

# A novel next-generation sequencing capture-based strategy to report somatic hypermutation status using genomic regions downstream to immunoglobulin rearrangements

Neil McCafferty,<sup>1,2</sup> James Peter Stewart,<sup>1</sup> Nikos Darzentas,<sup>3</sup> Jana Gazdova,<sup>1</sup> Mark Catherwood,<sup>4</sup> Kostas Stamatopoulos,<sup>5</sup> Anton W. Langerak<sup>6</sup> and David Gonzalez<sup>1</sup>

<sup>1</sup>Patrick G Johnston Centre for Cancer Research, Queens University Belfast, Belfast, UK;

<sup>2</sup>Barts Cancer Institute, Queen Marys University, London, UK; <sup>3</sup>Department of Hematology, University Hospital Schleswig-Holstein, Kiel, Germany; <sup>4</sup>Belfast City Hospital, Belfast, UK;

<sup>5</sup>Institute of Applied Biosciences, Centre for Research and Technology Hellas, Thessaloniki, Greece and <sup>6</sup>Laboratory of Medical Immunology, Department of Immunology, Erasmus University Medical Center, Rotterdam, The Netherlands

**Correspondence:** D. Gonzalez  
[D.GonzalezdeCastro@qub.ac.uk](mailto:D.GonzalezdeCastro@qub.ac.uk)

**Received:** August 15, 2022.

**Accepted:** December 15, 2022.

**Early view:** December 29, 2022.

<https://doi.org/10.3324/haematol.2022.281928>

©2023 Ferrata Storti Foundation

Published under a CC BY-NC license



## **Supplementary Data**

### **Supplementary Methods**

#### **Library preparation and sequencing**

An optimised and standardised library preparation and sequencing protocol was followed according to the EuroClonality-NDC protocol<sup>28</sup>. The additional 61 CLL samples matched the EuroClonality-NDC strategy except 500 ng gDNA was used for library preparation and no pre-hybridisation PCR cycles. These samples were pooled and sequenced on the NextSeq 500 (Illumina, Cambridge, UK) using a 75 bp paired-end strategy on a 150-cycle NextSeq 500/550 Mid Output Kit (Illumina, Cambridge, UK).

#### **Distinction of SHM from non-somatic variants**

Polymorphic variants in IGHJ-E were filtered out using the population database GnomAD (v3) and several minor allele frequency (MAF) filters for the global population (>2%, >1%, >0.2%, >0.1%, >0.02% and >0.01%, Figure S4)<sup>46</sup>. Recurrent polymorphic and artefactual variants were also identified from the 331 tumour cases if they were present above specified thresholds, using several tumour frequency (TF) filters (>50%, >40%, >30%, >20% and >10%, Figure S5). Minimum variant allele frequency (VAF) filters were also examined to account for variants below the limit of detection (>5%, >4%, >3%, >2% and >1%, Figure S6) since these could represent low-level artefactual variants. All remaining SNVs were visually assessed using IGV to ensure variants were somatic mutations (accounting for remaining alignment artefacts, Figure S7). Variant filters were sequentially optimised, and the best performing filter according to ROC/AUC analysis was selected in each analysis for a combined final filtering strategy containing the following parameters: (i) GnomAD MAF>0.02%, (ii) TF>10%, (iii) VAF>5%, (iv) IGV assessment (Figure S3).

## Supplementary Tables

<i>Diagnosed Disease</i>	<i>n (%)</i>	<i>B/T cell</i>
<i>Chronic lymphocytic leukaemia (CLL)<sup>a</sup></i>	<b>98 (29.6)</b>	<b>B</b>
<i>Mantle cell lymphoma (MCL)</i>	<b>34 (10.3)</b>	<b>B</b>
<i>Follicular lymphoma (FL)</i>	<b>34 (10.3)</b>	<b>B</b>
<i>Plasma cell myeloma (PCM)</i>	<b>24 (7.3)</b>	<b>B</b>
<i>Burkitt lymphoma (BL)</i>	<b>21 (6.3)</b>	<b>B</b>
<i>Diffuse large B-cell lymphoma (DLBCL)</i>	<b>21 (6.3)</b>	<b>B</b>
<i>B Acute Lymphoblastic leukaemia (B-ALL)</i>	<b>16 (4.8)</b>	<b>B</b>
<i>Marginal zone lymphoma (MZL)</i>	<b>7 (2.1)</b>	<b>B</b>
<i>Mucosa-Associated Lymphoid Tissue (MALT)</i>	<b>3 (0.9)</b>	<b>B</b>
<i>T Acute Lymphoblastic leukaemia (T-ALL)</i>	<b>40 (12.1)</b>	<b>T</b>
<i>Anaplastic large-cell lymphoma (ALCL)</i>	<b>12 (3.6)</b>	<b>T</b>
<i>CD30+ anaplastic large cell lymphoma (C-ALCL)</i>	<b>1 (0.3)</b>	<b>T</b>
<i>Angioimmunoblastic T-cell lymphoma (AITL)</i>	<b>9 (2.7)</b>	<b>T</b>
<i>Peripheral T-cell lymphoma, NOS (PTCL)</i>	<b>1 (0.3)</b>	<b>T</b>
<i>Intestinal T-cell lymphoma, NOS (ITCL)</i>	<b>1 (0.3)</b>	<b>T</b>
<i>Mycosis fungoides (MF)</i>	<b>3 (0.9)</b>	<b>T</b>
<i>Sézary syndrome (SS)</i>	<b>3 (0.9)</b>	<b>T</b>
<i>T-lymphoblastic lymphoma (T-LBL)</i>	<b>1 (0.3)</b>	<b>T</b>
<i>Enteropathy-associated T-cell lymphoma (EATL)</i>	<b>2 (0.6)</b>	<b>T</b>
<b>Total</b>	<b>331</b>	<b>258 B 73 T</b>

**Supplementary Table 1.** Diagnosed disease types of 331 samples. <sup>a</sup>95/98 CLL samples with Sanger sequencing data and in-frame IGH VDJ rearrangements identified by IMGT. CLL cases made up of 37 EuroClonality-NDC study cases and an additional 61 diagnostic CLL cases.

<i>Sanger Sequencing IGHV</i>	<i>n (%)</i>	<i>Detected by IGV (n)</i>	<i>Detected by ARResT/Interrogate (n)</i>
<b>V1-2</b>	1 (1.05%)	1	1
<b>V1-3</b>	2 (2.11%)	1	2
<b>V1-46</b>	1 (1.05%)	1	1
<b>V1-69</b>	13 (13.68%)	12	13
<b>V1-8<sup>a</sup></b>	3 (3.16%)	0	2
<b>V2-26</b>	1 (1.05%)	0	1
<b>V2-5</b>	3 (3.16%)	1	2
<b>V3-11</b>	2 (2.11%)	0	2
<b>V3-15</b>	2 (2.11%)	0	2
<b>V3-21</b>	9 (9.47%)	6	6
<b>V3-23</b>	7 (7.37%)	4	7
<b>V3-30</b>	7 (7.37%)	3	6
<b>V3-33</b>	3 (3.16%)	2	3
<b>V3-48</b>	6 (6.32%)	5	6
<b>V3-49</b>	1 (1.05%)	1	1
<b>V3-53</b>	1 (1.05%)	1	1
<b>V3-7</b>	5 (5.26%)	4	5
<b>V3-74</b>	4 (4.21%)	2	4
<b>V3-9<sup>b</sup></b>	1 (1.05%)	0	1

<i>Sanger Sequencing IGHV</i>	<i>n (%)</i>	<i>Detected by IGV (n)</i>	<i>Detected by ARResT/Interrogate (n)</i>
<b>V4-31</b>	1 (1.05%)	0	1
<b>V4-34</b>	9 (9.47%)	2	7
<b>V4-38<sup>a</sup></b>	1 (1.05%)	0	1
<b>V4-39</b>	5 (5.26%)	4	4
<b>V4-4</b>	1 (1.05%)	0	1
<b>V4-61</b>	2 (2.11%)	0	2
<b>V5-10</b>	1 (1.05%)	1	1
<b>V5-51</b>	3 (3.16%)	2	3
<b>Total</b>	95	53	86

**Supplementary Table 2.** IGHV genes identified by Sanger Sequencing in 95 CLL samples. IGV concordance: 55.79% (95% CI: 45.23-65.98%). ARResT/Interrogate concordance: 89.47% (95% CI: 81.49-94.84%). Sanger Sequencing SHM data showed 50 M-CLL and 45 U-CLL. <sup>a</sup>IGHV genes not present on GRCh38 or <sup>b</sup>not captured as part of the EuroClonality-NDC panel, and therefore not assessable.

## **Supplementary Figure Legends**

**Supplementary Figure 1.** Comparison of Sanger sequencing IGHV Identity in 95 CLL cases where the reported clonal IGHV gene could be detected by: A) Visual assessment with IGV (genomic alignment-based), where 53/95 (55.8%) cases were concordant with Sanger sequencing. B) ARResT/Interrogate, where 86/95 (90.5%) cases were concordant. Wilcoxon signed-rank: ns (not significant):  $p > 0.05$ ; \*\*\*\*:  $p \leq 0.0001$ .

**Supplementary Figure 2.** Reporting SNVs in IGHV using an NGS genomic alignment-based strategy. IMGT output of Sanger sequences reported on left (SNVs numbered in blue), IGV track of the rearranged IGHV on right (corresponding SNVs numbered with successful variant calling in green and unsuccessful in red). Example 1. Borderline case where the rearranged IGHV was detected by IGV and ARResT/Interrogate: NGS reads are aligned to the genome, 7/7 SNVs are detected. Example 2. Mutated case where the clonally rearranged IGHV was detected by ARResT/Interrogate but not IGV: NGS reads are aligned to the genome, however, only 1/6 shown SNVs are detected (1/17 in total).

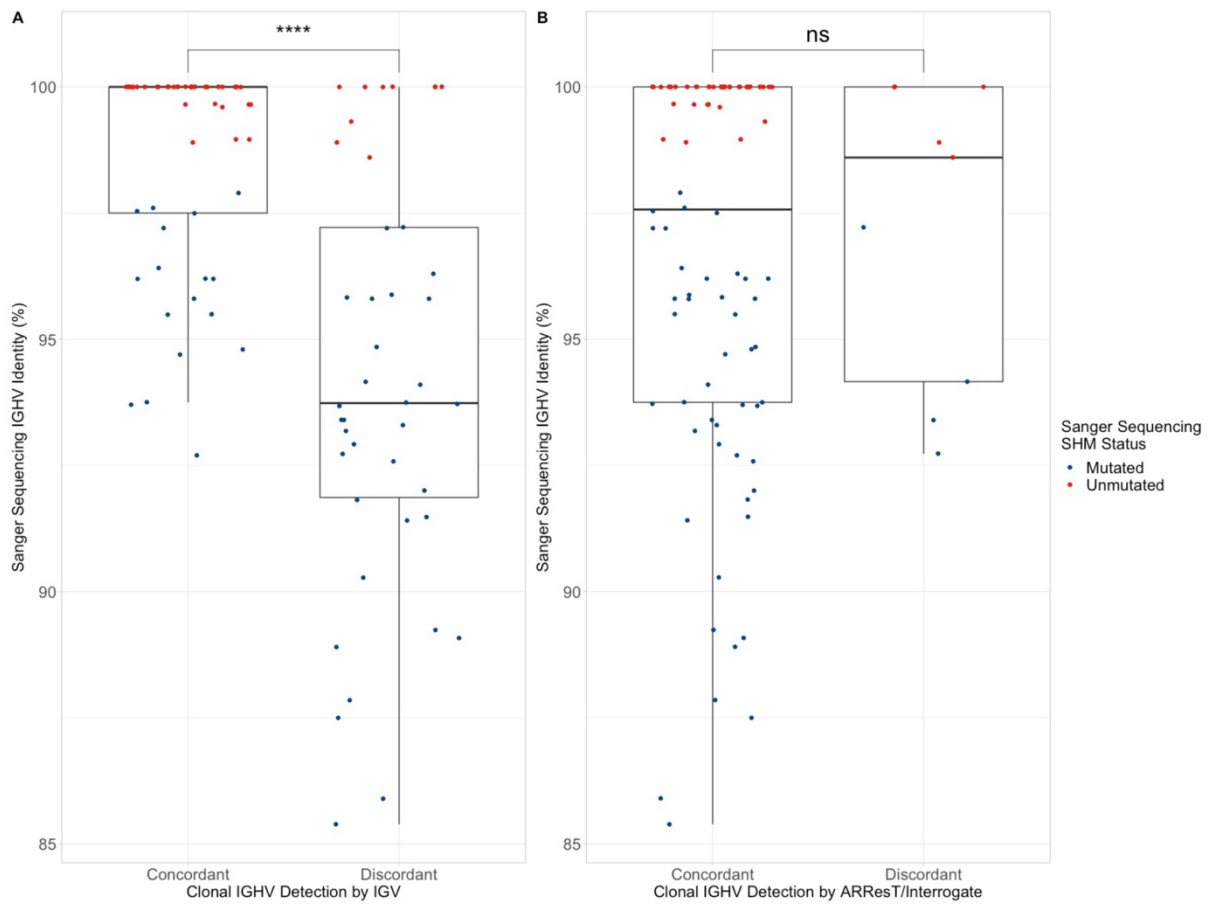
**Supplementary Figure 3.** ROC curves of IGHJ-E with Sanger Sequencing SHM status, selected variant filters and final IGHJ-E following SNV assessment for alignment artefacts. Somatic variant filters were selected sequentially, e.g., the best performing GnomAD filter (MAF > 0.02) was selected and carried through to the following analyses.

**Supplementary Figure 4.** ROC curves of IGHJ-E with Sanger sequencing SHM status. Population variants were selected based on minor allele frequency (MAF) 0.01-2% applied to filter SNPs. Population SNPs were sourced from GnomAD using all available GRCh38 sequences.

**Supplementary Figure 5.** ROC curves of IGHJ-E with Sanger sequencing SHM status. Recurrent SNPs with a tumor frequency (TF) 10-50% were used. A list of tumor SNPs were generated using all available EuroClonality-NDC study and diagnostic CLL cases (n=331)

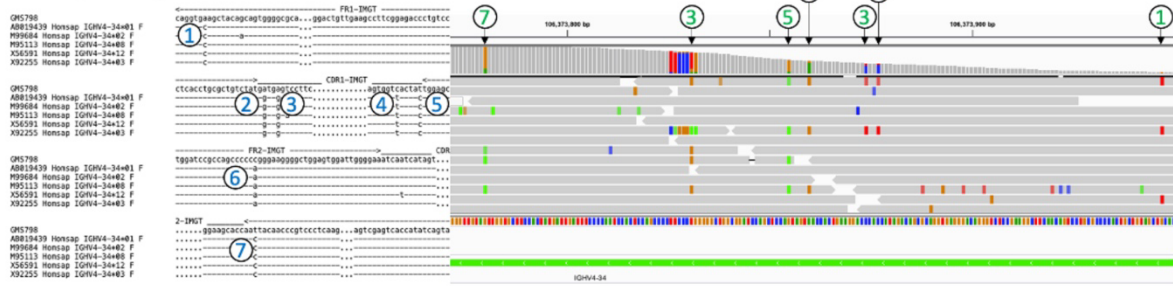
**Supplementary Figure 6.** ROC curves of IGHJ-E with Sanger sequencing SHM status. Minimum variant allele frequencies (VAF) 1-5% were applied to filter SNVs.

**Supplementary Figure 7.** A) IGV screenshot of IGHJ-E analysed region. This sample was reported as unmutated by Sanger sequencing. B) Zoomed region highlights sequencing/alignment artefact resulting in several false somatic variants. Variants were removed from the analysis and IGHJ-E identity was updated.

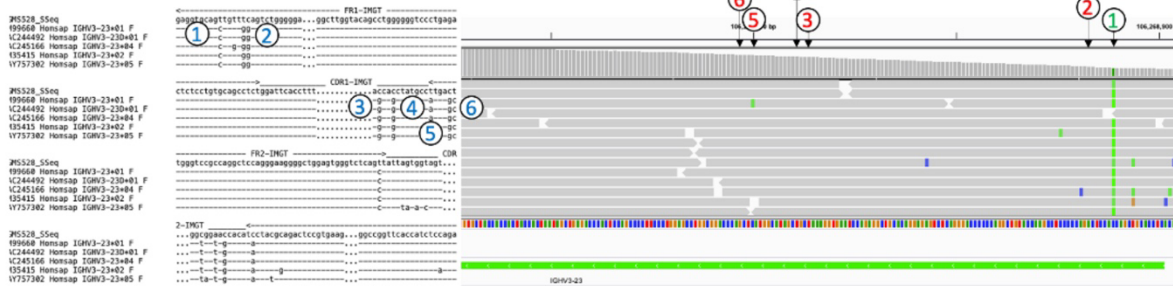


**Supplementary Figure 1.**

Example 1. GMS798 (CLL)  
 SSeq/IMGT: IGHV4-34 (97.54%)  
 ARResT/Interrogate: ✓ IGV: ✓

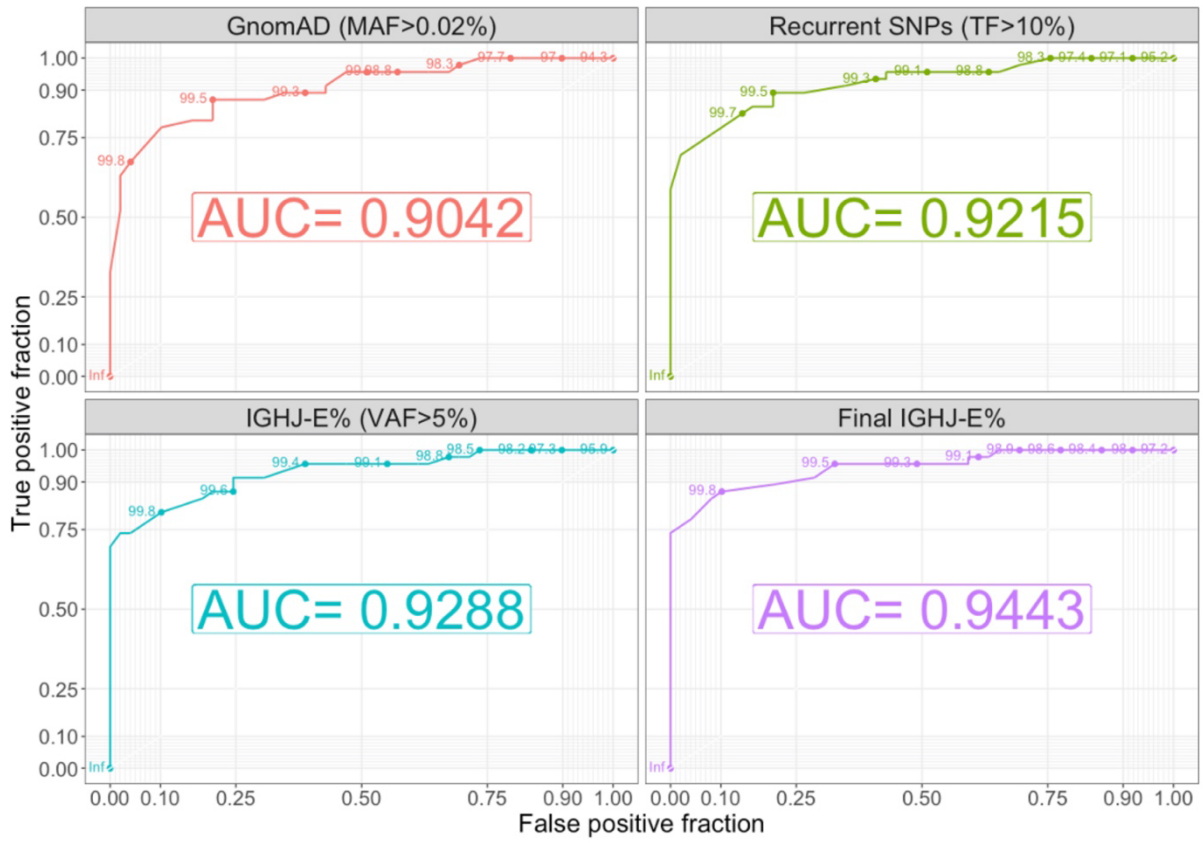


Example 2. GMS528 (CLL)  
 SSeq/IMGT: IGHV3-23 (94.10%)  
 ARResT/Interrogate: ✓ IGV: X

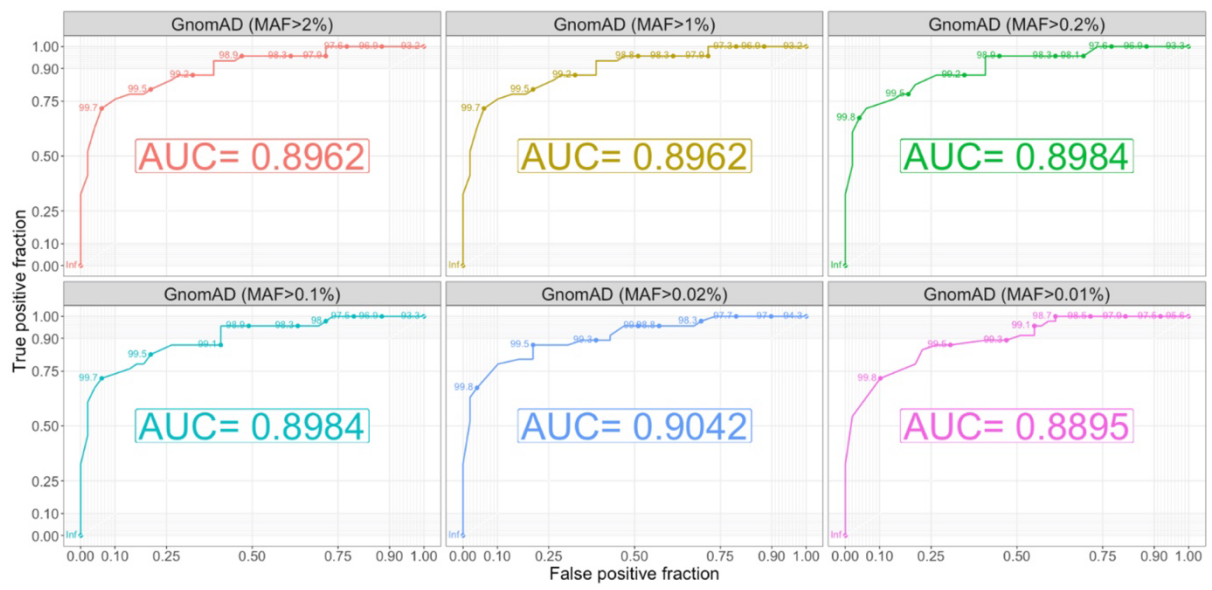


Supplementary Figure 2.

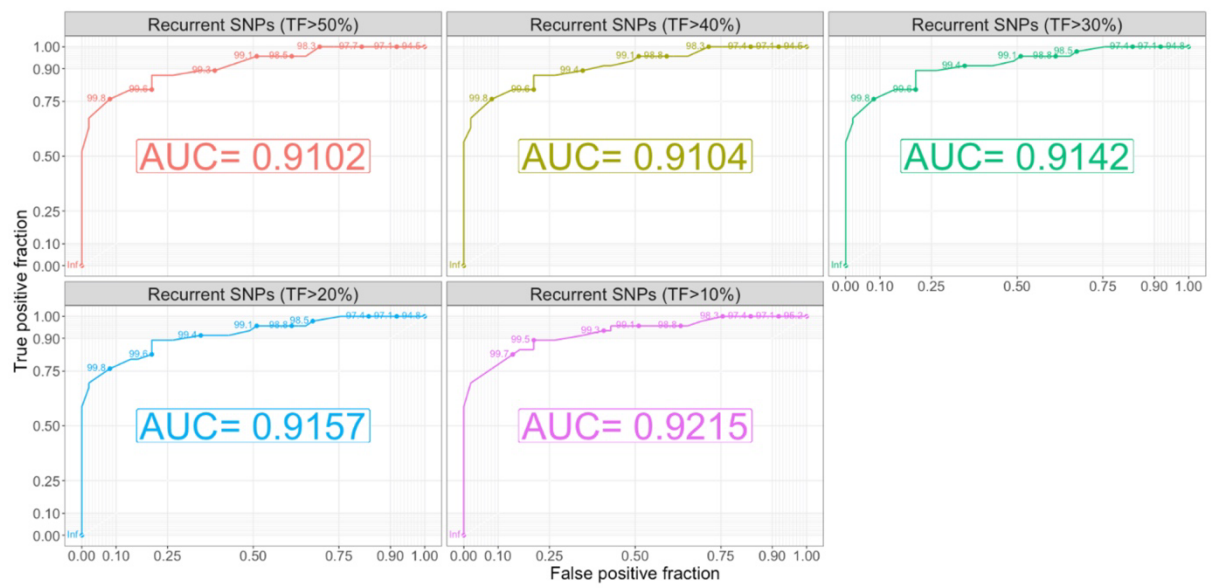




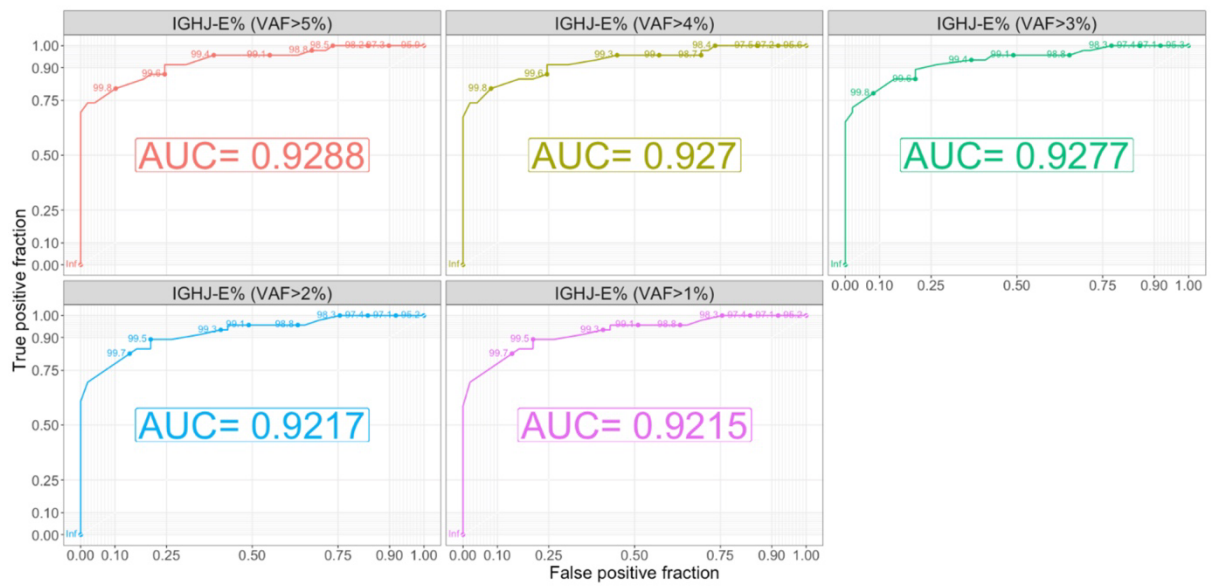
Supplementary Figure 3.



Supplementary Figure 4.

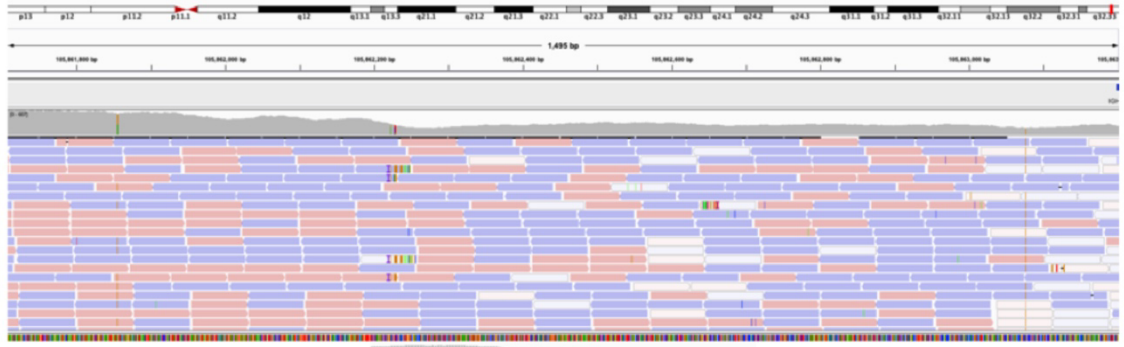


Supplementary Figure 5.

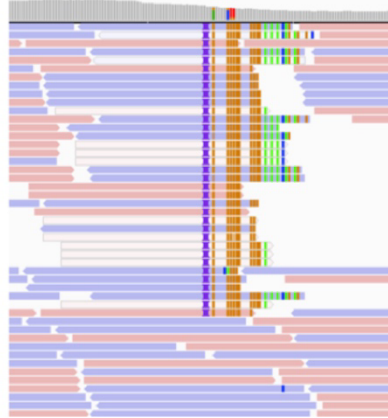


Supplementary Figure 6.

A.



B.



Supplementary Figure 7.