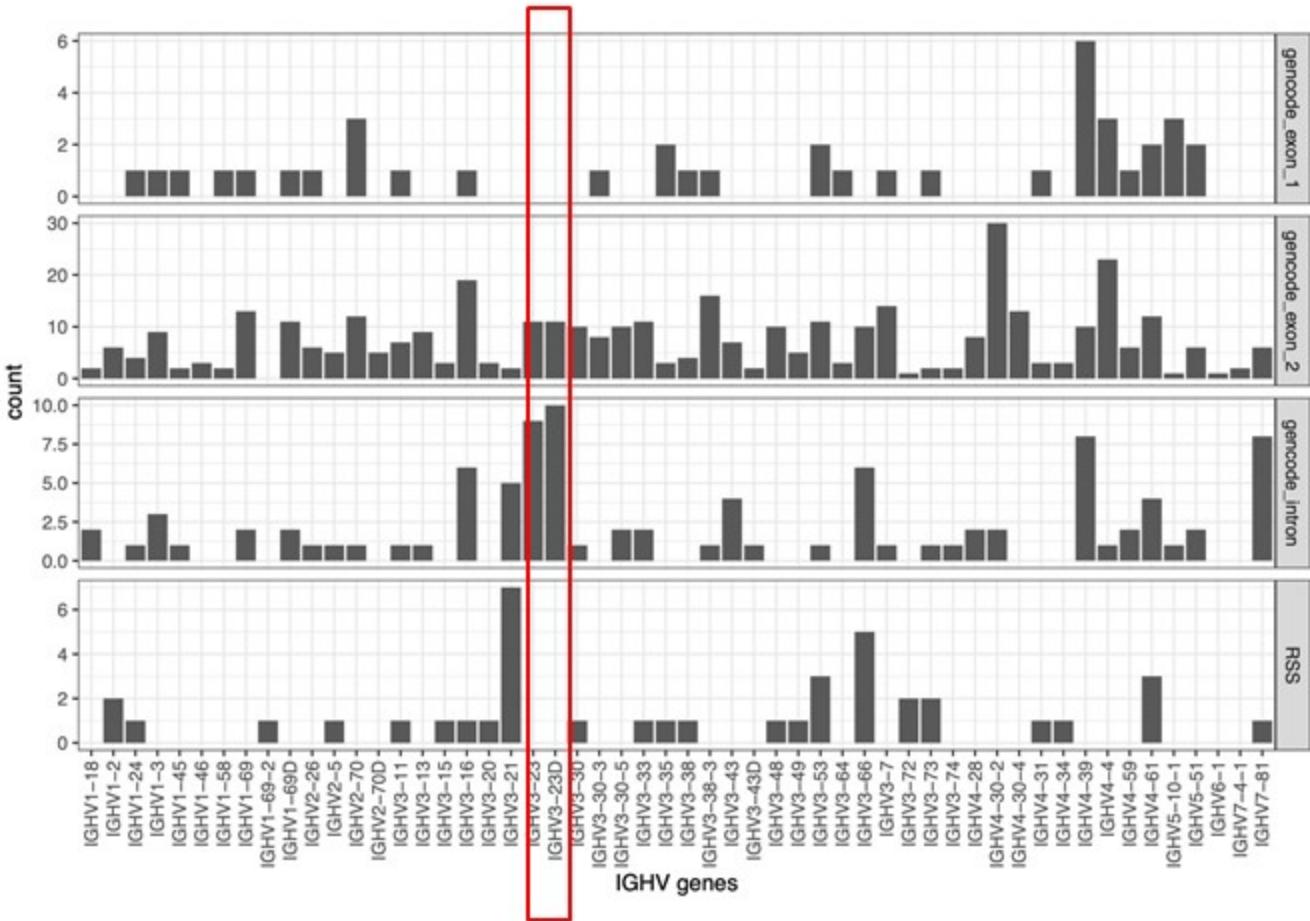
L-PART1

Gene alleles

Introns

RSS

SNVs are found in different gene components



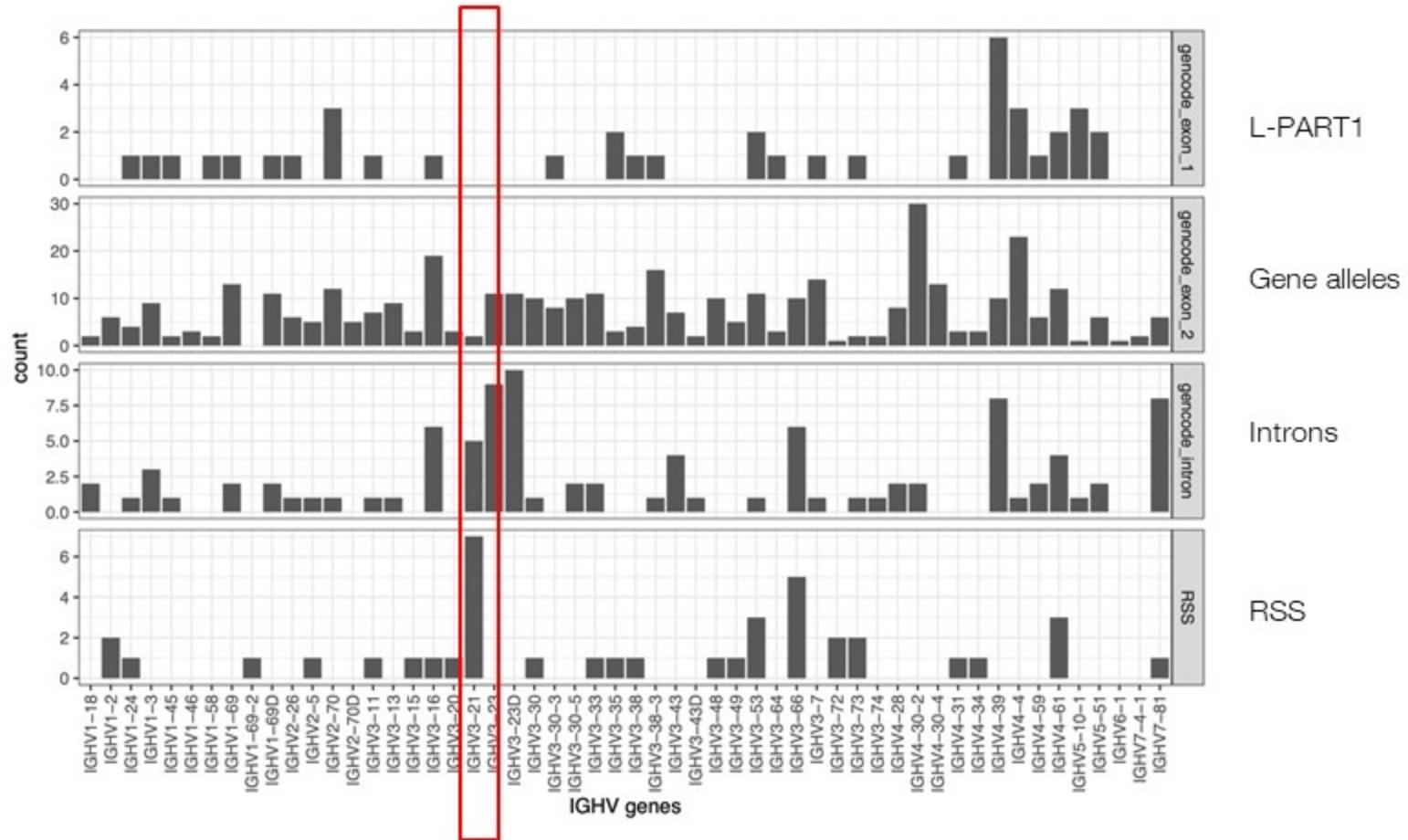
L-PART1

Gene alleles

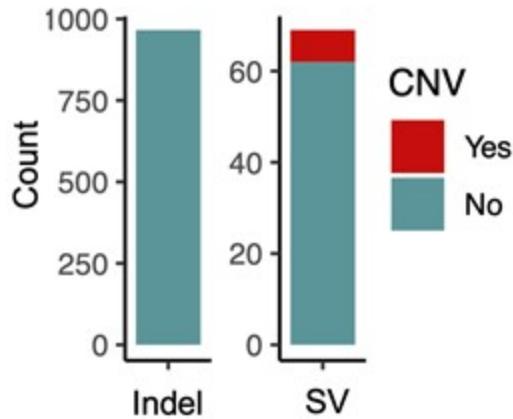
Introns

RSS

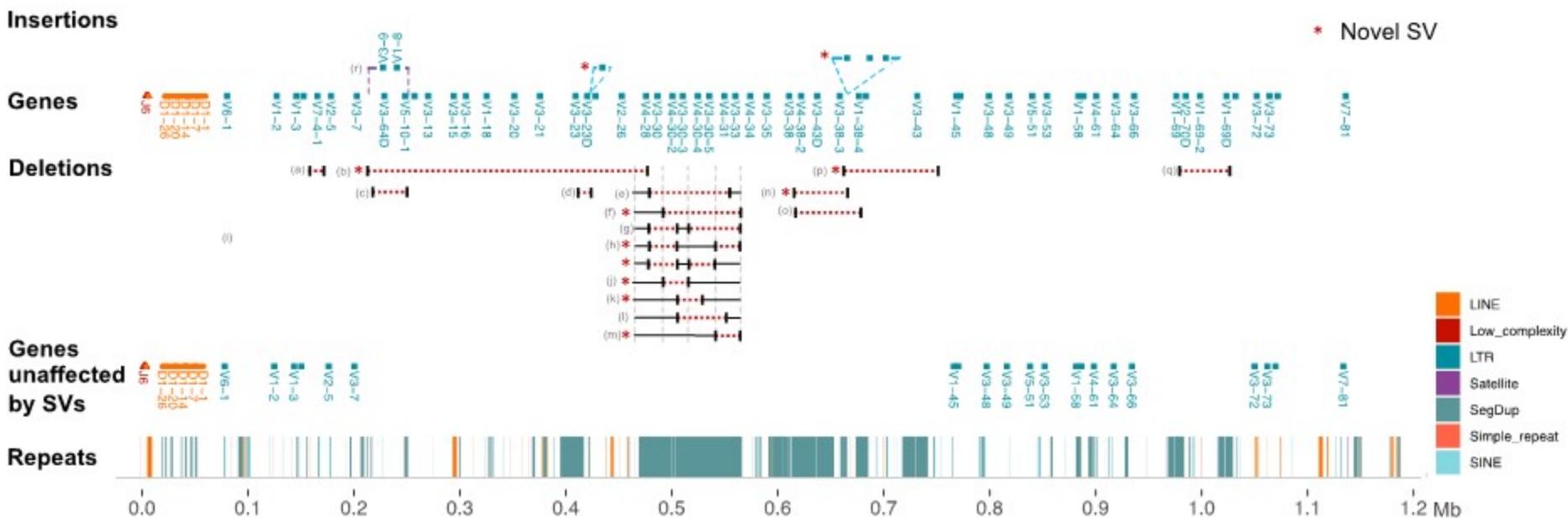
SNVs are found in different gene components



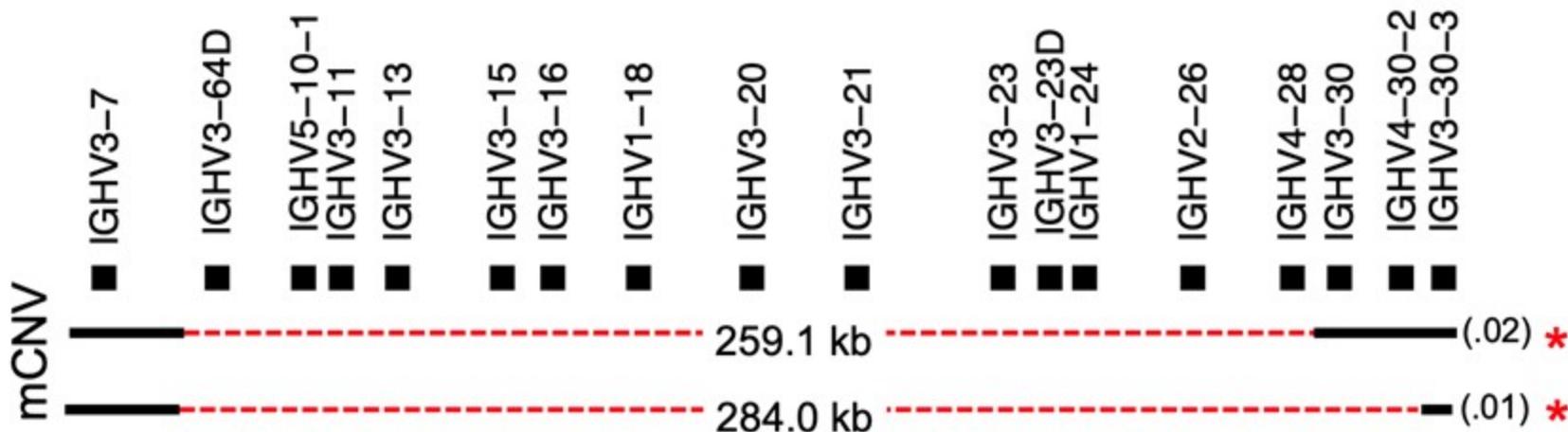
Large number of indels (966) and SVs (71) were also identified



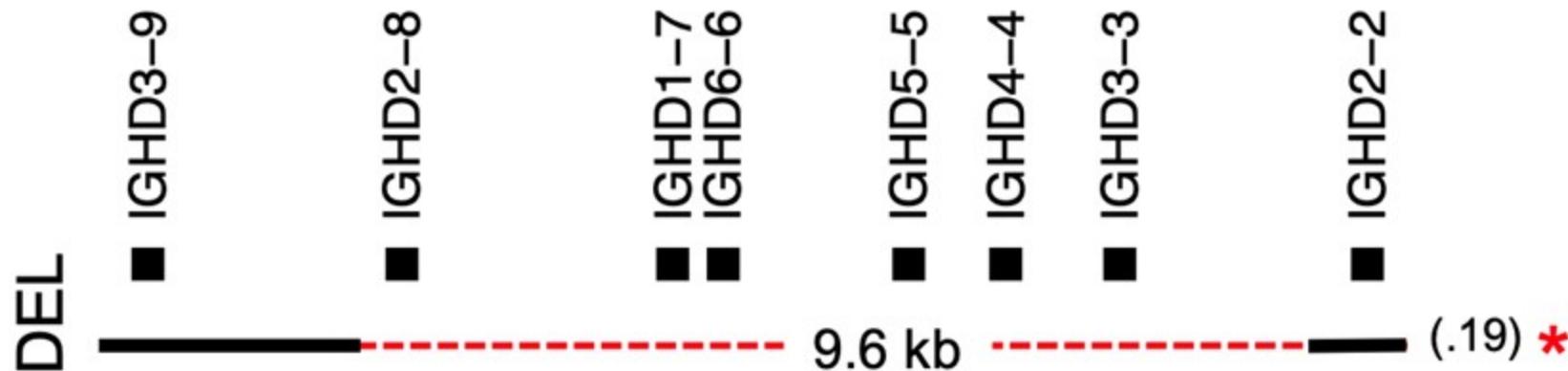
8 large SVs genotyped that affect gene copy number



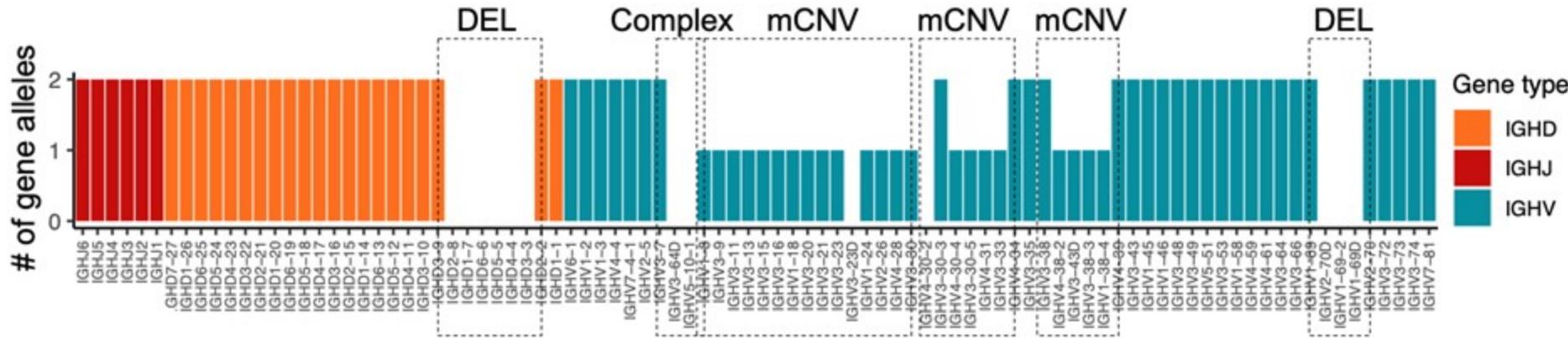
SVs affecting gene copy number were also resolved



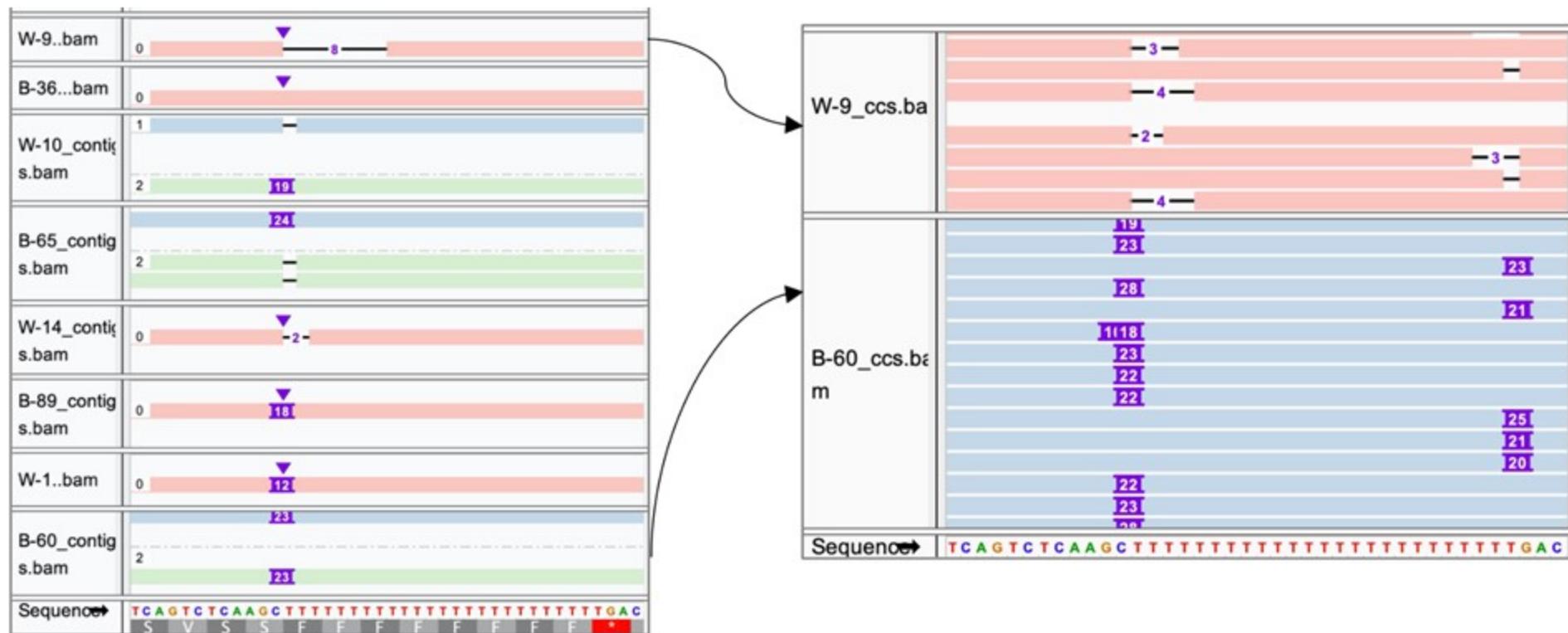
SVs affecting gene copy number were also resolved



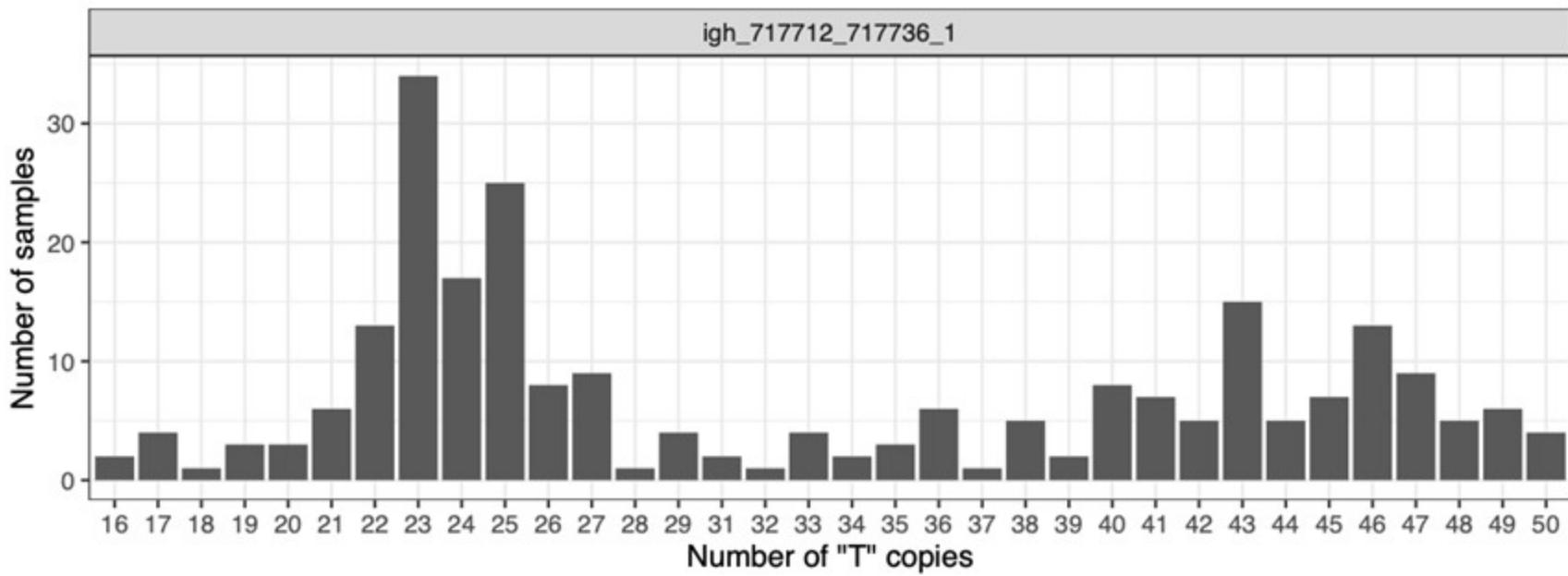
SVs affecting gene copy number were also resolved



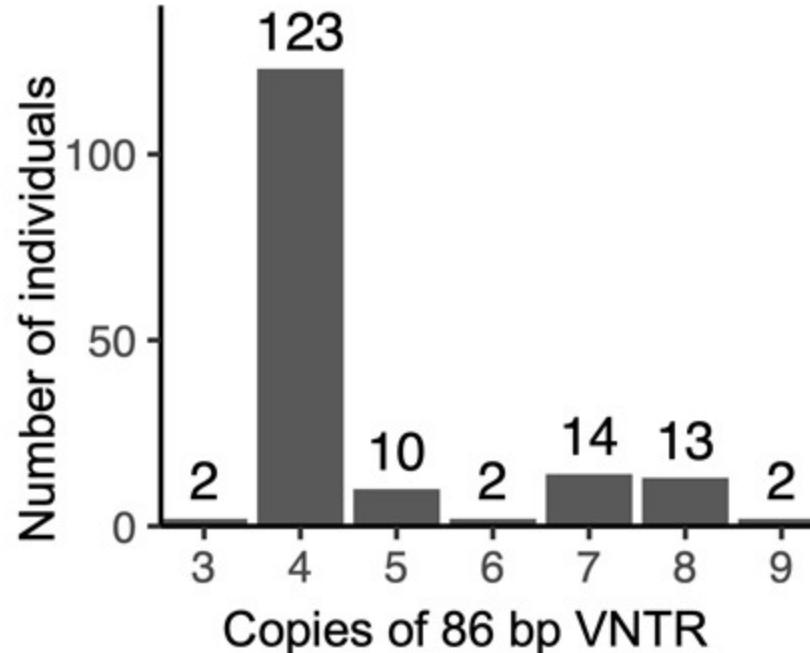
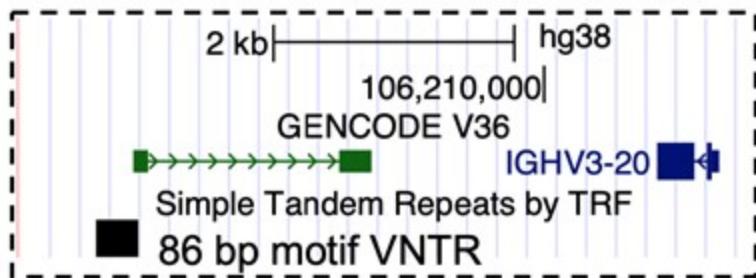
Large number of indels (966) and SVs (71) were also identified



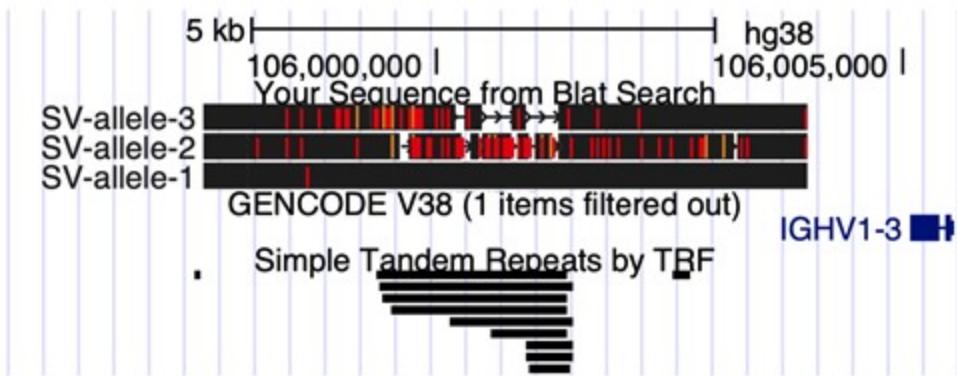
Large number of indels (966) and SVs (71) were also identified



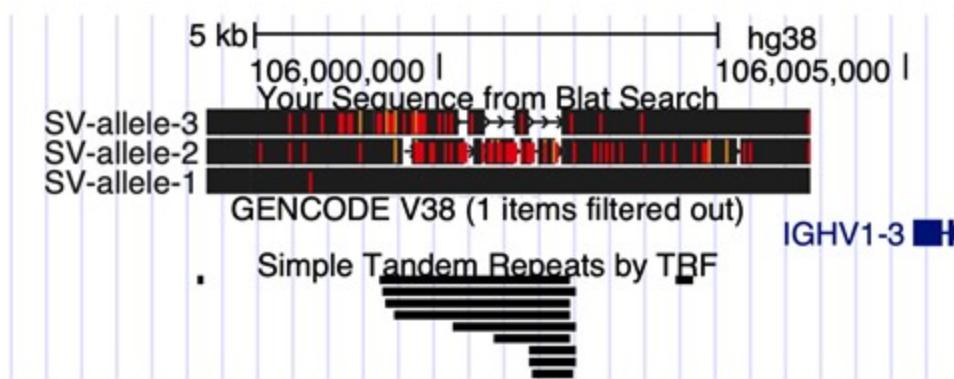
Large number of indels (966) and SVs (71) were also identified



Large number of indels (966) and SVs (71) were also identified



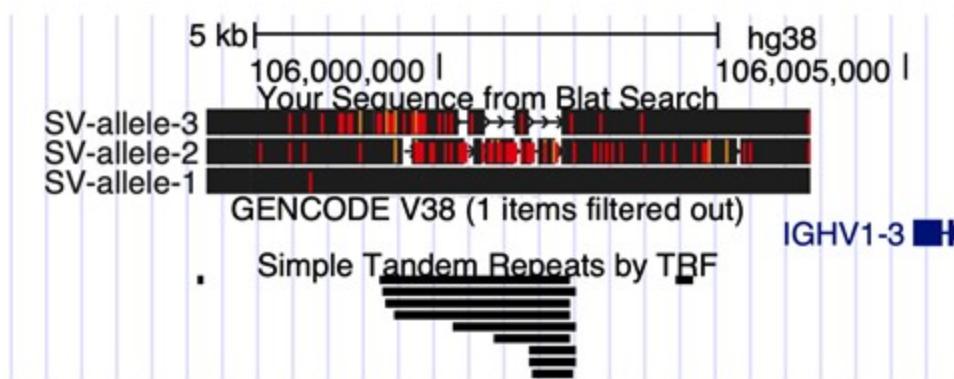
Large number of indels (966) and SVs (71) were also identified



Consensus motif sequence

- (3) **CTGGTGGTTC-TGAGCG-CCCC**
- (2) **-TGGTGGTTCCTGAGC**A**CCCCC**
- (1) **-TGGTG--TC**C**TGAGCG**C**CCCCC**

Large number of indels (966) and SVs (71) were also identified

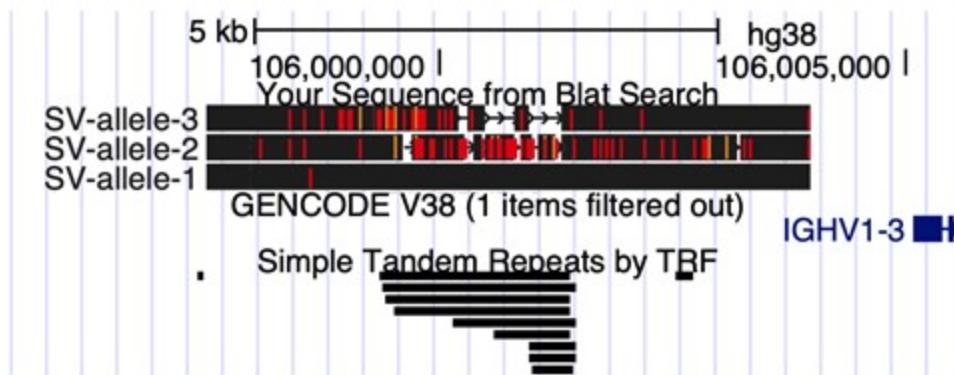


Consensus motif sequence

- (3) CTGGTGGTTC-TGAGCG-CCCC
- (2) -TGGTGGTTCCTGAGC**A**CCCCC
- (1) -TGGTG--TC**C**TGAGCG**C**CCCC

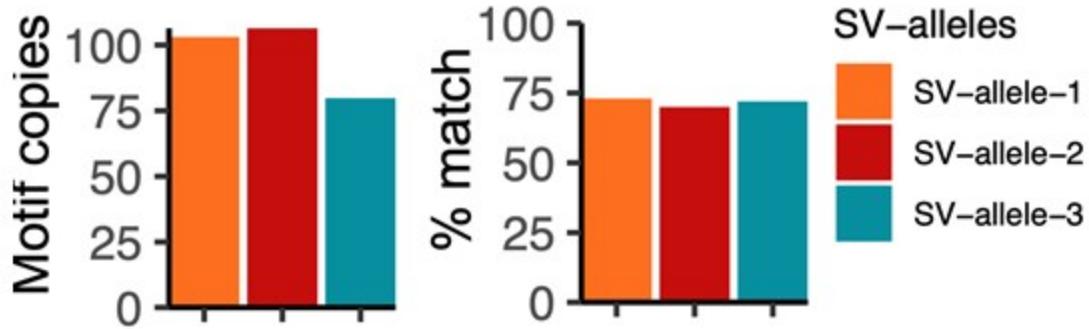


Large number of indels (966) and SVs (71) were also identified

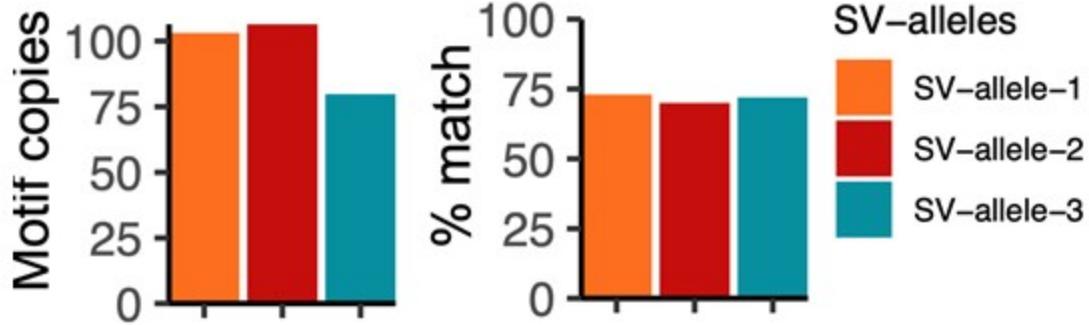


Consensus motif sequence

- (3) CTGGTGGTTC-TGAGCG-CCCC
- (2) -TGGTGGTTCCTGAGC~~A~~CCCCC
- (1) -TGGTG--TCCTGAGCGC~~C~~CCC

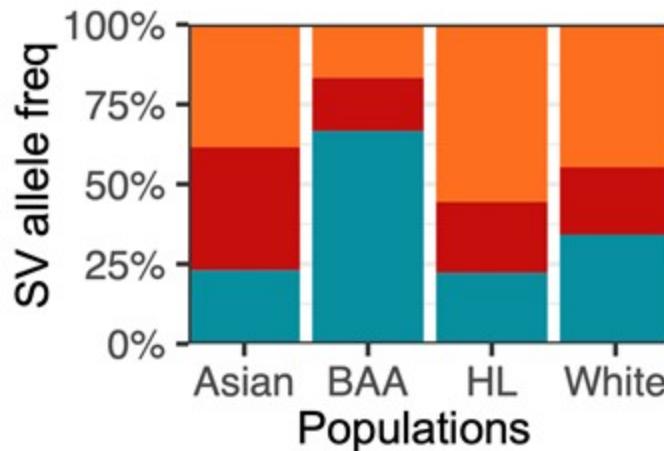


Large number of indels (966) and SVs (71) were also identified

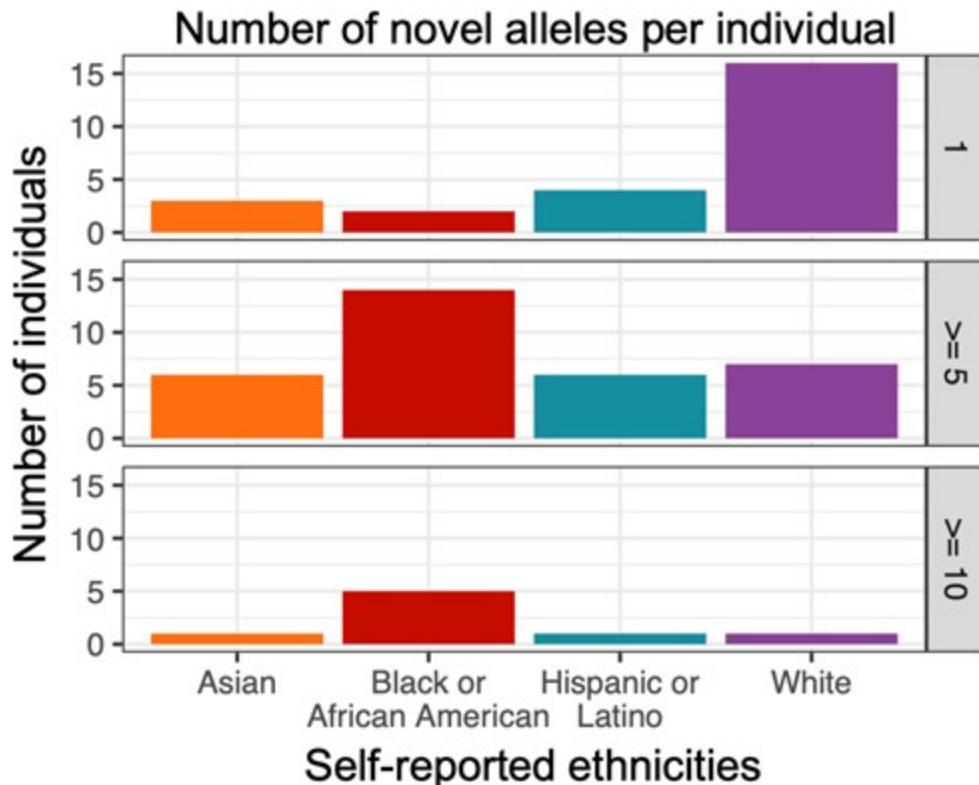
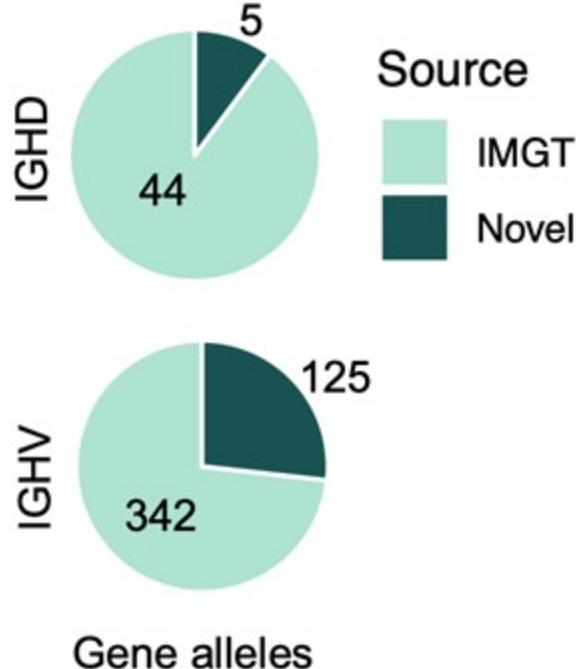


Consensus motif sequence

- (3) CTGGTGGTTC-TGAGCG-CCCC
- (2) -TGGTGGTTCCTGAGC~~A~~CCCCC
- (1) -TGGTG--TCCTGAGCGC~~C~~CCCCC



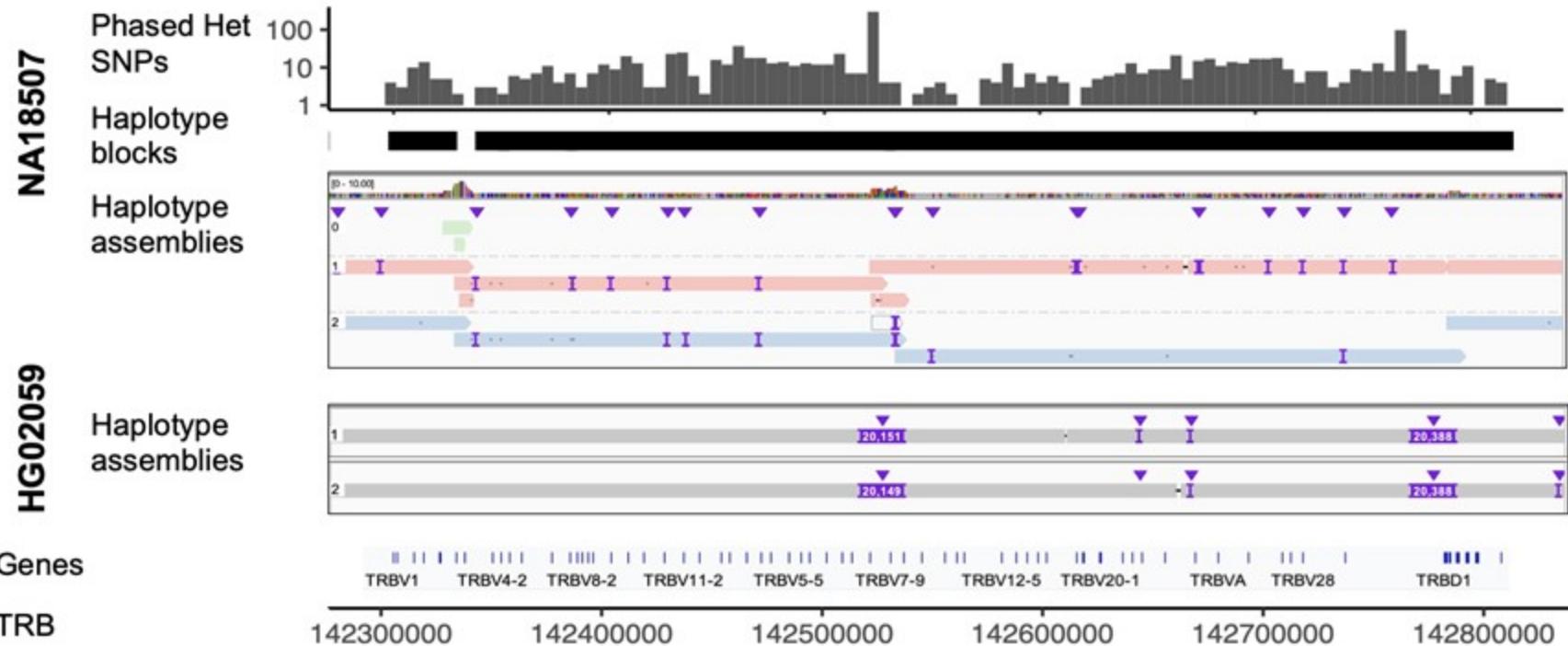
A large number of novel alleles were identified



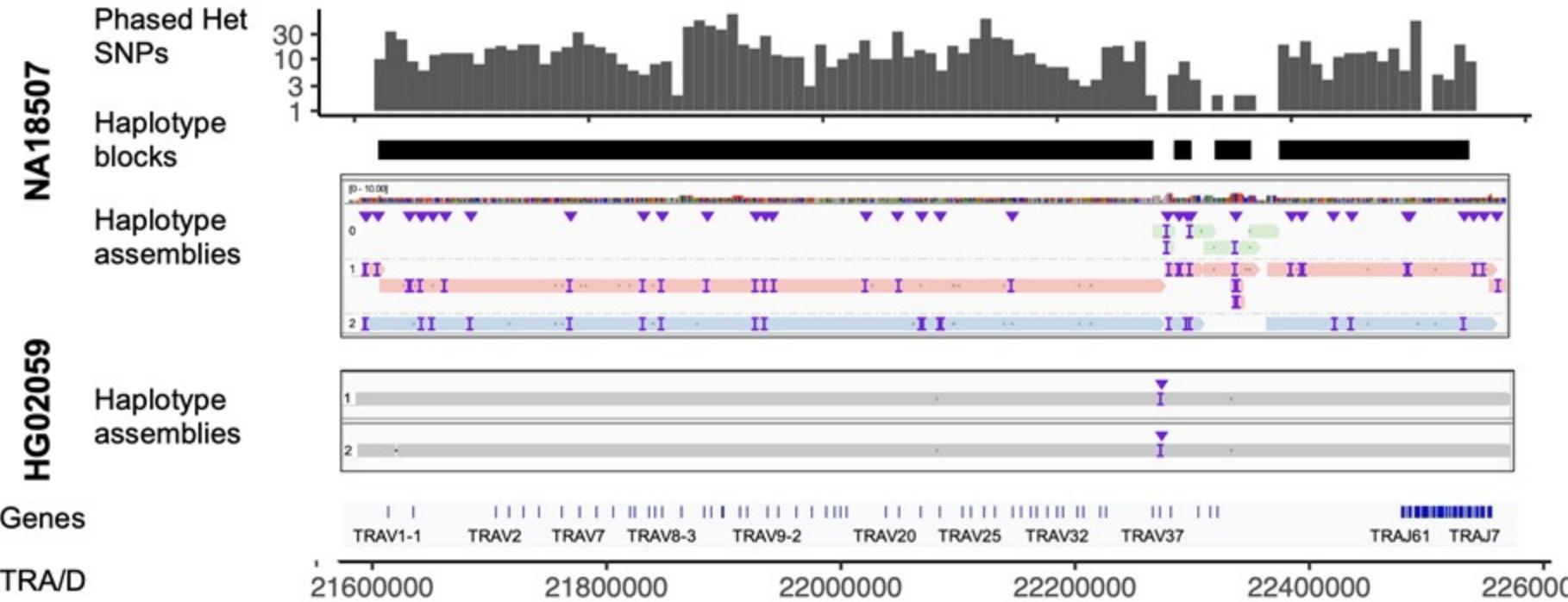
Outline of presentation

- 1) Introduction
 - a) Genetic variation in the immunoglobulin heavy chain locus (IGH)
 - b) Advantages of long-read sequencing over short-read sequencing
- 2) Framework for resolving the immunoglobulin heavy locus in a high throughput fashion using long read sequencing
- 3) Resolving the immunoglobulin heavy chain locus in a cohort with adaptive immune repertoire sequencing data
- 4) Resolving the T-cell receptor locus using long read sequencing
- 5) Application of framework to the immunoglobulin heavy chain locus in rhesus macaque
- 6) Conclusion

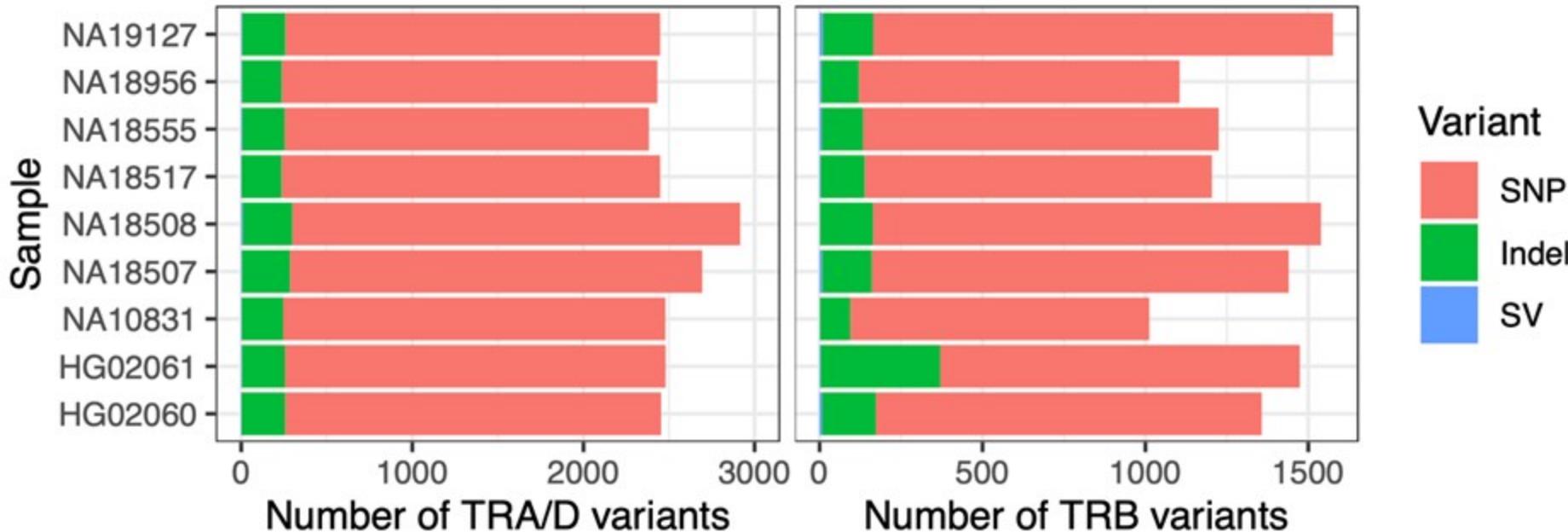
Resolving the TCR loci using long-read sequencing



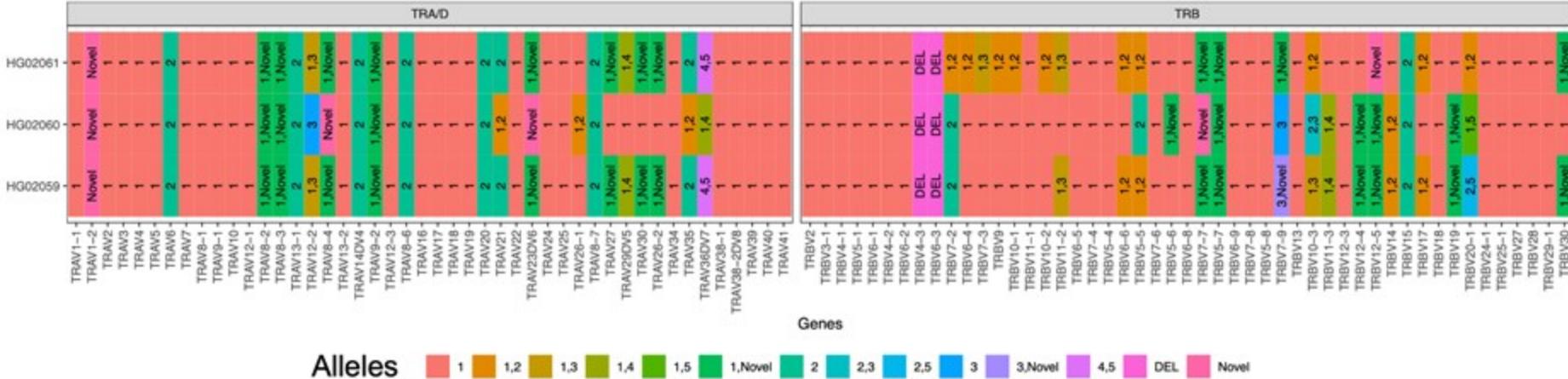
Resolving the TCR loci using long-read sequencing



SNPs, indels and SVs detected from the assemblies



Genes in TRA/D and TRB genotyped



Outline of presentation

- 1) Introduction
 - a) Genetic variation in the immunoglobulin heavy chain locus (IGH)
 - b) Advantages of long-read sequencing over short-read sequencing
- 2) Framework for resolving the immunoglobulin heavy locus in a high throughput fashion using long read sequencing
- 3) Resolving the immunoglobulin heavy chain locus in a cohort with adaptive immune repertoire sequencing data
- 4) Resolving the T-cell receptor locus using long read sequencing
- 5) Application of framework to the immunoglobulin heavy chain locus in rhesus macaque
- 6) Conclusion

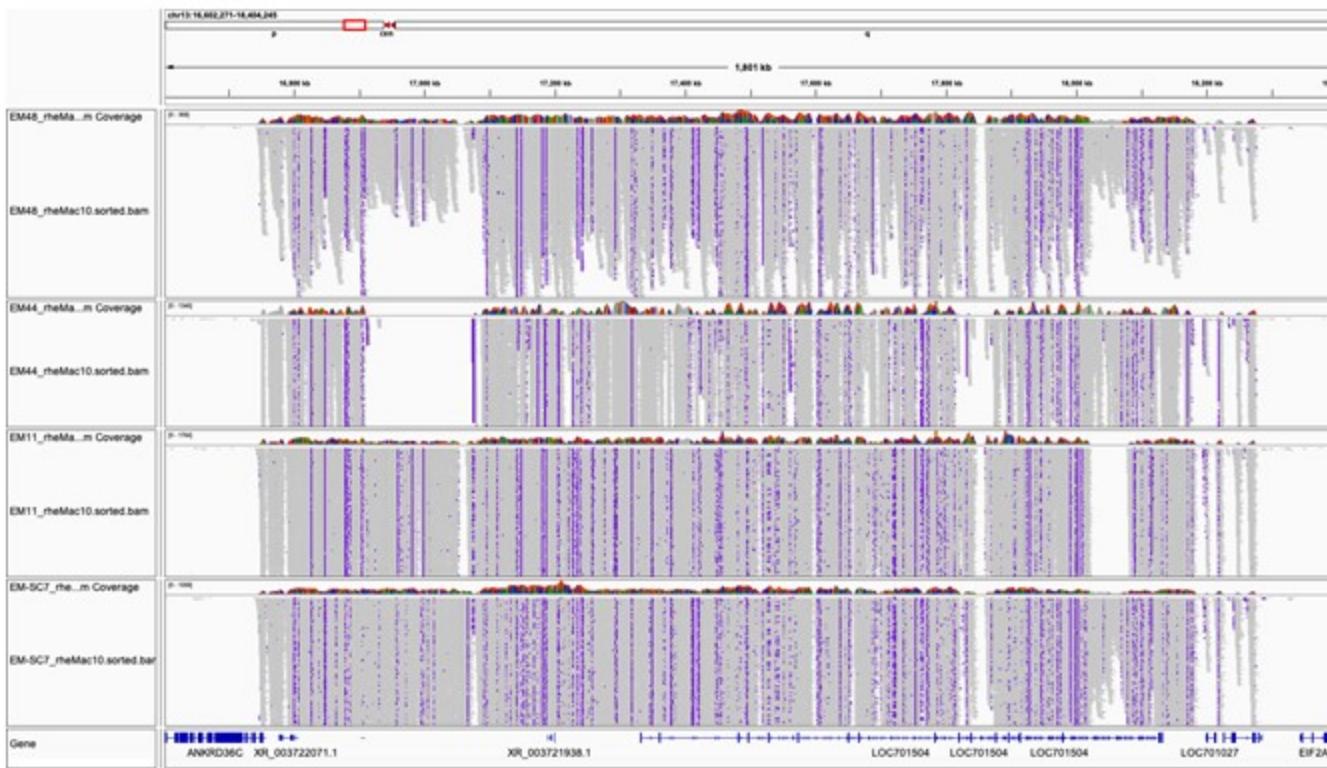
Added non-human functionality to IGenotyper



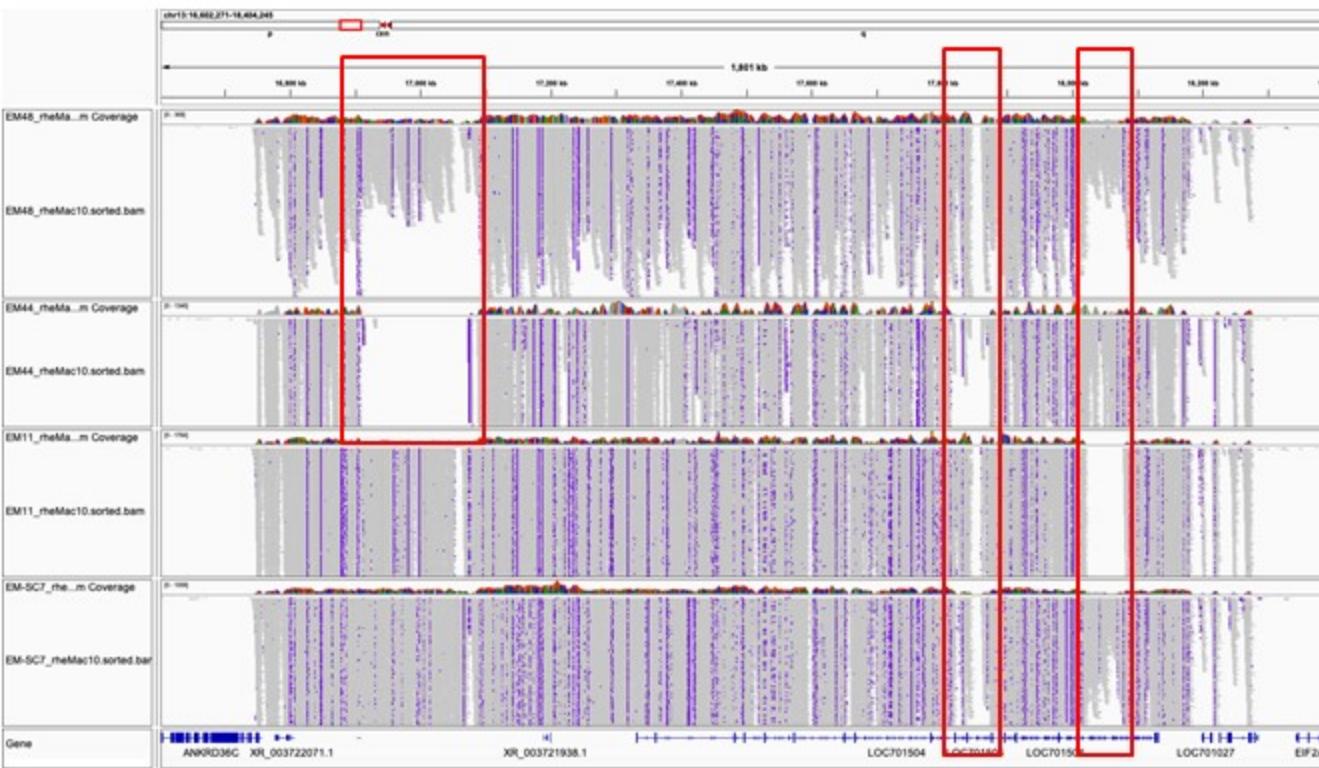
```
IG phase \
  --sample EM48 \
  --rhesus \
demultiplex.bc1009--bc1009.bam \
EM48_igenotyper
```

Collaboration:
Melissa Smith, PhD
Steve Bosingher, PhD

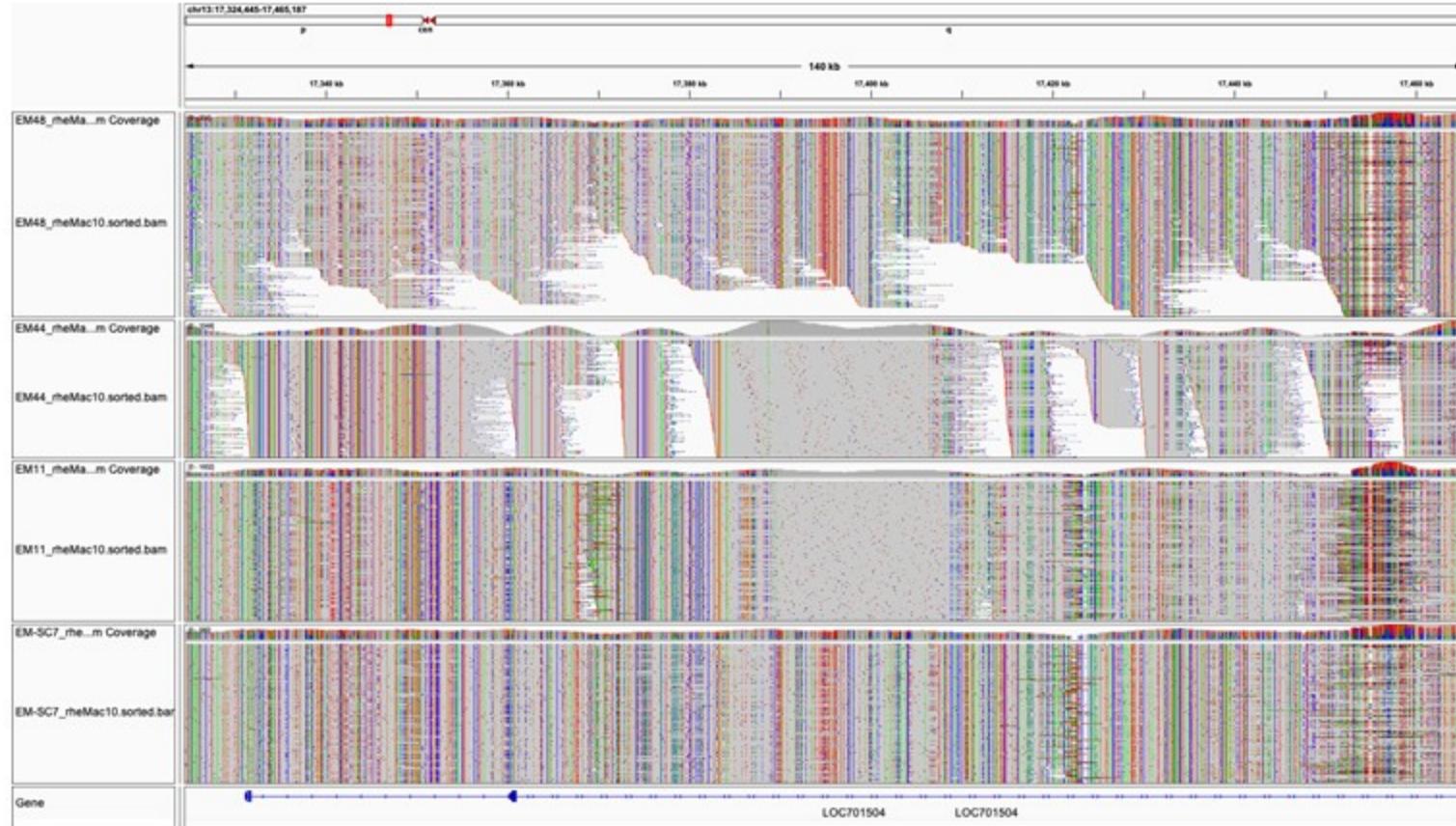
Preliminary targeted long-read sequencing coverage on 4 out 34 rhesus macaque individuals



Large structural variants detected in IGK from coverage profile



Large amount SNVs present



Conclusions

- Using a targeted approach and long read sequencing allows the resolution of the IG and TCR loci in humans and extendable to rhesus macaque
- Large amount of novel genetic variation in IGH was identified
 - ◆ Initial application to two samples identified 2 novel SVs and 16 novel alleles
 - ◆ SNVs were more accurately detected compared to NGS data
- Multiplexing allows for the application in a large cohort
- Large cohort identified even more novel genetic variation and alleles

Acknowledgments

University of Louisville

Corey Watson, PhD

Allison Silver, BS

Kaitlyn Shields, BS

William Gibson, BS

Melissa Smith, PhD

Harvard University

Wayne A. Marasco, MD, PhD

La Jolla Institute for Immunology

Shane Crotty, PhD

Emory University

Steve Bosinger, PhD

Adam Ericksen, PhD

Amit Upadhyay, PhD

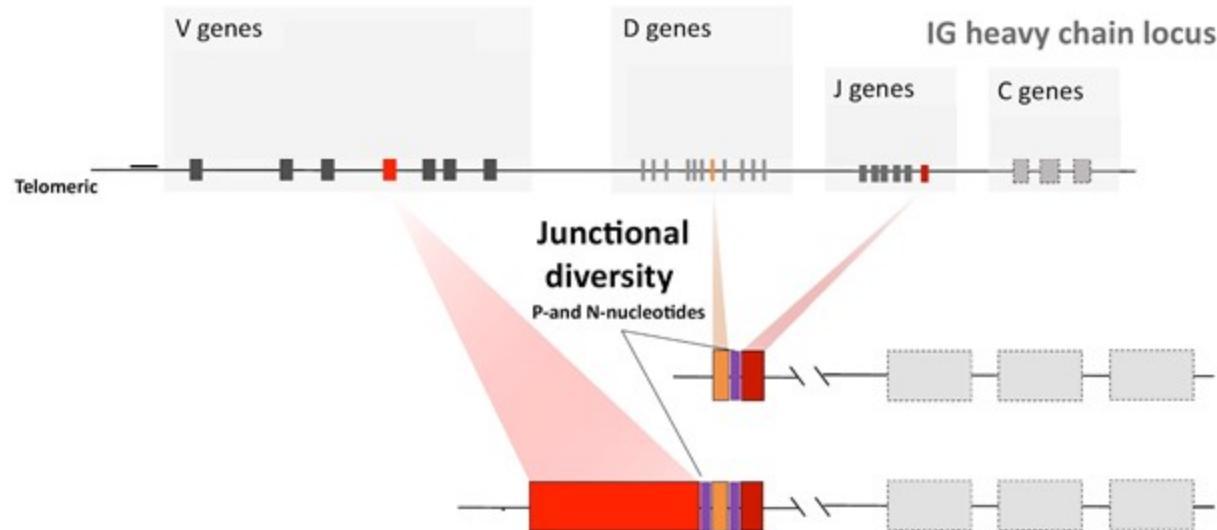


github.com/oscarlr/IGenotyper

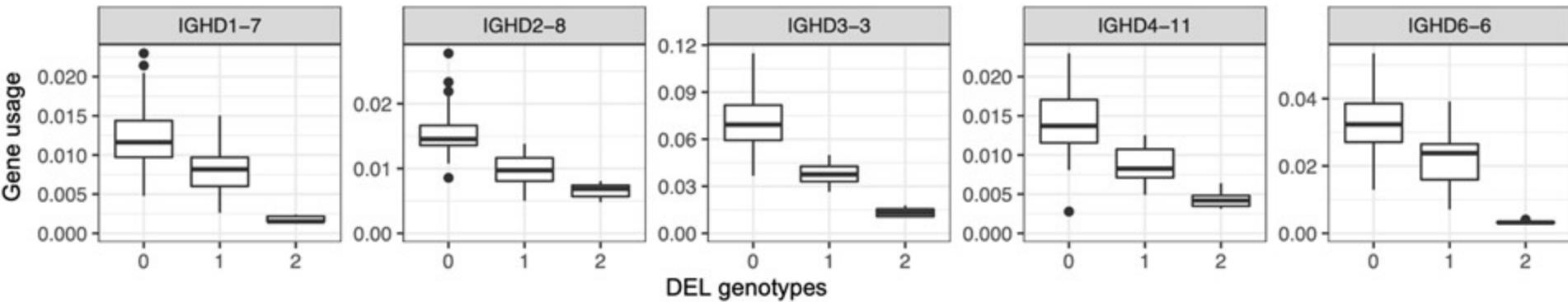
oscarlr.github.io

louisville.edu/research/watsonlab

DNA rearranges during the development of B cells



Using Ab repertoire data we correlated gene usage with genetic genotypes



Large number of structural variants in IG

Table 1. CNVs Identified from BAC and Fosmid Clones

Individual	Population	CNV Type	IGHV Genes Included in CNV ^a	GRCh37 Outer-Start (Breakpoint) ^b	GRCh37 Outer-End (Breakpoint) ^b	Event Size (-kbp)
CH17	nd	Insertion	V1-69D, V1-f, V3-h, V2-70D (gain)	107174927	107174941	46.6
CH17	nd	Complex event	V4-30-2 (gain) V4-31 (loss)	106804332	106810878	6.5 ^c /48.8 ^d
CH17	nd	Complex event	V5-a, V3-64D (gain) V3-9, V1-8 (loss)	106531320	106569343	38 ^c /37.7 ^e
CH17	nd	Insertion	V7-4-1 (gain)	106483362	106484225	9.5
NA12156	CEPH	Deletion	V4-39, V3-38 (loss)	106866357	106899042	32.7
NA15510 and NA19240	nd and Yoruba	Insertion	V1-c, V3-d, V3-43D, V4-b (gain)	106877146	106877535	61.1
NA18555 (haplotype A)	Han Chinese	Complex event	V3-30-5, V4-30-4, V3-30-3, V4-30-2 (gain)	106804332	106804333	49.2
NA18555 (haplotype B)	Han Chinese	Deletion	V4-31, V3-30 (loss)	106786254	106811213	24.9 ^c /73.9 ^d
NA18507	Yoruban	Complex event ^g	V4-30-4, V3-30-3 (gain) V3-30 (loss)	106784242	nd	25.2
NA18502	Yoruban	Complex event ^g	V3-30-5 (gain) V3-33, V4-31 (loss)	nd	106820685	24.7
NA18956 and NA12156	Japanese and CEPH	Duplication	V3-23D (gain)	106716650	106727861	10.8
NA19240 and NA12878	Yoruban and CEPH	Insertion	V7-4-1 (gain)	106483362	106484225	9.5

PICTURE

Testing capture plus IGenotyper on sample from Rheumatic heart disease GWAS study

ARTICLE

Received 15 Sep 2016 | Accepted 15 Feb 2017 | Published 11 May 2017

DOI: 10.1038/ncomms14946

OPEN

Association between a common immunoglobulin heavy chain allele and rheumatic heart disease risk in Oceania

Tom Parks¹, Mariana M. Mirabel², Joseph Kado^{3,4}, Kathryn Auckland¹, Jaroslaw Nowak⁵, Anna Rautanen¹, Alexander J. Mentzer¹, Eloi Marijon^{2,6}, Xavier Jouven^{2,6}, Mai Ling Perman⁴, Tuliana Cua⁷, John K. Kauwe⁸, John B. Allen⁸, Henry Taylor⁹, Kathryn J. Robson¹⁰, Charlotte M. Deane⁵, Andrew C. Steer^{11,12,*}, Adrian V.S. Hill^{1,*} & for the Pacific Islands Rheumatic Heart Disease Genetics Network[†]

GWAS SNVs were a combination of array, short-read data and imputed SNVs

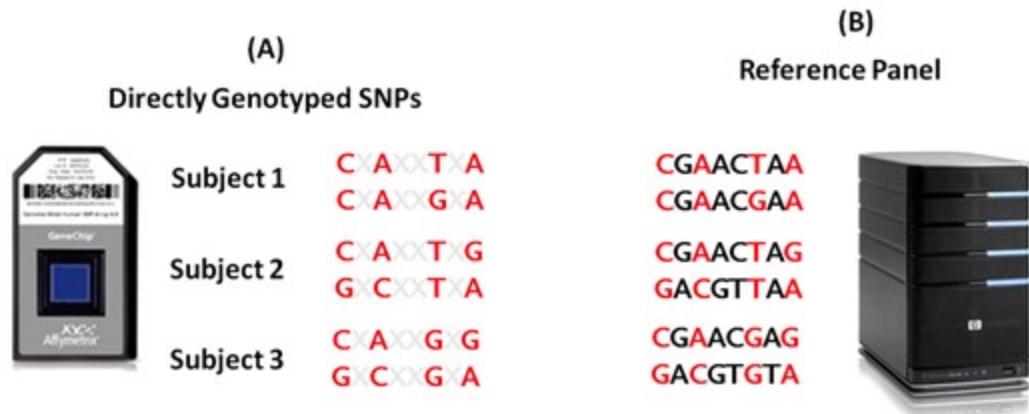
(A)

Directly Genotyped SNPs

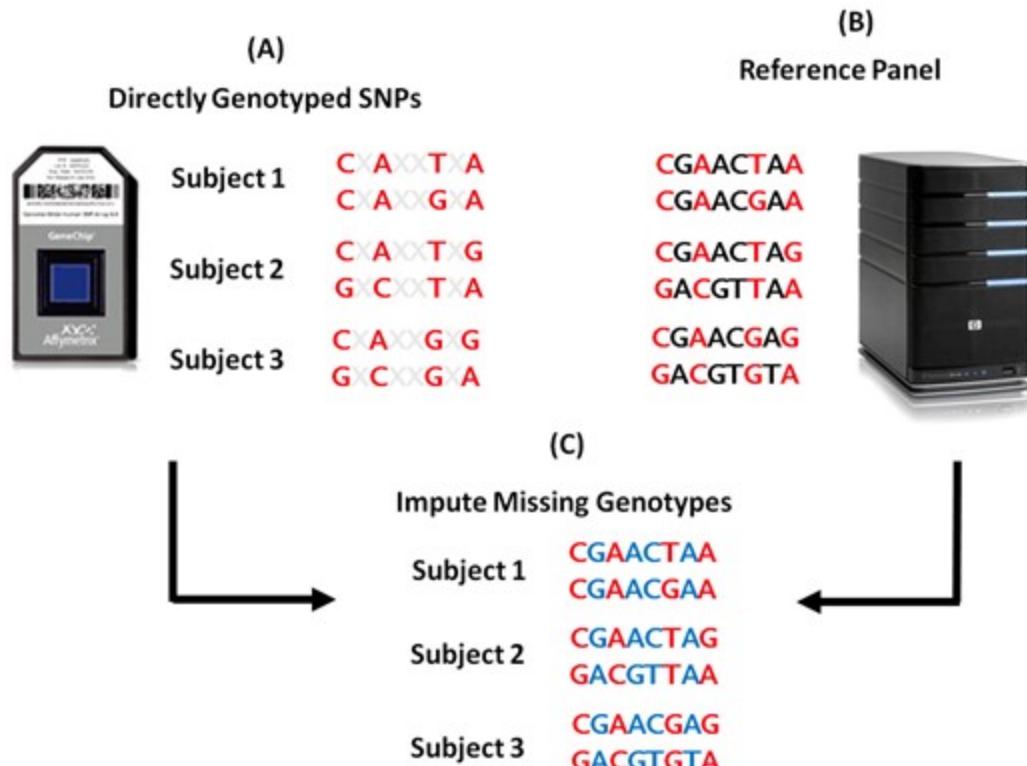


Subject 1	C A X X T A C A X X G A
Subject 2	C A X X T X G G C X X T X A
Subject 3	C A X X G G G C X X G A

GWAS SNVs were a combination of array, short-read data and imputed SNVs

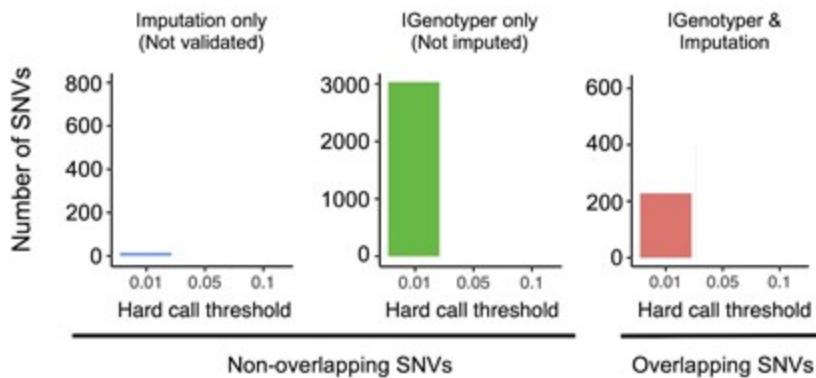


GWAS SNVs were a combination of array, short-read data and imputed SNVs



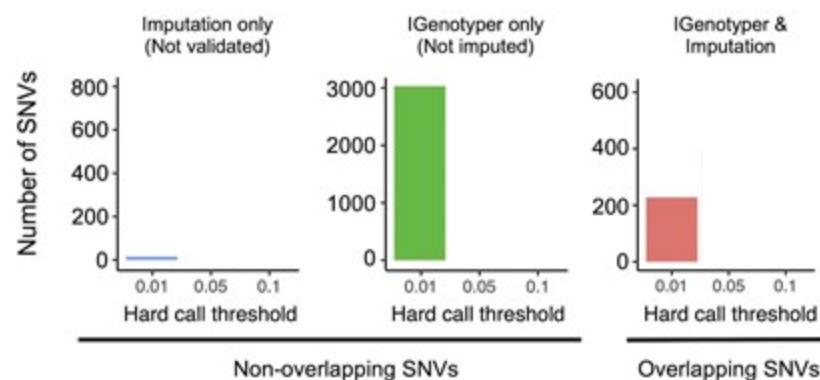
Reducing imputation threshold to allow more imputed SNVs increased number of FP SNV

a

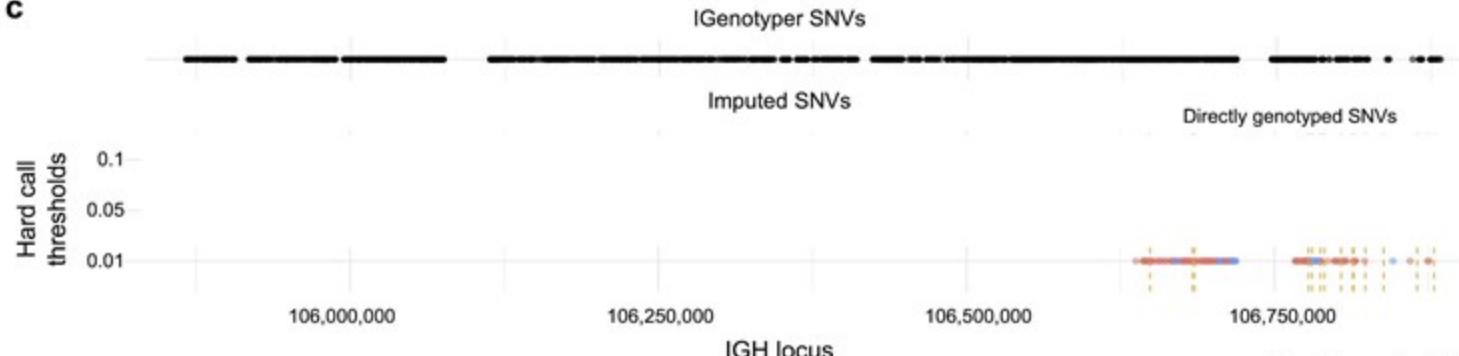


Reducing imputation threshold to allow more imputed SNVs increased number of FP SNV

a

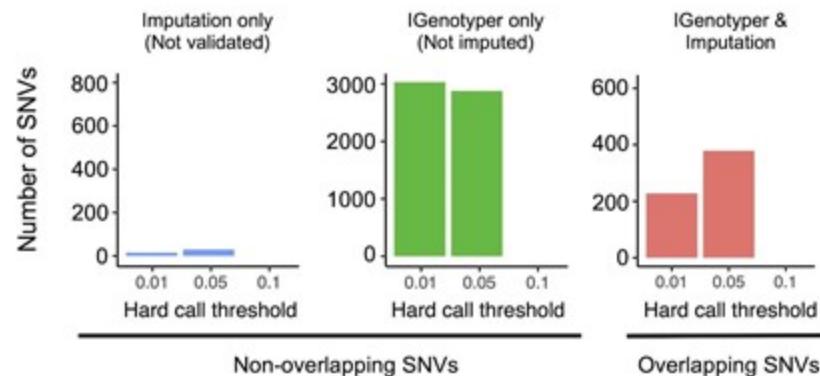


c

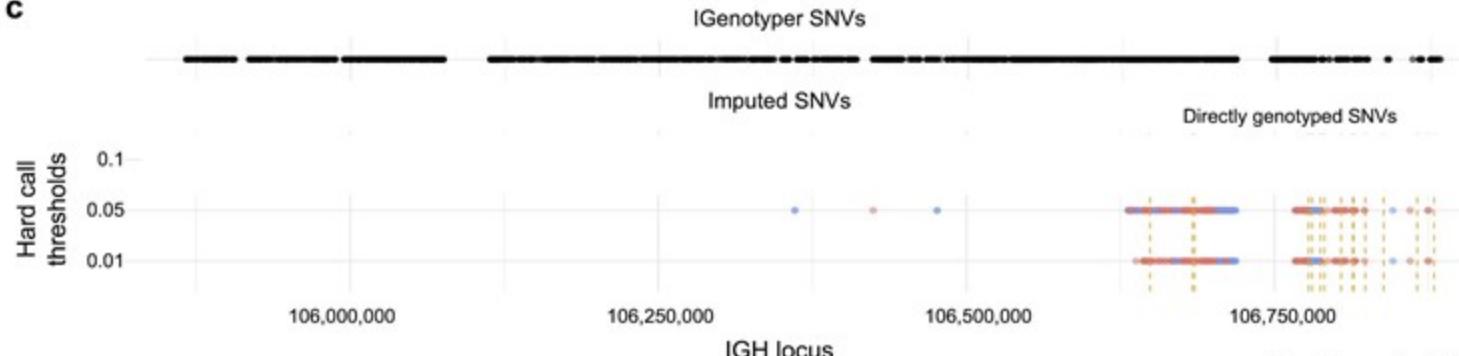


Reducing imputation threshold to allow more imputed SNVs increased number of FP SNV

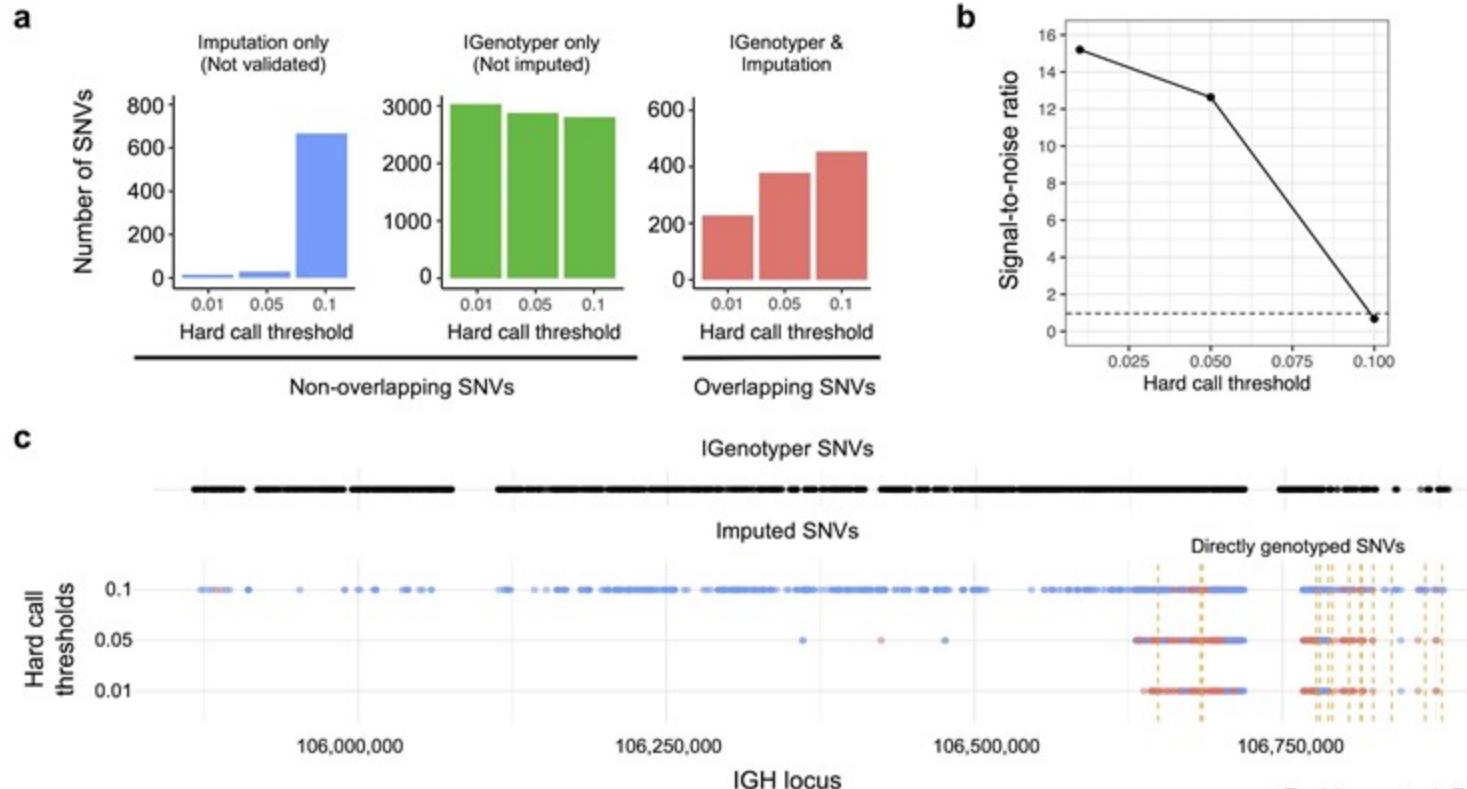
a



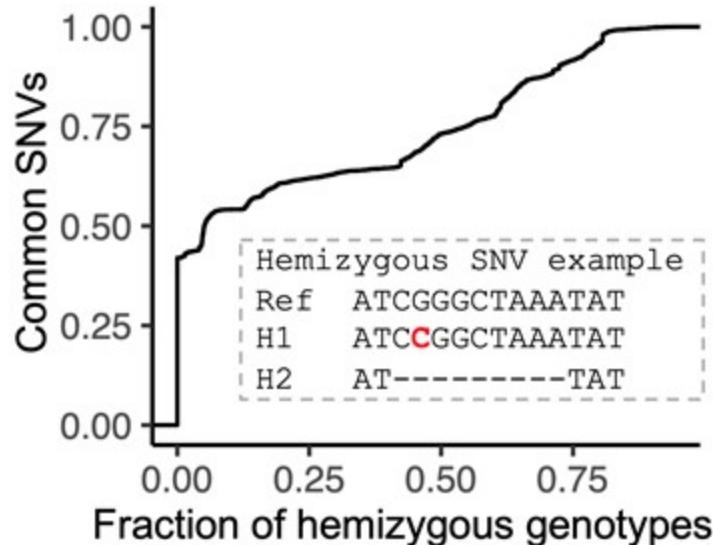
c



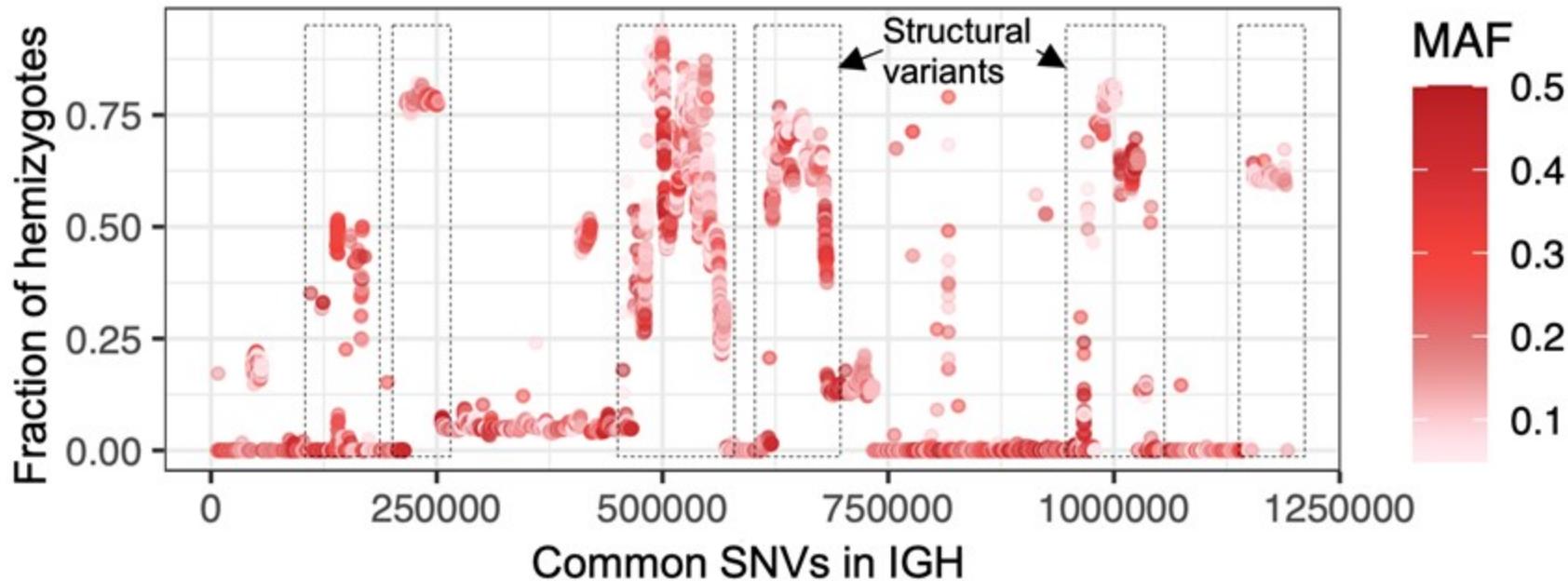
Reducing imputation threshold to allow more imputed SNVs increased number of FP SNV



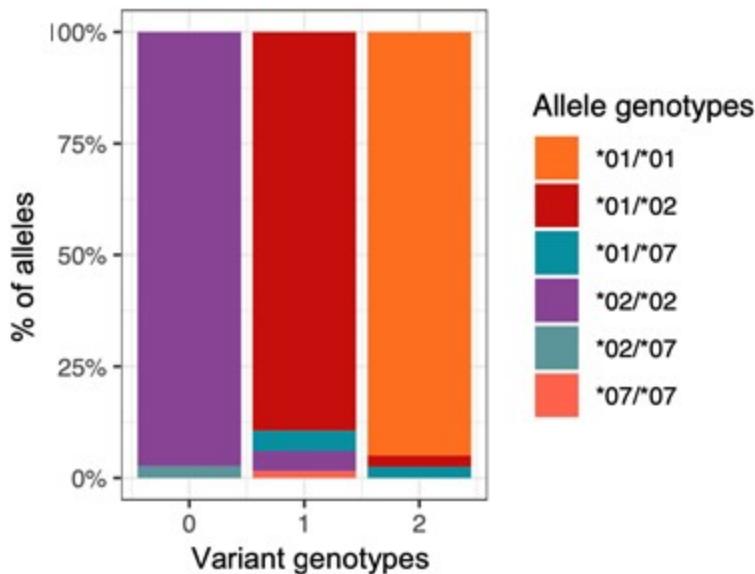
Large amount of hemizygous genotypes



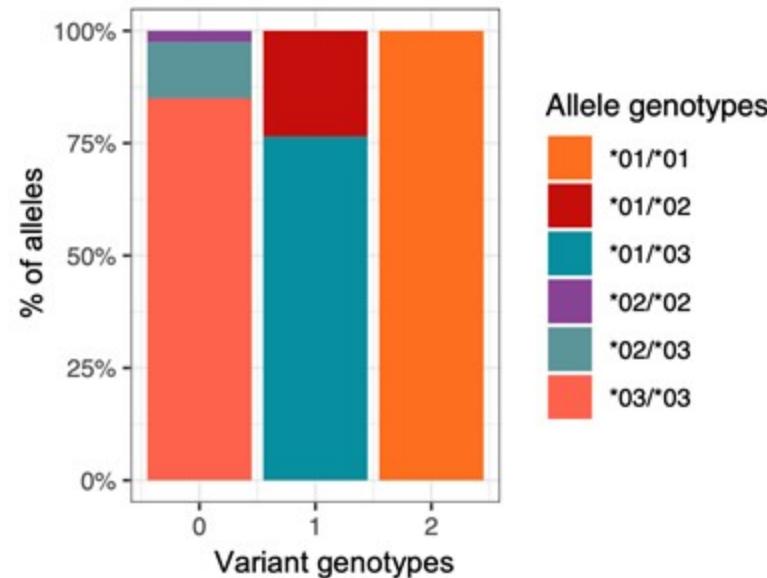
Large amount of hemizygous genotypes



Variants associated with gene usage contained different alleles across genotype groups

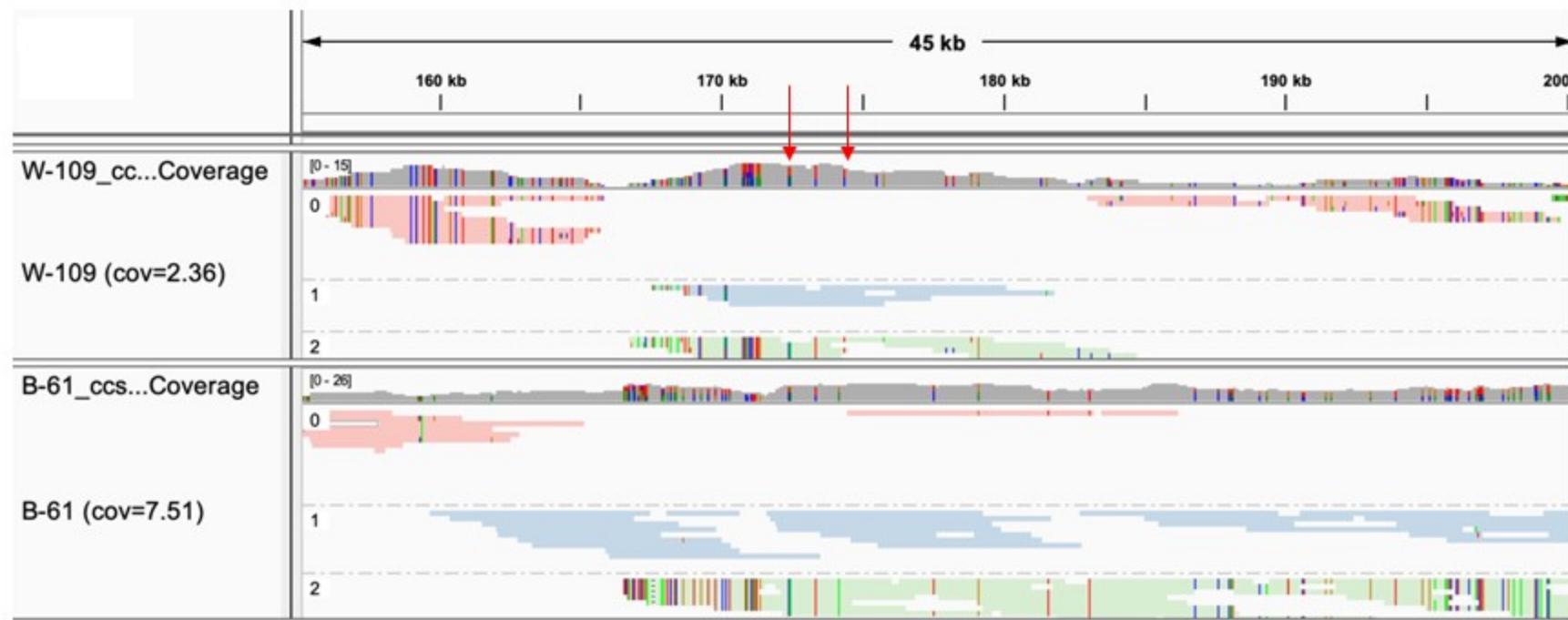


IGHV3-64

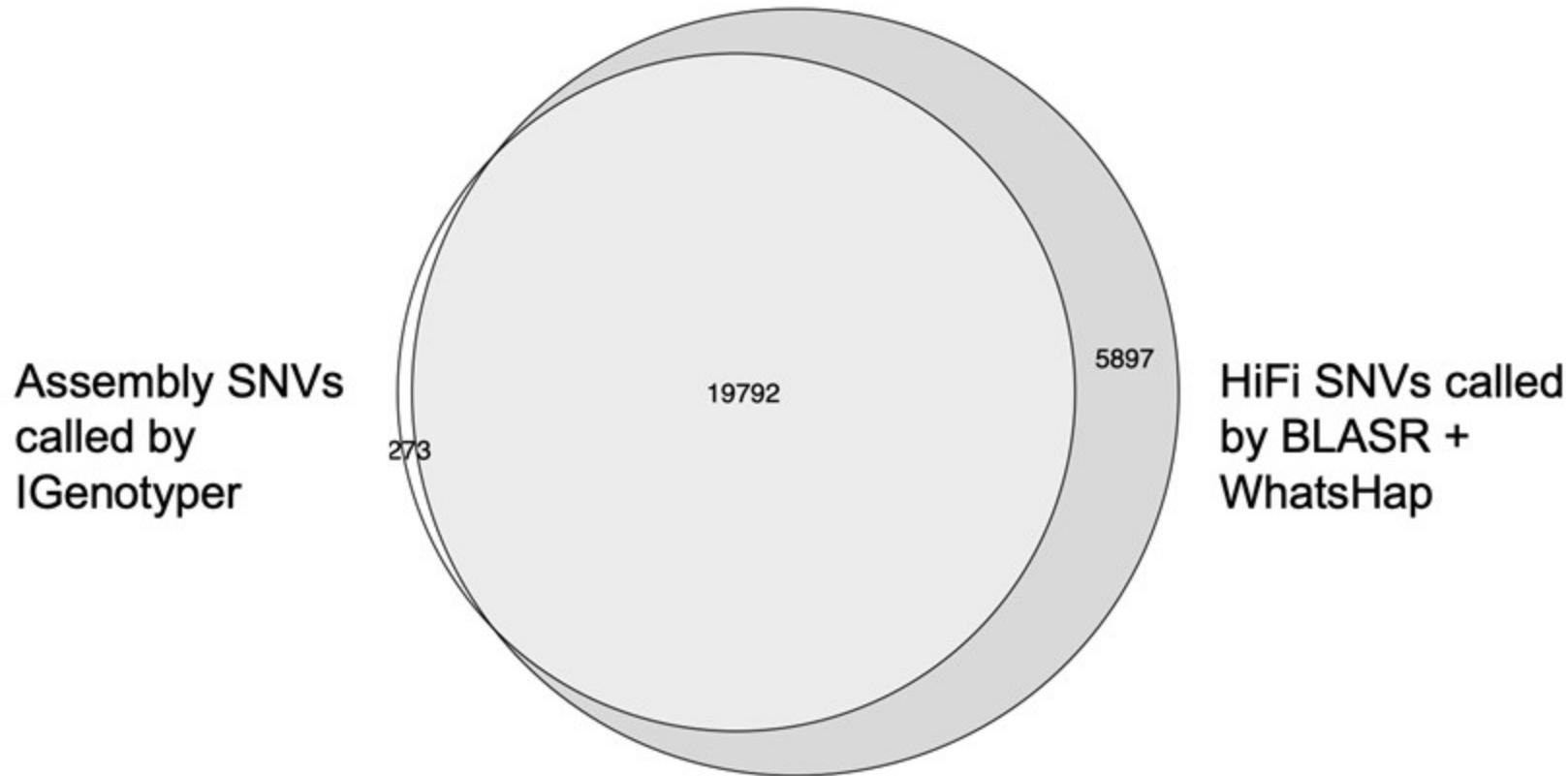


IGHV3-66

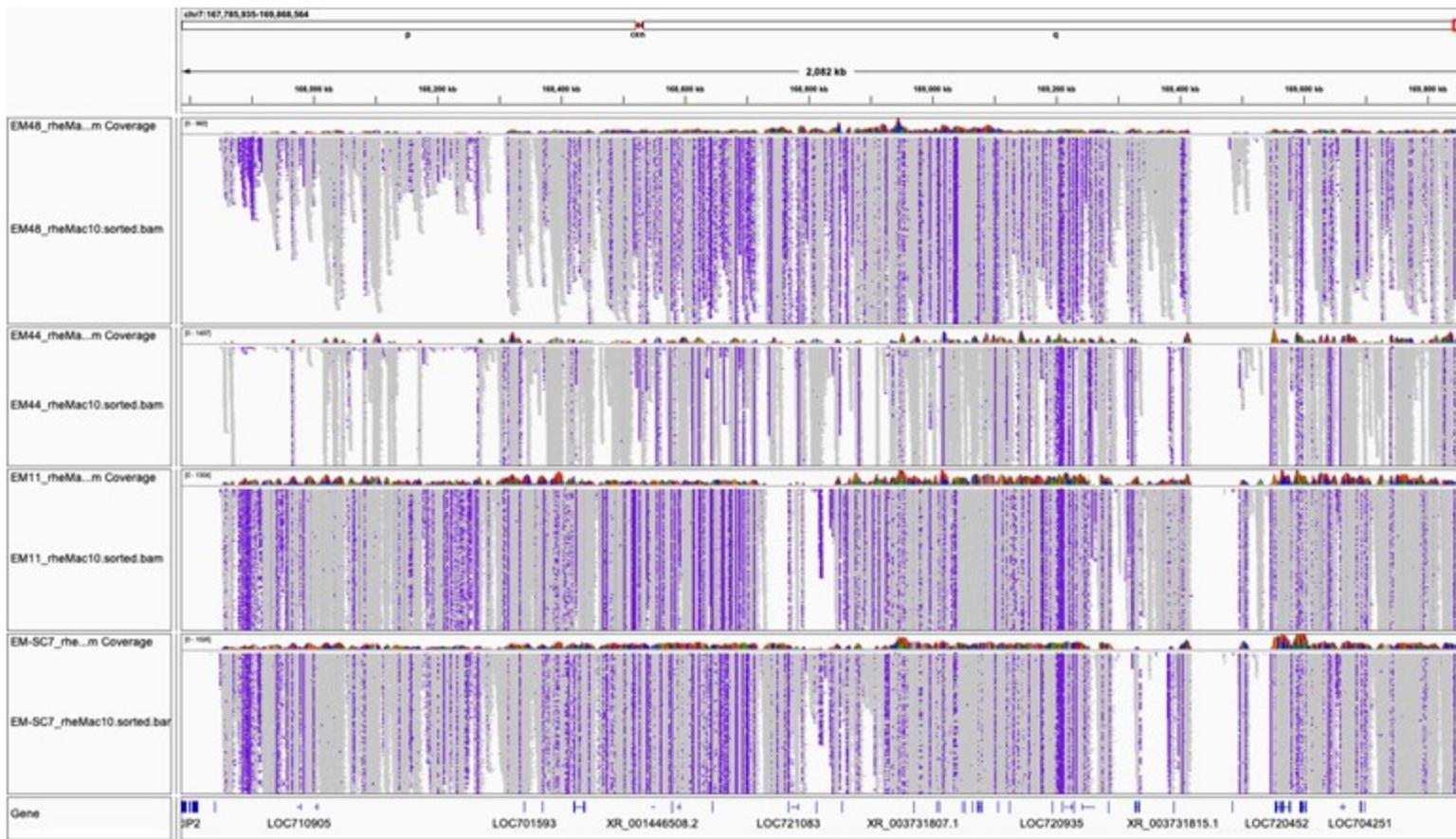
SNVs genotyped in samples with lower locus-wide coverage



Almost all SNVs were supported by CCS reads



IGH locus-wide coverage of long-read data



SVs present in IGK

