

# Ultrasensitive allele inference from immune repertoire sequencing data with MiXCR

Artem Mikelov<sup>5†\*</sup>, George Nefediev<sup>1†</sup>, Alexander Tashkeev<sup>2</sup>, Oscar L. Rodriguez<sup>3</sup>, Diego A. Ortman<sup>2</sup>, Valeriia Skatova<sup>1</sup>, Mark Izraelson<sup>1</sup>, Alexey N. Davydov<sup>1,4</sup>, Stanislav Poslavsky<sup>1</sup>, Souad Rahmouni<sup>2</sup>, Corey T. Watson<sup>3</sup>, Dmitriy Chudakov<sup>1,4</sup>, Scott D. Boyd<sup>5</sup>, Dmitry Bolotin<sup>1</sup>

† - These authors contributed equally to this work.

\*Corresponding authors: [amikelov@stanford.edu](mailto:amikelov@stanford.edu), [bolotin@milaboratories.com](mailto:bolotin@milaboratories.com)

1 - MiLaboratories Inc, San Francisco, CA, USA

2 - Unit of Animal Genomics, WELBIO, GIGA-R & Faculty of Veterinary Medicine, University of Liège (B34), Liège, Belgium

3 - Department of Biochemistry and Molecular Genetics, University of Louisville School of Medicine, Louisville, KY, USA

4 - Central European Institute of Technology, Masaryk University, Brno, Czech Republic

5 - Department of Pathology, Stanford University, Stanford, CA, USA

## Abstract

Allelic variability in the adaptive immune receptor loci, which harbor the gene segments that encode B cell and T cell receptors (BCR/TCR), is of critical importance for immune responses to pathogens and vaccines. Adaptive immune receptor repertoire sequencing (AIRR-seq) has become widespread in immunology research making it the most readily available source of information about allelic diversity in immunoglobulin (IG) and T cell receptor (TR) loci. Here we present a

novel algorithm for extra-sensitive and specific variable (V) and joining (J) gene allele inference, allowing reconstruction of individual high-quality gene segment libraries. The approach can be applied for inferring allelic variants from peripheral blood lymphocyte BCR and TCR repertoire sequencing data, including hypermutated isotype-switched BCR sequences, thus allowing high-throughput novel allele discovery from a wide variety of existing datasets. The developed algorithm is a part of the MiXCR software.

We demonstrate the accuracy of this approach using AIRR-seq paired with long-read genomic sequencing data, comparing it to a widely used algorithm, TIgGER. We applied the algorithm to a large set of IG heavy chain (*IGH*) AIRR-seq data from 450 donors of ancestrally diverse population groups, and to the largest reported full-length TCR alpha and beta chain (*TRA*; *TRB*) AIRR-seq dataset, representing 134 individuals. This allowed us to assess the genetic diversity within the *IGH*, *TRA* and *TRB* loci in different populations and to establish a database of alleles of V and J genes inferred from AIRR-seq data and their population frequencies with free public access through [vdj.online](https://vdj.online) database.

## Introduction

Adaptive immune repertoire diversity plays a crucial role in shaping the immune response and forming immunological memory. Most immune repertoire research has focused primarily on somatically derived immune receptor diversity, namely V(D)J recombination and somatic hypermutation (SHM) diversity. In recent years, however, the extent of population diversity has begun to be appreciated at both the immunoglobulin (IG) (Gidoni et al., 2019; Mikocziova et al. 2021; Corcoran et al.

2023; Rodriguez et al. 2023; Gibson et al. 2023) and T cell receptor (TCR) loci (Omer et al., 2022; M. Corcoran et al., 2023; Rodriguez et al. 2022). The functional significance of allelic variation in adaptive immune loci has also been recognized in the context of influenza, HIV and COVID-19 immunity and vaccination (Avnir et al. 2016; Lee et al. 2021; Leggat et al. 2022; Pushparaj et al., 2022).

Sequencing repertoires of adaptive immune receptors encoded by recombined germline V, D and J genes have become a major source of information about adaptive immune functions in health and disease. In recent years, AIRR-seq has been utilized to discover many novel alleles in TR and IG loci, becoming one of the major sources of information of the allelic diversity of TR and IG genes in different populations. However, the major obstacle for utilizing AIRR-seq datasets for genotyping and allelic discovery is the presence of somatically hypermutated sequences in most available immunoglobulin AIRR-seq datasets, along with the PCR and sequencing errors which affect both BCR and TR repertoire datasets. Hot-spot hypermutations and sequence errors have significantly hindered the ability to clearly detect individual polymorphisms. We aimed to overcome these issues with the algorithm described in this paper. However, there are two other challenging obstacles for accurate genotyping and haplotyping of TR and IG loci using AIRR-seq data only. Common structural variants (SVs) in IG loci (Rodriguez et al. 2023), especially gene duplications, in some cases make it hard to unequivocally map a sequence from AIRR-seq data to a particular germline gene without an additional source of information. Further, some alleles exhibit low usage levels, precluding their detection with AIRR-seq. Despite these limitations, AIRR-seq data remain valuable for applications focused on functional adaptive immune repertoires and their fluctuations in different conditions.

The ability to precisely call known allelic variants and infer novel ones from the same AIRR-seq data could enable new analyses of germline variation contribution to immune responses, and also improve the accuracy of many existing downstream approaches. There are several published methods for genotyping and allelic inference of V and J genes from AIRR-seq data (Table 1, numbers 2-5), however, each has important limitations. TIgGER (Gadala-Maria, Yaari, Uduman, & Kleinstein, 2015; Gadala-Maria et al., 2019) and Partis (Ralph & Matsen, 2019) are based on the idea that allelic sequence variants show a distinctive pattern over the background of SHMs. On the other hand, IgDiscover (M. M. Corcoran et al., 2016), a very robust and reliable tool for novel allele inference, requires data without hypermutations, thus excluding much published immunoglobulin repertoire data. The recently developed PIgLET software (Peres et al. 2023) has enhanced genotype inference capabilities through the use of *IGHV* allele similarity clustering, although it is not designed to infer novel alleles.

**Table 1. Tools for novel allele variants inference from AIRR-seq data and their characteristics.**

#	Tool name	Year	Supported chain type(s)	Supported gene type(s)	Programming language(s)	Suitable for inference from hypermutated repertoires
1	MiXCR	2023	<i>IGH, IGK, IGL, TRA, TRB</i>	V, J	Java, Kotlin	Yes
2	TlgGER	2015	<i>IGH, IGK, IGL</i>	V	R	Yes
3	IgDiscover	2016	<i>IGH, IGK, IGL, TRA, TRB</i>	V, D, J	Python	No
4	Partis	2019	<i>IGH, IGK, IGL</i>	V	C, C++, Perl, Python	Yes
5	ImPre	2016	<i>IGH, IGK, IGL, TRA, TRB</i>	V, J	C, Perl	Yes

Existing tools also require considerable depth of AIRR-seq data for reliable allele inference (for example, IgDiscover recommends at least 750,000 sequencing reads per individual library). Such sequencing depth is costly and not available for most publicly available AIRR-seq datasets. Here we present an algorithm for allelic inference and genotyping from both hypermutated and non-hypermutated repertoires, with low sequencing depth requirements. The algorithm performs well starting with a minimalistic gene reference library of only one allele for each gene, and even with some genes missing. These features make the tool especially useful for studying allelic diversity in non-model species where reference gene libraries are sparse and incomplete. The developed approach is integrated in MiXCR software and is available as the `findAlleles` command. Starting with version 4.0, MiXCR can process immune repertoire data directly from raw sequencing reads in FASTQ format. The MiXCR upstream pipeline supports all commercially available library preparation kits, as well as any custom protocols. It handles pre-processing, sequence alignment, and clonal

assembly based on customizable region of interest, such as the CDR3, the whole VDJ region, or any user-defined region. The output of the upstream pipeline can be generated in several formats: a highly efficient binary format, a tabular format with customizable fields, or the AIRR format. The `findAlleles` command can be executed on clonesets in binary format. For each individual `findAlleles` outputs a personalized reference allele library in either FASTA, tabular or json format.

The International ImMunoGeneTics Information System (IMGT®), established in 1989, is the oldest widely available source of information about immune receptors, including alleles. Recent advancements in high-throughput adaptive immune receptor repertoire sequencing (AIRR-seq) methods enabled a broader view of alleles, and many tools were developed to infer allelic variants from such data. In 2017, the AIRR-Community (a network of over 300 practitioners in the field of AIRR-seq, [www.airr-community.org](http://www.airr-community.org)) and IMGT® agreed on a process for adding new alleles inferred from AIRR-seq data to the IMGT® database (Ohlin et al., 2019). The AIRR Community also introduced the Open Germline Receptor Database (OGRDB, <https://ogrdb.airr-community.org/>, Lees et al., 2020), to track the addition of new alleles. Although being the most recognized source of germline immunoglobulin sequence data, IMGT® lacks information on population allele frequencies and harbors sequences ‘mapped’ to the identified genes at the specific genomic locations. However, structural variation is quite common (Rodriguez et al. 2023) in the IG and TR loci, while most new IG and TR sequence data is coming from AIRR-seq experiments, and can be hard to map to a particular germline locus position. Other databases of immune receptor gene alleles have been introduced,

such as pmTR ([Dekker, van Dongen, Reinders, & Khatri, 2022, https://pmtrig.lumc.nl/](https://pmtrig.lumc.nl/)), IgPdb (<https://cgi.cse.unsw.edu.au/~ihmmune/IgPdb>) and Karolinska Institutet human T cell receptor database (M. Corcoran et al., 2023, <https://gkhlab.gitlab.io/tcr/>). A comprehensive and well-maintained database of immune receptor gene alleles, including allelic variants inferred from AIRR-seq, is VDJbase (Omer et al. 2020, <https://vdjbase.org/>). However, VDJbase is not seamlessly integrated with any of the analysis tools, and using it for AIRR-seq data analysis requires conversion of sequence data formats. To accompany the MiXCR software, we have developed VDJ.online (<https://vdj.online/library>), a free and open database of immune receptor allelic sequences that enables examination, comparison, and downloading of sequences. The VDJ.online reference library is supplied with the MiXCR, allowing seamless AIRR-seq data processing with accurate V and J gene annotation, genotyping, and novel allele inference.

## Results

### **Novel approach to V and J gene allele variants inference and genotyping**

The main challenge of allelic inference from AIRR-seq data is the presence of hypermutations, PCR and sequencing errors, with a large fraction of them being hot-spot mutations occurring simultaneously in unrelated clones. We have overcome this challenge by consecutively applying several filters based on two major measures. The first one is the lower diversity bound, estimated as the number of unique combinations of J and V genes and CDR3-lengths of clonotypes. The second measure is based on the number of clonotypes with unmutated J and V genes. Filters are applied both at individual mutation and at mutation set levels (see

Methods for the detailed description). The mutations at germline-encoded positions in CDR3 are recovered, when possible, by using non-mutated clonotypes matching the inferred variants in the rest of the sequence. This approach allows both to infer novel (undocumented) V and J gene alleles and to perform genotyping with high sensitivity and precision.

### **Benchmarking of the V and J gene allele variants inference and genotyping**

To assess the performance of the developed algorithm we utilized publicly available datasets (Rodriguez et al. 2023) containing both AIRR-seq data and highly reliable genotyping data of the IGH locus reconstructed using Pacific Biosciences HiFi long-read sequencing from the same individuals. For the sake of comparison we utilized 33 AIRR-seq data sets of sufficient sequencing depth (> 500,000 sequencing reads) and at least 3,000 unique full-length clonotypes. Targeted long-read sequencing of genomic DNA (as described in Rodriguez et al., 2023) allowed us to observe the non-rearranged *IGH* locus, containing germline, unmutated *IGHV* and *IGHJ* genes. Since these genes were not rearranged, and did not contain somatic hypermutation, they provide a reliable ground truth for allele identification in our comparison. However, in some individuals, not all genes of interest were captured by the long-read sequencing. Consequently, we excluded the allele calls for these genes from our comparison. AIRR-seq data, which was derived from peripheral blood mononuclear cell samples containing both naive and antigen-experienced B cells, expressing either unmutated or somatically hypermutated BCR sequences, was used to infer the allele variants using our approach, with the PacBio germline DNA sequences as the gold standard for the true alleles present in each individual



(Rodriguez et al., 2020b). We also compared the performance of our algorithm to TlgGER, the most widely cited tool for this task.

Upstream analysis, including sequence alignment to reference V and J gene libraries and defining the full-length clonotypes, was performed using the tools' recommended pipeline, MiXCR's `analyze` module (Bolotin et al., 2015, <https://mixcr.com/mixcr/reference/mixcr-analyze/>), and Presto (Vander Heiden et al., 2014) and Change-O (Gupta et al., 2015) from Immcantation framework (<https://immcantation.readthedocs.io>). For further details please see the Methods section. For the alignment step and V and J gene annotation we used a custom minimalistic gene set library with only one allelic variant per V and J gene, derived from a custom public genome reference to match the one used for the long-read based genotyping (Rodriguez et al. 2020). Then we performed allele variant inference and genotyping with both tools for all datasets containing more than 3,000 unique full-length clonotypes and compared the resulting individualized V and J gene libraries with the accurate genotype inferred with the next generation long-read sequencing (Rodriguez et al. 2023), comparing nucleotide sequences of the genes. We also excluded poorly expressed allelic variants as determined by aligning the reads to the individualized gene reference libraries. Thus, in our benchmarking we focused on the question of detection of particular V or J gene allele sequences in the participants' AIRR-seq *IGH* data for the subsequent accurate clonotype annotation, which is crucial for many downstream applications of such data (e.g. lineage trees analysis). Importantly, we compared the abilities of both approaches using the sparse reference libraries and AIRR-seq data, containing varying amounts of errors and sequencing noise, including datasets incorporating unique molecular identifiers

and not. Therefore, we consider our benchmarking relevant to real world applications where the data quality is typically far from ideal in many aspects.

MiXCR on average detected 98% of the allelic variants of the V genes supported by the long-read based genotyping, while TlgGER detected 81% of the V gene alleles (**Fig. 1A**). MiXCR produced on average 1 allele call not supported by long-read based genotyping, while TlgGER yielded 2 potential false positive calls (**Fig. 1B**). The recall of TlgGER improved up to 94% on average when the upstream analysis and allele inference was performed utilizing the full built-in reference library containing all of the known alleles (**Supplemental Fig. S1A**). However, the number of the allele calls not supported by the long-read sequencing also increased, up to 5 potential false-positive calls on average. (**Supplemental Fig. S1B**). We assume that the TlgGER algorithm may exhibit improved performance when utilizing the full reference library, due to the algorithm's inherent design. Even with the most recent implementation utilizing the dynamic window as described by Gadala-Maria et al. (2019), there may be instances where a novel allele could be obfuscated by another novel allele that is more similar to the one in the minimal reference library. For MiXCR, transition to full reference library resulted in only minor changes in performance (**Supplemental Fig. S1A, S1B**).

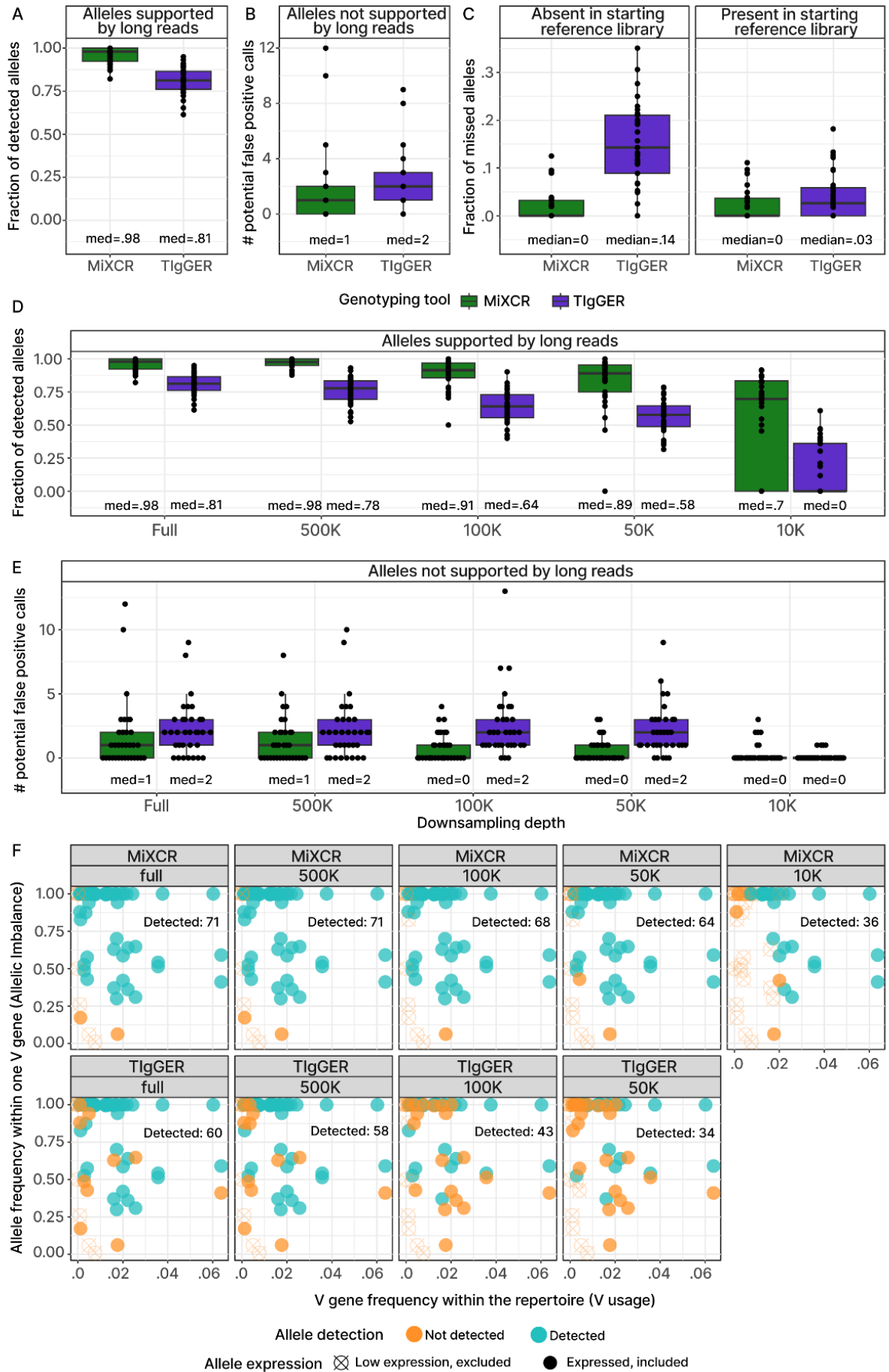
The difference in the number of called alleles between the two algorithms was also apparent when we compared rates of detection of the *de novo* inferred alleles. TlgGER did not detect on average 14% of alleles absent in the starting reference gene library, while MiXCR missed none of the alleles (**Fig. 1C**).

To test the sensitivity of the approaches we also downsampled the dataset to 500,000, 100,000, 50,000 and 10,000 raw sequencing reads. MiXCR allele detection rates decreased by 9 percentage points down to 89% on average when

downsampled to 50,000 reads, which is more than 10x downsampling for all of the datasets. TIgGER detection rates also deteriorated by 23 percentage points, detecting on average 58% of alleles with 50,000 reads. At the extreme level of downsampling by 10,000 sequencing reads MiXCR was able to detect 70% of alleles on average, while TIgGER yielded an error for 21 of the samples due to the low number of clones assigned to any of the V genes (**Fig. 1D**). The number of potential false-positive calls did not increase with lower downsampling depth. For MiXCR, there was a slight decrease in false positives at each downsampling step. Expectedly, both tools produced no potential false-positive calls at extreme downsampling depths (**Fig. 1E**). For MiXCR, the detection of the alleles clearly depended on the two variables - the frequency of the V gene in a particular repertoire and the imbalance in usage between different alleles for a particular V gene. For TIgGER, these parameters appeared to have little influence on detection rates (**Fig. 1F**). Sequencing quality influenced inference for both tools during upstream processing. TIgGER filters sequences based on average Phred quality scores, while MiXCR uses an adjustable threshold for each position. MiXCR's stringent default criteria resulted in no alleles being recovered for one sample at a 50,000 reads downsampling depth and for several more samples at 10,000 reads due to extensive read filtering. For the task of detecting J gene allelic variants, for which could not be performed with TIgGER, MiXCR yielded 100% sensitivity and specificity even with the datasets downsampled to 50,000 reads (**Supplemental Fig. S2**).

Furthermore, we compared the runtime of both tools and found that, on average, there was no significant difference. However, the runtime variability for TIgGER was considerably greater (**Supplemental Fig. S3**).

To assess the influence of sequencing and PCR errors on MiXCR inference performance, we compared results from data generated with and without unique molecular identifiers (**Supplemental Fig. S4 A, B**), which are known to eliminate these errors (Shugay et al., 2014). Additionally, we evaluated the impact of somatic hypermutation (SHM) load on allele inference performance (**Supplemental Fig. S4 C, D**). The number of potential false positive calls was unaffected by SHM frequency or by the presence of sequencing errors in data generated without using unique molecular identifiers (**Supplemental Fig. S4 B, D**). However, the fraction of false negative calls was influenced by both parameters (**Supplemental Fig. S4 A, C**).



**Figure 1. Detection of allelic variants of V genes by inference tools.** a, Fraction of allele calls supported by long-read based genotyping. b, Number of allele calls not supported by long-read based genotyping. c, Fraction of alleles, missed by MiXCR or TigGER, by presence in the initial reference library. d, e Sensitivity and specificity testing by downsampling each sample in the benchmarking dataset by 500,000, 100,000, 50,000 or 10,000 reads. d, fraction of identified allele calls supported by long-read based genotyping. e, number of identified allele calls not supported by long-read based genotyping. f, detection of the allele variants of V genes depending on V usage and allelic imbalance. Each dot represents a V gene allele present in the donor's genotype confirmed by long-read sequencing. The upper row represents detection by the developed algorithm; the lower, allele detection by the comparison tool TigGER. Columns represent different depths of downsampling by number of aligned reads, from right to left: full set of reads, 500,000, 100,000, 50,000, 10,000. V gene, and allele frequencies for each facet were calculated using the full set of reads and allele-resolved V and J gene reference library. Alleles excluded due to low expression (<10 clonotypes), are represented as empty crossed points. N=33 for a - e panels.

## Numerous *IGH*, *TRA* and *TRB* novel alleles detected using MiXCR allele inference

To investigate allelic diversity in human populations we applied the developed algorithm to a large collections of *IGH* (450 individuals) and full-length *TRA* and *TRB* (134 individuals) AIRR-seq datasets. The MiXCR allele inference and genotyping pipeline resulted in identification of both known and previously undocumented alleles, 384 *IGHV*, 128 *TRAV*, 144 *TRBV*, 14 *IGHJ*, 64 *TRAJ*, and 14 *TRBJ* in total.

Numerous previously undocumented alleles, absent from major databases mentioned above (OGRDB or IMGT), were detected: 183 *IGHV*, 33 *TRAV*, 7 *TRAJ*, 41 *TRBV* (**Fig. 2 A-D, F**). Of note, we did not detect any novel variant for any of the *IGHJ* and *TRBJ* genes (**Fig. 2E**). All of the novel alleles sequences were contributed to the public database of allelic variants and are available for download at <https://vdj.online/library>.

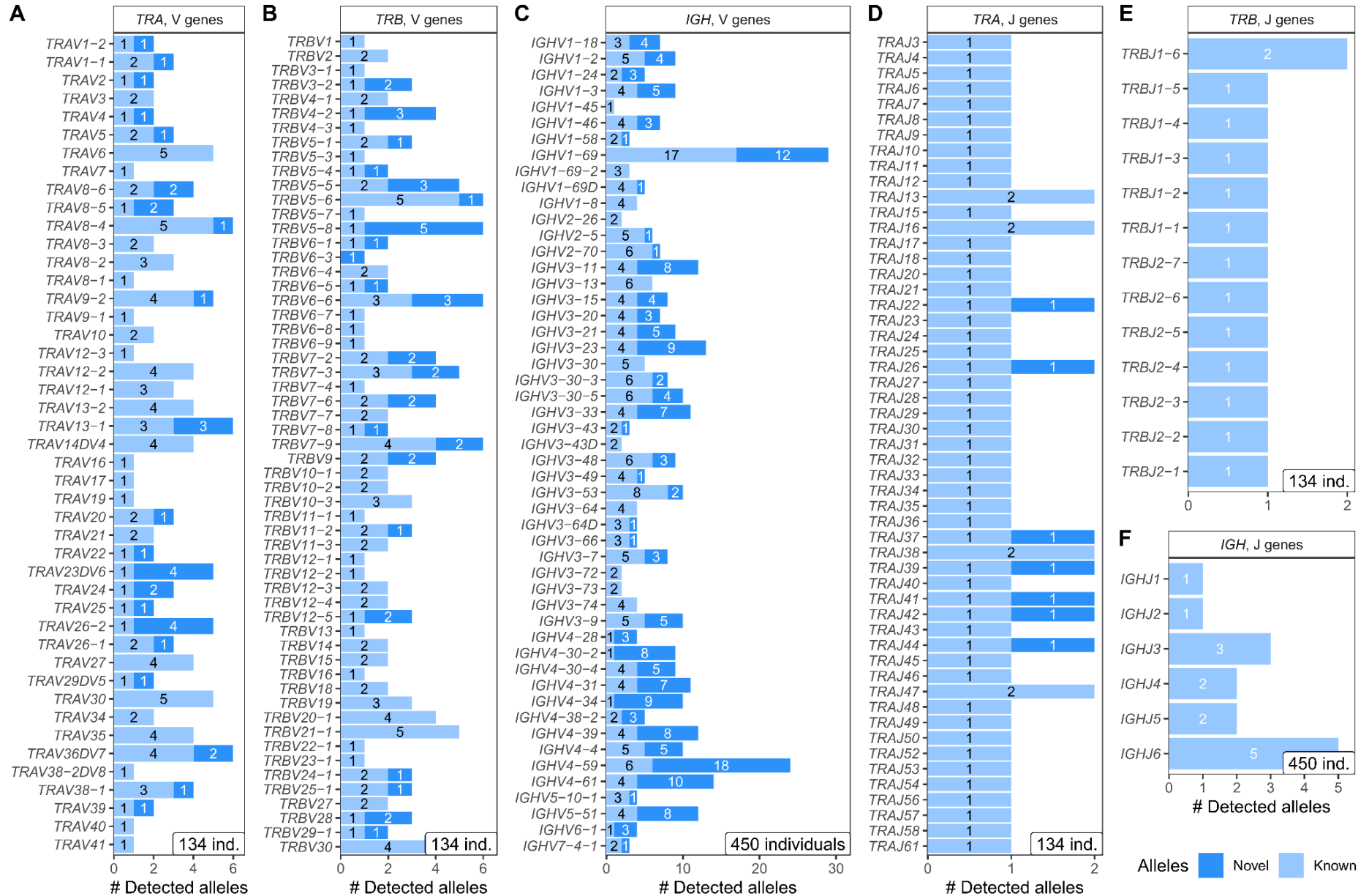
## Divergent allele frequency distribution in *IGHV* genes in African population

The considered *IGH* AIRR-seq datasets included repertoires from African, Asian, European, and Hispanic/Latino individuals (**Fig. 3A**). We did not observe significant differences in the number of detected novel *IGHV* alleles per donor between these groups (**Fig. 3B**). The sufficient sample sizes in the European and African populations allowed us to investigate differences in the number of *IGHV* and *IGHJ* alleles and allelic distributions between these two groups. The number of detected alleles was similar for all of the J genes (**Fig. 3C**) and most of the V genes with the exception of *IGHV1-3*, *IGHV1-69*, *IGHV3-53* and *IGHV4-30-2* (**Fig. 3D**).

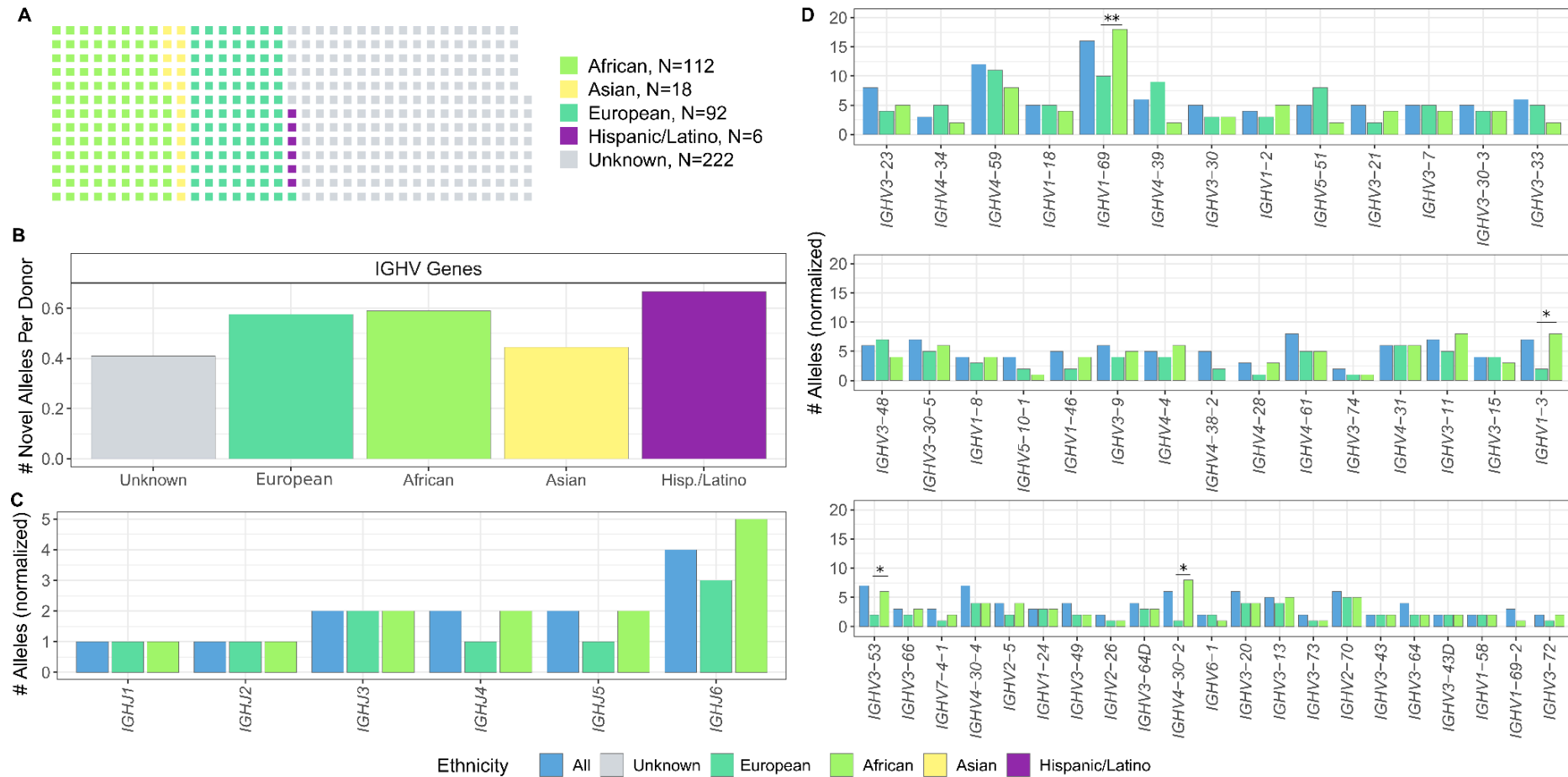
On the other hand, the allele frequency distribution in the African population was significantly different than that of other populations for 38 *IGHV* genes (**Supplemental Fig. S5,6**). Some V genes, even those with similar frequencies in a typical human repertoire, showed very distinct allele distributions. For example, several alleles of *IGHV1-69* and *IGHV3-48* appear at intermediate frequencies in the populations studied. In contrast, *IGHV3-23* had only one predominantly represented allele, while *IGHV3-7* had several alleles at the level of nucleotide sequence differences, but encoding the same amino acid sequence (**Fig. 4**). For these V genes, where the alleles were more evenly distributed, the allele distributions also showed greater differences between ethnic groups (**Fig. 4, Supplemental Fig. S5,6**).

The same difference was also observed for two of the *IGHJ* genes: *IGHJ3* and *IGHJ6* (**Supplemental Fig. S7**). In *TRBV* and *TRAV* loci for most of the gene frequencies, distributions were heavily skewed towards particular single allele variants (**Supplemental Fig. S8, S9**), which may be attributed to a more homogeneous cohort composition by ethnicity, with the predominant majority of participants being of European descent.





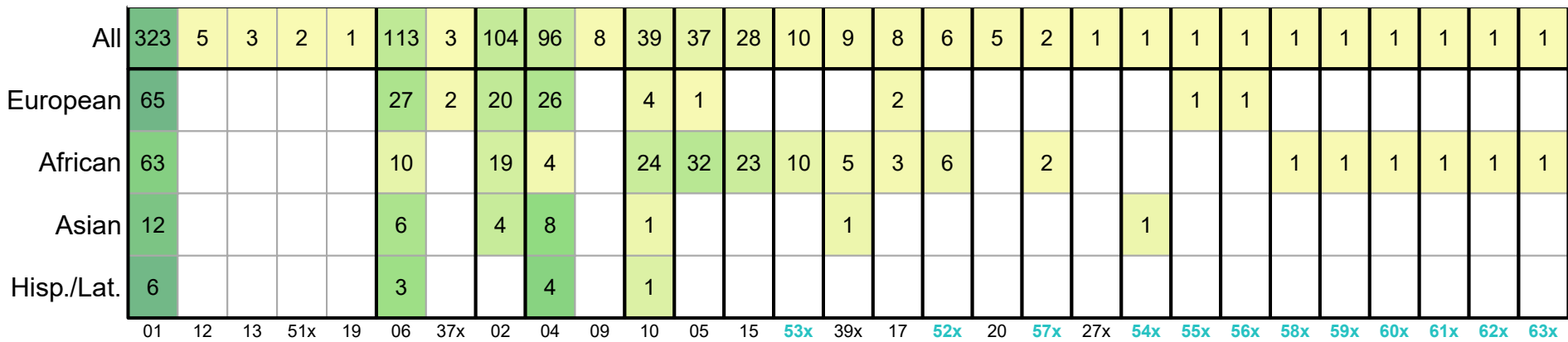
**Figure 2. Number of observed novel and known alleles.** a, *TRAV*. b, *TRBV*. c, *IGHV*. d, *TRAJ*. e, *TRBJ*. f, *IGHJ*



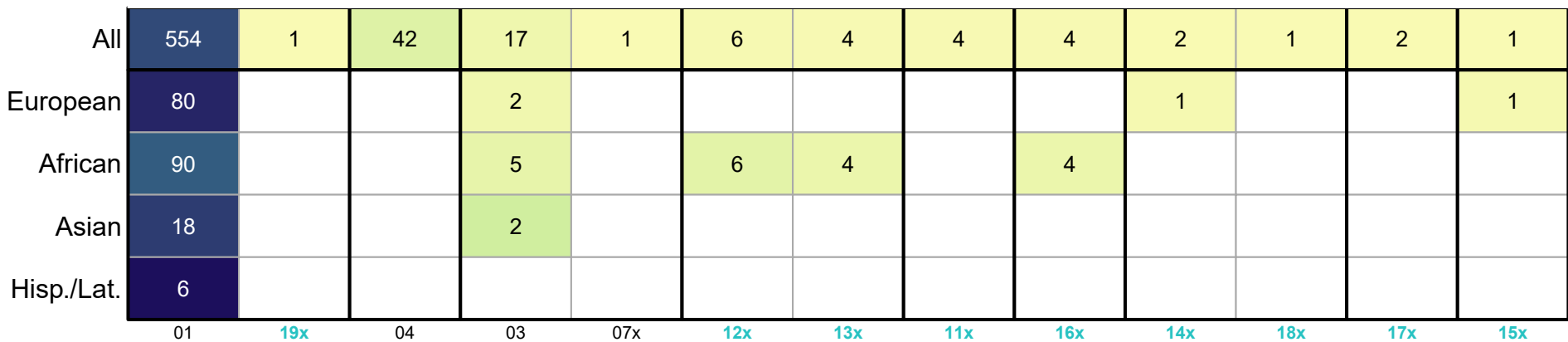
**Figure 3. *IGHV* and *IGHJ* allelic diversity by major ethnic groups.** a, Cohort composition. b, Number of detected novel alleles, normalized per number of individuals. c, Total number of detected alleles by *IGHJ* gene in European, African and general population, normalized by downsampling to a fixed number of individuals (N=92). d, Total number of detected alleles by *IGHV* gene in European, African and general population, normalized by downsampling to a fixed number of individuals (N=92). Comparison

between ethnicities in each b,c, d was performed using permutation test (1000 permutations, \* =  $p \leq 0.05$ , \*\* =  $p \leq 0.01$ , non-significant not shown)

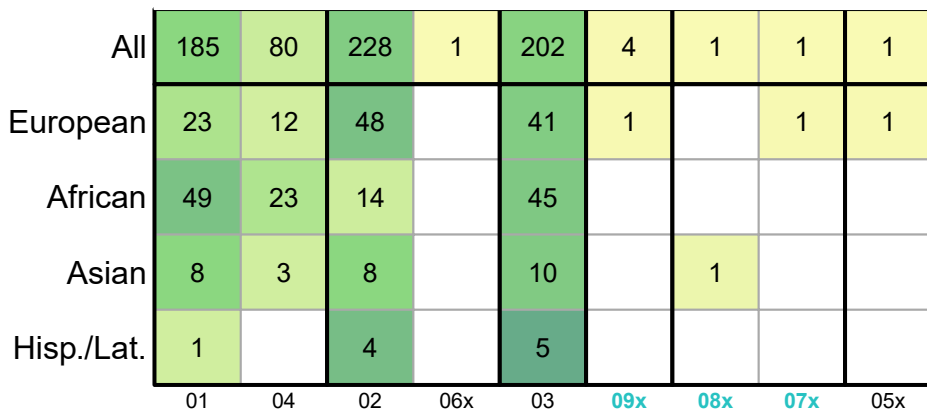
### IGHV1-69



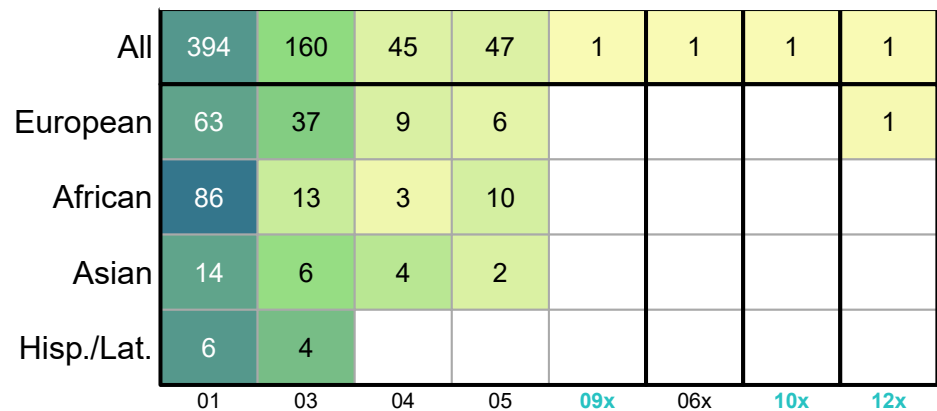
### IGHV3-23



### IGHV3-48



### IGHV3-7



**Figure 4. *IGHV* gene allele frequencies in major ethnic groups for selected *IGHV* genes.** Each column in heatmaps represents a particular allele; numbers for novel alleles, first reported in this study, are colored in green; the letter 'x' designates alleles inferred from AIRR-seq data, either in this study or previously, with the same sequences already present in OGRDB; bold lines separate groups of alleles with different amino acid sequences; groups of alleles with the same amino acid sequence are ordered by the aggregated frequencies of alleles; alleles within groups are order by allele frequency in the general population. Color represents the allele frequency within the ethnic group; numbers in cells represent the number of occurrences of the corresponding allele.

## Discussion

Immune receptor repertoire sequencing datasets have become a valuable source of information for studying immune responses across different health conditions, tissues and cell subsets. Recently developed specialized algorithms (Gadala-Maria et al., 2019; M. M. Corcoran et al., 2016; Zhang et al., 2016; Ralph & Matsen, 2019) allow inference of allelic variants of V and J genes of adaptive immune receptors from AIRR-seq data, scaling up the process of novel allele discovery and allowing AIRR-seq data analysis using individualized gene reference libraries, which significantly increases the accuracy and quality of many types of downstream repertoire analyses. However, all of the current approaches demand high sequencing depth and a significant number of unique receptor sequences for the analysis. Moreover, many prior approaches do not allow allelic inference from both hypermutated and non-hypermutated repertoires. The most comprehensive approach for precise genotyping and allelic inference, utilizing long-read sequencing of the immune receptor gene loci (Gibson et al., 2023; Rodriguez et al. 2023; Rodriguez et al., 2020b), has the greatest accuracy, and also allows for the investigation of structural variants. Although being the most desirable and accurate way to obtain donor-specific V and J gene genotypes, this methodology is costly and requires special experimental procedures. Here, we addressed this unmet need by developing an alternative approach for inferring allelic variants of V and J genes directly from AIRR-seq data, offering improved sensitivity and accuracy compared to existing tools.

Our method allows for successful allelic inference from datasets downsampled to as few as 50,000 sequencing reads. Moreover, the algorithm applicability is not

restricted to a particular type of AIRR-seq data; it can be applied to both repertoires containing hypermutated sequences (e.g., *IGH* repertoires generated from any isotype) as well as datasets containing only non-hypermutated sequences (e.g. TCR-repertoires). Multiple filtering steps integrated into our pipeline prevent false-positive polymorphism calling which typically arises due to the presence of hot-spot hypermutations and PCR and sequencing errors. Furthermore, we demonstrate high sensitivity and specificity of the approach utilizing a very sparse starting reference gene library, containing only one allelic variant per gene, which makes it even more useful for studying allelic diversity in non-model species for which V and J gene reference libraries are incomplete and lack allelic variants. We assume that the improved sensitivity compared to the comparison tool is due to our algorithm bypassing the regression modeling component, which requires a certain number of sequences for a specific V gene to reliably infer alleles. The developed approach is integrated within the MiXCR (Bolotin et al., 2015, <https://mixcr.com>) pipeline for immune-repertoire analysis, and allows seamless allelic inference and re-aligning repertoires to a personalized reference library.

Applying the developed approach to large collections of *IGH*, *TRA* and *TRB* repertoire datasets, we were able to identify a large number of previously undocumented V and J gene alleles. The number of novel *IGHV* alleles, normalized per donor, did not significantly differ among the different population groups. This finding suggests that the genetic diversity of *IGHV* genes even in the relatively better-studied European populations is still not fully characterized. Each additional sampling in the different population groups we studied continues to reveal novel alleles at similar rates. To facilitate sharing and usage of the discovered allele sequences we have established a



database of allelic variants integrated with MiXCR and publicly available at <https://vdj.online/library>.

Differences in allele frequency distributions may have major implications for susceptibility of different populations to diseases and vaccination outcomes (Avnir et al., 2016). Large sample sizes (450 individuals for IGH and 134 for TRA/TRB) allowed us to estimate allele frequencies for most of the studied genes in the population. For IGHV and J gene allelic variants we identify striking differences in allele frequency distributions between African donors and other major population groups. We also contributed the information on V and J gene allele frequencies to VDJ.online, making it a valuable public resource of such information. Having incorporated this database of allelic variants into the MiXCR platform, we hope that it will facilitate further advancement in the immune repertoire analysis field, adding the dimension of allele analysis with little additional effort and cost to many further studies.

# Methods

## Allele variants detection algorithm

The algorithm utilizes alignment and clonotype assembly information from the upstream AIRR-seq data processing, specifically mutation calls from reference V and J gene reference library for BCR or TCR clonotypes and V and J gene annotations, readily available after running the ‘analyze’ command in the MiXCR software (Bolotin et al. 2015). The clonotype definition for the purpose of allele inference may vary depending on the region covered by sequencing.

Using these defined sets of mutations which differentiate the particular clonotype sequences from the corresponding reference V or J gene, the algorithm then separately infers alleles for V and J genes. For simplicity, we describe the algorithm steps for V genes only, the J gene inference follows the same logic:

1. Clonotypes are grouped by the V genes. For the data without unique molecular barcodes only clonotypes with read count greater than one are utilized for subsequent analysis.
2. For each mutation within the group, including insertions and deletions, we define a set of clonotypes which contain this mutation.
3. The mutations are filtered based on the lower diversity bound, estimated as the number of unique combinations of J genes and CDR3-lengths of clonotypes containing that mutation. The mutations that don't exceed a predefined threshold for the value are removed from each of the clonotype's mutation sets.

4. Clonotypes are grouped by filtered mutation sets, including “empty” mutation sets, containing no mutations. The lower diversity bound is calculated for each of the groups as described above. Additionally, the number of clonotypes containing no mutations in J gene after filtering as described in step 3 is calculated. Mutation sets are then filtered by thresholds of these two parameters, resulting in a list of allele candidates.
5. Clonotypes are then assigned to the closest allele candidates. Clonotypes which can not be unambiguously assigned are filtered out. Lower bound of naive diversity is calculated as the number of unique combinations of J genes and CDR3-lengths for clonotypes with unmutated J gene sequences. Candidates are sorted by the score which represents the weighted sum of the lower bound of diversity and lower bound of naive diversity, calculated as described above. Formula for the score:

$$score = D_{all} + 2 \cdot D_{naiveByJgene}$$

Where  $D_{all}$  is the lower bound of diversity for all clonotypes;  $D_{naiveByJgene}$  - lower bound diversity, calculated only for clonotypes with no mutations in J gene.

6. Candidates with the score not lower than 0.35 of the maximum score are then selected for the subject-specific gene set library.
7. Mutations at germline-encoded positions in CDR3 are recovered using the non-mutated clonotypes, which totally match the inferred variants by the rest of the sequence excluding CDR3. Each position is considered if it has at least 5 clonotypes covering it and 70% nucleotide concordance. The right-most

position in CDR3, which meets these criteria, is reported by MiXCR (`reliableRegion` field in tabular output). The rest of CDR3 is picked from the closest allele in the database.

The process for inferring J gene alleles is the same, however the initial grouping is performed by J genes and V genes are used for all of the filtering steps.

This stepwise approach based sequential filtering first on the level of individual mutation and then on the level of mutation sets reduces noise introduced by SHM and sequencing and PCR-errors. The threshold of 0.35 for the final allele filtering was initially chosen from a theoretical consideration of possible distributions of expressed alleles for a V gene allowing the presence of three allelic variants due to possible V gene duplications. This was then corroborated by examining empirical score distributions for alleles in sequencing of *IGH* repertoire of a healthy donor with known genotype; in this case the donor was different from the one in the benchmarking of the algorithm.

In case of a significant difference between the reference library and a particular individual's genotype, the algorithm repeats the steps described above twice, with two different sets of parameters. The first step generates preliminary allele calls, which allows more precise estimation of numbers of clonotypes with unmutated gene sequences.

The algorithm always utilizes only one allele variant per gene as starting reference, preventing potential biases towards particular known sequences. In case there is a weak signal in a particular gene (usually represented by less than 20 clonotypes), the algorithm falls back to assigning one of the known alleles.

Finally, for all of the allele calls, the allele names are looked up in a reference database (the same as available at <https://vdj.online/library>) by exact match of nucleotide sequence. If there is no match the new name is derived from concatenation of the closest allele and sequence hash.

The described algorithm is integrated into MiXCR as the `findAlleles` command.

### **Data collection and repertoire sequencing**

For the benchmarking purposes we utilized *IGH* repertoire sequencing data, accompanied by a targeted long-read *IGH* locus sequencing from Rodriguez et al. 2023, selecting samples which had at least 500,000 sequencing (N=40), which was necessary for compatibility in downsampling experiments. *IGH* locus assembly and variant detection characterizing novel alleles were performed using iGenotyper (<https://github.com/oscarlr/IGenotyper>, Rodriguez et al. 2020b) as previously described (Rodriguez et al. 2023). Briefly, IGenotyper utilizes BLASR (Chaisson, M. J., & Tesler, G., 2012), WhatsHap (Martin et al., 2023), MsPAC (Rodriguez et al., 2020a), and Canu (Koren et al., 2017) for read alignment, calling and phasing single nucleotide variants, phasing reads, and assembling phased reads, respectively. Using the genotypes generated with IGenotyper, we constructed reference allele libraries in FASTA format for each participant. These libraries were then used to match with AIRR-seq derived allele sequences in the benchmarking.

For calculating population allele frequencies we used publicly available *IGH* AIRR-seq data from 6 published studies (total N=450) (Gidoni et al. 2019; Nielsen et al. 2020; Nielsen et al. 2019; Roskin et al. 2015; Davis et al. 2019; Rodriguez et al. 2022).

For generating high-quality full-length TCR repertoires, peripheral blood was collected from 134 individuals without major chronic immunological conditions at CHU of Liège, including COVID-19 patients and individuals after vaccination. 2.5 mL of blood was collected on PAXGene RNA tubes from each participant and stored at -80°C until use, RNA was extracted using the PAXgene Blood RNA Kit (Qiagen). cDNA libraries were generated using SMARTer Human TCR a/b Profiling Kit v2 (Takara Bio USA, San Jose, California, USA). Briefly, a rapid amplification of cDNA ends (RACE) approach with a template-switch effect was used to introduce 5' adaptors during cDNA synthesis. cDNA corresponding to *TRA* and *TRB* transcripts was further amplified and prepared for sequencing, which was performed on a MiSeq instrument with paired-end 2×300 bp reads using the NovaSeq 6000 SP Reagent Kit v1.5 (500 cycles) (Illumina, San Diego, California, USA). The protocol was approved by the ethics committee of Liège University Hospital (approval numbers 2021-54 and 2020/107).

### **Benchmarking of allele variants detection and genotyping**

Processing of the AIRR-seq data was performed using MiXCR v4.4.0 (<https://mixcr.com>, Bolotin et al. 2015) upstream pipeline 'analyze' command, parallelized using GNU Parallel (Tange 2018). Importantly, for the alignment step and V and J gene annotation we used a custom minimalistic gene set library with only one allelic variant per V and J gene, derived from a custom public genome reference to match the one used for the long-read assembly (Rodriguez et al. 2023). After processing we excluded samples with the resulting number of full-length *IGH* clonotypes less than 3,000 (N=7), which probably related to samples either with low cell counts or with low RNA yield. Then the allelic variants were inferred and

individual genotypes were reconstructed for each individual sample (N=33) with the algorithm described above integrated into MiXCR pipeline as `findAlleles` command.

To infer the alleles with the comparison tool, TlgGER (Gadala-Maria et al., 2019) we used the same set of samples (N=33). For initial AIRR-seq data processing we utilized tools pRESTO (Vander Heiden et al., 2014) and Change-O (Gupta et al., 2015), which are the part of the Immcantation framework along with TlgGER (<https://immcantation.readthedocs.io>), using commands and settings, recommended by the documentation. TlgGER v1.0.1 functions `findNovelAlleles` and `inferGenotype` were used for inferring novel alleles and reconstructing genotypes.

To test the sensitivity of the approaches we downsampled the dataset to 500,000, 100,000, 50,000 and 10,000 raw sequencing reads using `seqtk` (<https://github.com/lh3/seqtk>) v1.3 and applied the same upstream processing, allele inference and genotyping pipelines as for the full datasets.

The resulting sets of allele sequences were exported from both tools in FASTA format and matched with the sequences of the alleles present in the genotype of the donor, previously recovered with iGenotyper (Rodriguez et al. 2023), and the number of matches was determined. Comparison was performed on the sequences remaining after removing primers, 5' untranslated regions, and leader sequences. Importantly, due to the fact that *IGH* repertoire sequencing data utilized for comparison was derived using RNA-based technology, inference could be performed only for expressed V and J gene alleles. Thus, we excluded non-functional alleles and also those alleles from comparison which had less than 10 total clonotypes or

less than 3 “naive” clonotypes with no mutation calls assigned to these alleles, when utilizing the same MiXCR v4.3.2 upstream pipeline, but with the individual allele-resolved V and J gene reference libraries constructed from long-read based genotypes. Also we have excluded from comparison genes which were not captured by the long-read sequencing. In particular, *IGHJ* genes were covered only for 9 of 33 considered individuals. For the benchmarking purposes, we excluded alleles for low abundance genes with too low abundance, as defined by each of the tools. TlgGER could not infer novel alleles for genes with less than 50 clonotypes assigned to it, so we excluded such alleles from comparison for TlgGER. MiXCR reports the low abundance genes for which the analysis is impossible with the parameters described above. The average number of clonotypes assigned to those genes was less than 10, we excluded such alleles from comparison for MiXCR too. Finally, we have not taken into account false negative and false positive polymorphism calls in the whole CDR3 region for TlgGER; for MiXCR we applied stricter criteria and have not considered false negative and false positive polymorphism calls outside `reliableRegion`, defined by the tool as described above. The runtime for each of the samples was benchmarked with the R package “bench” v. 1.1.3 (Hester & Vaughan, 2023).

### **Novel allele inference and population frequencies**

Processing of the AIRR-seq data for both BCR and TCR repertoires was performed using MiXCR v4.3.2 `analyze` command. Repertoires containing less than 3000 unique clonotypes were not used for downstream analysis. The algorithm described above for inferring novel alleles and genotyping was used by invoking `findAlleles` MiXCR command under default settings. Alleles lacking designated



names by the International Union of Immunological Societies were given interim names composed of a number continuing the existing sequence, with 'x' letter added after the number. Those not found in OGRDB, the most up-to-date database of alleles inferred from AIRR-seq, were labeled as undocumented.

To compare the total number of alleles per V gene, we selected only European and African populations due to their sufficient size. We downsampled these populations to match the number of participants in the smallest group (N=92) and performed a permutation test to statistically validate the findings. For both TCR and immunoglobulin V and J gene alleles the number of haplotypes with these alleles were estimated using output tables from `findAlleles` command, utilizing only those alleles for which inference could be reliably performed as mentioned in the generated reports. Each case where the only one allele per gene was indicated was treated as a gene in homozygous state, thus not taking into account possible deletions of the genes on one of the chromosomes. We also limited our analysis with the genes detected in at least 15% of the donors. Allele frequencies were then calculated by dividing the number of haplotypes for a particular allele by the total number of haplotypes for this gene in the population. For the *IGH* data we also were able to calculate allele frequencies for four major ethnic groups - African, Asian, European, Hispanic/Latino (self-reported by participants, where missing assigned to “unknown”).

To evaluate pairwise similarity between *IGH* allele frequency distributions in different populations, we utilized Hellinger distance (Hellinger, 1909), calculated using the following formula:

$$H(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^k (\sqrt{p_i} - \sqrt{q_i})^2}$$

where  $P$  and  $Q$  represent the distributions of alleles of a particular V or J gene in two populations, and  $p_i$  and  $q_i$  represent frequencies of individual member  $i$  (one particular allele) of total number of alleles for the gene  $k$ .

We utilized a permutation test to statistically validate differences in the number of novel alleles (**Fig. 3B**), and the total number of alleles per gene (**Fig. 3C, D**), and Hellinger distance between allele distributions (**Supplementary Fig. S6**), reshuffling ethnicity labels 1,000 times, and then calculating the fraction of permutations where we observed the same absolute difference (or Hellinger distance) between groups or greater.

### **Software and packages**

All downstream data analyses and visualizations were conducted using R version 4.3.2 (R Core Team, 2023) with the following packages: bench v. 1.1.3 (Hester & Vaughan, 2023), Biostrings v. 2.70.3 (Pagès et al., 2024), ComplexHeatmap v. 2.15.4 (Gu, 2022; Gu, Eils, & Schlesner, 2016), cowplot v. 1.1.3 (Wilke, 2024), fuzzyjoin v. 0.1.6 (Robinson, 2020), ggpubr v. 0.6.0 (Kassambara, 2023), spgs v. 1.0.4 (Hart & Martínez, 2023), tidyverse v. 2.0.0 (Wickham et al., 2019), waffle v. 1.0.2 (Rudis & Gandy, 2023).

## Data access

Sequencing data generated in this study have been deposited in the ArrayExpress database ([www.ebi.ac.uk/arrayexpress](http://www.ebi.ac.uk/arrayexpress)) under accession number E-MTAB-13593.

MiXCR software including `findAlleles` module is free for academic nonprofit research and could be obtained at <https://github.com/milaboratory/mixcr>. Code required to reproduce the work is included as Supplemental Code.

## Competing Interests

AM consulted MiLaboratories Inc. GN, MI, AD, SP, DC, DB are employed by MiLaboratories Inc. SP, DC, DB are co-founders and share-holders of MiLaboratories Inc. SDB has consulted for Regeneron, Sanofi, Novartis, Genentech and Janssen on topics unrelated to this study and owns stock in AbCellera Biologics.

## Acknowledgments

We thank the study participants for their contributions to this research.

*Funding:* AM and SDB were partially supported by NIH/NIAID grants U19AI104209, R01AI127877, R01AI125567, R01AI130398, U19AI167903, U19AI057229.

*Author contributions:* AM, GN, SP and DB conceived the study. GN and AM developed and benchmarked the described algorithm. OLR and CTW performed contig assembly and allele calling from long-read sequencing data. AM, DAO, MI, VS, AD, performed data analysis for the large population cohort. AT performed

experimental work and SR supervised data collection for the TCR cohort. DC, SDB and DB supervised the study. AM wrote the manuscript with inputs from all authors.

## References

- Avnir, Y., Watson, C. T., Glanville, J., Peterson, E. C., Tallarico, A. S., Bennett, A. S., Qin, K., et al. (2016). IGHV1-69 polymorphism modulates anti-influenza antibody repertoires, correlates with IGHV utilization shifts and varies by ethnicity. *Scientific Reports*, *6*, 20842.
- Bolotin, D. A., Poslavsky, S., Mitrophanov, I., Shugay, M., Mamedov, I. Z., Putintseva, E. V., & Chudakov, D. M. (2015). MiXCR: software for comprehensive adaptive immunity profiling. *Nature Methods*, *12*(5), 380–381.
- Chaisson, M. J., & Tesler, G. (2012). Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR): application and theory. *BMC bioinformatics*, *13*, 238. <https://doi.org/10.1186/1471-2105-13-238>
- Corcoran, M., Chernyshev, M., Mandolesi, M., Narang, S., Kaduk, M., Ye, K., Sundling, C., et al. (2023). Archaic humans have contributed to large-scale variation in modern human T cell receptor genes. *Immunity*, *56*(3), 635-652.e6.
- Corcoran, M. M., Phad, G. E., Vázquez Bernat, N., Stahl-Hennig, C., Sumida, N., Persson, M. A. A., Martin, M., et al. (2016). Production of individualized V gene databases reveals high levels of immunoglobulin genetic diversity. *Nature Communications*, *7*, 13642.
- Dekker, J., van Dongen, J. J. M., Reinders, M. J. T., & Khatri, I. (2022). pmTR database: population matched (pm) germline allelic variants of T-cell receptor (TR) loci. *Genes and Immunity*, *23*(2), 99–110.
- Fitipaldi, H., & Franks, P. W. (2023). Ethnic, gender and other sociodemographic biases in genome-wide association studies for the most burdensome

non-communicable diseases: 2005-2022. *Human molecular genetics*, 32(3), 520–532. <https://doi.org/10.1093/hmg/ddac245>

Hart, A., & Martínez, S. (2023). spgs: Statistical patterns in genomic sequences. Retrieved from <https://CRAN.R-project.org/package=spgs>

Hellinger, E. (1909). Neue Begründung der Theorie quadratischer Formen von unendlichvielen Veränderlichen.. *Journal für die reine und angewandte Mathematik*, 1909(136), 210-271. <https://doi.org/10.1515/crll.1909.136.210>

Hester, J., & Vaughan, D. (2023). bench: High precision timing of r expressions. Retrieved from <https://CRAN.R-project.org/package=bench>

Gadala-Maria, D., Gidoni, M., Marquez, S., Vander Heiden, J. A., Kos, J. T., Watson, C. T., O'Connor, K. C., et al. (2019). Identification of Subject-Specific Immunoglobulin Alleles From Expressed Repertoire Sequencing Data. *Frontiers in Immunology*, 10, 129.

Gadala-Maria, D., Yaari, G., Uduman, M., & Kleinstein, S. H. (2015). Automated analysis of high-throughput B-cell sequencing data reveals a high frequency of novel immunoglobulin V gene segment alleles. *Proceedings of the National Academy of Sciences of the United States of America*, 112(8), E862-70.

Gibson, W. S., Rodriguez, O. L., Shields, K., Silver, C. A., Dorgham, A., Emery, M., Deikus, G., Sebra, R., Eichler, E. E., Bashir, A., Smith, M. L., & Watson, C. T. (2023). Characterization of the immunoglobulin lambda chain locus from diverse populations reveals extensive genetic variation. *Genes and immunity*, 24(1), 21–31. <https://doi.org/10.1038/s41435-022-00188-2>

Gidoni, M., Snir, O., Peres, A., Polak, P., Lindeman, I., Mikocziova, I., Sarna, V. K., et al. (2019). Mosaic deletion patterns of the human antibody heavy chain gene locus shown by Bayesian haplotyping. *Nature Communications*, 10(1), 628.

- Gu, Z. (2022). Complex heatmap visualization. *iMeta*. <https://doi.org/10.1002/imt2.43>
- Gu, Z., Eils, R., & Schlesner, M. (2016). Complex heatmaps reveal patterns and correlations in multidimensional genomic data. *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/btw313>
- Gupta, N. T., Vander Heiden, J. A., Uduman, M., Gadala-Maria, D., Yaari, G., & Kleinstein, S. H. (2015). Change-O: a toolkit for analyzing large-scale B cell immunoglobulin repertoire sequencing data. *Bioinformatics*, *31*(20), 3356–3358.
- Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., & Phillippy, A. M. (2017). Canu: scalable and accurate long-read assembly via adaptive *k*-mer weighting and repeat separation. *Genome research*, *27*(5), 722–736. <https://doi.org/10.1101/gr.215087.116>
- Kassambara, A. (2023). ggpubr: “ggplot2” based publication ready plots. Retrieved from <https://CRAN.R-project.org/package=ggpubr>
- Lees, W., Busse, C. E., Corcoran, M., Ohlin, M., Scheepers, C., Matsen, F. A., Yaari, G., et al. (2020). OGRDB: a reference database of inferred immune receptor genes. *Nucleic Acids Research*, *48*(D1), D964–D970.
- Martin, M., Ebert, P., & Marschall, T. (2023). Read-Based Phasing and Analysis of Phased Variants with WhatsHap. *Methods in molecular biology (Clifton, N.J.)*, *2590*, 127–138. [https://doi.org/10.1007/978-1-0716-2819-5\\_8](https://doi.org/10.1007/978-1-0716-2819-5_8)
- Mikocziova, I., Greiff, V., & Sollid, L. M. (2021). Immunoglobulin germline gene variation and its impact on human disease. *Genes and Immunity*, *22*(4), 205–217.
- Ohlin, M., Scheepers, C., Corcoran, M., Lees, W. D., Busse, C. E., Bagnara, D., Thörnqvist, L., et al. (2019). Inferred allelic variants of immunoglobulin receptor genes: A system for their evaluation, documentation, and naming. *Frontiers in*

*Immunology*, 10, 435.

Omer, A., Peres, A., Rodriguez, O. L., Watson, C. T., Lees, W., Polak, P., Collins, A. M., et al. (2022). T cell receptor beta germline variability is revealed by inference from repertoire data. *Genome Medicine*, 14(1), 2.

Pagés, H., Aboyou, P., Gentleman, R., & DebRoy, S. (2024). Biostrings: Efficient manipulation of biological strings. Retrieved from <https://bioconductor.org/packages/Biostrings>

Peres, A., Lees, W. D., Rodriguez, O. L., Lee, N. Y., Polak, P., Hope, R., Kedmi, M., Collins, A. M., Ohlin, M., Kleinstein, S. H., Watson, C. T., & Yaari, G. (2023). IGHV allele similarity clustering improves genotype inference from adaptive immune receptor repertoire sequencing data. *Nucleic acids research*, 51(16), e86. <https://doi.org/10.1093/nar/gkad603>

Pushparaj, P., Nicoletto, A., Sheward, D. J., Das, H., Castro Dopico, X., Perez Vidakovics, L., Hanke, L., et al. (2022). Immunoglobulin germline gene polymorphisms influence the function of SARS-CoV-2 neutralizing antibodies. *Immunity*, 56(1), 193-206.e7.

Ralph, D. K., & Matsen, F. A. (2019). Per-sample immunoglobulin germline inference from B cell receptor deep sequencing data. *PLoS Computational Biology*, 15(7), e1007133.

Robinson, D. (2020). fuzzyjoin: Join tables together on inexact matching. Retrieved from <https://CRAN.R-project.org/package=fuzzyjoin>

Rodriguez, O. L., Ritz, A., Sharp, A. J., & Bashir, A. (2020). MsPAC: a tool for haplotype-phased structural variant detection. *Bioinformatics (Oxford, England)*, 36(3), 922–924. <https://doi.org/10.1093/bioinformatics/btz618>

Rodriguez, Oscar L, Gibson, W. S., Parks, T., Emery, M., Powell, J., Strahl, M.,

- Deikus, G., et al. (2020). A novel framework for characterizing genomic haplotype diversity in the human immunoglobulin heavy chain locus. *Frontiers in Immunology*, *11*, 2136.
- Rodriguez, Oscar L., Silver, C. A., Shields, K., Smith, M. L., & Watson, C. T. (2022). Targeted long-read sequencing facilitates phased diploid assembly and genotyping of the human T cell receptor alpha, delta, and beta loci. *Cell Genomics*, 100228.
- Rodriguez, O. L., Safonova, Y., Silver, C. A., Shields, K., Gibson, W. S., Kos, J. T., Tieri, D., Ke, H., Jackson, K. J. L., Boyd, S. D., Smith, M. L., Marasco, W. A., & Watson, C. T. (2023). Genetic variation in the immunoglobulin heavy chain locus shapes the human antibody repertoire. *Nature communications*, *14*(1), 4419. <https://doi.org/10.1038/s41467-023-40070-x>
- Rudis, B., & Gandy, D. (2023). waffle: Create waffle chart visualizations. Retrieved from <https://CRAN.R-project.org/package=waffle>
- Shugay, M., Britanova, O. V., Merzlyak, E. M., Turchaninova, M. A., Mamedov, I. Z., Tuganbaev, T. R., Bolotin, D. A., Staroverov, D. B., Putintseva, E. V., Plevova, K., Linnemann, C., Shagin, D., Pospisilova, S., Lukyanov, S., Schumacher, T. N., & Chudakov, D. M. (2014). Towards error-free profiling of immune repertoires. *Nature methods*, *11*(6), 653–655. <https://doi.org/10.1038/nmeth.2960>
- Tange, O. (2018). GNU Parallel 2018. B GNU Parallel 2018 (c. 112). Ole Tange. <https://doi.org/10.5281/zenodo.1146014>
- Vander Heiden, J. A., Yaari, G., Uduman, M., Stern, J. N. H., O'Connor, K. C., Hafler, D. A., Vigneault, F., et al. (2014). pRESTO: a toolkit for processing high-throughput sequencing raw reads of lymphocyte receptor repertoires. *Bioinformatics*, *30*(13), 1930–1932.



Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ...

Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686>

Wilke, C. O. (2024). cowplot: Streamlined plot theme and plot annotations for “ggplot2”. Retrieved from <https://CRAN.R-project.org/package=cowplot>

Zhang, W., Wang, I.-M., Wang, C., Lin, L., Chai, X., Wu, J., Bett, A. J., et al. (2016).

IMPre: An Accurate and Efficient Software for Prediction of T- and B-Cell Receptor Germline Genes and Alleles from Rearranged Repertoire Data. *Frontiers in Immunology*, 7, 457.



## Ultrasensitive allele inference from immune repertoire sequencing data with MiXCR

Artem Mikelov, George Nefedev, Aleksandr Tashkeev, et al.

*Genome Res.* published online October 21, 2024

Access the most recent version at doi:[10.1101/gr.278775.123](https://doi.org/10.1101/gr.278775.123)

---

**P<P** Published online October 21, 2024 in advance of the print journal.

**Accepted Manuscript** Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

**Open Access** Freely available online through the *Genome Research* Open Access option.

**Creative Commons License** This manuscript is Open Access. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

**Email Alerting Service** Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

---



---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---